## Artificial Intelligence Techniques for Soil Erosion Mapping and Risk Assessment in Almería Province, Southeast Spain

**Kevin Goldsmith** 

A Thesis Submitted for the Degree of Doctor of Philosophy

School of Earth Sciences and Geography Kingston University

Submitted:

December 2005

KING	STON UNIN	/ERSI	TY L	BRARY
Acc No.	01497	70		PREF
	THESES	PHD	16	
	MOOAL	30	11/	d

## **IMAGING SERVICES NORTH**



ł

Boston Spa, Wetherby West Yorkshire, LS23 7BQ www.bl.uk

# THESIS CONTAINS CD

## Acknowledgements

I would like to take the opportunity to thank a number of people for their support and encouragement over the past four years, without whom this work would not have been completed.

Thanks to Dr. Stuart Downward for his undoubted enthusiasm for the subject, and the advice that he has provided along the way. Dr. Doreen Boyd, for sharing her knowledge and experience of machine vision and data-mining techniques, and Dr. Hazel Faulkner, for her extensive knowledge of the subject and study area, and her passion for geomorphology. I would like to acknowledge Dr. Peter Hooda, for his assistance regarding laboratory methods and techniques, and his proof reading of various chapters.

I would also like to express appreciation to my postgraduate colleagues, without whom I would have lost more hair than is currently the case. In no particular order these include; Simon Kiley, and his uncanny dry wit, Leigh Truelove, with whom I have held many GIS discussions, my namesake Paul Goldsmith and, Carolina Sanchez-Hernandez and her assistance with the sometimes amusing Spanish language. Also, Annabelle Boulay, the coffee morning crew; Mark Thomas and Al Garcia, the gruesome twosome, Paul Helps and John 'where your ashes gone' Groves, John Lignum, Bev Coldwell, Jenny Halliday, and finally the old School, Lorraine Yearon and Des Leech.

Thanks also to Joe and Lindy at Urra, and the many enjoyable nights that I have spent in Shady Grove, Caymars and Fats (not necessarily in that order). I look forward to enjoying a few more in the future. I also acknowledge the help and assistance that José Ruiz has provided, and look forward to West Ham playing Real Madrid in the not so distant future! I would also like to thank the British Geomorphological Research Group (BGRG) for their financial support.

Furthermore, I would like to acknowledge the welcomed distractions provided by my team-mates at Mitcham Cricket Club (both on and off the field), and other friends beyond Kingston University. There are far too many individuals to mention, you know who you are.

Special thanks go to my parents for their support and understanding, particularly when I was stuck in the 'trough of despair' (more than once!). These special thanks extend to my brother and sister, and my grandparents.

I would also like to thank my girlfriend's family who have over the past two years also had to deal with the traumas that PhD bring! Finally, thanks to my girlfriend Vicky, who has offered support and advice over the past two years.

## Abstract

This thesis provides an alternative method to for mapping soil erosion. The method is conducted in a small study area of 40 km<sup>2</sup> in the Sorbas Basin, Almería province, Southeast Spain.

Soil erosion is one of the most destructive land degradation processes and can often lead to serious environmental problems. It is important to implement appropriate management strategies to meet these challenges at a range of scales. However, prior knowledge of erosion processes and the extent to which they operate spatially is often limited and, traditional methods of soil erosion mapping are often time and labour intensive. This thesis explores the use of two Artificial Intelligence (AI) techniques for soil erosion mapping; Artificial Neural Networks (ANNs) and Decision Tree Classifiers (DTCs). The opportunities for employing such methods relate in part to their non-linear capabilities, their ability to learn in an inductive manner and incorporate multi-source data sets.

AI training and test data were collected from 520 individually sampled locations within the study area. At each site the dependent variable erosion was estimated, as were a range of independent variables through field study. Two Digital Elevation Models were developed. Laboratory analysis was also undertaken to explore the physico-chemical processes relating to soil dispersion and to determine the applicability of a soil sodicity meter developed by the Co-operative Research Centre for Soil and Land Management in Adelaide, Australia.

Results demonstrate that classification accuracy and overall performance is strongly dependent on the independent and dependent variables used, with the more expensive field collected data providing improved variables to those extracted from the Digital Elevation Models. Discriminant Analysis (DA) classifications were also employed to provide a linear comparison to the AI techniques, and performed comparably well. In the Artificial Neural Network classifications the composition of the training set was seen to exert significant bias, leading to poor performance and often misleading results. Laboratory analysis highlights the complex physico-chemical relationships associated with soil dispersion. The findings also indicate that no discernible relationship exists between the sodicity meter and standard laboratory procedures employed to measure the sodic properties of a soil.

The thesis demonstrates the potential for employing these methods for erosion risk analysis and the ability of inductive approaches to formulate rules that may enhance current levels of understanding associated with soil erosion processes. Mapped outputs produced by these methods may prove valuable in the management of landscapes susceptible to soil erosion.

## III

## Contents

ACKNOWLEDGEMENTS	I
ABSTRACT	II
CONTENTS	III
LIST OF FIGURES	VII
LIST OF TABLES	X
<ul> <li>1 INTRODUCTION TO THE THESIS</li> <li>1.1 INTRODUCTION</li> <li>1.2 RATIONALE FOR THE STUDY</li> <li>1.3 AIMS AND OBJECTIVES OF THE STUDY</li> <li>1.4 THESIS STRUCTURE</li> </ul>	1 1 2 4 6
2 SOIL EROSION AND EROSION MAPPING CHALLENGES IN SOUTHEAST SPAIN 2.1 INTRODUCTION	<b>8</b> 8
<ul> <li>2.2 LAND DEGRADATION AND SOIL EROSION</li> <li>2.3 CURRENT EROSION AND RISK MAPPING METHODS</li> <li>2.4 STUDY AREA <ul> <li>2.4.1 Introduction</li> <li>2.4.2 Regional Geological Setting</li> <li>2.4.3 Lithological Units</li> <li>2.4.4 Climate, Vegetation and Land-use</li> <li>2.4.5 The Geomorphological Sensitivity of the Almería Landscape</li> </ul> </li> <li>2.5 THE NEED FOR EROSION AND RISK MAPS AND THEIR APPLICATION</li> <li>2.6 CONCLUSIONS</li> </ul>	8 9 11 12 14 16 17 20 21
3 GEOMORPHOLOGICAL THEORY AND SOIL EROSION PROCESSES,	22
<ul> <li>MODELLING AND MAPPING</li> <li>3.1 INTRODUCTION</li> <li>3.2 SCALES OF INVESTIGATION AND THRESHOLDS</li> <li>3.3 BADLAND AND SEMI-ARID GEOMORPHOLOGY</li> <li>3.4 SOIL EROSION PROCESSES IN SOUTHEAST SPAIN</li> <li>3.4.1 Surface Erosion</li> <li>3.4.2 Subsurface Erosion (Piping)</li> <li>3.5 SOIL EROSION MODELLING AND MAPPING TECHNIQUES</li> <li>3.5.1 Erosion Process Models</li> <li>3.5.2 Mapping Techniques</li> <li>3.5.3 Synergistic Methods</li> <li>3.6 SOIL EROSION RISK, HAZARD AND POTENTIAL</li> <li>3.7 CONCLUSIONS</li> </ul>	22 22 27 29 30 32 42 43 45 47 48 49
<ul> <li>4 ARTIFICIAL INTELLIGENCE CLASSIFIERS</li> <li>4.1 INTRODUCTION</li> <li>4.2 ARTIFICIAL NEURAL NETWORKS (ANNs) <ul> <li>4.2.1 Introduction</li> <li>4.2.2 Activation Functions</li> <li>4.2.3 Training Methods and Algorithms</li> <li>4.2.4 Current Uses of Artificial Neural Networks</li> <li>4.2.5 Accuracy of Artificial Neural Networks</li> <li>4.2.6 Disadvantages of the Neural Network Approach</li> <li>4.2.7 Summary</li> </ul> </li> </ul>	<b>50</b> 50 54 54 56 59 63 68 69 75

4.3 DECISION TREE CLASIFIERS (DTCs)	75
4.3.1 Introduction	75
4.3.2 Decision Tree Algorithms	76
4.3.3 Decision Tree Structures	82
4.3.4 Current Uses of Decision Tree Classifiers	84
4.3.5 Disadvantages of the Decision Tree Approach	87
4.3.6 Summary	89
4.4 CONCLUSIONS	89
5 RESEARCH FRAMEWORK AND DATA METHODS	91
5 1 INTRODUCTION	91
5 2 METHODOLOGICAL SETTING	91
5 3 RESEARCH FRAMEWORK	93
5.3.1 Training Data Sets	94
5.3.2 Validation of Variables	95
5 4 SAMPLING STRATEGY	98
5 5 ATTRIBUTE ACOUIREMENT	102
5.5.1 Field Acquired Variables	102
5.5.2 DFM Acquired Variables	110
5 6 DATA TRANSCRIPTION	114
5.6.1 Data Transcription to Artificial Neural Networks	115
5.6.2 Data Transcription to Decision Tree Classifiers	117
5.6.3 Implementation of Independent and Dependent Variables	118
5.7. ACCLIDACY ASSESSMENT	120
5 & DEVELOPMENT OF AN FROSION RISK SCHEDULE	120
5.8.1 Methodology behind the Risk by Association Schedule	122
5.8.2 Justification for Risk Values	125
5.8.2 Summary	120
5.0.5 Summary	127
J.Y CONCLUSIONS	129
CONTERDOSION OF ASSIETCATIONS DISK SCHEDULES AND DUE	131
0 SOLL EROSION CLASSIFICATIONS, RISK SCHEDULES AND RULE	131
	121
0.1 INTRODUCTION 6.2 SOIL EDOSION OF ASSIERCATIONS USING ADTIERCIAL NEURAL	131
0.2 SOIL EROSION CLASSIFICATIONS USING ARTIFICIAL NEURAL	122
NETWORKS 6.2.1 Two Classifications Using Artificial Neural Networks	102
6.2.2 Three Class Classifications Using Artificial Neural Networks	137
6.2.2 Three Class Classifications Using Artificial Neural Networks	141
0.2.5 NINE CLASS CLASSIFICATIONS USING ARTIFICIAL NEURAL NELWORKS	145
6.3 SOIL EROSION CLASSIFICATIONS USING DECISION TREE CLASSIFIERS	145
6.3.1 Two Class Classifications Using Decision Tree Classifiers	150
6.3.2 Three Class Classifications Using Decision Tree Classifiers	154
6.3.3 Nine Class Classifications Using Decision Tree Classifiers	150
6.4 SOIL EROSION CLASSIFICATIONS USING DISCRIMINANT ANALISIS	150
6.4.1 Two Class Classifications Using Discriminant Analysis	161
6.4.2 Three Class Classifications Using Discriminant Analysis	167
6.4.3 Nine Class Classifications Using Discriminant Analysis	167
6.5 COMPARISON OF CLASSIFICATION TECHNIQUES	107
6.5.1 Comparison of Techniques for Two Class Classifications	172
6.5.2 Comparison of Techniques for Three Class Classifications	170
6.5.3 Comparison of Techniques for Nine Class Classifications	1/9
6.6 THE SELECTION AND INFLUENCE OF INDEPENDENT VARIABLES ON	100
CLASSIFIER PERFORMANCE	100
6.6.1 Independent Variables Extracted from the Digital Elevation Models	183
6.6.2 Independent Variables Acquired from the Field	186
6.6.3 Comparison of Field and DEM derived Independent Variables	189
6.6.4 Summary	191
6.7 RULE EXTRACTION USING RESPONSE SURFACES AND SPLITTING	

IV

CRITERIA	191
6.7.1 Artificial Neural Network Response Surfaces	194
6.7.2 Decision Tree Classifier Splitting Criteria	198
6.8 THE DETERMINATION OF OPTIMAL ARTIFICIAL NEURAL NETWORK	
ARCHITECTURES AND A REVIEW OF DECISION TREE TOPOLOGIES	204
6 9 FROSION RISK SCHEDULES AND POTENITAL	204
6.0.1 Erosion Probability Mans	208
6.0.2 Erosion Rick by Association Mans	214
6.9.2 Elosion Nisk by Association Maps	
0.10 CONCLUSIONS	
7 FIELD INVESTIGATIONS	215
7.1 INTRODUCTION	215
7.2 FIELD AND LABORATORY METHODS	217
7.2.1 Field Sodicity Meter	217
7.2.2.1 aboratory Methods	219
7 3 RESULTS AND DISCUSSION	222
7.3.1 Analysis of Laboratory Results	222
7.3.2 The Relationship between the Laboratory Analysis and the Field	234
Sodicity Meter Readings	234
7 4 SUMMARY	237
7.5 CONCLUSIONS	240
8 DISCUSSION	241
8.1 INTRODUCTION	241
8.2 THE DEPENDENT VARIABLE	241
8.2.1 The Effects of the Number of Classes Incorporated into the Dependent	242
Variable	
8.2.2 The Influence of the Training Data Set Composition on Classifier	245
Performance	
8.2.3 Interpretation of the Classified Soil Erosion Maps	250
8 3 THE SELECTION AND SELECTION OF INDEPENDENT VARIABLES	252
8.3.1 Independent Variables and their roles within the Classifications	253
8.3.2 Implications of the Field Sodicity Meter as an Independent Variable	256
8.3.3 Overall Performance of the Various Training Sets	250
8.5.5 Overall renormance of the various framing sets	259
8.4 THE USE OF ARTIFICIAL INTELLIGENCE TECHNIQUES FOR KNOWLEDGE GAIN AND RULE FYTRACTION	209
8 4 1 Knowledge Extraction through Decision Tree Classifiers	260
8.4.1 Knowledge Extraction through Artificial Neural Networks	209
8.4.2 Kinowiedge Extraction through Artificial Neural Networks	270
0.4.5 SUMMARY	212
8.5 ACCURACY ASSESSMENTS OF THE CLASSIFIED SOIL EROSION MAPS	2/3
8.6 UVERALL PERFORMANCE OF THE AT APPROACH AND A TRADITIONAL	278
CLASSIFICATION METHOD	206
8.7 SUMMARY	283
8.8 CONCLUSIONS	289
9 CONCLUSIONS	290
9 1 INTRODUCTION	290
9.2 SUMMARY OF THE MAIN FINDINGS	291
0.2 ELITIDE DESEADOU	294
0 A CONCLUDING COMMENTS	295
	_,,
	307

REFERENCES

APPENDIX 1 – Detailed Laboratory Analysis	328
APPENDIX 2 – Correlation Matrices	331
APPENDIX 3 – Decision Tree Structures	358
APPENDIX 4 - Neural Network Architectures, Errors and Accuracy and Decision Tree Relative Cost and Topology	382
<b>APPENDIX 5 – Sensitivity Analysis and Variable Importance</b>	397
APPENDIX 6 – Soil Erosion and Risk Maps	409

## VII

## **List of Figures**

## Page

Figure 2.1	Location map of the study area (Adapted from Mather and Stokes, 1996)	11
Figure 2.2	Detailed geological map of the Sorbas region (Adapted from Weijermars, 1991)	13
Figure 2.3	Geological map of the study area (Adapted from Mather, 2000a)	16
Figure 3.1	Spatial-temporal domains for different research scales (Delcourt and Delcourt, 1988)	23
Figure 3.2	The changes in channel gradient during cyclic, graded and steady time (Schumm and Lichty, 1965)	24
Figure 3.3	Geomorphological thresholds and reaction and relaxation	26
Figure 3.4	The location of total badlands, partial badlands and dissected softrock	29
-	areas susceptible to badland development (Harvey and Calvo, 1989)	
Figure 3.5	Simple representation of the Horton hypothesis (Adapted from Ward and Robinson, 2000)	31
Figure 3.6	Process dominance domains for the role of site regulators (Faulkner, 2003b)	35
Figure 3.7	A simple diagram highlighting the piping potential in convex morphologies	36
Figure 3.8	The Mocatán badlands demonstrating extensive piping activity leading to the creation of an extensive gully system	36
Figure 3.9	Cation exchange capacity against exchangeable sodium percentage for soil dispersivity classification (Gerber and Harmse, 1987)	39
Figure 3.10	Factors affecting the 'threshold concentration' curve (Quirk and Schofield, 1955) adapted by Sumner (1995)	40
Figure 3.11	A classification for the prediction of dispersive behaviour of red- brown earths (Rengasamy <i>et al.</i> , 1984) adapted by Faulkner <i>et al.</i> (2000)	42
Figure 4 1	The contrasting responses between two different systems	51
Figure 4.1	An example of a neural network with a 5.3.1 architecture	54
Figure 4.3	The internal process of the activation unit (Russell and Norvig, 1995)	56
Figure 4.4	Common activation functions and their accompanying mathematical functions	58
Figure 4.5	Error in the training and test sets during training (Picton, 2000)	74
Figure 4.6	A decision tree classifier (Friedl and Brodley, 1997)	76
Figure 4.7	The entropy function relative to a boolean classification (Mitchell, 1997)	80
Figure 4.8	Axis-parallel decision boundaries of a univariate decision tree (Adapted from Pal and Mather, 2003)	83
Figure 4.9	Decision boundaries for a multivariate decision tree classifier (Adapted from Pal and Mather, 2003)	83
Figure 5.1	The 520 sampling locations draped onto a DEM of the study area	102
Figure 5.2	The potential over-estimation of vegetation cover when using photography	104
Figure 5.3	The sliding scale classification for the dependent variable erosion	107
Figure 5.4	Examples of the erosion classes attributed using the erosion classification scale	108
Figure 5.5	The DEM generated with grid cell size of 10 metres	111
Figure 5.6	The DEM generated with grid cell size of 20 metres	112
Figure 5.7	DEM calculation of slope angle (Adapted from Longley et al., 2001)	113
Figure 5.8	The attributes determined from the DEMs (10 metre cell size)	114
Figure 5.9	The amalgamation of different erosion classes for use within different classifications	119

Figure 5.10	A simple diagrammatic representation of the independent and dependent variables	119
Figure 5.11	The progressive development of the erosion risk schedule	124
Figure 5.12	The logical progression of the risk by association schedule and the	124
I Iguie 5.12	assigned risk values based on the scenario	120
Figure 6.1	ROC curve for the two class ANN classification trained using the field acquired attributes and the 10 metre DEM independent variables.	135
Figure 6.2	Classified erosion map drape derived from the ANN trained using 10 metre DEM variables for a two class classification	136
Figure 6.3	Classified erosion map drape derived from the ANN trained using 10 metre DEM variables for a three class classification	140
Figure 6.4	Classified erosion map drape derived from the ANN trained using 20 metre DEM variables for a nine class classification	144
Eigura 6 5	An example of an error curve used to determine the optimum DTC	145
Figure 0.5	All example of all error curve used to determine the optimum DTC	145
Figure 0.0	ROC curve for the two class DTC classification trained using the	147
<b>F</b> :	The desiries frequencies for the true should be if the initial of the	
Figure 6. /	The decision tree grown for the two class classification using the field	148
	acquired independent variables	
Figure 6.8	Classified erosion map drape derived from the DTC trained using 10 metre DEM variables for a two class classification	149
Figure 6.9	The decision tree grown for the three class classification using the	152
	field acquired attributes and the classified vegetation independent variables	
Figure 6.10	Classified erosion map drape derived from the DTC trained using 10	153
	metre DEM variables for a three class classification	
Figure 6.11	The decision tree grown for the nine class classification using the field acquired independent variables	156
Figure 6.12	Classified erosion map drape derived from the DTC trained using 10	157
0	metre DEM variables for a nine class classification	157
Figure 6.13	Classified erosion man drape derived from the DA trained using 10	160
8	metre DEM variables for a two class classification	100
Figure 6 14	Classified erosion man drane derived from the DA trained using 10	162
I Igure 0.1 /	metre DEM variables for a three class classification	105
Figure 6.15	Classified erosion man drane derived from the DA trained using 10	144
rigure 0.15	metre DEM variables for a nine alogs alogsification	100
Figure 6 16	Overall accuracies achieved for the two class classification	1/0
Figure 6.10	BOC surves for both the ADDI and DTC should be discussed in the	168
rigure 0.17	the field acquired independent variable, classifications trained using 10 metre DEM attributes	172
Figure 6 18	Overall accuracies achieved for the three class classifications	177
Figure 6 10	Overall accuracies achieved for the nine class classifications	175
Figure 0.19	DOC outries derived from each of the eight ADDI two shows	1/0
rigure 0.20	classifications	189
Figure 6.21	ROC curves derived from each of the eight DTC two class classifications	190
Figure 6.22	Response surface for (A) slope angle against estimated vegetation and (B) slope angle against flow length, for the two class classification using the ANN trained with the field variables and 10 metre DEM variables	192
Figure 6.23	Response surface for (A) slope angle against estimated vegetation and (B) slope angle against sodicity, for the three class classification using the ANN trained with the field variables and 10 metre DEM variables	193
Figure 6.24	Response surface for (A) slope angle against slope aspect and (B) slope angle against classified vegetation, for the nine class classification using the ANN trained with the field variables and classified vegetation	194

Figure 6.25	Some of the rules extracted from the DTC grown for the two class classification using field acquired independent variables	195
Figure 6.26	The splitting criteria determined by the DTC trained using the 10 metre DEM data set	196
Figure 6.27	Some of the rules extracted from the DTC grown for the three class classification using field acquired and 10 metre DEM independent variables	197
Figure 6.28	Error verses accuracy graph for the ANNs trained using the field acquired independent variables (two class)	199
Figure 6.29	Error verses accuracy graph for the ANNs trained using the field acquired and 10 metre DEM independent variables (two class)	200
Figure 6.30	Error verses accuracy graph for the ANNs trained using the field acquired independent variables (three class)	200
Figure 6.31	Error verses accuracy graph for the ANNs trained using the field acquired independent variables and classified vegetation (three class)	201
Figure 6.32	Error verses accuracy graph for the ANNs trained using the 10 metre DEM independent variables (nine class)	201
Figure 6.33	Error verses accuracy graph for the ANNs trained using the field acquired independent variables (nine class)	202
Figure 6.34	Relative cost and terminal nodes for the field acquired independent variables (two class)	203
Figure 6.35	Relative cost and terminal nodes for the field acquired independent variables and classified vegetation (two class)	203
Figure 6.36	Erosion probability map drape produced from the DTC trained with the 10 metre DEM variables for a two class classification	206
Figure 6.37	Gully erosion probability map drape produced from the DTC trained with the 20 metre DEM variables for a three class classification	207
Figure 6.38	Risk by association map drape produced from the DTC trained with the 10 metre DEM variables for a two class classification	209
Figure 6.39	Risk by association map drape of surface erosion produced from the DTC trained with the 10 metre DEM variables for a three class classification	210
Figure 6.40	Risk by association map drape of subsurface erosion produced from the DTC trained with the 10 metre DEM variables for a nine class classification	211
Figure 6.41	Risk by association map drape for surface and subsurface erosion produced from the DTC trained with the 10 metre DEM variables for a nine class classification	212
Figure 6.42	Risk by association map of surface and subsurface erosion subsection	213
Figure 7.1	The field sodicity meter designed by the Australian Co-operative Research Centre for Soil and Land Management	217
Figure 7.2	The laboratory results for CEC and ESP plotted on the dispersivity classification graph identified by Gerber and Harmse (1987)	225
Figure 7.3	The relationship between ESP and pH	227
Figure 7.4	The relationship between ESP and EC using the power function	228
Figure 7.5	The relationship between SAR and EC plotted on log-transformed axes using the power function	229
Figure 7.6	The relationship between ESP and the field sodicity meter	235
Figure 7.7	The relationship between the Gerber and Harmse (1987) dispersivity index and the field sodicity meter	235
Figure 8.1	Comparison between slope and aspect measurements collected in the field and those extracted from the 10 and 20 metre DEMs	267
Figure 8.2	Response surfaces detailing the relationships between slope angle, vegetation cover and aspect for the two class classifications	272
Figure 8.3	The procedures required for the development of the AI techniques	287

## List of Tables

Table 3.1	The status of drainage basin variables during time spans of decreasing duration (Schumm and Lichty, 1965)	25
Table 3.2	Some studies investigating badland geomorphology and processes	27
Table 3.3	Summary of various soil erosion models	45
Table 4.1	Advantages and disadvantages of various classification techniques	53
Table 5.1	The variables collected for the study area and their sources	97
Table 5.2	A brief description of the attributes calculated from the derived	98
	DEMs	
Table 5.3	Common spatial sampling strategies	100
Table 5.4	The eight different data sets and the independent variables that they	118
	contain	
Table 5.5	An example of a correlation (error) matrix	121
Table 5.6	The producers accuracy	121
Table 5.7	The users accuracy	121
Table 6.1	Summary table of 'best' networks for two class classifications using	133
	ANNs	
Table 6.2	Correlation matrix for the two class ANN trained using the field	133
	acquired attributes and the 10 metre DEM independent variables	
Table 6.3	Summary table of the 'best' networks for the three class	137
	classifications using ANNs	
Table 6.4	Correlation matrix for the three class ANN trained using the field	138
T-bla ( 5	acquired attributes	130
Table 0.5	Correlation matrix for the three class ANN trained using the field	138
Table 6 6	Summary table of the 'best' networks for nine class classifications	141
	using ANNs	141
Table 67	Correlation matrix for the nine class ANN trained using the field	142
10000.7	acquired attributes and the classified vegetation independent variable	172
Table 6.8	Summary table of decision trees grown for two class classifications	146
	using DTCs	1.0
Table 6.9	Correlation matrix for DTCs trained using the field acquired	146
	independent variables	
Table 6.1	Summary table of the 'best' tree for three class classifications using	150
	DTCs	
Table 6.1	1 Correlation matrix for the three class DTC trained using the field	151
	acquired attributes and the classified vegetation independent variable	
Table 6.1	2 Summary table of the 'best' networks for nine class classifications	154
	using DTCs	
Table 6.1	3 Correlation matrix for the nine class DTC trained using the field	155
	acquired independent variables	1.00
Table 6.1	4 Summary table of the two class classifications using discriminant	159
	analysis	160
Table 6.1	5 Correlation matrix for the two class DA using the field acquired	139
<b>T</b> 11 (1	independent variables	161
Table 6.1	5 Summary table of the three class classifications using discriminant	101
T-bla ( 1	analysis Completion metric for three class DA using the field acquired and 10	162
Table 0.1	metre DEM independent variables	102
Tabla 6 1	Summary table of the nine class classifications using discriminant	164
	analysis	104
Table 6.1	9 Correlation matrix for DA using field acquired and 20 metre DEM	165
14010 0.1	independent variables	

Table 6.20	Correlation matrix for DA using the field acquired, 10 metre DEM independent variables and classified vegetation	165
Table 6.21	Area under curve comparison table for the ANNs and DTCs for the two class classifications	171
Table 6.22	Composition of erosion maps for all three techniques for the DEM classifications using each of the three class classifications	174
Table 6.23	Composition of erosion maps for all three techniques for the DEM classifications for the nine class classifications	177
Table 6.24	Sensitivity analysis for the ANN trained using the 10 metre data set for the two class classification	180
Table 6.25	Sensitivity analysis for the ANN trained using the 20 metre DEM data set fort he two class classifications	180
Table 6.26	Variable importance for the DTC trained using 10 metre DEM data set for the two class classification	181
Table 6.27	Variable importance for the DTC trained using the 20 metre DEM data set for the two class classification	182
Table 6.28	Sensitivity analysis for the ANN trained using the filed acquired data set for the two class classification	184
Table 6.29	Sensitivity analysis for the ANN trained using the filed acquired data set for the three class classification	184
Table 6.30	Sensitivity analysis for the ANN trained using the filed acquired data set for the nine class classification	184
Table 6.31	Variable importance for the DTC trained using the field acquired data set for the two class classification	185
Table 6.32	Variable importance for the DTC trained using the field acquired data set for the three class classification	185
Table 6.33	Variable importance for the DTC trained using the field acquired data set for the nine class classification	185
Table 6.34	Sensitivity analysis for the ANN trained using the field acquired data set and classified vegetation for the two class classification	186
Table 6.35	Variable importance for the DTC trained using the field acquired data set for the nine classification	186
Table 6.36	Sensitivity analysis for the ANN trained using the field acquired data set and the 10 metre DEM variables for the two class classification	187
Table 6.37	Variable importance for the DTC trained using the field acquired data set and the 10 metre DEM variables for the three class classification	187
Table 6.38	Sensitivity analysis for the ANN trained using the field acquired data set, classified vegetation and the 10 metre DEM variables for the three class classification	188
Table 6.39	Sensitivity analysis for the DTC trained using the field acquired data set, classified vegetation and the 10 metre DEM variables for the three class classification	188
Table 7.1	Results obtained from the laboratory analysis and the field sodicity meter	223
Table 7.2	Slope angle, aspect and vegetation cover values for each of the 55 sites sampled along with the erosion and sodicity meter values	233
Table 8.1	Accuracy of a two class classification interpreted from the three class correlation matrices	243
Table 8.2	Accuracy of a two and three class classification interpreted from the nine class correlation matrices	243
Table 8.3	Composition of the training data set	246
Table 8.4	Composition of the test data set	275
Table 8.5	Probability of scoring no errors in various size samples from a population with a range of real error proportions (van Genderen and Lock, 1977)	276

## 1

## Introduction to the Thesis

#### **1.1 INTRODUCTION**

Soil erosion, caused by wind and water is considered one of the most important and destructive land degradation processes. It is an increasing global problem (Vrieling *et al.*, 2002), largely due to its inherent ability of creating severe on and off-site impacts. These impacts can take many forms; both environmentally and economically. It is important to understand both the underlying processes that drive erosional activities and the impacts that these processes have to be able to implement appropriate management strategies and policies at a range of scales. The number of soil erosion investigations and studies has risen drastically in recent years as these concerns have grown, leading to the continual development and furthering of current knowledge.

Landscapes continually change, as does the extent of anthropogenic influences upon environments, resulting in the need for continual monitoring and increased process knowledge. It is estimated that during the past 50 years, human land use and other activities associated with it have resulted in the degradation of some 5 billion ha of land globally (Brady and Weil, 2002). In many regions erosion rates are significantly exceeding soil formation rates, (e.g. Australia) (see Edwards, 1991), highlighting a number of issues relating to long term sustainability.

#### **1.2 RATIONALE FOR THE STUDY**

The Mediterranean region is particularly sensitive and vulnerable to soil erosion. It affects different parts of the region to varying degrees (López-Bermúdez *et al.*, 1998). Rojo (1990) suggests that more than 22 million ha of land in Spain is affected by erosion rates in excess of 12 t ha<sup>-1</sup> yr<sup>-1</sup>, exceeding the estimated tolerable limits for soil formation of between 2 and 12 t ha<sup>-1</sup> yr<sup>-1</sup> in Mediterranean environments, a total of approximately 44 percent of the land. Furthermore, the National Institute for the Conservation of Nature (ICONA, 1988) estimated that within the Mediterranean region, more than 9 million ha of land are affected by very intense erosion rates of more than 50 t ha<sup>-1</sup> yr<sup>-1</sup> highlighting the need for management strategies and policies.

One area that is at present experiencing rapid change and increased risk of serious erosion, is the province of Almería, located in Southeast Spain. In recent years the interior regions of Almería have witnessed an extensive and rapid change in both the appearance and stability of the natural landscape. This is occurring largely as a consequence of changing agricultural regimes, from traditional dry farming (secano) methods, towards extensive plantation arboriculture (Faulkner *et al.*, 2003b). These areas are described as geomorphologically sensitive prior to any such agricultural changes, due to both the climatic characteristics and the sensitive lithologies found in the province (see Alexander *et al.*, 1996; Spivey, 1997; Faulkner *et al.*, 2000, 2003b). With the added implications of agricultural clearances, the risk of erosional activity is greatly increased.

Management of these environments, in the light of potential human developments, is important if these developments are to be economically and environmentally sustainable. Traditional techniques that have aided these management decisions include ground-based studies, remotely sensed imagery and aerial photography, and physical process models. The use of Artificial Intelligence (AI) based classifiers provides an approach that can utilise multi-source data sets (de Carvalho *et al.*, 2004), offer the ability to work in a non-linear manner, and an extensive array of studies have successfully incorporated them within environmental investigations. For example, Ermini *et al.* (2005) used Artificial Neural Networks (ANNs) to determine landslide susceptibility in the northern Apennines in Italy, and Benediktsson *et al.* (1990) compared ANNs to traditional statistical techniques when using multi-source data sets in remote sensing. Mahiny and Turner (2003) explored the use of ANNs for the determination of vegetation change in Australia, and de Carvalho *et al.* (2004) used both Decision Tree Classifiers (DTCs) and ANNs to map forests in Brazil. Fitzgerald and Lees (1993) compared the two techniques for use with remotely sensed data and as with all of the aforementioned studies, recognised the advantages they hold.

Understanding soil erosion processes and the extent to which they may be operating spatially in the environment is important to a range of different groups including regional governments, local governments and farmers. At smaller scales, land managers may have an intuitive empirical understanding of the landscape that aids their management decisions. However, erosion maps can aid these decisions and assist in the successful implementation of appropriate management strategies and policies. Nonetheless, they are only useful if they have been developed with a specific end-user in mind, and are at a spatial resolution that is relevant for a specific application. At present a soil erosion map exists for the study area, however, its resolution is too coarse to be of use at a 'local' scale. Therefore, an opportunity exists to develop and validate a methodology for mapping soil erosion processes spatially using multi-source data sets and AI techniques, at a scale deemed viable to assist in the implementation of management practices at a variety of spatial scales.

This study focuses on the development and implementation of a suitable methodology for spatially mapping erosion processes. At present, few methods can be implemented to determine the spatial extent of soil erosion processes using pre-existing or low-cost data sets in combination with one another that are relatively rapid in execution. Such a method would present itself as an applicable tool for assisting a range of management decisions at various scales. This is largely due to the fact that erosion can be considered a primary indicator of the sustainability of land use. The method presented here maps erosion in a visually qualitative manner, identifying extreme cases through to stable, non-affected sites using a simple sliding scale.

The following sections within this chapter outline the aims and objectives of the study and provide a brief overview of each of the remaining chapters.

#### **1.3 AIMS AND OBJECTIVES OF THE STUDY**

Through the use of the two AI classification techniques, Artificial Neural Networks (ANNs) and Decision Tree Classifiers (DTCs), soil erosion processes and the extent to which they are operating in a small study area located in Almería are mapped. Unlike some process models, the temporal dimension is not a factor as the classification process undertaken is in essence mapping the current spatial extent of erosion in the region: The technique maps 'form' and does not implicitly measure

change. Incorporating temporal parameters is beyond the scope of the study. The aims and objectives of the thesis are as follows:

## • Aim 1: Evaluate and compare the performance of Artificial Neural Networks and Decision Tree Classifiers as soil erosion classifiers.

Developing and constructing ANNs and DTCs for various classification problems, using different independent and dependent variables allows their performances to be compared and contrasted. This will highlight potential problems and advantages associated with each technique.

• Aim 2: To determine the ability of Artificial Neural Networks and Decision Tree Classifiers to further our current understanding of soil erosion processes and how the selection of dependent and independent variables influences the classifiers performance.

The two classifiers used possess the ability to use multi-source data sets, providing the opportunity to use variables in combination that was previously unachievable using traditional techniques. The influence that different combinations of independent and dependent variables have upon classifier accuracy can be determined as can their ability to enhance our understanding of the erosion processes.

Traditional modelling techniques are driven in a deductive manner by prior knowledge. However the use of AI classifiers offers the ability to work largely inductively by presenting the independent variables and allowing rules and parameters to be determined by the classifier alone. A particularly important objective within this aim involves the determination of the ability of a field sodicity meter to correctly, measure and classify soil sodicity. The presence of sodium as an exchangeable cation adversely affects extensive areas across the world (Sumner, 1995). Physical processes such as swelling and dispersion are responsible for soil degradation including surface crusting and hardsetting (Sumner, 1993), a result of poor aggregate stability (Greene *et al.*, 2002), a consequent reduction in hydraulic conductivity, and erosion. The determination of sodic soils is important so as to allow for appropriate management strategies, and has traditionally been undertaken through standard laboratory techniques. However, the determination through a simple field meter possesses a number of inherent advantages, such as the rapid measurement of levels of sodicity under field conditions. Therefore, the determination of the meters applicability is important within this thesis.

#### **1.4 THESIS STRUCTURE**

Chapter Two, in providing an introduction to the problem of soil erosion and soil erosion mapping in Southeast Spain, provides the general scope for the study. This chapter documents the need and requirement for erosion mapping detailing current erosion and risk mapping methods. The chapter also introduces the study area, detailing the geological and geomorphological setting, and highlights the perceived geomorphological susceptibility of the landscape. Chapter Three details geomorphological issues including scale and threshold concepts, erosion processes, as well as current and past soil erosion modelling methods and techniques. This provides the reader with the geomorphological process knowledge which aids the understanding of the erosion processes operating in Southeast Spain. Chapter Four provides a review of the artificial intelligence classification techniques used. It details the workings and theory behind artificial neural networks and decision tree classifiers. It also outlines their current uses, the advantages and disadvantages associated with each of them. Chapter Five details the methodology and research framework upon which this study is based. The chapter includes the determination, validation and justification of the dependent and independent variables, their acquirement and finally their implementation using both classification procedures (ANNs and DTCs). The development of an erosion risk schedule is also discussed and outlined based on the relatively simple concept of risk by association.

Chapter Six presents the results obtained through the various classifications undertaken, with particular emphasis upon the aims and objectives set out previously in section 1.3. Chapter Seven documents and discusses the results obtained through field and laboratory techniques employed to determine soil sodicity and the subsequent dispersivity potential. The chapter explores the relationships between various measured parameters, including physico-chemical characteristics, and concentrates particularly on the relationship between the field sodicity meter and the laboratory tests.

Chapter Eight presents a discussion of the results obtained, relating important findings to the appropriate literature. The results are analysed and their wider implications discussed, with particular reference and emphasis made to the aims and objectives of this work. Chapter Nine outlines the major findings, highlights the associated shortfalls and limitations as well as proposing where future work may be focused.

## 2

## Soil Erosion and Erosion Mapping Challenges in Southeast Spain

#### **2.1 INTRODUCTION**

This chapter provides an introduction to the problem of soil erosion and soil erosion mapping as well as the scope, emphasis and rationale for the research in this thesis. The importance of understanding the nature and magnitude of soil erosion is important for many management decisions, and erosion mapping must be appropriate to the end-user and for the application for which it is to be applied. Thus, land degradation, soil erosion and the multitude of associated problems are initially outlined, followed by a review of the methods and approaches used to produce erosion maps and erosion risk maps. The study area is introduced, including an overview of the regional geology, the lithological units, climate, vegetation and land-use as it affects the region's geomorphological sensitivity, exacerbated by the recent landscape changes.

### **2.2 LAND DEGRADATION AND SOIL EROSION**

Land degradation can be defined in many ways; largely involving any change in the land that reduces its condition or quality resulting in the deterioration of the physical and chemical properties of the soil (Imeson and Emmer, 1992), rendering the land less useful and of limited economic value. It can take many forms, including land clearance and deforestation, the agricultural mining of soil nutrients, pollution and poor agricultural practices (Brady and Weil, 2002; USDA, 2005). However, soil erosion is perhaps the most serious form of land degradation (Jayasuriya, 2003) and is

becoming an increasing global problem (Vrieling *et al*, 2002), creating severe on and off-site impacts. Consequently, it is important to understand both the underlying processes that drive erosional activities and the extent to which they operate, so as to provide a suitable position to implement appropriate management strategies and policies at a range of scales.

#### 2.3 CURRENT EROSION AND RISK MAPPING METHODS

Three traditional methods exist that have been used extensively in the determination and mapping of soil erosion. Firstly; ground based fieldwork incorporating methods such as erosion pins and sediment traps offer the ability to investigate erosion processes at the field plot scale. Secondly, remote sensing and aerial photography offer the ability to remotely determine erosion processes and the extent to which they are operating at a greater scale. Remote sensing consists of the interpretation of electromagnetic energy reflected by a target object observed by a sensor that is not in contact with the object (Mather, 1999). These approaches have a range of distinct advantages for studying geographical phenomena, including the ability to analyse problems at a range of temporal and spatial resolutions. Spatially, remote sensing exceeds other methods of data collection as vast areas of a landscape can be covered by either satellite imagery or photography which permit a range of investigations at spatial and temporal scales that would otherwise be impossible to achieve. Furthermore, many remote sensing instruments or platforms are hyperspectral and have the ability to measure varying levels of reflectance in different parts of the electromagnetic spectrum through narrowly defined spectral channels (Campbell, 2002). This allows for the rapid distinction and separation of different land coverages or atmospheric variations. King and Delpont (1993) discussed the usefulness and applicability of remote sensing for assessing the spatial and temporal variability of various factors that influence the susceptibility of soils to erosion. Signs of degradation can be recorded on bare ground, such as crusting and scouring which indicate poor surface conditions, they can also be provided by vegetation and land-use, both of which are important indicators of soil conditions. Finally the morphology of the landscape will assist in the determination of slopes at risk.

Modelling approaches have been used extensively with varying degrees of success, developed for investigations at a range of scales from the plot to the national and continental scales. Such approaches attempt to simulate the complex processes involved in soil erosion to improve our understanding of them. The process modelling approach also offers the ability to work at a range of scales and can provide both qualitative and quantitative outputs.

Although three traditional approaches have been identified here for mapping or predicting soil erosion, a fourth approach has been applied in a number of limited studies. Artificial Intelligence (AI) approaches offer the unique ability of pattern recognition and the identification of subtle relationships between dependent and independent variables using a range of data types. They can use relatively low-cost data from a range of sources and use small data sets. Here each technique has been outlined separately, however it is often the case that such methods are used in a synergistic manner in order to complement one another and provide more useful tools. It is important however to be aware of the view of soil erosion produced, either actual or predicted (Ellis, 1997) in order to determine the usefulness and applicability of each approach for any given problem.

## **2.4 STUDY AREA**

## **2.4.1 Introduction**

The study area is located in the province of Almería, near the small municipality of Sorbas in Southeast Spain (Figure 2.1). The total annual precipitation in the province is highly variable, ranging from less than 200mm a year in coastal areas such as the Cabo de Gata (130mm) (Tout, 1987), increasing to 400-500mm on the mountains (Sierras) (García-Latorre *et al.*, 2001). As a consequence of topography distinctive local micro-climates exist (Burke and Thornes, 1998) and in localities with the lowest rainfall, precipitation events largely take the form of short duration but high-intensity storms which occur in autumn and spring. These are intervened by midwinter droughts and hot dry summers (Mather, *et al.*, 2001a; Harvey, 1982). The mean annual temperature in the region is *c.*  $18^{\circ}$ C, averaging 23°C in the summer, 13°C in the winter (Wheeler, 1996) with maximums of *c.* 40°C in July and August (Mather *et al.*, 2001a).



Figure 2.1: Location map of the study area (Adapted from Mather and Stokes, 1996).

The study area is located in the Sorbas Basin, and enclosed by the Sierra de Bedar and Sierra de los Filabres to the north, the Sierra Alhamilla and the Sierra Cabrera to the south and the south-east respectively. The following sections within this chapter will review the geological and lithological setting of the region and their reported geomorphological sensitivities.

#### 2.4.2 Regional Geological Setting

Almería province is occupied within the easternmost part of the Betic Cordillera, which evolved in response to the relative motions of the European and African plates from the late Jurassic to the early Miocene (Bourrouilh and Gorsline, 1979; Smith and Woodcock, 1982; García *et al.*, 2003; Viseras *et al.*, 2003) or the early Cretaceous to the Miocene (Alonso-Chaves *et al.*, 2004). The Betic Cordillera, located in southern Spain, is an ENE-WSW trending thrust belt and Alpine fold (Keller *et al.*, 1995), and as stated by Lonergan *et al.* (1994) is split into Internal and External Zones as with other Alpine systems. The External Zone can be subsequently split into the Pre-Betic Zone, comprising of platform and shelf sequences of marginal facies (García-Hernández *et al.*, 1980), and the Sub-Betic Zone, which contains deep-water Cretaceous to early Tertiary sedimentary sequences with minor basaltic volcanics. In contrast, the Internal Zone, simply known as the Betic Zone, is composed dominantly of metamorphosed Paleozoic to Triassic rocks, resulting from the convergence between Africa and Europe (Lonergan, 1993).

The Betic Cordillera is comprised of a series of uplifted Paleozoic to Triassic metamorphic rocks separated by east-west orientated small sedimentary basins (Braga *et al.*, 2003), such as the Sorbas Basin, filled with post-orogenic detritus of Neogene

age (Harvey and Wells, 1987; Scotney *et al.*, 2000). The Sorbas Basin is located within the Internal zone of the Betic Cordillera and is well defined within the Sierra de los Filabres and Sierra de Bedar in the north, and the Sierra Cabrera and Alhamilla to the South, comprising metamorphic rocks from the Internal Betics (Mather *et al.*, 2001b). However, the eastern and western margins are less well defined topographical highs (Mather, 1993). The geological setting of the Basin can be seen in Figure 2.2. Within the Basin, compression has been dominantly north-south, with associated eastwest extension during the Quaternary (Mather and Westhead, 1993). Much of the compressional movement has taken place along the major left-lateral strike-slip faults (Weijermars, 1991; Bousquet, 1979) which forms part of the left-lateral, Trans-Alboran shear zone (de Larouzíere *et al.*, 1988).



Figure 2.2: Detailed geological map of the Sorbas region (Adapted from Weijermars, 1991).

As can be seen in Figure 2.3, Plio/Pleistocene conglomerates dominate the southern end of the Sorbas Basin. These conglomerates have been divided into two main units by Mather and Stokes (1996, 2001), Mather (1993) and Mather (2000a, 2000b) consisting of a Triassic-metacarbonate-rich unit (TRU), and a Messinian carbonaterich-unit (MRU), the former being overlain by the latter. The two units show a number of differences with regards to their composition, all of which are consequences of alterations in the fluvial systems from which they were developed, which can either be put down to autocyclic or allocyclic controls on the basin. Autocyclic controls may include a river capture event, consequently changing the source area, or it may be the exhaustion of a specific source area. The possible Allocyclic controls responsible for such an event to take place may include tectonic activity, consequently changing the development of the system, or climatic influences may play a role.

#### 2.4.3 Lithological Units

The TRU consists of Triassic limestone and Tortonian sandstone clasts, whereas the upper unit, the MRU consists of Messinian limestone, with the two units being separated by a weakly developed unconformity (Mather, 1993). The TRU is the dominant sedimentary sequence of the two units in the area, and as seen in Figure 2.3, is prolific in the west of the basin, where it forms 95-100 percent of the total sequence (Mather, 1993). However, Mather and Stokes (2001) infer that it becomes less dominant towards the east where it is overlain by the MRU

Mather (1993) suggests that the TRU was most probably sourced from an area on the northern side of the Sierra Alhamilla as the area is rich in Triassic limestone and

would contain an abundant supply of sandstone for transport. The lateral extent and thickness of the sequence suggests that the TRU was deposited through a number of coalescing alluvial fans, all of which were sourced from the southern Sierras at fairly rapid rates due to the lack of any observable well developed soils. Furthermore, paleocurrent data of the clasts identified in the unit also suggest the Alhamilla as the derivation of the material found within the unit.

The MRU however, was not likely to have been sourced from the Sierra Alhamilla area, primarily due to the lack of sandstone found within the unit, and the clast assemblage being dominated by Messinian limestone (Mather, 2000a). The spatial extent of the MRU tends to indicate that the potential source area was much smaller, and much closer, than the one supplying the TRU. With the development of relatively well developed paleosols, the system consisted of overall lower supply rates of sediment to the system through smaller channels (Mather, 2000a 2000b). Using paleocurrent data, the source area is therefore believed to be located in the region around Cantona (Mather, 2000a, 2000b, Mather and Stokes, 2001), southeast of the unit.



Figure 2.3: Geological map of the study area (Adapted from Mather, 2000a).

## 2.4.4 Climate, Vegetation and Land-use

The study area outlined in Figure 2.1 is commonly identified as being located within Mediterranean semi-arid and arid areas of the sub-tropical climatic belt (Bryan and Yair, 1982). During the winter months, precipitation events are associated with fronts coming from the Atlantic Ocean to the west. However after summer and during the autumn, they are largely a consequence of Mediterranean fronts (Lázaro *et al.*, 2001). Such events may provide storms and torrential rainfall during the months of maximum vegetation stress and minimum coverage, thus increasing their erosive potential. Moreover, Zukowskyj *et al.* (2005), Lázaro *et al.* (2001), Thornes (1996), López-Bermúdez and Romero-Diaz (1989) have stated that high inter-annual rainfall variations exist and the recurrence of drought means that the semi-arid area has similar conditions to arid areas for at least one in ten years (Cerdà, 1997).

Vegetation in the region is relatively sparse with a predominance of sclerophylous shrubs (García and Chuvieco, 2004), including *Stipa tenacissima* (Cerdà, 1997), *Retama sphaerocarpa* (Haase *et al.*, 1996) and *Anthyllis cytisoides* (Haase *et al.*, 2000; Blackburn and Steele, 1999). Nonetheless, substantial areas of the region are simply bare surfaces, particularly on south-facing slopes where water stress is at its most extreme (Zukowskyj *et al.*, 2005). Even under thin vegetation cover the semi-arid climate increases the effectiveness of runoff during storm events (Harvey *et al.*, 2001).

Due to low rainfall, farming in the region has traditionally been restricted to dry farming and non-irrigated cultivation (secano) (Spivey, 1997). However, since 1986 and Spain's succession to the European Union, the Common Agricultural Policy (CAP) promoted agricultural activity, largely in the form of irrigated plantation aboriculture (e.g. Olives) in the region. Furthermore, in the coastal regions large areas have been devoted to intensive horticulture under plastic greenhouses (Orgaz *et al.*, 2005). Approximately 27000 hectares of land is now covered by greenhouses in Almería (Molina-Aiz *et al.*, 2004), taking advantage of the high radiation levels, mild winters, the availability of underground water and the development of economic exchanges within the European Union (Gary, 2000).

### 2.4.5 The Geomorphological Sensitivity of the Almería Landscape

The region of the Sorbas Basin has undergone two major river capture and rejuvenation events. The first occurred during the early Pleistocene, and re-routed approximately 15 percent of the original Sorbas Basin drainage to the Carboneras basin (Mather, 2000a, 2000b). The second occurred during the late Pleistocene (ca.

100 ka; Harvey *et al.*, 1995) and re-routed 73 percent of the Carboneras basin in the south, to the east (Mather, 2000a). The Rambla Mocatán is 13 km above the capture site, and has consequently undergone significant moderation by river capture. This consequently led to base level changes of around 90 metres (Harvey *et al.*, 1995), and thus created a situation of intense incision. This led to the oversteepening of slopes, in an area that is highly sensitive to base level changes due to the extremely weak lithologies. Alexander *et al.* (1996), Spivey (1997) and Faulkner *et al.* (2000, 2003b) recognised the sodicity and consequent dispersive nature of both the TRU and MRU making them highly susceptible to erosion. As a result significantly sized badland landscapes have developed, one of particular importance here are the Mocatán badlands which have been the subject of numerous research investigations due largely to its apparent susceptibility to subsurface piping and erosion (see Faulkner *et al.*, 2000, Spivey, 1997; Alexander *et al.*, 1996).

Further intensifying the sensitivity of the Almería landscape are the extensive agricultural clearances in the interior regions developed as a direct response to the 1992 reforms of the CAP. Zukowskyj *et al.* (2005) suggests that the expansion of plantation arboriculture, through such clearances, has been ongoing for at least a decade and the resultant erosion risks have been documented. Faulkner *et al.* (2003b) stressed the heightened erosion risks associated with rapid and extensive clearances on the susceptible lithologies, whereby slopes are reshaped, old trees are removed and the land ploughed.

The region also suffers from great spatial and temporal annual and inter-annual variations in rainfall events (Geeson and Thornes, 1996; Faulkner *et al.*, 2003b).

Thus, frequent drought periods occur creating the potential for vegetation aridization (de la Rosa et al., 1999). The threat of desertification increases as the already unreliable winter rainfall may become more extensive and shift northwards as a product of climate change (Imeson and Emmer, 1992). It has been estimated through the use of Global Circulation Models that the temperature in the Mediterranean basin could increase by 1°C by 2030 (Perry, 1997). However, projecting the extent of precipitation change in the region is more difficult (Wigley, 1992). Nonetheless, Imeson and Emmer (1992) envisaged a range of short-term (50 years) impacts of climate change on sensitive Mediterranean soils. Firstly, the transport and distribution of salts and the salt balance of the soil can lead to an increase in potential erodibility and a subsequent decrease in the stability of soil aggregates in semi-arid regions, induced by a general decrease in precipitation or an increased evapotranspiration rate. Secondly, Imeson and Emmer (1992) proposed that the precipitation of calcium and magnesium carbonates could lead to a distinct caliche layer, often connected with desertification and can impede plant growth and productivity. Finally, as the organic content of a soil is closely correlated to precipitation, and to soil moisture, any variation that may occur as a result of climate change will strongly influence soil aggregate stability.

It is readily evident therefore that the landscape is geomorphologically sensitive as a result of the dispersive nature of the lithological units as discussed previously, coupled with the topographic nature of the region. Moreover, the climatic setting of the area means that extreme precipitation events can occur during summer months, when vegetation cover is at its lowest and thus erosion potential at its maximum. Combined with the agricultural renaissance of the region (Harvey *et al.*, 2001) and the

potential negative impacts associated with climate change, it is easy to understand the susceptibility and serious erosion risk.

## 2.5 THE NEED FOR EROSION AND RISK MAPS AND THEIR APPLICATION

Section 2.3 of this chapter has briefly documented the main methods and techniques employed to produce erosion maps and predicted erosion maps for studies at a range of scales. Such practices have been used to develop erosion maps for various regions of the world as they have numerous practical applications (e.g. for landscape managers, local government and individual landowners). Through the production of erosion maps it is possible to readily identify and determine the spatial extent of the operative processes. This enables the subsequent recognition of susceptible areas (for example, lithologies, soils, geologies and slopes) and may therefore assist in the understanding of specific processes. In addition to this, mapping and predicting risk allows environmental managers to precisely identify areas where intervention must be of high priority (Haboudane *et al.*, 2002).

It has been demonstrated here that the Almería province of southern Spain is a highly sensitive environment and the threat of erosion is of great concern. At present large badland landscapes occur in varying locations controlled largely by topographic and lithological features, however, the situation at present is such that management strategies and practices are required across the region as the sensitivity of the landscape is further increased. Erosion maps exist, but they are at a resolution that is too general to be of 'local' use, potentially inhibiting and restricting management strategies. Erosion maps serve a range of different purposes, dictated largely by the end user(s). Such maps must be at an appropriate scale and resolution to answer the question being posed. For example, the individual land owner (farmer) will require a map at a much finer resolution than would a local government who require a general overall view of the situation. Thus, through the development and implementation of a suitable methodology, the main aim of this thesis is to determine the applicability of an AI methodology with which such a map could be produced relatively quickly through the incorporation of a range of different data sets.

#### **2.6 CONCLUSIONS**

This chapter has provided an introduction to the multifaceted problem of land degradation and soil erosion. The various methods and techniques used to map erosion have also been discussed, from the traditional modelling and remote sensing approaches to the more recent and less conventional AI methods. Through the detailed introduction of the study area the geomorphological sensitivity of the region has been illustrated and the subsequent need for an erosion map proposed.

The following chapter introduces the subject of semi-arid and badland geomorphology, and related scale and threshold concepts. Furthermore, it details the erosion processes operating in southern Spain including an in-depth discussion of the physical and chemical processes related to soil dispersion, and provides a comprehensive review of erosion studies undertaken in the region.

## Geomorphological Theory and Soil Erosion Processes, Modelling and Mapping

## **3.1 INTRODUCTION**

This chapter investigates badland and semi-arid geomorphology. Badland and semiarid studies are important because the principle operative factor, soil erosion, is detrimental to land productivity and often requires some level of management. This chapter begins by attempting to comprehensively review the fundamental underpinnings of geomorphology. It also highlights the importance of recognising and working at the appropriate spatial and temporal scale. Badland geomorphology is briefly introduced and reviewed, followed by a review of surface and subsurface erosion processes. Particular reference is given to an area in the Almería province. Soil dispersivity is an important element of badland geomorphology and is extensively discussed to provide a better all-round understanding of the physical and chemical processes associated with subsurface erosion.

Soil erosion modelling is discussed with a review of a number of the key models along with techniques that have been used by various authors who have determined, quantified or mapped soil erosion. Finally, soil erosion *risk*, *hazard* and *potential* are reviewed to allow a better understanding of an often vague terminology.

## **3.2 SCALES OF INVESTIGATION AND THRESHOLDS**

The determination of a suitable scale to observe and understand the landscape is crucially important in the development of a successful geomorphological investigation. As a consequence, scale has formed a key element in geomorphological debate. Schumm and Lichty (1965) outlined distinctions between cause and effect in landscape evolution are dependent upon the time span and spatial extent of the geomorphic system. The resolution used in any research/investigation should depend upon the problem being investigated (Schumm, 1991). Delcourt and Delcourt (1988) stated that a successful research design first determines the scale at which the phenomenon of interest occurs and then defines appropriate methods of analysis to determine the patterns and processes operating. Based on this statement, a scale paradigm possessing a range of domains can be seen in Figure 3.1, highlighting four broad scales of investigation. The scales range from the largest, the mega-scale, to the macro-scale, the meso-scale and finally the smallest, the micro-scale. The scales are simply generalised domains and any given study can be within one of them or in some instances cross scale boundaries to ensure an appropriate scale of investigation.



Figure 3.1: Spatial-temporal domains for different research scales (Delcourt and Delcourt, 1988).

Schumm and Lichty (1965) recognised the need to identify appropriate timescales in which to study geomorphic systems. They classified time into three broad time spans; cyclic, graded and steady time (see Figure 3.2). The longest time span is cyclic time,
or geologic time, and spans an entire cycle of erosion. Graded time refers to a short span of cyclic time when a dynamic equilibrium exists for a short period within components of the system or a small area within the system. Steady time span can occur when over, very brief periods, none of the variables associated with a system change. Here, the relationship between space and time scale is important because phenomenon that do not appear to show change over time at one scale can be observed to change over time significantly at another. Table 3.1 demonstrates the status of drainage basin variables as an example during decreasing duration time spans. Drainage basin variables are arranged in a hierarchical manner of increasing degrees of dependence and attempts to demonstrate the influence of differing time spans. The status of each variable is determined as independent, dependent or not relevant, based upon the variable under consideration. Thus, Schumm and Lichty (1965) proposed that depending on the time span involved, time may either be an extremely important variable or of little significance when attempting to understand and study landforms.



Figure 3.2: The changes in channel gradient during cyclic, graded and steady time (Schumm and Lichty, 1965).

The concept of thresholds is closely related to that of scale, as they will vary significantly in both space and time. Bull (1980) regarded thresholds as a balance

between opposing tendencies. Campbell and Honsaker (1982) proposed that that the limits between equilibrium and disequilibrium within a system (e.g. periods of change and non-change) are defined and determined by thresholds. Schumm (1973) divided geomorphic thresholds into two separate categories; intrinsic and extrinsic. Intrinsic thresholds are regarded as inherent to the system and may be exceeded due to the cumulative internal stresses and are not triggered by external events (e.g. pore-water pressure changes on a slope). Extrinsic thresholds however occur as a result of an increased force or stress arising from an external event outside of the system (e.g. seismic activity or flooding). The determination of thresholds is often difficult (Campbell and Honsaker, 1982) and will be influenced by the sensitivity of the system (the propensity of the landscape to change), implying instability in the system and the consequent possibility of sudden irreversible change taking place (Thomas, 2001).

Drainage Basin Variables	Status of variables during designated time spans		
	Cyclic	Graded	Steady
1. Time	Independent	Not relevant	Not relevant
2. Initial Relief	Independent	Not relevant	Not relevant
3. Geology (lithology, structure)	Independent	Independent	Independent
4. Climate	Independent	Independent	Independent
5. Vegetation (type and density)	Dependent	Independent	Independent
6. Relief or volume of system above base level	Dependent	Independent	Independent
7. Hydrology (runoff and sediment yield per unit area within system)	Dependent	Independent	Independent
8. Drainage network morphology	Dependent	Dependent	Independent
9. Hillslope morphology	Dependent	Dependent	Independent
10. Hydrology (discharge of water and sediment from system)	Dependent	Dependent	Dependent

**Table 3.1:** The status of drainage basin variables during time spans of decreasing duration (Schumm and Lichty, 1965).

It is also important to acknowledge the time dependent nature of thresholds, threshold exceedance and system recovery, and the associated implications for system stability (Ritter *et al.*, 1999). Some systems therefore may be able to handle instabilities up to a critical point, as they may not exceed threshold boundaries (low sensitivity landscape units). If the recovery time is sufficiently long then the system will remain in a state of equilibrium, indefinitely (transient). However, if a further disruptive event occurs before the system can fully recover then the internal resistance may not be able to oppose change and will develop a new equilibrium (inter-transient). Figure 3.3 represents these concepts.



Figure 3.3: Geomorphological thresholds and reaction and relaxation. (A) Event causes system to react but completely recovers. (B) Event causes system to react and does not recover and reaches a new equilibrium. (C) Transient form where the system has sufficient time to recover from threshold events. (D) Inter-transient form where the system does not have sufficient time to fully recover from threshold events.

# 3.3 BADLAND AND SEMI-ARID GEOMORPHOLOGY

The term 'badland' is commonly used to describe areas of densely gullied landscapes, where vegetation is sparse or absent and useless for agricultural purposes (Bryan and Yair, 1982). They usually develop as a consequence of contributing factors that include climate, geology and lithology, and lead to the development of a highly dissected and complex landscape. Extensive badlands occur globally and in susceptible areas such as Alberta, Canada. Hodges and Bryan (1982) highlight the role of overland flow and runoff in their development and in particular the importance of individual soil surfaces and lithology and their relationship to overland flow. Numerous investigations have been undertaken to better understand badland geomorphology, some of which are detailed in Table 3.2.

Author	Location	Investigating
Campbell and	Alberta,	Morphological change occurring in the
Honsaker (1982);	Canada	badlands over a period of time.
Campbell (1982; 1989)		
Imeson et al. (1982)	Northeast Morocco	The influence of physical and chemical aspects of soil properties in badland development.
Drew (1982)	Saskatchewan, Canada	Subsurface pipe erosion in the Big Muddy badlands.
Torri and Bryan (1997)	Tuscany, Italy	Badland evolution in Tuscany through micropiping processes.
Sirvent et al. (1997)	Ebro Basin, Northeast Spain	Rates of erosion in badland areas in the Monegros region, Spain.
Boardman et al. (2003)	Great Karoo, South Africa	Development of badlands and gullies and overall land degradation in an area within the Great Karoo.

Table 3.2: Some studies investigating badland geomorphology and processes.

This thesis concentrates on environments in Southeast Spain (see Chapter Two) and in particular the Almería province that has been the subject of numerous investigations due to the coverage of badland areas, where they occupy extensive areas of the landscape (Harvey and Calvo, 1989; Harvey et al., 2001; Spivey, 1997). López-Bermúdez and Romero-Díaz (1989) identified and mapped the locations of susceptible marls in Southeast Spain. Harvey and Calvo (1989) identified areas within the province that were *total* badlands, *partial* badlands and deeply dissected softrock areas that are potentially susceptible to badland development (Figure 3.4). Total badlands are those where entire drainage systems are eroded as are the divides between the gullies. Partial badlands in contrast still have the divides between gully systems intact. Due to the complex nature of the region a number of different processes occur independently or in combination, leading to the development of badland landscapes with differing morphological characteristics. As discussed previously in section 2.4.5, this contributes to the geomorphological sensitivity of the region, as the lithologies are highly susceptible to erosion due largely to their low erosional resistance and poor structure.

Numerous investigations have highlighted the importance of piping in the development of badlands within the region. However, this contrasts with other global regions where literature indicates that Hortonian processes dominate. Although numerous investigations have highlighted the importance of piping in the region, including Harvey (1982), Imeson and Verstraten (1985), Alexander *et al.* (1996), Calvo-Cases and Harvey (1996), Faulkner *et al.* (2000), Faulkner *et al.* (2003a) and Faulkner *et al.* (2003b), the extent to which the piping process occurs varies spatially. Faulkner *et al.* (2000) commented upon the variation in the dominant process

operating in three different badland sites. These and other papers concerned with the Almería province are discussed in more detail in the following sub-section.



Figure 3.4: The location of total badlands, partial badlands and dissected softrock areas susceptible to badland development (Harvey and Calvo, 1989).

# 3.4 SOIL EROSION PROCESSES IN SOUTHEAST SPAIN

Soil erosion processes occur in a number of different environments and are often highly complex. Soil erosion is simply the process of detachment of individual soil particles, their transport and subsequent deposition (Rosewell *et al.*, 1991). Erosion can be split into three major categories; surface erosion, subsurface erosion and mass movements. Surface and subsurface erosion (piping) largely operate at the same spatial scales as one another, whereas mass movements tend to occur on a much larger scale. As a result mass movements have not been considered in this study.

### **3.4.1 Surface Erosion**

As mentioned previously, erosion is simply the detachment of soil particles and in its most common form operates at the soil surface. The principle controlling factor of surface erosion is overland flow derived from a precipitation event. It occurs once the infiltration capacity of the soil mass is exceeded during or after a prolonged precipitation event. Over time, overland flow becomes channelised and concentrated within small rills in the soil surface determined and controlled by factors such as slope angle, length and roughness (Stiegeler, 1979). However, sheet flow or wash precedes the development of rills and is simply the movement of water across a slope surface generated after the onset of precipitation. The term 'sheet' implies a smooth planar surface, however, the water is seldom of uniform depth due to the microtopography of the hillslope surface (Summerfield, 1991). The erosion potential of surface wash is largely a factor of the characteristics of the soil surface, vegetation cover, the slope gradient and the routing of water at the micro-level.

There are two general bases for our understanding of gully and badland development. The first outlined by Horton (1933, 1945) and elaborated by Strahler (1958) relating gully and badland morphology to surface erosion, and the second involves subsurface erosion through piping (Harvey, 1982). Horton (1933) proposed that precipitation is partitioned so that one part infiltrates the soil and the other part goes rapidly as overland flow (Figure 3.5). The Horton hypothesis has been dominant as the traditional process responsible for gully erosion and development. Bryan and Yair (1982) highlighted the fact that the extremely high drainage densities of badland areas often are regarded as evidence of the dominance of overland flow. Ideal conditions for the generation of the Hortonian process are often found in arid and semi-arid regions, where a bare soil is often present combined with the development of surface crusts (Ward and Robinson, 2000). Furthermore, a combination of sparse vegetation, steep slopes and surfaces with relatively low infiltration rates often associated with badland environments are often assumed to contribute to Hortonian overland flow (Bryan and Yair, 1982).



Figure 3.5: Simple representation of the Horton hypothesis (Adapted from Ward and Robinson, 2000).

### Rainsplash

With respect to surface erosion, Chorley and Schumm (1984) highlighted the importance of combined processes in soil detachment and subsequent erosion. Rainsplash is an important element of surface erosion processes occurring in areas where vegetation cover is limited or not present at all (Summerfield, 1991). As numerous studies have shown (Morgan, 1978; Pedersen and Hasholt, 1995), rainsplash is a significant erosive agent moving soil particles both upslope and downslope. However, when combined with overland flow the combination of water on the slope plus raindrop impact an intermediate level of erosion occurs. The water present on the surface dissipates the energy of the falling raindrops and thus reduces their ability to dislodge particles from the soil mass, however the particles that are

removed from the soil are easily entrained and transported by the surface water. Evans (1980) suggested that the two processes operating together are more efficient at moving soil particles than when they occur independently.

The effectiveness of rainsplash can be significantly reduced on clay materials where it can lead to particle compaction and surface sealing (Kuhn and Bryan, 2004). Therefore, rainsplash would seem to be a limited erosion process in badland environments as surface crusting makes material resistant to detachment (Bryan and Yair, 1982) yet contributes indirectly as compaction leads to surface sealing, a reduction in infiltration and a subsequent increase in overland flow.

# 3.4.2 Subsurface Erosion (Piping)

Subsurface erosion has received little literary attention until relatively recently (e.g. in Spain from the early 1980s). Subsurface erosion is known by different names in different places. European literature typically refers to piping, whereas Australian and New Zealand literature may refer to tunnelling (Boucher, 2002; Hosking, 1967). Piping is relatively common in semiarid environments (Parker and Higgins, 1990), where it is related to rilling and gullying and the development of badlands (Bryan and Yair, 1982). The process largely relies upon soil geochemistry factors, especially the amount of swelling clays present in sodic soils (soils containing high levels of exchangeable sodium) (Parker and Jenne, 1967) and topographic influences (See Chapter Two). If appropriate conditions suit, subsurface processes can seriously influence the shape of the landscape and can become the dominant operative erosive process.

Piping is a complex process and arises when substantial volumes of subsurface lateral throughflow pass through a dispersive subsoil (Crouch, 1976; Baillie et al., 1986). Piping is generally best developed in the presence of swelling clays, desiccation cracks, seasonal high rainfalls, differentially permeable layers of soil, steep hydraulic gradients along with a base level of erosion, and a suitable outlet (Jones, 1981). Torri et al. (1994) highlighted the dominance of pipe erosion over surface erosion processes in a study investigating the mechanisms of erosion in a badland area of Tuscany, Italy. They found that the materials were dispersive in nature: readily broken down on contact with water. Sediment rates were found to be higher near pipes as a result of physico-chemical slaking (disintegration of soil) and subsequent erosion. Subsurface erosion has also received much attention in the Dinosaur Park badlands of Alberta; De Boer and Campbell (1990), Bryan et al. (1984) and Campbell (1989) identified the importance of piping as an erosive process in the Albertan badlands and many subsequent studies have attempted to quantify it. Furthermore, Harvey (1982) and López-Bermúdez and Romero-Díaz (1989) investigated the role of piping in the development of badlands and gully systems in Southeast Spain. The studies recognised the importance of piping in the erosional development of some gully systems in the region. However, details relating to the physico-chemical processes involved have not always been fully documented due largely to a lack of understanding and confusion relating to the issue. When viewing a gully the assumption is often made that overland flow is the responsible process, however, detailed inspection of gully form can distinguish subtle differences that can be indicative of a collapsed pipe (e.g. the long profile is seldom uniform as would be expected from those developed through surface process alone). However, more recent research that takes into consideration soil science factors has revealed new insights

into the process and subsequently provides a more rounded understanding (Harvey, 1982; Alexander and Calvo, 1990; Alexander *et al.*, 1994; Sumner, 1995; Faulkner *et al.*, 2000).

Such is the sensitivity of a landscape to piping, dependent on local geo-chemical soil factors, that the nature and extent of pipes can be highly localised. The importance of site geochemical factors is highlighted by Faulkner et al. (2000) for the understanding of badland development in three areas in the Almería province of Southeast Spain. The study concluded that physico-chemical properties of three different badland sites are a useful tool for characterising piping behaviour. The study revealed the importance of subsurface erosion in the development of the Mocatán badlands of Southeast Spain, and how the dominant erosive process differs in two other badland areas studied (Vera and Tabernas). It was observed that the clay content is likely to strongly influence the propensity for which subsurface processes may occur. This is further illustrated by Faulkner et al. (2003b) in Figure 3.6, showing the variation in materials with both high and low clay contents. The difference in process occurs as a direct result of the fact that whilst the dispersal of clays normally slakes and seals the subsurface horizons (Naidu et al., 1995), in soils with low clay percentages this may not occur. This is simply a result of the fact that the material will disperse, yet will not render the material impermeable, as there is too little clay to do so. Therefore, the soil mass will merely deflocculate (break down) and because the clay is the only binding agent this can cause the complete destruction of the soil structure (Alexander et al., 1996).



(a) Materials with high clay content (b) Materials with a low clay content **Figure 3.6:** Process dominance domains for the role of site regulators. Domain controls R: resculpting; O: organic amendments; G: gypsum amendments (Faulkner *et al.*, 2003b).

Figure 3.6 also highlights the importance of the relief ratio, which is an important controlling factor for understanding subsurface erosion. Where hydraulic gradients allow, large pipes can develop and subsequently collapse to form large gullies as discussed by Faulkner *et al.* (2000, 2003b). The point is visualised in Figure 3.7, emphasising the importance of landscape morphology for subsurface processes. In surfaces with convex morphologies with an infiltrating surface and a substantial hydraulic head, subsurface erosion can occur, particularly in materials of low-bulk density. Therefore, it is relatively commonplace to see subsurface erosion processes operating behind terrace walls and other similar features in landscapes where physicochemical situations suit. In the Mocatán badlands of Southeast Spain, the capture and rejuvenation of the Rio Aguas (see section 2.4.5) has led to the over-steepening of slopes, coupled with the susceptibility of the dispersive marls provides the ideal setting for subsurface erosion processes to operate. Figure 3.8 highlights the effects that piping can have on the landscape, creating a dense gully network in a slope face.



Figure 3.7: A simple diagram highlighting the piping potential in convex morphologies.



**Figure 3.8:** The Mocatán badlands demonstrating extensive piping activity leading to the creation of an extensive gully system.

### Soil Dispersivity

The discussion here highlights the importance of a range of soil characteristics with regards to the soil dispersion process. As a consequence, problem soils and associated clay dispersion cannot be distinguished simply in terms of a particular Exchangeable Sodium Percentage (ESP) value or Sodium Adsorption Ratio (SAR) value, and thus relationships are sought between various characteristics and are discussed here.

Soil chemistry can have a profound influence upon a soil's characteristics and behaviour when in contact with water, such as during a precipitation event or irrigation practices, and understanding soil chemistry is important for understanding piping. One such characteristic is the result of the presence of excess sodium as an exchangeable cation, and there are large areas of the world where soils are consequently adversely affected (Sumner, 1995). The primary processes responsible for the degradation of such soils are slaking and dispersion that can lead to both the physical and chemical break down of a soil (So and Woodhead, 1987). From a landuse perspective this is a most undesirable characteristic as it can lead to hardsetting surfaces, rapid surface ponding or runoff, severe erosion, cloddy cultivation surfaces, poor crop establishment and shallow water and root penetration (Powell et al., 1995). A number of factors however can influence and control the extent to which a soil may be dispersive. Dispersion occurs when sodium cations cannot satisfy the exchange complex, a characteristic usually found in soils lacking in magnesium and calcium, and the following discussion will attempt to aid the understanding of the various processes involved.

It is well known and documented that the presence of some double-layer clay minerals are particularly sensitive to sodium on the exchange complex as it causes them to swell and disperse (Faulkner *et al.*, 2000). The dispersion of clays in soils is strongly influenced by the nature of the exchangeable cations and the amount of electrolyte present (Quirk and Schofield, 1955; Shainberg *et al.*, 1981). However, other factors can strongly influence soil dispersion and are consequently used to indicate the feature such as pH, ESP (Equation 1), exchangeable Ca:Mg ratio, bulk density (Powell *et al.*, 1995) and soil organic matter (Churchman *et al.*, 1995). In very

### KINGSTON UNIVERSITY LIBRARY

simplistic terms, soil dispersion occurs in cohesive soils when the repulsive forces between clay particles exceed the attractive forces (Bell and Walker, 2000). This is the case when sodium is the prominent adsorbed ion (Brady and Weil, 1999). As previously stated, a number of soil characteristics can influence the soils ability to flocculate and thus render it more susceptible to the dispersion process.

Prior to a full in-depth review of soil dispersion, it is important to understand how a sodic soil can be identified. The determination of a soils Cation Exchange Capacity (CEC) (Equation 2) is seen as one of the most reliable methods of identifying sodic soils (Murphy, 1995). In addition to this, the ESP, is widely used as an indicator of a soils susceptibility to dispersion (Rycroft *et al.*, 2002). Further to the CEC and ESP, the SAR (Equation 3) has been used to infer the equilibrium relation between soluble and exchangeable cations (Richards, 1954). The SAR is simply the proportion of water soluble sodium to calcium and magnesium in the soil (Davis *et al.*, 2003) and is often the parameter of choice when dealing with the sodicity of irrigation water or soil solution as the presence of salts can cause problems when determining the sum of the cations as in the case for ESP (Sumner, 1995).

$$ESP = \frac{(100 \times ExchangeableNa)}{CEC}$$
(Equation 1)

$$CEC = \sum (ExchangeableCa + Mg + K + Na + Al)$$
 (Equation 2)

$$SAR = \frac{Na^{+}}{\sqrt{\frac{Ca^{++} + Mg^{++}}{2}}}$$
(Equation 3)

Through the use and determination of the parameters discussed above, a great deal of research has been undertaken in an attempt to more fully understand soil sodicity and the associated problems. Gerber and Harmse (1987) produced a chart for assisting the determination of dispersive potential using the CEC and ESP meq/100 g clay (Figure 3.9). The chart is recognised as one of the most reliable of the chemical methods used for assessing a soil's potential dispersivity (Bell and Walker, 2000). From Figure 3.9, soils with an ESP above 15%, the critical limit, are classified as dispersive; the extent to which is controlled by the CEC. Elges (1985) proposed a threshold of 10%, above which soils that have had their free salts leached are prone to dispersion. ESP values as low as 6% have been identified as critical (Northcote and Skene, 1972). However, 15% is recognised as the standard level used for determining saline-sodic soils according to the USDA (Richards, 1959) who have carried out extensive work on the subject area. In practice obtaining a critical value above or below which a soil can be regarded as sodic is arbitrary, as a number of other factors have to be taken into consideration.



Figure 3.9: Cation exchange capacity against exchangeable sodium percentage for soil dispersivity classification (Gerber and Harmse, 1987). Where: VD, very dispersive; HD, highly dispersive; D, dispersive; M, marginally dispersive; ND, non-dispersive; CD, completely non-dispersive.

Electrical Conductivity (EC) also has a strong influence on a soils dispersion. EC is based on the concept that the electrical current carried by a salt solution under standard conditions increases as the salt concentration of a solution increases, and is a common way to measure soil salinity (Sparks, 1995). Quirk and Schofield (1955) investigated the effect of electrolyte concentration on soil permeability following up on work undertaken by Christiansen (1947) who had found that water of low electrolyte content resulted in soil surface sealing. Quirk and Schofield (1955) highlighted the fact that the electrolyte concentration must be increased in soils with an increasing ESP in order to maintain their flocculation. This is represented in Figure 3.10. A similar relationship highlighting the importance of electrolyte in the understanding of soil deflocculation was proposed by Kamphorst and Bolt (1976).



Figure 3.10: Factors affecting the 'threshold concentration' curve (Quirk and Schofield, 1955) adapted by Sumner (1995).

Further influential factors include soil pH and organic matter, and as with those factors discussed above, a great deal of research has been undertaken so as to understand their role. Organic matter can influence the susceptibility of potentially dispersive clays, as organic carbon as well as root and microbial filaments can contribute to the stabilisation of aggregates (Sumner, 1995). However, there is considerable disagreement concerning the effects of organic matter on the dispersion process in sodic soils. Loveland *et al.* (1987) found that dispersion ratios were strongly inversely correlated with organic matter in sodic soils with relatively low ESP values, yet Gupta *et al.* (1984) found that organic additions actually increased the dispersion of soils at high SARs.

A final characteristic of soils, the pH, can also have a profound influence upon soil processes. In soils possessing a high pH, the ability of organic matter to bind and assist in soil flocculation can be reduced, even to the extent where the organic matter begins to dissolve as in many saline-sodic soils.

A great deal of work has been carried out since the recognition of sodic soils and the associated problems that they encompass. Nonetheless, the problem of defining what characteristics a sodic soil should possess has not yet been resolved satisfactorily so as to give a universally accepted definition (Sumner, 1995). The boundaries defining the dispersive nature of a soil are not entirely clear. It was stated earlier that a value of 15% is the generally accepted ESP level above which soils are considered dispersive depending upon the level of electrolyte present. However, Rengasamy *et al.* (1984) identified slightly different boundaries, using the SAR and EC. A range of classes were distinguished relating to their relevant dispersive nature and can be seen in Figure 3.11. The boundaries were drawn after work was carried out on 138 samples of red-brown earths in Australia and, as can be seen from the graph, soils are potentially dispersive at extremely low SAR values, and the EC has to be relatively high so as to maintain flocculation. However, mechanical shaking was required for dispersion to

occur in class 2a soils. In soils with a relatively high EC level, SAR values above 3 are identified as dispersive and do so without mechanical shaking spontaneously, therefore making the management of such conditions difficult.



Figure 3.11: A classification scheme for the prediction of dispersive behaviour of red-brown earths (Rengasamy *et al.*, 1984) adapted by Faulkner *et al.* (2000).

### 3.5 SOIL EROSION MODELLING AND MAPPING TECHNIQUES

In order to manage badland and semi-arid environments it is an advantage to be able to represent the spatial distribution of soil erosion. The ability to determine whether areas are at potential risk from erosion and to map these at an appropriate spatial scale is valuable. A number of soil erosion modelling methods and mapping techniques exist and have been widely used in a range of environments. As outlined in section 2.3 the spatial extent of soil erosion has been determined using one of three main approaches, or a combination of them, process models, remote sensing and field based investigations.

# **3.5.1 Erosion Process Models**

Erosion process models are designed around independent variables or predictors that are believed to influence the extent to which erosion processes operate. They are either site specific (developed for a specific area/region), or generic and can be applied across a range of environments. This section will briefly outline some of the models in global use and some of studies where they have been applied.

The most widely used erosion model applied over large areas is the Universal Soil Loss Equation (USLE) (USDA, 1978). The equation takes the form:

### A = RKLSCP

(Equation 4)

43

where A: Average annual soil loss in tons per acre;

- R: Rainfall erosivity factor;
  K: Soil erodibility factor;
  L: Slope length;
  S: Slope;
  C: Cropping factor;
- P: Conservation practice factor.

The USLE has since been adapted to become the Revised Universal Soil Loss Equation (RUSLE) which contains several improvements to the original approach. The RUSLE uses the same basic factors as the USLE, but improvements have been made through the way of the computing of the soil erosion factors (Shi *et al.*, 2004). Millward and Mersey (1999) modified the RUSLE further for use within a mountainous region in Mexico and to incorporate GIS. The research demonstrated the potential of the approach in such environments as well as its wider applicability. However, the need for more research into the modelling requirements within different environments was recognised as a fundamental objective for future research.

The US Department of Agriculture (USDA) also developed the Water Erosion Prediction Project (WEPP), a physically based model developed for the quantitative prediction of erosion in small to medium sized basins (Flanagan and Nearing, 1995; Nearing *et al.* 1989; Wright and Webster, 1991; Hairsine and Rose, 1992a, b). A number of investigations have been undertaken to evaluate the predictive ability of the model such as Savabi *et al.* (1995), Tiscareno-Lopez *et al.* (1995), Ghidley and Alberts (1996), Zhang *et al.* (1996b) and Soto and Díaz-Fierros (1998). The latter study found that in general the predictions showed reasonable agreement with the observed values in the study area.

As with WEPP, the European Soil Erosion Model (EUROSEM) is process and event based, considering fundamental hydrologic and erosion processes (de la Rosa *et al.*, 1999). The model has a modular structure simulating erosion by using a water and sediment routing scheme and has been widely evaluated (Folly *et al.*, 1999). Cai *et al.* (2005) used the model in the Three Gorges Dam area and highlighted its ability to simulate runoff but encountered problems in determining sediment concentration and soil loss in a single event.

A number of further models exist and have been briefly summarised in Table 3.3.

Model	Summary	Reference
TOPMODEL	<ul> <li>Catchment Scale.</li> <li>Calculates hydrological processes in order to predict sediment yield.</li> <li>Runoff only.</li> </ul>	Beven and Kirkby (1979) Beven <i>et al.</i> (1984)
CREAMS	<ul> <li>Field/Plot scale.</li> <li>Calculates runoff and chemical transport from agricultural systems.</li> </ul>	Knisel (1980)
ANSWERS	<ul> <li>Catchment scale.</li> <li>Simulates hydrological processes during and after a rainfall event.</li> </ul>	De Roo (1993)
EROSION- 2D/3D	<ul> <li>2D - Slope scale.</li> <li>3D - Catchment scale.</li> <li>Calculates rainfall induced soil erosion and deposition for single storm events.</li> </ul>	Schmidt <i>et al</i> . (1999)
CORINE	<ul> <li>Continental Scale – Mediterranean Region.</li> <li>Based on a simplification of the USLE.</li> <li>Determines erosion risk at a 1km resolution.</li> </ul>	Briggs and Giordano (1995)
LISEM	<ul> <li>Catchment scale.</li> <li>Simulates the hydrology and sediment transport during and immediately after a rainfall event.</li> </ul>	De Roo et al. (1996)

Table 3.3: Summary of various soil erosion models.

# **3.5.2 Mapping Techniques**

A contrasting approach to soil erosion modelling is to map soil erosion spatially, and a range of techniques and approaches exist. The spatial resolution and extent of mapping investigations can vary significantly, from very simplistic ground based surveys where erosion is physically mapped at relatively small scales, to more complex studies using remote sources such as remote sensing and aerial photography over more extensive scales.

## Remote Sensing and Aerial Photography

Remote sensing offers the ability to study phenomena at scales that would otherwise be impossible. Various sensors possess different spatial, temporal and spectral resolutions, and the associated advantages with such have led to its growth and development in erosion studies.

Servenay and Prat (2003) explored the application of remotely sensed imagery and black and white aerial photography for the determination of erosion in Mexico. The study found that erosion in the region was strongly related to land use changes that had occurred, and that Spot satellite imagery allowed the discrimination between different types of erosion, which would not have been possible with black and white photography. Metternicht and Zinck (1998) developed an erosion map for a small study area in Bolivia using the Synthetic Aperture Radar (SAR) JERS-1 and Landsat TM imagery, and highlighted the benefits of the synergistic use of remotely sensed imagery.

Pickup and Nelson (1984) mapped soil erosion status using Landsat MSS in an arid region of central Australia, and found that the methodology provided a quick and simple means of mapping erosion status and landscape instability. Finally, Singh *et al.* (In Press) analysed the use of the NOAA/AVHRR data for the calculation of environmental degradation in central Brazil. As erosive processes change both the physical and chemical properties of a soil, monitoring such changes through time help identify and analyse the processes. It was determined that using such data provides a useful means of calculating soil colour and vegetation indexes and thus is quite helpful for the determination of soil erosion.

# Field Mapping

The most traditional approach to mapping and assessing soil erosion is through fieldbased investigations. In comparison to remote sensing, field mapping has a range of advantages and disadvantages. The advantages are that they can provide levels of spatial resolution that would otherwise not be possible from remote sensing or photographic techniques, however, field mapping is very time consuming and may be intensive as outlined by Evans (1992).

Field based assessments of erosion are variable. For example, erosion pins can be used to accurately measure soil loss from individual locations. Saynor *et al.* (2003) used erosion pins to determine the extent of soil loss at 49 different sites in a catchment in northern Australia, and Arnáez and Larrea (1995) at 118 sampling points in La Rioja, Spain. Although such work provides very detailed investigations typically related to specific study sites, it may be difficult to infer erosion characteristics for wider areas.

### 3.5.3 Synergistic Methods

A synergistic methodology recognises that no one method may be entirely appropriate to provide information concerning erosion at a suitable scale and resolution. A synergistic methodology may incorporate two or more different approaches or techniques in an attempt to better understand and determine soil erosion. For example, de Jong *et al.* (1999) used the soil erosion model for the Mediterranean (SEMMED) to simulate soil loss in two experimental locations. The model incorporates the use of a DEM, Landsat TM remotely sensed imagery, a soil map and field data, and accurately produced erosion maps at the regional scale. The use of remote sensing imagery to provide inputs for process models is relatively commonplace. Mongkolsawat *et al.* (1994) used Landsat TM data to provide a land use coverage to be used with the USLE in Thailand, and Lu *et al.* (2004) used imagery in combination with the RUSLE in Brazil. Haboudane *et al.* (2002) combined the use of remote sensing and DEMs to determine areas susceptible to erosion and degradation in the Guadelentin Basin, Southeast Spain. Floras and Sgouras (1999) mapped erosion in central Greece and incorporated field-based studies on individual plots, DEMs and Landsat data, and found the methodology successful.

### 3.6 SOIL EROSION RISK, HAZARD AND POTENTIAL

The application of soil erosion information to applied scenarios requires that we can determine soil erosion parameters that are meaningful to landscape managers. The application of soil erosion *risk*, *hazard* and *potential* are three such parameters that can be used. They are defined as follows:

- Soil erosion *risk* is a product of probability and loss, and is the conventional method by which it is calculated (Bondi, 1985; Middleton, 1999; Smith, 2001).
- A *hazard* is a naturally occurring or human-induced process with the potential to create loss, and *risk* is the actual exposure of something of human value to a hazard (Smith, 2001; Wright, 2003).
- Soil erosion *potential* is the inherent risk of erosion irrespective of current land use or vegetation cover whereas *actual risk* is the risk of erosion under current conditions. This approach is supported by Ellis (1997), taking it a stage further proposing that if only physical factors are taken into account when modelling erosion then *potential risk* is close to the *actual risk*. However, if social factors are considered, then the output is one of soil erosion *hazard*.

A number of previous investigations have attempted to develop erosion risk models or erosion risk maps in order to better understand and determine the associated hazards with erosion processes. However, the definition of risk is highly contentious within both the physical and social sciences and therefore requires the basic outlining of their fundamental definitions.

The end-users choice of parameter will be largely dependent on the questions asked; for example, a soil erosion hazard map may not provide the necessary information on which to base management decisions because social factors (e.g. risk) may be more important.

# **3.7 CONCLUSIONS**

As outlined in the introduction to this chapter, a number of topics have been introduced and comprehensively discussed regarding the underpinnings of geomorphology as a science, largely based on spatial and temporal concepts. The soil erosion processes responsible for land degradation in the Almería province of Spain have been documented, the numerous modelling approaches that have been developed have been reviewed, and a discussion attempting to clarify the common confusion between the terms *soil erosion risk, hazard* and *potential* used within this study has been clarified.

# 4

# **Artificial Intelligence Classifiers**

#### **4.1 INTRODUCTION**

This chapter will briefly introduce the concept of Artificial Intelligence (AI) and provide a comprehensive review of both Artificial Neural Networks (ANNs) and Decision Tree Classifiers (DTCs). The review will include the general theory behind the techniques, their respective advantages and disadvantages and their past and current uses. The chapter will draw to a close with a general overview of the two classifiers and some brief concluding remarks.

With the vast improvements and developments of data sources and collection techniques over recent years, substantial amounts of geographical information covering extensive areas have become readily available for use by academics and commercial organisations. This information comes from various sources, such as remotely sensed satellite derived data, obtained through various sensors which detect radiation from the Earth for different parts of the electromagnetic spectrum (Heywood *et al.*, 1998). Aerial photographic imagery is now also readily attained for many areas of the Earth's surface. Even imagery for those which are not covered, can be acquired relatively easily and cheaply. These data can now be stored, managed and manipulated readily with the development of technology and software applications. Coupled with this has been the development of Geographic Information Systems (GIS), for processing, analysing, and visualising the growing amounts of digital spatial data (Goodchild *et al.*, 1996). This has allowed scientific investigations to

proceed that were hitherto impossible over a relatively short period of time, and at a spatial scale of study that was previously unpractical.

With the development of data sets and tools such as GIS it is relatively simple to integrate data from one or more sources and subsequently perform some form of statistical analysis or model application to describe landscape form and process. Landscape processes and responses, change, can be linear and/or non-linear. A linear response is one where for a graded and measured input (A) the landscape has an equal and proportionally graded output (B). A non-linear response is where a graded and measured input provokes a functional response that may not necessarily be linearly proportional; such that the response, at first sight, may appear to be somewhat random and unexplained (see Figure 4.1).



**Figure 4.1:** The contrasting responses between two different systems where; A is an input and B is an output, with a linear process on the left and a non-linear process on the right.

Conventional statistical techniques tend to make prior assumptions with regards to the nature of the data, and tend to assume a linear response between inputs and outputs, and are coupled with parametric restrictions (Spellman, 1999). However, because many environmental processes are apparently non-linear, the use of linear

classification systems is inappropriate when attempting to fully describe landscape behaviour, and tend to perform poorly.

In order to combat this, an approach that would allow the use of multi-source data sets in research activities concerned with investigating non-linear patterns and processes is required. This has largely been achieved in the form of machine learning techniques, a sub-discipline of AI, such as ANNs and DTCs. Machine learning implies that we use the 'machine' (computer) to process the input data in such a way as to determine the functional relationships that provide a set of observable outputs. These types of classifiers have a number of distinct advantages over more traditional counterparts, some of which have been outlined in Table 4.1. This chapter will introduce and comprehensively review both classification techniques. It will highlight their applicability in determining landscape change for use in this thesis and discuss their associated limitations.

Data Classifier	Advantages	Disadvantages
Artificial Neural Networks	<ul> <li>Ability to weight significance of variables.</li> <li>No prior assumptions.</li> <li>Ability to make supervised and unsupervised classifications.</li> <li>Easily integrate multi-source data.</li> <li>Ability to map non linear functions.</li> <li>Each component within the structure of the network is responsible for only one small part of the input-output mapping operation.</li> <li>Ability to recognise subtle patterns in training data that may be missed by conventional statistical analysis</li> </ul>	<ul> <li>Neural networks are opaque and understanding the internal workings is complex.</li> <li>The identification of an appropriate architecture.</li> <li>The determination of the training algorithm to be used.</li> <li>Determining the optimum learning parameters for the algorithm.</li> <li>The identification of appropriate stopping conditions for the training process.</li> </ul>
Decision Tree Analysis	<ul> <li>Provide a simple representation for propositional knowledge that can be used for decision making and classification.</li> <li>Robust to errors.</li> <li>Relatively easy to implement.</li> <li>Can assist in our understanding of various problems.</li> </ul>	<ul> <li>Problems of over-fitting.</li> <li>In many domains, not all of the attribute data will be known for every example.</li> <li>When an attribute has a large number of possible values it may potentially be a problem (Multivalued attributes).</li> <li>Continuous valued attributes may also be a problem as they may contain a large or infinite set of values</li> </ul>
Linear Regression	<ul> <li>Easy to implement.</li> <li>Highlights potential relationships between two variables.</li> </ul>	<ul> <li>Only incorporates two variables.</li> <li>Assumption of linearity.</li> </ul>
Multiple Regression	• Reveal relationships between several independent and one dependent variable.	<ul> <li>Assumes a normal distribution.</li> <li>Assumption of linearity.</li> </ul>
Discriminant Analysis	<ul> <li>Multivariate technique.</li> <li>Limited use of categorical dependent variables permitted.</li> </ul>	• Assumption of linearity.
Logistic Regression	• Categorical dependent variable can be used.	• Can only be used in cases of dichotomous dependent variable.
Cluster Analysis	<ul> <li>Can be used as an unsupervised classification technique.</li> <li>Highlights potential relationships between two variables.</li> </ul>	Only incorporates two variables.

Table 4.1: Advantages and disadvantages of various classification techniques.

### 4.2 ARTIFICIAL NEURAL NETWORKS (ANNs)

### **4.2.1 Introduction**

A neural network is a series of connections linking individual nodes (neurons) composing input neurons, 'function determining' neurons, and output neurons (Figure 4.2) A neural network can be conceptualised as an information-processing system that has certain performance characteristics in common with biological neural networks (Fausett, 1994; Campbell, 2002). They are a computational (statistical) mechanism that is able to represent and compute a "mapping" from one multivariate space of information to another, given a set of inputs and outputs but without necessarily knowing any details of their relationship (Goh, 1995; Yang *et al.*, 2003). They have a number of distinct advantages over more traditional statistical techniques making them appealing for a wide range of operations, and have consequently been used in a variety of applications where traditional statistical methods are traditionally employed (e.g. maximum likelihood classifiers) (Warner and Misra, 1996).



Figure 4.2: An example of an artificial neural network with a 5:3:1 architecture.

ANNs attempt to imitate some of the functions of the mammalian brain. The first attempts to model the fundamental cell of the brain, the neuron, were made by McCulloch and Pitts (1943), who described a simple neural network. It was only from the mid- to late 1980s that the realisation of the abilities of neurocomputing came

about. Neural networks work on the same principle as the human brain, but on a much smaller and simple scale. They typically comprise a number of simple processing units, or neurons, linked by weighted connections, in specified layers (Tveter, 1998), collectively termed the architecture (Figure 4.2).

The input layer simply accepts the input data, with one node representing each element (independent variables), in the case of Figure 4.2 there are five inputs, thus representing five independent variables. Perhaps the most popular network architecture in use today is that of the Multilayer Perceptron (MLP) (Hontoria et al., In Press). In situations where neural networks have been implemented, it is usually to a large extent MLPs which have been used (Brown et al., 1998; Boyd et al., 2002; Spellman, 1999; Gong, 1996), usually as a result of their ability to model functions of almost arbitrary complexity, with the number of layers, and the number of units in each layer, determining the function complexity. In a MLP the input layer distributes the data to the nodes within the 'hidden' layer; in our example here there are three nodes. Here, the weighted sum of the inputs are computed, passed through a transfer function, and forwarded to the node(s) in the output layer (Brown et al., 1998). The transfer function enables the network to perform in a non-linear manner, using either a sigmoid or stepped function (see Figure 4.4) The node(s) in the output layer then calculate the sum of the incoming weighted values, passes them through a transfer function and generates a result. The nodes are arranged in a layered feed-forward topology, and the network has a simple interpretation as a form of input-output model, with the weights and thresholds (biases) the free parameters of the model. A 'threshold' is the point that defines one decision outcome from another.

# **4.2.2** Activation Functions

The basic operation of an artificial neuron involves a number of steps so as to compute an output signal to be passed on to the neurons in the next layer of the network. Each unit, node, receives signals from its input links and computes a new activation level (from the activation function) that it sends along each of its output links (German and Gahegan, 1996; Russell and Norvig, 1995). As can be seen in Figure 4.3 the process can be split into two components. Firstly, a linear component or input function, net<sub>i</sub>, calculates the weighted sum of the input values of that unit. This is followed by a transfer function, or activation function g, which simply transforms the net input to a neuron into its output value (Fausett, 1994). Different models are obtained by using different mathematical functions for g, some of which are illustrated in Figure 4.4.



Figure 4.3: The internal process of the activation unit (Russell and Norvig, 1995).

A range of activation functions exist, and are explained briefly below.

*Linear Function*: The application of a simple linear activation function will constrain the performance of a network and will not be able to learn non-linear patterns. Step Function: The incorporation of the step function allows the hidden layers to learn in a non-linear manner, usually between the limits of 0 and 1. Working on a simple threshold concept.

Sine Function: Similar to the step function, but operates between the limits of 1 and -1. In such cases where positive and negative values are produced the networks tend to be trained faster.

*Binary Sigmoid Function*: Can be scales to have any range of values appropriate for a given problem.

Bipolar Sigmoid Function: A type of binary sigmoid function, but with the specific limits of 1 and -1.

The activation function is part of what determines the dynamic state of a neural network system, and the ability of the network to change the weights enables the system to adapt and change (McCord Nelson and Illingworth, 1991). Furthermore, it acts as a range limiter, allowing the discrimination of inputs if they are arranged within some appropriate range (McCord Nelson and Illingworth, 1991). Without these limits the input data could produce outputs at the extremes of the range and would create problems when determining one from another. A good example would be that of the binary sigmoid function seen in Figure 4.4 (c) where inputs of  $\chi$ <-1 or  $\chi$ >+1 would yield similar outputs in either the positive or negative direction. Therefore,  $\chi$ =10 would be very similar to  $\chi$ =100.



Figure 4.4: Common activation functions and their accompanying mathematical functions.

#### 4.2.3 Training Methods and Algorithms

Deciding whether or not to use an AI approach to describe landscape change behaviour is dependent on the questions the user (e.g. the landscape manager) is asking. For example, a small specific scale problem (e.g. a localised landslide failure) may be readily managed through quick and empirical field observations. However, more extensive areas and problems may render simple field/empirical relationships difficult, inefficient (e.g. non-linear) and/or unproductive. In these cases a neural network (AI) approach may be appropriate. The decision to use a neural network model should therefore be made on the basis of employing an appropriate tool for a task (Abrahart and White, 2001). The user must be satisfied that an ANN can outperform alternative methods in its use of resources (e.g. time, money and efficiency) and deemed relevant to a study.

The user must carry out a number of vitally important steps before the model can be applied to a specific problem. Once the type of network has been chosen the architecture of the network must be selected. A great deal of literature has been written regarding the determination of an optimal network topology for a problem, such as the number of layers and neurons within each layer. Blum (1992) suggested that the number of hidden nodes should fall be between the number of input and output nodes. Moreover, Berry and Linoff (1997) proposed a general rule of thumb whereby the hidden nodes should be no more than twice that of the input layer, and Bourquin *et al.* (1997) applied Kolmogorov's theorem stating for *n* inputs a hidden layer of 2n + 1 nodes is sufficient. However, many studies have successfully used more hidden nodes than in the input and output layer. Unfortunately determining the optimal number of layers and the number of neurons within those layers is still largely
believed to be a process of trial and error (Ghiassi and Saidane, 2005; Spellman, 1999), problem specific (Wang *et al.*, 1994) and of all the configuration issues has the fewest theoretical guidelines (Jarvis and Stuart, 1996).

Once an architecture has been established, the network requires training. The ANN learning process can be grouped into two specific techniques, 'supervised' and 'unsupervised'. An example of a supervised classification would be where information relating to dependent variables is available (e.g. an erosion map); an unsupervised classification however does not contain any knowledge of the dependent variable, and using the independent variables outputs are clustered into similar groups to one another. Supervised networks are presented with both dependent and independent variables. Unsupervised neural networks are essentially data classifiers (Openshaw and Openshaw, 1997) and are largely employed when there is no available training set on which the network could be trained. Maybe perhaps little is known about their functional relationships in the first place. Therefore, it is the job of the neural network to find results and relationships with little prior knowledge of the process in question.

There are numerous training algorithms, all of which are used in conjunction with varying network types. The back-propagation neural network (BPNN) trained by the generalised delta rule (Rumelhart *et al.*, 1986) (Equation 5) has been successfully utilised in many fields, especially for pattern recognition due to its learning ability (Dai and MacBeth, 1997). Back-propagation is probably the best known training algorithm for ANNs (Tveter, 1998) and comprises a number of distinct advantages and disadvantages over other training algorithms, such as Conjugate Gradient

Descent, Levenberg-Marquardt, Kohonen, and Probabilistic training algorithms. During the training of the ANN, the learning algorithm aims to estimate optimal values of the weights by minimising an error function, usually the sum of the squared error between the target and network predicted output (Wang *et al.*, 1994).

The generalised delta rule commonly used to train back-propagation ANNs can be written as:

$$\Delta w_{iik}(n+1) = \eta \delta_{ik} o_{ii} + \alpha \Delta w_{iik}(n)$$
 (Equation 5)

Where Dai and MacBeth (1997) state that  $\eta$  is the learning rate and  $\alpha$  is the momentum rate.  $\Delta w_{ijk}$  is the change of the weighted connection between  $n_{ij}$  ( $n_{ij}$  is the *j*th node in the *i*th layer) and  $n_{i+l\,k}$ .  $o_{ij}$  is the output of  $n_{ij}$ , and  $\delta_{ik}$  is the change in error as a function of the change in the network input to the  $n_{i+l\,k}$ . Finally, (n+1) indicates the (n+1)th step.

Back-propagation works by simply back-propagating the error of an output based upon the actual, and a targeted output. The algorithm works in two distinct stages, a feed-forward phase, and a back-propagation phase. The performance of a BPNN trained by the generalised delta rule is influenced by such factors as the size and structure of the network as well as training parameters, such as the learning rate and momentum. Defining such parameters can be problematic and few rules actually exist. The convergence speed (iteration number) of the training procedure is dependent upon both the learning rate  $\eta$ , and the momentum rate a (Dai and MacBeth, 1997). The learning rate controls the amount by which weights are adjusted during training. The momentum on the other hand considers the weight changes associated with previous learning steps and uses the knowledge to slow down ineffective oscillations (Salomon and van Hemmen, 1996).

The difficult process of selecting both a learning rate and a momentum rate that will allow an ANN to perform at its optimum has been recognised in the literature, for example Dai and MacBeth (1997) and Maier and Dandy (1998). The goal is to identify both a learning and momentum rate that will allow the network to converge relatively quickly, and within some specified error threshold. The learning rate should be selected to be as large as possible so that the network reaches its error threshold quickly. However if it is too large, then the learning process may become unstable and oscillate (Pao, 1988) or in some cases the network may fail to converge at all (Maier and Dandy, 1998). The way to increase the learning rate without leading to oscillation is to use a momentum term, where step size is increased or decreased if error rates are reduced or increased respectively (McClelland and Rumelhart, 1988). This also needs to be relatively large so as to increase the speed of the algorithm when a number of consecutive steps are made in the right direction. However, as in the case of the learning rate, if the chosen value is too large then the network will not converge and too small a value would mean that the time taken to converge would be substantial. Dai and MacBeth (1997) found the learning parameters  $\eta$  and  $\alpha$  were highly influential in relation to the speed of convergence when using a BPNN to pick seismic arrivals. A relationship was found between the iteration number and the two training parameters, and can be written as follows:

Iteration\_Number = 
$$IN_0 \times \frac{(1-\alpha)}{10 \times \eta}$$
 (Equation 6)

where:  $\alpha$  is the momentum term,  $\eta$  is the learning rate and  $IN_0$  is the iteration number.

Dai and MacBeth (1997) stated that the optimum value for  $\eta$  is likely to fall between 0.6 and 0.7 and  $\alpha$  between 0.8 and 0.9. McClelland and Rumelhart (1988) also reported in most of their simulations that  $\alpha$  should be around 0.9 and  $\eta$  around 0.7. In reality however, these findings can only act as guides, as the ideal values for the two parameters are likely to be problem specific (Maier and Dandy, 1998, 2001).

# 4.2.4 Current Uses of Artificial Neural Networks

The revival of interest in neural networks in the 1980s came largely as a result of the increasing awareness of their many abilities. This awareness however, was only in a very small number of research areas (e.g. in medical use), and even with the great advances in data collection sources, computing power and the networks themselves, their adoption has been limited. In the few cases where neural networks have been employed, the results have tended to be positive to such an extent that they have, in many circumstances, replaced the more traditional techniques (e.g. in remote sensing neural networks are widely preferred to more traditional classifiers when using spectral unmixing models (see Liu and Wu, 2005). These uses have ranged from simple pattern recognition, such as the automatic recognition of hand-written characters (see Le Cun et al., 1990) to medical problems, where symptoms are identified and a diagnosis offered (see Anderson, 1986a;b). Refenes et al. (1994) examined the use of an ANN as an alternative to classical statistical techniques for forecasting stock market activities. They found that even simple neural learning procedures such as the back-propagation algorithm far outperform current statistical techniques in forecasting accuracy terms. This comes as a result of the limitations of classical statistical techniques that reach their limitations in applications with nonlinearities in the data set (Refenes et al., 1994). These are just a small number of examples where ANNs have provided promising advances in the understanding of highly complex processes. However, their adoption has been less noticeable in geography and related sub-disciplines, despite the attempts to promote their benefits by a small number of researchers (Openshaw and Openshaw, 1997). As a result, the possibility that ANNs may hold to geographical issues is immense, and this promise has only recently become apparent.

One research area that has realised this potential is remote sensing, where ANNs have been successfully used in numerous applications (Paola and Schowengerdt, 1993, 1995; Luoto and Hjort, 2005). For example, ANNs have been used in both supervised and unsupervised classifications, including land cover classifications (Civco, 1993; Downey et al., 1992; Jarvis and Stuart, 1996), image inversion (Smith, 1993) estimating soil physical properties (Chang and Islam, 2000) and cloud classifications (Lee et al., 1990). This potential was largely recognised due to a networks ability to learn and classify data. In remote sensing, various sensors either on satellites or aircraft operate at relatively large spatial scales, producing vast amounts of data relating to the reflectance of various surfaces of the Earth. In more traditional techniques, the wavebands that were deemed the best for a given classification were used in a statistical analysis, such as regression, with the others being disregarded. However, ANNs have the ability to weight the significance of the independent variables used (Boyd et al., 2002), and thus allow the integration of all available wavebands. In theory this allows for better, more accurate analysis, without the added implications of choosing the correct wavebands to use.

Foody *et al.* (1995a; 1995b) recorded significantly higher classification accuracies when using an ANN classification than from the discriminant analysis when attempting to classify crop type on fields in Feltwell, in the UK. The outputs gained from the neural approach gave accuracies as high as 98 percent, but it was only during the use of non-normally distributed data that the differences between the two techniques became apparent as the neural approach continued to produce highly accurate outputs.

Moreover, Boyd *et al.* (2002) tested a number of neural networks in comparison with regression analysis to estimate forest cover in the Pacific Northwest USA. It was found that the neural network approach was most attractive in this role as a result of the advantages they hold over more traditional techniques and not simply as a result of an improved accuracy. ANNs make no assumptions about the data, can easily integrate multi-source data and can weight the significance of the discriminating variables. Overall, ANNs have been used to classify remotely sensed data to accuracies that are generally comparable to or higher than those derived from conventional statistical classifications (Hepner *et al.*, 1990; Medina and Vasquez, 1991; Short, 1991).

Other uses of neural networks in geographical fields have been extremely limited, despite the potential that they offer. As seen in the remote sensing applications, the advantages that they hold have been sufficient to warrant their wider applicability in the field of remote sensing. Boyd *et al.*, (2002) detailed that one of the major attractions is that they offer a powerful means to analyse complex data sets without making assumptions about them, such as linearity's. A great deal of real-world

systems and processes are non-linear and thus, simple linear models fail to capture the essence of the underlying phenomenon (Spellman, 1999).

ANNs are in principle capable of solving any non-linear classification problem, provided the network contains a sufficiently large number of free parameters (i.e. hidden units and/or connections) (Wu, 1997). Abrahart and White (2001) found this most beneficial in modelling sediment transfers, where it was apparent that the neural solution provided a tighter fit to the data with a pronounced reduction in outliers when compared to a multiple linear regression technique. The neurocomputing technique was also found to have further advantages for such an application. As described earlier, each component within the overall structure of the network is responsible for only one small part of the total input-output mapping operation. Therefore, each individual data input can have no more than a marginal influence with respect to the complete solution, thus allowing for substantial fault tolerance (Abrahart and White, 2001). This consequently allows the model to generate reasonable results, even in cases of incomplete data, and/or data containing substantial 'noise'. This statement is supported by Gong (1996) who also found networks to be, to some extent, tolerant of noise when using them for geological mapping and Hepner et al. (1990) who were using a minimal training set for classification purposes. Furthermore, Skidmore et al. (1997) highlighted their ability to identify subtle patterns in input training data, many of which would be missed by conventional statistical analyses. This is the advantage of the training process, whereby the network actually learns and therefore does not require a prior knowledge of the functional relationships between the input and output variables primarily due to their non-linear abilities at handling complex data patterns (de la Rosa et al., 1999). Hence, any underlying relationships, no matter how small,

will have a good chance of being picked up and acknowledged by the network, thus allowing for the development of a model more closely representing the underlying patterns and processes being investigated.

One area in which ANNs have only recently attracted attention in geographical and environmental analysis is looking at the activities involved in land degradation, especially that of soil erosion, and subsequently attempting to produce accurate, representative models. To date, this attention has at best been limited, and consequently while the suitability of ANNs for this type of application is still not fully understood, it appears promising. De la Rosa et al. (1999, 2000) used neural network applications in the development of a model to assess agricultural soil erosion, potential vulnerability, and the impacts upon crop productivity in Spain. In these studies ANNs along with DTCs were used to formulate, calibrate and perform a validation analysis on the ImpelERO model. Within this model, networks were applied to capture the interactions between the land and management qualities in order to produce one output, a vulnerability index to soil erosion (de la Rosa et al., 1999). Variables such as runoff erosivity, relief hazard, soil erodibility, crop protection, tillage translocation and productivity influence were used for training purposes. The results of the first of the two studies were promising, and highlighted the ability of ANNs in recognising the main interrelationships of the input parameters, and subsequently could accurately reproduce the soil erosion vulnerability patterns that were observed in the field. The model was later applied to 20 selected benchmark sites in western Europe to quantify the soil erosion vulnerability under several crops, the impact of soil erosion on crop production and the optimum management strategies (de la Rosa et al., 2000). The model was found to give much more accurate results

than the CORINE model (see Chapter Three, section 3.5.1), but would require further studies if the model were to be put into practice over larger geographical areas because the model was developed and trained with regards to one specific geographical area and its specific characteristics, and these may differ in other areas.

#### 4.2.5 Accuracy of Artificial Neural Networks

The accuracy of the output(s) of the neural network is dependent upon the input variables and the complexity of the questions being asked: The more complex the problem (e.g. landscape classifications over large geographical area encompassing many different land types or 'detailed' mapping at smaller scales) the less likely that the accuracy of the output will represent the 'real world'. This is clearly scale dependent and involves a trade-off between the goals and expectations of the analysis and exactly what the practical application of the analysis is. Some users may be willing to accept a lesser accuracy of output in order to minimise input complexities and provide required spatial coverage, whereas other users may demand higher accuracies that necessitate more detailed inputs and greater training. Therefore, acceptable accuracy is subjective and there is no one level at which the neural networks and decision trees can universally be said to be right or wrong; success must be measured specifically on a case-by-case basis.

Ellis (1997, 2002) incorporated the use of multi-source data sets along with neural networks and GIS when investigating the application of machine learning techniques for erosion modelling. The studies identified the current inability of traditional mathematical erosion models when applied over complex regions and therefore investigated the usefulness of both neural networks and decision trees in the

modelling of areas prone to land degradation through the process of soil erosion. The study took place in an area of approximately 93km<sup>2</sup> in New South Wales, Australia. The data inputs for the models were Landsat Thematic Mapper bands 1-7 for use as a surrogate for both vegetation cover and broad land management practices, a soil map, tree cover map and a DEM produced from 1:25000 topographic maps. The output data used to train the models was obtained from a 1:25000 soil erosion map of the area. It was found that even using only a 5 percent training set, the overall accuracy of the network ranged from 91 to 92 percent. In this case, although this may seem relatively high it was deemed poor as the networks classified only one or two prominent classes, largely ignoring the smaller classes, and thus obtaining high accuracy levels.

Harris and Boardman (1990) applied an expert system approach for the prediction of soil erosion in the South Downs in Sussex, England. The expert system was trained using a total of 334 erosion events recorded during a six year monitoring exercise, and the results were comparable to those obtained using more traditional process-based models such as the USLE and CREAMS. In an extension of the research, Harris and Boardman (1998) increased the data set to 450 events and applied both an ANN and an expert system to the problem. The results indicated that although the accuracy of the AI techniques matched that of the process models, the accuracy of the expert system was not improved considerably with the addition of further data.

# 4.2.6 Disadvantages of the Neural Network Approach

As with any modelling technique the ANN approach has several drawbacks and related problems, some of which are relatively easy to overcome and others that tend to be more challenging. One of the major problems is that the vast majority of those who could possibly apply neural networks to their research are unfamiliar with the technology and thus would not incorporate it into their work. Secondly, but not such a problem today, was the slow processing speeds and technological hindrances that made such an approach unappealing. However, the vast majority of the evidence against using neural networks comes as a result of their associated challenges for their meaningful application to real-world problems, rather than from the technique itself. There are a number of such points, but the main argument against using neural networks are frequently based on the premise that they are nothing more than blackbox models which provide no scientific explanation or theoretical understanding of underlying fundamental processes of the real-world (Abrahart and White, 2001).

Networks are often described as *opaque*; and it is not easy to look inside them so as to ascertain how they produce their results (Minsky, 1991). This is a distinct disadvantage of connectionist systems when taken in comparison with other more traditional, rule-oriented approaches to artificial intelligence (Alexander and Mozer, 1999). Results gained through ANNs in any study can therefore only be validated statistically as the network output and the desired output can be compared in situations when the actual or true output is known. In order for the user to gain confidence in the technique, it is important to be able to obtain an understanding of the internal workings of the network and, of course, to make inferences into the processes or activities being studied. For example, an ANN may produce an accurate model representation of the real-word for a simple, small-scale problem, but be unsuitable (without development) when scaled-up to a larger, more geographically complex landscape. This is a well documented drawback within the literature, for

example Boyd *et al.* (2002) found opacity to be a problem when attempting to identify the relationship between remotely sensed data and the phenomenon of interest. There is a great deal of research written contributing to determining possible ways of countering the problem of opacity (see Alexander and Mozer, 1999; McMillan *et al.*, 1991 and Hinton, 1989). The solutions used are largely based on rule extraction techniques, approaches that provide visibility into a networks operation by casting its weights in the form of symbolic rules (Alexander and Mozer, 1999). Such methods involve the long and drawn out process of removing variables one by one so as to determine the influence of a given variable upon the workings of the network as a whole.

A second concern when using neural networks is the identification of an appropriate architecture. The specification of an appropriate network topology is a key issue because it governs the capability of the network to provide an adequate approximation of the input-output relationships (Benediktsson *et al.*, 1990; Hepner *et al.*, 1990; Lee *et al.*, 1990; Wang *et al.*, 1994). However, there is no assurance that the optimum topology is identified and the network may therefore perform below its capabilities. Too small a network will be incapable of representing the desired function, yet too big and it will be able to memorise all of the examples by forming a large lookup table but is unable to generalise (Russell and Norvig, 1995). Therefore, a network with a large number of hidden nodes may perform well reclassifying the learning sample, but would perform less well when faced with new cases (Lees, 1996). Thus, it is critical to develop an architecture large enough to perform a given task, but small enough as to not hinder its ability to generalise (Hinton, 1990).

Capability not only depends upon the network structure but also on the learning algorithm that is used to determine the weights of the interconnections. There is a great deal written with regards to how to combat this problem in terms of rules and laws (see Wang *et al.*, 1994). Many believe the network architecture is problem specific and extremely difficult to handle more efficiently than through a time-consuming process of trial and error (Lengellé and Denœux, 1996; Brown *et al.*, 1998).

The final problems to be discussed here all concern the ANN training process. Before the training of a network is to take place, a sufficiently large sample needs to be obtained, which is also unbiased towards the population that it represents. Once again no laws as such exist, and the sample size is largely deduced by the user. It is vitally important that the sample used to train the network is representative of the population from which it has been taken, given the particular sensitivity of machine learning algorithms to unrepresentative learning samples (Lees, 1996). Various studies have incorporated data that has not been representative and have consequently had difficulties when using neural networks. Spellman (1999) for example used ANNs for the prediction of surface ozone concentrations in the UK. The results from the study revealed that although they predicted ozone concentrations more accurately than the conventional regression-based model, they tended to underestimate high air pollution episodes. This problem could be due to the fact that very poor air quality events are rare and as a result there are few examples in the training dataset (Spellman, 1999) and the model therefore performs poorly when predicting such events. Such a problem was also documented by Ellis (2002) when using machine learning techniques for erosion modelling. Although the model reached prediction accuracies as high as 92

percent, small classes were often poorly predicted. This was put down to the fact that such classes were not well represented in the training samples, and thus some form of stratified sampling would best be employed to enable the prediction of the underrepresented classes.

The choice of training algorithm is also vitally important in the training of a neural network. A great deal of literature has been written regarding the many various algorithms, their advantages and their drawbacks, where they have been successful and where they have not (see McCord Nelson and Illingworth, 1991; Russell and Norvig, 1995). However, some algorithms have added complications in that they require learning parameters to be set before they can begin to train a given network. A prime example of this is the back-propagation algorithm, which has been documented in some detail throughout this chapter. The learning parameters, such as the learning rate and the momentum will influence how the network learns. The learning procedure might be unstable and the weights of the back-propagation neural network may become trapped in a local minimum and fails to discover lower values from the global minimum. However, Rumelhart et al. (1986) discovered that very rarely did networks get stuck in poor local minima that were significantly worse than the global minimum, and was only encountered in networks possessing just enough connections to perform a task. It usually occurs when the momentum causes the algorithm to increase its speed if a number of consecutive steps change the weights in the right direction. As with the majority of these issues, a great deal of work has been concerned with identifying optimal values for these parameters (see Gorman and Sejnowski, 1988; Dowla et al., 1990; Fausett, 1994; Dai and MacBeth, 1997), yet it is largely a process of trial and error which is often time consuming and inconsistent.

Finally, it is important to determine at which point training should stop. The importance of a large enough sample size on which to train a network has been well emphasised here, but it can be the case that an ANN becomes over-trained. The traditional method of control involves the implementation of fixed stopping conditions based on out-of-sample performance (Bishop, 1995). These stopping conditions usually work on the premise of the associated errors within the network, and the user can specify an error level below which training will be suspended, or when the error fails to improve by a given amount over a given number of epochs. It is important that stopping conditions are set otherwise the predictive capabilities of the network will suffer and the performance will be poor. This point is further illustrated in Figure 4.5, where it is evident that at a given point the error level in the test data stops falling, and may even begin to rise. Thus at this point training should cease and over-training can be avoided.



Figure 4.5: Error in the training and test sets during training (Picton, 2000).

# 4.2.7 Summary

In summary, ANNs have an enormous potential as a classification tool possessing a number of inherent advantages. However, as has been demonstrated within the previous sections of this chapter, they also have a series of shortfalls that may limit their applicability. Therefore, prior to implementing an ANN, the user must be satisfied that it is the most appropriate tool and be aware of all of the associated limitations.

#### **4.3 DECISION TREE CLASSIFIERS (DTCs)**

# **4.3.1 Introduction**

A Decision Tree Classifier (DTC) is an attempt to use a 'stratified' or 'layered' approach to the problem of distinguishing between different groups or classes (Mather, 2001). Friedl and Brodley (1997) defined DTCs as a classification process that subdivides a data set into smaller subsets through a series of questions or tests at each branch in the tree. This works through the simple partition of the space  $\chi$  of possible observations into sub-regions corresponding to the leaves, since each example will be classified by the label of the leaf that it finally reaches (Ripley, 1996). Figure 4.6 shows the workings of a very basic DTC. The tree can be separated into three principle components; firstly, the root node at the very top of the tree, followed by a number of internal nodes known as splits, and finally terminal nodes or leaves. Unlike ANNs, each node in a DTC may have two or more descendant nodes but only a single parent node from which they are grown (Friedl and Brodley, 1997). The process is considered to be a chain or union of basic decisions based on the results of a series of questions as opposed to a single, complex decision, making the

classification process more comprehensible (Pal and Mather, 2003; Safavian and Landgrebe, 1991).



**Figure 4.6:** A decision tree classifier. Each box is a node at which tests (T) are applied to recursively split the data. A, B and C are classes assigned to each observation (Friedl and Brodley, 1997).

# **4.3.2 Decision Tree Algorithms**

The main objectives of DTCs are; to correctly classify as much of the training sample as possible, to generalise beyond the training sample so that unseen samples could be classified with as high of an accuracy as possible and be easy to update and have as simple a structure as possible (Safavian and Landgrebe, 1991). In order to satisfy these general criteria, a number of steps have to be taken. Firstly, a DTC structure has to be selected. In the past this done manually, using spectral plots (Pal and Mather, 2003), or by experts with substantial knowledge of the subject (Cawsey, 1998), however, recent developments have led to automatic design methods (e.g. the computer may create many decision trees and select the most appropriate for the classification task). Secondly, the choice of feature subsets to be used at each internal node has to be determined. Finally, the choice of the decision rule to be used within the network at each node has to be selected. These choices should be made on the premise that each classification task is unique and therefore requires individual considerations when they are being developed.

Constructing a DTC is simple when there is an exact partition of  $\chi$ , however, this is rarely the case, and the classes are found to overlap, often known as a noisy classification problem (Ripley, 1996). In order to account for this, there exist two possible strategies; stopping tree construction early, so monitoring its progress, and secondly pruning the tree after construction. Pruning simply involves the identification and removal of the least reliable nodes (Bradley and Lovell, 1995), thus reducing the complexity of the tree and making it more comprehensible (Pal and Mather, 2003). Ockham's razor insists upon economy or simplicity when choosing between different hypothesis to explain a fact (Rodríguez-Fernández, 1999). It is therefore generally assumed, following Ockham's economy principle, that the simpler a DTC is, in terms of its structure, the more general its applicability will be, and thus its performance on unseen cases will be better than that of a more complex tree. It is commonly agreed that finding new ways to build smaller decision trees is a desirable goal (Berzal et al., In Press), and pruning is one way in which decision trees can be drastically simplified.

Most algorithms that have been developed for decision trees are variations of a core algorithm that employs a top-down, greedy search through the space of possible DTCs (Mitchell, 1997). Such algorithms have been used, and have been the centre of decision tree induction research for a number of years. Probably the most well-known is the ID3 algorithm (Quinlan, 1986) and its successor C4.5 (Quinlan, 1993). There

are a number of alternatives, such as the ACLS algorithm (Patterson and Niblett, 1983) and the ASSISTANT algorithm (Kononenko *et al.*, 1984), both of which are generalisations of the ID3 algorithm.

There are two possible approaches to the induction task of creating a DTC. The first approach would be to generate all of the possible decision trees that correctly classify the training set and simply select the simplest one. This is a viable approach when the task is relatively simple, as the number of trees is indeed finite, but potentially very large. The ID3 in contrast, attempts to classify data when there are many attributes and the training set contains a number of objects, and thus computationally becomes more efficient. In general the ID3 has been found to construct simple decision trees, but the approach that it uses cannot guarantee that better trees have not been overlooked (Quinlan, 1986).

The central choice in the ID3 algorithm is selecting which attribute to test at each node in the tree, as it is important to select the attribute that is most useful for classifying the data first (Mitchell, 1997), allowing for as small a tree as possible. The ID3 algorithm uses an information gain criterion, which can be interpreted as the decrease in conclusion uncertainty resulting from the last test at a given node (Pomorski and Perche, 2001). This relatively simply process is known as *information gain*, which measures how well a given attribute separates the training examples according to their target classification (Mitchell, 1997). The process by which this operation can be undertaken is known in information theory as *entropy*. This is simply a method of determining the purity or impurity of a collection of examples, or data

set. Therefore, given a collection S, which contains both positive and negative examples of a given target, the entropy of S relative to this boolean classification is:

$$Entropy(S) = -p_{\oplus} \log_2 p_{\oplus} - p_{\Theta} \log_2 p_{\Theta}$$
 (Equation 7)

where:  $p_{\oplus}$  is the proportion of positive examples in S and  $p_{\Theta}$  is the proportion of negative examples in S. It is important to know that  $0 \log 0$  is 0 here.

It is often useful to use an example to illustrate this rule. Therefore, suppose S is a collection of 16 examples, of some boolean concept, 11 of which are positive, and 5 are negative. The entropy of S relative to this boolean classification is therefore:

$$Entropy(S) = -\frac{p_{\oplus}}{p_{\oplus} + p_{\Theta}} \log \frac{p_{\oplus}}{p_{\oplus} + p_{\Theta}} - \frac{p_{\Theta}}{p_{\oplus} + p_{\Theta}} \log \frac{p_{\Theta}}{p_{\oplus} + p_{\Theta}}$$
(Equation 8)

For example, if there are a total of 16 classes, with a split of 11 and 5, then;

$$Entropy(S) = (11+,5-) = -\frac{11}{16}\log_2\frac{11}{16} - \frac{5}{16}\log_2\frac{5}{16}$$
$$= 0.896$$

Therefore, the entropy is 0 if all of the members of S belong to the same class, or in other words it returns a low entropy value for subsets with high homogeneity (Kervahut and Potvin, 1996), and the entropy is 1 if the set contains an equal number of positive and negative examples. It is therefore likely to be the case (as in the example following Equation 8), that the entropy will fall somewhere between 0 and 1. The form of the entropy function can be seen in Figure 4.7 highlighting this point.



Figure 4.7: The entropy function relative to a boolean classification (Mitchell, 1997).

The above discussion has been formed around a boolean classification, without regarding a classification with multiple values. In such circumstances, a target attribute that can take on c different values, has an entropy of S that can be defined as:

$$Entropy(S) \equiv \sum_{i=1}^{c} -p_i \log_2 p_i$$
 (Equation 9)

where  $p_i$  is the proportion of S belonging to class *i*. Once calculated, the entropy can be used as a measure to determine the ability of any given attribute to correctly classify the training set. The ID3 system uses a simple rule whereby the attribute that gains the most information is selected (Quinlan, 1986). This measure is known as information gain, and is the expected reduction in entropy after a partition via a specific attribute has been made. ID3 examines all of the candidate parameters and chooses the one that maximises the gain (Pomorski and Perche, 2001), thus selecting the attribute for the root of the decision tree. This process would then be repeated time and again until a tree has been fully built from a given training set. The information gain can be written as:

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$
 (Equation 10)

where Values(A) is the set of all possible values for the attribute A, and  $S_v$  is the subset of S for which attribute A has value v (i.e.  $S = \{s \in S | A(s) = v\}$ ). The equation is merely the entropy of the subset S, and the expected entropy after it has been partitioned by A. Gain (S, A) is therefore the information provided about the *target* function value, given the value of some other attribute A (Mitchell, 1997).

Unfortunately, this approach has been found to be bias, as Kononenko *et al.* (1984) suggested, stating that the gain criterion tends to favour attributes with many values. There have been a few steps taken in order to prevent this undesirable effect, such as Kononenko *et al.* (1984), whose ASSISTANT algorithm took active steps to tackle this problem. It achieves this by ensuring that all tests only have two possible outcomes. This will not be discussed in detail here, but for a more comprehensive discussion see Konenenko *et al.* (1984) and Quinlan (1986). An alternative method however is the *gain ratio* criterion (Quinlan, 1986), which penalises attributes with numerous values by incorporating a *split information*, as follows

$$SplitInformation(S, A) = -\sum_{i=1}^{c} \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$
(Equation 11)

where  $S_1$  through  $S_c$  are the *c* subsets of examples from the partitioning of *S* by the *c*-valued attribute *A*. We can therefore determine the *gain ratio* by using the following equation.

$$GainRatio(S, A) = \frac{Gain(S, A)}{SplitInformation(S, A)}$$
(Equation 12)

The gain ratio criterion will therefore act against the natural bias and deter the selection of attributes with uniformly distributed values. Quinlan (1987) found that when all the attributes are binary, the gain ratio procedure found a considerably smaller decision tree than that produced otherwise. The C4.5 algorithm uses the same technique as ID3, but its main contribution consists of introducing continuous parameters (Quinlan, 1996).

# **4.3.3 Decision Tree Structures**

DTCs algorithms can either have a uniform or heterogeneous set of algorithms to estimate the splits at the internal nodes (Friedl and Brodley, 1997). More traditional approaches tend to use homogeneous hypothesis spaces, such as Univariate Decision Trees (UDTs) and Multivariate Decision Trees (MDTs). A UDT is limited to testing a single feature (Swain and Hauska, 1969), whereas a MDT allows testing of multiple features at a node (Brodley and Utgoff, 1995). As each test in a UDT is based on a single independent variable, splits or thresholds can only be created at right angles to the axis that represents the selected variable (Pal and Mather, 2003). Therefore, as can be seen in Figure 4.8, the decision boundaries created for a UDT can only be parallel to either of the axis as only the outcome of a test applied to a single feature or independent variable at each internal node can be determined (Swain and Hauska, 1969).



**Figure 4.8:** Axis-parallel decision boundaries of a Univariate Decision Tree (Adapted from Pal and Mather (2003)).

Unlike a UDT, a MDT is not restricted to splits of the instance space that are only at right angles to the features' axes (Brodley and Utgoff, 1995). The two are similar, yet the splitting test at each node in a MDT may be based on a number of independent variables (Friedl and Brodley, 1997). Breiman *et al.* (1984) and Utgoff and Brodley (1990) suggest that UDTs may be limited in situations where data sets can only be split and thresholds identified using combinations of features (tests), as opposed to unions of single features, and that MDTs may be better employed. As can be seen in Figure 4.9, in such cases the feature space may be split through linear combinations of features, so as to correctly classify the training data.



**Figure 4.9:** Decision boundaries for a multivariate decision tree classifier (Adapted from Pal and Mather (2003)).

# 4.3.4 Current Uses of Decision Tree Classifiers

Decision trees have a number of advantages over more traditional classification algorithms (Hansen *et al.*, 1996), as well as their machine learning counterparts. This has made them an attractive tool for which they have been employed for a number of classification problems in a number of fields. To begin, they are strictly nonparametric and therefore do not require any assumptions regarding the distributions of the training data (independent variable) (Friedl and Brodley, 1997; Pal and Mather, 2003). This advantage has been recognised within the remote sensing community and is of particular use when a single land cover type is represented by more than one cluster in the spectral space (DeFries and Cheung-Wai, 2000). The efficiency of DTCs, when compared to maximum likelihood methods has been attributed to the non-parametric nature that they possess (Borak and Strahler, 1999; McLachlan, 1992). A maximum likelihood classifier applies a single classification operator to an entire data set assuming that it is normally distributed, and thus performs poorly on nonnormally distributed data.

The opportunity that DTCs have developed has been taken on board in many remote sensing research areas, contributing largely to their rapid growth and development. For example, Yang *et al.* (2003) indicated that by using hyperspectral data as input, they could outperform logistic regression quite considerably when classifying agricultural plots (in Canada). In addition, Rogan *et al.* (2002) found decision trees outperformed a Maximum Likelihood Classification (MLC) approach by about 10 percent in an investigation into various methods for monitoring multi-temporal vegetation change using Thematic Mapper imagery. Friedl and Brodley (1997) found that decision trees consistently classified to higher accuracies than either the Linear

Discriminant Function (LDF) or a MLC when used to produce land cover classifications. Goel *et al.* (2003) stated that DTC algorithms have potential in the classification of remotely sensed spectral data.

Such an advantage can also be extrapolated into further areas of research, where more traditional statistical approaches tend to perform poorly, particularly when the entity being studied is believed to be non-linear. In such cases, the implementation of a decision tree can often reveal non-linear and hierarchical relationships between input variables and the dependent variable (DeFries and Cheung-Wai, 2000; Han and Kamber, 2001), as well as uncovering any structure in the data (Breiman *et al.*, 1984; Safavian and Landgrebe, 1990; Sethi *et al.*, 1990; Brown *et al.*, 1993). Ellis (1997) evaluated the use of DTCs to model soil erosion and found them easy to interpret, offering the ability with which to better understand the relationship between variables. Unfortunately however, the implementation and use of DTCs for use within similar research areas has been extremely limited. However, DTCs have been used to determine soil properties in Australia (Henderson *et al.*, 2005), predict and model pasture productivity in New Zealand (Zhang *et al.*, 2005a, 2005b) and to determine vegetation distributions (Moore *et al.*, 1991).

Unlike other non-linear techniques, in particular ANNs, a decision tree is not a 'blackbox' (Borak and Strahler, 1999) and the internal workings and theory (rules and boundaries) behind the classifier is open to view (Pal and Mather, 2003). This allows our understanding of certain non-linear entities to be furthered in cases where decision trees tend to classify well. This advantage can be seen in numerous studies and in numerous fields of research where such an approach has been employed. In contrast to conventional single-stage classifiers where all of the sample data is tested against all of the classes (a very long drawn out process), a decision tree becomes more computationally efficient as it only tests a sample against certain subsets of classes, thus reducing unnecessary computations (Safavian and Landgrebe, 1991). Decision trees can be trained relatively quickly and are rapid in execution (Gahegan and West, 1998; Fayed and Irani, 1992; Hampson and Volper, 1986). Other techniques, such as the ANN approach, take much longer to train and finally develop a classifier that can be used, as the method has to continually learn and adapt to the training data that it is presented with. The time taken to train a decision tree is therefore significantly faster than neural networks (Wierczorkowska, 2000). Furthermore, there are fewer parameters involved with creating and developing a tree, and the procedure is relatively straightforward in comparison to a neural network methodology.

Another important advantage that DTCs possess in comparison to more traditional classification techniques, unlike parametrically based classifiers, is that they can function without difficulty on data of any statistical type, so the individual attribute dimensions can be a mixture of nominal, ordinal and quantitative data (Gahegan and West, 1998). This allows far greater freedom and a much wider range of applications for such a classifier. This is especially true in cases where more than one data set is considered to be the optimal approach for classifying any particular entity.

In addition to those advantages already discussed, Friedl and Brodley (1997) also found decision trees to be very efficient and robust, and insensitive to noise, making them highly appealing for remote sensing applications. However, German *et al.*  (2002) found that they were noise intolerant, as they depend on the specified sets of data being available at each node for rule resolution. Nonetheless, in their comparison study of two bayesian classifiers (namely the minimum distance-to-mean (MDM) and the MLC, as well as a linear discrimination analysis (LDA) and an ANN), it was found that the decision tree offered the best all-round choice.

#### 4.3.5 Disadvantages of the Decision Tree Approach

As with all classification techniques, decision trees possess a number of inherent disadvantages. A relatively large proportion of these problems are related to the fact that they learn from a training data set presented to them by the user. The characteristics of any data set used to train a supervised classifier has a considerable influence on the accuracy of the resulting classification (Campbell, 1981). To begin, a set of training data needs to be representative of the entity with which it is being used to study. If the training data however is not representative, the decision tree will not perform well when presented with new inputs. Furthermore, in a situation where a small data set is used, the tree may encounter the "Hughes phenomenon" (Hughes, 1968), where the classification accuracy may decrease as the number of features increases for a constant training set size. Such a problem must therefore be considered, but training set size is problem specific, and therefore usually difficult to determine.

It is often the case that the training data will possess noise of some description. The description of objects may include attributes based on measurements or subjective judgements, both of which may give rise to errors in the values of attributes, and some of the objects may even have been misclassified. Many of the branches of a DTC may

reflect chance occurrences in the training data set rather than portraying true underlying relationships (Kim and Koehler, 1995). Quinlan (1986) suggests that in order for a decision tree to handle noise, the algorithm must be able to;

- (i) work with inadequate attributes, because noise can cause a comprehensive set of attributes to appear inadequate.
- decide that testing further attributes will not improve the predictive accuracy of the decision tree. Therefore refraining from increasing the complexity of the decision tree to accommodate noise generated cases.

A similar problem to that discussed above, is that of overfitting. As the training data examples are only samples of all possible instances, it is possible to add branches to a tree that would improve performance on the training set, whilst reducing the performance and limiting the general applicability of the tree on unseen cases (Mitchell, 1997). It is often the case that decision trees constructed from a collection of examples, although accurate and efficient, often suffer the disadvantage of excessive complexity and are therefore incomprehensible (Quinlan, 1987). In such cases, where the trees become complex and opaque, it is useful to attempt to simplify the classifier, usually through a pruning method. Quinlan (1987) investigated the use of four pruning methods to assist in the simplification of complex decision trees, and found that they all achieved a significant simplification, and often an actual increase in classification improvements on unseen cases. This finding is supported by the work of Pal and Mather (2003) who also found that although pruning increased the error on the classification of the training set, it reduced the tree size considerably (from 713 to 231 nodes in this example), and was better at classifying on unseen examples.

Other problems with decision tree classifiers largely relate to the actual design and implementation process. Firstly, the user has to determine which type of decision tree would best classify a given problem. Once a tree has been grown, a post growth procedure such as pruning can then be adopted to simplify the tree and make it more comprehensible, but again which one to choose can be a difficult task. Finally, the number of training patterns has to be determined. This is very often determined by cost, as data acquisition in many instances can be expensive. The size of the training data set however is important, as the classifier accuracy will improve as the set size increases, but only up until a point. Again, such a scenario is problem specific and no general rules can be applied.

#### 4.3.6 Summary

The usefulness and wide-scale applicability of DTCs for a number of research questions has been documented and discussed within section 4.3 of this chapter. DTCs as with other non-parametric classifiers, offer the ability to handle non-linear data sets and have frequently been shown to yield higher classification accuracies than MLCs (Rogan *et al.*, 2002). Moreover, they have the ability to work with multi-source data-sets and provide useful and coherent outputs. Nevertheless, they are reliant upon training data and the growth and determination of an optimal tree is not straightforward.

## **4.4 CONCLUSIONS**

This chapter has introduced the concept of machine learning and has documented and reviewed the two classification techniques to be used in this thesis. Both techniques have unique advantages and disadvantages which are highlighted and supported by the considerable literature. Although both have broadly similar advantages over more traditional techniques, they have been compared and contrasted for a range of applications. Goel *et al.* (2003) for example proposed that although better results were gained using ANNs over DTCs, the formulation of precise rules using the latter approach was both an interesting and useful feature. Tu (1996) supported this suggesting that ANNs may be particularly useful when the primary goal is prediction. However, in situations where this is not the case and an important goal is to determine possible relationships between independent and dependent variables, other methods may better be employed.

These are potentially innovative methods by which soil erosion can be classified and mapped using multi-source data sets. These techniques may offer the ability to better understand operative processes and assist in our all-round geomorphological knowledge.

Prior to the employment of the machine learning techniques discussed here, the development and creation of a data set has to be fulfilled with which the classifiers can be trained. The following chapter, Chapter Five, will document and detail the building of the data set and provide a comprehensive methodology on which the thesis is built.

# 5

# **Research Framework and Data Methods**

#### **5.1 INTRODUCTION**

Previous chapters have considered the issues of soil erosion, modelling and mapping, and Artificial Intelligence (AI) methods to determine landscape elements that are eroding or vulnerable to erosion. This is seen as valuable as an applied tool for landscape managers in order to aid landscape management decisions. This chapter provides descriptions of the research framework and the data methods used to meet the research aims and objectives outlined in Chapter One. A number of components are required to provide output types (as discussed in Chapter Four). The components (independent variables) that determine the propensity for a landscape element to erode are reviewed along with their validation (why any particular component is needed; its validity to the output type) and how they have been calculated and collected. The chapter also discusses the sampling strategy used to provide a training data set for both the Artificial Neural Networks (ANNs) and Decision Tree Classifiers (DTCs), and the implementation of the data used in the training procedures.

## 5.2 METHODOLOGICAL SETTING

The use of AI techniques has been demonstrated in recent years with the fruition of a number of successful studies concerning a wide range of environmental issues (see Chapter Four). This has largely occurred as awareness has increased and the realisation that such techniques can be applied using low-cost data sets, but still provide useful and acceptable results. Furthermore, AI approaches can often handle

multi-source data sets and establish non-linear relationships that may otherwise go unseen between independent and dependent variables and/or undetected by more conventional classifiers.

A primary aim of this study concerns the applicability of two AI techniques for the classification of erosion processes and their spatial extent in the semi-arid landscape of the Almería province, Southeast Spain. To determine the applicability of AI techniques a number of classifications are undertaken using combinations of both dependent and independent variables of landscape erosion as training data. Different combinations of independent variables are used to determine landscape erosion. The methodology investigates what combinations of these are most influential in determining soil erosion with accuracy. This is valuable in a practical sense because the high costs associated with field sampling are frequently acknowledged. Lunetta et al. (1991) stressed the difficult balance in deciding what field data to choose between when sample data is expensive when obtaining a statistically valid data set. The search for methods and techniques that efficiently use low-cost training data as opposed to expensive field collected data is an important debate. An extensive array of data sources are now available that may offer potential replacements or substitutes for variables previously collected in the field. Such data sources have come to fruit largely as a result of the development of Geographical Information Systems (GIS) and associated technologies. There is no universal answer because different applications will have different spatial resolution requirements. For example, a locally-based, small scale study may still discover that field-collection and analysis is a more appropriate method (it may be cheaper and quicker), while an extensive regional

investigation may conclude that AI-based techniques provide a valuable alternative to exhaustive and expensive detailed field mapping.

In order to provide a focus, a number of basic research questions are investigated. The questions are as follows:

- How do ANNs and DTCs perform as classifiers of soil erosion for different levels of classification (e.g. simple (binary) classes to detailed (interval) classes)?
- How do ANNs and DTCs compare with one another as classifiers of soil erosion and where appropriate how do they compare with a more traditional statistical technique (Discriminant Analysis)?
- To what extent does the selection of independent variables influence classifier performance?
- Can ANNs and DTCs further our current knowledge and understanding of soil erosion processes?
- What are the physico-chemical relationships between various soil parameters and what is the applicability of the field sodicity meter developed by the Australian Co-operative Research Centre for Soil and Land Management?

# **5.3 RESEARCH FRAMEWORK**

In order to meet the stated aims and objectives (Chapter One) and to answer the questions above (section 5.2), a solid research framework detailing rationale and justification for the major methodological steps undertaken is required. The subsections here will detail the development of the training data sets, the validation and justification of the dependent and independent variables and the sampling strategies used to collect the data.

#### 5.3.1 Training Data Sets

As outlined in Chapter Four, both ANNs and DTCs learn through the presentation of training data. It is therefore vitally important that the data set is representative (sample of landscape information attributes that adequately describe the population) (Friedl and Brodley, 1997; Lillesand and Kiefer, 2000) so as to allow the classifiers the best possible chance of accurately classifying unseen cases. Campbell (2002) and Hixson *et al.* (1980) suggest that the selection of an appropriate training data set may even be more important than the classifier used.

To produce a supervised classification, both independent and dependent variables are required. The independent variables were derived through a range of sources and can be seen in Table 5.1. As outlined previously, independent variables were compiled from primary and secondary information sources. Primary sources were obtained in the field, and in some cases (e.g. geology) this field evidence is used to corroborate secondary (mapped) information. Secondary data sources are those that are 'remote' from the site (e.g. geology map).

However, the dependent variable (estimate of erosion based on a visual scale) was compiled and determined through fieldwork alone because currently comprehensive soil erosion mapping does not exist for the area (at an appropriate resolution). If one had existed it may have provided the data set from which the dependent variable could have been derived.

#### 5.3.2 Validation of Variables

Independent and dependent variables must be individually validated to establish that (as suggested from established literature) they are important components of landscape change prior to soil erosion modelling (Le Bissonnais *et al.*, 2001). Therefore, the justification for their inclusion will be made here based on prior knowledge of erosion processes.

In this study nine independent variables and one dependent variable were chosen (Table 5.1). Of these, slope angle, aspect and vegetation cover (estimated and classified) can be derived in different ways. They are as follows:

- Slope Angle: The angle of a slope influences its level of stability as well as controlling factors such as runoff velocity and the potential of soil detachment. It has thus been incorporated within numerous soil erosion models, including the Universal Soil Loss Equation (USLE), the Revised Universal Soil Loss Equation (RUSLE) and the Water Erosion Prediction Project (WEPP). Moreover Desmet and Govers (1996) and Mcgregor (1957) recognised the importance of slope in erosion processes.
- Slope Aspect: Slope aspect strongly influences factors that will in turn limit or exacerbate erosional processes. The aspect will control soil moisture retention and will consequently control vegetation growth on that slope. Shrimali *et al.* (2001) incorporated slope aspect for the determination of erosion susceptible areas for a catchment in northern India.
- Vegetation Cover: The extent of vegetation cover is a vitally important factor limiting erosional processes and subsequent land degradation (Cyr et al., 1995, Cerdá, 1999). Plant roots will bind the underlying soil mass and the above ground
matter will protect the soil from raindrop impact and the detachment of soil particles from surface runoff.

- Geology (Lithology): Vrieling et al. (2002) included lithology when developing a methodology for erosion risk mapping, as did Rafaelli et al. (2001) when developing an erosion process model for use within the a catchment in Argentina.
- Sodicity: Sodic soils contain a higher than desirable amount of sodium attached to the clay particles, and when in contact with water tend to swell and disperse (Rengasamy and Bourne, 2001). Shainberg (1990) highlighted that sodic clay soil surfaces can be easily slaked leading to a reduction in surface infiltration rates as pore spaces become clogged with particles. However, in soils containing low clay percentages and organic content, infiltration rates may remain unaffected whilst the soil is destructured due to the loss of its only flocculating agent (Faulkner et al., 2000). Therefore, the determination of soil sodicity allows inferences to be made with regards to its potential erosion characteristics. Furthermore, the determination of soil sodicity at each site is useful as although lithology maps detail known soils sodicity can vary greatly over small distances. López-Bermúdez and Romero-Díaz (1989) discussed the extent of piping and badland development in Southeast Spain, and mapped the extent of marls associated with subsurface erosion. Unfortunately, the resolution of the map is too coarse to be used as a further independent variable here. For a more detailed discussion refer to section 3.4.2.
- *Plan and Profile Curvature*: Plan and profile curvature refers to the curvature of the surface in the direction of the slope and perpendicular to the slope respectively. Profile curvature can assist in the identification of zones of erosion

and deposition, and plan curvature areas of convergent and divergent flow (Pallaris, 2000; Gallant and Wilson, 1996; Haboudane et al., 2002).

- Flow Accumulation: Flow accumulation is an important variable as topography controls the accumulation of water and energy in the landscape (MacMillan *et al.*, 2004) and will thus exert controls on erosive processes.
- Flow Length: Flow length refers to the calculated distance of flow upstream or downstream for any given cell. This parameter has also been used by Shrimali et al. (2001) in the aforementioned study.

Source	Independent Variables	Dependent Variable
DEM Generated	Slope Angle	
(Grid cells of 10 and 20m)	Slope Aspect	
	Plan Curvature	
	Profile Curvature	
	Flow Accumulation	
	Flow Length	
Ground Survey	Slope Angle	Erosion
	Slope Aspect	
	Vegetation Estimate	
Photographs	Vegetation Classification	
Sodicity Meter	Sodicity	
Geology Map	Geology (Lithology)	
Table 5.1: The variables col	lected for the study and their	sources.

As detailed in Table 5.1 a range of independent variables were obtained from DEMs with grid cell sizes of 10 and 20 metres, providing a range of spatial attributes of hydrologic characters (Zhang *et al.*, 1996a). Table 5.2 provides a brief description of the range of attributes acquired to assist in their validation.

Attribute	Description
Slope Angle	Slope identifies the maximum rate of change in value from each cell to its neighbours.
Slope Aspect	Aspect identifies the down-slope direction of the maximum rate of change in value from each cell to its neighbours.
Plan Curvature	Calculates the curvature of the surface perpendicular to the slope direction.
Profile Curvature	Calculates the curvature of the surface in the direction of the slope.
Flow Accumulation	Creates a grid of accumulated flow to each cell, by accumulating the weight for all cells that flow into each downslope cell.
Flow Length	Calculates upstream or downstream distance along a flow path for each cell.

Table 5.2: A brief description of the attributes calculated from the derived DEMs.

#### **5.4 SAMPLING STRATEGY**

The aim of the training stage (see Chapter Four) within supervised classifications is to derive a representative sample of characteristics for each class (Chen and Stow, 2002). In order to achieve this, an appropriate sampling strategy is required to ensure that the sample is representative of the population of information for the study area. This has long proved a challenge in numerous investigations and has subsequently been discussed at great length. A number of difficulties still exist. Chapter Four discussed at length the problem of identifying appropriate training set sizes for any supervised classification technique and identified that only some general rules of thumb exist. For example, conventional probabilistic classifiers require large training sets, typically for each class in the region of 10 to 30 times the number of discriminating variables (Piper, 1992; Swain, 1978). Campbell (2002) suggests between 5 and 10 examples of each class during supervised classifications. Foody (1995) recognised the limitation with classifiers for remote sensing purposes, relevant here nonetheless, suggesting that obtaining extensive data sets is contrary to the overall goal of scaling-up from limited ground data, and furthermore they rarely exist

(Foody *et al.*, 1995a). In this context this is an important point; if the number of individual training cases becomes too large the user may be better served by simply conducting a detailed field survey and dispensing with the AI approach altogether.

Foody *et al.* (1995a; 1995b) and Chen and Stow (2002) discuss the issues relating to training set composition for use within supervised classifications. One of the most important factors is inter-class variations within the training set: equal numbers of each dependent variable may improve classification accuracy particularly in smaller classes. Thus it is important to be aware of the composition of the training set, the characteristics of which should be carefully outlined. Nevertheless, it is not always possible to determine the location of various classes spatially, as for example may be the case in remote sensing investigations, and it is therefore difficult to control training set composition. The problem is highlighted by Welsh *et al.* (1996) when modelling rare species (Leadbeater's Possum) and associated swamping effects.

A range of spatial sampling strategies exist in order to meet this criteria, including random, stratified and random stratified sampling strategies (see Table 5.3). The use of a rigid sampling strategy was decided against. The reasons for this are listed and discussed below:

- (i) It is extremely expensive to spend long periods of time in the field (Chen and Stow, 2002). Therefore, data collection had to be relatively quick and efficient.
- (ii) Some areas within the study area have limited accessibility. For example, an extensive portion of the study area is inaccessible due to the existence of a gypsum quarry (see Figure 5.1). Furthermore, only a limited number of roads exist within the study area, and thus the majority of the sampling had to be

undertaken on foot; and the extremely variable terrain made some locations impossible to reach and others potentially dangerous.

- (iii) The distribution of erosion processes varies spatially, and a sampling technique whereby predetermined sites are sampled disregards this. For example, subsurface processes may only be operating in certain locations and may be missed if sites are randomly selected prior to fieldwork.
- (iv) Areas of extreme erosion are less prominent than are areas showing no appreciable sign of erosion. The use of a rigid sampling strategy could therefore potentially miss or limit the number of cases incorporated into the training set relating to the less geographically extensive processes.
- (v) Training data sets are required to be representative of the entire population being sampled (see Chapter Four). Therefore, through site selection in the field it is possible to control the contents of the training set to some extent. For example, sites with limited vegetation cover, moderate vegetation cover and extensive vegetation cover would all be incorporated. This would be impossible to control using when predetermined locations are selected.

Sampling Strategy	Description
Simple random sampling	Samples selected totally at random within the study area with no predetermined assumptions made.
Stratified sampling	Samples collected systematically based upon some predetermined spatial framework such as a grid.
Random stratified sampling	Samples collected at random within some predetermined locations, such as a grid.
Clustered sampling	Samples taken around a number predetermined locations within the study area.

Table 5.3: Common spatial sampling strategies.

As previously stated, the classifiers require data relating to as many of the field locations as possible accompanied by the range of different outputs (dependent variable). Prior to the extensive fieldwork season, the determination of a suitable measurement of the dependent variable (erosion) was required. This was achieved through a preliminary field excursion where the study area was covered extensively. This enabled the determination of different erosion processes as well as the extent to which they were seen to be operating upon the landscape. This enabled the construction of an erosion risk scale. The scale is discussed in more detail later.

Fieldwork was undertaken during the months of September and October in 2003. The approach that was taken was as follows:

- (i) An area of 40 km<sup>2</sup> (8 km x 5 km) was selected within the DEM within which sampling was to be concentrated. This can be seen in Figure 5.1.
- Using geology (lithology) as a primary variable, an extensive range of sampling points were visited located on different units within the study area. Of the independent variables geology is the only one that completely changes spatially. For example, slope angles of 45 degrees are likely to exist in numerous locations throughout the study area, however, geological units are specifically constrained. Thus, to develop a representative training set samples were required on the different units.
- (iii) As many sampling points were visited as possible within the specified fieldwork period. A total of 520 sites were sampled, the locations of which can be seen in Figure 5.1.
- (iv) Sites were chosen based upon their sampling viability. As previously discussed, extensive areas were not viable to visit either due to safety issues or

simple logistics. However, if a site was deemed viable and could be identified on the DEMs it was chosen and various data collection methods and techniques were subsequently employed to extract the required data.



**Figure 5.1:** The 520 sampling locations draped onto a DEM of the study area (study area is enclosed within the red box, and yellow box is the approximate location of the gypsum quarry).

#### 5.5 ATTRIBUTE ACQUIREMENT

Attribute acquirement concerns how the variables were collected (primary and secondary), and how they are extracted (e.g. DEMs). As indicated by Table 5.1, independent variables were acquired through field-based methods and two DEMs. The following sub-sections detail the methods through which each variable has been collected.

#### 5.5.1 Field Acquired Variables

A number of independent variables along with the dependent variable erosion were measured and determined in the field. Sites were determined using the sampling strategy outlined in section 5.4, however, it was important to ensure that each site fulfilled a range of simple criteria prior to its selection and inclusion within the training data.

A description of how each variable was collected and/or measured is detailed below.

- *Slope Angle*: Recorded using a compass clinometer to the nearest degree. The measurement was taken on any part of the slope that was deemed to be representative of the entire slope (i.e. not an extreme point).
- Slope Aspect: Using a compass-clinometer a bearing of the predominant direction of the slope was recorded to the nearest degree ranging from. 1° to 360°. In cases where the study site was flat, and thus had no dominant aspect, a value of -1 was assigned. In such cases the slope angle was said to be 0° (i.e. flat).
- *Vegetation Cover*: The extent of vegetation cover at each site was estimated through simple visual inspection, providing an essentially subjective but nonetheless useful assessment. Vegetation coverage was deemed to involve only the extent of ground cover at or very near the soil surface. Tree canopies therefore do not accommodate the same extensive coverage, as would low-level scrub. The reason for this is twofold; firstly very few trees grow in the semi-arid conditions, and they usually have limited canopies with little leaf matter. Secondly, overland flow, a strong erosive force is constrained by vegetation cover at or very near the soil surface.

To provide a less subjective view of vegetation coverage a number of photographs were taken at each site, attempting to incorporate the entire slope. The photographs provided the base for simple classifications to be undertaken in order to accurately calculate the vegetation cover as a percentage. However, a number of practical issues arise from this method and have to be taken into consideration. The main concern is that such photographs will generally tend to lead to an overestimation in vegetation cover, as the image is oblique and will rarely be perpendicular to the slope surface, leading to the impression that more vegetation cover is present than may actually be the case (Figure 5.2).



Figure 5.2: The potential over-estimation of vegetation cover when using photography.

Using a simple supervised classification technique within ERDAS Imagine 8.6, a single photograph chosen for each site was classified using the maximum likelihood algorithm. The chosen image was deemed the optimum of all of those taken at a specific site (total of three), incorporating the entire slope and from the best possible angle (in relation to the point mentioned previously). The maximum likelihood classifier was chosen (it quantitatively evaluates both the variance and covariance of the category spectral response when classifying unknown pixels) (Lillesand and Kiefer, 2000). Furthermore, maximum likelihood classifiers have been used successfully in numerous remote sensing investigations where the assumption of normality is applicable to response patterns as is the case within this application. As in any supervised classification technique it is important to determine an appropriate training method so as to incorporate spectral response

patterns of each class to be classified, namely vegetation and non-vegetation. Through a small number of trial classifications it was determined that 20 different training points each consisting of five pixels for both vegetation and nonvegetation produced good results. The classifications were monitored using the histograms to monitor the normality of the spectral responses.

- *Geology (Lithology)*: With the aid of a number of different geology and lithology maps the geological unit on which each of the training sites was located was recorded. Previous field excursions allowed the visual identification of distinctive geologies, whilst those that were not so distinctive were classified using the maps.
- Sodicity: The field sodicity meter was used to determine the potential of soil sodicity and its dispersive properties. The method involved collecting a subsurface soil sample at each site at a depth of 20 centimetres. This was done largely to avoid sampling surface crusts that may vary significantly and thus should be sampled separately (Richards, 1954). Each sample was tested using the field sodicity meter using the methodology detailed in section 7.2.1, and also in the laboratory with standard procedures for the determination of dispersive properties so as to provide a basis for the comparison of the results.
- *Erosion*: The dependent variable erosion was determined at each site through the use of a simple soil erosion index. Prior field excursions were undertaken to develop the index, estimating the most extreme areas where various erosion processes operate and areas less affected by erosion. The index is a simple sliding scale ranging from 0, no appreciable erosion, to 8, severe subsurface gully erosion and can be seen in Figure 5.3 and Figure 5.4. The presence of surface wash was limited in the study area and was seen to be operating in combination with minor rill erosion when it did occur as runoff was concentrated into small channels, and

consequently they have been combined in Class 1. In addition to minor rilling two further classes, moderate and severe rill erosion were incorporated into the classification scheme. Rills were defined within the classification as being no deeper than around 30 centimetres. Erosion above 30 centimetres was defined as a gully. A number of definitions have been used in the literature to distinguish between rills and gullies (Øygarden, 2003). For example, Bocco (1991) stated that a depth of around 0.5 metres should differentiate between rills and gullies. Furthermore, the Food and Agriculture Organisation (FAO) (1965) defined gullies as channels whose depth and width do not allow normal tillage. However, the classification here used a 30-centimetre rule to distinguish between the two, and as long as it was implemented consistently it was felt that it would have little negative effect on the overall output.

This approach is subjective and is dependent on the surveyor's interpretation. In defence, it can be stated that:

- (i) The surveyor must be "trained" and have a degree of expertise in the field that is sufficient to distinguish between classes.
- (ii) The "system" or "scheme" used to classify the landscape is implemented consistently.

Through the field excursions, it was felt that slopes possessing gullies developed through the process of subsurface erosion were the most severely eroded slopes in the study area as opposed to those containing gullies developed through Hortonian processes. Moreover, gullies developed through piping were seen to range more in terms of the extent to which they operate and as such are contained within three classes as opposed to the two classes dedicated to Hortonian gully erosion.



Figure 5.3: The sliding scale classification for the dependent variable erosion.

While developing the classification scheme, it was important to keep in mind the issues associated with training sets used within supervised classifications. The number of different classes was the main concern as it could severely limit the classifiers performance. It was important to produce a classification that was representative of the processes identified in the field as well as the extent to which they were seen to be operating. Too many classes would make it very difficult for the classifier to perform well as the class separability would be reduced due to underlying causal relationships becoming blurred or even lost. Thus, the classification scheme produced is a balance between too many classes, (reducing classifier accuracy), and too few classes (resulting in the loss of diversity and the able to be representative).



Class 0: No Appreciable Erosion: No visible evidence of any significant soil erosion.





**Class 5: Moderate Gully Erosion:** Well developed gully(s) eroded through surface processes.



Class 1: Surface Wash/Minor Rill Erosion: Limited evidence of wash processes and/or minor rill erosion... Class 6: Minor Subsurface Gully Erosion: Small gully(s) developed through subsurface processes.



Class 2: Moderate Rill Erosion: Well developed rill networks potentially limiting farming activities.



Class 3: Severe Rill Erosion: Extensively developed rill networks with potential to develop into gully. erosion.





Class 7: Moderate Subsurface Gully Erosion: Well developed gully network eroded through subsurface processes.



Class 8: Severe Subsurface Gully Erosion: Very severe and well developed gully networks produced through subsurface processes.

Class 4: Minor Gully Erosion: Small gully(s) developed through surface processes.

Figure 5.4: Examples of the erosion classes attributed using the erosion classification scale.

Taking all of the above into consideration it was decided after the completion of the preliminary fieldwork to use a total of nine classes. The classes contained the entire range of processes that were seen to be operating in the region without over complicating the schedule.

The independent variables were inputted into the classifiers as continuous data, as opposed to classed or grouped data where appropriate. Geology for example is naturally a classed data set, as is the field sodicity meter. However, the variables extracted from the DEMs and those collected in the field (slope, aspect and vegetation) are inherently continuous and have not been adapted or grouped into classes. The independent variable aspect however, is often grouped and manipulated into predetermined classed intervals, avoiding the potential confusion of using the data in a continuous format where the two extremities, namely 1° and 360° are different but in reality very similar. Therefore, it might seem natural to group the data into simple classes, such as Northeast, Southeast, Southwest and Northwest. The reason for not taking this approach is twofold. Firstly, Brouwer (2002) suggested that continuous data should be used in preference of categorical data, as the latter will tend to cause a discontinuous relationship between independent variables and the dependent variable. Secondly, decision tree classifiers are naturally suited to either continuous or categorical inputs, and there is little supporting evidence to suggest that changing from one to another offers any potential benefit. Due to their non-linear capabilities, decision trees possess the ability to ask more than one question of an independent variable (see Chapter Four). Thus, even if a strong trend existed on slopes with aspects ranging from Northeast to Northwest (315° to 45°), it should be identified by asking two simple questions. Firstly, is the point of interest greater than

315°? If so, is it also less than 45°? ANNs also have this ability, however the method by which it undertakes the task is somewhat different, and in the interest of comparability and fairness, both classifiers used the same data sets.

#### **5.5.2 DEM Acquired Variables**

Two DEMs were developed for the study area produced by digitising contours from the Sorbas and Polopos 1:25 000 topographic maps, with 10 and 20 metre grid-cell resolutions. The topographic maps were scanned and geocorrected using ERDAS Imagine 8.6, and then exported to ArcMap where the contours were digitised onscreen using height as the attribute for each. This work was carried out by the Landscape and Ecology Research Group at the University of Hertfordshire as part of an ongoing investigation into geomorphological processes in and around the Sorbas Basin. At the time of study this was the only DEM of an appropriate resolution encompassing the entire study area.

The accuracy of a DEM and the variables derived from it is dependent upon the source of the elevation data, the methods by which the model is created, the structure of the elevation data, the vertical and horizontal resolution of the data and the topographic complexity of the landscape being modelled (Reuter *et al.*, 2006, Thompson *et al.*, 2001). Such issues are likely to heavily constrain DEM accuracy, and the DEMs developed for use within this investigation may suffer as a result. The topographic maps themselves are likely to contain some degree of error and this will be propagated further through the geocorrection process and the subsequent contour digitisation. Furthermore, the interpolation algorithm used will exert a strong

influence over the final DEM, and therefore the purpose of the DEM should be determined prior to algorithm selection.

The contour spacing on the topographic maps is 10 metres and thus constrains the grid-cell sizes that can realistically be created. Moreover, the scale at which the erosive processes operate must also be taken into consideration. Generally, the grid cell size for the DEM developed from the digitised contours with 10 metre spacing should be 20 metres (Livingstone pers. com.). However, it is important to recognise the fact that the DEM is the strongest link to the real field environment attributes, and therefore should be used to its maximum potential, and although this may incorporate errors, a DEM with a grid cell size of 10 metres was developed.



Figure 5.5: The DEM generated with grid cell size of 10 metres.

The DEMs were created using the TOPOGRID command in ARC/INFO. The TOPOGRID command is an interpolation method specifically designed for the creation of hydrologically correct DEMs (simulates natural movement of water over the surface). It is based on the ANUDEM programme developed by Hutchinson (1988; 1989) using a drainage enforcement algorithm and the DEMs generated can be

seen in Figures 5.5 and 5.6. The variables derived from the DEMs were acquired through the use of the GRID function in ARC/INFO; namely slope angle, aspect, plan and profile curvature, and flow length and flow accumulation (Figure 5.7).



Figure 5.6: The DEM generated with grid cell size of 20 metres.

Each grid was developed and the data extracted for each of the 520 training set points using Hawth's Analysis Tools (Beyer, 2003), designed to perform spatial analysis functions largely for ecological applications. Nonetheless, the tools are simply extensions to ArcMap and allow the extraction of data from known location points, a useful function currently unavailable in ARCGIS.

Resolution exerts a strong influence upon the quality of a DEM and the accuracy of the model. It is likely therefore that the two DEMs constructed with 10 and 20 metre grid cells will vary in terms of quality. As stated previously, the DEM is the strongest link to the real field environment when not incorporating the field collected data, and as such should be used to its maximum potential, and thus a finer resolution model was created. Previous studies have shown that as DEM resolution decreases, slope gradients decreased, particularly in areas of relatively steep slopes (Reuter *et al.*, 2006; Thompson *et al.*, 2001; Wolock and Price, 1994). Furthermore, Thieken *et al.* (1999) found that flow length and drainage density decreased as DEM resolution reduced. These findings are generally associated with the smoothing effect resulting from a reduction in DEM resolution, a consequence of the method by which slope angle and aspect are calculated. The slope angle and aspect of any given cell within a DEM is a function of itself and its eight neighbours. A detailed explanation of the calculation of slope angle is provided below and in Figure 5.7.



Figure 5.7: DEM calculation of slope angle (Adapted from Longley et al., 2001).

As stated previously, when using a grid-based DEM the common approach is to use a moving 3x3 window to derive slope angle (Zhou and Liu, 2004). Figure 5.7 details the steps required for calculating the slope angle of cell 5. Such an approach is highly dependent upon the resolution of the DEM, or the grid cell spacing. For example, to measure slope at a 10 metre resolution, a landscape profile of 30 metres is required. This issue is further exacerbated by the 20 metre DEM where slope angle is measured over a distance of 60 metres, resulting in further smoothing of the landscape compared to that of the 10 metre elevation model. Therefore, as suggested by Longley *et al.* (2001), slope is a function of resolution.

#### 5.6 DATA TRANSCRIPTION

In this instance data transcription is the process of transferring information that has been acquired (section 5.5) into the AI classifiers. As with any supervised classification technique, both training and testing examples are required to both calibrate and validate the classifier (Muchoney and Strahler, 2002). The data set developed was thus split into two subsets, namely training and testing, using a ratio of 3:1, creating a training set of 390 cases and a test set comprising of 130 cases. Figure 5.8 shows the range of different independent variables used within the classifications. The implementation of the data for each of the classification techniques is outlined in the following sub-sections.



Figure 5.8: The attributes determined from the DEMs (10 metre cell size).

#### 5.6.1 Data Transcription to Artificial Neural Networks

Neural Network models were constructed using Version 4.0 of the TRAJAN Neural Network Simulator software (Trajan, 1999). The software allows for a number of networks with various structures to be created and trained using an assortment of algorithms. The programme also allows the incorporation of both continuous and discrete variables within the training and test data sets.

As mentioned previously, the classifiers require both a training data set and a test data set. However, ANNs also require a verification data set. The verification set is used to track the training and ensure that the network is learning correctly and to avoid overlearning (see Chapter Four) (Sebastiá *et al.*, 2003). Furthermore, the verification set is used to determine the network error and thus the best network prior to any testing. It is recommended by the TRAJAN software that a ratio of 2:1 training to verification is used. Thus, the networks used within this research consisted of 260 training cases, 130 verification cases and a further 130 test cases.

Prior to the training of any neural network, a number of parameters and criteria have to be set. One of the most important considerations is the identification and implementation of an appropriate architecture as it governs the networks capability (Benediktsson *et al.*, 1990; Hepner *et al.*, 1990; Wang *et al.*, 1994; Fitzgerald and Bean, 2001). The type of network chosen for use here was a multilayer perceptron (MLP) consisting of an input layer, a hidden layer and an output layer. The number of hidden layers has been kept constant at one within all of the networks used within this research as many environmental applications have used only one layer (Ellis, 1997; Dedecker *et al.*, 2004; Tiira, 1999) and should allow enough flexibility to determine underlying non-linear patterns.

The back-propagation training algorithm was chosen largely due to their extensive use in a range of classification problems. The training algorithm, devised independently by Robbins and Monro (1951), Rumelhart *et al.* (1986), Werbos (1994) and Parker (1985) has been used successfully with many network types and in particular multilayer perceptrons (MLPs) and has become the single most useful neural networking algorithm (Tveter, 1998). However, as outlined in Chapter Three, the back-propagation algorithm requires training parameters to be selected. The learning rate for the training of the networks was set at 0.6 and the momentum term at 0.8 based upon prior knowledge identified in the literature. Dai and MacBeth (1997) identified that learning rate should be between 0.6 and 0.7 and the momentum should be between 0.8 and 0.9. In support of this McClelland and Rumelhart (1988) found a learning rate of 0.7 and momentum of 0.9 produced the optimum results within their research, however Tveter (1998) noted that it is likely to be problem specific.

To ensure that the optimum network topology is achieved using only a single hidden layer, each classification was run using only one node up to a maximum of 25 nodes, a technique used by Leane *et al.* (2003), Malhotra and Malhotra (2003) and Moatar *et al.* (1999), to identify the ideal architecture for a given problem. A maximum of 25 nodes in the hidden layer was chosen as Blum (1992) suggests the number should fall between the number of input and outputs, and Berry and Linoff (1997) proposed that there should be no more than twice that of the input layer. Furthermore, each classification was replicated 20 times as the outcome varies each time a network is trained. In an investigation into the classification performance of ANNs, Sánchez *et al.* (1996) trained networks with varying topologies ten times with the same data, as did Gupta and Sexton (1999) when undertaking a comparison between neural network training algorithms. Thus, for any classification a total of 500 neural networks were built and trained. A single optimum network was chosen for each of the different 25 network topologies chosen based upon the verification error.

#### 5.6.2 Data Transcription to Decision Tree Classifiers

The decision tree growing software used was CART 5.0 (Brieman *et al.*, 1984; Salford Systems, 2004). The data mining software allows the incorporation of both continuous and categorical data for both the independent and dependent variables and consists of a range of splitting methods. CART 5.0 uses binary splits that divide each parent node into two child nodes by posing simple yes or no questions to the data. Other DTCs such as the CHAID programme allows multiple splits at each node, however this can lead to less accurate splits.

The CART 5.0 programme requires the training set file and a test set file, incorporating 390 and 130 examples respectively. The programme works by growing the largest possible tree consisting of a number of terminal nodes. The tree is then pruned removing small sections that have little or no influence on the classification accuracy creating a number of smaller trees. The best tree is selected by testing for error rates of costs, achieved by using the test set to calculate the rate at which cases are misclassified.

#### 5.6.3 Implementation of Independent and Dependent Variables

To answer the questions posed in section 5.2 regarding the ability of ANNs and DTCs to classify soil erosion and the influence that different independent variables have upon levels of classification accuracy, a range of different classifications incorporating different combinations of dependent and independent variables are undertaken. Three broad classifications were run; a two class classification, three class classification and a further classification including all (nine) of the classes. For each of these, eight different sets of independent variables were used to train the classifiers, detailed in Table 5.4. The amalgamation of different classes for the two, three and nine class classifications can be seen in Figure 5.9, and a graphical representation of the different layers incorporated within a classification (independent and dependent variables) can be seen in Figure 5.10.

Data Set					Inde	ependen	t Varia	ables				
	Field Slope Angle	Field Slope Aspect	Estimated Vegetation	Field Sodicity Meter Value	Geology	Classified Vegetation	DEM Slope Angle	DEM Slope Aspect	Flow Length	Flow Accumulation	Plan Curvature	Profile Curvature
10 Metre DEM										alu:		
20 Metre DEM												
Field Variables												
Field Variables and Classified Vegetation												
Field Variables and 10 metre DEM												
Field Variables and 20 metre DEM												R. Mark
Field Variables, 10 metre DEM and classified Vegetation												
Field Variables, 20 metre DEM and classified Vegetation												

Table 5.4: The eight different data sets and the independent variables that they contain.



Figure 5.9: The amalgamation of different erosion classes for use within different classifications.



Figure 5.10: A simple diagrammatic representation of the independent and dependent variables.

#### 5.7 ACCURACY ASSESSMENT

The calculation of accuracy is vitally important when any type of classification has been undertaken. It is important to remember that errors are likely to be present in classifications and the maps they produce as they are simply generalisations of reality (Brown *et al.*, 1999; Dicks and Lo, 1990). Thus, before their value can be determined it is important to question their ability to meet the needs and requirements of the intended application.

The level of accuracy that is acceptable is highly subjective and dependent on the questions being asked. Consider two scenarios, firstly that of an individual landowner and secondly that of a regional or provincial government. It is unlikely that the level of accuracy acceptable at the regional scale (for the government) is sufficiently detailed for the individual landowner. Furthermore, the resolution of the product developed will dictate its wider applicability and usefulness.

Numerous methods of accuracy assessment for classification procedures exist and have been discussed widely in the literature. However, the confusion matrix, otherwise known as the error matrix or the correlation matrix, is the standard form for representing site-specific error (Campbell, 2002) and thus currently at the core of accuracy assessment literature (Foody, 2002). An error matrix simply compares information from reference sites to information on predicted sites for a number of sampled areas (Congalton and Green, 1999). The overall accuracy of a classification is simply calculated through the sum of the major diagonal of the matrix divided by the total number of samples classified. An example of this can be seen in Table 5.5. An error matrix is an effective way to represent accuracy in that the accuracies of each

#### Research Framework and Data Methods

category are plainly described (Congalton, 1991), and thus provides an obvious foundation for accuracy assessment (Campbell, 2002; Canters, 1997).

	Ground Truth Data								
		1	2	3	4	5	6	Total	
	1	34	0	1	0	0	0	35	
	2	2	35	1	2	2	1	43	
	3	1	0	44	1	0	1	47	
Classified Image	4	5	1	4	40	2	0	52	
	5	8	11	0	7	45	1	72	
	6	0	3	0	0	1	47	51	
	Total	50	50	50	50	50	50	245	
	Accuracy	0.8167							

Table 5.5: An example of a correlation (error) matrix.

			Grou	nd Truth	Data		
		1	2	3	4	5	6
	1	0.68	0	0.02	0	0	0
Classified	2	0.04	0.7	0.02	0.04	0.04	0.02
	3	0.02	0	0.88	0.02	0	0.02
Image	4	0.1	0.02	0.08	0.8	0.04	0
, mugo	5	0.16	0.22	0	0.14	0.9	0.02
	6	0	0.06	0	0	0.02	0.94
	Total	1	1	1	1	1	1

Table 5.6: The producer's accuracy.

	Ground Truth Data									
		1	2	3	4	5	6	Total		
Classified Image	1	0.9714	0	0.0286	0	0	0	1		
	2	0.0465	0.814	0.0233	0.0465	0.0465	0.0232	1		
	3	0.0213	0	0.9362	0.0212	0	0.0212	1		
	4	0.0962	0.0192	0.0769	0.7692	0.0385	0	1		
	5	0.1111	0.1528	0	0.0972	0.625	0.0138	1		
	6	0	0.0588	0	0	0.0196	0.9215	1		

Table 5.7: The user's accuracy.

The correlation matrix provides further valuable insights into classification performance in addition to overall accuracy. All of the non-diagonal elements within a matrix provide errors of commission and omission (Lillesand and Kiefer, 2000), otherwise known as the producers and users accuracy. The producer's accuracy is a simple measure of omissions, indicating areas of a specific coverage that have been omitted from the classification. An example of the producer's accuracy can be seen in Table 5.6 derived from the correlation matrix in Table 5.5. In contrast, the user's accuracy determines the extent of commission errors, highlighting the probability of examples classified within a specific class actually belonging to that class in reality. Table 5.7 shows the user's accuracy derived from the error matrix in Table 5.5.

In an attempt to further determine the usefulness of the AI approach to the classification of soil erosion it is compared to the more traditional statistical approach of Discriminant Analysis (DA). Discriminant analysis is a classification technique used to classify unknown cases from a training set rather like the ANNs and DTCs previously discussed. It makes the assumption that the independent variables are normally distributed, encouraging the use of continuous variables rather than discrete variables (Blackard and Dean, 1999). Nonetheless, McLachlan (1992) supports the use of linear discriminant analysis in cases where the assumption may be violated and Blackard and Dean (1999) suggest that this is relatively common and has a limited affect on results.

#### **5.8 DEVELOPMENT OF AN EROSION RISK SCHEDULE**

An erosion risk schedule comprises a set of relatively simple rules, applies them to an erosion map and produces an erosion risk map. It would have been possible to produce an erosion risk map rather than one of actual soil erosion by simply replacing the erosion scale (Figure 5.3) with a risk scale. However, the approach taken here is

such that an erosion risk output is derived through the incorporation of the developed erosion map in combination with simple topographic variables.

The risk schedules developed here are based upon the maps created for the two class and the nine class classifications and simply a logical progression from them as seen in Figure 5.11.

The erosion risk schedule is based on a suite of logical geomorphological principles using current erosion as the basis for the determination of risk along with some basic topographical attributes and only concerns risk by association. As outlined in Chapter Three, risk can relate to *actual* or *potential* depending upon whether land use is taken into consideration or other potentially changeable variables such as vegetation cover. However, the map created here simply considers the current conditions and therefore is a classification of *actual* risk.



Figure 5.11: The progressive development of the erosion risk schedule.

#### 5.8.1 Methodology behind the Risk by Association Schedule

The development of the risk schedule can be seen diagrammatically in Figure 5.11. The 10 metre DEM is used with the slope angle grid and an erosion map derived through the AI techniques. Through the use of the two topographical parameters it is possible to determine the concavity or convexity of a cell in relation to one of its eight neighbours. The DEM determines whether a cell is above or below an adjacent cell based on height, and the slope angle will allow the calculation of profile shape to be made. For example, if A is above an adjacent cell B and the slope angle is greater for cell A, then a concavity exists between the two cells A and B. Using concavities and convexities in relation to surrounding cells that may be eroding, some simple rules can be constructed producing a model for soil erosion risk by association.

The risk for each cell is a function of its eight neighbouring cells. The basic rules vary slightly for the two class schedule and the nine class schedule, however, the general principles are the same for both. The models were developed using a combination of programmes including ArcMap, ArcInfo, and Idrisi. The models are simply developed via a series of simple logical steps gradually building into a complex grid.

Risk by association only takes into consideration the perceived risk to any given slope posed by neighbouring slopes that under current conditions are eroding (according to the classification). It does not therefore determine the erosion risk for slopes as single entities as it is determining *actual* risk and the erosion map used is thus definitive. It is inherently difficult to quantitatively determine levels of risk, and the approach outlined here can be said to be semi-quantitative. The rules developed are simply outlined in Figure 5.12 whereby a single value is attributed to a given cell for each of its eight neighbouring cells if they are currently eroding. If they are not eroding then the risk through association is zero and no value is assigned. Four possible scenarios exist for the two class risk schedule, labelled A, B, C and D, and eight scenarios for the nine class schedule A, B, C, D for both no erosion (0) and surface erosion (1) seen in Figure 5.11. A simple risk through association map is thus created by applying the relevant risk factor for each morphology and calculating the total for each individual cell, such that if a cell were adjacent to 3 eroding cells with risk factors of 2, 1 and 3 a value of 6 would be assigned to that cell. The following sub-section details the justification for the risk values chosen.

#### 5.8.2 Justification for Risk Values

In order to determine the level of erosion risk using a simple set of rules, a risk factor has to be attributed to each topographical morphology. Young and Mutchler (1969) found that erosion was less than half on concave slopes than on convex slopes under otherwise similar conditions, a finding supported further by Hadley and Toy (1977). This is largely a consequence of the fact that the steepest gradient is at the top of the profile, and reducing gradient as both the slope length and runoff increase. Moreover, Faulkner *et al.* (2003b) highlighted the potential of coupling on convex sites further increasing the potential risk through association.

Scenario A in the two class schedule, seen in Figure 5.12 shows a convex morphology with erosion occurring on the flatter top section of the profile only. It could therefore be proposed that the eroding cell could potentially capture the cell below as it possesses a steeper gradient and would have a larger contributing area. Scenario C is in fact the reversal of A, and due to the potential risk posed by headward erosion or undercutting at the slope base a risk factor of 3 is applied, as opposed to 2 for scenario A. The concave morphologies, scenarios B and D, are assigned values of 1 and 2 respectively. As outlined previously such morphologies tend to reduce surface erosional processes, however risk through association is such that both scenarios could potential be at risk through runoff progression in B and headward erosion in D.

The risk schedule, derived from the nine class erosion map, includes the same rules to account for the risk by association of surface erosion but also has some additional rules regarding the subsurface erosion process. Furthermore, the same risk values are attributed to slopes at risk from subsurface erosion regardless of whether they are currently not eroding or eroding under surface processes. The main issue to be confronted are discussed in Chapter Three and by Faulkner *et al.* (2003b) concerning subsurface erosion on convex morphologies. Pipes preferentially develop in convex slopes with an infiltrating surface and a potential outlet. Consequently the convex morphologies A and C have been allocated values of 5 and 4 respectively. Scenario A is at high risk from piping as the slope above it is eroding and may provide potential outlets. Scenario C may be at risk from headward erosion, undercutting and indeed piping if an adjacent slope is eroding under such a process. The risk associated with concave morphologies however is not as great, and as such have been attributed values of 2 and 3 respectively.



**Figure 5.12:** The logical progression of the risk by association schedule and the assigned risk values based on the scenario.

#### 5.8.3 Summary

Risk through association is a very simple concept. The risk values assigned are inherently subjective and based on prior knowledge of the erosion processes operating. Nonetheless, relatively complex outputs can be achieved through the implementation of some simple rules in combination with only a limited number of input variables. Thus, in order to improve the risk by association map further variables could be incorporated such as vegetation cover and soil sodicity.

#### **5.9 CONCLUSIONS**

This chapter has outlined the research framework and data methods upon which the thesis is based, using a series of questions in an attempt to provide a focus and structure. A review of the training data sets has been undertaken, including the methods through which they were derived and the justification and validation of the different variables incorporated within them. The methods used to train and develop the different classifiers have been comprehensively outlined, as has the development of the risk by association schedule. In order to answer the question posed, regarding the physico-chemical relationships and the applicability of the field sodicity meter, a range of soil samples were collected and a number of laboratory analysis methods were undertaken. The methods, results and a detailed discussion relating to this can be seen in Chapter Seven, Field Investigations.

The following chapter reviews and details the results obtained through the research framework and data methods outlined in this chapter, reporting the main findings and relating them to the research questions aforementioned. Following the documentation of the results in Chapter Six and the Field Investigations in Chapter Seven, the allround implications are discussed in Chapter Eight and related to relevant literature where appropriate.

### 6

# Soil Erosion Classifications, Risk Schedules and Rule Extraction

#### **6.1 INTRODUCTION**

This chapter provides details of the results obtained from the implementation of the research framework outlined in Chapter Five. The results are documented with particular reference to the research questions stated in section 5.2 (fulfilling the stated aims and objectives detailed in Chapter One). The results of the soil erosion classifications are presented for both the artificial neural networks and the decision tree classifiers, followed by the classification results using discriminant analysis (DA). The classifications are compared and contrasted, with particular emphasis regarding the influences and effects that using different independent and dependent variables has upon classifier performance. The ability to extract rules or knowledge regarding the various processes from the ANNs and DTCs is reviewed, along with a discussion regarding neural network and decision tree topology. Finally, potential risk maps and erosion risk schedule maps for the study area are presented and discussed.

## 6.2 SOIL EROSION CLASSIFICATIONS USING ARTIFICIAL NEURAL NETWORKS

A total of 500 artificial neural networks were trained using the eight different independent variable data sets outlined in Table 5.4. Classifications were undertaken for a simple binary problem (two class), an intermediate three class problem, and a detailed (interval) nine class problem incorporating all of the erosion classes outlined in Figures 5.3 and 5.4. For each of the different classifications undertaken, a single
optimum or 'best' network is selected, based on their verification error (the root mean square (RMS) of the errors on each individual case).

### 6.2.1 Two Class Classifications Using Artificial Neural Networks

The two class classifications involved the amalgamation of all of the erosion classes into a single class (1), and the no appreciable erosion class (0), creating a basic binary output. Such a classification highlights the ability of the classifier to delineate between eroding cells, and those portraying no appreciable signs of erosion.

The results obtained are summarised in Table 6.1, highlighting the overall accuracy, the verification error and the network topology for each of the eight different classifications trained using different data sets. The level of accuracy (see Chapter Five, section 5.7) ranged from 66.2 (3.s.f.) percent using the field collected variables to 75.4 percent for both the field collected data used in conjunction with the 10 metre DEM variables, and the data set using the same variables and in addition the classified vegetation data. Accuracy's of 77.7 percent were attained for the classifications trained with the field attributes, classified vegetation and the 10 metre DEM, and the field attributes and the 20 metre DEM, but these had slightly higher error rates than their counterparts documented in Table 6.1.

Verification errors (RMS) ranged from 0.362 to 0.459, highlighting a general trend whereby classifications undertaken using primarily DEM independent variables suffered from an increased error rate compared to those that incorporated field collected variables.

To have under shad dhay hay part	Overall Accuracy	Verification Error	Network Topology
DEM 10 metre	0.692308	0.4526253	6-21-1
DEM 20 metre	0.684615	0.459258	6-13-1
Field collected	0.661538	0.3618585	5-22-1
Field collected and DEM 10 metre	0.753846	0.3692149	9-18-1
Field collected and DEM 20 metre	0.707692	0.3701264	9-5-1
Field collected and classified vegetation	0.669231	0.3938281	5-6-1
Field collected, classified vegetation and DEM 10 metre	0.753846	0.3930671	9-12-1
Field collected, classified vegetation and DEM 20 metre	0.723077	0.3932958	9-3-1

Table 6.1: Summary table of 'best' networks for two class classifications using ANNs.

The correlation matrix for the optimum network identified for the two class classification can be seen in Table 6.2, trained using the field collected variables and the 10 metre DEM attributes. The users and producers accuracy in this classification can also be seen in the table, revealing valuable insights into the classification performance. It is evident from the results that only 60 percent of the areas which were not eroding have been successfully classified, in contrast to 83.5 percent of those that were eroding. This indicates a good classification for the latter. Moreover, 79.8 percent of cells classified as eroding were actually seen to be doing so, compared with 65.9 percent of the cells that were not eroding.

	And Statistics	Actu	al			
		0	1	Total	Overall	
Predicted	0	27	14	41	Accuracy	
	1	18	71	89		
	Total	45	85	130		
	Producers Accuracy		Users A	Accuracy	0.753846	
	0	60%	0	65.9%		
	1	83.5%	1	79.8%		

**Table 6.2:** Correlation matrix for the two class ANN trained using the field acquired attributes and the 10 metre DEM independent variables.

To better understand classifier performance when undertaking binary classifications, receiver operating characteristic (ROC) curves can be determined. Zweig and Campbell (1993) highlighted the ability and usefulness of ROC curves. These summarise the performance of a two class classifier across the range of thresholds by plotting the sensitivity (class two true positives) against one minus the specificity (class one false negatives). The true positive rate is defined as the number of positives correctly classified divided by the total number of positives. The false positive rate is the number of negatives incorrectly classified, and divided by the total number of negatives. The perfect classifier would therefore produce an ROC curve that hugs the left and top sides of the graph, and thus the total area under the graph would be one (perfect). The ROC curve is a useful tool for comparing different classifications as it takes into account the performance of the classifier across all possible thresholds.

Figure 6.1 shows the ROC curve for the classification. The Area Under the Curve (AUC) is 0.88 indicating a relatively strong classifier, as a value of 1.0 would indicate a perfect classifier and a value of 0.5 a purely random classifier (Tseng-Chung and Li-Chiu, 2005). Pearce and Ferrier (2000) suggested that classifiers with AUC between 0.5 and 0.7 are poor, between 0.7 and 0.9 are good, and above 0.9 are excellent.



Figure 6.1: ROC curve for the two class ANN classification trained using the field acquired attributes and the 10 metre DEM independent variables.

Due to data constraints, erosion maps could not be created using the independent variables collected in the field as these were not known at every location within the study area. Therefore, erosion maps have been created using classifiers trained using the DEM data only (as this is known at every cell). Figure 6.2 shows the binary soil erosion map produced from the 10 metre DEM data draped on the topographical map in order to provide some spatial reference.

Erosion maps can be seen for the ANNs trained using the 10 and 20 metre DEM data in Appendix Six (Figures A6.1 and A6.3), and draped over the topographical maps (Figures 6.2 and A6.4). Correlation matrices for each of the two class ANN classifications undertaken can be seen in Appendix Two, Tables A2.1 to A2.8 inclusive.



**Figure 6.2:** Classified erosion map drape derived from the ANN trained using 10 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).

### 6.2.2 Three Class Classifications Using Artificial Neural Networks

In an attempt to determine the ability of neural networks to differentiate between areas of no appreciable erosion, rill erosion and gully erosion, classifications incorporating the three classes were undertaken (see Figure 5.9). The three rill erosion classes were amalgamated, as were the five gully classes, and incorporated into a classification with the 'no appreciable erosion' class. Once again a range of different networks were trained using the different data sets, the results of which can be seen in Table 6.3.

As was the case with the two class ANN classifications, the levels of overall accuracy appear to be highly influenced by the incorporation of the field collected independent variables. The networks trained using independent variables derived from the elevation models not only have the lowest overall accuracy but do so in accordance with the highest verification error of all eight classifications.

alenar jurit	Overall Accuracy	Verification Error	Network Topology
DEM 10 metre	0.469231	0.4432813	6-13-3
DEM 20 metre	0.546154	0.4377942	6-16-3
Field collected	0.569231	0.334454	5-13-3
Field collected and DEM 10 metre	0.584615	0.3856856	9-21-3
Field collected and DEM 20 metre	0.523077	0.3803771	9-23-3
Field collected and classified vegetation	0.6	0.3502805	5-10-3
Field collected, classified vegetation and DEM 10 metre	0.584615	0.4196299	9-10-3
Field collected, classified vegetation and DEM 20 metre	0.576923	0.3868505	9-20-3

Table 6.3: Summary table of 'best' networks for three class classifications using ANNs.

The matrices for both the classifications undertaken using the field acquired attributes, and the field acquired attributes used in combination with the classified vegetation can be seen in Tables 6.4 and 6.5 respectively. It becomes apparent that the overall accuracy in both classifications suffers as a result of the particularly poor determination of rill erosion (class 1), when compared to the no erosion and gully erosion classes. The classification of rill erosion is poorer in the latter correlation matrix (Table 6.5), but has a better overall accuracy than that achieved using only the field acquired variables due to the improved classification of both other classes.

Statistics		i in klasti	Actual			1201 127.0
		0	1	2	Total	Overall
Predicted	0	32	13	19	64	Accuracy
1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1	1	4	4	1	9	
the lines of	2	9	10	38	57	
	Total	45	29	56	130	0.56923
	Producers Accuracy			Users Accuracy		
	0	71.1%	0		50%	
	1	13.8%	1		44.4%	
	2	67.9%	2	2	67%	

**Table 6.4:** Correlation matrix for the three class ANN trained using the field acquired attributes.

			Actual		Sec. 10.95	
		0	1	2	Total	Overall
Predicted	0	35	13	17	65	Accuracy
alente ferret	1	1	2	0	3	
	2	9	12	41	62	
	Total	45	29	56	130	0.6
	Producers Accuracy		Users Accuracy			
problem to	0	78%	0		53.8%	
	1	6.9%	1		66.7%	
23.25	2	73.2%	2		66.1%	

**Table 6.5:** Correlation matrix for the three class ANN trained using the field acquired attributes and the classified vegetation independent variable.

The point is emphasised further through the determination and analysis of both the users and producers accuracy. The producers accuracy for rill erosion (class 1) within both of the aforementioned classifications highlights the large associated omission errors. Of the 29 test samples, only four have been correctly classified for the field acquired attribute classification and only two for the field and classified vegetation classification. The misclassifications appear to have been shared by the other two classes relatively evenly. This is indicated by the similar users accuracies for each in both classifications.

An interesting point to note that is evident in both classifications concerns the misclassification of the 'no appreciable erosion' and 'gully erosion' classes; both networks appear to have difficulty in differentiating between the two in certain instances. For example, in the correlation matrix seen in Table 6.4, there are 45 cases of class 0 (no appreciable erosion), of which 32 have been correctly classified, yet of the remaining ten misclassified cases, nine have been attributed to class 3 (gully erosion) and only four to class 1 (rill erosion).

Figure 6.3 emphasises the problem further, showing the three class erosion map produced from a neural network trained using the 10 metre DEM data. The classification highlights the ANNs inability to determine cases of rill erosion, as only a very limited number of cells in the study area have been attributed to the class. This is an issue more prominent in those networks trained only using the DEM data. The problem is evident to a lesser or greater degree throughout all of the three class classifications, the matrices of which can be seen in Appendix Two.



**Figure 6.3:** Classified erosion map drape derived from the ANN trained using 10 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).

## 6.2.3 Nine Class Classifications Using Artificial Neural Networks

The final classification to be undertaken using neural networks involved all nine classes detailed in the erosion classification scheme (Figure 5.3). Using the same procedures and techniques as those used in the two and three class classifications, networks were once again trained using the eight different data sets, the results of which are detailed in Table 6.6.

The overall performance of networks for the classification of all nine classes is much lower than those achieved in the two and three class classifications, with overall accuracies ranging from 30 to 39.2 percent. As with the classifications detailed in sections 6.2.1 and 6.2.2, the highest accuracies are produced in accordance with field acquired independent variables. Table 6.7, the correlation matrix for the network trained using the field collected variables and the classified vegetation, highlights the apparent inability in this instance of the neural network to differentiate between the different classes of erosion.

	Overall Accuracy	Verification Error	Network Topology
DEM 10 metre	0.346154	0.2892779	6-5-9
DEM 20 metre	0.3	0.2840651	6-14-9
Field collected	0.376923	0.2710848	5-8-9
Field collected and DEM 10 metre	0.392308	0.2727597	9-3-9
Field collected and DEM 20 metre	0.361538	0.2751743	9-3-9
Field collected and classified vegetation	0.392308	0.2708853	5-3-9
Field collected, classified vegetation and DEM 10 metre	0.392308	0.2803297	9-2-9
Field collected, classified vegetation and DEM 20 metre	0.361538	0.2717551	9-14-9

Table 6.6: Summary table of 'best' networks for nine class classifications using ANNs.

Of the nine classes incorporated in the classification, the network was only able to attribute unseen examples to three of the classes, namely 'no appreciable erosion', 'moderate gully erosion' or 'minor subsurface gully erosion'. As a consequence, the overall accuracy is particularly poor, as are both the producers and users accuracies. This problem appears to have been replicated in some form throughout all of the nine class classifications undertaken using the ANNs, whereby all of the unknown test cases have been incorrectly classified into two or three classes. This problem may be attributed to a swamping effect in the training data whereby the number of training cases of class 0 (no appreciable erosion), exceeded those in the other eight classes subsequently causing the network to learn such examples well whilst neglecting less affluent classes.

						Actua	1					
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	38	10	3	4	3	9	3	3	0	73	Accuracy
	1	0	0	0	0	0	0	0	0	0	0	
	2	0	0	0	0	0	0	0	0	0	0	
	3	0	0	0	0	0	0	0	0	0	0	
Predicted	4	0	0	0	0	0	0	0	0	0	0	0.392308
	5	3	1	0	0	0	0	2	0	0	6	
	6	4	6	2	3	4	9	13	6	4	51	
	7	0	0	0	0	0	0	0	0	0	0	
	8	0	0	0	0	0	0	0	0	0	0	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	roduc	ers Ac	curac	y		Users Accuracy					
	0			84.4	%		0 52.1%					
	1			-			1			-		
	2			-			2		-			
	3			-			3					
	4			-			4		14			
	5			-			5					
	6			72.29	%		6		25.5%			
	7			-			7		-			
	8			-			8			-		

**Table 6.7:** Correlation matrix for the nine class ANN trained using the field acquired attributes and the classified vegetation independent variable.

The potential swamping effect can be seen in Figure 6.4, the classified erosion map derived using the 20 metre data set. The vast majority of the study area has been classified as not eroding to any appreciable extent, however, this classification has incorporated five other classes, although they are relatively infrequent except for the 'moderate subsurface gully erosion' class. Moreover, the classified map produced for the 10 metre classification has not been shown as the network classified every cell in the image as being in the same class, 'no appreciable erosion'.



**Figure 6.4:** Classified erosion map drape derived from the ANN trained using 20 metre DEM variables for a nine class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).

# 6.3 SOIL EROSION CLASSIFICATIONS USING DECISION TREE CLASSIFIERS

As with the neural networks, decision trees were grown for each of the three different classifications using the eight training sets. Unlike the ANNs however, only a single tree is grown for each classification problem as the optimum solution for any given problem is found every time and thus there is no need for replications to be made. The CART 5.0 decision tree software grows the largest possible decision tree and recursively prunes backwards developing a set of trees ranging in size. A suitable tree can be selected from the set based either upon the number of terminal nodes or the cross-validated relative cost (error), and in order to provide a consistent methodology the latter option was preferred. Figure 6.5 shows an error curve used to determine the optimum tree based on the relative cost (misclassifications), and in this example it is the tree containing 16 terminal nodes.





### 6.3.1 Two Class Classifications Using Decision Tree Classifiers

Table 6.8 provides a summary of the results obtained for the binary classification incorporating 'erosion' and 'no appreciable erosion'. It is evident that accuracies are generally lower for the decision trees trained using only variables acquired from elevation models. Moreover, the relative cost for these classifications is also considerably higher than for the other classifications, with the 20 metre DEM

classification performing particularly poorly. However, accuracies exceeded 77 percent in two of the eight classifications, but the trees with the lowest errors were those grown only from the field acquired variables and the same data used in combination with the 10 metre DEM variables.

	Overall Accuracy	Relative Cost	Terminal Nodes
DEM 10 metre	0.676923	0.682 +/- 0.088	6
DEM 20 metre	0.569231	0.899 +/- 0.091	11
Field collected	0.761538	0.480 +/- 0.079	12
Field collected and DEM 10 metre	0.761538	0.480 +/- 0.079	16
Field collected and DEM 20 metre	0.738462	0.505 +/- 0.079	10
Field collected and classified vegetation	0.776923	0.550 +/- 0.081	2
Field collected, classified vegetation and DEM 10 metre	0.746154	0.535 +/- 0.082	15
Field collected, classified vegetation and DEM 20 metre	0.776923	0.550 +/- 0.081	2

**Table 6.8:** Summary table of decision trees grown for two class classifications using DTCs.

S. College		Actu	al		
		0	1	Total	Overall
Predicted	0	34	20	54	Accuracy
	1	11	65	76	
	Total	45	85	130	0.761538
and the second	Producers Accuracy		Users A	Accuracy	
	0	75.6%	0	63%	
Elline r V. S.	1	76.5%	1	85.5%	

**Table 6.9:** Correlation matrix for DTCs trained using the field acquired independent variables.

Table 6.9 details the correlation matrix produced for the tree grown using only the field acquired variables, as this is the optimum tree grown for the problem. The producers accuracy for both classes exceed 75 percent with the users accuracy exceeding 85 percent for the 'erosion' class. However, a high omission rate regarding

the 'no appreciable erosion' class culminates in a relatively low users accuracy of 63 percent. All of the correlation matrices for the DTCs can be seen in Appendix Two.

As with the neural networks, ROC curves can be determined for the decision trees trained for a binary classification. Figure 6.6 shows the rate of true positives plotted against the rate of false positives. The AUC is 0.86 indicating a good classification performance according to Pearce and Ferrier (2000), with good separation of the frequency curves.



Figure 6.6: ROC curve for the two class DTC classification trained using the field acquired independent variables.

The decision tree developed for this classification problem can be seen in Figure 6.7 detailing the various rules developed through the training stage. The implications associated with the grown tree will be discussed in detail later, however, the classified map derived from the DTC grown using the 10 metre DEM data can be seen in Figures 6.8.



Figure 6.7: The decision tree grown for the two class classification using the field acquired independent variables.



**Figure 6.8:** Classified erosion map drape derived from the DTC trained using 10 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).

## 6.3.2 Three Class Classifications Using Decision Tree Classifiers

The details of the decision tree classifiers grown for the three class classification problem are outlined in Table 6.10. The 'optimum' or 'best' tree was that trained using the field collected variables used in association with the classified vegetation independent variable. The overall accuracy produced from this classification was in excess of 63 percent, and also had the lowest relative cost. The tree possessed 14 terminal nodes, and the architecture can be seen in detail in Figure 6.9 detailing the rules governing the splits.

As with the previous classifications, the weakest classifications appear to have been derived using the training data sets provided by the digital elevation models, culminating in both the lowest overall accuracies and the highest errors.

And the main sector of the sec	Overall Accuracy	Relative Cost	Terminal Nodes
DEM 10 metre	0.5	0.826 +/- 0.064	5
DEM 20 metre	0.515385	0.798 +/- 0.064	19
Field collected	0.569231	0.691 +/- 0.068	4
Field collected and DEM 10 metre	0.607692	0.653 +/- 0.066	16
Field collected and DEM 20 metre	0.607692	0.660 +/- 0.065	17
Field collected and classified vegetation	0.630769	0.597 +/- 0.065	14
Field collected, classified vegetation and DEM 10 metre	0.623077	0.608 +/- 0.065	12
Field collected, classified vegetation and DEM 20 metre	0.623077	0.608 +/- 0.065	12

Table 6.10: Summary table of 'best' tree for three class classifications using DTCs.

The correlation matrix for the DTC trained using the field acquired attributes and classified vegetation variable can be seen in Table 6.11. It is evident that the decision tree performs reasonably well at classifying two of the three classes, namely 'no

# Soil Erosion Classifications, Risk Schedules and Rule Extraction

appreciable erosion' and 'gully erosion', producing relatively high producers and users accuracies. However, the decision tree appears to perform poorly when classifying class 1, 'rill erosion'. Only 11 of the actual 29 cases have been correctly attributed, culminating in a producers accuracy of 37.9 percent. Moreover, a high error of commission due to misclassifications has also led to a low users accuracy.

			Actual			
		0	1	2	Total	Overall
Predicted	0	35	9	13	57	Accuracy
1.15-161	1	7	11	9	27	
an and the	2	3	7	36	46	0.630769
	Total	45	29	56	130	
1.	Producers Accuracy		Users Accuracy		uracy	
-	0	77.8%	0		61.4%	
	1	37.9%	1		40.7%	
	2	64.3%	2	1	78.3%	

**Table 6.11:** Correlation matrix for the three class DTC trained using the field acquired attributes and the classified vegetation independent variable.

The classifications derived from the decision trees trained with the 10 metre DEM variables for the three class classification can be seen in Figure 6.10.



Figure 6.9: The decision tree grown for the three class classification using the field acquired attributes and the classified vegetation independent variables.



**Figure 6.10:** Classified erosion map derived from a DTC trained using 10 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).

### 6.3.3 Nine Class Classifications Using Decision Tree Classifiers

The overall accuracies achieved for the nine class classifications using decision trees ranged from 13.1 to 29.2 percent, and occurred in combination with high error rates. Of the eight different training sets used, only four produced overall accuracies that exceeded 25 percent. A further four of the decision trees grown possessed fewer terminal nodes than the eight dependent variable erosion classes indicating the apparent inability of the trees to distinguish between and separate different levels of erosion.

	Overall Accuracy	Relative Cost	Terminal Nodes
DEM 10 metre	0.223077	0.889 +/-0.020	10
DEM 20 metre	0.130769	0.867 +/- 0.047	20
Field collected	0.292308	0.801 +/- 0.051	31
Field collected and DEM 10 metre	0.238462	0.826 +/- 0.039	6
Field collected and DEM 20 metre	0.269231	0.858 +/- 0.043	51
Field collected and classified vegetation	0.253846	0.844 +/- 0.036	6
Field collected, classified vegetation and DEM 10 metre	0.238462	0.826 +/- 0.039	6
Field collected, classified vegetation and DEM 20 metre	0.292308	0.845 +/- 0.033	5

**Table 6.12:** Summary table of 'best' networks for nine class classifications using DTCs.

Table 6.12 summarises the decision trees grown from the different training sets. Continuing the apparent trend outlined in the majority of the two and three class classifications, the DTCs grown using the DEM training data not only produced the lowest overall accuracy but also did so in combination with the highest relative cost. The trees trained using the field derived independent variables, and the field derived independent variables used in combination with the classified vegetation produced the highest overall accuracy. However, the latter of the two classifications possessed a larger error. Table 6.13 is the correlation matrix for the DTC trained using the field derived training data set. The decision tree grown for this classification can be seen in Figure 6.11.

				1		Actua	1					
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	19	2	0	2	1	2	1	0	0	27	Accuracy
Predicted	1	0	4	1	2	1	3	3	1	0	15	
	2	15	4	1	1	0	2	0	0	0	23	
	3	2	0	0	0	0	2	0	0	0	4	
	4	5	1	0	0	3	0	4	2	0	15	0.292308
	5	2	4	3	2	0	4	3	1	0	19	
	6	0	0	0	0	1	1	2	1	0	5	
	7	1	2	0	0	1	2	4	2	1	13	
	8	1	0	0	0	0	2	1	2	3	9	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	rodu	cers Ac	curac	y							
	0		42.2%				0			70.4%		
	1	1			23.5%					26.7%		
	2	2			20%					4.3%		
	3	3			-				÷			
	4		42.9%				4		20%			
	5	5			22.2%				21.1%			
	6	6			11.1%				40%			
	7			22.2	%		7		15.4%			
	8	8			ó		8			33.3%		

**Table 6.13:** Correlation matrix for the nine class DTC trained using the field acquired independent variables.

The high misclassification rates are reflected in the low producers and users accuracies seen in Table 6.13. The only exceptions were the producers accuracy for the 'severe subsurface gully erosion' class and the users accuracy for the 'no appreciable erosion' class. 'Severe rill erosion', class 3, is the only class not to be attributed a single correct case and thus no users or producers accuracies have been calculated. The erosion map produced for the nine class classifications based on the 10 metre resolution DEM data can be seen in Figures 6.12.

Soil Erosion Classifications, Risk Schedules and Rule Extraction



Figure 6.11: The decision tree grown for the nine class classification using the field acquired independent variables.



**Figure 6.12:** Classified erosion map drape derived from the DTC trained using 10 metre DEM variables for a nine class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).

# 6.4 SOIL EROSION CLASSIFICATIONS USING DISCRIMINANT ANALYSIS

Discriminant Analysis (DA) classifications were undertaken using the StatistiXL statistical add-in for Microsoft Excel. Using the same training data splits as those used in the ANN and DTC classifications, 390 training examples and 130 test examples, DA classifications were trained in order to provide an comparative traditional statistical and linear baseline from which the AI approaches can be analysed. DA was selected over other techniques as it is multivariate and allows the limited use of categorical variables, unlike many of its counterparts, a important aspect required to allow direct contrasts and comparisons to be made between the different techniques (see Table 4.1).

#### 6.4.1 Two Class Classifications Using Discriminant Analysis

Using discriminant analysis eight classifications were undertaken using the same training sets as those used for the training of the Artificial Neural Networks and Decision Tree Classifiers. Overall accuracies ranged from 63.1 percent up to a maximum of 78.5 percent depending on the training set used within the classification (Table 6.14). Unlike the classifications undertaken using the AI techniques, where the optimal solution was determined based upon some validation error or cost function, DA provides a single solution and its associated overall accuracy.

Discriminant analysis performance is highly variable depending upon the training data used. The DA possessing the optimum performance was that trained using the field, classified vegetation and 10 metre DEM independent variables. The correlation matrix for this classification can be seen in Table 6.15. The producers accuracy for both classes exceed 75 percent and the users accuracy for class 1 'erosion' exceeds 88

## Soil Erosion Classifications, Risk Schedules and Rule Extraction

percent. However, due to the high commissions error the users accuracy for the 'no appreciable erosion' class is only 64.2 percent as a further 17 examples of 'no appreciable erosion' were classified as currently eroding.

	Overall Accuracy
DEM 10 metre	0.630769
DEM 20 metre	0.676923
Field collected	0.769231
Field collected and DEM 10 metre	0.769231
Field collected and DEM 20 metre	0.753846
Field collected and classified vegetation	0.769231
Field collected, classified vegetation and DEM 10 metre	0.784615
Field collected, classified vegetation and DEM 20 metre	0.769231

Table 6.14: Summary table of the two class classifications using discriminant analysis.

		Actu	al		
		0	1	Total	Overall
Predicted	0	34	17	53	Accuracy
	1	11	68	77	
	Total	45	85	130	0.784615
	Producers A	curacy	Users A	Accuracy	
	0	75.6%	0	64.2%	
	1	80%	1	88.3%	

Table 6.15: Correlation matrix for the two class DA using the field acquired independent variables.

All of the correlation matrices relating to the DA classifications can be seen in Figures A2.49 to A2.72 inclusive in Appendix Two. The erosion map produced for the discriminant analysis using the 10 metre resolution DEM training data can be seen in Figure 6.13.



**Figure 6.13:** Classified erosion map drape derived from the DA trained using 10 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).

## 6.4.2 Three Class Classifications Using Discriminant Analysis

Table 6.16 details the overall accuracies using DA for the three class classifications. The DA based on the training data sets derived solely from the digital elevation models produced the weakest classifiers, and the strongest used the field-acquired data in combination with the 10 metre resolution DEM data. The correlation matrix produced from this classification can be seen in Table 6.17.

Large errors of omission and commission result in the discriminant analysis classification, with reasonably low producers and users accuracy's, using the field collected independent variables in combination with the DEM data with 10 metre spatial resolution. The matrix highlights the fact that omission errors occur in the 'no appreciable erosion' class and the 'rill erosion' class, and substantial commission errors are apparent in the 'rill erosion' class.

	Overall Accuracy
DEM 10 metre	0.523077
DEM 20 metre	0.523077
Field collected	0.6
Field collected and DEM 10 metre	0.615385
Field collected and DEM 20 metre	0.592308
Field collected and classified vegetation	0.523077
Field collected, classified vegetation and DEM 10 metre	0.6
Field collected, classified vegetation and DEM 20 metre	0.607692

Table 6.16: Summary table of the three class classifications using discriminant analysis.

			Actual					
		0	1	2	Total	Overall		
Predicted	0	26	4	8	38	Accuracy		
	1	11	12	8	31			
	2	8	11	42	61	0.615385		
	Total	45	29	56	130			
	Producers	s Accuracy		Users Accuracy				
	0	57.8%	(	)	68.4%			
Kin ( Starta	1	41.4%	1	1	38.7%			
and the second	2	75%	2	2	68.9%			

**Table 6.17:** Correlation matrix for three class DA using the field acquired and 10 metre DEM independent variables.

Figure 6.14 shows the erosion map produced for the classification incorporating 'no appreciable erosion', 'rill erosion' and 'gully erosion', using the independent variables extracted from the 10 metre elevation model.



Figure 6.14: Classified erosion map drape derived from a DA trained using the 10 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).

### 6.4.3 Nine Class Classifications Using Discriminant Analysis

The summary table for the classifications using all nine erosion classes highlights the low overall accuracies of all of the eight different DA classifications (Table 6.18). The maximum overall accuracy attained was 27.7 percent. These were obtained for the classifications which were trained using the field collected variables accompanied by the coarser 20 metre DEM data and that trained using the field variables, 20 metre DEM variables and classified vegetation. The correlation matrix for both classifications can be seen in Tables 6.19 and 6.20 respectively.

	Overall Accuracy
DEM 10 metre	0.2
DEM 20 metre	0.253846
Field collected	0.2
Field collected and DEM 10 metre	0.215385
Field collected and DEM 20 metre	0.276923
Field collected and classified vegetation	0.261538
Field collected, classified vegetation and DEM 10 metre	0.2
Field collected, classified vegetation and DEM 20 metre	0.276923

Table 6.18: Summary table of the nine class classifications using discriminant analysis.

						Actua	1				100000	1. Carlos and
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	18	2	1	0	1	2	2	0	0	26	Accuracy
	1	8	4	2	1	0	0	0	1	0	16	
	2	6	4	0	2	1	3	2	2	0	20	
	3	7	0	0	0	1	2	1	1	0	12	0.276923
Predicted	4	0	0	0	2	2	0	1	1	0	6	
	5	4	2	1	2	1	4	1	0	0	15	
	6	0	3	1	0	0	4	4	0	0	12	
	7	2	2	0	0	1	2	3	3	3	16	
	8	0	0	0	0	0	1	4	1	1	7	
	Total	45	17	5	7	7	18	18	9	4	130	
	1	Produc	cers Ac	curacy			-					
	0		40%				0				69.2	
	1	1			23.5%					25%		
	2	2			-					-		
	3	3			-							
	4		28.6%				4		33.3%			
	5	22.2%				5		26.7%				
	6	22.2%				6		33.3%				
	7			33.39	%		7					
	8	8			5		8			6		

Table 6.19: Correlation matrix for DA using the field acquired and 20 metre DEM independent variables.

					Actual						
	0	1	2	3	4	5	6	7	8	Total	Overall
0	19	2	2	0	1	3	1	0	0	28	Accuracy
1	8	3	1	0	0	0	1	1	0	14	
2	4	3	0	2	0	2	2	3	0	16	
3	6	0	0	0	1	2	0	0	0	9	
4	0	2	0	2	1	0	1	1	1	8	0.276923
5	5	2	1	2	1	4	1	0	0	16	
6	0	3	1	1	1	4	6	1	0	17	
7	1	2	0	0	2	2	3	2	2	14	
8	2	0	0	0	0	1	3	1	1	8	
Total	45	17	5	7	7	18	18	9	4	130	
1	Produc	ers Ac	curacy								
0	0			42.2%					67.9%	0	
1	1			17.6%					21.4%		
2					2			-			
3	-				3						
4		14.3%				4		12.5%			
5		22.2%				5		25%			
6	33.3%				6		35.3%				
7	7			22.2%				14.3%			
		-	250/			0	-	12.50/			
	0 1 2 3 4 5 6 7 8 Total 0 1 2 3 4 5 6 7 8	0         19           1         8           2         4           3         6           4         0           5         5           6         0           7         1           8         2           Total         45           Product         0           1         2           3         4           5         6           7         8	0         1           0         19         2           1         8         3           2         4         3           3         6         0           4         0         2           5         5         2           6         0         3           7         1         2           8         2         0           Total         45         17           Producers Ac         0         1           2         3         3           4         5         6           7         8         2			$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	0         1         2         3         4         5         6         7         8           0         19         2         2         0         1         3         1         0         0           1         8         3         1         0         0         0         1         1         0           2         4         3         0         2         0         2         2         3         0           2         4         3         0         2         0         2         2         3         0           2         4         3         0         2         0         2         2         3         0           3         6         0         0         1         1         1         0         0           4         0         2         0         2         1         1         1         1         0         0           6         0         3         1         1         1         1         0         0         1         3         1         1         1         1         0         0         1         3         1	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$

**Table 6.20:** Correlation matrix for DA using the field acquired, 10 metre DEM independent variables and classified vegetation.

Finally, the erosion map produced by the DA classification technique trained using the 10 metre DEM independent variables can be seen in Figure 6.15.



**Figure 6.15:** Classified erosion map drape derived from the DA trained using 10 metre DEM variables for a nine class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).

166

### **6.5 COMPARISON OF CLASSIFICATION TECHNIQUES**

The presentation of the results in for all of the classifications undertaken using ANNs, DTCs and DA allows for contrasts and comparisons to be made between the different techniques (sections 6.2, 6.3 and 6.4). An important aim outlined in Chapter One is to better understand and determine the usefulness and applicability of AI techniques for the spatial mapping of erosion processes. In order to fulfil this aim, a detailed analysis of the results is required, and this has been undertaken within the following subsections.

As expected, the overall accuracy of the classifications deteriorates as the complexity of the task increases. The best results are achieved for the two class classifications, followed by the three class and then the nine class classifications. Although this general trend is readily identifiable throughout the results obtained and presented in the previous sections of this chapter, the classification technique and the data set used to train the classifier does influence overall performance.

### 6.5.1 Comparison of Techniques for Two Class Classifications

The first point of interest concerning the two class classifications is the significantly reduced overall accuracies achieved using the digital elevation model training data for both the DTCs and DA. However, the ANN classifications do not appear to improve drastically through the incorporation of the field acquired variables, but do outperform both the DTCs and DA using both resolutions of training data provided by the elevation models. The DTCs and DA by contrast do appear to improve with the incorporation of the independent variables collected in the field. Neither of these two techniques appears to greatly outperform the other in any of the eight models
produced, however, they both tend to produce improved results over the neural networks.



Figure 6.16: Overall accuracies achieved for the two class classifications.

As discussed previously, only classified erosion maps can be produced for each of the classifications based upon the varying resolution DEM data, as other independent variables incorporated within other data sets are not known throughout the entire study area. Nonetheless, comparing and contrasting the maps produced from these classifications provides important inferences into the potential strengths and weaknesses of each method and the associated advantages and disadvantages of the different techniques. Some of the developed erosion maps have been shown in the previous sections of this chapter, and the remainder can be seen in Appendix Six. In addition, a Digital Versatile Disk (DVD) is located on the inside of the rear cover of this thesis containing all of the erosion maps in PDF format for more detailed inspection.

Figure 6.16 incorporates the overall accuracies for each of the three classification techniques using all eight training data sets. It appears that in terms of overall accuracy, the ANN classifications are not highly variable across the range of training data sets. This is particularly evident in the classifications using the remotely obtained data from the DEMs. On both occasions ANNs produce the best classification of all three methods in terms of overall accuracy, yet do not appear to drastically improve with the incorporation of the independent variables collected in the field. However, the potential problems associated with these classifications are highlighted in the erosion maps produced from the 10 and 20 metre DEM data (Figures 6.2 and A6.3). It is readily evident that the two erosion maps have classified the vast majority of the cells within the study area as currently eroding, as opposed to showing no appreciable signs of erosion. The correlation matrices for the two neural networks from which the maps have been produced highlight the fact that of the 130 test cases presented, 113 were classified as eroding in the 10 metre DEM network (Table A2.1) and 122 in the 20 metre DEM network (Table A2.2), implying their apparent inability to distinguish between the two classes. However, the ROC curves for the classifications reflect the poor performances with AUC values of 0.67 for the higher resolution data and 0.63 for the more coarse resolution data, indicating poor classification performances based on the parameters set out by Pearce and Ferrier (2000).

In contrast, the DTCs and DA do not appear to suffer from the same inability to separate the two classes for the 10 and 20 metre DEM classifications, and although they produce lower overall accuracies in both cases they may in fact be better classifiers. Both the DTCs and DA erosion maps produced from the 10 and 20 metre training sets are extremely similar, and in strong contrast to those of the ANNs. As

outlined previously, the neural network solutions appear to classify the vast majority of the study area as eroding with limited patches showing no appreciable erosion. Thus, although the ANNs in fact appear to produce superior results using the DEM training data, the DTCs and DA may actually be better employed. It is important to bear in mind the simple fact that if every case in the test set were classified as eroding then the overall accuracy would exceed 65 percent. Furthermore, the area under the ROC curve for the decision tree DEM classifications reflect far superior classifiers with values of 0.74, using the 10 metre variables, and 0.76 for the 20 metre variables. Unfortunately the receiver operating characteristics cannot be calculated using the DA technique within the StatistiXL software, however, the similarities between the two methods can be seen visually through the classified erosion maps (Figures A6.11 and A6.23). It is however clear that the DTCs have classified more of the study area as eroding than the DA for both DEM classifications, a point further emphasised within the correlation matrices. Within the test data set, 45 cases of 'no appreciable erosion' exist along with 85 'erosion' cases, making up the total of 130 cases. Within the DA classifications 15 and 17 of these were misclassified as eroding for the 10 and 20 metre DEM classifications respectively (Tables A2.49 and A2.50), compared with 18 and 23 for the DTCs (Tables A2.25 and A2.26), indicating the tendency for the latter method to potentially classify more of the 'no appreciable erosion' cells as eroding.

As indicated in Figure 6.16, DTCs and DA appear to outperform the ANNs for the two class classifications in terms of overall accuracy (with the exception of the previously discussed DEM classifications). Moreover, the comparison of the different ROC curves derived from the classifications tends to support the superiority of the decision trees for the two class classifications in the majority of cases. All eight

# Soil Erosion Classifications, Risk Schedules and Rule Extraction

classifications undertaken using DTCs, possessed ROC curves where the AUC exceeded 0.7, three of which exceeded 0.9, indicating excellent classifiers. Table 6.21 details the AUC statistics derived from the respective ROC curve for each of the eight different two class classifications for both ANNs and DTCs. The statistics are comparable to one another implying little variation between the two techniques for all of the classifications trained incorporating the independent field variables. The best classifier based on the ROC curves was the decision tree trained using the field acquired variables with the classified vegetation and the 10 metre DEM data, with 92 percent of the area under the curve. Figure 6.17 shows both the ANN and DTC ROC curves for the classifications emphasising the strong performance of both, but in particular the excellent decision tree classifier produced for the two class problem.

	Area Under ROC Curve		
	ANN	DTC	
DEM 10 metre	0.67	0.74	
DEM 20 metre	0.63	0.76	
Field collected	0.88	0.86	
Field collected and DEM 10 metre	0.88	0.91	
Field collected and DEM 20 metre	0.88	0.91	
Field collected and classified vegetation	0.88	0.77	
Field collected, classified vegetation and DEM 10 metre	0.88	0.92	
Field collected, classified vegetation and DEM 20 metre	0.84	0.77	

**Table 6.21:** Area under curve comparison table for the ANNs and DTCs for the two class classifications.



Figure 6.17: ROC curves for both the ANN and DTC classifications trained using the field acquired independent variables, classified vegetation and the 10 metre DEM attributes.

#### 6.5.2 Comparison of Techniques for Three Class Classifications

Of the different classification techniques, DTCs and DA produced better results for the three class classification and, as with the two class problem, appear to outperform the neural networks. The overall accuracies achieved through each of the three different techniques can be readily compared and contrasted with one another in Figure 6.18. It is clear that the weakest classifiers were produced from the two digital elevation model training sets of differing resolutions. The maps produced from these classifications can, as with the two class classifications assist in the understanding and analysis of different classifier performances ad highlight advantages and potential shortfalls associated with the different techniques.



Figure 6.18: Overall accuracies achieved for the three class classifications.

The foremost point of interest associated with the three class classification maps concerns the neural networks. The maps derived from the ANN 10 and 20 metre DEM classifications (Figures A6.5 and A6.7), highlight the extreme bias towards the 'no appreciable erosion' and 'gully erosion' classes, at the expense of the 'rill erosion' class. As discussed previously in section 6.5.1, the ANNs overall performance tend to suffer, largely as a consequence of not classifying any reasonable number of cases for one of the erosion groups, possibly due to the learning procedure and the training set used. The training set comprises a total of 390 cases, of which 152 are 'no appreciable erosion', 66 are 'rill erosion' and 172 are 'gully erosion'. As can be seen in Figures A6.5 and A6.7, rill erosion is virtually unclassified, with a total of 10 cells attributed in the 10 metre DEM classification and none in the 20 metre classification. This problem is not evident in the DEM classifications for either the DTCs (Figures A6.15 and A6.17) or the DA (Figures A6.27 and A6.29) which classify significant portions of the study area as eroding through rilling. Table 6.22 documents the composition of

the training set along with the composition of the erosion maps produced (10 and 20 metre) to aid in the determination of any possible detrimental influences caused by the training data.

	Class	1	2	3
Training Data Set Composition	Cases	152	66	172
	%	39	16.9	44.1
	Class	1	2	3
Artificial Neural Networks	10 Metre DEM	242103	10	561937
	%	30.11	0.0012	69.88
	20 Metre DEM	49828	0	151347
	%	24.8	0	75.2
A PERSONAL WARRANT	Class	1	2	3
	10 Metre DEM	281296	158354	364400
Decision Tree Classifiers	%	34.98	19.69	45.32
	20 Metre DEM	58232	59433	83510
	%	28.95	29.54	41.51
	Class	1	2	3
	10 Metre DEM	411303	140622	252125
Discriminant Analysis	%	51.15	17.49	31.36
	20 Metre DEM	71853	69703	59619
	%	35.72	34.65	29.63

**Table 6.22:** Composition of erosion maps for all three techniques for the DEM classifications using each of the three class classifications.

The results outlined in Table 6.22 emphasise a distinct weakness associated with the neural network technique in comparison with the DTCs and DA. The rill erosion class makes up the smallest group of the three within the training data set and it would appear that this has a significant impact upon neural network performance, signifying the potential swamping effect that may be a causal factor. In both classifications undertaken using the DEM training sets, not only do the ANNs appear to ignore the 'rill erosion' class but also classify the majority of the study area as 'gully erosion', the largest group in the training set. This suggests that their performance is strongly influenced by the training set composition. In contrast, neither the DTCs or DA appear to be constrained in such a manner, and do have a far more even spread of classified classes. As discussed previously, the overall accuracy results may appear somewhat misleading as the ANN classification trained with the coarser of the two

DEM data sets appear to outperform both of the other two techniques. It is evident that it fails to classify a single case out of a total of 201175 cells as eroding through rill processes, and as with the two class classifications, produces superior overall accuracies compared with DA or DTCs. This is achieved by classifying cases in the test set in the most prominent classes, resulting in high errors of omission and commission.

A similar trend is also evident within the six classifications developed with the independent variables collected in the field. Through examination of the correlation matrices, it is apparent that the neural networks perform poorly for the three class problem. The overall accuracies suffer as a result of the minor inclusion of rill erosion cases.

#### 6.5.3 Comparison of Techniques for Nine Class Classifications

All of the classifiers trained and produced for the nine class classification problem performed poorly irrespective of the data set used to train them. A neural network containing three hidden nodes produced the best overall accuracy classification, 39.2 percent. Figure 6.19 shows the overall accuracies attained by the ANNs, DTCs and DA, highlighting the poor performances associated with each. Of the three methods, the neural networks produce the highest overall accuracies, followed by the DTCs in the majority of the cases with the DA being the weakest. However, as with the two and three class classifications, the overall accuracies produced from the ANNs are somewhat misleading because they tend to classify the vast majority of unknown cases into one or two erosion classes.



Figure 6.19: Overall accuracies achieved for the nine class classifications.

The problem is further highlighted through the erosion maps produced from the DEM classifications. The map produced from the ANN 10 metre DEM classification has not been shown as every cell was classified in the same class, namely 'no appreciable erosion'. Figure 6.4, the erosion map for the 20 metre DEM classification, shows the same problem but to a lesser extent whereby the vast majority of the coverage is classified as 'no appreciable erosion', and only six of the eight classes have had cases attributed to them. Table 6.23 contains the composition of the training data set and the number of cells that have been classified within each class for both the 10 and 20 metre DEM classifications for each of the three different techniques.

interesting the extension when the fetering we concerns to feel to not be presented to a signed the life checkle during at the LIFE's however it is in redshift to compare the marked take to a second of the scale of the schuld or present take or every call writer 0 of marked takes to the scale of the scale of the schuld or present takes the peterin "for or or interesting the scale of the scale of the schuld or present to the scale of the sc

Training	Class	0	1	2	3	4	5	6	7	8
Data Set	Cases	152	44	12	12	22	54	59	27	8
Composition	%	38.97	11.28	3.08	3.08	5.64	13.85	15.13	6.92	2.05
READ TO THE	Class	1	2	3	4	5	6	7	8	9
Artificial Neural Networks	10 Metre DEM	804050	0	0	0	0	0	0	0	0
and the second se	%	100	0	0	0	0	0	0	0	0
li enne en	20 Metre DEM	83732	104901	19	0	0	5	1611	1896	9011
	%	41.621	52.144	0.009	0.000	0.000	0.002	0.801	0.942	4.479
	Class	1	2	3	4	5	6	7	8	9
Decision Tree Classifiers	10 Metre DEM	239639	139580	76487	5740	76346	77354	57372	67789	63743
Chastine	%	29.80	17.36	9.51	0.71	9.50	9.62	7.14	8.43	7.93
Calaling Setter	20 Metre DEM	38831	46163	16351	2097	19850	19812	19936	29183	8952
and the	%	19.30	22.95	8.13	1.04	9.87	9.85	9.91	14.51	4.45
	Class	1	2	3	4	5	6	7	8	9
Discriminant Analysis	10 Metre DEM	317146	74105	35255	59551	30919	10309 7	29068	34069	120840
	%	39.44	9.22	4.38	7.41	3.85	12.82	3.62	4.24	15.03
	20 Metre DEM	43314	49382	8186	23274	6613	27653	14532	10606	17615
	%	21.53	24.55	4.07	11.57	3.29	13.75	7.22	5.27	8.76

**Table 6.23:** Composition of erosion maps for all three techniques for the DEM classifications for the nine class classifications.

Through analysis and interpretation of the results outlined in Table 6.23, it would appear that the composition of the training data set exerts strong influences upon the neural network classifications undertaken for the nine class problem, as evident in the cases of those trained using the elevation model data. It is clearly evident that the ANNs attribute the vast majority of the unknown cases in the study area to the largest constituent classes within the training set. The same effect is not as pronounced in either the DA classifications or the DTCs, however it is impossible to comment extensively upon these results as the actual or ground truth of every cell within the study area is unknown. Thus, it may in fact be the case that certain classes are inherently less prominent than others irrespective of the composition of the training set. Nonetheless, the figures presented in Table 6.23 are such that it is readily apparent that the neural networks appear to be strongly influenced by the training process to such an extent that some classes have no cases attributed to them. For example, class six, 'minor subsurface gully erosion', of which the training data comprises nearly 14 percent, had no cells attributed within the 10 metre classification and only five in the 20 metre classification.

With the incorporation of the independent variables collected in the field within the training process, it is clear that all three classification methods produced improved results. The ANNs achieved higher overall accuracies than either the DTCs or the DA in any of the eight classifications. However, the analysis of the correlation matrices further emphasises the potential swamping effect associated in particular with the neural networks. The correlation matrix for the optimum neural network trained for the nine class problem seen in section 6.2.3 (Table 6.7) grouped all 130 test cases into only three of the possible eight erosion classes. Likewise, none of the other seven classifications attributed unknown test cases to all nine classes and in the worst instance only a single erosion class was used in the 10 metre DEM classification, up to a maximum of five (20 metre DEM classification).

In contrast to the ANNs, the DTCs and DA techniques do not appear to suffer from the same swamping problems. In both instances, either the majority of classes have been attributed unknown cases or all of them have, particularly in the DA classifications.

# 6.6 THE SELECTION AND INFLUENCE OF INDEPENDENT VARIABLES ON CLASSIFIER PERFORMANCE

An extensive array of independent and dependent variables have been used in combination with one another so as to aid the understanding of their influence and effect upon classification outputs. An important aim set out in Chapter One concerned the issue of selecting appropriate dependent and in particular independent variables for such classification problems, in order to provide some insights into the costbenefits associated with different techniques.

Overall accuracy and classifier performance is variable between those classifications using only independent variables extracted from the two digital elevation models and those incorporating the field collected variables within the training stage. This general trend is evident throughout each classification problem, namely two class, three class and nine class classifications. However, to fulfil the stated aim it is important to determine the extent to which the trend holds and the importance and influence of individual independent variables as opposed to the overall data set.

Discussed in Chapter Four, ANNs and DTCs can determine the usefulness and applicability of individual predictors and weight them accordingly, so that independent variables, believed to have little significance upon classifier performance, can largely be ignored, in contrast to those that may be heavily weighted due to their perceived value. Thus, through the analysis of such weightings, proposals and suggestions can be forwarded to assist and aid future research.

#### 6.6.1 Independent Variables Extracted from Digital Elevation Models

A total of six independent variables were derived from the two DEMs, including slope angle, aspect, flow length and accumulation, and plan and profile curvature (see Table 5.1). Sensitivity analysis can be carried out to determine the contribution of each independent variable to the dependent variable in a neural network. They have been used extensively in studies where neural networks have been applied, in an attempt to better understand the neural solution identified (Park and Chung, In Press; Pastor-Bárcenas *et al.*, 2005; Olden and Jackson, 2002; Jaiswal *et al.*, 2005). Table 6.24 and 6.25 detail the sensitivity analysis for the neural network trained for the two class classification using the 10 metre and 20 metre resolution data respectively. Each variable is given a rank and an error value for both the training and verification stages. The rank is the most important variable with regards to overall network performance (1 being the most important), and the error is the error rate that would occur if the variable were not included in the model.

	011-81	Slope Angle	Aspect	Flow Length	Flow Accumulation	Plan Curvature	Profile Curvature
Train	Rank	1	2	4	6	5	3
	Error	0.5095	0.4971	0.4858	0.4831	0.4846	0.4910
Verify	Rank	1	2	4	3	6	5
	Error	0.4678	0.4609	0.4571	0.4592	0.4560	0.4570

**Table 6.24:** Sensitivity analysis for the ANN trained using the 10 metre DEM data set for the two class classification.

		Slope Angle	Aspect	Flow Length	Flow Accumulation	Plan Curvature	Profile Curvature
Train	Rank	1	6	2	5	4	3
	Error	0.4978	0.4825	0.4892	0.4843	0.4844	0.4844
Verify	Rank	2	5	1	4	3	6
	Error	0.4630	0.4594	0.4682	0.4594	0.4619	0.4593

**Table 6.25:** Sensitivity analysis for the ANN trained using the 20 metre DEM data set for the two class classification.

It can be concluded from the sensitivity analysis that slope angle would appear to be the most important independent variable derived from the elevation models for the neural network classifications. Of the other five variables, the importance varies between both training and verification data sets as well as between the different resolutions, yet the margin of error between them is extremely small and in some instances is negligible. Nonetheless, the analysis does show that all six of the inputs in the networks assist, to some degree, in the classification process.

The variable importance can be calculated for the DTCs, producing a scored index based upon the contribution of each independent variable, taking into account its role as a primary splitter and as a surrogate to any of the primary splitters in a tree (Salford Systems, 2004). The most important variable is assigned a value of 100, whilst a value of zero indicates that the variable played no role in the analysis. The variable importance index however, concerns all trees grown during a classification procedure, not just the tree selected. Therefore, caution must be exercised when analysing the results. Unlike ANNs however, DTCs are explicit in their topology and through viewing the tree itself inferences can be made regarding the individual influence of various independent variables.

elevele energy of a set	Variable Importance
Aspect	100
Flow Length	82.17
Slope Angle	73.07
Profile Curvature	15.34
Flow Accumulation	3.87
Plan Curvature	3.58

**Table 6.26:** Variable importance for the DTC trained using the 10 metre DEM data set for the two class classification.

he ter the last lines for the black of	Variable Importance
Flow Length	100
Flow Accumulation	81.70
Slope Angle	75.34
Aspect	35.08
Profile Curvature	25.71
Plan Curvature	3.38

**Table 6.27:** Variable importance for the DTC trained using the 20 metre DEM data set for the two class classification.

Table 6.26 and 6.27 are the variable importance values for the two class problem using the two DEM data sets. The weakest predictor appears to be plan curvature, adding little to the overall performance in either of the classifications, and flow length can be seen to be an extremely useful predictor. Through visual analysis of tree structure, it is apparent that in the 10 metre DEM classification only slope angle, aspect and flow length are used as predictors, and slope angle, flow length, flow accumulation, aspect and profile curvature are used in the 20 metre DEM classification. In both cases, slope angle is the root node indicating its high entropy factor, but flow length and aspect provide a number of data splits.

The ANN and DTC classifications undertaken for the three class problem suggest, through the sensitivity analysis and variable importance, that the most influential independent variables are slope angle, aspect and flow accumulation. The neural network sensitivity analysis for both DEM resolutions can be seen in Appendix Five, Tables A5.9 and A5.10, which indicate that the least influential variables are flow accumulation, plan and profile curvature. For the same classifications, the decision trees also highlight flow length, slope angle and aspect as the most important independent variables, with the exception of the 20 metre classification whereby flow accumulation was in fact an important contributor as shown in Tables A5.33 and A5.34. These findings are reinforced by the sensitivity analysis for the neural

networks trained for the nine class classification (Tables A5.17 and A5.18). However, the variable importance analysis for the 10 metre decision tree trained to classify all nine classes inverts the trend, with flow accumulation, plan curvature and profile curvature being the most significant variables (Table A5.41). Nonetheless, the 20 metre classification appears to comply with the overall trend with flow length being the most influential variable and plan and profile curvatures being the least influential (Table A5.42).

### 6.6.2 Independent Variables Acquired from the Field

The advantages of using digital elevation model (DEM) data for the study of various geographical phenomena are well documented and understood. The advantages of DEMs are based largely on a combination of cost and time issues. Therefore, their incorporation within numerous investigations has been extensive, with varied degrees of success. To determine the ability of such data sets within this work and beyond, it is important to provide a comparison. Therefore, as outlined in Chapter Five, a range of independent variables were collected in the field and used to train the various classifiers. The following discussion details the apparent ability of each of the field collected independent variables as predictors of soil erosion.

The sensitivity analysis undertaken for the three neural network classifications, two class, three class and nine class, can be seen in Tables 6.28, 6.29 and 6.30 respectively. Taking all of the classifications into account using field data, geology would seem to be the variable that exerts the strongest controls on the network performances followed by the slope angle. The remaining three independent variables, aspect, estimated vegetation cover and sodicity meter, have inconsistent levels of

influence throughout the three classifications, revealing no significant underlying patterns. Interestingly the DTCs support these findings, and in all three cases slope angle is ranked as the most important. Geology is the second most important in two out the three classifications (Tables 6.31, 6.32 and 6.33). Of the less important variables, aspect, estimated vegetation and the results of the sodicity meter, the latter is the weakest. In the three class classification the sodicity meter has no influence whatsoever upon the decision tree construction.

		Slope Angle	Aspect	Estimated Vegetation	Field Sodicity Meter	Geology
Train	Rank	5	4	3	2	1
	Error	0.3857	0.4078	0.4103	0.4163	0.4283
Verify	Rank	1	4	5	3	2
	Error	0.4208	0.3823	0.3799	0.3840	0.4202

**Table 6.28:** Sensitivity analysis for the ANN trained using the field acquired data set for the two class classification.

		Slope Angle	Aspect	Estimated Vegetation	Field Sodicity Meter	Geology
Train	Rank	3	2	5	4	1
	Error	0.3932	0.3950	0.3851	0.3911	0.4083
Verify	Rank	2	5	4	3	1
	Error	0.3782	0.3416	0.3587	0.3632	0.3965

**Table 6.29:** Sensitivity analysis for the ANN trained using the field acquired data set for the three class classification.

alus th	and the	Slope Angle	Aspect	Estimated Vegetation	Field Sodicity Meter	Geology
Train	Rank	1	4	3	5	2
	Error	0.2766	0.2744	0.2763	0.2723	0.2766
Verify	Rank	1	3	4	5	2
	Error	0.2845	0.2797	0.2741	0.2738	0.2825

**Table 6.30:** Sensitivity analysis for the ANN trained using the field acquired data set for the nine class classification.

	Variable Importance
Slope Angle	100
Geology	59.61
Aspect	58.89
Estimated Vegetation	40.53
Sodicity Meter	32.56

**Table 6.31:** Variable importance for the DTC trained using the field acquired data set for the two class classification.

	Variable Importance
Slope Angle	100
Geology	69.2
Estimated Vegetation	57.21
Aspect	44.22
Sodicity Meter	0

**Table 6.32:** Variable importance for the DTC trained using the field acquired data set for the three class classification.

orado, var oekerning in I	Variable Importance
Slope Angle	100
Aspect	75.38
Estimated Vegetation	63.39
Sodicity Meter	60.20
Geology	54.80

**Table 6.33:** Variable importance for the DTC trained using the field acquired data set for the nine class classification.

Table 6.34 details the sensitivity analysis for the two class problem incorporating the classified vegetation independent variable within the network. The results are highly comparable to those using the estimated vegetation, with slope angle and geology being the most important predictors and sodicity being the least important. This trend holds for both the three and nine class classifications undertaken using the same training data. The classified vegetation variable therefore does not appear to have a significantly increased role in the class separation process within any of the classifications. However, it is important to highlight the fact that all three classifications incorporating the field variables and the classified vegetation have

		Slope Angle	Aspect	Classified Vegetation	Field Sodicity Meter	Geology
Train	Rank	1	3	4	5	2
	Error	0.4377	0.4101	0.3998	0.3927	0.4195
Verify	Rank	1	4	3	5	2
	Error	0.4537	0.3992	0.4185	0.3963	0.4210

**Table 6.34:** Sensitivity analysis for the ANN trained using the field acquired data set and classified vegetation for the two class classification.

The decision tree classifications rank the slope angle as the most important variable in all three problems. The classified vegetation however is one of the lesser important variables in the two and three class problems, attaining importance values of 0.67 and 58.1 respectively. Nonetheless, as can be seen in Table 6.35, the classified vegetation variable was ranked third in the list of importance.

	Variable Importance
Slope Angle	100.00
Aspect	69.22
Classified Vegetation	65.19
Sodicity Meter	61.69
Geology	56.12

**Table 6.35:** Variable importance for the DTC trained using the field acquired data set for the nine class classification.

### 6.6.3 Comparison of Field and DEM derived Independent Variables

Artificial neural networks and decision tree classifiers were also trained using a combination of both field collected independent variables and those extracted from DEMs. The results of the sensitivity analysis for the network trained using the field data and the 10 metre DEM variables are shown in Table 6.36. The sensitivity analysis for the other classifications are not shown here because a strong trend exists throughout the networks trained for all three problems using both the field data in combination with the 10 and 20 metre DEM independent variables. Generally

throughout all of these classifications the four predicting variables extracted from the DEMs are ranked lowest out of the nine. Thus, the five variables determined through fieldwork were deemed the most important with geology and slope angle being the most influential. The trend is apparent throughout all of the classifications, with the exact rank for each variable changing slightly but not significantly. It is interesting to point out that the error is higher for the field variables in contrast to the DEM variables, indicating that the former has a far more important role in the network.

	VISIDARS	Slope Angle	Aspect	Est. Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	2	3	4	5	1	9	7	8	6
	Error	0.4266	0.4184	0.4138	0.4135	0.4493	0.3760	0.3784	0.3769	0.3808
Verify	Rank	1	3	2	5	4	6	9	8	7
	Error	0.4546	0.4070	0,4074	0.3991	0.3996	0.3989	0.3942	0.3949	0.3968

**Table 6.36:** Sensitivity analysis for the ANN trained using the field acquired data set and the 10 metre DEM variables for the two class classification.

The decision trees grown for the classification problems using the combined data sets also appear to rank the field variables as more important than those of the DEM. The slope angle is generally ranked as the most important and the plan and profile curvatures as the least useful. A typical example of the variable importance for these classifications can be seen in Table 6.37, which indicates once again that the remotely collected DEM variables provide less assistance for the classification procedure when using decision trees.

	Variable Importance
Slope Angle	100
Geology	80.58
Aspect	60.89
Estimated Vegetation	53.61
Sodicity Meter	51.15
Flow Length	34.53
Flow Accumulation	10.4
Profile Curvature	8.7
Plan Curvature	1.64

**Table 6.37:** Variable importance for the DTC trained using the field acquired data set and the 10 metre DEM variables for the three class classification.

Finally, the training data set comprising of field collected variables, classified vegetation cover and those extracted from the DEMs produced neural networks and decision trees with similar results to those using the estimated vegetation cover. In addition, the variable importance and sensitivity also varies little from those classifications, whereby the field collected data is generally ranked as more useful independent variables than are those attained from the elevation models. The general point is illustrated through Tables 6.38 and 6.39, the sensitivity analysis and the variable importance using the field data, classified vegetation and the 10 metre DEM variables for the three class problem.

		Slope Angle	Aspect	Class. Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	2	3	4	5	1	7	9	8	6
	Error	0.3847	0.3824	0.3691	0.3568	0.3875	0.3422	0.3405	0.3406	0.3433
Verify	Rank	2	3	9	4	1	7	6	5	8
	Error	0.4405	0.4301	0.4198	0.4241	0.4428	0.4199	0.4204	0.4219	0.4198

**Table 6.38:** Sensitivity analysis for the ANN trained using the field acquired data set, classified vegetation and the 10 metre DEM variables for the three class classification.

	Variable Importance
Slope Angle	100
Geology	96.17
Aspect	65.03
Estimated Vegetation	63.56
Sodicity Meter	61.15
Flow Length	36.17
Flow Accumulation	18.22
Profile Curvature	10.4
Plan Curvature	6.42

**Table 6.39:** Sensitivity analysis for the DTC trained using the field acquired data set, classified vegetation and the 10 metre DEM variables for the three class classification.

### 6.6.4 Summary

In general, it is apparent throughout the many classifications that the independent variables measured in the field are more useful than those derived from the DEMs. Both the variable importance rankings and the sensitivity analysis clearly show this to be the case, to a lesser or greater degree, depending upon the classification in question. This is further supported through the use of the ROC curves created from the two class problems and seen in Figures 6.20 and 6.21 for the ANNs and the DTCs respectively.



Figure 6.20: ROC curves derived from each of the eight ANN two class classifications.

The ROC curves derived from the neural networks highlight the distinct improvement in the classifiers that incorporate the field data in their training, developing weak classifiers into good classifiers as suggested by Pearce and Ferrier (2000). A similar trend exists within the DTCs grown. However, in this instance good classifiers have been transformed into excellent ones through the incorporation of the field variables.



Figure 6.21: ROC curves derived from each of the eight DTC two class classifications.

In summary the overall accuracies for the classifications indicate generally improved performances when incorporating the field data sets into the training set as opposed to purely using those variables from the DEMs. This generalisation holds true for the ANNs, DTCs and DA techniques. However, the extent to which it operates varies from one classification to another.



# 6.7 RULE EXTRACTION USING RESPONSE SURFACES AND SPLITTING CRITERIA

An advantage of using ANNs and DTCs is that they work inductively, formulating rules, parameters and thresholds based on the training data they receive. Therefore, an opportunity exists to remove and review knowledge from the classifiers, and determine the ability of such techniques to produce informative rules relating to the processes. However, as discussed in Chapter Four, it is not easy to view the internal workings of artificial neural networks due to their black-box nature. Decision trees are perhaps more straightforward in their interpretation as the splitting rules are stated at each individual node. Response surfaces offer the ability to visualise the behaviour of a network by plotting two variables against one another whilst holding all others equal. The technique allows the user to view potential thresholds identified within the training data by the networks.

#### 6.7.1 Artificial Neural Network Response Surfaces

Figure 6.22 provides response surfaces for the two class classification trained using the 10 metre DEM data in combination with the field acquired variables. The first surface plots slope angle against estimated vegetation (A), and the second plots slope angle against flow length (B). Response surface A demonstrates the increased slope angle required to cause erosion as the vegetation cover increases, up until a specific point beyond which the angle of the trend appears to invert, with cells possessing lower slope angles and increased vegetation cover eroding. The second surface, plotting slope angle against flow length, shows a linear relationship whereby an increasing slope angle is required to cause erosion as the flow length decreases.



**Figure 6.22:** Response surface for (A) slope angle against estimated vegetation and (B) slope angle against flow length, for the two class classification using the ANN trained with the field variables and 10 metre DEM variables (Note: z-axis is erosion where 1 = no appreciable erosion and 2 = erosion).

The response surfaces derived from the three class classifications emphasise the apparent difficulty the ANNs have distinguishing between rill erosion and the two other classes with a single threshold boundary. This can be seen in both response surfaces shown in Figure 6.23. The surfaces do highlight the non-linear capabilities of ANNs. The first of the two surfaces plots slope angle against estimated vegetation cover (A). As may be expected, an increase in vegetation cover limits and restricts the onset of erosion even as slope angle increases. However, even slopes possessing maximum vegetation cover (100 percent), are vulnerable to erosion on relatively steep slopes. The second response surface possesses a highly non-linear form, incorporating the field sodicity meter data and slope angle (B). The response surface implies that even on very low angled slopes, if sodicity levels are significantly high then subsurface erosion may occur. It is important however to stress the limitations

associated with such graphs, largely as they only incorporate two variables whilst holding all others equal.



**Figure 6.23:** Response surface for (A) slope angle against estimated vegetation and (B) slope angle against sodicity, for the three class classification using the ANN trained with the field and 10 metre DEM variables (Note: z-axis is erosion where 1 = no appreciable erosion, 2 = rill erosion and 3 = gully erosion).

Two further response surfaces have been generated, seen in Figure 6.24 for a network trained for the nine class classification. As with the those produced for the three class problem, the response surfaces demonstrate the neural networks extremely poor class separation. The graphs have been produced for the ANN trained using the field and classified vegetation independent variables. Each contain two threshold boundaries to distinguish between three erosion classes. The two relationships plotted between slope and aspect (A), and slope and classified vegetation (B), were the only variables to portray a relationship. Of the other independent variables no threshold boundaries existed, that is to say that a single class existed in the z-axis across the entire decision region.



**Figure 6.24:** Response surface for (A) slope angle against slope aspect and (B) slope angle against classified vegetation, for the nine class classification using the ANN trained with the field variables and classified vegetation (Note: z-axis is erosion where 0 = no appreciable erosion and 7 = minor subsurface gully erosion).

#### 6.7.2 Decision Tree Classifier Splitting Criteria

The concept of entropy has been discussed in Chapter Four (section 4.3.2), where the data split offering the maximum knowledge gain is the root node, and usually the most important variable. In contrast to ANNs, DTCs produce easily interpretable rules with explicit structures. Splitting criteria is a simple parameter on which to extract rules and theories that the decision trees may have discovered existing within the training data.

The decision trees grown for the two and three class classifications with training data that incorporates the field acquired independent variables all contain slope angle as the root node (see Appendix Two). A common split value of 19 appears to provide the maximum knowledge gain. Slopes with angles equal to or below 19 degrees therefore have been classified as non-eroding regardless of other factors, however, above this further questions are generally asked depending upon the tree in question. The most common questions relate to vegetation cover, geology and sodicity. Figure 6.25 gives an example of the some of the rules extracted from the two class problem trained using the field derived data only.

RULE 1 IF Slope > 19 AND Estimated Vegetation Cover > 55 AND Sodicity Meter  $\leq$  5.5 AND Slope > 43 THEN Terminal Node 11 = Erosion RULE 2 IF Slope > 19 AND Estimated Vegetation Cover  $\leq$  55 AND Geology = Gypsum AND Aspect > 185 THEN Terminal Node 6 = Erosion

Figure 6.25: Some of the rules extracted from the DTC grown for the two class classification using field acquired independent variables.

The majority of trees grown using the field acquired data found the slope angle to be the most important variable with a threshold value of 19 degrees assigned. The trees generally split the data such that a value of 19 or less resulted in 'no appreciable erosion' and above 19 further questions would be posed. It is not so easy to determine such simple premises as the tree progresses down to its terminal nodes as these rules are based upon a number of previous questions that have been answered. The rules determining the splits are useful and can reveal important insights into the soil erosion phenomenon. For example, as with the slope angle, estimated vegetation in excess of 55 percent or equal to or less than 55 percent has been identified as a general split criteria in many of the trees. Figure 6.26 demonstrates the non-linear abilities of DTCs, detailing the rules produced from the tree grown for a binary classification using the 10 metre DEM independent variables. The two graphs detail the recursive partitioning of the data by the decision tree and simply assist in the visualisation of the decision boundaries identified. The first two rules are outlined in the first of the two graphs (top), which determining the majority of the feature space. However, a second graph is required (bottom) to further differentiate between examples plotted in this region.



Figure 6.26: The splitting criteria determined by the DTC trained using the 10 metre DEM data set.

Figure 6.27 details some of the rules extracted from a tree grown for a three class classification, using both the 10 metre DEM data as well as the field collected data. The trees grown using data incorporating the field attributes, portray subtle differences with one another. These trees largely consist of the same rules. However, it is not so easy to extract simple rules or patterns within the training data from the trees grown for the nine class classifications. The main reason for this is that due to the increased complexity of the problem. The trees grown are generally larger and far more complicated. Even in circumstances where they are relatively small, they do not distinguish between the classes well and extracting knowledge from them would provide little relevant information.

RULE 1 IF Slope ≤ 19 THEN

Terminal Node 1 = No Appreciable Erosion

RULE 2 IF Slope > 19 AND Geology = Gypsum THEN

Terminal Node 2 = Rill Erosion

RULE 3 IF Slope > 19 AND Geology  $\neq$  Gypsum AND Estimated Vegetation Cover > 55 AND Sodicity Meter  $\leq$  5.5 AND Flow Length  $\leq$  42 THEN

Terminal Node 8 = Gully Erosion

Figure 6.27: Some of the rules extracted from the DTC grown for the three class classification using field acquired and 10 metre DEM independent variables.

## 6.8 THE DETERMINATION OF OPTIMAL ARTIFICIAL NEURAL NETWORK ARCHITECTURES AND A REVIEW OF DECISION TREE TOPOLOGIES

The difficulties associated with the determination of optimum ANN topologies have been outlined in Chapter Four (section 4.2.3). This is a problem not encountered through the use of DTCs. In order to improve the understanding of such issues, neural networks were trained using a single hidden layer comprising a single node up to a maximum of 25 nodes (see Chapter Five). This is largely a process of trial and error, as suggested by Spellman (1999), and a method employed by Chen *et al.* (2002), Plumb (2002), Maier and Dandy (1998), Jaiswal *et al.* (2005) and Hussain *et al.* (1991).

The general trend that overall accuracy increases in accordance with a reduced error rate as the network topology becomes more complex (i.e. as the number of hidden nodes increases) would be expected until some specific point. However, results suggest that the identification of ideal network architectures is a difficult task. Error versus accuracy graphs can be seen in Appendix Four (Figures A4.1 to A4.24) for each of the neural networks trained for the two, three and nine class classifications. Close inspection of the graphs reveals a number of interesting points, particularly when considering each of the three classifications individually (two, three and nine classes).

The two class classifications generally indicate that the optimum number of nodes in the hidden layer is case specific. Throughout the eight classifications, the verification error appears to reach a significant low threshold, beyond which oscillation tends to occur. However, it may be the case that a lower error threshold is identified beyond

# Soil Erosion Classifications, Risk Schedules and Rule Extraction

this point. Figures 6.28 and 6.29 show the error versus accuracy graphs produced for the ANNs trained using the field acquired independent variables, and the field acquired and 10 metre DEM independent variables respectively. In both cases it is possible to identify the point at which verification error rapidly decreases, and subsequently oscillates around this level. The graphs also highlight the fact that overall accuracy tends to vary with little evidence of any trend or pattern as the number of hidden nodes increases. It is worthwhile mentioning that the networks presented with data collected in the field, the 10 metre DEM data, or a combination of the two, tend to reach a more distinguishable optimisation point.



Figure 6.28: Error versus accuracy graph for the ANNs trained using the field acquired independent variables (two class).

The error versus accuracy graphs derived for the three class neural network classifications also tend to indicate that the number of nodes in the hidden layer produces highly variable error and accuracy rates. The verification error rate within all of the classifications, with the exception of those trained with only the DEM data, tends to rapidly reduce as the networks become more complex. However, there is little evidence within the graphs revealing any distinct patterns regarding overall accuracy, and thus optimum topology is largely dictated (at least in these classifications) by the verification error. The networks that used training data incorporating the data collected in the field, tend to show more distinctively where the optimum topology appears to lie. Figures 6.30 and 6.31 for example, clearly show the optimum network architecture.



Figure 6.29: Error versus accuracy graph for the ANNs trained using the field acquired and 10 metre DEM independent variables (two class).



Figure 6.30: Error versus accuracy graph for the ANNs trained using the field acquired independent variables (three class).

200



Figure 6.31: Error versus accuracy graph for the ANNs trained using the field acquired independent variables and classified vegetation (three class).

The optimum network topology for the most complex classification which incorporates all nine erosion classes, should also be determined by verification error as overall accuracy is highly variable. The majority of the networks tend to reach their lowest verification error with relatively small architectures, largely consisting of less than ten hidden nodes. Verification error beyond this point tends to generally increase in most of the classifications, highlighted in Figures 6.32 and 6.33.



Figure 6.32: Error versus accuracy graph for the ANNs trained using the 10 metre DEM independent variables (nine class).



Figure 6.33: Error versus accuracy graph for the ANNs trained using the field acquired independent variables (nine class).

As suggested previously, DTCs tend not to suffer from the same architecture issues as ANNs. Nonetheless, examination of tree structure in relation to error (relative cost) can reveal interesting insights into the tree growth procedure and overall classifier ability. As discussed in Chapter Five (section 5.6.2), extensive decision trees are grown and subsequently pruned through the removal of terminal nodes, thereby ensuring that the optimum tree is found and not overlooked.

Graphs plotting the relative cost against the number of terminal nodes within a tree can be seen for all classifications in Appendix Four, Figures A4.25 to A4.48 inclusive. The relative cost or error would be expected to decrease as the number of terminal nodes in the decision tree increases or as the tree becomes more complex. However, this should begin to plateau at some specific point, and in some instances may begin to rise as the generalisation ability of the tree is exceeded.

Such a trend is clearly distinguishable within the error curves plotted for the various classifications. Optimum tree topologies are readily identifiable, however, few

apparent trends or patterns are evident. For example, Figures 6.34 and 6.35 show the error curves for the two class classification trees grown using the field acquired data and the field acquired data with classified vegetation respectively. The trees are trained using largely similar data sets but the architecture for the tree grown incorporating classified vegetation is far simpler.



Figure 6.34: Relative cost and terminal nodes for the field acquired independent variables (two class).



Figure 6.35: Relative cost and terminal nodes for the field acquired independent variables and classified vegetation (two class).

The number of terminal nodes within the trees grown for the three class classification also appears highly variable and ranges from four in the tree grown using the field data only to 19 in that trained using the 20 metre DEM data. The more complex nine class problem produced a range of DTCs that were grown significantly larger than any involved in the aforementioned classifications. For example, using the field collected data and the field collected data and 20 metre DEM variables, trees incorporated 31 and 51 terminal nodes respectively.
#### **6.9 EROSION RISK SCHEDULES AND POTENTIAL**

As outlined in Chapter Five, it is important to develop and produce outputs (maps) that are meaningful and useful to landscape managers. The importance or value of current erosion maps (actual) to people or organisations at a range of levels to assist and aid decision-making is understood. However, varying assortments of erosion risk maps can be developed to complement the actual erosion maps.

The following sub sections present the various erosion risk maps developed through a range of different methods. The methods include:

- Erosion probability maps constructed through decision tree growth.
- Erosion risk by association maps produced through the implementation of the methodology detailed in section 5.8.

#### **6.9.1 Erosion Probability Maps**

The different decision trees grown each classify the training data by recursively partitioning the example cases into smaller subsets until a terminal node (leaf) is reached. Due to the nature of decision tree structure – in particular the data split at each terminal node - it is possible to produce erosion probability maps as opposed to actual or classified erosion maps. For example, if a data set presented to a decision tree is split in such a way that terminal node  $\chi$  contains 100 examples, 90 of which are attributed to class A and the remainder to class B, the probability of unknown cases at the same node actually belonging to class A and class B is 90 percent and 10 percent respectively.

Using this simple concept, erosion probability maps have been produced from trees grown using the two, three and nine class classifications. As with the erosion maps however, risk maps have only been produced from the trees trained using either the 10 metre or 20 metre DEM data as this is available for every cell within the study area, whereas the field collected data are not.

Figure 6.36 is the erosion probability map produced using the DTC trained using 10 metre DEM data for the binary classification. According to the erosion map, white areas are currently eroding, and the coloured areas relate to cells that are at present not eroding but may be in reality. For example, based on the training data, blue areas on the map have a 33 percent chance of eroding. Figure 6.37 shows an erosion risk map detailing the potential for gully erosion produced from the tree trained using the 20 metre DEM data. The probabilities range from zero, to a maximum of 50 percent.

Such maps can be perceived as erosion risk maps: highlighting erosion potential in a quantitative manner. Similar maps have been constructed based upon various trees grown using either resolution DEM data, and can be seen in Appendix Six (Figures A6.35 to A6.41 inclusive) and on the DVD attached (DTCRiskmaps).



**Figure 6.36:** Erosion probability map drape produced from the DTC trained with the 10 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure 6.37: Gully erosion probability map drape produced from the DTC trained with the 20 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).

207

#### 6.9.2 Erosion Risk by Association Maps

Erosion risk created as a result of neighbouring or adjacent slopes (cells) eroding has been calculated using the methodology detailed in Chapter Five (section 5.8). Risk by association is simply a function of current soil erosion processes and the topographical nature of the landscape. Using the rules presented in Figure 5.12, risk maps have been created incorporating different degrees of complexity. Figure 6.38 is the risk by association map developed using a two class classification, where all erosion classes (see Figure 5.3) are amalgamated into a single class and the risk calculated using the stated rules.

However, Figures 6.39 and 6.40 show the risk by association maps for the more complex three class schedule, incorporating no appreciable erosion, surface, and subsurface erosion. A nine class classification map was used, and the appropriate classes amalgamated to create a new three class map (no appreciable erosion, surface erosion and subsurface erosion). It is important to note that a map produced from a three class classification could not be used. While these did differentiate between rilling and gullying they did not differentiate between surface and subsurface processes. Figure 6.39 shows the risk associated with surface erosion and Figure 6.40 shows the risk associated with subsurface erosion. The total risk by association maps shown in Figures 6.41 and 6.42 are simply the summation of Figures 6.39 and 6.40.



**Figure 6.38:** Risk by association map drape produced from the DTC trained with the 10 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



**Figure 6.39:** Risk by association map drape of surface erosion produced from the DTC trained with the 10 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



**Figure 6.40:** Risk by association map drape of subsurface erosion produced from the DTC trained with the 10 metre DEM variables for a nine class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure 6.41: Risk by association map drape of surface and subsurface erosion produced from the DTC trained with the 10 metre DEM variables for a nine class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



**Figure 6.42:** Risk by association map of surface and subsurface erosion subsection (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).

#### 6.10 CONCLUSIONS

The results obtained and collated throughout the course of the study and the implementation of the research methodology have been presented in the various sections of this chapter. The data has been presented in such a way so as to allow the implications associated with the aims and objectives set out in the introduction to be comprehensible.

The implications, regarding the results obtained and the data presented, have not been discussed here as a comprehensive discussion is undertaken in Chapter Eight. The following chapter documents the data collected and results obtained from the field investigations and contains a discussion relating solely to the findings. The discussion chapter concerns the results presented here as well as those from Chapter Seven and highlights potential causes or reasons for the findings. The data is discussed in relation to both current and previous literature in an attempt to provide a better all-round understanding of the issues and culminates in the conclusions in Chapter Nine.

## 7

### **Field Investigations**

#### 7.1 INTRODUCTION

Soil sodicity can have major impacts upon soil structure which affects soil permeability and infiltration (Sparks, 1995), and can lead to a number of undesirable physical properties (see Chapter Three). This can include poor soil structure, the breakdown of soil aggregates, surface crusting and the consequent reduction of infiltration rates (Rengasamy *et al.*, 1984). The determination of soil sodicity is an important factor along with other soil characteristics in the estimation of dispersivity potential. However, such analysis can be particularly time consuming and reliant upon laboratory facilities. Consequently it has long been desirable to produce or formulate a method by which soil dispersivity and sodicity can be readily measured or estimated in the field and which reduces the need for extensive laboratory analysis.

The Co-operative Research Centre for Soil and Land Management in Australia designed and produced a sodicity meter. The meter provides a method by which water turbidity can be measured to provide a useful indication of soil sodicity whilst in the field. However, to determine the usefulness and applicability of the meter it is important to validate it using traditional laboratory based techniques. This chapter explores potential relationships between a range of laboratory techniques and the field sodicity meter and attempts to better understand the chemical and physical processes involved in soil dispersion. The following sections detail the methods by which this was achieved, presents the results and discusses the wider implications of the findings.

#### 7.2 FIELD AND LABORATORY METHODS

Shallow subsurface soil samples were collected and analysed at each of the 520 sites visited in the field. Ideally each of the samples would be transported for laboratory testing to provide a more comprehensive analysis. This was impractical because soils should not be transported across international borders (Spain to the UK in this case) and thus a smaller subsection of the samples were collected. Fifty-five samples were chosen from the total of 520 training sites, with particular emphasis on the dispersive marls. The chosen samples were collected in and around the Mocatán catchment as the local lithologies, the TRU and MRU (see Chapter Two) appear heavily piped and provides a good location for the testing of the sodicity meter.

All soil samples were collected from the subsurface at a depth of 20 centimetres, which was sufficiently deep to avoid the thick, hard-setting crusts present at some sites. Ideally surface samples would also have been collected and analysed to assist in the determination of variations between surface and subsurface soils. However, analysing such a large number of samples during the fieldwork period would have been impractical due to time constraints. Samples to be further analysed in the laboratory were chosen at random (using computer generated random numbers), in order to compile a set of data that would allow unbiased comparisons to be made between results obtained using the field and laboratory procedures. The following sub-sections detail the methods used within this study.

#### 7.2.1 Field Sodicity Meter

Soil samples at each site were tested for sodicity using a simple field sodicity meter. The sodicity meter used in this study has been developed by the Australian Cooperative Research Centre for Soil and Land Management (CRCSLM), and calibrated using a total of sixty Alfisols and Vertisols from South Australia and Victoria (Rengasamy pers. com). The meter measures water turbidity, and is a plastic tube with a white disc at one end and a scale running down its length from non-sodic to highly sodic, and can be seen in Figure 7.1. The sodicity meter works on the simple premise that the clays present in sodic soils swell and if there is a sufficient excess of sodium then deflocculation will occur when mixed with water. Thus, the level of turbidity within the water increases, as the clays become suspended as a direct result of deflocculation.



Figure 7.1: The field sodicity meter designed by the Australian Co-operative Research Centre for Soil and Land Management.

The use of such an inexpensive meter has been recognised within this research, as it allows the characterisation of sodic soils relatively quickly and easily in the field environment, thus reducing the need for extensive laboratory analysis. However, it is important to determine the relationship between the sodicity meter observations and standard sodicity measurements. In order to assist in the determination of a relationship, the original scale was adapted to incorporate a wider range of values. This was achieved by taking the scale of non-sodic, sodic and highly sodic and substituting it for a scale of 0 to 14, where 0 is non-sodic and 14 highly sodic.

The field method for using the meter is as follows:

1. Weigh approximately 100g of soil into a clean 600 ml glass jar.

2. Add 500 ml of rainwater or distilled water to the jar to give a 1:5 ratio of soil to water. Bottled drinking water used here as it was not possible to obtain distilled water and the local drinking water may contain gypsum, influencing the electrical conductivity and thus potentially skewing the results.

3. Gently pour this water down the side of the jar without disturbing the soil at the bottom. Invert the jar slowly once and then return to its original position allowing to stand for 4 hours.

4. Lower the meter with the white disc at the bottom of the plastic tube into the suspension, until the disc is no longer visible (when viewed from above).

5. Place a moistened finger over the top of the tube and remove the meter from the suspension, with a level of liquid in the tube. Read the level against the coloured scale.

#### 7.2.2 Laboratory Methods

A number of soil analysis methods were undertaken on the samples that were collected and returned from the field. These tests included pH, Electrical Conductivity (EC), Exchangeable Sodium Percentage (ESP), Cation Exchange Capacity (CEC), Sodium Adsorption Ratio (SAR) and organic matter content. The following section details the soil analysis methods used and the procedures undertaken. All of the samples were prepared by air-drying for three days, disaggregated and passed through a 2mm sieve.

#### Soil pH and Electrical Conductivity (EC)

The following procedure was carried out on the soil samples in order to determine both the pH level and EC:

- 1. Weigh out 15 g <2mm soil into a small glass beaker.
- 2. Add 30 ml of de-ionised water to create a 1:2 soil to water ratio.
- 3. Stir each sample intermittently and leave to stand for thirty minutes.
- 4. Using a pre-calibrated pH and EC meter determine the pH and conductivity values.

#### Cation Exchange Capacity (CEC)

CEC of the soils was determined by using a standard method (Bower et al., 1952).

- 1. Weigh 4 g of soil into a 50 ml centrifuge tube.
- 2. Using a measuring cylinder, add 33 ml of 1M sodium acetate solution to the centrifuge tube. Seal the tube and shake for 10 minutes on a shaker.
- 3. Centrifuge and decant the supernatant.

- 4. Treat the sample with 2 additional 33 ml aliquots of sodium acetate solution. Each time re-suspend the sample before putting on to the shaker, and discarding the supernatant after each centrifuge.
- 5. Suspend the sample in 33 ml ethanol and shake for 5 minutes.
- 6. Centrifuge and discard the supernatant. Repeat this washing procedure another two times, re-suspending the sample each time before adding the ethanol.
- 7. Add 33 ml of 1M ammonium acetate and shake for 10 minutes.
- 8. Centrifuge and decant the supernatant into a 100 ml volumetric flask.
- 9. Repeat the extraction procedure twice more, and carefully make the contents of the flask to the 100ml mark with de-ionised water.
- 10. Determine the Na content of the solution in the flask by flame emission spectrometry. This solution usually requires dilution before analysis. Typically a 10 times dilution should bring down the Na concentration to the readable level on the flame photometer. That is 10 ml solution diluted to 100 ml. This dilution must be taken into account when making subsequent calculations.

#### Exchangeable Sodium (ES)

The procedure for calculating the exchangeable Na of the soil samples is as follows:

- 1. Weigh 4 g of soil into a 50 ml centrifuge tube.
- 2. Using a measuring cylinder, add 33 ml of de-ionised water.
- 3. Shake for 10 minutes on a shaker.
- 4. Centrifuge and discard the supernatant, removing the water soluble Na.
- 5. Add 33 ml of 1M ammonium acetate. Re-suspend the sample and shake for 10 minutes.
- 6. Centrifuge and decant the supernatant into a 100 ml volumetric flask.

 Treat the sample with 2 additional aliquots of ammonium acetate solution. Each time re-suspending, shaking and centrifuging the sample. Make the volume up to 100 ml with de-ionised water.

It is important to note here that the samples are treated with de-ionised water to begin, to remove any water soluble sodium that may otherwise skew the results. Once the exchangeable sodium had been determined it was used to calculate the Exchangeable Sodium Percentage (ESP) and the Exchangeable Sodium Ratio (ESR) using Equations 13 and 14 respectively.

$$ESP = \frac{ES}{CEC} \times 100$$
 (Equation 13)

$$ESR = \frac{ES}{CEC - ES}$$
(Equation 14)

#### Sodium Adsorption Ratio (SAR)

The SAR is calculated using the Equation determined by Richards (1954) based on the relationship between the ESR and the SAR (r = 0.923,  $r^2 = 0.852$ , n = 59), and is as follows:

$$SAR = \frac{ESR + 0.0126}{0.01475}$$
 (Equation 15)

#### Determination of Soil Organic Matter (SOM)

In order to determine the percentage of soil organic matter (SOM) present in each soil sample the following loss-on-ignition (LOI) procedure was followed (Avery and Bascomb, 1987).

1. Weigh a clearly labelled porcelain crucible.

- 2. Weigh out approximately 5 g <2mm air-dried soil directly into a crucible. Record the weight of the crucible plus the soil.
- Place the crucible containing soil in an oven at 105°C and leave for 16 hours (or overnight).
- 4. Remove the crucible from the oven using long-handled tongs and place in a dessicator to cool. Record the oven-dry weight of the soil and place in a muffle furnace at 850°C for one hour.
- 5. Remove the crucible and place in a dessicator to cool. Re-weigh the crucible plus soil sample and calculate the percentage weight loss using Equation 16.

$$LOI = \frac{furnace\_weight\_loss}{oven\_dry\_weight} \times 100$$
 (Equation 16)

#### 7.3 RESULTS AND DISCUSSION

The results obtained from these various methods allow a number of inferences to be made regarding the soil characteristics as well as the determination of any relationships that may exist between the different soil properties.

#### 7.3.1 Analysis of Laboratory Results

Using the data detailed in Table 7.1, the results obtained from the laboratory analysis and field sodicity meter, a number of relationships have been examined to determine their strength (see Appendix One Tables A1.1, A1.2 and A1.3 for detailed data). The data determined within this study was analysed with reference to domains proposed by Gerber and Harmse (1987) and Rengasamy *et al.* (1984) to determine their dispersive potential and is discussed with reference to the wider literature concerning the subject area. Finally, the relationships between the field sodicity meter and the various laboratory determined variables is explored.

SAMPLE	pН	EC (µS cm <sup>-1</sup> )	CEC	ESP	SAR	LOI	SOIL	SODICITY
	-		(cmol/kg)	%		%	ORGANIC MATTER %	METER
1	8.05	4420	13.342	13.330	11.282	3.915	2.361	0
2	8.61	8220	10.472	10.487	8.797	3.805	2.293	0
3	6.02	90	9.744	1.696	2.024	3.599	2.165	9
4	8.18	310	8.462	4.030	3,702	4.159	2.512	9
5	5.99	10580	14.239	2.540	2.621	4.183	2.527	0
6	7.38	500	13.657	44.676	55.603	4.208	2.542	0
7	7.98	6330	9.710	18.577	16.322	3.521	2.118	0
8	8.12	2620	12.691	9.742	8.172	3.984	2.404	0
9	6.08	240	8.659	3.919	3.620	4.390	2.655	10
10	8.11	80	7.104	2.320	2.465	4.896	2.969	9
11	8.78	1310	7.721	22.230	20.233	3.717	2.239	9
12	7.98	1280	17.401	5.528	4.821	4.250	2.569	0
13	7.97	1650	14.387	4.723	4.215	4.027	2.431	0
14	8.19	90	9.704	1.740	2.055	4.042	2.440	9
15	7.96	11820	9.691	20.037	17.843	3.956	2.386	0
16	8.03	130	13.855	0.555	1,233	6.996	4,269	9
17	8.35	90	7 715	2 159	2 350	6 596	4 021	
18	7 74	120	9 693	1 300	1 747	5 010	3 039	
19	8.05	110	6 796	2 433	2 545	4 474	2 707	
20	7.56	12070	22 838	14 786	12 618	3 337	2.004	
21	7.88	360	13 733	1 227	1 696	3 931	2.004	8
22	8 35	4060	15 288	18 203	16.033	3 635	2.571	
22	7 99	9460	7 623	25 310	23 830	3.000	1 096	
23	8 36	490	13 653	5 950	<u>23.039</u> 5.143	3 205	1.900	
25	5.00	9560	16,807	2 149	2 343	3 801	2 346	
25	7 78	160	5 424	2 380	2.543	5.069	2.540	
20	8.00	100	9 214	1.635	1 091	4 797	2 001	2
21	83	110	8 454	0.037	1.901	4.707	2.901	
20	8.8	7600	15 027	6 584	5.633	3 800	2.093	
29	83	7000	7 601	1 602	1 059	4 714	2.250	
30	0.5	10240	16 964	12 222	11 275	2 970	2.030	
32	82	200	10.004	1 857	2 127	3.070	2.333	
32	<u> </u>	200	12 212	0.716	1 3/3	4.410	2.072	
24	77	790	11 921	0.710	1 247	4.009	2.500	
- 34	70	1560	9 972	5.620	4 900	2.015	3.043	
30	1.0	550	12 200	4 352	3 030	4 965	2 040	0
	0.4		9 126	4.332	1 966	4.005	2.949	12
	77	170	0.120	0.927	1.000	4.010	2.734	8
30	7 5	140	14.400	0.021	1.420	7 107	180.0	<u> </u>
	1.5 E 0	260	12 255	0.521	1.404	1.107	4.33/	<u>8</u>
40	<u>9.0</u>	140	10.000 Q AEQ	1 200	1.290	4.000	2.110	<u></u>
	76	160	11 202	2 800	2 207	6.001	2.500	0
42	7.0	150	12 225	1 206	1 692	6 257	2 914	9
43		100	0.555	1 609	2 025	5 251	2 252	
	0.4	120	7 004	2 400	2.023	1 010	2 083	
40	<u>0.1</u> 7 6	70	5 550	1 570	1 026	4 660	2.503	
40	77		14 141	1 442	1.930	4.005	2.020	
4/	- 1.1 - Q A	2120	25 000	13 406	11 251	3 217	4.01	
40	0.4	4020	20.000	14 072	11 057	2 770	1.331	
49	0.1	4030	21.437	0.640	1 200	2.110	1,007	
	0.3	100	17 220	0.010	1.200	2.000	1./ 12	
	0.3	100	20.070	1 204	1.492	2.014	1.184	
52	0.5	180	47 550	12045	1./30	4.014	2.422	10
53		4390	11.000	12.015	10.112	4.342	2.625	
54	<u>- 1.9</u>	120	11.309	1.0//	1.592	4.003	2.787	8
55	7.9	1/0	14.204	0.881	1.457	5.893	3.585	9

Table 7.1: Results obtained from the laboratory analysis and the field sodicity meter.

Of the 55 soil samples analysed, only six had ESP values in excess of 15 percent: the recognised level above which soils are classed as sodic. This does not solely determine the soils dispersivity as other factors are also important but it is an important indicator in the understanding of the deflocculation process. The relationship identified by Gerber and Harmse (1987) between the ESP and CEC has been discussed in Chapter Three, and the results obtained have been plotted on the dispersivity classification graph in Figure 7.2. Based on the parameters and domains set out by Gerber and Harmse (1987), the majority of the samples analysed fall within the 'non-dispersive' and completely 'non-dispersive' categories. Nonetheless, it is evident that a relatively small number of the samples fall within the dispersive and highly dispersive categories. However, using the domains set by Rengasamy et al. (1984) based upon soil EC and SAR, the vast majority of the samples are either dispersive (class 1) or potentially dispersive (class 2a and 2b). This is due largely to the fact that the EC is not sufficiently high to ensure the flocculation of the samples, and thus, in samples where the SAR exceeds 3 the soil is dispersive, and below 3 potentially dispersive (Figure 3.11).

As can be seen from Figure 7.2, the majority of the samples appeared to be eroding through subsurface processes in the field (yellow points) and surface processes (blue points) and only a small number were not eroding under any visible process (orange points), based upon the field classifications used within this study (see Figures 5.3 and 5.4). As would be expected, the soils that are not eroding in the field plot within the 'completely non-dispersive' and 'non-dispersive' fields within the diagram. However, as stated previously a large number of the samples, including ones that have eroded through piping, also plot within these classes implying some confusion with the

classification scheme. Three samples plot within the dispersive and highly dispersive classes that were eroding, but through surface processes.



**Figure 7.2:** The laboratory results for CEC and ESP plotted on the dispersivity classification graph identified by Gerber and Harmse (1987). Where: VD, very dispersive; HD, highly dispersive; D, dispersive; MD, marginally dispersive; ND, non-dispersive and CD, completely non-dispersive (NB. 3 samples have not been plotted as their ESP was in excess of 20 percent, thus exceeding the range of the graph).

Typical sodic soils have a pH in excess of 8.5, as well as an ESP above 15 percent (Brady and Weil, 2002; Sparks, 1995; Mzezewa *et al.*, 2003). However, Naidu *et al.* (1995) highlighted that soils with ESPs as low as 5 can display sodic soil characteristics in situations whereby associated variables such as EC and organic matter suit. The samples analysed in this study have pH values between 7.5 and 8.5, with a small number both exceeding this and others below. Figure 7.3 shows the relationship between ESP and pH. As can be seen, the relationship is weak (r = 0.09,  $r^2 = 0.01$ ), however, there is a general trend where an increased pH value coincides

with an increased ESP. This is generally expected due to the fact that an increase in sodium will render a soil more alkaline and thus increase the pH level. Nonetheless, soil mineralogy and buffering capacity will exert strong controls on the level of pH. While the general trend can be seen to increase, a number of samples cluster about the x-axis at a range of pH levels and thus would appear to reduce the level of correlation between the two parameters. Mzezewa et al. (2003) found whilst identifying sodic soils in Zimbabwe for reclamation and improvement strategies, that observed high pH values were generally associated with high ESP values. However, the difficulty associated with using pH to indicate the presence and extent of sodium within a soil was highlighted by Fireman and Wadleigh (1951). It was recognised that soil pH can be influenced by numerous other factors, including adsorbed cations, soil-to-water ratio, texture, carbon-dioxide pressure, insoluble carbonates, gypsum, soluble salts, organic matter and the type of clay mineral. However, it is a useful indicator as McBride (1994) suggested that; with all other things equal clays with a given sodicity became more dispersible as pH increased.

Faulkner *et al.* (2000) and Alexander *et al.* (1999) identified relationships between pH and SAR for three badland sites, Vera, Tabernas and Mocatán. A strong relationship was found where r = 0.89 (n = 10) at the Mocatán site, and r = 0.82 (n = 12) at the Vera site. However, it was determined that 'signatures' exist for each site and significantly vary. Generally, the Mocatán samples had SAR values well in excess of those from Vera producing a different relationship between the two parameters. The ESP and SAR ranges obtained here are generally comparable to the Vera data set and not the Mocatán site where SAR values reached a maximum around 400. This may be due to the fact that Faulkner *et al.* (2000) sampled in and around pipe systems and on

recently cleared locations on the TRU (see Section 2.4), known to be highly sodic and discussed in detail in section 2.4. Furthermore, Faulkner *et al.* (2003a) identified a series of different relationships between SAR and pH for a single gully in the Vera badlands. The study identified significant variations between surface and subsurface SARs, and differences in the relationship between SAR and pH at the base, middle and top of the slope. It was concluded that the high level of complexity suggested by the spatial variation identified within the results, made it difficult to formulate simple conclusions.



Figure 7.3: The relationship between ESP and pH.

Figure 7.4 shows the relationship between ESP and EC. A wide range of values exist for both the parameters, but the majority of the samples analysed had low levels of exchangeable sodium in combination with low ECs. A general trend of increasing ESP with EC can be seen as indicated by an  $r^2$  value of 0.56. Faulkner *et al.* (2000)

and Alexander *et al.* (1999) identified a strong log-transformed relationship between SAR and EC (r = 0.75) for the three different badland sites combined. The Mocatán badlands in particular had a significant relationship (r = 0.71, n = 14), and all except one sample plotted in the 'dispersive' domain determined by Rengasamy *et al.* (1984). Imeson *et al.* (1982) also identified dispersive soil types in Morocco based upon the relationship between ESP and EC, and SAR and EC, both studies highlighting the usefulness of the parameters.



Figure 7.4: The relationship between ESP and EC using the power function.

The relationship between SAR and EC can be seen in Figure 7.5, plotted on logtransformed axes. It was found that the EC could account for 54 percent of the variability found in the SAR values ( $r^2 = 0.54$ ). As outlined previously, nearly all samples tested are classified as dispersive or potentially dispersive based on the domains outlined by Rengasamy *et al.* (1984). High levels of electrolyte concentration removes the tendency for clay minerals to swell and disperse (Kamphorst and Bolt, 1978; Walker, 1997; Sparks, 1995) and thus enhancing and promoting flocculation. Even at levels of ESP as low as 5, accompanying low levels of electrolyte can make the soil structure weak and increase the risk of erosion (Qadir and Schubert, 2002).



Figure 7.5: The relationship between SAR and EC plotted on log transformed axes using the power function.

Organic matter content in the 55 soils analysed ranged from 2.01 to 7.11 percent, with a mean of 4.45 and a standard deviation of 1.1 using the LOI method to approximate organic content. Using a relationship, identified by Hooda (1992), between loss-onignition and soil organic matter based on 43 soil samples from various locations within the UK, actual values of organic matter have been estimated. The relationship for calculating soil organic matter (SOM) can be seen in Equation 17 and possesses an  $r^2$  value of 0.96. This has been used as a correction factor as the LOI method overestimates organic content (Marx *et al.*, 1999) as soils containing appreciable quantities of clay lose 'structural' water, and CaCO<sub>3</sub> loses CO<sub>2</sub> to form calcium oxide at temperatures around 770°C (Rowell, 1994).

$$SOM = -0.0622 + (0.619 \times LOI)$$
 (Equation 17)

The correction factor has been applied to all of the LOI results and can be seen, along with all of the other results obtained, for each of the 55 samples in Table 7.1. Whilst often at the same time, organic matter can both suppress swelling and enhance dispersion (Churchman *et al.*, 1995; Oades, 1984), it is generally agreed that it is useful for counteracting the unfavourable effects of exchangeable sodium in soils (Richards, 1954). It has been proposed that some organic compounds, especially low molecular weight humic substances, can destabilise soil aggregates (D'Acqui *et al.*, 1999), and these organic anions increase the negative charge of mineral colloids, thereby increasing the density of the diffuse layer of cations thus promoting dispersion (Oades, 1984). Nonetheless, low levels of soil organic matter (around 3 percent or below), may indicate a potential problem for sodic soils as it is a useful flocculating agent binding aggregates (Singer and Le Bissonnais, 1998).

A number of the samples analysed, after the LOI correction factor has been applied, have low SOM contents (below 3 percent). The maximum SOM was 4.34 percent and the minimum 1.18 percent, with a mean of 2.69 and a standard deviation of 0.68. These values suggest that a number of the samples therefore may be poorly structured and weakly aggregated. Not only may these soils be susceptible to deflocculation when other variables suit, but they will also be susceptible to surface erosion processes. As highlighted in Figure 7.2, the majority of the samples appeared to be eroding through subsurface processes. However, according to the parameters set by Gerber and Harmse (1987), only a small number of the samples were dispersive or moderately dispersive. The soil organic content appears to be relatively low in the majority of the samples, with only 12 of the 55 samples actually having in excess of 3 percent SOM. The EC of the samples does not appear to be considerably high and soils with low ECs and low soil organic contents will have light textures, and may not maintain flocculation when wet.

Rengasamy et al. (1984) found that dispersion may occur even in soils at very low SARs and that may be evident in some cases here. Faulkner et al. (2000) highlighted the susceptibility of the Mocatán badlands to subsurface processes as a result of the inherent deflocculation of the sensitive lithological units. It is clearly evident however, that the results obtained here vary somewhat from those presented by Faulkner et al. (2000). The most distinguishable variation being that the SAR values appear much lower in this investigation. The importance of the landscape morphology may go some way to explaining the subsurface processes that are evidently operating on soils but which may not appear highly susceptible. Pipes preferentially develop where large hydraulic gradients exist (see Figure 3.7), such as behind terrace walls. The Mocatán badlands offer the ideal setting for these processes, and it was proposed by Faulkner et al. (2000) that in soils with low clay contents deflocculation may not cause the reduction of hydraulic conductivity as not all pore spaces are filled, and thus erosion can continue unabated. Usually dispersion will lead to the movement of clay particles into a region of 0.1-0.5 mm depth where they clog conducting pores (Mamedov et al., 2002).

Table 7.2 gives details including the slope angle and aspect and estimated vegetation cover regarding the location from which each sample was collected. Furthermore, the process and level of erosion, if any, is given based on the erosion scale presented in Chapter Five (see Figure 5.3), and finally the sodicity meter readings are also given. The Table has been produced in an attempt to better understand the influence and controls that the surrounding physical features may have upon the erosion processes and the relationship with the soil physico-chemical parameters.

It becomes readily apparent that topography is strongly influential with regards to the extent of erosion processes (through visual inspection of Table 7.2). For example, 14 of the 55 samples were documented as eroding in classes 7 and 8, severe subsurface gully erosion. The majority of the slope angles associated with these classes are in excess of 40 degrees, and by and large, the steepest slopes of all the samples. The aspect of these 14 sites varies, as does the extent of vegetation cover. In addition to this, almost all of the sites that do not appear to be eroding under any visible process are flat and have varying degrees of vegetation cover. Therefore it appears that topography, and indeed the slope angle, is highly influential in controlling the extent of erosion. This may offer a potential explanation as to why so many of the soils appeared to be dispersive in the field, yet do not display drastically high ESP values or other critical soil parameters in the laboratory. In environments such as Mocatán, where topography allows steep hydraulic gradients to exist, it is possible to suggest that even soils containing relatively low levels of exchangeable sodium may have weak soil structures leading to deflocculation and extensive pipe development, thus contributing to its highly distinctive morphology described by Faulkner et al. (2000).

	SAMPLE	SLOPE	SLOPE	VEGETATION	EROSION	SODICITY
1 45 360 60 8 0   2 50 340 10 7 0   3 18 30 60 6 9   4 0 -1 20 0 9   5 44 20 10 7 0   6 40 270 10 8 0   7 38 200 10 6 0   9 36 120 40 5 10   10 0 -1 10 0 9   11 38 90 40 6 9   12 40 120 30 6 0   13 34 50 30 6 0   14 12 10 80 5 9   15 50 340 30 7 0   16 30 20 60 4 9   17 42 230 30 6 8   20 <		ANGLE (°)	ASPECT (°)	COVER %	(SCALE)	METER
2 50 340 10 7 0   3 18 30 60 6 9   4 0 -1 20 0 9   5 44 20 10 7 0   6 40 270 10 8 0   7 38 200 10 6 0   9 36 120 40 5 10   10 0 -1 10 0 9   11 38 90 40 6 9   12 40 120 30 6 0   13 34 50 30 7 0   16 30 20 60 4 9   17 42 230 30 6 8   18 30 340 80 6 9   20 38 180 30 6 0   21 50 20 90 4 8   22	1	45	360	60	8	0
3 18 30 60 6 9   4 0 -1 20 0 9   5 44 20 10 7 0   6 40 270 10 8 0   7 38 200 10 6 0   9 36 120 40 5 10   10 0 -1 10 0 9   11 38 90 40 6 9   12 40 120 30 6 0   13 34 50 30 6 0   14 12 10 80 5 9   15 50 340 30 7 0   16 30 20 60 4 9   17 42 230 30 6 9   19 0 -1 50 5 9   20 38 180 30 6 2   22 <	2	50	340	10	7	0
4 0 -1 20 0 9   5 44 20 10 7 0   6 40 270 10 8 0   7 38 200 10 6 0   9 36 120 40 5 10   10 0 -1 10 0 9   11 38 90 40 6 9   11 38 90 40 6 9   12 40 120 30 6 0   13 34 50 30 6 9   15 50 340 30 7 0   16 30 20 60 4 9   17 42 230 30 6 0   21 50 20 90 4 8   22 46 180 30 6 2   23 46 320 50 7 0   24	3	18	30	60	6	9
5   44   20   10   7   0     6   40   270   10   8   0     7   38   200   10   6   0     9   36   120   40   5   10     10   0   -1   10   0   9     11   38   90   40   6   9     12   40   120   30   6   0     13   34   50   30   6   0     14   12   10   80   5   9     15   50   340   30   6   8     18   30   20   60   4   9     17   42   230   30   6   9     20   38   180   30   6   2     21   50   20   90   4   8     22   46   180   40   5   9	4	0	-1	20	0	9
6   40   270   10   8   0     7   38   200   10   6   0     9   36   120   40   5   10     10   0   -1   10   0   9     11   38   90   40   6   9     12   40   120   30   6   0     13   34   50   30   6   0     14   12   10   80   5   9     15   50   340   30   7   0     16   30   20   60   4   9     17   42   230   30   6   8     19   0   -1   50   5   9     20   38   180   30   6   0     21   50   20   90   4   8     22   46   180   40   5   0	5	44	20	10	7	0
7 38 200 10 6 0   8 52 40 70 6 0   9 36 120 40 5 10   10 0 -1 10 0 9   11 38 90 40 6 9   12 40 120 30 6 0   13 34 50 30 6 0   14 12 10 80 5 9   15 50 340 30 7 0   16 30 20 60 4 9   17 42 230 30 6 9   19 0 -1 50 5 9   20 38 180 30 6 0   21 50 20 90 4 8   22 46 180 40 5 0   23 46 30 30 6 2   24	6	40	270	10	8	0
8   52   40   70   6   0     9   36   120   40   5   10     10   0   -1   10   0   9     11   38   90   40   6   9     12   40   120   30   6   0     13   34   50   30   6   0     14   12   10   80   5   9     15   50   340   30   7   0     16   30   20   60   4   9     17   42   230   30   6   8     18   30   340   80   6   9     19   0   -1   50   5   9     20   38   180   40   5   0     23   46   320   50   7   0     24   40   30   30   6   2	7	38	200	10	6	0
9   36   120   40   5   10     10   0   -1   10   0   9     11   38   90   40   6   9     12   40   120   30   6   0     13   34   50   30   6   0     14   12   10   80   5   9     15   50   340   30   7   0     16   30   20   60   4   9     17   42   230   30   6   8     18   30   340   80   6   9     19   0   -1   50   5   9     20   38   180   30   6   0     21   60   20   90   4   8     22   46   180   40   5   0     23   46   320   50   7 <td< td=""><td>8</td><td>52</td><td>40</td><td>70</td><td>6</td><td>0</td></td<>	8	52	40	70	6	0
10   0   -1   10   0   9     11   38   90   40   6   9     12   40   120   30   6   0     13   34   50   30   6   0     14   12   10   80   5   9     15   50   340   30   7   0     16   30   20   60   4   9     17   42   230   30   6   8     18   30   340   80   6   9     19   0   -1   50   5   9     20   38   180   30   6   0     21   50   20   90   4   8     22   46   180   40   5   0     23   42   30   30   6   2   2     25   46   330   10	9	36	120	40	5	10
11 38 90 40 6 9   12 40 120 30 6 0   13 34 50 30 6 0   14 12 10 80 5 9   15 50 340 30 7 0   16 30 20 60 4 9   17 42 230 30 6 8   18 30 340 80 6 9   19 0 -1 50 5 9   20 38 180 30 6 0   21 50 20 90 4 8   22 46 180 40 5 0   23 46 320 50 7 0   24 40 30 30 6 2   25 46 330 10 8 0   28 28 10 10 8 0   29	10	0	-1	10	0	9
12   40   120   30   6   0     13   34   50   30   6   0     14   12   10   80   5   9     15   50   340   30   7   0     16   30   20   60   4   9     17   42   230   30   6   8     18   30   340   80   6   9     19   0   -1   50   5   9     20   38   180   30   6   0     21   50   20   90   4   8     22   46   180   40   5   0     23   46   320   50   7   0     24   40   30   30   6   2   2     25   46   330   10   8   0   30   37   9   38   30 <t< td=""><td>11</td><td>38</td><td>90</td><td>40</td><td>6</td><td>9</td></t<>	11	38	90	40	6	9
13 34 50 30 6 0   14 12 10 80 5 9   15 50 340 30 7 0   16 30 20 60 4 9   17 42 230 30 6 8   18 30 340 80 6 9   20 38 180 30 6 0   21 50 20 90 4 8   22 46 180 40 5 0   23 46 320 50 7 0   24 40 30 30 6 2   25 46 320 50 7 9   28 20 30 60 6 8   29 38 180 10 8 0   30 34 160 50 6 8   31 38 20 20 6 0   32 <td>12</td> <td>40</td> <td>120</td> <td>30</td> <td>6</td> <td>0</td>	12	40	120	30	6	0
14   12   10   80   5   9     15   50   340   30   7   0     16   30   20   60   4   9     17   42   230   30   6   8     18   30   340   80   6   9     19   0   -1   50   5   9     20   38   180   30   6   0     21   50   20   90   4   8     22   46   180   40   5   0     23   46   320   50   7   0     24   40   30   30   6   2     25   46   330   10   8   0     26   0   -1   80   5   2     27   34   200   30   7   9     28   20   30   60   6 <td< td=""><td>13</td><td>34</td><td>50</td><td></td><td>6</td><td>0</td></td<>	13	34	50		6	0
15503403070163020604917422303068183034080669190-150592038180306021502090482246180405023463205070244030306225463301080260-18052273420030792820306068293818010803034160506831382202060320-140593334260407934282404059353810050563630607048420-1100844801001047453426030594620206058470-110094830300305	14	12	10	80	5	9
16302060491742230306818303408069190-150592038180306021502090482246180405023463205070244030306225463301080260-18052273420030792820306068293818010803034160506831382202060320-1405933342604079342824040593538100505636306060612374040507840400270507844801001047453426030594620206058470-11009483030030 <td< td=""><td>15</td><td>50</td><td>340</td><td>30</td><td>7</td><td>0</td></td<>	15	50	340	30	7	0
1742230306818303408069190-150592038180306021502090482246180405023463205070244030306225463301080260-18052273420030792820306068293818010803034160506831382202060320-140593334260407934282404059353810050563630606061237404050784040290307541423607048420-1105943103607008448010010474534260305946202060 <td< td=""><td>16</td><td>30</td><td>20</td><td>60</td><td>4</td><td>9</td></td<>	16	30	20	60	4	9
18 30 340 80 6 9   19 0 -1 50 5 9   20 38 180 30 6 0   21 50 20 90 4 8   22 46 180 40 5 0   23 46 320 50 7 0   24 40 30 30 6 2   25 46 330 10 8 0   26 0 -1 80 5 2   27 34 200 30 7 9   28 20 30 60 6 8   30 34 160 50 6 8   31 38 220 20 6 0   32 0 -1 40 5 9   33 34 280 40 7 9   34 28 240 40 5 9   35	17	42	230	30	6	8
190-150592038180306021502090482246180405023463205070244030306225463301080260-18052273420030792820306068293818010803034160506831382202060320-14059333426040793428240405935381005056363060606123740405068380-10084040290307541423607048420-1100944801001047453426030594620206058470-110094830300305 <td>18</td> <td>30</td> <td>340</td> <td>80</td> <td>6</td> <td>9</td>	18	30	340	80	6	9
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	19	0	-1	50	5	9
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	20	38	180	30	6	0
22 $46$ $180$ $40$ $5$ $0$ $23$ $46$ $320$ $50$ $7$ $0$ $24$ $40$ $30$ $30$ $6$ $2$ $25$ $46$ $330$ $10$ $8$ $0$ $26$ $0$ $-1$ $80$ $5$ $2$ $27$ $34$ $200$ $30$ $7$ $9$ $28$ $20$ $30$ $60$ $6$ $8$ $29$ $38$ $180$ $10$ $8$ $0$ $30$ $34$ $160$ $50$ $6$ $8$ $31$ $38$ $220$ $20$ $6$ $0$ $32$ $0$ $-1$ $40$ $5$ $9$ $33$ $34$ $260$ $40$ $7$ $9$ $34$ $28$ $240$ $40$ $5$ $9$ $35$ $38$ $100$ $50$ $5$ $6$ $36$ $30$ $60$ $60$ $6$ $12$ $37$ $40$ $40$ $50$ $6$ $8$ $38$ $0$ $-1$ $0$ $0$ $8$ $39$ $40$ $270$ $50$ $7$ $8$ $40$ $40$ $290$ $30$ $7$ $5$ $41$ $42$ $360$ $70$ $4$ $8$ $42$ $0$ $-1$ $10$ $5$ $9$ $43$ $10$ $360$ $70$ $4$ $8$ $44$ $80$ $100$ $10$ $4$ $7$ $44$ $80$ $100$ $10$ $4$ <t< td=""><td>21</td><td>50</td><td>20</td><td>90</td><td>4</td><td>8</td></t<>	21	50	20	90	4	8
23463205070244030306225463301080260-18052273420030792820306068293818010803034160506831382202060320-14059333426040793428240405935381005056363060606123740405068380-1008394027050784040290307541423607048420-110094310300305044801001047453426030594620206058470-1700751283405079520-1200105312360605<	22	46	180	40	5	0
24 40 30 30 6 2   25 46 330 10 8 0   26 0 -1 80 5 2   27 34 200 30 7 9   28 20 30 60 6 8   29 38 180 10 8 0   30 34 160 50 6 8   31 38 220 20 6 0   32 0 -1 40 5 9   33 34 260 40 7 9   34 28 240 40 5 9   35 38 100 50 5 6   36 30 60 60 6 12   37 40 40 50 6 8   38 0 -1 0 0 8   39 40 270 50 7 8   40	23	46	320	50	7	Ō
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	24	40	30	30	6	2
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	25	46	330	10	8	0
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	26	0	-1	80	5	2
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	27	34	200	30	7	9
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	28	20	30	60	6	8
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	29	38	180	10	8	0
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	30	34	160	50	6	8
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	31	38	220	20	6	0
33 $34$ $260$ $40$ $7$ $9$ $34$ $28$ $240$ $40$ $5$ $9$ $35$ $38$ $100$ $50$ $5$ $6$ $36$ $30$ $60$ $60$ $6$ $12$ $37$ $40$ $40$ $50$ $6$ $8$ $38$ $0$ $-1$ $0$ $0$ $8$ $39$ $40$ $270$ $50$ $7$ $8$ $40$ $40$ $290$ $30$ $7$ $5$ $41$ $42$ $360$ $70$ $4$ $8$ $42$ $0$ $-1$ $10$ $5$ $9$ $43$ $10$ $360$ $70$ $0$ $8$ $44$ $80$ $100$ $10$ $4$ $7$ $45$ $34$ $260$ $30$ $5$ $9$ $46$ $20$ $20$ $60$ $5$ $8$ $47$ $0$ $-1$ $10$ $0$ $9$ $48$ $30$ $300$ $30$ $5$ $0$ $49$ $38$ $10$ $40$ $7$ $0$ $50$ $0$ $-1$ $70$ $0$ $7$ $51$ $28$ $340$ $50$ $7$ $9$ $52$ $0$ $-1$ $20$ $0$ $10$ $53$ $12$ $360$ $60$ $5$ $1$ $54$ $26$ $300$ $70$ $6$ $8$ $55$ $0$ $-1$ $20$ $0$ $9$	32	0	-1	40	5	9
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	33	34	260	40	7	9
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	34	28	240	40	5	9
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	35	38	100	50	5	6
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	36	30	60	60	6	12
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	37	40	40	50	6	8
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	38	0	-1	0	0	8
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	39	40	270	50	7	8
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	40	40	290	30	7	5
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	41	42	360	70	4	8
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	42	0	-1	10	5	9
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	43	10	360	70	0	8
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	44	80	100	10	4	7
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	45	34	260		5	9
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	46	20	20	60	5	8
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	47	0	-1	10	0	9
49   38   10   40   7   0     50   0   -1   70   0   7     51   28   340   50   7   9     52   0   -1   20   0   10     53   12   360   60   5   1     54   26   300   70   6   8     55   0   -1   20   0   9	48	30	300	30	5	0
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	49	38	10	40	7	0
51   28   340   50   7   9     52   0   -1   20   0   10     53   12   360   60   5   1     54   26   300   70   6   8     55   0   -1   20   0   9	50	0	-1	70	0	7
52   0   -1   20   0   10     53   12   360   60   5   1     54   26   300   70   6   8     55   0   -1   20   0   9	51	28	340	50	7	9
53   12   360   60   5   1     54   26   300   70   6   8     55   0   -1   20   0   9	52	0	-1	20	0	10
54   26   300   70   6   8     55   0   -1   20   0   9	53	12	360	60	5	1
55 0 -1 20 0 9	54	26	300	70	6	8
	55	0	-1	20	0	9

**Table 7.2:** Slope angle, aspect and vegetation cover values for each of the 55 sites sampled along with the erosion and sodicity meter values (N.B. -1 is flat).

# 7.3.2 The Relationship between the Laboratory Analysis and the Field Sodicity Meter Readings

An important additional aim of the field and laboratory investigations was to determine the relationship between the field sodicity meter and actual measured sodicity parameters. The relationship between ESP values and the sodicity meter measurements can be seen in Figure 7.6, and it is evident that there is no discernable positive correlation between the two parameters (quite the opposite). Soils with high ESP values, in many cases, have been classified as 'non-sodic' with the field meter, and soils with low ESPs have in some circumstances been classed as 'sodic' and 'highly sodic'. However, a strong relationship between the two variables is not necessarily expected, as ESP is not individually responsible for the deflocculation of soils as other factors will also exert influences. Nonetheless, it would be expected that some general trend would exist if indeed the field sodicity meter could correctly classify varying levels of soil sodicity. A relationship between ESP and the sodicity meter (Figure 7.6) where r = -0.61, suggests that the meter is actually working inversely, thus higher ESP values are being classified in the lower sodicity ranges.

Using relationships identified by commentators in the literature, such as Gerber and Harmse (1987) and Rengasamy *et al.* (1984) to identify the dispersive nature of any given soil, provides a useful baseline from which the sodicity meter can be compared. The correlation between the field sodicity meter and soil dispersivity based on Gerber and Harmse's (1987) deflocculation categories can be seen in Figure 7.7. It is evident that no relationship exists between the field and laboratory methods, and in fact a negative correlation can be discerned where an increase in dispersivity class is associated with a decrease in sodicity.



Figure 7.6: The relationship between ESP and the field sodicity meter.



Figure 7.7: The relationship between the Gerber and Harmse (1987) dispersivity index and the field sodicity meter. Where: 0 is completely non-dispersive; 1 is non-dispersive; 2 is moderately dispersive; 3 is dispersive and 4 is highly dispersive (NB. 4 samples have not been plotted as their ESP was in excess of 20 percent and thus exceeded the range of the graph).

The sodicity meter works on the simple premise that clay minerals present in sodic soils swell on contact with water, thus causing the deflocculation of the soil. Therefore, measuring the turbidity of the water after a soil sample has been suspended, and then left to stand for a given period of time, should provide a general indication of the level of sodicity. However, when a soil sample is mixed with water, the settling rates of different particles will vary. For instance, when using the hydrometer method to determine particle size analysis larger particles will settle first. Sand and silt will settle prior to clay particles which may take in excess of 8 hours. Hence, the method used to determine the level of sodicity for a given sample using the field meter may be responsible for the apparent poor contradictions.

Outlined in section 7.2.1, the field sodicity meter method requires the sample to stand for 4 hours before the measurement is taken. Clay particles therefore would still be suspended after only 4 hours in soils containing only very small amounts of exchangeable sodium where the soil aggregates are still flocculated, creating a large degree of error associated with the meter. In order to counter this, the method instructs that the water is gently poured down the side of the jar and should be inverted gently once only. This in theory minimises the mechanical breakdown of the soil aggregates and only spontaneously dispersible clays should deflocculate. However, considerable discrepancies may occur as some flocculated clays will still be suspended even with minimal shaking and there is also the risk that the sample is not sufficiently mixed creating the potential scenario whereby clays that may disperse do not do so as they have not made contact with the water.

#### 7.4 SUMMARY

The swelling and dispersion of sodic soils is an important process in the understanding of soil erosion and degradation in the Almería province, Southeast Spain. An extensive array of soil physico-chemical parameters have been documented as playing an important role in the deflocculation process. Many of the parameters which have been determined in this study, using a range of samples, provide a useful insight into the erosional processes operating. A number of relationships, largely based on those identified in the extensive literature regarding sodic soils and their associated problems and implications, have been explored within this chapter. Furthermore, the 55 samples analysed have been classified using two well-known and reliable dispersivity classification domains to determine the level of agreement between them and the field observations. Finally, the extent of relation between the field sodicity meter and the laboratory analysis has been tested in order to determine its usefulness and applicability for future research.

A number of weak relationships were identified between the different variables, such as ESP and pH, where a general increase in pH was associated with an increase in the ESP. ESP and SAR produced an  $r^2$  value when correlated against EC of 0.56 and 0.54 respectively, and it must be noted here that SAR was determined from ESP, thus similar correlations will occur. The relationships identified in this study are not as strong as some that have been documented in the literature, such as Faulkner *et al.* (2000), however, there are some general trends and the results have enabled a valuable insight into the chemistry of the soils. Furthermore, it has provided a reliable baseline from which the field sodicity meter can be accurately assessed.

It is evident that the majority of the samples analysed here plot in the completely nondispersive domain identified by Gerber and Harmse (1987) (Figure 7.2). However, some of those plotted within that domain were seen to be eroding through subsurface processes, indicating some disagreement. Using the domains set by Rengasamy et al. (1984), all samples plot within the 'dispersive' or 'potentially dispersive' classes. The variability between different classification domains is therefore highlighted and indicates the complexity involved in classifying levels of dispersivity and the subsequent caution that should be exercised when using them. Such classification schemes fail to take into consideration all associated physical or chemical variables. and they do not account for field conditions such as slope angle, vegetation cover or landscape morphology. Such parameters have been shown to influence the susceptibility of a soil to the swelling and deflocculation processes and therefore require consideration when attempting to understand subsurface or pipe erosion, and at best such parameters should only be used as broad guidelines. Although ESP and SAR levels do not appear to be extremely high it is clearly evident that subsurface processes are operating extensively. Faulkner et al. (2000) and Faulkner et al. (2003b) suggested that the soils in and around Mocatán, developed in the Gochar Formation (see section 2.4) are low in clay content. As a consequence the deflocculation of clays may not cause a reduction in hydraulic conductivity as would normally occur when clays slake and pore spaces become clogged (Naidu et al., 1995; Irvine and Reid, 2001). Rather, infiltration rates continue unaffected and the clay fraction becomes dispersed leading to the extensive development of subsurface pipes. The apparent susceptibility of the material to erode is further enhanced as a result of the landscape morphology, the associated hydraulic gradients in the region, and to a lesser extent,

the degree of vegetation cover, particularly in locations such as Mocatán documented previously by Faulkner *et al.* (2000) (see section 2.4.5).

An objective of this research was to determine and quantify the relationship between the laboratory analysis of the samples and the field sodicity meter. The results presented in this chapter have highlighted the poor relationship identified between the field meter and ESP as well as with the classification scheme of Gerber and Harmse (1987). The results suggest that the methodology by which the meter works may indeed be responsible, at least in part, for the poor correlations with standard sodicity methods. As suggested previously the method may overestimate sodicity levels as flocculated clays may still be suspended after 4 hours, and it may underestimate sodicity levels as the sample may not be fully mixed with water thus inhibiting the dispersion process.

The use of such field meters should therefore be undertaken with appropriate caution and awareness of the associated implications, especially those associated with the methodology. In such cases where field analysis is required, it may be better to use an ion-specific electrode field meter which can accurately predict exchangeable sodium (Irvine and Reid, 2001), in combination with electrical field EC and pH meters that can accurately, quickly and easily measure such parameters.

Furthermore, it is important to be aware of potential limitations and associated problems incorporated within the field investigations undertaken within this study. Only a relatively small number of samples have been analysed, all of which have been collected and sampled at the same depth, 20cm. It would have been useful to collect
surface samples and compare and contrast them to the subsurface ones, as Faulkner *et al.* (2000) identified significant variations between the two using SAR. In addition, Faulkner *et al.* (2000) and Alexander *et al.* (1999) identified site signatures for three different badland locations in Almería, and Faulkner *et al.* (2003a) highlighted significant variations in sample geochemistry within a single gully in Vera. However, this study has incorporated all samples together and it would therefore be expected to decrease the level of correlation identified within the aforementioned studies.

## 7.5 CONCLUSIONS

This chapter has fulfilled the objectives that were set out initially, exploring the physico-chemical relationships between various parameters and providing a critique of the sodicity meter. It has provided a valuable insight into the operative erosional processes and the caution that should be exercised when using field meters as a means of assessing levels of sodicity. The following chapter provides a comprehensive discussion of all of the results obtained throughout this study with reference to the wider literature where appropriate.

# 8

# Discussion

#### **8.1 INTRODUCTION**

This chapter provides an in-depth discussion relating to the results provided from the soil erosion classifications in Chapter Six and the Field Investigations in Chapter Seven. The chapter discusses the results and makes particular reference to the stated aims and objectives set out in Chapter One.

At the conclusion of this chapter the advantages and disadvantages of employing the two AI techniques for the mapping of soil erosion processes will be better understood. The chapter will explore the influence of the dependent and independent variables (sections 8.2 and 8.3 respectively), the ability of these classifiers to further current understanding of soil erosion processes and identify possible rules or thresholds of erosion processes (section 8.4). This is followed by an assessment of the erosion classification accuracies (section 8.5) and a discussion of the overall performance of the classifications (section 8.6). Finally, a summary is provided in section 8.7. This will allow conclusions to be drawn in Chapter Nine and highlight opportunities for further research.

## **8.2 THE DEPENDENT VARIABLE**

The dependent variable throughout is 'erosion'. A number of issues relating to the dependent variable are worthy of discussion. The following sub-sections describe; the effects associated with the number of classes within the dependent variable (section

8.2.1), the influence of training set composition (section 8.2.2), and the interpretation of the classified soil erosion maps in section 8.2.3.

# 8.2.1 The Effects of the Number of Classes Incorporated into the Dependent Variable

The number of erosion classes incorporated within a classification affects classifier performance, and exerts influences over a range of further issues that will in turn be discussed in the following sub-sections. The varying degrees of complexity are outlined in Figure 5.9, indicating the amalgamation of the different classes culminating in the dependent variables for the two, three and nine class classifications. As expected, the overall accuracy achieved through the various classification techniques varies according to the complexity of the problem.

Overall accuracies reduce, as the classification problem becomes more complex. The extent to which this occurs is important, particularly in relation to practical applications. It may not be important to some practitioners that the classifier has not distinguished between every type of erosion and as long as it has distinguished between certain critical classes the 'product' may satisfactorily serve their particular purpose. For example, it may be important for one practitioner that the classifier can distinguish between all nine classes, as they require a very detailed output. However, another user may simply require knowledge relating to areas that are currently eroding and not eroding. Therefore, it is important to determine the extent to which classifiers can accurately predict unknown test cases, and the associated misclassifications.

Tables 8.1 and 8.2 provide accuracy results derived from the various classification techniques, based upon the three class and nine class classifications respectively. The

#### Discussion

results indicate that, particularly in the three class DTCs, the accuracies obtained actually exceed those produced specifically for the two class problem (Table 8.1). Thus, from a practical point of view, it may appear logical to produce a three class classification and simply degrade the classification if the practitioner simply requires knowledge relating to the presence or absence of erosion. However, such an assumption may incorporate significant implications.

	ANNs	DTCs	DA
10 Metre DEM	56.2%	66.9%	65.4%
20 Metre DEM	66.9%	68.5%	68.5%
Field Variables	65.4%	77.7%	75.4%
Field Variables and 10 metre DEM	65.4%	76.9%	76.2%
Field Variables and 20 metre DEM	71.5%	75.4%	74.6%
Field Variables and Classified Vegetation	69.2%	75.4%	65.4%
Field Variables, 10 metre DEM and classified Vegetation	68.5%	75.4%	76.2%
Field Variables, 20 metre DEM and classified Vegetation	63.1%	75.4%	74.6%

**Table 8.1:** Accuracy of a two class classification interpreted from the three class correlation matrices. (Values in blue indicate a superior overall accuracy to those produced from the corresponding two class classifications).

COMPANY AND	ANNs		DTCs		DA	
	2 Class	3 Class	2 Class	3 Class	2 Class	3 Class
10 Metre DEM	34.6%	34.6%	67.7%	46.2%	67.7%	41.5%
20 Metre DEM	39.2%	36.9%	65.4%	41.5%	70%	48.5%
Field Variables	71.5%	61.5%	73.8%	55.4%	73.8%	51.5%
Field Variables and 10 metre DEM	70%	57.7%	73.1%	53.8%	73.8%	53.1%
Field Variables and 20 metre DEM	69.2%	59.2%	73.8%	57.7%	73.1%	52.3%
Field Variables and Classified Vegetation	67.7%	58.5%	78.5%	57.7%	76.9%	56.2%
Field Variables, 10 metre DEM and classified Vegetation	70.8%	56.9%	73.1%	53.8%	73.8%	51.5%
Field Variables, 20 metre DEM and classified Vegetation	57.7%	53.8%	69.2%	55.4%	73.1%	51.5%

**Table 8.2:** Accuracy of a two and three class classification interpreted from the nine class correlation matrices. (Values in blue indicate a superior overall accuracy to those produced from the corresponding two or three class classifications, and those in red indicate the same overall accuracy).

Similar scenarios occur when interpreting the correlation matrices produced for the nine class classifications. Table 8.2 details the accuracies derived through the interpretation of the nine class classifications for each of the three techniques applied to the two and three class problems. The table highlights a number of cases whereby a classifier produced for the most complex classification outperforms, or reaches the same level of accuracy, as those produced specifically for either the two or three class problem. Interestingly however, only the neural networks have produced a classifier trained for the nine class classification which actually distinguishes between 'no appreciable erosion', 'rill erosion' and 'gully erosion' better than does its three class classification counterpart. The likely cause for this is a swamping effect during the training procedure, that produces in erroneous misclassifications into one or two classes, resulting in significant commission errors. This problem is discussed fully sub-section 8.2.2.

It is likely to be the case that in the decision trees and discriminant analysis the training undertaken for a nine class classification causes the classifier to identify subtle underlying relationships required to distinguish between classes. In the less complex classification problems the DTCs and DA can generalise more by ignoring or overlooking less obvious patterns or trends located in the training data. However, this does not necessarily imply that by training with an increased number of classes in the dependent variable, an improved classifier will be identified for a less complex task. Due to the increased complexities involved, the classifier has to search more extensively for rules or parameters and, in the case of the decision trees, results in more splits (nodes) which subsequently increases the size of the tree. Consequently, this will lead to higher error rates and decreased generalisation abilities beyond the

training data. Therefore, although it may appear cost-effective and time-efficient to produce a more detailed classification and simply decrease the associated complexities as required, it is likely that higher error rates will be encountered, combined with a reduction in the applicability of the classifier.

#### 8.2.2 The Influence of Training Data Set Composition on Classifier Performance

It would appear that the composition of the training data exerts considerable influence over a number of the trained classifiers and lead to potentially poor and often misleading results. The issue is related to inter-class variability; the difference in the number of training examples presented to the classifier for each class of the dependent variable.

Table 8.3 details the number of examples comprised within the training set for the two, three and nine class classifications. It is evident that the training data possesses inter-class variabilities within each of the different classifications, particularly in the more detailed nine class problem. Unfortunately, this is to a large extent unavoidable due to the unknown distribution of erosion processes and the opportunistic method with which sampling was undertaken. This significantly reduces the ability to control the composition of the training data as the erosion processes occurring at the sites are unknown prior to data recording.

In the case of the two class classifications, where the model attempts to distinguish between cells that are either showing 'no appreciable erosion' or 'erosion', the training data consisted of 152 and 238 cases respectively. With the exception of the two ANNs trained with the DEM data, the reasonable large number of cases representing each class resulted in the development of classifiers that did not appear to suffer as a consequence of training set composition. The classified erosion maps produced from these classifications can be seen in Figures 6.2 (10 metre) and A6.3 (20 metre), show a strong bias towards the erosion class which is particularly evident in the latter where the vast majority of the cells within the map are classified as eroding. However, it is worthwhile highlighting the fact that the networks used to derive the maps produced poor overall accuracies when compared to other training sets, higher error rates and worst ROC curves.

Class	2 Class		3 Class		9 Class	
	Cases	Percentage	Cases	Percentage	Cases	Percentage
0	152	39 %	152	39 %	152	39 %
1					40	10.3 %
2			68	17.4 %	16	4.1 %
3		1			12	3.1 %
4	238	61 %			22	5.6 %
5					54	13.8 %
6			170	43.6%	59	15.1 %
7	7			[	27	6.9 %
8					8	2.1 %

Table 8.3: Composition of the training data set.

A similar problem is evident throughout all eight neural networks trained for the three class classification as the two most prominent classes in the training data, namely 'no appreciable erosion' and 'gully erosion', comprise a total of 82.6 percent of the entire data set. This appears to generally encourage the allocation of unknown cases presented to the networks into one of these two classes, at the expense of the 'rill erosion' class. Evidence of this can be seen in the error matrices constructed for the classifications and is supported further by the erosion maps produced from the 10 and 20 metre DEM training data (Figures 6.3 and A6.7 respectively). The matrices possess extensive commission errors as a number of unknown cases of 'rill erosion'

have been attributed to either the 'no appreciable erosion' class or the 'gully erosion' class - a problem that not only leads to poor overall prediction and performance but also to misleading and ambiguous results. However, this problem does not appear to be replicated in either the decision trees or the discriminant analysis classifications.

As highlighted by Table 8.3, the inter-class variability is most extreme in the training data used within the nine class classifications, where the training data composition varies from 2.1 percent in the case of 'extreme subsurface gully erosion' up to 39 percent in the 'no appreciable erosion class'. Such variation appeared to negatively influence the neural networks ability to classify unknown examples into classes that were less well represented in the training data. This is highlighted in the correlation matrices. The problem was evident in the respective decision trees and discriminant analysis, yet it did not occur to the same degree. The matrices confirm that in the majority of the classifications undertaken using the decision trees or discriminant analysis, unknown cases in the test data set were attributed to every class, and only in a small number of cases were classes missed or ignored. Moreover, in some cases where this occurred it was in fact the most prominent class 'no appreciable erosion' that was missed, indicating perhaps the difficulty in distinguishing between classes, rather than training set composition.

Some authors have expressed the problem of training set composition and the issues involved particularly when undertaking any classification involving environmental phenomenon. Foody *et al.* (1995a, 1995b) and Foody and Mathur (2004) highlighted the potential effects associated with inter-class variability's, and made the point that if a class were more abundant in the training stage, it generally encouraged the allocation to that class incurring significant commission errors, supporting the findings outlined here. Furthermore, Tourenq *et al.* (1999) and Manel *et al.* (1999) found that when the composition of the training set was uneven in a simple binary classification, neural networks struggled to reach acceptable levels of accuracy and delivered better predictions for the largest occurring class. This can lead to substantial commission errors and subsequent misleading results. For example, within this study the ANNs produced superior overall accuracies for the nine class classifications compared to both the DTCs and DA. However, the networks tended to simply assign all unknown test cases into one or two classes, and are in fact particularly poor classifiers.

The difficulties associated with the production and development of training sets are numerous and are largely concerned with issues relating to cost. When studying physical and environmental phenomena, the problems are further increased as the variable of interest does not generally occur evenly and its presence or absence is often unknown. A similar problem has been discussed by Welsh *et al.* (1996) when modelling the abundance of rare species in south-eastern Australia. If the training data collected here were to be evened out in an attempt to counter such problems information would be lost by reducing the number of cases associated with the largest class in the training set (Tourenq *et al.*, 1999). Furthermore, Ellis (1997) found that creating proportionally selected training sets simply led to the over-prediction of smaller classes and further reduced overall accuracies. Nonetheless, the problems associated with misleading results, particularly in the case of the ANNs, can have serious negative implications. The misclassification of severely eroding cells in the

field into the class of 'no appreciable erosion', could have detrimental consequences if such a method were employed as a land management tool.

Taking the discussion into account regarding the composition of the training set, it is important to consider the means by which the classifiers are tested or at least be aware of the limitations involved. For example, it is possible to propose that a classifier (ANN, DTC or DA) is capable of attributing unknown cases to each of the dependent variable classes but has simply not received such examples within the test set. This leads to the conclusion that inter-class variability within training data has significant detrimental affects upon classifier performance. Nevertheless, the erosion maps developed for the classifications undertaken using the neural networks provide further evidence to support the findings outlined in the matrices. The 10 and 20 metre DEMs contain a significantly increased number of cells, compared with the test set. Therefore, if a classifier possessed the ability to allocate unknown cases to a particular class, it is likely to identify them within such data sets. However, through visual inspection of the maps and the percentage of cells classified for each class for the three and nine class classifications, detailed in Tables 6.22 and 6.23 respectively, the networks do not appear to possess the ability to classify the less prominent classes. Nonetheless, it must not be forgotten that the networks trained using the DEM data sets produced the lowest overall accuracies, possessed the highest error rates and occurred in combination with the lowest AUC statistic for their respective ROC curve.

In summary, the training set composition exerts a significant influence over ANN performance, leading to the allocation of unknown cases into the most affluent classes

within the data set. This results in significant commission errors and misleading overall performance statistics. However, the influence of training set composition can be assessed through simple visual inspection of decision tree topology. Once a tree has been grown, its ability or inability to allocate unknown cases into various classes is determined by the explicit tree structure which allows the user to highlight any potential weakness and, if desired, an appropriate alternative tree may be selected prior to any classification.

#### **8.2.3 Interpretation of Classified Soil Erosion Maps**

Definitions of soil erosion risk, hazard and potential have been provided in Chapter Three (section 3.6). The classifications undertaken in this thesis incorporate the dependent variable 'erosion' that has been derived through ground surveys where current or actual erosion processes were recorded. The maps produced from the various classifications attempt to replicate the processes operating on the ground, and are considered to be erosion maps. As stated in this thesis, the 'acceptable' level of overall accuracy varies depending upon the practitioners requirements; be that the individual landowner (farmer) at the micro-scale, to local, and national government at the meso and macro-scales. The level of acceptable accuracy is likely to change significantly between these users and determining this level is highly subjective with few guidelines detailing such parameters.

Taking these issues into consideration, the soil erosion maps developed and produced through the various classifications may be better interpreted as maps of 'potential' erosion, rather than 'actual' erosion. The maps could be considered to infer soil erosion risk. Actual risk and potential risk are closely related and the maps can be treated as either because the independent variables are all physical parameters (see section 3.6 for a fuller discussion).

A distinct advantage of incorporating supervised classifiers, such as ANNs and DTCs, is that they learn through the presentation of training cases, where both the independent and dependent variables are known. Therefore, when presented with an unknown case, the classifier attempts to identify similar scenarios within the training data set, based upon the independent variables, and will assign the case to the appropriate predetermined class. In reality, this means that the neural network or decision tree may assign an unknown case to the incorrect class (where the correct class would be what is actually occurring on the ground at the present moment in time or the 'truth'). This would imply that there is some level of error present within the classification. However, it does suggest that the classifier has identified some similarities between the unknown case and the data presented to it within the training stage. A similar approach was incorporated by Ermini et al. (2005), Yesilnacar and Topal (2005) and Gómez and Kavzoglu (2005), where ANNs were used to produce landslide susceptibility maps for the northern Apennines, Italy, the Hendek region, Turkey, and the Jabonosa River Basin in Venezuela respectively.

The importance associated with the development and production of soil erosion maps has been stressed throughout this thesis. Nevertheless, it is likely that maps describing soil erosion risk will provide an invaluable tool for environmental managers and policy makers, assisting in the identification of areas where intervention should be sought (Haboudane *et al.*, 2002; Shrestha *et al.*, 2004). In theory, the only method by which such maps could be validated is to adopt a "wait and see" approach (Ermini *et*  *al.*, 2005). Unlike conventional approaches to risk or susceptibility mapping, the approach used here is based upon the classifier's 'experience' or 'recognition' of erosion processes within the study area, as opposed to rules or laws.

#### **8.3 THE SELECTION AND INFLUENCE OF INDEPENDENT VARIABLES**

An integral aim associated with the work undertaken within the thesis is to better understand and determine the issues relating to independent variable selection and their ability to define different soil erosion processes. The costs associated with obtaining large training sets are widely documented (see Foody and Mathur, 2004; Foody, 2002; Muchoney and Strahler, 2002; Jackson and Landgrebe, 2001; Buckheim and Lillesand, 1989). Through the use of remotely obtained independent variables (e.g. DEMs), such costs can be significantly reduced. However, the issues related to cost-benefit and trade-offs are highly complex and involve a number of parameters that require careful consideration.

The primary issue to be addressed concerns the identification of the end-user, as this will determine the specification of the end product. Once this has been achieved, it is possible to state the appropriate level of accuracy that is required, the monetary funds available for the research, and relevant time scales needed for research. At this point the appropriate trade-offs can be sought, between accuracy and resources, along with the identification of the best possible approach.

The following discussion concerns; the role of the different independent variables incorporated into the various classifications, the implications of the field sodicity meter and the overall performance of the different data sets.

### 8.3.1 Independent Variables and their roles within the Classifications

Prior to the development and construction of training data sets, it was necessary to identify through the literature, independent variables that are believed to influence and determine soil erosion processes (see section 5.3.2). Independent variables were subsequently identified and collected through either field techniques or remote sources. Within the different classifications, it would be expected that different variables appear more influential or important than others, as the information they provide assists in the discrimination between various dependent variable classes.

The digital elevation model data sets comprised a total of six independent variables (see Tables 5.1 and 5.2). Within the ANN classifications the most important predicting variable for the training procedure is slope angle and this supports the view that it is in general the most influential factor in soil erosion (Faulkner *et al.*, 2003b; Nearing *et al.*, 1991), and increasing it positively increases soil erodibility (Cerdà and García-Fayos, 1997). Of the six networks trained using the DEM data only that produced for the three class classification calibrated with the 20 metre data did not rank slope angle as the most important variable and instead ranked flow length ahead of it. Flow length was particularly useful for the three class and nine class classifications, where it was important to distinguish between varying types of erosion. The sensitivity analysis therefore tends to support the wider literature whereby gully erosion is strongly controlled by the angle of the slope and the contributing area (Martínez-Casasnovas *et al.*, 2004; Desmet *et al.*, 1999).

As the contributing area is a strong controlling factor it would seem reasonable to assume that flow accumulation would also rank highly. However, it was discovered that this is not the case and flow accumulation has little influence on the networks performance. This is highlighted by the small change in verification error if flow accumulation were to be left out. It may be the case that where two independent variables offer very similar information, namely flow length and flow accumulation, one of them will be dismissed and largely ignored in favour of the other. Of the three remaining variables, aspect is most useful for classifying erosion, followed by profile and finally plan curvature.

The variable importance data produced for the decision trees trained with the DEM data indicate that flow length is one of the most important independent variable throughout the two, three and nine class classifications. However, individual tree topology provides valuable insights into the ability of different variables to differentiate between the various dependent variable classes. Within the two class classifications, slope angle provides the root node and the maximum entropy (see Chapter Four) for both the 10 and 20 metre resolutions.

When attempting to distinguish between three classes using the DEM training sets, slope angle was replaced with flow length as the root node in both trees. This would suggest that flow length could be interpreted as a surrogate of erosive potential or flow accumulation, an obvious parameter by which rill and gully erosion may be delineated. Nevertheless, slope angle was the second most important splitting variable in both trees, followed by slope aspect. Interestingly, flow accumulation provided the highest entropy within the nine class classification when a tree was grown from the 10 metre DEM data. As in the case of the neural networks, this may suggest that flow

length and accumulation provide very similar data, as they could both be seen as a surrogate for erosion potential.

The neural networks and decision trees, trained using the field collected data, generally support the wider literature with regards to the dominant factors controlling soil erosion. As in the case with the DEM data, slope angle featured highly in both the sensitivity analysis and within individual tree topologies within each of the three classifications (two, three and nine classes). This would not only support the wider literature concerning the issue, but also emphasise a point made in section 7.3.1 where slope angle appeared highly influential in subsurface processes, even in soils containing low levels of exchangeable sodium.

Geology is also highly ranked throughout the classifications, and this is likely to be attributed in part, to the relatively poor performance associated with the field sodicity meter (see Chapter Seven). The independent variable geology provides the classifier with data associated with the lithological units that are susceptible to deflocculation processes, and may provide a surrogate in place of the meter. This point is discussed in sub-section 8.2.3.

Finally, the perceived importance of the independent variables of vegetation cover and field aspect varied throughout the classifications. Throughout the sensitivity analyses the two variables were closely ranked, and their role within the decision trees was also varied. This would suggest that neither of the variables offers significantly improved data over the other, and slope aspect is perhaps a good surrogate for vegetation cover. Nevertheless, all of the neural networks trained with both variables incorporated them within the final trained network and, in some of the decision trees, both variables provided splits. This implies that each of the variables may provide subtlety improved splitting criteria over the other.

In summary, it might have been expected that different independent variables are useful for distinguishing between sites of no appreciable erosion and erosion, or indeed different types of erosion. However, the results indicate that within the neural networks all of the independent variables in each classification provide some means of determining between dependent variable classes. Moreover, the sensitivity analysis shows that the importance of individual variables changes little between the three different classification problems. A similar trend is evident using the decision trees, whereby slope angle appears to be the most important variable, containing the highest entropy in most trees, and it is not apparent that any variables are particularly useful at splitting the data when varying levels of complexity are involved.

### 8.3.2 Implications of the Field Sodicity Meter as an Independent Variable

An objective stated in the introduction to the thesis concerned the determination of the usefulness and applicability of the field sodicity meter developed by the Co-operative Research Centre for Soil and Land Management in Adelaide, Australia. As described in Chapter Seven, a range of laboratory techniques were undertaken in addition to the work carried out in the field in order to ascertain the extent of the relationship between the meter and laboratory based sodicity measurements. In addition various physico-chemical relationships were also explored.

The discussion in the previous chapter highlighted the apparent limitations associated with the sodicity meter. There is little evidence within the results to suggest that any discernible relationship exists between the results gained through the meter and those derived through more conventional laboratory techniques. A range of proposals were considered to account for the disparities. These included; particle size which may influence turbidity and settling rates, differing soil textures, and general flaws associated with the methodology (see sections 7.3.2 and 7.4). Taking these findings into account, it is important to determine the meter's role within the various neural networks and decision tree classifiers.

The results obtained through the sensitivity analysis tend to indicate that the sodicity meter is a relatively poor predicating variable of soil erosion. The sensitivity analysis generally ranks the sodicity meter as one of the least influential variables of the five field acquired predictors when the data set is used to train a network. However, it tends to outperform the DEM data when used in combination in most cases. Furthermore, the meter does appear to present at least some useful information to the network, as its removal from each of the networks that it has been used, leads to a subsequent increase in verification error based upon the sensitivity analyses.

The variable importance statistics relating to the DTCs show the meter to have a marginal influence upon classifier performance in some cases but none in others. As outlined previously, the variable importance statistics relate to every tree grown irrespective of the one chosen. Therefore, it is important to inspect the individual tree, in order to determine the ability of individual variables to assist in class separation. The sodicity meter is incorporated within the trees grown for the binary classification

in four of the six trained with data sets incorporating the variable. Interestingly, the sodicity meter appears to split the data in contradicting ways. For example, in the two class classification trained using only field acquired variables (Figure A3.3), nodes 10 and 11 use the same splitting criteria ( $\leq 0.5$  and > 0.5). In the former node, data is split in such a way that the outcome is 'no appreciable erosion' and 'erosion' respectively, as may be expected. In contrast, node 11 splits it such that low levels of dispersion (according to the mater), leads to 'erosion' and higher levels 'no appreciable erosion'.

As discussed in Chapters Two and Three, subsurface erosion is strongly controlled by various physico-chemical relationships. The independent variable geology provides a generalised surrogate for such parameters, as sensitive lithologies have been highlighted in the literature (TRU and MRU) (see Alexander *et al.*, 1996; Spivey, 1997; Faulkner *et al.*, 2000, 2003b). The majority of both the ANN and DTC classifications indicate that geology provides superior data than that provided by the sodicity meter. However, the spatial heterogeneity of soil physico-chemical properties has been highlighted in Chapter Seven and by Corwin and Lesch (2005), Corwin *et al.* (2003) and Horney (2005). Ardahanlioglu *et al.* (2003) also demonstrated that the spatial distribution of ESP, EC and pH are highly spatially variable, particularly in sodic soils (Samra *et al.*, 1988). Thus, geology is not an ideal surrogate for an accurate sodicity meter. This may provide an explanation for the apparent inability of the classifiers to distinguish between various erosion classes and in particular within the three and nine class classifications.

As highlighted in Chapter Seven, Irvine and Reid (2001) proposed the benefit of using a sodium-specific electrode meter for accurately measuring the Exchangeable Sodium Percentage (ESP) in the field. Furthermore, Nuttall *et al.* (2003) demonstrated that an ion-specific electrode could accurately predict soil sodicity through ESP. These approaches, like the sodicity meter, are site specific. Therefore, although it is possible to produce a classification using such variables, it is not possible to incorporate them within the production of a map. As such, further investigations may assess the ability of extensive mapping techniques, including remote sensing, as it offers the ability to work at extensive spatial and spectral scales. Remote sensing has been used successfully for the identification and mapping of saline soils (see Goossens and Van Ranst, 1998; Farifteh *et al.*, In Press; Tóth *et al.*, 1991), and the technique may lend itself to the identification of sodic soils.

#### **8.3.3 Overall Performance of the Various Training Sets**

The end-user and the purpose of the developed 'product' largely dictates the accuracy level required for a tool such as an erosion map. This inherently incorporates issues relating to cost and time, and the following discussion attempts to review such issues in relation to the results presented in Chapter Six.

It could be argued that three different end products have been produced through the implementation of the research methods outlined in Chapter Five. These comprise a two class, three class and nine class soil erosion map. Acceptable levels of accuracy, or error will, undoubtedly, vary between the different products. This reiterates the difficulty associated with their identification. Nevertheless, it is important to examine the performance of the classifiers trained using the various independent variables to

improve the understanding of the potential issues involved in cost-benefits and tradeoffs between the more expensive field collected data, and the less expensive DEM extracted data.

The results achieved through the AI approaches for the two class classification suggest that the digital elevation model data alone fails to provide adequate definable boundaries within the training set provided and specifically between 'no appreciable erosion' and 'erosion'. The limited ability of both the neural networks and decision trees trained with the data is readily apparent when comparing the overall accuracies with those achieved in the majority of the other classifications, where field collected variables were incorporated. This is by no means surprising and would be expected to occur to a certain degree, as the DEM data is simply a grid-based generalisation of the 'real' landscape. Nevertheless, the issue is further compounded within both the ANNs and DTCs when the classifiers trained with the DEM data not only produce the lowest overall accuracies, but do so in combination with the highest error rates and the lowest AUC values identified from the ROC curves.

These results outlined above regarding the two class classifications are largely replicated in the neural networks and decision trees trained for the more complex three class and nine class classifications. The DEMs produce the lowest overall accuracy's and the highest error rates, compared to the classifiers trained with field acquired data. The results suggest that the independent variables acquired in the field provide improved data and subsequently enhance classifier performance. However, to understand the complexities of these findings fully, an in-depth discussion is required.

Although the independent variables collected through field methods would appear to provide improved training data, the results indicate that the variables extracted from the DEM, or at least some of them, assist in the classifications. Thus, if the decision were made to use the more expensive approach - obtain data from the field - it would make economic sense to assess the potentially enhanced classifier performance that may occur if the cheaper DEM data were included. Such improvements however, appear to be case specific, and involve potential areas of concern. For example, the networks and trees trained with data from both the field and the DEM, appear in some cases, to produce higher error rates compared to those trained using only the field data. This is a somewhat surprising and unforeseen anomaly, as the two AI techniques offer the ability to remove or, at least limit, the influence of, variables that appear to reduce classifier performance (see Chapter Four). Therefore, the presentation of the DEM data to the classifiers trained with field data should, in the worst case, obtain the same level of error (no improvement). The reasons for the apparent increase in error rates are not therefore readily explicable.

The method by which neural networks determine the importance of individual independent variables within the TRAJAN software is by means of a sensitivity analysis. This technique assesses the relative contribution of each independent variable by successively testing the network with each input excluded. Therefore, if a variable is excluded and the verification error subsequently increases, it is reasonable to assume that the missing variable provides some useful information enhancing network performance, and is thus reinstated. Although the assumptions made by this technique are reasonable, it does possess certain limitations (Abrahart *et al.*, 1999) which may in part explain such findings. These may include:

- Sensitivity analysis is a post-training procedure, and can only be assessed on completion of neural network training. Therefore, although the internal weights will have been altered automatically by the network, giving less weighting to less important variables, sensitivity analysis has no effect on the overall network unless it has identified an independent variable that when removed leads to a reduction in verification error.
- Only individual variables are removed, and the influence assessed, as opposed to the extraction of two or three variables. This means that the network may overlook or miss the optimum set of training variables.
- Neural networks will reach different solutions each time they are trained, even when using identical parameters (i.e. training data, architecture, momentum and learning rate). Therefore, the optimum network produced may differ from that produced at an alternative moment in time.

In contrast to sensitivity analysis, the decision tree software, CART, grows the largest tree possible with the optimum separation and performance on the training data. However, such trees (as stressed in Chapter Four) possess little generalisation ability beyond the training data, and thus a technique of recursively pruning the tree until an optimum error point is reached is implemented. However, this method would appear to lead to problems, such as those encountered here. The initial objective is for a tree to be grown that classifies as many of the training set cases correctly, irrespective of error rates. The tree uses as many of the independent variables as required, even if the information they may offer is minimal. From this point the tree is cut-back (pruned), removing terminal nodes and sections of the tree that overfit the data. However, the resultant tree, although possessing the smallest relative cost (error) of all grown trees,

has in fact been subject to bias as a result of the initial growth process. It may therefore, include variables that have little more than a marginal influence on overall performance.

In an attempt to determine the potential benefit of incorporating detailed vegetation cover classifications as independent variables, rather than estimations made in the field, the AI techniques were trained with data sets incorporating the information derived through ground-based photography (see section 5.5). However, the results achieved through the various classifications suggest that the determination of the potential benefits is not as straightforward as might be expected. Using the neural network classification approach, the field acquired training set incorporating the classified vegetation produced higher overall accuracies, compared to the training set using the estimated vegetation cover, but did so in combination with larger rates of error. Nevertheless, the optimal networks using the classified vegetation did appear to be smaller in size, and are likely therefore to have more generalisation ability beyond the training data. This trend was seen throughout the two, three and nine class classifications, also occurring in the classifications where DEM data was used in combination. However, in a number cases where the DEM data was incorporated, the verification error increased, further highlighting the points made previously. While the sensitivity analysis statistics show a marginal increase in the influence of the classified vegetation variable over the estimated vegetation variable, it is important to be aware that the error values are not directly comparable as the overall verification error rates are network specific.

The decision trees trained using the field data and the classified vegetation however, have to be viewed with caution as a result of their composition. For example, the tree produced for the simple two class classification only incorporates the independent variable slope angle and disregards the classified vegetation, producing a marginally better tree in terms of overall accuracy than that produced using the estimated vegetation (see Figure A3.6). Furthermore, as noted with the neural networks, the trees trained with the classified vegetation and the DEM data often possess higher error rates than those that only use the field and classified vegetation data alone.

Therefore there is little evidence to suggest, using the overall accuracy's, error rates and the ROC curves where appropriate, that the incorporation of the classified vegetation as an independent variable offers significant benefits to classifier performance in comparison to the estimated vegetation. Although the inclusion of the classified vegetation has produced trees and networks that may have better generalisation abilities beyond the training data (i.e. smaller architectures), the time, and subsequent indirect costs associated with its production suggest that they may outweigh the potential benefits.

The independent variables compiled within the different training data sets have a strong influence upon the performance of the classifiers. The results tend to suggest that the independent variables extracted from the DEMs do not provide the same quality data as that incorporated within the field collected data. This is signified by the increased error rates produced from the classifications, reduced overall accuracies and the ROC curves where appropriate. Furthermore, the results indicate that the finer resolution DEM data outperforms the coarser 20 metre data. As outlined in Chapter

Five, the generation of a DEM derived through the digitisation of contours with 10 metre spacing is 20 metres. However, a point made then and emphasised further now is that the DEM is the closest link to the landscape when the field data is not used. Thus, it is vitally important to use the data set to its maximum potential and results indicate that the 10 metre DEM provides the classifiers with better independent variable information than does the 20 metre DEM. This may also reflect potential problems associated with mixed pixels. As the size of the cell increases, it is likely that more than one erosive process may be operating. Therefore, the results may not necessarily indicate that the 10 metre data provides superior data to that provided by the 20 metre DEM, but may involve issues relating to process scale and mixed pixels.

Figure 8.1 consists of seven graphs comparing the slope angle and aspect data measured in the field with that extracted from the 10 metre and 20 metre DEMs, and also comparisons between the two DEMs. The graphs clearly show that the finer resolution DEM bears a stronger resemblance to the field collected slope angle and aspect than the 20 metre resolution model. However, the relationships do not appear to be particularly strong, supporting the fact that the classifications using field data outperformed those using the elevation model data.

The comparisons highlight some important scale issues associated with DEMs, some of which have been discussed by Reuter *et al.* (2006) and Thompson *et al.* (2001). The 20 metre DEM is a generalisation of the finer resolution 10 metre DEM, and as such lower slope gradients are produced due to the smoothing effect. Decreasing the horizontal resolution of a DEM will have similar implications on all of the other extracted variables such as aspect and profile curvature, and may therefore go some way to explaining the reduced classifier performances when incorporating such data sets. This issue is clearly evident in Figure 8.1c where the 10 and 20 metre DEM slope angles are plotted against one another. The data points generally sit closer to the x-axis than the y-axis, inferring that slope angles are further underestimated in the 20 metre DEM due to the smoothing effect outlined in Chapter Five. Furthermore, it is likely that the smoothing effect also reduces the overall correlation of the slope data in Figures 8.1a and 8.1b. However, it does not appear to effect the relationships of the aspect data to the same degree, and this comes as a result of the fact that it will not influence the direction of a slope. It is therefore important to be aware of such scale issues and the implications associated with using such data sets.

The graphs detailing the relationships of the various aspect data measurements also highlight a point worthy of mention. Firstly, Figures 8.1e and 8.1f contain a significant number of outliers, in particular along the y-axis. This has occurred as a result of the fact that flat areas were assigned a slope value of -1 (no slope) and an aspect value of -1 (no aspect). Secondly, the issue associated with circular data is responsible for the two significant data clusters in the top-left and bottom-right of the graphs, particularly evident in Figure 8.1g. Due to the inherent circular nature of aspect data, it is possible that slopes that are general north facing will possess significantly different aspects. For example, a slope with an aspect of 350° is similar to a slope with an aspect of 10°, as the difference is merely 20°, however the numerical difference is actually 340. The two data clusters are instances where the 10 and 20 metre DEMs have attributed an aspect value either side of 360°. Figure 8.1d shows the elevation of the 10 metre DEM against that of the 20 metre DEM.

Discussion





**Figure 8.1:** Comparison between slope and aspect measurements collected in the field (520 points) and those extracted from the 10 and 20 metre DEMs. Slope angles and aspect are in degrees, and elevation is in metres (NB the red line is the 1:1, and the DEM against DEM plots are for the same 520 points).

Taking these issues into account, it is likely that DEM quality will influence classifier performance. However, it is important to remember the non-linear capabilities associated with the two AI techniques incorporated within this study (see Chapter Four). In particular, the issue associated with aspect data, caused by its circular nature, should not inhibit the classifiers ability to discriminate between various classes. That is to say that the AI techniques possess the ability to overcome this problem. An important consideration however, is that in order to give the classifiers the best opportunity to learn from the training set, data that resembles most closely that of the

real world should be incorporated.

In summary, the cost-benefit complexities associated with the identification of appropriate levels of accuracy will determine the independent variables to be used within the classifications. As stated previously, if the decision were made to utilise expensive field collected data, then it may make sense to also incorporate the cheaper DEM acquired data. However, this comes with associated limitations that the user has to be aware of, such as potential increases in error. As a consequence, the best approach may be to undertake some preliminary classifications in order to aid the decisions regarding the optimum cost-benefit. For example, if the DEM data were used, it may be beneficial to incorporate geology from a map, as this is likely to significantly improve classifier performance. However, it is important to determine whether an erosion map, or an erosion risk map (as they may be interpreted), is going to be used individually or in association with other sources. The maps produced here are indicative; that is to say that they provide information relating to potential areas of concern, and are by no means definitive. Therefore, it is likely that such maps are to be used in combination with other methods, techniques or data sets, to formulate and implement appropriate management strategies. It is crucial to identify the optimum outputs (benefits) in association with a range of different inputs (cost), and these are likely to be case specific.

# **8.4 THE USE OF ARTIFICIAL INTELLIGENCE TECHNIQUES FOR KNOWLEDGE GAIN AND RULE EXTRACTION**

One of the aims outlined in Chapter One was to determine the ability of artificial neural networks and decision tree classifiers to operate as inductive tools. A distinct advantage of inductive learning algorithms is their ability to generate interesting rules and parameters (Bobbin and Recknagel, 2001). This reveals underlying patterns and process and furthers current knowledge and understanding (Bui *et al.*, In Press). The following sub-sections discuss the potential of each classification technique to provide this information.

# 8.4.1 Knowledge Extraction through Decision Tree Classifiers

Due to their explicit nature, decision tree classifiers are easy to interpret in terms of knowledge gain and rule extraction. Due to the number of trees grown for the two, three and nine class classifications, it is not possible to discuss at length the various rules or data splits identified. Nevertheless, a number of rules are consistently used to split the training data, implying the existence of some general underlying patterns.

Slope angle is generally attributed as the root node containing the maximum knowledge gain. Using the field acquired data, the decision trees identify a threshold of 19 degrees from both the two and three class classifications. Examples below this threshold are assigned as 'no appreciable erosion' by the trees. The same classifications undertaken using the DEM data identifies a threshold at 14.5 degrees. These values are slightly larger than that identified by Kosmas *et al.* (2000), where slopes in excess of 13 degrees were seen to be eroding significantly, in a semi-arid environment in Greece.

A number of trees have identified a threshold for vegetation cover, again within the two and three class classifications. A number of trees split the data at 55 percent, and this would support work undertaken by López-Bermúdez *et al.* (1998) in Murcia, where vegetation covers in excess of 50 percent provided significant protection against erosion. Also, Thornes (1990) identified that erosion increased rapidly when vegetation cover was below 30 percent. However, the trees trained with the classified vegetation cover tend to split the data at 73.05 percent. The increase between the estimated and classified vegetation is likely to occur as a result of the potential over-estimation associated with the latter, highlighted in Figure 5.2.

It is important to be aware that further splits are often required before a terminal node is reached and the splitting criteria discussed above are not definitive decisions. Nevertheless, tree structures are such that it is relatively straightforward to develop and construct simple, openly interpretable diagrams such as that seen in Figure 6.30. Such methods could be employed to assign values to process dominance domains, such as that outlined by Faulkner *et al.* (2003b) (Figure 3.6). Furthermore, rules and thresholds identified through such techniques could be incorporated within simple risk schedules, where a layered GIS approach is adopted, such as Faulkner *et al.* (2003b).

# 8.4.2 Knowledge Extraction through Artificial Neural Networks

The difficulties associated with extracting data from ANNs, and their 'black-box' nature has been discussed at length in Chapter Four. However, through the use of response surfaces, it is possible to visualise the behaviour of an artificial neural network.

The network trained with the field acquired variables for the two class classifications suggest that slopes of any angle will erode if vegetation cover is sparse (below 6 percent) (Figure 8.2a). However, as vegetation cover increases, soil erosion decreases ('erosion' to 'no appreciable erosion') in a non-linear manner. Nevertheless, even with 100 percent vegetation cover, erosion occurs on slopes in excess of approximately 45 degrees. Figure 8.2c shows the response surface for the classified vegetation and slope angle, and a similar trend can be seen to exist. Slopes at low angles with no vegetation cover show 'no appreciable erosion'. A largely linear trend can be seen with an increase in vegetation cover controlling erosion as slope angle increases up until a specific point. Beyond which erosion occurs irrespective of vegetation. As with the previous example, this occurs at around 45 degrees. Finally, Figure 8.2b shows the relationship between slope angle and aspect. As might be expected, slopes with south facing aspects are more prone to erosion than those facing north (e.g. decreased relative vegetation cover).

Unfortunately, sensitivity analysis does not allow the user to fully determine and understand the role of individual variables. It also fails to acknowledge the entire network and, whilst viewing two variables, all others are held constant. Furthermore, sensitivity analysis is unable to be used when categorical variables are employed, such as geology in this work. Discussion





#### 8.4.3 Summary

An important aim associated with the work undertaken here, relates to the ability of neural networks and decision trees to enhance our current understanding of soil erosion processes in the Almería province. The ability of inductive learning techniques for knowledge discovery has been recognised in various applications, and has been proposed as a means of generating explicit knowledge in an understandable structure, that is potentially useful to a user/practitioner and also provides new and interesting concepts (Bradley *et al.*, 1998; Kusiak *et al.*, 2006). The two AI approaches adopted here provide very different means by which knowledge discovery can be fulfilled, and if the central aim of any research concerns the determination of

272

rules and thresholds, it is important to be aware of the implications associated with each approach. The knowledge extracted through such techniques is only as good as the data used to train the classifier (Bologna, 2004). If inadequate data is used to formulate rules in an inductive manner, then erroneous knowledge will be extracted which would comprise substantial error levels. Furthermore, the user must take into account the overall performance of the classifier, and be aware of potential shortcomings.

In summary, the decision tree classifiers offer the more straightforward and easily comprehensible rule extraction, through simple explicit splitting criteria. Furthermore, when non-numeric or categorical data are incorporated within the classifications, the difficulties associated with extracting knowledge from neural networks have been highlighted. Mak and Munakata (2002) provide further support for this issue; suggesting that when categorical data are involved, neural networks should be avoided, if the aim of the work is rule extraction. In addition, issues associated with training time should also be considered as should the possibility of adopting a combined approach to data mining (Mak and Munakata, 2002; Hashemi *et al.*, 1998).

# **8.5 ACCURACY ASSESSMENTS OF THE CLASSIFIED SOIL EROSION MAPS**

An important point to be raised for discussion involves to the determination of accuracy and the methods by which it is assessed. The issues related to accuracy assessment have been briefly outlined in section 5.7, where the correlation matrix was introduced. However, to understand and appreciate the potential drawbacks and limitations associated with the classifications undertaken, it is important to be fully aware of a range of important issues.

Largely as a result of the growth in remote sensing investigations, a great deal of literature has been written on the subject of accuracy assessment (see Foody, 2002; Congalton, 1991; Congalton *et al.*, 1983; Hord and Brooner, 1976; Story and Congalton, 1986; van Genderen and Lock, 1977; Ginevan, 1979; Hay, 1979; Stehman, 2000). The literature describes; issues relating to sampling strategies, acceptable levels of accuracy, errors, and sample size, all of which will affect the confidence levels which the end-user can place in the final product.

These considerations would have little impetus on accuracy assessment if time and cost issues were irrelevant. As discussed previously however, a range of cost-benefit issues apply when undertaking a study, and undoubtedly exert influences upon the accuracy assessment itself.

An important issue of concern related to the work undertaken here relates to the number of samples incorporated within the test set. A total of 130 cases were presented to each of the trained classifiers and the results detailed in correlation matrices. The number of cases attributed to each of the dependent variable classes is reduced as the classifications become more complex except in the case of the 'no appreciable erosion' class, as this was included within each classification. Table 8.4 details the number of examples comprised in the test set for each of the three classifications.

Class	2 Class	3 Class	9 Class
	Cases	Cases	Cases
0	45	45	45
1			17
2		29	5
3			7
4	85		7
5		56	18
6	7	1 5	18
7	7	1 [	9
8	]	Ι Γ	4

Table 8.4: Composition of the test data set.

The majority of the literature written regarding accuracy assessment has been undertaken with remote sensing applications in mind. This is largely a result of the fact that classifications have traditionally been carried out in few other research areas. Consequently, it is important to be aware of factors that should be taken into consideration and discussed. Within the extensive array of literature a range of general 'rules of thumb' and 'best practices' have been proposed. However, the methods by which remote sensing classifications are carried out are not necessarily replicable in other applications where classifications have been undertaken. For example, the determination of the dependent variable is often derived through sources that are remote to the subject(s) of interest, including maps and aerial photography. This allows both the identification and collection of an extensive number of training and testing examples relatively easily and in a time-efficient manner, and few complications arise regarding classes with low spatial coverages and/or limited accessibility. This not only holds true for remote sensing investigations, but any classification whereby an existing source can be used to determine the dependent and independent variables. However when this is not the case and field data collection methods have to be employed, serious time and cost issues require consideration. It is extremely expensive to spend extensive periods of time in the field and as such, cost-
benefit decisions have to be made. Once again, these issues come back to the end-user specifications and the expectation placed upon the end-product.

Van Genderen and Lock (1977) produced a table detailing the probability of obtaining no errors in samples of varying sizes, reproduced here in Table 8.5. The table is based on the binomial expansion  $(p + q)^x$ , where q = 1 - p. By using this approach it is possible to determine the range of the true errors based upon the 95 percent confidence limits. The table is divided into two by a stepped line, indicating the minimum sample size required to statistically validate an acceptable level of accuracy. The value documented above the line is the probability of obtaining an error free sample, which is low even when true errors are present in appreciable numbers. However, below the line, there is a high probability that the acquired results were achieved using a method that was relatively error free (van Genderen and Lock, 1977). For example, if an acceptable level of accuracy were identified as 80 percent, then 15 samples per class would be required to be confident of the accuracy of the output.

	X	5	10	15	20	25	30	35	40	45	50	60
q												
0.99												<u>0.5472</u>
0.95							<u>0.2146</u>	<u>0.1661</u>	<u>0.1285</u>	<u>0,0994</u>	<u>0.0769</u>	0.0461
0.90				0.2059	<u>0.1216</u>	<u>0.0718</u>	0.0424	0.0250	0.0148	0.0087	0.0052	0.0461
0.85				<u>0.0874</u>	0.0388	0.0172						
0.80	Į		<u>0.1074</u>	0.0352								
0.70	[ (	0.1681	0.0282									
0.60	<u>ا</u>	0.07 <b>78</b>										
0.50		0.0313										

**Table 8.5:** Probability of scoring no errors in various size samples from a population with a range of real error proportions, where q is the specified interpretation accuracy and  $\chi$  is the sample size (N.B. the stepped line indicates approximate 0.05 level of probability) (van Genderen and Lock, 1977).

It is important therefore, to acknowledge the limitations associated with the determination of accuracy levels derived within this study. It is also important to

highlight the differences in the methods by which test data is collected here, and that often used in remote sensing investigations, as this may influence overall performance.

A variety of errors exist and can be encountered when undertaking any classification. These can have profound effects on classification results and should be carefully considered when assessing classifier performance. Foody (2002) considered the accuracy of the dependent variable, and potential errors in the source from which it has been determined. For example, all maps possess an inherent degree of error and it is often the case that any disagreement between the map and the classified output is assumed to indicate error in the latter (Fitzgerald and Lees, 1994). Furthermore, to avoid potential confusion and indirectly improve classifier performance, sampling is often carried out in large homogeneous areas in order to avoid mixed pixels (Foody, 2002; Richards, 1996; Foody and Arora, 1996). Consequently, classifier performance is likely to be overestimated and exaggerated.

Unlike many remote sensing land cover classifications, the dependent variable within this work varies significantly over short distances; a slope showing little erosion may be adjacent to one that it heavily eroded. Moreover, all of the training and testing sites were individually visited and ground truthed, ensuring that little bias occurred in the determination of the dependent variable, other than that inherent in the classification scheme (Figure 5.3). In some instances, a similar problem to that of mixed pixels in remote sensing occurred, whereby more than one erosion process was seen to be operating on a slope. Unfortunately, this is unavoidable, and the dominant erosive process was attributed in such cases. Nevertheless, these could potentially lead to errors in the final classifications.

Classification accuracy is an important consideration, and the potential errors involved must be acknowledged. The number of test cases varies from the two, three and nine class classifications, inhibiting the confidence that can be placed in the recorded accuracies of the final outputs. However, each of the 130 test sites were visited and the errors associated with the data should therefore be minimal. Thus, the quality of the test set may be superior to that used in many remote sensing investigations, even if the number of examples is comparably small in some instances.

In summary, the determination of an acceptable level of accuracy is user-defined, and the manner in which it is recorded varies. At present, no standard method exists by which accuracy should be determined or indeed reported (Foody, 2002). It is generally agreed however, that the correlation matrix offers the most comprehensive way in which results can be presented but a range of further techniques and statistics exist. For example, ROC curves present an ideal tool to analyse the ability of a classifier to determine between two classes and requires no test data as the graph is developed through the separability of the training data.

#### 8.6 OVERALL PERFORMANCE OF THE AI APPROACH AND A TRADITIONAL CLASSIFICATION METHOD

An integral research aim proposed in Chapter One was to evaluate and compare the performance of artificial neural networks and decision trees as soil erosion classifiers. The classifiers were selected based on their perceived abilities detailed in the literature (see Chapter Four), and in particular the advantages that they hold over

more conventional classification techniques. In order to explore such advantages, a detailed research question was posed in section 5.2, proposing that the AI approaches were compared to a traditional method, Discriminant Analysis (DA).

The extensive array of classifications undertaken and results obtained tend to suggest that ANNs and DTCs possess varying abilities for the classification of soil erosion processes. It is fair to say that both techniques applied to the two class classification performed reasonably well. The best neural network classified 75.4 percent of the test set correctly and the best decision tree 77.7 percent. Furthermore, the networks and decision trees obtained AUC statistics of 88 percent and 92 percent respectively, implying a good classification for the former and an excellent classification for the latter according to Pearce and Ferrier (2000). However, the results derived through the discriminant analysis are not only comparable to those achieved through the AI techniques but obtain an overall accuracy of 78.5 percent in one of the classifications, which actually exceeds them. This would suggest therefore, that the same level of distinction between sites that are currently eroding in the field and sites that show no appreciable erosion can be achieved using a linear technique. Unfortunately however, ROC curves cannot be extracted through the DA approach limiting the extent to which comparisons can be made. Nevertheless, a point worthy of mention is the fact that two of the eight decision trees grown for the two class classification are linear. splitting the data based only on the slope angle (Figures A3.6 and A3.8).

The classification results produced for the three class classification also tend to suggest that the discriminant analysis method is superior to the ANNs and comparable to the DTCs. The neural networks appear to produce the weakest classifiers of the

three techniques. As discussed previously in section 8.2.2, this is likely to occur as a result of training set composition and occurs to a greater degree in the nine class classifications. In contrast to the neural networks, decision trees and discriminant analysis do not appear to assign unknown test cases to the most prominent classes in the training set. This leads to highly misleading results, and appears to be a serious weakness associated with ANNs.

As detailed in Chapter Four, ANNs and DTCs have often been found to outperform traditional statistical classifiers such as DA. Due to the non-linear capabilities associated with neural networks and decision trees, it would be expected that they would either perform to comparable levels as the discriminant analysis classifications or exceed them. However, the results outlined in Chapter Six show that the DA performs to comparable levels as the AI techniques and even outperforms them in some instances. It is important to bear in mind the issues previously discussed in section 8.2.2 regarding training set composition and the influence that this has upon the ANN classifications in the three and nine class problems. As discussed previously, the decision tree classifiers produced linear solutions within some of the two class classifications. In such instances the comparable performance from DA would be expected. However, this is not the case within the three or nine class classifications. Back et al. (1996) found a similar problem whereby DA occasionally produced superior results compared with neural network for predicting bankruptcy. Altman et al. (1994) found DA performed as well as an ANN approach for distinguishing between strong and weak industrial organisations in Italy.

These findings may suggest that the classification problem is linear and, non-linear classifiers over-complicate the task. However, both AI techniques possess the ability to work in a linear manner and, a linearly discernible relationship should not be a limiting factor (i.e. should be able to classify irrespective of linearities and/or non-linearities). Moreover, through the visual analysis of scatter-plots derived from the various independent variables, it is evident that the problem is largely non-linear. Therefore, it is possible to suggest that it may occur as a result of limitations associated with the training data. The data may not provide a sufficient number of examples to delineate between classes in the non-linear manner required. Furthermore, the training parameters may limit classifier performance and the orthogonal splits used by the univariate decision trees in the classifications may result in more splitting criteria as more questions are required.

A number of factors may influence the performance associated with the various classifications (ANNs, DTCs). These may include the following.

• Limitations associated with the training data.

Composition of training data appears to cause significant bias within the neural networks, leading to reduced overall accuracies in the three class classification, and spuriously increased accuracies in the nine class classifications when compared with the other two techniques. Furthermore, the training and test data sets are relatively small. However, due to temporal constraints and the extreme topographic nature of the landscape, this is largely unavoidable. Similar issues may be encountered when employing alternative practical applications (e.g. mapping).

Limitations associated with the independent variables.

The independent variables identified within the literature may not be able to adequately provide the information required by the classifiers to determine between the various classes within the dependent variable. Furthermore, errors may be present within the inputs.

• Errors present within the dependent variable.

The dependent variable may contain errors that lead to confusion within the various classifiers and consequently result in misclassifications. These errors may be inherent in the erosion classification scheme implemented and occur largely as a result of scale. When employing grid based or raster analysis, as is the case here, it is assumed that each cell or pixel is occupied by a single homogeneous class (Campbell, 2002). However, as is often the case in remote sensing investigations, mixed pixels exist, producing composite signatures particularly when high variation occurs over very short distances (Campbell, 2002; Mather, 2001). Therefore, when hard classifiers are used, such as ANNs, DTCs and DA, mixed pixels may cause confusion. Due to the scale of investigation, it is likely that more than one erosion process may be operating on any individual slope and, irrespective of the consistency with which the schedule is applied, hard classifiers such as those used here may become confused. Furthermore, due to the inherent subjectivity of the erosion scale used to determine the dependent variable in the field, errors may be incorporated into the classifications. Coupled with the previous point, error propagation will potentially limit classifier performance further. However, there is no alternative and it highlights the value of trained investigators (e.g. geomorphologists) working in close association with planners.

Training methods and parameters associated with the classifiers.

A valid factor that may limit the ability of the AI classifiers is the training methods employed and various associated parameters incorporated within them. Both ANNs and DTCs require the selection of a training algorithm and it is possible that those used here, although based on previous studies undertaken and detailed in the literature, do not provide the best solutions for the given problem. Furthermore, the use of the back-propagation algorithm for the training of the neural networks requires that both momentum and learning-rate terms are set. As with the choice of algorithm, these were based on previous studies and it is likely that alternative settings would have produced in different results. This point is emphasised by Jarvis and Stuart (1996). Neural network architecture will result in varying classification performance. Networks were trained with a single hidden layer comprising a single node, up to a maximum of 25 nodes. Again, this is likely to be problem specific and more hidden nodes, or indeed hidden layers, may have been led to improved network performance (see Maier and Dandy, 1998).

#### • Test data

There may be a number of factors involved in the apparent misclassifications, resulting in the interpreted risk map. These would involve the previously outlined issues regarding errors inherent in the independent and dependent variables, the ability of the independent variables to discriminate between different erosion processes, and classifier parameters. Furthermore, temporal issues may be important, where cells may be classified as eroding under a specified condition or extent, but may not be seen to be doing so at present, although they might do so over the course of time.

Taking the aforementioned points into consideration it may be possible to suggest where enhancements may be made to the AI approach documented within this study. consequently improving classifier performance. For example, the findings presented and discussed in Chapter Seven highlight the limitations associated with using the sodicity meter as an independent variable. Therefore, if an improved method by which sodicity can be accurately measured were used, or laboratory analysis were carried out on every sample, then overall classifier performance may be improved. Furthermore, training data could comprise of variables known to influence sodicity rather than attempt to assign a given level of sodicity based on some pre-requisite knowledge, the difficulties of which have been discussed extensively in the previous chapter. For example, training data could include variables such as soil pH, EC, ESP, CEC and organic content, thus allowing the classifier to identify potentially dispersive soils. A number of issues have been highlighted associated with DEM quality, and it may be possible to increase the classifier performance through incorporating improved DEM data. However, it is important to remember that the resolution of an elevation model has to represent the scale of the subject of interest, that of erosion in this instance. Therefore, it is not as straightforward as simply increasing DEM resolution. That is not to say that an improved method of DEM creation could not be employed for future studies. The number of training examples could also be increased and thereby present the classifiers with more examples from which to learn. The extent to which this may enhance classifier performance is unclear. However it is important to be aware that although it is reasonable to assume that by increasing the knowledge base from which the classifiers learn will lead to an increase in overall accuracy, Harris and Boardman (1998) found little improvement in levels of accuracy when increasing the number of training examples from 334 to 450 when classifying soil erosion. The

reasons proposed for this apparent anomaly were twofold. Firstly, the new data may have introduced more noise to the classifier and therefore limiting its effect, and secondly the classifier may have reached its optimum level of performance. Either way, this example highlights the difficulties associated with identifying potential areas of improvement when using AI techniques. It is impossible to know for certain, yet it may be the case that little is achieved through significant improvements to the quality and quantity of training data.

#### 8.7 SUMMARY

The methods described in this thesis have widespread practical applications for determining actual and/or potential erosion. Applying these methods requires careful consideration. Perhaps the most important issues for practitioners to determine relate to time scales, from project design to implementation, and associated resource requirements (cost). If a high-resolution investigation is required for a relatively small area, it may be better to undertake a detailed approach with a fieldwork based bias. In contrast, if knowledge of an extensive area is needed (e.g. regional scale) then alternative approaches may be appropriate.

Through the implementation of the research framework detailed in Chapter Five, the thesis has evaluated the applicability of two AI approaches to soil erosion mapping and risk assessment. It is possible therefore to state that decision tree classifiers offer the better approach of the two AI methods employed and, possess distinct advantages over discriminant analysis. As discussed previously, the applicability of an approach is vitally important and needs to be time efficient. In comparison to ANNs, decision trees require fewer parameters to be set prior to training. Furthermore, they perform

poorly when inter-class variabilities are present. Figure 8.3 outlines the required steps undertaken in the development of both an ANN and DTC classifier.

DTCs simply require the modeller to select the training algorithm. In contrast, ANNs trained using the back-propagation algorithm require that both momentum and learning rate parameters are set. Furthermore, network topology has to be chosen. As discussed in Chapter Four, this is largely a process of trial and error (Ghiassi and Saidane, 2005; Spellman, 1999), although some general guidelines have been proposed (see Blum, 1992; Berry and Linoff, 1997 and Bourquin et al., 1997). This is undoubtedly a time consuming process, with little insurance that the optimum network has been identified. The results of the classifications carried out in this study tend to suggest that the proposed guidelines offer little assistance in terms of optimal architecture, and the 'best' network is indeed case specific (see section 6.8). Nevertheless, the graphs indicate that the verification error relating to the networks trained using the field acquired variables does reduce drastically up to some specific point, beyond which oscillation may occur. Furthermore, in the majority of the these networks, results indicate that a minimum of five nodes are required in the hidden layer, as error rates appear to reduce significantly up to this point. Overall accuracies on the other hand reveal little in relation to potential patterns or trends.

In contrast, decision tree classifiers determine their own topology based upon relative cost. Therefore, the implications associated with employing neural networks require more consideration than that of DTCs. In practical terms, it is likely that for any given problem a trial and error procedure is required to ensure that a suitable network has

been identified. Even then the learning process may produce a relatively poor network.



Figure 8.3: The procedures required for the development of the two AI techniques.

While the implementation of DTCs allows the construction of erosion probability maps, such as those detailed in 6.9.1, artificial neural networks are not as viewable in terms of classifications and misclassifications associated with the training data. Therefore, the production of such maps would be a far more complicated timeconsuming process. Decision tree classifiers are also more efficient with training data. By using cross validation, as opposed to having a specific validation data set as is the case in the neural networks, each example presented to the tree is employed for training. Section 8.2.3 discussed the concept of interpreting the classified soil erosion maps as 'potential' erosion maps, inferring erosion risk. This approach has been adopted elsewhere and may present a vitally important management tool assisting in the implementation of appropriate policies and strategies. An erosion risk-by-association methodology has also been presented; the maps from which can be seen in section 6.9.2, Appendix 6 and on the accompanying DVD. The methodology outlined can be applied elsewhere and has the ability to highlight potentially susceptible locations based upon some simple topographical and process based rules.

The soil erosion maps presented within this thesis are indicative rather than definitive: they highlight susceptible areas where management strategies should be targeted and resources directed. It is suggested therefore, that the method presented be used to supplement and support alternative approaches and techniques.

The difficulties associated with determining between different erosion processes, and the degree to which they are operating, has been highlighted. Several factors have been proposed that may be responsible for the difficulties associated with discriminating between classes. However, even in the best case scenario; with good quality training data, optimum training parameters and useful independent variables, the differentiation between dependent variable classes will be small. As outlined by Foody (2002) and Felix and Binney (1989), classification errors often occur between highly similar classes. Moreover, attempts to classify on the basis of discrete classes may exacerbate the issue. The method presented within the thesis is only viable if it is easier to obtain the required variables than it is to physically map soil erosion. Nevertheless, the advantages associated with employing classification techniques include; the determination of rules, particularly in the case of DTCs, potential non-linear capabilities and, the ability to interpret classifications as 'risk' rather than 'actual' erosion.

#### **8.8 CONCLUSIONS**

This chapter has discussed, with reference to the aims and objectives, the results obtained through the implementation of the analytical framework and research methods set out in Chapter Five. With reference to the wider literature, the discussion has evaluated and compared the performance of the two AI techniques for soil erosion mapping, determined the influences associated with the independent and dependent variables and assessed the ability of ANNs and DTCs to inductively formulate useful rules and thresholds.

The following chapter summarises the main findings and highlights avenues where future research may focus. This is followed by some concluding comments drawing the thesis to a conclusion.

# 9

## Conclusion

#### 9.1 INTRODUCTION

The methodology presented in this thesis provides an alternative approach to the more traditional techniques employed for soil erosion mapping. Through the realisation of the potential benefits offered by Artificial Intelligence techniques, soil erosion maps and risk maps have been produced. These maps have a range of applications at varying spatial scales, with value to practitioners ranging from the individual landowner, to local and regional levels of government.

The thesis evaluates and compares two AI techniques for the mapping of soil erosion processes in Almería province, Southeast Spain. The approach aims to evaluate and compare the performance of the artificial neural networks and decision tree classifiers when trained, using readily available and low-cost data sets (e.g. digital elevation models), and additionally, more expensive field collected data. Furthermore, their ability to enhance our current understanding of soil erosion processes and how the selection of dependent and independent variables influence classifier performance has been assessed.

The following sections of this chapter summarise the main findings with reference to the stated aims outlined in Chapter One. As a result, it is possible to highlight potential avenues for future research, particularly in light of limitations and shortcomings associated with the approach employed here.

#### 9.2 SUMMARY OF THE MAIN FINDINGS

Several key findings, principally relating to the aims set out in this thesis, are identified.

- The composition of the training data set appears to detrimentally affect the performance of the artificial neural networks in a number of the classifications. The training data set comprises an uneven number of examples for each class within the dependent variable due largely to the opportunistic methods by which sampling was undertaken in the field. The variations were at their most extreme within the three and nine class classifications and, led to extensive commission errors and misleading results. In contrast, neither the decision trees nor the discriminate analysis classifiers appear susceptible to this problem.
- The inclusion of independent variables extracted from the digital elevation models
  into the field acquired training set often led to an increased verification error or
  relative cost in the neural network and decision tree classifications respectively.
  An attractive advantage often used for promoting the wider implementation of the
  ANNs and DTCs, is their ability to determine the importance of individual
  independent variables, and ignore those that detrimentally affect classifier
  performance (see Chapter Four). However, results indicate that error rates often
  increase when field acquired data is incorporated along with DEM data, as
  opposed to when the field data is used to train the classifier alone.
- Soil erosion maps produced from the classifications may be better interpreted as 'potential' erosion maps that infer soil erosion risk. The supervised AI techniques

employed here learn through the presentation of training cases, where the independent and dependent variables are known. The classifier assigns unknown cases based upon its learning experience. Thus, misclassifications or the difference between actual and the predicted erosion, can be interpreted as erosion risk.

- Discriminant analysis classifies soil erosion to comparable levels of accuracies as both AI techniques. DA consistently outperforms the neural networks in the three and nine class classifications because the linear technique does not appear to suffer as a result of training set composition and associated inter-class variability.
- Classifier performance is strongly controlled by the dependent and independent variables. As may be expected, classifier performance generally reduces as the complexity of the problem increases. The independent variables acquired in the field provided superior data to that extracted through either of the digital elevation models. In addition, the classifiers trained using variables extracted from the higher resolution digital elevation model proved superior to those trained using the coarser model. This is evident throughout the two, three and nine class classifications.
- A number of weak relationships were identified between various soil characteristics measured under laboratory conditions. The trends generally follow those identified by Faulkner *et al.* (2000) and Alexander *et al.* (1999). However, the Sodium Adsorption Ratios within these investigations tend to exceed those measured here. The likely causal factor is that samples were collected primarily

from the Mocatán badland site, in and around pipe entrances that appear highly dispersible. In contrast, samples collected for use within this study were located on slopes containing varying degrees of erosion, on the TRU and MRU (see Chapter Two).

- The difficulties associated with the determination and identification of dispersive soils has become readily evident. Although various dispersivity indices exist (e.g. Gerber and Harmse (1987) (Figure 3.9) and Rengasamy *et al.* (1984) (Figure 3.11)), it is difficult to classify or group soils into deflocculation classes. The results highlight the issue of conflicting boundaries identified by different authors producing different indices. Sites that were seen to be eroding in the field did not appear to be 'at risk', based primarily upon the laboratory findings. Due to the fact that soil dispersion is a highly complex process, involving a range of different chemical (e.g. pH, ESP, CEC, organic content etc.) and physical (e.g. slope angle) variables, it is important to consider all soil characteristics without basing presumptions upon one or two variables.
- There is no discernible relationship between the field sodicity meter and soil characteristics calculated through extensive laboratory work. Possible reasons for this include:
  - (i) The method

It is possible that the method itself may be flawed. To avoid the mechanical breakdown of the soil structure through shaking, the method advises that the jar (comprised of a 1:5 soil to water ratio) is inverted slowly once. The sample may not be sufficiently mixed and thus limit the

potential dispersion process. In addition, the mixture is allowed to settle for four hours, but it may take in excess of eight hours for clay particles that are still flocculated to settle.

- (ii) Soils may vary from those on which the meter has been calibrated
   The meter has been designed and produced in Australia, and it may be the case that the soils on which it has been calibrated vary in some way to those with which it has been applied here.
- A number of limitations were identified. Consideration of these limitations has been a valuable lesson in determining the applicability of these approaches to practical applications.

### 9.3 FUTURE RESEARCH

Through the course of this work, a range of potential research avenues have been identified. These are outlined below.

The interpretation of the classified erosion maps as erosion 'risk' maps can be explored further. The inclusion of land-use as an additional independent variable would provide an important social aspect to the models that will enhance their applicability. This will also allow the determination and assessment of landscape responses to various land-use changes. The trained classifiers offer the ability to run various scenarios based purely on physical parameters. For example, if a landscape manager required information concerning erosion under varying degrees of vegetation cover, it would be very easy to run a range of simple

scenarios. If land-use were incorporated, it would be possible to determine the impact that any change may have.

- Potential advantages of using classification ensemble approaches can be assessed.
   Various classification ensemble methods exist and, have been proven to improve classification accuracy (Gislason *et al.*, In Press). Such methods include boosting and bagging. Boosting involves the aggregation of models through voting (McBratney *et al.*, 2003) and the re-training of poorly classified samples. In contrast, bagging involves training many classifiers using bootstrapping methods where subsets of the training data are created to further enhance classifier performance (Gislason *et al.*, In Press; Lawrence *et al.*, 2004).
- Soil physico-chemical relationships can be explored further. The findings described in Chapter Seven highlight the complex relationships between various soil characteristics and in particular the inability and conflicting nature of the dispersivity classifications set out by Gerber and Harmse (1987) and Rengasamy et al. (1984). Furthermore, a method by which soil sodicity can be measured quickly and accurately is required and will provide a valuable means of producing a more rapid assessment of erosion risk in highly dispersive landscapes.

### 9.4 CONCLUDING COMMENTS

This thesis presents an alternative method by which soil erosion processes can be mapped, and the extent to which they are operating determined. Furthermore, the ability of such methods to contribute to current knowledge regarding soil erosion processes has been explored, as have the predictive abilities of different data sets and the influence of the number of classes in the dependent variable.

Mapping soil erosion processes remains a global challenge to geomorphologists. The methodology presented within this thesis provides an alternative approach to mapping soil erosion and highlighting susceptible areas where intervention should be prioritised. However, it is important to be aware of the associated limitations and in particular the dependence of such techniques upon good quality training data. Nevertheless, decision tree classifiers present a useful tool that does not appear to suffer as a consequence of inter-class variations in the data, which is often unavoidable when studying environmental phenomena.

In conclusion, the method presented for mapping actual soil erosion may be used to determine potential erosion. Both outputs provide valuable tools for landscape managers and may provide an alternative method by which intervention and resources may be directed.

# References

Abrahart, R. J. and White, S. M. (2001) Modelling sediment transfer in Malawi: Comparing backpropagation neural network solutions against a multiple linear regression benchmark using small data sets. *Physics and Chemistry of the Earth, Part* B: Hydrology, Oceans and Atmosphere, 26, pp. 19-24.

Abrahart, R. J., See, L. and Kneale, P. E. (1999) Applying saliency analysis to neural network rainfall-runoff modelling. Fourth International Conference on Geocomputation, Washington.

Alexander, J. A. and Mozer, M. C. (1999) Template-based procedures for neural network interpretation. *Neural Networks*, 12, pp. 479-498.

Alexander, R. W. and Calvo, A. (1990) The influence of lichens on slope processes in some Spanish badlands. In: Thornes, J. B. (Ed.) *Vegetation and Erosion*. Wiley, Chichester. pp. 385-398.

Alexander, R. W., Harvey, A. M., Calvo, A., James, P. A. and Cerda, A. (1994) Natural stabilisation mechanisms on badland slopes: Tabernas, Almería, Spain. In: Millington, A. C. and Pye, K. (Eds.) *Environmental Change in Drylands: Biogeographical and Geomorphological Perspectives*. Wiley, Chichester. pp. 85-111.

Alexander, R. W., Spivey, D. B., Faulkner, H. and Willshaw, K. (1996) Badland morphology and geoecology: Mocatán system: Processes and patterns. In: Mather, A. and Stokes, M. (Eds.) Second Cortijo Urra Field Meeting Southeast Spain. University of Plymouth.

Alexander, R. W., Spivey, D. B., Faulkner, H. and Willshaw, K. (1999) Badland morphology and geoecology: Mocatán system: processes and patterns. In: Mather, A.E. and Stokes, M. (Eds.) BSRG/BGRG SE Spain Field Meeting Guide Book. University of Plymouth, England. pp 134-151.

Alonso-Chaves, F. M., Soto, J. I., Orozco, M., Kilias, A. A. and Tranos, M. D. (2004) Tectonic evolution of the Betic Cordillera: An overview. *Bulletin of the Geological* Society of Greece, **36**, pp. 1598-1607.

Altman, E. I., Marco, G. and Varetto, F. (1994) Corporate distress diagnosis: Comparisons using linear discriminant analysis and neural networks (the Italian experience). *Journal of Banking and Finance*, 18, pp. 505-529. Anderson, J. A. (1986a) Cognitive capabilities of a parallel system. In: Bienenstock, E., Fogelman-Souli, F. and Weisbuch, G. (Eds.) *Disordered Systems and Biological Organisation*. NATO ASI series, F20. Spinger-Verlag, Berlin.

Anderson, J. A., Golden, R. M. and Murphy, G. L. (1986b) Concepts in distributed systems. In: Szu, H. H. (Ed.) *Optical and Hybrid Computing*. 634. Bellington, WA. Society of photo-optical instrumentation engineers. pp.260-272.

Ardahanlioglu, O., Oztas, T., Evren, S., Yilmaz, H. and Yildrim, Z. N. (2003) Spatial variability of exchangeable sodium, electrical conductivity, soil pH and boron content in salt- and sodium-affected areas of the Igdir plain (Turkey). *Journal of Arid Environments*, 54, pp. 495-503.

Arnáez, J. and Larrea, V. (1995) Erosion processes and rates on road-sides of hillroads (Iberian System, La Rioja, Spain). *Physics and Chemistry of the Earth*, 20, pp. 395-401.

Avery, B. W. and Bascomb, G. (1987) Soil Survey Laboratory methods. Prentice-Hall, London.

Back, B., Laitinen, T. and Sere, K. (1996) Neural networks and genetic algorithms for bankruptcy predictions. *Expert Systems with Applications*, 4, pp. 407-413.

Baillie, I. C., Faulkner, H., Espin, G. D. and Levett, M. J. (1986) Problems of protection against piping and surface erosion in central Tunisia. *Environment and Conservation*, 13, pp. 27-40.

Bell, F. G. and Walker, D. J. H. (2000) A further examination of the nature of dispersive soils in Natal, South Africa. *Quarterly Journal of Engineering Geology* and Hydrogeology, 33, pp. 187-199.

Benediktsson, J. A., Swain, P. H. and Ersoy, O. K. (1990). Neural network approaches versus statistical methods in classification of multi-source remote sensing data. *IEEE Transactions on Geoscience and Remote Sensing*, 28, pp. 540-551.

Berry, M. J. A. and Linoff, G. (1997) Data Mining Techniques. Wiley, Chichester.

Berzal, F., Cubero, J. C., Marín, N. and Sánchez, D. (In Press) Building multi-way decision trees with numerical attributes. *Information Sciences*.

Beven, K. and Kirkby, M. (1979) A physically based variable contributing area model of basin hydrology. *Hydrological Science Bulletin*, 24, pp. 43-64.

Beven, K. J., Kirkby, M. J., Schoffield, N. and Tagg, A. (1984) Testing a physicallybased flood forecasting model TOPMODEL for three UK catchments. *Journal of Hydrology*, **69**, pp. 119-143.

Beyer, H. L. (2003) Hawth's Analysis Tools for ArcGIS. http://www.spatialecology.com/htools (Accessed 20/1/2003). Bishop, C. M. (1995) Neural Networks for Pattern Recognition. Clarendon press, Oxford.

Blackard, J. A. and Dean, D. J. (1999) Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables. *Computers and Electronics in Agriculture*, **24**, pp. 131-151.

Blackburn, G. A. and Steele, C. M. (1999) Towards the remote sensing of Matorral vegetation physiology: Relationships between spectral reflectance, pigment, and biophysical characteristics of semiarid bushland canopies. *Remote Sensing of Environment*, **70**, pp. 278-292.

Blum, A. (1992) Neural networks in C++. Wiley, Chichester.

Boardman, J., Parsons, A. J., Holland, R., Holmes, P. J. and Washington, R. (2003) Development of badlands and gullies in the Sneeuberg, Great Karoo, South Africa. *Catena*, **50**, pp. 165-184.

Bobbin, J. and Recknagel, F. (2001) Knowledge discovery for prediction and explanation of blue-green algal dynamics in lakes by evolutionary algorithms. *Ecological Modelling*, 146, pp. 253-262.

Bocco, G. (1991) Gully erosion. Processes and models. Progress in Physical Geography, 15, pp. 392-406.

Bologna, G. (2004) Is it worth generating rules from neural network ensembles? *Journal of Applied Logic*, **2**, pp. 325-348.

Bondi, H. (1985) Risk in Perspective. In: Cooper, M. G. (Ed.) Risk: Man-Made Hazards to Man. Clarendon Press, Oxford. pp. 8-17.

Borak, J. S. and Strahler, A. H. (1999) Feature selection and land cover classification of a MODIS-like data set for semiarid environment. *International Journal of Remote Sensing*, **20**, pp. 919-938.

Boucher, S. C. (2002) The Initiation and Development of Tunnel Erosion near Costerfield, Victoria. Unpublished PhD. Thesis, Monash University, Australia.

Bourquin, J., Schmidi, H., van Hoogevest, P. and Leuenberger, H. (1997) Basic concepts of artificial neural networks (ANN) modelling in the application to pharmaceutical development. *Pharmaceutical Development and Technology*, **2**, pp. 95-109.

Bourrouilh, R. and Gorsline, D. S. (1979) Pre-Triassic fit and Alpine tectonics of continental blocks in the western Mediterranean. *Geological Society of America Bulletin*, **90**, pp. 1074-1083.

Bousquet, J. C. (1979) Quaternary strike-slip faults in southeastern Spain. *Tectonophysics*, **52**, pp. 277-286.

Bower, C. A., Reitemeier, R. F. and Fireman, M. (1952) Exchangeable cation analysis of saline and alkali soils. *Soil Science*, **73**, pp. 251-261.

Boyd, D. S., Foody, G. M. and Ripple, W. J. (2002) Evaluation of approaches for forest cover estimation in the Pacific Northwest, USA using remote sensing. *Applied Geography*, 22, pp. 375-392.

Bradley, A. P. and Lovell, B. C. (1995) Cost-sensitive Decision Tree Pruning: Use of the ROC Curve. Proceedings of the Eighth Australian Joint Conference on Artificial Intelligence. Canberra, ACT, pp. 1-8.

Bradley, P. S., Fayyad, U. M. and Mangasarian, O. L. (1998) Mathematical programming for data mining: formulation and challenges. *Informs Journal on Computing*, 3, pp. 217-238.

Brady, N. C. and Weil, R. R. (1999) The Nature and Properties of Soils (Twelfth Edition). Prentice-Hall, New Jersey.

Brady, N. C. and Weil, R. R. (2002) The Nature and Properties of Soils (Thirteenth Edition). Prentice-Hall, New Jersey.

Braga, J. C., Martín, J. M. and Quesada, C. (2003) Patterns and average rates of late Neogene-Recent uplift of the Betic Cordillera, SE Spain. *Geomorphology*, **50**, pp. 3-26.

Breiman, L., Friedman, J. H., Olshen, R. A. and Stone, C. J. (1984) Classification and Regression Trees. Wadsworth, Belmont, California.

Briggs, D. J. and Giordano, A. (1995) CORINE Soil Erosion Report. European Commission.

Brodley, C. E. and Utgoff, P. E. (1995) Multivariate decision trees. *Machine Learning*, 19, pp. 45-77.

Brouwer, R. K. (2002) A feed-forward network for input that is both categorical and quantitative. *Neural Networks*, **15**, pp. 881-890.

Brown, D. E., Corruble, V. and Pittard, C. L. (1993) A comparison of decision tree classifiers with back-propagation neural networks for multimodal classification problems. *Pattern Recognition*, **26**, pp. 953-961.

Brown, D. G., Lusch, D. P. and Duda, K. A. (1998) Supervised classification of glaciated landscapes using digital elevation data. *Geomorphology*, **21**, pp. 233-250.

Brown, J. F., Loveland, T. R., Ohlen, D. O. and Zhu, Z. (1999) The global land-cover characteristics database: the user's perspective. *Photogrammetric Engineering and Remote Sensing*, **65**, pp. 1069-1074.

Bryan, R. and Yair, A. (Eds.) (1982) Perspectives on studies of badland geomorphology. In: *Badland Geomorphology and Piping*. Geobooks, Norwich. pp. 1-12.

Bryan, R. B., Imeson, A. C. and Campbell, I. A. (1984) Solute release and sediment entrainment on microcatchments in the Dinosaur Park badlands, Alberta, Canada. *Journal of Hydrology*, **71**, pp. 79-106.

Buckheim, M. P. and Lillesand, T. M. (1989) Semi-automated training field extraction and analysis for efficient digital image classifications. *Photogrammetric Engineering and Remote Sensing*, **55**, pp. 1347-1355.

Bui, E. N., Henderson, B. L. and Viergever, K. (In Press) Knowledge discovery from models of soil properties developed through data mining. *Ecological Modelling*.

Bull, W. B. (1980) Geomorphic thresholds as defined by ratios. Coates, D. R. and Vitek, J. D. (Eds.) *Thresholds in Geomorphology*. Allen and Unwin, London. pp. 259-263.

Burke, S. and Thornes, J. B. (1998) Actions taken by National and Non-Governmental Organisations to Mitigate Desertification in the Mediterranean. Concerted Action on Mediterranean Desertification. Concerted Action Report 1.

Cai, Q. G., Wang, H., Curtin, D. and Zhu, Y. (2005) Evaluation of the EUROSEM model with single event data on Steeplands in the Three Gorges Reservoir Areas, China. *Catena*, **59**, pp. 19-33.

Calvo-Cases, A. and Harvey, A. M. (1996) Morphology and development of selected badlands in SE Spain: Implications of climatic change. *Earth Surface Processes and Landforms*, 21, pp. 725-735.

Campbell, I. A. (1982) Surface morphology and rates of change during a ten year period in the Alberta badlands. In: Bryan, R. B. and Yair, A. (Eds.) Badland Geomorphology and Piping. Geobooks, Norwich. pp. 221-238.

Campbell, I. A. (1989) Badlands and badland gullies. In: Thomas, D. S. G. (Ed.) Arid Zone Geomorphology. Belhaven, London. pp. 159-185.

Campbell, I. A. and Honsaker, J. L. (1982) Variability in badlands erosion; problems of scale and threshold identification. In: Thorn, C. E. (Ed.) Space and Time in Geomorphology. George Allen and Unwin, London. pp. 59-79.

Campbell, J. B. (1981) Spatial correlation effects upon accuracy of supervised classification of land cover. *Photogrammetric Engineering and Remote Sensing*, 47, pp. 355-363.

Campbell, J. B. (2002) Introduction to Remote Sensing (Third Edition). Taylor and Francis, London.

Canters, F. (1997) Evaluating uncertainty of area estimates derived from fuzzy landcover classification. *Photogrammetric Engineering and Remote Sensing*, **63**, pp. 403-414.

Cawsey, A. (1998) The essence of artificial intelligence. Prentice Hall, Harlow.

Cerdà, A. (1997) The effect of patchy distribution of *Stipa tenacissima* L. on runoff and erosion. *Journal of Arid Environments*, **36**, pp. 37-51.

Cerdá, A. (1999) Parent material and vegetation affect soil erosion in Eastern Spain. Soil Science Society of America Journal, 63, pp. 362-368.

Cerdà, A. and García-Fayos, P. (1997) The influence of slope angle on sediment, water and seed losses on badland landscapes. *Geomorphology*, 18, pp. 77-90.

Chang, D. H. and Islam, S. (2000) Estimation of soil physical properties using remote sensing and artificial neural network. *Remote Sensing of Environment*, 74, pp. 534-544.

Chen, D. and Stow, D. (2002) The effect of training strategies on supervised classifications at different spatial resolutions. *Photogrammetric Engineering and Remote Sensing*, 68, pp. 1155-1161.

Chen, Y, Jiao, T, McCall, T. W., Baichwal, A. R. and Meyer, M. C. (2002) Comparison of four artificial neural network software programs used to predict the in vitro dissolution of controlled-release tablets. *Pharmaceutical Development and Technology*, 7, pp. 373-379.

Chorley, R. J. and Schumm, S. A. (1984) Geomorphology. Menthuen, London.

Christiansen, J. E. (1947) Some permeability characteristics of saline and alkali soils. Agric. Eng., 28, pp. 147-150.

Churchman, G. J. and Weissmann, D. A. (1995) Particle mobility as a parameter in soil dispersibility. In: Naidu, R., Sumner, M. E. and Rengasamy, P. (Eds.) Australian Sodic Soils: Distribution, Properties and Management. pp. 191-194. CSIRO publication.

Churchman, G. J. Skjemstad, J. O. and Oades, J. M (1995) Effects of clay minerals and organic matter on sodicity. In: Naidu, R., Sumner, M. E. and Rengasamy, P. (Eds.) Australian Sodic Soils: Distribution, Properties and Management. CSIRO, Melbourne. pp. 107-119.

Civco, D. L. (1993) Artificial neural networks for land-cover classification and mapping. Int. J. Geographic Information Systems, 7, pp. 173-186.

Congalton, R. G. (1991) A review of assessing the accuracy of classifications of remotely sensed data. *Remote Sensing of Environment*, 37, pp. 35-46.

Congalton, R.G. and Green, K. (1999) Assessing the Accuracy of Remotely Sensed Data: Principles and Practices. Lewis Publishers, New York.

Corwin, D. L. and Lesch, S. M. (2005) Characterising soil spatial variability with apparent soil electrical conductivity I. Survey protocols. *Computers and Electronics in Agriculture*, **46**, pp. 103-133.

Corwin, D. L., Kaffka, S. R., Hopmans, J. W., Mori, Y., van Groenigen, J. W., van Kessel, C., Lesch, S. M. and Oster, J. D. (2003) Assessment and field-scale mapping of soil quality properties of a saline-sodic soil. *Geoderma*, **114**, pp. 231-259.

Crouch, R. J. (1976) Field tunnel erosion – a review. Journal of Soil Conservation Service, 32, pp. 98-111.

Cyr, L., Bonn, F. and Pesant, A. (1995) Vegetation indices derived from remote sensing for an estimation of soil protection against water erosion. *Ecological Modelling*, **79**, pp. 277-285.

D'Acqui, L. P. D., Churchman, G. J., Janik, L. J., Ristori, G. G. and Weissmann, D. A. (1999) Effect of low organic matter removal by low-temperature ashing on dispersion of undisturbed aggregates from a tropical crusting soil. *Geoderma*, 93, pp. 311-324.

Dai, H. and MacBeth, C. (1997) Effects of learning parameters on learning procedure and performance of a BPNN. *Neural Networks*, **10**, pp. 1505-1521.

Davis, J. G., Waskom, R. M., Bauder, T. A. and Cardon, G. E. (2003) Managing sodic soils. http://www.ext.colostate.edu/pubs/crops/00504.pdf (Accessed 10/11/2003).

De Boer, D. H. and Campbell, I. A. (1990) Runoff chemistry as an indicator of runoff sources and routing in semi-arid, badland drainage basins. *Journal of Hydrology*, **121**, pp. 379-394.

de Carvalho, L. M. T., Clevers, J. G. P. W., Skidmore, A. K. and de Jong, S. M. (2004) Selection of imagery data and classifiers for mapping Brazilian semideciduous Atlantic forests. *International Journal of Applied Earth Observation and Geoinformation*, 5, pp. 173-186.

de Jong, S. M., Paracchini, M. L., Bertoli, F., Folving, S., Megier, J. and De Roo, A. P. J. (1999) Regional assessment of soil erosion using the distributed model SEMMED and remotely sensed data. *Catena*, **37**, pp. 291-308.

de la Rosa, D., Mayol, F., Moreno, J. A., Bonson, T. and Lozano, S. (1999) An expert system/neural network model (ImpelERO) for evaluating agricultural soil erosion in Andalucian region, southern Spain. Agriculture, Ecosystems and Environment, 73, pp. 211-226.

de la Rosa, D., Moreno, J. A., Mayol, F. and Bonsón, T. (2000) Assessment of soil erosion vulnerability in western Europe and potential impact on crop productivity due

to loss of soil depth using the ImpelERO model. Agriculture, Ecosystems and Environment, 81, pp. 179-190.

de Larouzíere, F. D., Bolze, J., Bordet, P., Hernandez, J., Montenat, C., and D'Estevou, P. (1988) The Betic segment of the lithospheric Trans-Alboran shear zone during the Late Miocene. *Tectonophysics*, 152, pp. 41-52.

De Roo, A. P. J. (1993) Modelling surface runoff and soil erosion in catchments using Geographical Information systems. Validity and applicability of the ANSWERS model in two catchments in the loess area of south Limburg and Devon. *Nederlandse Geografische Studies*, No. 157.

De Roo, A.P.J., Wesseling, C.G. and Ritsema C.J (1996) LISEM: a single-event physically based hydrological and soil erosion model for drainage basins. I: theory, input and output. *Hydrological Processes*, 10, pp. 1107-1117.

Dedecker, A. P., Goethals, P. L. M., Gabriels, W. and De Pauw, N. (2004) Optimization of Artificial Neural Network (ANN) model design for prediction of macroinvertebrates in the Zwalm river basin (Flanders, Belgium). *Ecological Modelling*, 174, pp. 161-173.

DeFries, R. S. and Cheung-Wai Chan, J. (2000) Multiple criteria for evaluating machine learning algorithms for land cover classification from satellite data. *Remote Sensing of Environment*, 74, pp. 503-515.

Delcourt, H. R. and Delcourt, P. A. (1988) Quaternary landscape ecology: Relevant scales in space and time. Landscape Ecology, 2, pp. 23-44.

Desmet, P. J. and Govers, G. (1996) A GIS procedure for the automated calculation of the USLE LS factor on topographically complex landscape units. *Journal of Soil and Water Conservation*, **51**, pp. 427-433.

Desmet, P. J. J., Poesen, J., Govers, G. and Vandaele, K. (1999) Importance of slope gradient and contributing area for optimal prediction of the initiation and trajectory of ephemeral gullies. *Catena*, **37**, pp. 377-392.

Dicks, S. E. and Lo, T. H. C. (1990) Evaluation of thematic map accuracy in a landuse and land-cover programme. *Photogrammetric Engineering and Remote Sensing*, 56, pp. 1247-1252.

Dowla, F., Taylor, S. R., and Anderson, R. W. (1990) Seismic discrimination with artificial neural networks: preliminary results with regional spectral data. *Bulletin of the Seismological Society of America*, **80**, pp. 1346-1373.

Downey, I. D., Power, C. H, Kanellopoulos, I. and Wilkinson, G. G. (1992) A performance comparison of Landsat TM land cover classification based on neural network techniques and traditional maximum likelihood and minimum distance algorithms. In: Cracknell, R. A. and Vaughan, R. A. (Eds.) Remote Sensing: From Research to Operation (Proceedings of the 18<sup>th</sup> Annual Conference of the UK Remote Sensing Society. Dundee, pp. 518-528.

Drew, D. P. (1982) Piping in the Big Muddy badlands, southern Saskatchewan, Canada. In: Bryan, R. and Yair, A. (Eds.) *Badland Geomorphology and Piping*. Geobooks, Norwich. pp. 293-304.

Edwards, K. (1991) Soil formation and erosion rates. In: Chairman, P. E. V. and Murphy, B. W. (Eds.) Soils, Their Properties and Management. Sydney University Press, Sydney. pp. 36-47.

Elges, H. F. W. K. (1985) Dispersive soils. The Civil Engineer in South Africa, 27, pp. 347-355.

Ellis, F. (1997) Evaluating Techniques for Soil Erosion Modelling: A role for Artificial Intelligence? Unpublished PhD. Thesis. Australian National University, Canberra.

Ellis, F. G. (2002) The application of machine learning techniques to erosion modelling. http://www.sbg.ac.at (Accessed 12/06/02).

Ermini, L., Catani, F. and Casagli, N. (2005) Artificial neural networks applied to landslide susceptibility assessment. *Geomorphology*, **66**, pp. 327-343.

Evans, R. (1980) Mechanics of water erosion and their spatial and temporal controls: an empirical viewpoint. In: Kirkby, M. J. and Morgan, R. P. C. (Eds.) *Soil Erosion*. Wiley, Chichester. pp. 109-128.

Evans, R. (1992) Assessing erosion in England and Wales. Proceedings 7<sup>th</sup> ISCO Conference. pp. 82-91.

Farifteh, J., Farshad, A. and George, R. J. (In Press) Assessing salt-affected soils using remote sensing, solute modelling, and geophysics. *Geoderma*.

Faulkner, H., Alexander, R. and Wilson, B. R. (2003a) Changes to the dispersive characteristics of soils along an evolutionary slope sequence in the Vera badlands, southeast Spain: implications for site stabilisation. *Catena*, **50**, pp. 243-254.

Faulkner, H., Ruiz, J., Zukowskyj, P. and Downward, S. (2003b) Erosion risk associated with rapid and extensive agricultural clearances on dispersive materials in southeast Spain. *Environmental Science and Policy*, **6**, pp. 115-127.

Faulkner, H., Spivey, D. and Alexander, R. (2000) The role of some site geochemical processes in the development and stabilisation of three badland sites in Almería, Southern Spain. *Geomorphology*, **35**, pp. 87-99.

Fausett, L. (1994) Fundamentals of neural networks: Architectures, Algorithms, and Applications. Prentice Hall, New Jersey.

Fayed, U. M. and Irani, K. B. (1992) Technical note on the handling of continuous values attributes in decision tree generation. *Machine Learning*, **8**, pp. 87-102.

Felix, N. A. and Binney, D. L. (1989) Accuracy assessment of a Landsat-assisted vegetation map of the coastal plain of the Arctic National Wildlife Refuge. *Photogrammetric Engineering and Remote Sensing*, **55**, pp. 475-478.

Fireman, M. and Wadleigh, C. H. (1951) A statistical study of the relation between pH and the exchangeable-sodium-percentage of western soils. *Soil Science*, **71**, pp. 273-285.

Fitzgerald, E. M. and Bean, C. J. (2001) Sub-basalt imaging problems and the application of artificial neural networks. *Journal of Applied Geophysics*, **48**, pp. 183-197.

Fitzgerald, R. W. and Lees, B. G. (1993) Assessing the classification accuracy of multisource remote sensing data. *Remote Sensing of Environment*, 47, pp. 362-368.

Flanagan, D. C. and Nearing, M. A. (Eds.) (1995) USDA water erosion prediction project (WEPP). Hillslope profile and watershed model documentation. NSERL Report No. 10, USDA-ARS, West Lafayette, IN.

Floras, S. A. and Sgouras, I. D. (1999) Use of geoinformation techniques in identifying and mapping areas of erosion in a hilly landscape of central Greece. *International Journal of Applied Earth Observation and Geoinformation*, 1, pp. 68-77.

Folly, A., Quinton, J. N. and Smith, R. E. (1999) Evaluation of the EUROSEM model using data from the Catsop watershed, The Netherlands. *Catena*, **37**, pp. 507-519.

Food and Agriculture Organisation (1965) Soil erosion by water. Some measures for its control on cultivated lands. Agricultural Paper 81, Rome.

Foody, G. M. (1995) Using prior knowledge in artificial neural network classification with a minimal training set. *International Journal of Remote Sensing*, 16, pp. 301-312.

Foody, G. M. (2002) Status of land cover classification accuracy assessment. *Remote Sensing of Environment*, **80**, pp. 185-201.

Foody, G. M. and Arora, M. K. (1996) Incorporating mixed pixels in the training, allocation and testing stages of supervised classifications. *Pattern Recognition Letters*, 17, pp. 1389-1398.

Foody, G. M. and Mathur, A. (2004) Toward intelligent training of supervised image classifications: directing training data acquisition for SVM classification. *Remote Sensing of Environment*, 93, pp. 107-117.

Foody, G. M., McCulloch, M. B. and Yates, W. B. (1995a) The effect of training set size and composition on artificial neural network classification. *International Journal of Remote Sensing*, **16**, pp. 1707-1723.

Foody, G. M., McCulloch, M. B. and Yates, W. B. (1995b) Classification of remotely sensed data by an artificial neural network: Issues related to training data characteristics. *Photogrammetric Engineering and Remote Sensing*, **61**, pp. 391-401.

Friedl, M. A. and Brodley, C. E. (1997) Decision tree classification of land cover from remotely sensed data. *Remote Sensing of Environment*, **61**, pp. 399-409.

Gahegan, M. and West, G. (1998) The classification of complex geographic datasets: An operational comparison of artificial neural networks and decision tree classifiers. *Third International Conference on GeoComputation*, Bristol, United Kingdom.

Gallant, J., C. and Wilson, J. P. (1996) TAPES-G: A grid-based terrain analysis program for the environmental sciences. *Computers and Geosciences*, 22, pp. 713-722.

García, A., Zhu, Z., Ku, T. L., Sanz de Galdeano, C., Chadwick, O. A. and Chacón Montero, J. (2003) Tectonically driven landscape development within the eastern Alpujarran corridor, Betic Cordillera, SE Spain (Almería). *Geomorphology*, **50**, pp. 83-110.

García, M. and Chuvieco, E. (2004) Assessment of the potential of SAC-C/MMRS imagery for mapping burned areas in Spain. *Remote Sensing of Environment*, **92**, pp. 414-423.

García-Hernández, M., Lopez-Garrido, A. C., Rivas, P., Sanz de Galdeano, C. and Vera, J. A. (1980) Mesozoic paleogeographic evolution of the external zones of the Betic Cordillera. *Geol. Mijnbouw*, **59**, pp. 155-168.

García-Latorre, J., García-Latorre, J. and Sanchez-Pícon, A. (2001) Dealing with aridity: socio-economic structures and environmental changes in an arid Mediterranean region. *Land Use Policy*, 18, pp. 53-64.

Gary, C. (2000) Book reviews. Scientia Horticulturae, 86, pp. 261-266.

Geeson, N. A. and Thrones, J. B. (1996) MEDALUS II. Executive Summary Phase II. Davies Wise, Bristol.

Gerber, A. and Harmse, H. J. Von M. (1987) Proposed procedure for identification of dispersive soils by chemical testing. *The Civil Engineer in South Africa*, **29**, pp. 397-399.

German, G. W. H. and Gahegan, M. N. (1996) Neural network architectures for the classification of temporal image sequences. *Computers and Geosciences*, 22, pp. 969-979.

German, G. W. H., West, G. and Gahegan, M. (2002) Statistical and AI techniques in<br/>classification:Acomparison.http://divcom.otago.ac.nz/sirc/webpages/99German.pdf (Accessed 12/8/04).

Ghiassi, M. and Saidane, H. (2005) A dynamic architecture for artificial neural networks. *Neurocomputing*, **63**, pp. 397-413.

Ghidley, F. and Alberts, E. E. (1996) Comparison of measured and WEPP predicted runoff and soil loss for midwest claypan soil. *Trans. ASAE*, **39**, pp. 1395-1402.

Gislason, P. O., Benediktsson, J. A. and Sveinsson, J. R. (In Press) Random forests for land cover classification. *Pattern Recognition Letters*.

Goel, P. K., Prasher, S. O., Patel, R. M., Landry, J. A., Bonnell, R. B. and Viau, A. A. (2003) Classification of hyperspectral data by decision trees and artificial neural networks to identify weed stress and nitrogen status of corn. *Computers and Electronics in Agriculture*, **39**, pp. 67-93.

Goh, A. T. C. (1995) Modelling soil correlation's using neural networks. Journal of Computing in Civil Engineering, 9, pp. 275-277.

Gómez, H. and Kavzoglu, T. (2005) Assessment of shallow landslide susceptibility using artificial neural networks in Jabonosa River Basin, Venezuela. *Engineering Geology*, **78**, pp. 11-27.

Gong, P. (1996) Integrated analysis of spatial data from multiple sources: Using evidential reasoning and artificial neural network techniques for geological mapping. *Photogrammetric Engineering and Remote Sensing*, **62**, pp. 513-523.

Goodchild, M. F., Steyaert, L. T. and Parks, B. O. (1996) Process and Research Issues. GIS world books, For Collins.

Goossens, R. and Van Ranst, E. (1998) The use of remote sensing to map gypsiferous soils in the Ismailia Province, Egypt. Geoderma, 87, pp. 47-56.

Gorman, R. P. and Sejnowski, T. J. (1988) Analysis of hidden units in a layered network trained to classify sonar targets. *Neural Networks*, 1, pp. 75-89.

Greene, R. S. B., Eggleton, R. A. and Rengasamy, P. (2002) Relationship between clay mineralogy and the hardsetting properties of soils in the Carnarvon horticultural district of Western Australia. *Applied Clay Science*, **20**, pp. 211-223.

Gupta, J. N. D. and Sexton, R. S. (1999) Comparing backpropagation with a genetic algorithm for neural network training. *Omega*, 27, pp. 679-684.

Gupta, R. K., Bhumbla, D. K. and Abrol, I. P. (1984) Effect of soil pH, organic matter and calcium carbonate on the dispersion behaviour of alkali soils. *Soil Science*, 137, pp. 245-251.

Haase, P., Pugnaire, F. I., Clark, S. C. and Incoll, L. D. (2000) Photosynthetic rate and canopy development in drought-deciduous shrub and *Anthyllis cytisoides* L. Journal of Arid Environments, **46**, pp. 79-91.

Haase, P., Pugnaire, F. I., Maria Fernández, E., Puigdefábregas, J., Clark, S. C. and Incoll, L. D. (1996) An investigation of rooting depth of the semiarid shrub *Retama* sphaerocarpa (L.) Boiss. By labelling of ground water with a chemical tracer. Journal of Hydrology, 177, pp. 23-31.

Haboudane, D., Bonn, F., Royer, A., Sommer, S. and Mehl, W. (2002) Land degradation and erosion risk mapping by fusion of spectrally-based information and digital geomorphometric attributes. *International Journal of Remote Sensing*, 23, pp. 3795-3820.

Hadley, R. F. and Toy, T. J. (1977) Relation of surficial erosion on hillslopes to profile geometry. *Journal of Research of US Geological Survey*, 5, pp. 487-490.

Hairsine, P. B. and Rose, C. W. (1992a) Modelling water erosion due to overland flow using physical principles: 1, sheet flow. *Water Resources Research*, 28, pp. 245-250.

Hairsine, P. B. and Rose, C. W. (1992b) Modelling water erosion due to overland flow using physical principles: 2, Rill flow. *Water Resources Research*, 28, pp. 237-243.

Hampson, S. E. and Volper, D. J. (1986) Linear function neurons: structure and training. *Biol. Cybern*, 53, pp. 203-217.

Han, J. and Kamber, M. (2001) Data Mining: Concepts and Techniques. Academic Press, San Diego.

Hansen, M., Dubayah, R. and DeFries, R. (1996) Classification trees: an alternative to traditional land cover classifiers. *International Journal of Remote Sensing*, 17, pp. 1075-1081.

Harris, T. and Boardman, J. (1990) A rule-based expert system approach to predicting waterborne soil erosion. In: Boardman, J., Foster, I. D. L. and Dearing, J. A. (Eds.) *Soil Erosion on Agricultural Land*. Wiley, Chichester. pp. 401-412.

Harris, T. M. and Boardman, J. (1998) Alternative approaches to soil erosion prediction and conservation using expert systems and neural networks. In: Boardman, J. and Favis-Mortlock, D. (Eds.) *Modelling Soil Erosion by Water*. Springer, Berlin. pp. 461-477.

Harvey, A. M. (1982) The role of piping in the development of badlands and gully systems in Southeast Spain. In: Bryan, R. and Yair, A. (Eds.) Badland Geomorphology and Piping. Geobooks, Norwich. pp. 317-335.

Harvey, A. and Calvo, A. (1989) Distribution of badlands in southeast Spain: Implications of climate change. In: Imeson, A. C. and De Groot, R. S. (Eds.) Landscape-Ecological Impact of Climate Change. Discussion Report on Mediterranean Region. pp. 14. Harvey, A. M. and Wells, S. G. (1987) Response of Quaternary fluvial systems to differential epeirogenic uplift: Aguas and Feos river systems, southeast Spain. *Geology*, 15, pp. 689-693.

Harvey, A. M., Alexander, R. W. and Spivey, D. B. (2001) Excursion: Badlands. In: Mather, A., Martín, J. M., Harvey, A. M. and Braga, J. C. (Eds.) *A Field Guide to the Neogene Sedimentary Basins of the Almería Province, South-East Spain.* Blackwell Science, Oxford. pp. 279-303.

Harvey, A. M., Miller, S. Y. and Wells, S. G. (1995) Quaternary soil and river terrace sequences in the Aguas/Feos river systems: Sorbas basin, southeast Spain. In: Lewin, J., Macklin, M. G. and Woodward, J. C. (Eds.) *Mediterranean Quaternary River Environments*. Balkema, Rotterdam. pp. 263-281.

Hashemi, R. R., Le Blanc, L. A., Rucks, C. T. and Rajaratnam, A. (1998) A hybrid intelligent system for predicting bank holding structures. *European Journal of Operational Research*, 109, pp. 390-402.

Henderson, B. L., Bui, E. N., Moran, C. J. and Simon, D. A. P. (2005) Austrial-wide predictions of soil properties using decision tress. *Geoderma*, 124, pp. 383-398.

Hepner, G. F., Logan, T., Ritter, N., Bryant, N. (1990) Artificial neural network classification using a minimal training set: Comparison to conventional supervised classification. *Photogrammetric Engineering and Remote Sensing*, **56**, pp. 469-473.

Heywood, I., Cornelius, S. and Carver, S. (1998) An introduction to Geographical Information Systems. Longman, Essex.

Hinton, G. E. (1989) Connectionist learning procedures. Artificial Intelligence, 40, pp. 185-234.

Hinton, G. E. (1990) Connectionist learning procedures. In: Kondratoff, Y. and michalski, R. (Ed.) *Machine learning III*. pp. 555-610. Morgan Kaufman, San Mateo.

Hixson, M., Scholz, D. and Fuhs, N. (1980) Evaluation of several schemes for classification of remotely sensed data. *Photogrammetric Engineering and Remote Sensing*, **46**, pp. 1547-1553.

Hodges, W. K. and Bryan, R. B. (1982) The influence of material behaviour on runoff initiation in the Dinosaur Badlands, Canada. In: Bryan, R. and Yair, A. (Eds.) Badland Geomorphology and Piping, Geobooks, Norwich. pp. 13-46.

Hontoria, L., Aguilera, J., and Zufiria, P. (In Press) An application of the multilayer perceptron: Solar radiation maps in Spain. *Solar Energy*.

Hooda, P. S. (1992) The Behaviour of Trace Metals in Sewage Sludge-Amended Soils. Unpublished PhD. Thesis, University of London.

Horney, R. D., Taylor, B., Munk, D. S., Roberts, B. A., Lesch, S. M. and Plant, R. E. (2005) Development of practical site-specific management methods for reclaiming salt-affected soil. *Computers and Electronics in Agriculture*, **46**, pp. 379-397.

Horton, R. E. (1933) The role of infiltration in the hydrologic cycle. *Trans. AGU*, 14, pp. 446-460.

Horton, R. E. (1945) Erosional development of streams and their drainage basins; hydrophysical approach to quantitative morphology. *Bulletin of the Geological Society of America*, **56**, pp. 275-370.

Hosking, P. L. (1967) Tunnelling erosion in New Zealand. Journal of Soil and Water Conservation, 22, pp. 149-151.

Hughes, G. F. (1968) On the mean accuracy of statistical pattern recognisers. *IEEE Transactions on Information Theory*, 14, pp. 55-63.

Hussain, A. S., Yu, A. and Johnson, R. D. (1991) Application of neural computing in pharmaceutical product development. *Pharmaceutical Research*, **8**, pp. 1248-1252.

Hutchinson, M. F. (1988) Calculation of hydrologically sound digital elevation models. Third International Symposium on Spatial Data Handling, Sydney. Columbus, Ohio: International Geographical Union

Hutchinson, M. F. (1989) A new procedure for gridding elevation and stream line data with automatic removal of spurious pits. *Journal of Hydrology*, **106**, pp. 211-232.

ICONA (1988) Agresividad de las lluvias en España. Ministerio de Agricultura, Pesca y Alimentación. ICONA, Madrid.

Imeson, A. C. and Emmer, I. M. (1992) Implications of climate change on land degradation in the Mediterranean. In: Jeftic, L., Milliman, J. D. and Sestini, G. (Eds.) *Climatic Change in the Mediterranean*. Arnold, London. pp. 95-128.

Imeson, A. C. and Verstraten, J. M. (1985) The erodibility of highly calcareous soil material from southern Spain. *Catena*, **12**, pp. 291-306.

Imeson, A. C., Kwaad, F. J. P. M. and Verstraten, J. M. (1982) The relationship of soil physical and chemical properties to the development of badlands in Morocco. In: Bryan, R. and Yair, A. (Eds.) *Badland Geomorphology and Piping*. Geobooks, Norwich. pp. 47-70.

Irvine, S. A. and Reid, D. J. (2001) Field prediction of sodicity in dryland agriculture in Central Queensland, Australia. *Australian Journal of Soil Research*, **39**, pp. 1349-1357.

Jackson, Q. and Landgrebe, D. (2001) An adaptive classifier design for highdimensional data analysis with a limited training data analysis with a limited training data set. *IEEE Transactions on Geoscience and Remote Sensing*, **39**, pp. 2664-2679.
Jaiswal, S., Benson, E. R., Bernard, J. C. and Van Wicklen, G. L. (2005) Neural network modelling and sensitivity analysis of a mechanical poultry catching system. *Biosystems Engineering*, **92**, pp. 59-68.

Jarvis, C. H. and Stuart, N. (1996) The sensitivity of a neural network for classifying remotely sensed imagery. *Computers and Geosciences*, 22, pp. 959-967.

Jayasuriya, R. T. (2003) Economic assessment of technological change and land degradation in agriculture: application to the Sri Lanka tea sector. *Agricultural Systems*, **78**, pp. 405-423.

Jones, J. A. A. (1981) The Nature of Soil Piping: A Review of Research (BGRG Research Monograph 3). Geobooks, Norwich.

Kamphorst, A. and Bolt, G. H. (1976) Adsorption of cations by soil. In: Bolt, G. H. and Bruggenwert, M. G. M. (Eds.) Soil Chemistry: Basic Elements. Elsevier, London, pp. 54-90.

Kamphorst, A. and Bolt, G. H. (1978) Saline and Sodic Soils. In: Bolt, G. H. and Bruggenwert, M. G. M. (Eds.) Soil Chemistry A. Basic Elements. Developments in Soil Science 5A. Elsevier, Amsterdam. pp. 171-191.

Keller, J. V. A., Hall, S. H., Dart, C. J. and McClay, K. R. (1995) The geometry and evolution of a transpressional strike-slip system: the Carboneras fault, SE Spain. *Journal of the Geological Society*, **152**, pp. 339-351.

Kervahut, T. and Potvin, J, Y. (1996) An interactive-graphic environment for automatic generation of decision trees. *Decision Support Systems*, 18, pp. 117-134.

Kim, H. and Koehler, G. J. (1995) Theory and practice of decision tree induction. Omega Int. J. Mgmt Sci., 23, pp. 637-652.

King, C. and Delpont, G. (1993) Spatial assessment of erosion: Contribution of remote sensing, a review. *Remote Sensing Reviews*, 7, pp. 223-232.

Knisel, W. G. (1980) CREAMS: a field scale model for chemicals, runoff and erosion from agricultural management systems. USDA Conservation Research Report 26.

Kononenko, I., Bratko, I. and Roskar, E. (1984) *Experiments in Automatic Learning* of Medical Diagnostic Rules (Technical Report). Jozef Stefan Institute, Ljubljana, Yugoslavia.

Kosmas, C., Donalatos, N. G. and Gerontidis, St. (2000) The effect of land parameters on vegetation performance and degree of erosion under Mediterranean conditions. *Catena*, **40**, pp. 3-17.

Kuhn, N. J. and Bryan, R. B. (2004) Drying, soil surface condition and interrill erosion on two Ontario soils. *Catena*, 57, pp. 113-133.

Kusiak, A., Caldarone, C. A., Kelleher, M. D., Lamb, F. S., Persoon, T. J. and Burns, A. (2006) Hypoplastic left heart syndrome: knowledge discovery with a data mining approach. *Computers in Biology and Medicine*, **36**, pp. 21-40.

Lawrence, R., Bunn, A., Powell, S. and Zambon, M. (2004) Classification of remotely sensed imagery using stochastic gradient boosting as a refinement of classification tree analysis. *Remote Sensing of Environment*, **90**, pp. 331-336.

Lázaro, R., Rodrigo, F. S., Gutiérrez, L., Domingo, F. and Puigdefábregas, J. (2001) Analysis of a 30-year rainfall record (1967-1997) in semi-arid SE Spain for implications on vegetation. *Journal of Arid Environments*, **48**, pp. 373-395.

Le Bissonnais, Y., Montier, C., Jamagne, M., Daroussin, J. and King, D. (2001) Mapping erosion risk for cultivated soil in France. *Catena*, **46**, pp. 207-220.

Le Cun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W. and Jackel, L. D. (1990) Handwritten digit recognition with a backpropagation network. In: Touretsky, D. S. (Ed.) Advances in Neural Information Processing Systems 2. pp. 396-404. Morgan Kaufman, San Mateo.

Leane, M. M., Cumming, I. and Corrigan, O. (2003) The use of artificial neural networks for the selection of the most appropriate formulation and processing variables in order to predict the in vitro dissolution of sustained release minitablets. *AAPS PharmSciTech*, 4, pp. 1-12.

Lee, J., Weger, R. C., Sengupta, S. K. and Welch, R. M. (1990) A neural network approach to cloud classification. *IEEE Transactions on Geoscience and Remote Sensing*, 28, pp. 846-855.

Lees, B. (1996) Sampling strategies for machine learning using GIS. In: Goodchild, M. F., Steyaert, L. T. and Parks, B. O. (Ed.) GIS and Environmental Modelling: Process and Research Issues. pp. 39-42. GIS world books, For Collins.

Lengellé, R. and Denœux, T. (1996) Training MLP's layer by layer using an objective function for internal representations. *Neural Networks*, 9, pp. 83-97.

Lillesand, T. M. and Kiefer, R. W. (2000) Remote Sensing and Image Interpretation (Fourth Edition). Wiley, Chichester.

Liu, W. and Wu, E. Y. (2005) Comparison of non-linear mixture models: sub-pixel classification. *Remote Sensing of Environment*, 94, pp. 145-154.

Lonergan, L. (1993) Timing and kinematics of deformation in the Malaguide Complex, internal zone of the Betic Cordillera, Southeast Spain. *Tectonics*, 12, pp. 460-476.

Lonergan, L., Platt, J. P. and Gallagher, L. (1994) The internal/external zone boundary in the eastern Betic Cordillera, SE Spain. *Journal of Structural Geology*, 16, pp. 175-188.

Longley, P. A., Goodchild, M. F., Maguire, D. J. and Rhind, D. W. (2001) Geographic Information Systems and Science. Wiley, Chichester.

López-Bermúdez, F. and Romero-Diaz, M. A. (1989) Piping erosion and badland development in South-East Spain. In: Yair, A. and Berkowicz, B. (Eds.) Arid and Semi-arid Environments – Geomorphological and Pedological Aspects. Catena Supplement, 14, pp. 59-73.

López-Bermúdez, F., Romero-Díaz, M. A., Martínez-Fernandez, J. and Martínez-Fernandez, J. (1998) Vegetation and soil erosion under a semi-arid Mediterranean climate: a case study from Murcia (Spain). *Geomorphology*, 24, pp. 51-58.

Loveland, P. J., Hazelden, J. and Sturdy, R. G. (1987) Chemical properties of saltaffected soils in north Kent and their relationship to soil instability. *Journal of Agricultural Science*, **109**, pp. 1-6.

Lu, D., Li, G., Valladares, G. S. and Batistella, M. (2004) Mapping soil erosion risk in Rhondonia, Brazilian Amazonia: Using RUSLE, remote sensing and GIS. Land Degradation and Development, 15, pp. 499-512.

Lunetta, R. S., Congalton, R. G., Fenstermaker, L. K., Jensen, J. R., McGwire, J. R. and Tinney, L. R. (1991) Remote sensing and geographical information system data integration: error sources and research issues. *Photogrammetric Engineering and Remote Sensing*, 57, pp. 677-687.

Luoto, M. and Hjort, J. (In Press) Evaluation of current statistical approaches for predictive geomorphological mapping. *Geomorphology*.

MacMillan, R. A., Jones, K. R. and McNabb, D. H. (2004) Defining a hierarchy of spatial entities for environmental analysis and modeling using digital elevation models (DEMs). *Computers, Environment and Urban Systems*, **28**, pp. 175-200.

Mahiny, A. S. and Turner, B. J. (2003) Modeling Past Vegetation Change Through Remote Sensing and G.I.S: A comparison of Neural Networks and Logistic Regression Methods. 7<sup>th</sup> International Conference on Geocomputation, University of Southampton.

Maier, H. R. and Dandy, G. C. (1998) The effect of internal parameters and geometry on the performance of back-propagation neural networks: an empirical study. *Environmental Modelling and Software*, 13, pp. 193-209.

Maier, H. R. and Dandy, G. C. (2001) Neural network based modelling of environmental variables: A systematic approach. *Mathematical and Computer Modelling*, 33, pp. 669-682.

Mak, B. and Munakata, T. (2002) Rule extraction of expert heuristics: A comparative study of rough sets with neural networks and ID3. *European Journal of Operational Research*, 136, pp. 212-229.

Malhotra, R. and Malhotra, D. K. (2003) Evaluating consumer loans using neural networks. Omega, 31, pp. 83-96.

Mamedov, A. I., Shainberg, I. and Levy, G. J. (2002) Wetting rate and sodicity effects on interill erosion from semi-arid Israeli soils. *Soil and Tillage Research*, **68**, pp. 121-132.

Manel, S., Dias, J. M. and Ormerod, S. J. (1999) Comparing discriminant analysis, neural networks and logistic regression for predicting species distributions: a case study with a Himalayan river bird. *Ecological Modelling*, **120**, pp. 337-347.

Martínez-Casasnovas, J. A., Ramos, M. C. and Poesen, J. (2004) Assessment of sidewall erosion in large gullies using multi-temporal DEMs and logistic regression analysis. *Geomorphology*, **58**, pp. 305-321.

Marx, E. S., Hart, J. and Stevens, R. G. (1999) Soil Test Interpretation Guide. EC 1478. Oregon State University, Corvallis.

Mather, A. (1993) Basin inversion: some consequences for drainage evolution and alluvial architecture. *Sedimentology*, 40, pp. 1069-1089.

Mather, A. (2000a) Impact of headwater river capture on alluvial system development: an example from the Plio-Pleistocene of the Sorbas Basin, SE Spain. *Journal of the Geological Society*, **157**, pp. 957-966.

Mather, A. (2000b) Adjustment of a drainage network to capture induced base-level change: an example from the Sorbas Basin, SE Spain. *Geomorphology*, **34**, pp. 271-289.

Mather, A. and Stokes, M. (2001) Marine to continental transition. In: Mather, A., Martín, J. M., Harvey, A. M. and Braga, J. C. (Eds.) A Field Guide to the Neogene Sedimentary Basins of the Almería Province, South-East Spain. Blackwell Science, Oxford. pp. 186-224.

Mather, A. and Stokes, M. (Eds.) (1996) Second Cortijo Urra Field Meeting Southeast Spain: Field Guide. University of Plymouth.

Mather, A. and Westhead, R. K. (1993) Plio/Quaternary strain of the Sorbas Basin, SE Spain: Evidence from sediment deformation structures. *Quaternary Proceedings*, 3, pp. 57-65.

Mather, A., Martín, J. M., Harvey, A. M. and Braga, J. C. (2001a) Introduction to the field guide. In: Mather, A., Martín, J. M., Harvey, A. M. and Braga, J. C. (Eds.) A Field Guide to the Neogene Sedimentary Basins of the Almería Province, South-East Spain. Blackwell Science, Oxford. pp. 1-8.

Mather, A., Martín, J. M., Harvey, A. M. and Braga, J. C. (2001b) Introduction to the Neogene geology of the Sorbas Basin. In: Mather, A., Martín, J. M., Harvey, A. M. and Braga, J. C. (Eds.) *A Field Guide to the Neogene Sedimentary Basins of the Almería Province, South-East Spain*. Blackwell Science, Oxford. pp. 9-28.

Mather, P. M. (1999) Computer Processing of Remotely-Sensed Images. Wiley, Chichestser.

Mather, P. M. (2001) Computer Processing of Remotely Sensed Images: An Introduction. (Second Edition). Wiley, Chichester.

McBratney, A. B., Mendonça Santos, M. L. and Minasny, B. (2003) On digital soil mapping. *Geoderma*, 117, pp.3-52.

McBride, M. B. (1994) Environmental Chemistry of Soils. Oxford university Press, Oxford.

McClelland, J. L. and Rumelhart, D. E. (1988) *Explorations in Parallel Distributed Processing*. The MIT Press, Cambridge.

McCord Nelson, M. and Illingworth, W. T. (1991) A Practical Guide to Neural Nets. Addison-Wesley, Massachusetts.

McCulloch, W. S. and Pitts, W. (1943) A logical calculus of the ideas immanent in nervous activity. *Bulletin of Mathematical Biophysics*, 5, pp. 115-133.

Mcgregor, D. (1957) Some observations on the geographical significance of slope. *Geography*, **42**, pp. 167-173.

McLachlan, G. J. (1992) Discriminant Analysis and Statistical Pattern Recognition. Wiley, Chichester.

McMillan, C., Mozer, M. C. and Smolensky, P. (1991) Rule induction through integrated symbolic and subsymbolic processing. In: Moody, J.E., Hanson, S. J. and Lippmann, R. P. (Eds.) Advances in Neural Information Processing Systems (Volume 4). Morgan Kaufmann, California. pp. 969-976.

Medina, F. I. and Vasquez, R. (1991) Use of simulated neural networks form aerial image classification. *ACSM-ASPRS Annual Convention*. ACSM and ASPRS, Bethesda, Maryland, pp. 268-274.

Metternicht, G. I. and Zinck, J. A. (1998) Evaluating the information content of JERS-1 SAR and Landsat TM data for the discrimination of soil erosion features. *ISPRS International Journal of Photogrammetry and Remote Sensing*, **53**, pp. 143-153.

Middleton, N. (1999) The Global Casino: An Introduction to Environmental Issues. Arnold, London.

Millward, A.A. and Mersey, J. E. (1999) Adapting the RUSLE to model soil erosion potential in a mountainous tropical watershed. *Catena*, **38**, pp. 109-129.

Minsky, M. (1991) Logical versus analogical or symbolic versus connectionist or neat versus scruffy. *AI Magazine*, Summer, pp. 35-51.

Mitchell, T. M. (1997) Machine Learning. McGraw-Hill, New York.

Moatar, F., Fessant, F. and Poirel, A. (1999) pH modelling by neural networks. Application of control and validation data series in the Middle Loire river. *Ecological Modelling*, **120**, pp. 141-156.

Molina-Aiz, F. D., Valera, D. L. and Álvarez, A. J. (2004) Measurement and simulation of climate inside Almería-type greenhouses using computational fluid dynamics. *Agricultural and Forest Meteorology*, **125**, pp. 33-51.

Mongkolsawat, C., Thirangoon, P. and Sriwongsa, S. (1994) Soil erosion mapping with universal soil loss equation and GIS. http://www.gisdevelopment.net/aars/acrs/1994/ts3/ts3001pf.htm (Accessed 26/2/05).

Moore, D. M., Lees, B. G. and Davey, S. M. (1991) A new method for predicting vegetation distributions using decision tree analysis in a geographic information system. *Environmental Management*, 15, pp. 59-71.

Morgan, R. P. C. (1978) Field studies of rainsplash erosion. Earth Surface Processes, 3, pp. 295-299.

Muchoney, D. M. and Strahler, A. H. (2002) Pixel- and site-based calibration and validation methods for evaluating supervised classification of remotely sensed data. *Remote Sensing of Environment*, **81**, pp. 290-299.

Murphy, B (1995) Relationship between the Emerson aggregate test and exchangeable sodium percentage in some subsoils from central New South Wales. In: Naidu, R., Sumner, M. E. and Rengasamy, P. (Eds.) Australian Sodic Soils: Distribution, Properties and Management. CSIRO, Adelaide. pp. 101-105.

Mzezewa, J., Gotosa, J. and Nyamwanza, B. (2003) Characterisation of a sodic soil catena for reclamation and improvement strategies. *Geoderma*, **113**, pp. 161-175.

Naidu, R., Sumner, M. E. and Rengasamy (Eds.) (1995) Australian Sodic Soils: Distribution, Properties and Management. CSIRO, Melbourne.

Nearing, M. A., Bradford, J. M. and Parker, S. C. (1991) Soil detachment by shallow flow at low slopes. *Soil Science Society of America Journal*, **55**, pp. 1532-1536.

Nearing, M. A., Foster, G. R., Lane, L. J. and Finkner, S. C. (1989) A process-based soil erosion model for USDA-water erosion prediction project technology. *Transactions of the American Society of Agricultural Engineering*, **32**, pp. 1587-1593.

Northcote, K. H. and Skene, J. K. M. (1972) Australian Soils with Saline and Sodic Properties. CSIRO Australian Soil Publication. no. 27.

Nuttall, J. G., Armstrong, R. D., Connor, D. J. and Matassa, V. J. (2003) Interrelationships between edaphic factors potentially limiting cereal growth on alkaline soils in north-western Victoria. *Australian Journal of Soil Research*, **41**, pp. 277-292.

Oades, J. M. (1984) Soil organic matter and structural stability: mechanisms and implications for management. *Plant Soil*, **76**, pp. 319-367.

Olden, J. D. and Jackson, D. A. (2002) Illuminating the "black-box": a randomization approach for understanding variable contributions in artificial neural networks. *Ecological Modelling*, **154**, pp. 135-150.

Openshaw, S. and Openshaw, C. (1997) Artificial Intelligence in Geography. John Wiley, Chichester.

Orgaz, F., Fernández, M. D., Bonachela, S., Gallardo, M. and Fereres, E. (2005) Evapotranspiration of horticultural crops in an unheated plastic greenhouse. *Agricultural Water Management*, **72**, pp. 81-96.

Øygarden, L. (2003) Rill and gully development during an extreme winter runoff event in Norway. *Catena*, **50**, pp. 217-242.

Pal, M. and Mather, P. M. (2003) An assessment of the effectiveness of decision tree methods for land cover classification. *Remote Sensing of Environment*, **86**, pp. 554-565.

Pallaris, K. (2000) Terrain modelling for risk assessment in the Cabuyal River catchment; Comparison of results with farmer perceptions. Advances in Environmental Monitoring and Modelling, 1, pp. 149-177.

Pao, Y. H. (1988) Adaptive Pattern Recognition and Neural Networks. Addison-Wesley, New York.

Paola, J. D. and Schowengerdt, R. A. (1993) A review and analysis of neural networks for classification of remotely sensed multi-spectral imagery. Research institute for Advanced Computer Science, NASA Ames Research Centre, Tech. Rep. 93.05 (NASA-CR-194291).

Paola, J. D. and Schowengerdt, R. A. (1995) A detailed comparison of backpropagation neural network and maximum-likelihood classifiers for urban land use classification. *IEEE Transactions on Geoscience and Remote Sensing*, 33, pp. 981-996.

Park, Y. S. and Chung, Y. J. (In Press) Hazard rating of pine trees from a forest insect pest using artificial neural networks. *Forest Ecology and Management*.

Parker, D. B. (1985) Learning logic. Technical report TR-47. Centre for Computational Research in Economics and Management Science. Cambridge, MA.

Parker, G. G. and Higgins, C. G. (1990) Piping and pseudokarstin drylands. In: Higgins, C. G. and Coates, D. R. (Eds.) Groundwater Geomorphology: The Role of Subsurface Water in Earth Surface Processes and Landforms. *Geological Society of America Special Paper*, 252, pp. 77-110. Parker, G. G. and Jenne, E. A. (1967) Structural failure of western US highways caused by piping. US Geological Survey Water Resources Division, pp. 27.

Pastor-Bárcenas, O., Soria-Olivas, Martín-Guerrero, J. D., Camps-Valls, G., Carrasco-Rodríguez, J. L. and del Valle-Tascón, S. (2005) Unbiased sensitivity analysis and pruning techniques in neural networks for surface ozone modelling. *Ecological Modelling*, **182**, pp. 149-158.

Patterson, A. and Niblett, T. (1983) ACLS User Manual. Intelligent Terminals Ltd, Glasgow.

Pearce, J. and Ferrier, S. (2000) Evaluating the predictive performance of habitat models developed using logistic regression. *Ecological Modelling*, 133, pp. 225-245.

Pedersen, H. S. and Hasholt, B. (1995) Influence of wind speed on rainsplash erosion. Catena, 24, pp. 39-54.

Perry, A. (1997) Mediterranean climate. In: King, R., Proudfoot, L. and Smith, B. (Eds.) *The Mediterranean: Environment and Society*. Arnold, London.

Pickup, G. and Nelson, D. J. (1984) Use of Landsat radiance parameters to distinguish soil erosion, stability, and deposition in arid Central Australia. *Remote Sensing of Environment*, 16, pp. 195-209.

Picton, P. (2000) Neural Networks. Palgrave, Hampshire.

Piper, J. (1992) Variability and bias in experimentally measured classifier error rates. *Pattern Recognition Letters*, 13, pp. 685-692.

Plumb, A. P., Rowe, R. C., York, P. and Doherty, C. (2002) The effect of experimental design on the modelling of a tablet coating formulation using artificial neural networks. *European Journal of Pharmaceutical Science*, 16, pp. 281-288.

Pomorski, D. and Perche, P. B. (2001) Inductive learning of decision trees: application to fault isolation of an induction motor. *Engineering Application of Machine Learning*, 14, pp. 155-166.

Powell, B., Ahern, C. R. and Baker, D. E. (1995) Soil dispersion in Queensland sodic soils and their classification. In: Naidu, R., Sumner, M. E. and Rengasamy, P. (Eds.) *Australian Sodic Soils: Distribution, Properties and Management*. CSIRO, Adelaide. pp. 81-87.

Qadir, M. and Schubert, S. (2002) Degradation processes and nutrient constraints in sodic soils. Land Degradation and Development, 13, pp. 275-294.

Quinlan, J. R. (1986) Induction of decision trees. Machine Learning, 1, pp. 81-106.

Quinlan, J. R. (1987) Simplifying decision trees. International Journal of Man-Machine Studies, 27, pp. 221-234. Quinlan, J. R. (1993) C4.5: Programs for Machine Learning. Morgan Kauffman, Los Altos, CA.

Quinlan, J. R. (1996) Improved use of continuous attributes in C4.5. Journal of Artificial Intelligence, 4, pp. 77-90.

Quirk, J. P. and Schofield, R. K. (1955) The effect of electrolyte concentration on soil permeability. *Journal of Soil Science*, **6**, pp. 163-178.

Rafaelli, S. G., Montgomery, D. R. and Greenberg, H. M. (2001) A comparison of thematic mapping of erosion intensity to GIS-driven process models in an Andean drainage basin. *Journal of Hydrology*, 244, pp. 33-42.

Refenes, A. N., Zapranis, A. and Francis, G. (1994) Stock performance modelling using neural networks: A comparative study with regression models. *Neural Networks*, 7, pp. 375-388.

Rengasamy, P. and Bourne, J. (2001) *Managing Sodic, Acidic and Saline Soils*. Cooperative Research Centre for Soil and Land Management. Victoria, Austalia.

Rengasamy, P., Greene, R. S. B., Ford, G. W. and Mehanni, A. H. (1984) Identification of dispersive behaviour and the management of red-brown earths. *Australian Journal of Soil Research*, 22, pp. 413-431.

Reuter, H. I., Wendroth, O. and Kersbaum, K. C. (2006) Optimisation of relief classification for different levels of generalisation. *Geomorphology*, 77, pp. 79-89.

Richards, L. A. (Ed.) (1954) *Diagnosis and Improvement of Saline and Alkaline Soils*. USDA Handbook 60, Washington DC.

Ripley, B. D. (1996) Pattern Recognition and Neural Networks. Cambridge University Press, Cambridge.

Ritter, D. F., Kochel, R. C. and Miller, J. R. (1999) The disruption of Grassy Creek: implications concerning catastrophic events and thresholds. *Geomorphology*, **29**, pp. 323-338.

Robbins, H. and Monro, S. (1951) A stochastic approximation method. Annals of Mathematical Statistics, 22, pp. 400-407.

Rodríguez-Fernández, J. L. (1999) Ockham's razor. Endeavour, 23, pp. 121-125.

Rogan, J., Franklin, J. and Roberts, D. A. (2002) A comparison of methods for monitoring multitemporal vegetation change using Thematic Mapper imagery. *Remote Sensing of Environment*, **80**, pp. 143-156.

Rojo, L. (1990) Plan Nacional de Restauración Hidrológico-Forestal y Control de la Erosión, Tomo I. Memoria. Tomo II. Mapas. ICONA, Madrid.

Rosewell C. J., Crouch, R. J., Morse, R. J., Leys, J. F., Hicks, R. W. and Stanley, R. J. (1991) Forms of Erosion. In: Charman, P. E. V. and Murphy, B. W. (Eds.) Soils - Their Properties and Management: A Soil Conservation Handbook for New South Wales. Sydney University Press, Sydney. pp. 12-35.

Rowell, D. L. (1994) Soil Science: Methods and Applications. Longman, Harlow.

Rumelhart, D. E., Hinton, G. E. and Williams, R. J. (1986) Learning representations by back-propagating errors. *Nature*, **323**, pp. 533-536.

Russell, S. and Norvig, P. (1995) Artificial Intelligence: A Modern Approach. Prentice Hall, London.

Rycroft, D. W., Kyei-Baffour, N. and Tanton, T. (2002) The effect of sodicity on the strength of a soil surface. *Irrigation and Drainage*, **51**, pp. 339-346.

Safavian, R. S. and Landgrebe, D. (1991) A survey of decision tree classifier methodology. *IEEE Transactions on Systems Machines and Cybernetics*, 21, pp. 660-674.

Salford Systems (2004) CART 5.0 Decision Tree Software. California.

Salomon, R. and Leo van Hemmen, J. (1996) Accelerating backpropogation through dynamic self-adaption. *Neural Networks*, 9, pp. 589-601.

Samra, J. S., Sharma, K. N. S. and Tyaki, N. K. (1988) Analysis of spatial variability in sodic soils. I. Structural analysis. *Soil Science*, **145**, pp. 180-187.

Sánchez, M. S., Swierenga, H., Sarabia, L. A., Derks, E. and Buydens, L. (1996) Performance of multi layer feedforward and radial base functional neural networks in classifications and modelling. *Chemometrics and Intelligent Laboratory Systems*, 33, pp. 101-119.

Savabi, M. R., Flanagan, D. C., Hebel, B. and Engel, B. A. (1995) Application of WEPP and GIS-GRASS to a small watershed in Indiana. *Journal of Soil and Water Conservation*, 50, pp. 477-483.

Saynor, M. J., Erskine, W. D. and Evans, K. G. (2003) Bank erosion in the Ngarradj catchment: Results of erosion pin measurements between 1998 and 2001. *Supervising Scientist Report*, **176**. Supervising Scientist, Darwin.

Schmidt J. V. Werner, M. and Michael, A. (1999) Application of the EROSION 3D Model to the Catsop Watershed, The Netherlands. In: De Roo, A. (Ed.) Modelling Soil Erosion by Water at the Catchment Scale. *Catena* 418.

Schumm, S. A. and Lichty, R. W. (1965) Time, space and causality in geomorphology. *American Journal of Science*, 263, pp. 110-119.

Schumm, S. A. (1973) Geomorphic thresholds and complex response of drainage systems. In: Morisawa (Ed.) *Fluvial Geomorphology*. George Allen and Unwin, London. pp. 299-310.

Schumm, S. A. (1991) To Interpret the Earth: Ten Ways to be Wrong. Cambridge University Press, Cambridge.

Scotney, P., Burgess, R. and Rutter, E. H. (2000) 40Ar/39Ar age of the Cabo de Gata volcanic series and displacements on the Carboneras fault zone, SE Spain. *Journal of the Geological Society*, **157**, pp. 1003-1008.

Sebastiá, M., Fernández, Olmo, I. and Irabien, A. (2003) Neural network prediction of unconfined compressive strength of coal fly ash – cement mixtures. *Cement and Concrete Research*, 33, pp. 1137-1146.

Servenay, A. and Prat, C. (2003) Erosion extension of indurated volcanic soils of Mexico by aerial photographs and remote sensing analysis. *Geoderma*, 117, pp. 367-375.

Sethi, I. K. (1990) Entropy nets: from decision trees to neural networks. *Proceedings* of the IEEE, 78, pp. 1605-1613.

Shainberg, I. (1990) Soil response to saline and sodic conditions. In: Tanji, K. K. (Ed.) Agricultural Salinity Assessment and Management. American Society of Civil Engineering, New York. pp. 91-112.

Shainberg, I., Rhoades, J. D. and Prather, R. J. (1981) Effect of low electrolyte concentration on clay dispersion and hydraulic conductivity of a sodic soils. *Soil Science Society of America Journal*, **45**, pp. 273-277.

Shi, Z. H., Cai, C. F., Ding, S. W., Wang, T. W. and Chow, T. L. (2004) Soil conservation planning at the small watershed level using RUSLE with GIS: a case study in the Three Gorge Area of China. *Catena*, **55**, pp. 33-48.

Short, N. (1991) A real-time expert system and neural network for the classification of remotely sensed data. *ACSM-ASPRS Annual Convention*. ACSM and ASPRS, Bethesda, Maryland, pp. 406-418.

Shrestha, D. P., Zinck, J. A. and Van Ranst, E. (2004) Modelling land degradation in the Nepalese Himalaya. *Catena*, 57, pp. 135-156.

Shrimali, S. S., Aggarwal, S. P. and Samra, J. S. (2001) Prioritising erosion-prone areas in hills using remote sensing and GIS – a case study of the Sukhna Lake catchment, Northern India. *International Journal of Applied Earth Observation and Geoinformation*, **3**, pp. 54-60.

Singer, M. J. and Le Bissonnais, Y. (1998) Importance of surface sealing in the erosion of some soils from a Mediterranean climate. *Geomorphology*, 24, pp. 79-85.

Singh, D., Meirelles, M. S. P., Costa, G. A., Herlin, I., Berrior, J. P. and Silva, E. F. (In Press) Environmental degradation analysis using NOAA/AVHRR data. Advances in Space Research.

Sirvent, J., Desir, G., Gutierrez, M. Sancho, C. and Benito, G. (1997) Erosion rates in badland areas recorded by collectors, erosion pins and profilometer techniques (Ebro Basin, NE Spain). *Geomorphology*, 18, pp. 61-75.

Skidmore, A. K., Turner, B. J., Brinkhof, W. and Knowles, E. (1997) Performance of a neural network: mapping forests using GIS and remotely sensed data. *Photogrammetric Engineering and Remote Sensing*, **63**, pp. 501-514.

Smith, A. G. and Woodcock, N. H. (1982) Tectonic syntheses of the Alpine-Mediterranean region: a review. In: Berckhemer, H. and Hsu, K. (Eds.) Alpine-Mediterranean Geodynamics. *American. Geophysical. Union*, Geodynamics Series, 7, pp. 15-38.

Smith, J. A. (1993) LAI inversion using a back-propagation neural network trained with a multiple scattering model. *IEEE Transactions on Geoscience and Remote Sensing*, 31, pp. 1102-1106.

Smith, K. (2001) Environmental Hazards: Assessing Risk and Reducing Disaster. Routeledge, London.

So, H. B. and Woodhead, T. (1987) Alleviation of soil physical limits to productivity of legumes in Asia. In: Wallis, E. S. and Blythe, D. E. (Eds.) *Food Legume Improvement of Asian Farming Systems*. ACIAR proceedings No. 18, pp. 112-120.

Soto, B. and Díaz-Fierros, F. (1998) Runoff and soil erosion from areas burnt scrub: comparison of experimental results with those predicted by the WEPP model. *Catena*, **31**, pp. 257-270.

Sparks, D. L. (1995) Environmental Soil Chemistry. Academic Press, London.

Spellman, G. (1999) An application of artificial neural networks to prediction of surface ozone concentrations in the United Kingdom. *Applied Geography*, **19**, pp. 123-136.

Spivey, D. (1997) Scale, process and badland development in Almería province. Unpublished PhD thesis. University of Liverpool.

Stiegeler, S. E. (1979) A Dictionary of Earth Sciences. Pan Books, London.

Strahler, A. N. (1958) Dimensional analysis of watershed geomorphology. *Geological Society of America Bulletin*, **60**, pp. 279-299.

Summerfield, M. A. (1991) Global Geomorphology. Pearson Education, Harlow.

Sumner, M. E. (1993) Sodic soils: new perspectives. Australian Journal of Soil Research, 31, pp. 683-750.

Sumner, M. E. (1995) Sodic Soils. In: Naidu, R., Sumner, M. E. and Rengasamy, P. (Eds.) Australian Sodic Soils: Distribution, Properties and Management. CSIRO, Adelaide. pp. 1-34.

Swain, P. H. (1978) Fundamentals of pattern recognition. In: Swain, P. H. and Davis, S. M. (Eds.) *Remote Sensing: The Quantitative Approach*. McGraw-Hill, New York. pp. 136-187.

Swain, P. H. and Hauska, H. (1969) The decision tree classifier: design and potential. *IEEE Transactions on Geoscience and Remote Sensing*, **15**, pp. 142-147.

Thieken, A. H., Lucke, A., Diekkruger, B. and Richter, O. (1999) Scaling input data by GIS for hydrological modeling. *Hydrological Processes*, **13**, pp. 611-630.

Thomas, M. F. (2001) Landscape sensitivity in time and space – an introduction. *Catena*, 42, pp. 83-98.

Thompson, J. A., Bell, J. C. and Butler, C. A. (2001) Digital elevation model resolution: effects on terrain attribute calculation and quantitative soil-landscape modeling. *Geoderma*, 100, pp. 67-89.

Thornes, J. B. (1990) The interaction of erosional and vegetational dynamics in land degradation: spatial outcomes. In: Thornes, J. B. (Ed.) Vegetation and Erosion: Processes and Environments. Wiley, Chichester. pp. 41-53.

Thornes, J. B. (1996) Introduction. In: Brandt, J. C. and Thornes, J. B. (Eds.) *Mediterranean Desertification and Land Use*. pp. 1-11.

Tiira, T. (1999) Detecting teleseismic events using artificial neural networks. Computers and Geoscience, 25, pp. 929-938.

Tiscareno-Lopez, M., Weltz, M. A. and Lopes, V. L. (1995) Assessing uncertainties in WEPP's soil erosion predictions on rangelands. *Journal of Soil Water and Conservation*, 50, pp. 512-516.

Torri, D. and Bryan, R. (1997) Micropiping processes in biancana evolution in southeast Tuscany, Italy. *Geomorphology*, 20, pp. 219-235.

Torri, D., Colica, A. and Rockwell, D. (1994) Preliminary study of the erosion mechanisms in a biancana badland (Tuscany, Italy). *Catena*, 23, pp. 281-294.

Tóth, T., Csillag, F., Biehl, L. L. and Michéli, E. (1991) Characterization of semivegetated salt-affected soils by means of field remote sensing. *Remote Sensing of Environment*, 37, pp.167-180.

Tourenq, C., Aulagnier, S., Mesléard, F., Durieux, L., Johnson, A., Gonzalez, G. and Lek, S. (1999) Use of artificial neural networks for predicting rice crop damage by greater flamingos in the Camargue, France. *Ecological Modelling*, **120**, pp. 349-358.

Tout, D. G. (1987) South-East Almería province, Spain – The driest region in Europe. *Weather*, **42**, pp. 242-247.

Trajan Neural Network Simulator Software (1999) User Manual for Version 4.0. Protaprint, Durham.

Tseng-Chung, T. and Li-Chiu, C. (2005) Predicting multilateral trade credit risks: comparisons of Logit and Fuzzy Logic models using ROC curve analysis. *Expert* Systems with Applications, 28, pp. 547-556.

Tu, J. V. (1996) Advantages and disadvantages of using artificial neural networks versus logistic regression for predicting medical outcomes. *Journal of Clinical Epidemiology*, **49**, pp. 1225-1231.

Tveter, D. R. (1998) The Pattern Recognition Basis of Artificial Intelligence. IEEE Computer Society, California.

United States Department of Agriculture (1978) Agricultural Research Service (ARS). Wischmeir, W. H. and Smith, D. D. *Predicting Rainfall Erosion Losses*. Handbook 537.

United States Department of Agriculture (2005) Human induced land degradation http://www.soils.usda.gov/use (Accessed 26/1/05).

Utgoff, P. E. and Brodley, C. E. (1990) An incremental method of finding multivariate splits for decision trees. *Machine Learning, Proceedings of the Seventh International Conference on Machine Learning*. pp. 58-65. Morgan Kaufmann, Austin, TX.

van Genderen, J. L. and Lock, B. F. (1977) Testing land-use map accuracy. *Photogrammetric Engineering and Remote Sensing*, **43**, pp. 1135-1137.

Viseras, C. Calvache, M. L., Soria, J. M. and Fernández, J. (2003) Differential features of alluvial fans controlled by tectonic or eustatic accommodation space. Examples from the Betic Cordillera, Spain. *Geomorphology*, **50**, pp. 181-202.

Vrieling, A., Sterk, G. and Beaulieu, N. (2002) Erosion risk mapping: A methodological case study in the Colombian Eastern Plains. *Journal of Soil and Water Conservation*, 57, pp. 158-163.

Walker, D. J. H. (1997) Dispersive Soils in KwaZulu-Natal. MSc Thesis, University of Natal, Durban, South Africa.

Wang, Z., Di Massimo, C., Tham, M. T. and Morris, A. J. (1994) A procedure for determining the topology of multilayer feedforward neural networks. *Neural Networks*, 7, pp. 291-300.

Ward, R. C. and Robinson, M. (2000) Principles of Hydrology (Fourth Edition). McGraw-Hill, Maidenhead.

Warner, B. and Misra, M. (1996) Understanding neural networks as statistical tools. *The American Statistician*, **50**, pp. 284-293.

Weijermars, R. (1991) Geology and tectonics of the Betic Zone, SE Spain. *Earth Science Reviews*, **31**, pp. 153-184.

Welsh, A. H., Cunningham, R. B., Donelly, C. F. and Lindenmayer, D. B. (1996) Modelling the abundance of rare species: statistical models for counts with extra zeros. *Ecological Modelling*, **88**, pp. 297-308.

Werbos, P. J. (1994) The Roots of Backpropagation. Wiley, New York.

Wheeler, D. (1996) Spanish climate: regions and diversity. *Geography Review*, 10, pp. 34-40.

Wieczorkowska, A. A. (2000) Application of decision trees to wavelet-based classification of musical instrument sounds. In: Kłopotek, M., Michalewicz, M. and Wierzchoń (Eds.) Advances in Soft Computing. Physica-Verlag, New York. pp.45-53.

Wigley, T. M. L. (1992) Future climate of the Mediterranean basin with particular emphasis on changes in precipitation. In: Jeftic, L., Milliman, J. D. and Sestini, G. (Eds.) *Climatic Change in the Mediterranean*. Arnold, London. pp. 15-44.

Wolock, D. M. and Price, C. V. (1994) Effects of digital elevation model map scale and data resolution on a topography based watershed model. *Water Resources Research*, **30**, pp. 3041-3052.

Wright, A. C. and Webster, R. (1991) A stochastic distributed model of soil erosion by overland flow. *Earth Surface Processes and Landforms*, 16, pp. 207-226.

Wright, J. A. (2003) Environmental Chemistry. Routledge, London.

Wu, C. H. (1997) Artificial neural networks for molecular sequence analysis. Computers and Chemistry, 21, pp. 237-256.

Yang, C. C., Prasher, S. O., Enright, P., Madramootoo, C., Burgess, M., Goel, P. K. and Callum, I. (2003) Application of decision tree technology for image classification using remote sensing data. *Agricultural Systems*, **76**, pp. 1101-1117.

Yesilnacar, E. and Topal, T. (2005) Landslide susceptibility mapping: A comparison of logistic regression and neural networks methods in a medium scale study, Hendek region (Turkey). *Engineering Geology*, **79**, pp. 251-266.

Young, R. A. and Mutchler, C. K. (1969) Effect of slope shape on erosion and runoff. *Trans. ASAE*, 12, pp. 231-239.

Zhang, B., Valentine, I. and Kemp, P. (2005a) Modelling the productivity of naturalised pasture in the North Island, New Zealand: a decision tree approach. *Ecological Modelling*, In Press.

Zhang, B., Valentine, I., Kemp, P. and Lambert, G. (2005b) Predictive modelling of hill-pasture productivity: integration of a decision tree and a geographic information system. *Agricultural Systems*, In Press.

Zhang, L., O'Neill, A. L. and Lacey, S. (1996a) Modelling approaches to the prediction of soil erosion in catchments. *Environmental Software*, **11**, pp. 123-133.

Zhang, X. C., Nearing, M. A., Risse, L. M. and McGregor, K. C. (1996b) Evaluation of WEPP runoff and soil loss predictions using natural runoff plot data. *Trans. ASAE*, **39**, pp. 855-863.

Zhou, Q. and Liu, X. (2004) Analysis of errors of derived slope and aspect related to DEM data properties. *Computers and Geosciences*, **30**, pp. 369-378.

Zukowskyj, P., Alexander, R., Teeuw, R., Faulkner, H., Sullivan, M. and Ruiz, J. (2005) Changing vegetation density and a comparison with agricultural clearances in Sorbas, south east Spain. http://www.nerc.ac/arsf/contents/zukowskyj\_p\_et\_al.doc. (Accessed 29/1/2005).

Zweig, M. H. and Campbell, G. (1993) Receiver-operating characteristic (ROC) plots: a fundamental evaluation tool in clinical medicine. *Clinical Medicine*, **39**, pp. 561-577.

# Appendix 1 – Detailed Laboratory Analysis

Sample	CEC Sample Weight (g)	Deflection Reading	Calibration Curve Equation	Calculation	Per Kg Soil	CEC cmol/kg
1	4.0074	27	y=2.1143x+1	12.3	3068.6	13.342
2	4.0018	23	y=2.3714x+0.1429	9.64	2408.6	10.472
3	4.0099	20	y=2.1143x+1	8.99	2241.1	9.744
4	4.0006	21	y=2.62x+0.6	7.79	1946.3	8.462
5	4.0083	29	y=2.1257x+1.0952	13.1	3275	14.239
6	4.0083	30	y=2.3714x+0.1429	12.6	3141.1	13.657
7	3.9993	24	y=2.62x+0.6	8.93	2233.2	9.710
8	4.0139	26	y=2.1257x+1.0952	11.7	2918.9	12.691
9	3.9932	18	y=2.1257x+1.0952	7.95	1991.5	8.659
10	4.0036	15	y=2.1257x+1.0952	6.54	1633.8	7.104
11	4.003	17	y=2.3714x+0.1429	7.11	1775.8	7.721
12	4.018	35	y=2.1143x+1	16.1	4002.2	17.401
13	4.0021	29	y=2.1143x+1	13.2	3309.1	14.387
14	4.0015	24	y=2.62x+0.6	8.93	2232	9.704
15	4.007	24	y=2.62x+0.6	8.93	2228.9	9.691
16	4.0074	28	y=2.1143x+1	12.8	3186.7	13.855
17	4.0059	17	y=2.3714x+0.1429	7.11	1774.5	7.715
18	4.0062	24	y=2.62x+0.6	8.93	2229.4	9.693
19	4.0049	17	y=2.62x+0.6	6.26	1563	6.796
20	4.0026	50	y=2.3714x+0.1429	21	5252.7	22.838
21	4.0071	28	y=2.1257x+1.0952	12.7	3158.6	13.733
22	4.001	31	y=2.1257x+1.0952	14.1	3516.2	15,288
23	4.0057	19	y=2.62x+0.6	7.02	1753.2	7.623
24	4.0095	30	v=2.3714x+0.1429	12.6	3140.2	13.653
25	4.0044	34	v=2.1257x+1.0952	15.5	3865.6	16.807
26	4.008	12	v=2.3714x+0.1429	5	1247.5	5 424
27	4,0057	17	v=2.1143x+1	7.57	1889.2	8 214
28	4.0044	21	y=2.62x+0.6	7.79	1944.4	8 454
29	4.0088	33	v=2.3714x+0.1429	13.9	3456.3	15.027
30	4.0106	16	y=2.1257x+1.0952	7.01	1748.3	7 601
31	4.007	37	v=2.3714x+0.1429	15.5	3878.8	16 864
32	3 9936	27	v=2.62x+0.6	10.1	2523.1	10.970
33	4.0039	25	v=2.1257x+1.0952	11.2	2808.7	12 212
34	4.0012	24	v=2.1143x+1	10.9	2718.8	11.821
35	4.0023	22	y=2.62x+0.6	8.17	2040.8	8 873
36	4.0036	27	y=2.3714x+0.1429	11.3	2828.8	12 299
37	4 0032	17	v=2 1257x+1 0952	7.48	1869	8 126
38	3 9962	29	v=2.1143x+1	13.2	3313.9	14 408
30	4 0065	29	v=2.1257x+1.0952	13.1	3276.5	14 246
40	4 0036	27	v=2.1143x+1	123	3071.5	13 355
40	4 0023	21	v=2.62x+0.6	7.79	1945.4	8 458
41	4 0004	25	v=2.3714x+0.1429	10.5	2620.2	11 392
42	4 0095	27	v=2 1143x+1	123	3067	13 335
45	4 0036	21	v=2 3714x+0 1429	88	2196.8	9.551
44	4 0092	15	v=2.1257x+1.0952	6.54	1631.6	7 094
45	4.0002	14	v=2 62x+0 6	5.11	1278.6	5.550
40	4 0008	31	y=2 3714x+0 1429	13	3252.4	14 141
47	4 0147	50	y=2 1143x+1	23.2	5772 7	25 099
40	3 9983	43	v=2 1257x+1 0952	19.7	4930.4	21.033
49	4 004	46	v=2 1143x+1	213	5315 B	23.111
50	4.0021	35	v=2 1257x+1 0952	15.0	3095 4	17 329
52	4.0021	42	v=2.1143x+1	10.8	4825.4	20.070
52	4.0105	43	v=2.62x+0.6	16.3	4020.1	17 559
53	4.0075	23	v=2 1143x+1	10.2	4036.2	11.000
54	4.0004	35	v=2 62x+0 6	10.4	2001.1	11,309
55	4.0022	55	1-2.02210.0	13.1	3280.6	14.264

Table A1.1: Detailed Cation Exchange Capacity laboratory results.

Sample	ESP Sample Weight (g)	ESP Calibration Curve Equation	Deflection Reading	Calculation	Per Kg Soil	ESP cmol/kg	ESP cmol/kg
1	4.089	y = 2.3242x + 1.1246	40	16.72635746	409	1.7785	0.500
2	4	y = 2.3611x + 2.1451	26	10.10329931	253	1.0982	0.308
3	4,1006	y = 2.5229x + 2.0683	6	1.55840501	38	0.1652	0.046
4	4.0078	y = 2.5229x + 2.0683	10	3.143882041	78.4	0.3411	0.095
5	4	v = 2.3611x + 2.1451	10	3.326796832	83.2	0.3616	0.100
6	4 0046	y = 2.273x + 1.2628	129	56.19762429	1403	6.1014	1.690
7	4.0062	y = 2.5229x + 2.0683	44	16.6204368	415	1.8038	0.498
8	4.0002	y = 2.3611x + 2.1451	29	11.37389352	284	1.2363	0.341
9	4 0281	y = 2.5229x + 2.0683	10	3.143882041	78	0.3393	0.093
10	4 1111	y = 2.5229x + 2.0683	6	1,55840501	37.9	0.1648	0.045
11	4.0095	y = 2.5229x + 2.0683	42	15.82769828	395	1.7163	0.469
12	4 1084	y = 2.5229x + 2.0683	25	9.089420905	221	0.9619	0.262
12	4.0404	y = 2.5229x + 2.0683	18	6.314836101	156	0.6795	0.185
10	4.0404	y = 2.5229x + 2.0683	6	1 55840501	38.8	0 1689	0.046
14	4.018	y = 2.5558x + 2.1368	48	17 94475311	447	1 9418	0.525
10	4.0002	y = 2.5500x + 2.1000	2	0 707983688	17.7	0.077	0.021
10	4.0002	y = 2.5002x + 0.1040	6	1 55840501	38.3	0.1666	0.045
17	4.0074	y = 2.5220x + 2.0000	5	1 162035753	29	0.126	0.034
10	4.0090	y = 2.5229x + 2.0000	6	1 55840501	38	0.1653	0.034
19	4.0903	y = 2.0223X + 2.0003	72	31 12063352	777	3 3768	0.044
20	4.007	y = 2.273x + 1.2020	12	1 55840501	387	0.1684	0.900
21	4.0220	y = 2.5229x + 2.0005	0	25 90405459	643	0,1004	0.045
22	4.0117	y = 2.5502x + 0.1945	00	17 00054457	045	2.7900	0.742
23	4.0121	y = 2.5229x + 2.0003	47	7 502042974	444	0.9400	0.511
24	4.0164	y = 2.5229x + 2.0083	21	1.003943074	10/	0.0123	0.214
25	4.004	y = 2.3611x + 2.1451	10	3.326796832	83.1	0.3612	0.095
26	4.0552	y = 2.273x + 1.2628	4	1.204223493	29.7	0.1291	0.032
27	4.0053	y = 2.3242x + 1.1246	4	1.23/1508/1	30.9	0.1343	0.033
28	4	y = 2.5558x + 2.1368	4	0.72900853	18.2	0.0792	0.020
29	4.0676	y = 2.3611x + 2.1451	24	9.2562365	228	0.9894	0.245
30	4	y = 2.5558x + 2.1368	5	1.120275452	28	0.1218	0.030
31	4.0039	y = 2.3611x + 2.1451	51	20.69158443	51/	2.2469	0.553
32	4.0215	y = 2.5502x + 0.1945	5	1.884362011	46.9	0.2037	0.050
33	4.0132	y = 2.3242x + 1.1246	3	0.806901299	20.1	0.0874	0.021
34	4.0017	y = 2.3611x + 2.1451	4	0.785608403	19.6	0.0854	0.021
35	4.0014	y = 2.3611x + 2.1451	13	4.597391047	115	0.4995	0.121
36	4.0789	y = 2.5502x + 0.1945	13	5.021370873	123	0.5352	0.130
37	4.0018	y = 2.5502x + 0.1945	3	1.100109795	27.5	0.1195	0.029
38	4.0119	y = 2.5502x + 0.1945	3	1.100109795	27.4	0.1192	0.029
39	4.0073	y = 2.3611x + 2.1451	5	1.209139808	30.2	0.1312	0.031
40	4.008	y = 2.3611x + 2.1451	4	0.785608403	19.6	0.0852	0.020
41	4.0677	y = 2.5502x + 0.1945	3	1.100109795	27	0.1176	0.028
42	4.0398	y = 2.273x + 1.2628	8	2.964012319	73.4	0.319	0.075
43	4.0339	y = 2.5502x + 0.1945	4	1.492235903	37	0.1608	0.038
44	4.0012	y = 2.5502x + 0.1945	4	1.492235903	37.3	0.1622	0.038
45	4.0048	y = 2.3611x + 2.1451	6	1.632671213	40.8	0.1773	0.042
46	4.0185	y = 2.3242x + 1.1246	3	0.806901299	20.1	0.0873	0.020
47	4.019	y = 2.5502x + 0.1945	5	1.884362011	46.9	0.2039	0.047
48	4.0436	y = 2.5502x + 0.1945	80	31.29382009	774	3.3648	0.770
49	4.0017	y = 2.5502x + 0.1945	71	27.76468512	694	3.0166	0.686
50	4.0136	y = 2.5502x + 0.1945	3	1.100109795	27.4	0.1192	0.027
51	4.0164	y = 2.5502x + 0.1945	4	1.492235903	37.2	0.1615	0.037
52	4.0031	y = 2.3611x + 2.1451	8	2.479734022	61.9	0.2693	0.061
53	4.0252	y = 2.5502x + 0.1945	50	19,53003686	485	2.1095	0.475
54	4	y = 2.5558x + 2.1368	5	1,120275452	28	0,1218	0.027
55	4.0219	y = 2.5229x + 2.0683	5	1,162035753	28.9	0,1256	0.028
00	1.0210	a support of the local data and the			2010	0.1800	0.020

Table A1.2: Detailed Exchangeable Sodium Percentage laboratory results.

Sample Number	ESR (ES/(CEC - ES))	$SAR y = -0.0126 \pm 0.01475y$
1	0 153805891	11 28175535
2	0 117153463	8 796844956
3	0.01725073	2 023778337
4	0.041997469	3.701523337
5	0.026056793	2 620799534
6	0.80754033	55.60273426
7	0.228156557	16.32247846
8	0.107931315	8 171614554
9	0.040788866	3.619584158
10	0.023752268	2.464560518
11	0.285839167	20.23316389
12	0.058513706	4.821268203
13	0.049573202	4.215132318
14	0.017712274	2.055069444
15	0.250578128	17.84258494
16	0.005585035	1,232883718
17	0.022068121	2.350381062
18	0.01317035	1 747142362
19	0.024935745	2 544796275
20	0.173515328	12 61798837
21	0.012417539	1,696104338
22	0.22388799	16.0330841
23	0.33902491	23.83897698
24	0.063261598	5.143159185
25	0.021965889	2 343450087
26	0.024384385	2 507415956
27	0.01662165	1.981128786
28	0.009461741	1.495711259
29	0.070479799	5 632528732
30	0.016280342	1,957989285
31	0.15371329	11.27547728
32	0.018922537	2.137121154
33	0.007210242	1.343067235
34	0.007273409	1.347349747
35	0.059657078	4,898784952
36	0.045498714	3.93889588
37	0.014927836	1 866293987
38	0.008343538	1.419900878
39	0.009294602	1.484379828
40	0.006422479	1.289659578
41	0.014097681	1.810012247
42	0.028807941	2.807318047
43	0.0122086	1.681938973
44	0.017269705	2.025064717
45	0.025627274	2 591679592
46	0.015954649	1.935908392
47	0.014626796	1.84588449
48	0.154820517	11.3505435
49	0.163767635	11 9571278
50	0.005183168	1 205638533
51	0.009410158	1 492214107
52	0.013004995	1 735931889
53	0.13655765	10 11238304
54	0.010884632	1 592178474
55	0.008885291	1 456629888
		1.400020000

 Table A1.3: Determination of Sodium Adsorption Ratio from Exchangeable Sodium Ratio.

# **Appendix 2 – Correlation Matrices**

**Correlation Matrices Produced for the Artificial Neural Networks** 

All and a state	- AN AND A	Actu	ıal		
		0	1	Total	Overall
Predicted	0	11	6	17	Accuracy
	1	34	79	113	
	Total	45	85	130	0.692308
	Producers A	ccuracy	Users A	Accuracy	
	0	24.4%	0	64.7%	
	1	92.9%	1	69.9%	10.25

**Two Class Classifications** 

Table A2.1: Correlation matrix for 10 metre DEM independent variables.

		Actı	ıal		
		0	1	Total	Overall
Predicted	0	0 6		8	Accuracy
	1	39	83	122	
	Total	45	85	130	0.684615
	Producers A	ccuracy	Users A	Accuracy	
	0	13.3%	0	75%	
	1	97.6%	1	68%	

Table A2.2: Correlation matrix for 20 metre DEM independent variables.

		Actı	ial	ERCO PAR	
N.S. MASIN		0	1	Total	Overall
Predicted	0	28	27	55	Accuracy
	1	17	58	75	
	Total	45	85	130	0.661538
	Producers A	ccuracy	Users A	Accuracy	
	0	62.2%	0	50.9%	
	1. Contra 1. Contra 1. Contra	74.1%	1	77.3%	

Table A2.3: Correlation matrix for field collected independent variables.

1.4.6.2.6.4.1		Actu	ıal		Selection of
- 2 - 1 - 9		0	1	Total	Overall
Predicted	0	27	14	41	Accuracy
	1	18	71	89	
	Total	45	85	130	
	Producers A	ccuracy	Users A	Accuracy	0.753846
	0	60%	0	65.9%	
	1	83.5%	1	79.8%	

Table A2.4: Correlation matrix for field collected and the 10 metre DEM independent variables.

2 breve Since	C Class No.	Actu	ıal		a de cara de la composición de la
and we had		0	1	Total	Overall
Predicted	0	30	23	53	Accuracy
	1	15	62	77	1.1.1.1.1.1.1
	Total	45	85	130	0.707692
	Producers A	ccuracy	Users A	Accuracy	
	0	66.6%	0	56.6%	
	1	72.9%	1	80.5%	

Table A2.5: Correlation matrix for field collected and 20 metre DEM independent variables.

	1. Stranger	Actu	ıal	9-3-9-5-1 (	
LEDIC DALLS	Condultation	0	1	Total	Overall
Predicted	0	34	32	66	Accuracy
	1	11	53	64	
	Total	45	85	130	0.669231
	Producers A	ccuracy	Users A	Accuracy	
	0	75.6%	0	51.5%	
	1	62.4%	1	82.8%	

Table A2.6: Correlation matrix for field collected and classified vegetation independent variables.

	And the second	Actu	ıal		Stiff Second and State
		0	1	Total	Overall
Predicted	0	28	15	43	Accuracy
Tribler and al	1	17	70	87	
in the star	Total	45	85	130	0.753846
	Producers A	ccuracy	Users A	Accuracy	
	0	62.2%	0	65.1%	
	1	82.4%	1	80.5%	

 Table A2.7: Correlation matrix for field collected, classified vegetation and 10 metre

 DEM independent variables.

	Contraction of the	Actı	ıal		Brene Lesse
		0	1	Total	Overall
Predicted	0	25	16	41	Accuracy
	1	20	69	89	14
	Total	45	85	130	0.723077
	Producers A	ccuracy	Users A	Accuracy	
	0	55.6%	0	61%	1.1.1.1
	1	81.2%	1	77.5%	

 Table A2.8: Correlation matrix for field collected, classified vegetation and 20 metre

 DEM independent variables.

#### Three Class Classifications

			Actual	1. 1. 1. 1. 1. 1.	Selle Maria	
		0	1	2	Total	Overall
Predicted	0	34	15	31	80	Accuracy
	1	0	0	0	0	
	2	11	12	27	50	
	Total	45	29	56	130	0.469231
	Producers Accuracy		Users Accuracy			
	0	75.6%	(	)	42.5%	
	1	- 11	1	1	18	
	2	48.2%	1	2	54%	

Table A2.9: Correlation matrix for 10 metre DEM independent variables.

Sectores 14		1	Actual	ne basi	Alleys up th	
		0	1	2	Total	Overall
Predicted	0	18	11	5	34	Accuracy
	1	0	0	0	0	
	2	27	16	53	96	0.546154
	Total	45	29	56	130	
	Producers	s Accuracy	1	Users Acc	uracy	1
	0	40%	0		52.9%	
	1	-	1		-	1
	2	94.6%	2		55.2%	

Table A2.10: Correlation matrix for 20 metre DEM independent variables.

			Actual		anal.	Charles and	
Stand La	2	0	1	2	Total	Overall	
Predicted	0	32	13	19	64	Accuracy	
	1	4	4	1	9		
	2	9	10	38	57	194	
	Total	45	29	56	130	0.569231	
	Producers	s Accuracy	1	Users Acc	uracy		
	0	71.1%	(	)	50%		
Sarah Sarah	1	13.8%	1		44.4%		
	2	67.9%	2	2	66.7%		

Table A2.11: Correlation matrix for field collected independent variables.

	T. S. BURKER		Actual			1.5.12 4.50	
		0	1	2	Total	Overall	
Predicted	0	34	15	19	68	Accuracy	
	1	4	4	1	9	Contraction of the	
	2	7	8	38	53		
	Total	45	29	56	130	0.584615	
101320	Producer	s Accuracy	1	Users Acc	uracy	a state of	
	0	75.6%	(	)	50%		
	1	13.8%	1		44.4%		
	2	67.9%	1	2	71.7%		

Table A2.12: Correlation matrix for field collected and the 10 metre DEM independent variables.

	A DE SUL I		Actual			
		0	1	2	Total	Overall
Predicted	0	12	2	2	16	Accuracy
	1	0	0	0	0	
Page 1	2	33	25	56	114	and such as
Section 1	Total	45	29	56	130	0.523077
	Producers	s Accuracy	1	Users Acc	euracy	
	0	26.7%	(	)	75%	Constantino de
	1		1		-	
	2	100%	2	2	49.1%	

Table A2.13: Correlation matrix for field collected and the 20 metre DEM independent variables.

AND CALLS	The second s		Actual			
		0	1	2	Total	Overall
Predicted	0	35	13	17	65	Accuracy
Contract of the second	1	1	2	0	3	
	2	9	12	41	62	0.6
	Total	45	29	56	130	
ELANDARY.	Producers	s Accuracy	1	Users Acc	uracy	
	0	77.8%	(	)	53.8%	
	1	6.9%	1		66.7%	
	2	73.2%	2	2	66.1%	

Table A2.14: Correlation matrix for field collected and classified vegetation independent variables.

		1	Actual	14.14 876.9		
		0	1	2	Total	Overall
Predicted	0	32	11	17	60	Accuracy
	1	6	4	1	11	
	2	7	12	40	59	
	Total	45	29	56	130	0.584615
and the state of	Producers	s Accuracy	1	Users Acc	uracy	
ANT HAR MARK	0	71.1%	0	)	53.3%	
	1	13.8%	1		36.4%	
Sale Mary Ing	2	71.4%	2		67.8%	

**Table A2.15:** Correlation matrix for field collected, classified vegetation and 10metre DEM independent variables.

			Actual			
		0	1	2	Total	Overall
Predicted	0	36	18	21	75	Accuracy
	1	3	2	0	5	
	2	6	7	37	50	0.576923
	Total	45	29	56	130	
	Producers	s Accuracy		Users Acc	uracy	
	0	80%	(	)	48%	
	1	6.9%	1		40%	1
58517 57-1	2	66.1%	1 2	2	74%	

 Table A2.16: Correlation matrix for field collected, classified vegetation and 20 metre DEM independent variables.

#### Nine Class Classifications

	1		1			Actua	1		61.1	6.97	1. 1. 1. 1.	And strains
19 march 1		0	1	2	3	4	5	6	7	8	Total	Overall
	0	45	17	5	7	7	18	18	9	4	130	Accuracy
	1	0	0	0	0	0	0	0	0	0	0	
	2	0	0	0	0	0	0	0	0	0	0	
	3	0	0	0	0	0	0	0	0	0	0	
Predicted	4	0	0	0	0	0	0	0	0	0	0	0.346154
	5	0	0	0	0	0	0	0	0	0	0	
	6	0	0	0	0	0	0	0	0	0	0	
	7	0	0	0	0	0	0	0	0	0	0	
	8	0	0	0	0	0	0	0	0	0	0	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	roduc	ers Ac	curac	y.							
	0			1009	6		0			34.6%	6	
	1	_		-			1		-			
	2			-			2			-		
	3			-			3			-		
	4		-	-			4					
	5			-			5			-		
	6						6			-		
	7								-			Second Second Second
	8			-			8			-		

Table A2.17: Correlation matrix for 10 metre DEM independent variables.

						Actua	1					Contractor of the
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	37	15	3	7	6	18	15	6	1	108	Accuracy
	1	0	1	0	0	0	0	0	0	0	1	
	2	0	0	0	0	0	0	0	0	0	0	
	3	0	0	0	0	0	0	0	0	0	0	
Predicted	4	0	0	0	0	0	0	0	0	0	0	0.3
	5	6	0	2	0	0	0	1	1	1	11	
	6	0	0	0	0	0	0	0	1	1	2	
	7	2	1	0	0	1	0	2	1	1	8	
	8	0	0	0	0	0	0	0	0	0	0	
	Total	45	17	5	7	7	18	18	9	4	130	
	I	roduc	ers Ac	curac	y		S MAL	User	s Accu	iracy		
	0			82.2	%		0			34.39		
	1			5.8%	6		1			100%	6	
	2			-			2			-		
	3						3			-		
	4			-			4			-		
	5	10.00		-			5			-		
	6			-			6					
	7	7		11.1%					12.5%			1000
	8			-			8					

Table A2.18: Correlation matrix for 20 metre DEM independent variables.

						Actua	1						1
		0	1	2	3	4	5	6	7	8	Total	Overall	ī
	0	32	7	2	3	2	6	3	1	0	56	Accuracy	
	1	5	3	1	1	1	0	0	0	0	11		1
6. 18 . 19	2	0	0	0	0	0	0	0	0	0	0		
	3	0	0	0	0	0	0	0	0	0	0		
Predicted	4	0	0	0	0	0	0	0	0	0	0	0.376923	
Part Part	5	1	0	0	1	1	1	2	1	0	7		
140-5 L	6	7	7	2	2	3	11	13	7	4	56		
Mar Martin	7	0	0	0	0	0	0	0	0	0	0		
12 10 10 10	8	0	0	0	0	0	0	0	0	0	0		
	Total	45	17	5	7	7	18	18	9	4	130		
Dr. Burger	P	roduc	ers Ac	curac	Y		1	User	s Accu	iracy			
E. S. E. L.	0			71.19	%		0			57.1%	6		
Ener Sale	1			17.6	%		1			27.3%	6		
	2			-			2			÷			
1028530536	3			-			3			-			
- Andrew Street	4			-			4			-			
	5			5.6%	6		5		_	14.3%	6		
	6			72.29	%		6			23.2%	6		
	7			-			7			-			
	8			-			8			-			

Table A2.19: Correlation matrix for field collected independent variables.

				8.28	2.11	Actua	1					
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	33	9	2	2	2	7	4	1	0	60	Accuracy
	1	0	0	0	0	0	0	0	0	0	0	1.1.1.1.1.1.1.1
	2	0	0	0	0	0	0	0	0	0	0	
	3	0	0	0	0	0	0	0	0	0	0	
Predicted	4	0	0	0	0	0	0	0	0	0	0	0.361538
	5	4	1	1	0	2	1	1	1	0	11	Contraction of the second
	6	8	7	2	5	3	10	13	7	4	59	
	7	0	0	0	0	0	0	0	0	0	0	
	8	0	0	0	0	0	0	0	0	0	0	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	roduc	ers Ac	curac	y			User	s Accu	racy	1 Million	1
	0			73.3	%		0			55%	)	
	1			-			1			-		
	2			-			2		-			
	3			-			3			-		
	4			-			4			-		
	5			5.6%	6		5			9.1%	, 0	
	6			72.29	%		6			22%		
	7						7		-			
	8			-			8 -					

Table A2.20: Correlation matrix for field collected and the 10 metre DEM independent variables.

E. CHA. L. CAR						Actua	1	15		S. Stall		STURING ST
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	35	10	2	1	4	6	4	2	0	64	Accuracy
	1	2	1	1	0	0	0	0	0	0	4	
	2	0	0	0	0	0	0	0	0	0	0	
	3	0	0	0	0	0	0	0	0	0	0	5 500000000
Predicted	4	0	0	0	0	0	0	0	0	0	0	0.369231
	5	8	6	2	6	3	12	14	7	4	62	
	6	0	0	0	0	0	0	0	0	0	0	
	7	0	0	0	0	0	0	0	0	0	0	
	8	0	0	0	0	0	0	0	0	0	0	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	roduc	ers Ac	curac	y		1 State					
	0			77.8	%		0			54.7%	6	
	1			5.9%	6		1			25%		
	2			-			2			2		
	3			-			3			-		
	4			-			4			×.		
	5			66.79	%		5			19.4%	6	
	6			-			6			-		
	7			-			7		-			
	8			-			8					

Table A2.21: Correlation matrix for field collected and the 20 metre DEM independent variables.

	1	1.8.		Coza-		Actua	1					
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	38	10	3	4	3	9	3	3	0	73	Accuracy
	1	0	0	0	0	0	0	0	0	0	0	
	2	0	0	0	0	0	0	0	0	0	0	
	3	0	0	0	0	0	0	0	0	0	0	
Predicted	4	0	0	0	0	0	0	0	0	0	0	0.392308
	5	3	1	0	0	0	0	2	0	0	6	
	6	4	6	2	3	4	9	13	6	4	51	
	7	0	0	0	0	0	0	0	0	0	0	
	8	0	0	0	0	0	0	0	0	0	0	
	Total	45	17	5	7	7	18	18	9	4	130	1
	F	roduc	ers Ac	curac	y		A.	User	s Accu	iracy		
	0		84.4%				0		52.1%			
	1			-	1.1		1			-		
	2			-			2			-		
	3			-			3			-		
	4			-			4			-		
	5		-				5		-			
	6	6		72.2%			6		25.5%			
	7		-				7		-			
	8			-			8			-		

Table A2.22: Correlation matrix for field collected and classified vegetation independent variables.

						Actua	1			1200	Section 1	
The second	10007	0	1	2	3	4	5	6	7	8	Total	Overall
	0	37	8	2	4	3	7	2	4	0	67	Accuracy
	1	2	1	1	0	0	2	3	0	0	9	
	2	0	0	0	0	0	0	0	0	0	0	
	3	0	0	0	0	0	0	0	0	0	0	
redicted	4	0	0	0	0	0	0	0	0	0	0	0.392308
	5	0	0	0	0	0	0	0	0	0	0	
	6	6	8	2	3	4	9	13	5	4	54	
	7	0	0	0	0	0	0	0	0	0	0	
	8	0	0	0	0	0	0	0	0	0	0	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	roduc	ers Ac	curac	y			User	s Accu	iracy		
	0	1112	82.2%				0		55.2%			
	1			5.9%	6		1			11.19	6	
	2			-			2		-			
	3			-	1.11	100	3		1.11	-		
	4	-		-			4					
	5			-	-		5			-		
	6		72.2%				6		24.1%			
	7	1.1		-			7		2			1
	9			-			8					

Table A2.23: Correlation matrix for field collected, classified vegetation and 10 metre DEM independent variables.

						Actua	1	1				
State of the		0	1	2	3	4	5	6	7	8	Total	Overall
	0	31	10	3	6	3	10	5	3	1	72	Accuracy
	1	5	3	1	1	0	0	0	0	0	10	
	2	0	0	0	0	0	0	0	0	0	0	
	3	0	0	0	0	0	0	0	0	0	0	2 3 50-0 D
Predicted	4	0	0	0	0	0	0	0	0	0	0	0.361538
	5	0	0	0	0	0	0	0	0	0	0	
	6	9	4	1	0	4	8	13	6	3	48	
	7	0	0	0	0	0	0	0	0	0	0	S
	8	0	0	0	0	0	0	0	0	0	0	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	roduc	ers Ac	curac	y		Users Accuracy				No. of the state	
	0			68.9	%		0			43.19	6	
	1		17.6%				1		30%			
	2			-			2					
	3	-					3		-			
	4			-		-	4			-		
	5			-			5					
	6			72.2	%		6			6.3%	6	
	7	7 -			7		-					
	8		I	-			8			-	1	

 Table A2.24: Correlation matrix for field collected, classified vegetation and 20 metre DEM independent variables.

**Correlation Matrices Produced for the Decision Tree Classifiers** 

Two Class Classifications

Stall and		Actu	ial	A STOR	
2511 124		0	1	Total	Overall
Predicted	0	27	24	51	Accuracy
	1	18	61	79	
	Total	45	85	130	and the
	Producers A	ccuracy	Users A	Accuracy	0.676923
	0	60%	0	52.9%	
	1	71.8%	1	77.2%	

Table A2.25: Correlation matrix for 10 metre DEM independent variables.

El Martin I	Serie and the series of the	Actu	ıal	ALC: MAL	
C.S.S.S. Contraction	Lintal	0	1	Total	Overall
Predicted	0	22	33	55	Accuracy
	1	23	52	75	
	Total	45	85	130	
	Producers A	ccuracy	Users A	Accuracy	0.569231
	0	48.9%	0	40%	
	1	61.2%	1	69.3%	

Table A2.26: Correlation matrix for 20 metre DEM independent variables.

I State State	1	Actu	ıal		a series and a series of
	1	0	1	Total	Overall
Predicted	0	34	20	54	Accuracy
	1	11	65	76	S. Standard
	Total	45	85	130	
	Producers A	ccuracy	Users A	Accuracy	0.761538
	0	75.6%	0	63%	
	1	76.5%	1	85.5%	

Table A2.27: Correlation matrix for field collected independent variables.

1446233		Actu	ıal		
The second second	1.1	0	1	Total	Overall
Predicted	0	34	20	54	Accuracy
	1	11	65	76	
	Total	45	85	130	1.21 79.77
	Producers A	ccuracy	Users A	Accuracy	0.761538
	0	60%	0	65.9%	
	1	83.5%	1	79.8%	

Table A2.28: Correlation matrix for field collected and the 10 metre DEM independent variables.

a part of an	(Classification	Actu	ıal		Spectrum and the second
		0	1	Total	Overall
Predicted	0	35	24	59	Accuracy
	1	10	61	71	
	Total	45	85	130	1 min
C. S. STANK	Producers A	ccuracy	Users A	Accuracy	0.738462
Star Harris	0	77.8%	0	59.3%	21
	1	71.8%	1	85.9%	- 25

Table A2.29: Correlation matrix for field collected and the 20 metre DEM independent variables.

		Actu	ıal			
Tables 2.3	<ul> <li>Conseliation</li> </ul>	0	1	Total	Overall	
Predicted	0	25	9	34	Accuracy	
	1	20	76	96		
	Total	45	85	130		
	Producers A	ccuracy	Users A	Accuracy	0.776923	
	0	55.6%	0	73.5%		
	1	72.9%	1	79.2%		

Table A2.30: Correlation matrix for field collected and classified vegetation independent variables.

The state of the	1.1.1	Actu	ial		
-4.7-2.10		0	1	Total	Overall
Predicted	0	31	19	50	Accuracy
	1	14	66	80	
	Total	45	85	130	
	Producers A	ccuracy	Users A	Accuracy	0.746154
	0	68.9%	0	62%	
	1	77.6%	1	82.5%	

 Table A2.31: Correlation matrix for field collected, classified vegetation and 10 metre DEM independent variables.

	res (Lotalise)	Actı	ial		
	D. Produktor	0	1	Total	Overall
Predicted	0	25	9	34	Accuracy
	1	20	76	96	
	Total	45	85	130	
	Producers A	ccuracy	Users A	Accuracy	0.776923
	0	55.6%	0	73.5%	
	1	72.9%	1	79.2%	

Table A2.32: Correlation matrix for field collected, classified vegetation and 20 metre DEM independent variables.

#### Appendices

## Three Class Classifications

Barrie State			Actual			S. Sector
The strange	0	0	1	2	Total	Overall
Predicted	0	18	4	12	34	Accuracy
	1	8	7	6	21	
	2	19	16	40	75	in the second
	Total	45	29	56	130	0.5
T I	Producers	Accuracy	1	Jsers Acci	uracy	
	0	40%	0		52.9%	
	1	24.1%	1		33.3%	
La la Conta	2	71.4%	2		53.3%	

Table A2.33: Correlation matrix for 10 metre DEM independent variables.

9. Stander	La grand h	1	Actual			Contraction of
E.C.		0	1	2	Total	Overall
Predicted	0	18	5	9	32	Accuracy
	1	11	8	8	27	
	2	16	14	41	71	0.515385
	Total	45	29	56	130	
	Producers	s Accuracy	1	Users Acc	uracy	1
	0	40%	0	)	56.3%	
	1	27.6%	1		29.6%	
	2	73.2%	2	2	57.7%	]

Table A2.34: Correlation matrix for 20 metre DEM independent variables.

	and the second		Actual		1 - Stin these	
		0	1	2	Total	Overall
Predicted	0	25	4	5	34	Accuracy
	1	9	11	15	35	1.00.00024
	2	11	12	38	61	0.569231
	Total	45	29	56	130	
	Producers	s Accuracy	1	Users Acc	uracy	
SUIT	0	55.6%	(	)	73.5%	
	1	37.9%	1		31.4%	
and the states	2	67.9%	2	2	62.3%	

Table A2.35: Correlation matrix for field collected independent variables.

	al se la se		Actual		NUL SUSA	Ann San Gel
		0	1	2	Total	Overall
Predicted	0	27	6	6	39	Accuracy
and a long	1	9	10	10	29	
	2	9	11	42	62	
	Total	45	29	56	130	0.607692
	Producers	s Accuracy		Users Acc	uracy	]
	0 60%		(	)	69.2%	1
	1	34.5%	1		34.5%	]
	2	75%	2		67.7%	]

Table A2.36: Correlation matrix for field collected and the 10 metre DEM independent variables.

			Actual			
Contraction of the		0	1	2	Total	Overall
Predicted	0	28	7	8	43	Accuracy
	1	7	9	8	24	
	2	10	11	42	63	
	Total	45	29	56	130	0.607692
	Producers	s Accuracy	1	Users Acc	uracy	]
	0 62.2%		(	)	65.1%	
	1	31%	1		37.5%	
	2 75%		2	2	66.7%	

Table A2.37: Correlation matrix for field collected and the 20 metre DEM independent variables.

		-	Actual			
		0	1	2	Total	Overall
Predicted	0	35	9	13	57	Accuracy
	1	7	11	9	27	
	2	3	7	36	46	
	Total	45	29	56	130	0.630769
	Producers	s Accuracy	1	Users Acc	uracy	
-	0	77.8%	(	)	61.4%	
	1	37.9%	1		40.7%	
	2	64.3%	2	2	78.3%	

Table A2.38: Correlation matrix for field collected and classified vegetation independent variables.

		1	Actual			
Constant of the		0	1	2	Total	Overall
Predicted	0	34	9	12	55	Accuracy
	1	8	11	10	29	
	2	3	7	36	46	
STREET.	Total	45	29	56	130	0.623077
	Producers	s Accuracy	1	Users Acc	uracy	
	0	75.6%	- (	)	61.8%	
	1	37.9%	1	L	37.9%	
-	2	64.3%	2	2	78.3%	

 Table A2.39: Correlation matrix for field collected, classified vegetation and 10 metre DEM independent variables.

	State States		Actual		19. 00.0-5 Mar	Carlo State
		0	1	2	Total	Overall
Predicted	0	34	9	12	55	Accuracy
	1	9	11	10	30	
	2	2	7	36	45	0.623077
Terbles 175	Total	45	29	56	130	
	Producers	s Accuracy	1	Users Acc	uracy	]
3409534000	0	75.6%	(	)	61.8%	
24. 19 States	1	37.9%	1		36.7%	
	2	64.3%	2		80%	

**Table A2.40:** Correlation matrix for field collected, classified vegetation and 20metre DEM independent variables.

## Nine Class Classifications

		18		10		Actua	1			No.		C. A DRAME U
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	15	2	1	1	1	2	1	4	0	27	Accuracy
	1	0	0	0	0	0	0	0	0	0	0	
	2	5	1	0	0	1	1	0	2	0	10	1/21/201
	3	0	0	0	0	0	0	0	0	0	0	
Predicted	4	3	3	0	1	0	4	1	0	0	12	0.223077
	5	6	3	1	2	3	6	8	1	0	30	
	6	7	3	0	1	1	2	4	1	0	19	
	7	0	0	0	0	0	0	0	0	0	0	
	8	9	5	3	2	1	3	4	1	4	32	
	Total	45	17	5	7	7	18	18	9	4	130	
	P	roduc	ers Ac	curacy	y							
	0			33.39	%		0			55.6%	6	
	1			-			1		-			
	2			-			2			-		
	3			-			3			-		
	4						4			-		
5				33.39	%		5			20%		
6			22.2%				6			21.19	6	
		-				7		-				
	8				6		8		12.5%			

Table A2.41: Correlation matrix for 10 metre DEM independent variables.

Falls a Constant					1.00	Actua	1	1-	1		Tem alter	
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	0	0	0	0	0	0	0	0	0	0	Accuracy
	1	12	2	0	2	0	2	1	2	0	21	
	2	11	5	4	1	1	0	5	2	1	30	1.000
	3	0	0	0	0	0	2	0	0	0	2	
Predicted	4	3	2	0	1	1	3	2	1	0	13	0.130769
	5	0	0	0	0	0	0	0	0	0	0	
	6	12	2	1	3	2	6	7	1	0	34	
	7	4	3	0	0	1	5	1	1	1	16	
	8	3	3	0	0	2	0	2	2	2	14	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	roduc	ers Ac	curacy	Y		hand a					
	0						0			-		
	1			11.89	%		1		9.5%			
	2			80%	ò		2			13.39	6	
	3			(im)			3			7.7%	ó	
	4			14.39	%		4			-		
	5			-	1		5					
	6	6		38.9%					20.6%			
	7		11.1%				7		6.3%			1993
	8			50%	6		8		14.3%			

Table A2.42: Correlation matrix for 20 metre DEM independent variables.

			a lot			Actua	1	1 3 1		1.15	1/10/11/17/1	2.2.5
THE REAL PROPERTY OF		0	1	2	3	4	5	6	7	8	Total	Overall
	0	19	2	0	2	1	2	1	0	0	27	Accuracy
Carl Start	1	0	4	1	2	1	3	3	1	0	15	
	2	15	4	1	1	0	2	0	0	0	23	
	3	2	0	0	0	0	2	0	0	0	4	
Predicted	4	5	1	0	0	3	0	4	2	0	15	0.292308
	5	2	4	3	2	0	4	3	1	0	19	
	6	0	0	0	0	1	1	2	1	0	5	
Star Carlo	7	1	2	0	0	1	2	4	2	1	13	
MEST	8	1	0	0	0	0	2	1	2	3	9	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	roduc	ers Ac	curac	y		100.00					
12-46	0			42.2	%	_	0		70.4%			
	1			23.5	%	_	1		26.7%			
	2			20%	ó		2			4.3%	ó	
	3			-			3			-		
	4			42.9	%		4			20%		
	5			22.2	%		5			21.19	/0	
	6			11.1%			6		40%			
	7		22.2%				7		15.4%			
	8			75%	, 0		8		33.3%			

Table A2.43: Correlation matrix for field collected independent variables.

THE REAL PROPERTY.				3.3		Actua	1			1		phil 15 Source
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	16	2	0	0	1	3	0	0	0	22	Accuracy
The second second	1	14	7	3	2	1	4	1	1	0	33	Accurrence
	2	6	5	1	3	1	2	3	1	0	22	
The second second	3	8	1	1	2	1	3	3	2	0	21	
Predicted	4	0	0	0	0	1	2	4	1	0	8	0.238462
Research -	5	0	0	0	0	0	0	0	0	0	0	A.3530-40
	6	0	0	0	0	0	0	0	0	0	0	
1. 1. 1. 1. 1. 1.	7	0	0	0	0	0	0	0	0	0	0	
	8	1	2	0	0	2	4	7	4	4	24	
	Total	45	17	5	7	7	18	18	9	4	130	
	P	roduc	ers Ac	curacy	y			User	s Accu	iracy	21.2.2	
	0			35.6	%		0		72.7%			
	1			41.29	%		1		21.2%			]
	2			20%	0		2		4.5%			
	3			28.6	%		3			9.5%	0	
	4			14 30	2/0		4			12.50	10	
				110	/0	-	5	-		1.4.10.1		
3				-		-	5			-		
	6			-			0	-+	-			
	7		-				7		-			
	8				6	_	8		_	16.79	10	

Table A2.44: Correlation matrix for field collected and the 10 metre DEM independent variables.

	1.			1		Actua	1		10-7		E Year	China Herb
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	18	3	0	1	1	1	1	0	0	25	Accuracy
	1	12	4	4	3	1	1	3	1	0	29	
	2	5	2	0	1	0	3	0	1	0	12	
	3	1	1	0	0	0	0	0	1	0	3	1.111.141
Predicted	4	1	0	1	0	1	3	4	1	0	11	0.269231
	5	0	2	0	1	2	2	1	2	0	10	
	6	2	0	0	1	1	4	5	1	0	14	
	7	6	5	0	0	1	3	3	2	1	21	
	8	0	0	0	0	0	1	1	0	3	5	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	roduc	ers Ac	curac	y			User	s Accu	iracy		
	0			40%	0		0			72%		
	1			23.5	%		1		13.8%			
	2			-			2			-		
	3			-	S		3			-		
	4			14.3	%		4			9.1%	5	
	5			11.19	%		5			20%		
	6			27.8	%		6			35.7%	6	
	7		22.2%				7		9.5%			
	8	-	1	75%	ó		8			60%		

Table A2.45: Correlation matrix for field collected and the 20 metre DEM independent variables.

The second	12.1.5				1.5	Actua	1		1.1		1 Acres	and the second
Distance in the		0	1	2	3	4	5	6	7	8	Total	Overall
	0	21	2	0	0	0	2	0	0	0	25	Accuracy
Predicted	1	13	5	3	3	1	4	3	3	0	35	0.253846
	2	6	5	1	3	1	2	3	1	0	22	
	3	4	3	1	1	2	4	1	0	0	16	
	4	0	0	0	0	1	2	4	1	0	8	
	5	0	0	0	0	0	0	0	0	0	0	
	6	0	0	0	0	0	0	0	0	0	0	
	7	0	0	0	0	0	0	0	0	0	0	
	8	1	2	0	0	2	4	7	4	4	24	
	Total	45	17	5	7	7	18	18	9	4	130	
	Producers Accuracy Users Accuracy											
	0	46.7%				0		84%				
	1		29.4%				1		14.3%			
	2		20%				2		4.5%			
	3		14.3%				3		6.3%			
	4		14.3%				4		12.5%			
	5		-				5		4			
	6		-				6		~			
	7						7		-			
	8		100%				8		16.7%			

Table A2.46: Correlation matrix for field collected and classified vegetation independent variables.
				- 10		Actua	1	1. S.	R	ST.U.S.		121-212112
Contraction of the log	1.1.1.1.1.1	0	1	2	3	4	5	6	7	8	Total	Overall
1. 1. 1. 1. 1.	0	16	2	0	0	1	3	0	0	0	22	Accuracy
and a second	1	14	7	3	2	1	4	1	1	0	33	
	2	6	5	1	3	1	2	3	1	0	22	
	3	8	1	1	2	1	3	3	2	0	21	
Predicted	4	0	0	0	0	1	2	4	1	0	8	0.238462
ME ROLLING	5	0	0	0	0	0	0	0	0	0	0	
and the second	6	0	0	0	0	0	0	0	0	0	0	
PAULT	7	0	0	0	0	0	0	0	0	0	0	
M.F. Barrelle	8	1	2	0	0	2	4	7	4	4	24	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	roduc	ers Ac	curac	y		Br. con	User	s Accu	iracy		
	0			35.6	%		0		<u>.</u>	72.7%	6	
	1		10	41.2	%	_	1	2	21.2%			
	2		61	20%	6		2			4.5%	Ó	
	3		-	28.6	%	112 13	3			9.5%	Ó	
174	4			14.39	%		4			12.5%	6	
E-WALLEN	5			-			5			-		
	6			-			6			-		
	7			-			7			-		
	8			1009	10		8			16.7%	10	

 Table A2.47: Correlation matrix for field collected, classified vegetation and 10 metre DEM independent variables.

THE REAL PROPERTY.						Actua	1					
The second second		0	1	2	3	4	5	6	7	8	Total	Overall
	0	30	9	3	2	2	7	1	1	0	55	Accuracy
	1	0	0	0	0	0	0	0	0	0	0	
	2	6	5	1	3	1	2	3	1	0	22	
	3	8	1	1	2	1	3	3	2	0	21	
Predicted	4	0	0	0	0	1	2	4	1	0	8	0.292308
	5	0	0	0	0	0	0	0	0	0	0	
	6	0	0	0	0	0	0	0	0	0	0	
	7	0	0	0	0	0	0	0	0	0	0	
	8	1	2	0	0	2	4	7	4	4	24	
	Total	45	17	5	7	7	18	18	9	4	130	
	F	roduc	ers Ac	curacy	y	1.	Users Accuracy					
	0			66.79	%	_	0			54.5%		
	1			-			1					
	2			20%	6		2 4			4.5%	0	
	3			28.6	%		3			9.5%	0	
	4			14.39	%		4			12.5%	10	
	5			-			5			-		
	6	6		-			6		annen -			
	7		-				7		-			
	8			100%	6		8		16.7%			

 Table A2.48: Correlation matrix for field collected, classified vegetation and 20 metre DEM independent variables.

**Correlation Matrices Produced for the Discriminant Analysis Classifiers** 

#### **Two Class Classifications**

	in as mererally	Actu	ıal		C.S. S.S.S.	
and a second		0	1	Total	Overall	
Predicted	0	30	33	63	Accuracy	
	1	15	52	67		
	Total	45	85	130	0.630769	
	Producers A	ccuracy	Users A	Accuracy		
	0	66.7%	0	47.6%		
	1	61.2%	1	77.6%		

Table A2.49: Correlation matrix for 10 metre DEM independent variables.

		Acti	ıal			
	Thread	0	1	Total	Overall	
Predicted	0	28	25	53	Accuracy	
	1	17	60	77	0.676923	
	Total	45	85	130		
	Producers A	ccuracy	Users A	Accuracy		
	0	62.2%	0	52.8%		
	1	70.6%	1	77.9%		

Table A2.50: Correlation matrix for 20 metre DEM independent variables.

	Provide States	Actu	ıal			
Phillips 1997		0	1	Total	Overall	
Predicted	0	33	18	51	Accuracy	
	1	12	67	79	- 2-78-612	
	Total	45	85	130	0.769231	
	Producers A	ccuracy	Users /			
	0	73.3%	0	64.7%		
	1	78.8%	1	84.8%		

Table A2.51: Correlation matrix for field collected independent variables.

name -	Sector States	Actu	ıal		74506696	
		0	1	Total	Overall	
Predicted	0	33	18	51	Accuracy	
	1	12	67	79		
	Total	45	85	130	0.769231	
	Producers A	ccuracy	Users A	Accuracy		
	0	73.3%	0	64.7%		
	1	78.8%	1	84.8%		

Table A2.52: Correlation matrix for field collected and the 10 metre DEM independent variables.

		Actu	ıal			
TCH Z ALAND		0	1	Total	Overall	
Predicted	0	33	18	51	Accuracy	
	1	12	67	79	Tetal. 1	
	Total	45	85	130	0.769231	
	Producers A	ccuracy	Users A	Accuracy		
	0	73.3%	0	64.7%		
	1	78.8%	1	84.8%		

Table A2.53: Correlation matrix for field collected and 20 metre DEM independent variables.

	2. Cartelatino	Actu	ıal		an a processi
		0	1	Total	Overall
Predicted	0	34	19	53	Accuracy
	1	11	66	77	
	Total	45	85	130	0.769231
	Producers A	ccuracy	Users A	Accuracy	The second second
	0	75.5%	0	64.2%	
	1	77.6%	1	85.7%	1.12

Table A2.54: Correlation matrix for field collected and classified vegetation independent variables.

		Actu	ıal	1. 1. 1. 23.	A BANGAR CON	
		0	1	Total	Overall	
Predicted	0	34	17	51	Accuracy	
	1	11	68	79		
	Total	45	85	130	0.784615	
	Producers A	ccuracy	Users A	Accuracy		
	0	75.6%	0	66.7%		
	1	80%	1	86.1%		

Table A2.55: Correlation matrix for field collected, classified vegetation and 10 metre DEM independent variables.

	A sen al sen	Actu	ıal			
		0	1	Total	Overall	
Predicted	0	34	19	53	Accuracy	
	1	11	66	77	0.769231	
	Total	45	85	130		
	Producers A	ccuracy	Users A	Accuracy		
	0	75.5%	0	64.2%		
	1	77.6%	1	85.7%		

 Table A2.56: Correlation matrix for field collected, classified vegetation and 20 metre DEM independent variables.

### Appendices

## Three Class Classifications

		1	Actual	Des des	1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1	C. Salaria
Save States		0	1	2	Total	Overall
Predicted	0	27	8	19	54	Accuracy
	1	8	7	5	20	1.
	2	10	12	34	56	1
	Total	45	29	56	130	0.523077
	Producers	s Accuracy	1	uracy		
	0	60%	(	)	50%	
	1	24.1%	1	L .	35%	
All and and a second	2	60.7%	2	2	60.7%	

Table A2.57: Correlation matrix for 10 metre DEM independent variables.

and the second		1.4.	Actual			A TRANSPORT
In a line of the	0	0	1	2	Total	Overall
Predicted	0	20	5	11	36	Accuracy
	1	14	9	8	31	
	2	11	13	39	63	25.20,000
	Total	45	29	56	130	0.523077
	Producers		Users Acc			
	0	44.4%	(	)	55.6%	
	1	31%	1	1	29%	]
	2	69.6%	1	2	61.9%	1000

Table A2.58: Correlation matrix for 20 metre DEM independent variables.

		1	Actual			
		0	1	2	Total	Overall
Predicted	0	25	4	8	37	Accuracy
	1	12	10	7	29	
Short Sur	2	8	13	43	64	]
	Total	45	29	56	130	0.6
	Producers	s Accuracy		Users Acc	uracy	
	0	55.6%	(	)	67.6%	]
	1	34.5%	1		34.5%	
	2	76.8%	2	2	67.2%	

Table A2.59: Correlation matrix for field collected independent variables.

		1	Actual			Calling Tally
		0	1	2	Total	Overall
Predicted	0	26	4	8	38	Accuracy
	1	11	12	8	31	
	2	8	11	42	61	
	Total	45	29	56	130	0.615385
	Producers	s Accuracy	1	Users Acc	uracy	
	0	57.8%	(	)	68.4%	
	1	41.4%	1		38.7%	
	2	75%	2		68.9%	1

Table A2.60: Correlation matrix for field collected and the 10 metre DEM independent variables.

	1		Actual	No. 1 in Shi	12 July 2 2 2 2	
		0	1	2	Total	Overall
Predicted	0	25	5	8	38	Accuracy
	1	11	10	8	29	
	2	9	12	42	63	1
	Total	45	29	56	130	0.592308
	Producers	s Accuracy	1	Users Acc	uracy	
	0	55.6%	(	)	65.8%	]
	1	34.5%	1		34.5%	
	2	75%	2	2	66.7%	

Table A2.61: Correlation matrix for field collected and the 20 metre DEM independent variables.

	million State		Actual	de la conte		WASSING.
		0	1	2	Total	Overall
Predicted	0	27	8	19	54	Accuracy
	1	8	7	5	20	
	2	10	12	34	56	1
	Total	45	29	56	130	0.523077
	Producers	s Accuracy	1	Users Acc	curacy	
	0	60%	(	)	50%	
	1	24.1%	1		35%	1
	2	60.7%	1	2	60.7%	

Table A2.62: Correlation matrix for field collected and classified vegetation independent variables.

			Actual			
		0	1	2	Total	Overall
Predicted	0	27	4	9	40	Accuracy
	1	12	9	7	28	and the second
	2	6	14	42	62	
	Total	45	29	56	130	0.6
	Producers	s Accuracy		Users Acc	uracy	
	0	60%		)	67.5%	
	1	31%		1	32.1%	
	2	75%		2	67.7%	

 Table A2.63: Correlation matrix for field collected, classified vegetation and 10 metre DEM independent variables.

10210		1	Actual			
		0	1	2	Total	Overall
Predicted	0	26	6	8	40	Accuracy
	1	11	10	7	28	
	2	8	11	43	62	
	Total	45	29	56	130	0.607692
Pable AL.6	Producers	s Accuracy	1	Users Acc	uracy	
	0	57.8%	(	)	65%	
	1	34.5%	1		35.7%	
	2	76.8%	2	2	69.4%	

 Table A2.64: Correlation matrix for field collected, classified vegetation and 20 metre DEM independent variables.

## Nine Class Classifications

						Actua		100				
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	19	4	0	0	3	5	3	1	0	35	Accuracy
	1	1	1	1	2	0	0	2	2	0	9	
	2	3	0	0	0	0	0	1	1	1	6	1.00
and the second	3	2	1	0	0	0	3	2	2	0	10	
Predicted	4	0	0	0	1	0	0	0	0	0	1	0.2
	5	10	1	1	2	1	4	5	1	0	25	
	6	2	1	0	0	0	1	0	0	1	5	
	7	3	3	0	1	1	3	2	0	0	13	
	8	5	6	3	1	2	2	3	2	2	26	
	Total	45	17	5	7	7	18	18	9	4	130	
	P	roduc	ers Ac	curac	y	1	6.14	User	s Accu	iracy	THE PARTY OF	
	0			42.2	%		0			54.3%	0	
	1			5.9%	0	_	1			11.19	6	
	2						2			-		
	3						3					
	4			4			4			-		
	5			22.2	%		5			16%		
	6			-			6			-		
	7			-			7			-		
	8			50%	0		8			7.7%	0	

Table A2.65: Correlation matrix for 10 metre DEM independent variables.

						Actua	1					
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	14	1	1	3	1	1	1	0	0	22	Accuracy
	1	12	3	1	1	0	1	0	2	0	20	
	2	1	2	1	0	0	0	0	2	0	6	
	3	5	1	1	1	1	4	3	1	1	18	
Predicted	4	0	0	0	0	0	0	0	0	0	0	0.253846
	5	6	4	0	1	1	7	7	1	0	27	
	6	4	4	1	1	3	5	4	1	0	23	
	7	2	2	0	0	1	0	2	1	1	9	
	8	1	0	0	0	0	0	1	1	2	5	2 3
	Total	45	17	5	7	7	18	18	9	4	130	
	P	roduc	ers Ac	curac	y		1	User	s Acci	iracy		
	0			31.1	%	-	0			63.69	1/0	
	1			17.6	%		1			15%	0	
	2			20%	ó		2			16.79	%	5 I - 15 - 1
	3			14.3	%		3			5.6%	6	
	4			-			4			-		
	5			38.9	%		5			25.99	%	
	6			22.2	%		6			17.49	%	
	7			11.19	%		7			11.19	%	in a start of the
and the	8			50%	6		8			40%	ó	ANTEN A

Table A2.66: Correlation matrix for 20 metre DEM independent variables.

CHARLES AND AND					1 Martin	Actua	1		1000	the set		
		0	1	2	3	4	5	6	7	8	Total	Overall
CARLES ST	0	18	2	1	0	1	2	1	0	0	25	Accuracy
A 7 15 742	1	6	2	1	1	0	1	2	1	0	14	
Store States	2	7	5	0	3	1	2	1	2	0	21	
	3	7	0	0	0	1	2	1	1	0	12	
Predicted	4	0	2	1	0	2	4	2	1	1	13	0.2
	5	5	3	2	2	1	1	1	1	0	16	
1. 2. 1. 1. 1. 1.	6	0	2	0	0	1	2	0	0	0	5	
	7	1	1	0	1	0	1	6	2	2	14	
	8	1	0	0	0	0	3	4	1	1	10	
	Total	45	17	5	7	7	18	18	9	4	130	
A STATE OF	F	roduc	ers Ac	curac	y		5.2	User	s Acci	iracy		
	0			40%	0		0			72%		
	1			11.8	%		1			14.3%	0	
	2			-			2			-		
	3			+			3			(+)		
	4			28.6	%		4			15.4%	6	
	5			5.6%	6		5			6.3%	, D	
	6			-			6			-		
	7			22.29	%		7			14.39	6	
	8			25%	ó		8			10%		

Table A2.67: Correlation matrix for field collected independent variables.

					101	Actua	1	1013	6			
		0	1	2	3	4	5	6	7	8	Total	Overall
A Carton	0	18	1	1	0	1	2	1	1	0	25	Accuracy
	1	6	3	1	1	0	1	2	1	0	15	
Rufe Instan	2	9	5	0	3	1	5	2	1	0	26	
	3	4	1	0	0	0	1	0	0	0	6	i National and the second
Predicted	4	0	0	0	1	1	0	1	1	0	4	0.215385
	5	6	3	2	2	2	2	2	0	0	19	
	6	0	2	1	0	1	1	0	0	0	5	
	7	1	2	0	0	1	3	4	3	3	17	
	8	1	0	0	0	0	3	6	2	1	13	
	Total	45	17	5	7	7	18	18	9	4	130	
1. 97. 47. 60	F	roduc	ers Ac	curac	y	1		User	s Acci	iracy	111.2.1.1	
	0			40%	6		0			72%	0	
	1			17.6	%		1		16	20%	0	
A STREET, # .	2			-			2			-		
in stands	3			-			3			-		
	4			14.3	%		4			25%	0	
in the state	5			11.19	%		5			10.59	%	
	6			-			6			-		
	7			33.3	%		7		1	17.6	%	
	8			25%	6		8			7.7%	6	

Table A2.68: Correlation matrix for field collected and the 10 metre DEM independent variables.

						Actua	1	1.5	2			
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	18	2	1	0	1	2	2	0	0	26	Accuracy
	1	8	4	2	1	0	0	0	1	0	16	- Canada
1.11	2	6	4	0	2	1	3	2	2	0	20	
	3	7	0	0	0	1	2	1	1	0	12	
Predicted	4	0	0	0	2	2	0	1	1	0	6	0.276923
Reservation	5	4	2	1	2	1	4	1	0	0	15	
Letter 2 al	6	0	3	1	0	0	4	4	0	0	12	
File Cont	7	2	2	0	0	1	2	3	3	3	16	
	8	0	0	0	0	0	1	4	1	1	7	
101 P 31 3 3	Total	45	17	5	7	7	18	18	9	4	130	
The second	P	roduc	ers Ac	curacy	y			User	s Acce	iracy	0.0497.1	
Section 1.	0			40%	0		0			69.2%	6	
	1			23.5	%		1			25%		
	2						2			-		
a set that	3						3			-		
	4			28.69	%		4			33.3%	6	
1415 25476	5			22.29	%		5			26.7%	6	
	6			22.29	%		6			33.3%	6	
	7			33.39	%		7			18.8%	6	
Star No.2	8			25%	0		8			14.3%	6	

Table A2.69: Correlation matrix for field collected and the 20 metre DEM independent variables.

		112				Actua	1				1999	The Wall
		0	1	2	3	4	5	6	7	8	Total	Overall
	0	21	2	2	0	1	0	1	0	0	27	Accuracy
	1	5	2	0	0	0	1	2	1	0	11	in the second
	2	6	4	0	4	0	1	0	3	0	18	
	3	7	0	0	0	1	2	1	0	0	11	
Predicted	4	0	5	0	1	1	3	3	1	1	15	0.261538
	5	3	3	2	1	1	5	0	0	0	15	1. 16. 24033
	6	0	0	1	0	2	1	3	1	0	8	
	7	1	1	0	1	1	2	5	1	2	14	
	8	2	0	0	0	0	3	3	2	1	11	
	Total	45	17	5	7	7	18	18	9	4	130	
	I	roduc	ers Ac	curac	у		in the	User	s Acci	iracy	1-	
	0			46.7	%		0			77.89	%	
	1			11.8	%		1			18.29	%	
	2			-			2			-		1.1.1
	3			-			3			-		
	4			14.3	%		4			6.7%	6	
	5			27.8	%		5			33.39	%	
	6			16.7	%		6			37.5	%	100
	7			11.19	%		7			7.1%	6	1.70.4
	8			25%	6		8		1.4	9.1%	6	

Table A2.70: Correlation matrix for field collected and classified vegetation independent variables.

		1.1211				Actua	1	1				
		0	1	2	3	4	5	6	7	8	Total	Overall
Predicted	0	20	1	2	0	1	3	1	1	0	29	Accuracy
	1	6	1	0	0	0	1	2	1	0	11	0.2
	2	8	5	0	4	1	2	1	2	0	23	
	3	4	1	0	0	0	3	1	0	0	9	
	4	0	0	0	1	0	0	2	1	0	4	
	5	5	6	2	2	2	3	3	1	0	24	
	6	0	1	1	0	2	0	0	0	0	4	
	7	1	2	0	0	1	3	3	1	3	14	
	8	1	0	0	0	0	3	5	2	1	12	
	Total	45	17	5	7	7	18	18	9	4	130	
	Producers Accuracy						19144					
	0	44.4%				0			69%			
	1	5.9%				1		9.1%			_	
	2	-				2			-			
	3	-				3		-				
	4					4						
	5	16.7%				5		12.5%			23.9	
	6		-				6			-	14 C. 1	
	7		11.1%				7			7.1%		
	8		25%				8			8.3%		

 Table A2.71: Correlation matrix for field collected, classified vegetation and 10 metre DEM independent variables.

- Sold Sold States		11120		31.7		Actual	1			y ST		S. S. Starting
		0	1	2	3	4	5	6	7	8	Total	Overall
Predicted	0	19	2	2	0	1	3	1	0	0	28	Accuracy
	1	8	3	1	0	0	0	1	1	0	14	0.276923
	2	4	3	0	2	0	2	2	3	0	16	
	3	6	0	0	0	1	2	0	0	0	9	
	4	0	2	0	2	1	0	1	1	1	8	
	5	5	2	1	2	1	4	1	0	0	16	
	6	0	3	1	1	1	4	6	1	0	17	
	7	1	2	0	0	2	2	3	2	2	14	
	8	2	0	0	0	0	1	3	1	1	8	
	Total	45	17	5	7	7	18	18	9	4	130	
	F											
	0		42.2%				0			67.99		
	1	17.6%				1			21.49			
	2	-				2			-			
	3	-				3						
	4		14.3%				4		12.5%			
	5	22.2%				5			25%			
	6		33.3%				6			35.39		
	7		22.2%				7			14.39		
	8		25%				8			12.59		

 Table A2.72: Correlation matrix for field collected, classified vegetation and 20 metre DEM independent variables.



Figure A3.1: The decision tree grown for the two class classification using the 10 metre DEM independent variables.



Figure A3.2: The decision tree grown for the two class classification using the 20 metre DEM independent variables.



Figure A3.3: The decision tree grown for the two class classification using the field acquired independent variables.



**Figure A3.4:** The decision tree grown for the two class classification using the field acquired and 10 metre DEM independent variables.



Figure A3.5: The decision tree grown for the two class classification using the field acquired and 20 metre DEM independent variables.



Figure A3.6: The decision tree grown for the two class classification using the field acquired and classified vegetation.



**Figure A3.7:** The decision tree grown for the two class classification using the field acquired, classified vegetation and 10 metre DEM independent variables.



**Figure A3.8:** The decision tree grown for the two class classification using the field acquired, classified vegetation and 20 metre DEM independent variables.



Figure A3.9: The decision tree grown for the three class classification using the 10 metre DEM independent variables.



Figure A3.10: The decision tree grown for the three class classification using the 20 metre DEM independent variables.



Figure A3.11: The decision tree grown for the three class classification using the field acquired independent variables.



**Figure A3.12:** The decision tree grown for the three class classification using the field acquired and 10 metre DEM independent variables.



**Figure A3.13:** The decision tree grown for the three class classification using the field acquired and 20 metre DEM independent variables.





W = 13.000

W = 21 000

Node 1 RELD\_SLOPE <= 19.000

W = 390.000 N = 390

Terrinal

Node 3

W = 48 000

Contraction of the local division of the loc

Node 7

N=32

Termolel

Node 4

W # 22.000

Termal

Node 1

W = 129.000

Termel

Node 2 W # 21.000



**Figure A3.15:** The decision tree grown for the three class classification using the field acquired, classified vegetation and 10 metre DEM independent variables.



**Figure A3.16:** The decision tree grown for the three class classification using the field acquired, classified vegetation and 20 metre DEM independent variables.



Figure A3.17: The decision tree grown for the nine class classification using the 10 metre DEM independent variables.



Figure A3.18: The decision tree grown for the nine class classification using the 20 metre DEM independent variables.



Figure A3.19: The decision tree grown for the nine class classification using the field acquired independent variables.

Appendices



**Figure A3.20:** The decision tree grown for the nine class classification using the field acquired and 10 metre DEM independent variables.



**Figure A3.21:** The decision tree grown for the nine class classification using the field acquired and 20 metre DEM independent variables.



**Figure A3.22:** The decision tree grown for the nine class classification using the field acquired independent variables and classified vegetation.



**Figure A3.23:** The decision tree grown for the nine class classification using the field acquired, classified vegetation and 10 metre DEM independent variables.



**Figure A3.24:** The decision tree grown for the nine class classification using the field acquired, classified vegetation and 20 metre DEM independent variables.

# Appendix 4 – Neural Network Architectures, Errors and Accuracy and Decision Tree Relative Cost and Topology

#### **ARTIFICIAL NEURAL NETWORKS**

Error versus Accuracy Graphs for the Two Class Classifications



Figure A4.1: 10 metre DEM independent variables.



Figure A4.2: 20 metre DEM independent variables.

#### Appendices



Figure A4.3: Field acquired independent variables.



Figure A4.4: Field acquired and 10 metre DEM independent variables.



Figure A4.5: Field acquired and 20 metre DEM independent variables.


Figure A4.6: Field acquired independent variables and classified vegetation.



Figure A4.7: Field acquired independent variables, classified vegetation and 10 metre DEM independent variables.



Figure A4.8: Field acquired independent variables, classified vegetation and 20 metre DEM independent variables.



# Error versus Accuracy Graphs for the Three Class Classifications

Figure A4.9: 10 metre DEM independent variables.



Figure A4.10: 20 metre DEM independent variables.



Figure A4.11: Field acquired independent variables.



Figure A4.12: Field acquired and 10 metre DEM independent variables.



Figure A4.13: Field acquired and 20 metre DEM independent variables.



Figure A4.14: Field acquired independent variables and classified vegetation.



Figure A4.15: Field acquired independent variables, classified vegetation and 10 metre DEM independent variables.



Figure A4.16: Field acquired independent variables, classified vegetation and 20 metre DEM independent variables.



Error versus Accuracy Graphs for the Nine Class Classifications

Figure A4.17: 10 metre DEM independent variables.



Figure A4.18: 20 metre DEM independent variables.



Figure A4.19: Field acquired independent variables.



Figure A4.20: Field acquired and 10 metre DEM independent variables.



Figure A4.21: Field acquired and 20 metre DEM independent variables.



Figure A4.22: Field acquired independent variables and classified vegetation.



**Figure A4.23:** Field acquired independent variables, classified vegetation and 10 metre DEM independent variables.



Figure A4.24: Field acquired independent variables, classified vegetation and 20 metre DEM independent variables.

# DECISION TREE CLASSIFIERS

### **Relative Cost Graphs for the Two Class Classifications**



Figure A4.25: Relative cost and terminal nodes for the 10 metre DEM independent variables.



Figure A4.26: Relative cost and terminal nodes for the 20 metre DEM independent variables.



Figure A4.27: Relative cost and terminal nodes for the field acquired independent variables.



Figure A4.28: Relative cost and terminal nodes for the field acquired and 10 metre DEM independent variables.



Figure A4.29: Relative cost and terminal nodes for the field acquired and 20 metre DEM independent variables.



Figure A4.30: Relative cost and terminal nodes for the field acquired independent variables and classified vegetation.



Figure A4.31: Relative cost and terminal nodes for the field acquired independent variables, classified vegetation and 10 metre DEM independent variables.



Figure A4.32: Relative cost and terminal nodes for the field acquired independent variables, classified vegetation and 10 metre DEM independent variables.



## **Relative Cost Graphs for the Three Class Classifications**

Figure A4.33: Relative cost and terminal nodes for the 10 metre DEM independent variables.



Figure A4.34: Relative cost and terminal nodes for the 20 metre DEM independent variables.



Figure A4.35: Relative cost and terminal nodes for the field acquired independent variables.



Figure A4.36: Relative cost and terminal nodes for the field acquired and 10 metre DEM independent variables.



Figure A4.37: Relative cost and terminal nodes for the field acquired and 20 metre DEM independent variables.



Figure A4.38: Relative cost and terminal nodes for the field acquired independent variables and classified vegetation.



Figure A4.39: Relative cost and terminal nodes for the field acquired independent variables, classified vegetation and 10 metre DEM independent variables.



Figure A4.40: Relative cost and terminal nodes for the field acquired independent variables, classified vegetation and 20 metre DEM independent variables.

### **Relative Cost Graphs for the Nine Class Classifications**



Figure A4.41: Relative cost and terminal nodes for the 10 metre DEM independent variables.



Figure A4.42: Relative cost and terminal nodes for the 20 metre DEM independent variables.



Figure A4.43: Relative cost and terminal nodes for the field acquired independent variables.



Figure A4.44: Relative cost and terminal nodes for the field acquired and 10 metre DEM independent variables.



Figure A4.45: Relative cost and terminal nodes for the field acquired and 20 metre DEM independent variables.



Figure A4.46: Relative cost and terminal nodes for the field acquired independent variables and classified vegetation.



Figure A4.47: Relative cost and terminal nodes for the field acquired independent variables, classified vegetation and 10 metre DEM independent variables.



Figure A4.48: Relative cost and terminal nodes for the field acquired independent variables, classified vegetation and 20 metre DEM independent variables.

# Appendix 5 - Sensitivity Analysis and Variable Importance

## Artificial Neural Network Sensitivity Analysis

		Slope Angle	Aspect	Flow Length	Flow Accumulation	Plan Curvature	Profile Curvature
Train	Rank	Rank 1		4	6	5	3
	Error	0.5095	0.4971	0.4858	0.4831	0.4846	0.4910
Verify	Rank	1	2	4	3	6	5
	Error	0.4678	0.4609	0.4571	0.4592	0.4560	0.4570

Table A5.1: Sensitivity analysis for the ANN trained using the 10 metre DEM data set for the two class classification.

		Slope Angle	Aspect	Flow Length	Flow Accumulation	Plan Curvature	Profile Curvature
Train	Rank	Rank 1		2	5	4	3
	Error	0.4978	0.4825	0.4892	0.4843	0.4844	0.4844
Verify	Rank	2	5	1	4	3	6
	Error	0.4630	0.4594	0.4682	0.4594	0.4619	0.4593

**Table A5.2:** Sensitivity analysis for the ANN trained using the 20 metre DEM data set for the two class classification.

		Slope Angle	Aspect	Estimated Vegetation	Field Sodicity Meter	Geology
Train	Rank	5	4	3	2	1
	Error	0.3857	0.4078	0.4103	0.4163	0.4283
Verify	Rank	1	4	5	3	2
	Error	0.4208	0.3823	0.3799	0.3840	0.4202

**Table A5.3:** Sensitivity analysis for the ANN trained using the field acquired data set for the two class classification.

andre the		Slope Angle	Aspect	Est. Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	2	3	4	5	1	9	7	8	6
	Error	0.4266	0.4184	0.4138	0.4135	0.4493	0.3760	0.3784	0.3769	0.3808
Verify	Rank	1	3	2	5	4	6	9	8	7
verny	Error	0.4546	0.4070	0.4074	0.3991	0.3996	0.3989	0.3942	0.3949	0.3968

**Table A5.4:** Sensitivity analysis for the ANN trained using the field acquired data set and the 10 metre DEM variables for the two class classification.

		Slope Angle	Aspect	Est. Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	1	4	3	5	2	7	8	6	9
	Error	0.4350	0.4134	0.4177	0.4086	0.4214	0.3721	0.3679	0.3724	0.3677
Verify	Rank	3	9	2	5	1	4	8	6	7
	Error	0.4382	0.3703	0.4459	0.3899	0.4492	0.3979	0.3765	0.3889	0.3832

**Table A5.5:** Sensitivity analysis for the ANN trained using the field acquired data set and the 20 metre DEM variables for the two class classification.

	Real Property in	Slope Angle	Aspect	Classified Vegetation	Field Sodicity Meter	Geology
Train	Rank	1	3	4	5	2
	Error	0.4377	0.4101	0.3998	0.3927	0.4195
Verify	Rank	1	4	3	5	2
	Error	0.4537	0.3992	0.4185	0.3963	0.4210

**Table A5.6:** Sensitivity analysis for the ANN trained using the field acquired and classified vegetation data set for the two class classification.

		Slope Angle	Aspect	Class. Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	1	4	3	6	2	5	9	7	8
	Error	0.4379	0.3615	0.3865	0.3389	0.4176	0.3508	0.3231	0.3291	0.3266
Verify	Rank	1	2	3	5	4	6	9	7	8
	Error	0.4753	0.4600	0.4503	0.4159	0.4461	0.4028	0.3931	0.3994	0.3935

**Table A5.7:** Sensitivity analysis for the ANN trained using the field acquired data set, classified vegetation and the 10 metre DEM variables for the two class classification.

	1 Cal	Slope Angle	Aspect	Class. Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	1	2	3	4	5	6	7	9	8
	Error	0.4448	0.4392	0.4203	0.4157	0.4146	0.4053	0.4042	0.4037	0.4040
Verify	Rank	2	4	3	9	1	8	7	6	5
	Error	0.4303	0.4023	0.4154	0.3940	0.4325	0.3941	0.3943	0.3955	0.3958

**Table A5.8:** Sensitivity analysis for the ANN trained using the field acquired data set, classified vegetation and the 20 metre DEM variables for the two class classification.

		Slope Angle	Aspect	Flow Length	Flow Accumulation	Plan Curvature	Profile Curvature
Train	Rank	Rank 1		2	5.	6	4
	Error	0.4578	0.4470	0.4511	0.4445	0.4439	0.4458
Verify	Rank	1	3	2	6	5	4
	Error	0.4527	0.4477	0.4492	0.4433	0.4442	0.4473

**Table A5.9:** Sensitivity analysis for the ANN trained using the 10 metre DEM data set for the three class classification.

		Slope Angle	Aspect	Flow Length	Flow Accumulation	Plan Curvature	Profile Curvature
Train	Rank	2	3	1	4	6	5
	Error	0.4458	0.4400	0.4544	0.4389	0.4383	0.4386
Verify	Rank	2	3	1	5	4	6
	Error	0.4489	0.4398	0.4494	0.4379	0.4386	0.4378

**Table A5.10:** Sensitivity analysis for the ANN trained using the 20 metre DEM data set for the three class classification.

		Slope Angle	Aspect	Estimated Vegetation	Field Sodicity Meter	Geology
Train	Rank	3	2	5	4	1
	Error	0.3932	0.3950	0.3851	0.3911	0.4083
Verify	Rank	2	5	4	3	1
	Error	0.3782	0.3416	0.3587	0.3632	0.3965

**Table A5.11:** Sensitivity analysis for the ANN trained using the field acquired data set for the three class classification.

		Slope Angle	Aspect	Est Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	2	5	4	3	1	6	8	9	7
	Error	0.3859	0.3636	0.3648	0.3832	0.4182	0.3453	0.3350	0.3348	0.3385
Verify	Rank	2	5	4	3	1	7	8	9	6
, cruy	Error	0.4163	0.3921	0.4110	0.4155	0.4216	0.3898	0.3859	0.3856	0.3902

**Table A5.12:** Sensitivity analysis for the ANN trained using the field acquired data set and the 10 metre DEM variables for the three class classification.

		Slope Angle	Aspect	Est. Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	1	2	4	5	3	9	6	7	8
	Error	0.3891	0.3803	0.3738	0.3645	0.3757	0.3514	0.3528	0.3526	0.3520
Verify	Rank	5	4	2	3	1	6	8	7	9
	Error	0.3878	0.3928	0.4010	0.3987	0.4208	0.3829	0.3804	0.3806	0.3804

**Table A5.13:** Sensitivity analysis for the ANN trained using the field acquired data set and the 20 metre DEM variables for the three class classification.

		Slope Angle	Aspect	Classified Vegetation	Field Sodicity Meter	Geology
Train	Rank	2	5	3	4	1
	Error	0.4028	0.3743	0.3954	0.3796	0.4161
Verify	Rank	1	5	4	3	2
1 U	Error	0.3984	0.3568	0.3625	0.3705	0.3798

**Table A5.14:** Sensitivity analysis for the ANN trained using the field acquired and classified vegetation data set for the three class classification.

John	1.6 %	Slope Angle	Aspect	Class. Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	2	3	4	5	1	7	9	8	6
	Error	0.3847	0.3824	0.3691	0.3568	0.3875	0.3422	0.3405	0.3406	0.3433
Verify	Rank	2	3	9	4	1	7	6	5	8
	Error	0.4405	0.4301	0.4198	0.4241	0.4428	0.4199	0.4204	0.4219	0.4198

**Table A5.15:** Sensitivity analysis for the ANN trained using the field acquired data set, classified vegetation and the 10 metre DEM variables for the three class classification.

1		Slope Angle	Aspect	Class. Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	2	5	4	3	1	6	8	7	9
	Error	0.3811	0.3699	0.3734	0.3762	0.4025	0.3685	0.3530	0.3532	0.3519
Verify	Rank	2	5	3	4	1	6	7	8	9
	Error	0.4081	0.3956	0.4044	0.3988	0.4127	0.3924	0.3879	0.3877	0.3870

**Table A5.16:** Sensitivity analysis for the ANN trained using the field acquired data set, classified vegetation and the 20 metre DEM variables for the three class classification.

Tank		Slope Angle	Aspect	Flow Length	Flow Accumulation	Plan Curvature	Profile Curvature
Train	Rank	1	3	2	5	6	4
	Error	0.2926	0.2913	0.2917	0.2906	0.2904	0.2911
Verify	Rank	1	2	3	5	6	4
	Error	0.2943	0.2919	0.2917	0.2895	0.2893	0.2897

**Table A5.17:** Sensitivity analysis for the ANN trained using the 10 metre DEM data set for the nine class classification.

		Slope Angle	Aspect	Flow Length	Flow Accumulation	Plan Curvature	Profile Curvature
Train	Rank	1	3	2	6	4	5
	Error	0.2920	0.2892	0.2907	0.2852	0.2874	0.2853
Verify	Rank	4	3	1	6	2	5
	Error	0.2847	0.2849	0.2872	0.2842	0.2863	0.2843

**Table A5.18:** Sensitivity analysis for the ANN trained using the 20 metre DEM data set for the nine class classification.

		Slope Angle	Aspect	Estimated Vegetation	Field Sodicity Meter	Geology
Train	Rank	1	4	3	5	2
	Error	0.2766	0.2744	0.2763	0.2723	0.2766
Verify	Rank	1	3	4	5	2
	Error	0.2845	0.2797	0.2741	0.2738	0.2825

**Table A5.19:** Sensitivity analysis for the ANN trained using the field acquired data set for the nine class classification.

		Slope Angle	Aspect	Est Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	1	4	3	5	2	8	9	7	6
	Error	0.2856	0.2821	0.2837	0.2787	0.2851	0.2776	0.2773	0.2777	0.2784
Verify	Rank	1	2	3	5	4	6	7	9	8
	Error	0.2884	0.2852	0.2766	0.2737	0.2760	0.2732	0.2732	0.2728	0.2730

Table A5.20: Sensitivity analysis for the ANN trained using the field acquired data set and the 10 metre DEM variables for the nine class classification.

		Slope Angle	Aspect	Est Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv,
Train	Rank	2	3	4	7	1	5	9	6	8
	Error	0.2786	0.2771	0.2766	0.2686	0.2814	0.2728	0.2685	0.2687	0.2686
Verify	Rank	2	1	3	5	4	6	7	9	8
	Error	0.2851	0.2865	0.2844	0.2768	0.2798	0.2762	0.2759	0.2755	0.2757

**Table A5.21:** Sensitivity analysis for the ANN trained using the field acquired data set and the 20 metre DEM variables for the nine class classification.

		Slope Angle	Aspect	Classified Vegetation	Field Sodicity Meter	Geology
Train	Rank	1	4	3	5	2
	Error	0.2879	0.2773	0.2812	0.2748	0.2825
Verify	Rank	2	4	3	5	1
	Error	0.2789	0.2755	0.2782	0.2709	0.2803

**Table A5.22:** Sensitivity analysis for the ANN trained using the field acquired and classified vegetation data set for the nine class classification.

	1.4.4	Slope Angle	Aspect	Class. Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	1	4	3	6	2	7	5	9	8
	Error	0.2852	0.2752	0.2777	0.2716	0.2806	0.2716	0.2718	0.2713	0.2713
Verify	Rank	2	4	3	7	1	5	9	6	8
	Error	0.2891	0.2821	0.2864	0.2812	0.2896	0.2819	0.2805	0.2818	0.2808

**Table A5.23:** Sensitivity analysis for the ANN trained using the field acquired data set, classified vegetation and the 10 metre DEM variables for the nine class classification.

		Slope Angle	Aspect	Class. Veg.	Field Sodicity Meter	Geology	Flow Length	Flow Acc.	Plan Curv.	Profile Curv.
Train	Rank	3	2	5	4	1	6	8	7	9
	Error	0.2733	0.2742	0.2687	0.2721	0.2788	0.2643	0.2624	0.2627	0.2621
Verify	Rank	4	3	5	2	1	6	7	9	8
	Error	0.2769	0.2775	0.2764	0.2778	0.2780	0.2732	0.2722	0.2720	0.2722

**Table A5.24:** Sensitivity analysis for the ANN trained using the field acquired data set, classified vegetation and the 20 metre DEM variables for the nine class classification.

### **Decision Tree Classifiers Variable Importance**

	Variable Importance
Aspect	100
Flow Length	82.17
Slope Angle	73.07
Profile Curvature	15.34
Flow Accumulation	3.87
Plan Curvature	3.58

**Table A5.25:** Variable importance for the DTC trained using the 10 metre DEM data set for the two class classification.

	Variable Importance
Flow Length	100
Flow Accumulation	81.70
Slope Angle	75.34
Aspect	35.08
Profile Curvature	25.71
Plan Curvature	3.38

**Table A5.26:** Variable importance for the DTC trained using the 20 metre DEM data set for the two class classification.

	Variable Importance
Slope Angle	100
Geology	59.61
Aspect	58.89
Estimated Vegetation	40.53
Sodicity Meter	32.56

**Table A5.27:** Variable importance for the DTC trained using the field acquired data set for the two class classification.

NOP WEET WRITE AL	Variable Importance
Slope Angle	100
Aspect	65.11
Geology	63.23
Flow Length	44.98
Estimated Vegetation	42.9
Sodicity Meter	33.17
Profile Curvature	5.59
Plan Curvature	4.93
Flow Accumulation	2.57

**Table A5.28:** Variable importance for the DTC trained using the field acquired data set and the 10 metre DEM variables for the two class classification.

angeneral Receivers	Variable Importance
Slope Angle	100
Geology	67.5
Aspect	54.42
Estimated Vegetation	47.29
Sodicity Meter	32
Flow Length	28.15
Flow Accumulation	20.49
Profile Curvature	5.01
Plan Curvature	0

**Table A5.29:** Variable importance for the DTC trained using the field acquired data set and the 20 metre DEM variables for the two class classification.

Variable Importance
100
51.583
3.048
0.671
0

**Table A5.30:** Variable importance for the DTC trained using the field acquired data set and classified vegetation for the two class classification.

	Variable Importance
Slope Angle	100
Aspect	82
Classified Vegetation	66.326
Flow Length	43.069
Geology	39.680
Sodicity Meter	20.289
Plan Curvature	15.519
Flow Accumulation	8.833
Profile Curvature	0

**Table A5.31:** Variable importance for the DTC trained using the field acquired data set, classified vegetation and 10 metre DEM data set for the two class classification.

	Variable Importance
Slope Angle	100
Aspect	51.583
Geology	3.048
Classified Vegetation	0.671
Flow Accumulation	0.192
Flow Length	0.003
Sodicity Meter	0
Plan Curvature	0
Profile Curvature	0

**Table A5.32:** Variable importance for the DTC trained using the field acquired data set, classified vegetation and 20 metre DEM data set for the two class classification.

Hard States	Variable Importance
Flow Length	100
Slope Angle	95.85
Aspect	45.67
Profile Curvature	36.87
Flow Accumulation	27.84
Plan Curvature	19.02

**Table A5.33:** Variable importance for the DTC trained using the 10 metre DEM data set for the three class classification.

	Variable Importance
Flow Length	100
Flow Accumulation	78.4
Slope Angle	60.35
Aspect	24.23
Profile Curvature	15.64
Plan Curvature	9.48

**Table A5.34:** Variable importance for the DTC trained using the 20 metre DEM data set for the three class classification.

	Variable Importance
Slope Angle	100
Geology	69.2
Estimated Vegetation	57.21
Aspect	44.22
Sodicity Meter	0

**Table A5.35:** Variable importance for the DTC trained using the field acquired data set for the three class classification.

	Variable Importance
Slope Angle	100
Geology	80.58
Aspect	60.89
Estimated Vegetation	53.61
Sodicity Meter	51.15
Flow Length	34.53
Flow Accumulation	10.4
Profile Curvature	8.7
Plan Curvature	1.64

**Table A5.36:** Variable importance for the DTC trained using the field acquired data set and the 10 metre DEM variables for the three class classification.

CINCAR MINING STREET	Variable Importance
Slope Angle	100
Geology	75.16
Aspect	59.61
Estimated Vegetation	55.02
Flow Length	49.2
Sodicity Meter	45.9
Flow Accumulation	13.53
Plan Curvature	6.85
Profile Curvature	3.05

**Table A5.37:** Variable importance for the DTC trained using the field acquired data set and the 20 metre DEM variables for the three class classification.

	Variable Importance
Slope Angle	100
Aspect	84.304
Geology	61.233
Classified Vegetation	58.099
Sodicity Meter	53.274

**Table A5.38:** Variable importance for the DTC trained using the field acquired data set and classified vegetation for the three class classification.

	Variable Importance
Slope Angle	100
Geology	96.17
Aspect	65.03
Classified Vegetation	63.56
Sodicity Meter	61.15
Flow Length	36.17
Flow Accumulation	18.22
Profile Curvature	10.4
Plan Curvature	6.42

**Table A5.39:** Variable importance for the DTC trained using the field acquired data set, classified vegetation and 10 metre DEM data set for the three class classification.

	Variable Importance
Slope Angle	100
Aspect	96.167
Classified Vegetation	65.025
Geology	63.560
Sodicity Meter	59.862
Flow Length	56.978
Profile Curvature	18.233
Flow Accumulation	16.533
Plan Curvature	7.858

**Table A5.40:** Variable importance for the DTC trained using the field acquired data set, classified vegetation and 20 metre DEM data set for the three class classification.

	Variable Importance
Flow Accumulation	100
Profile Curvature	55.69
Plan Curvature	41.48
Slope Angle	32.61
Aspect	26.36
Flow Length	13.67

**Table A5.41:** Variable importance for the DTC trained using the 10 metre DEM data set for the nine class classification.

	Variable Importance
Flow Length	100
Flow Accumulation	92.51
Aspect	81.3
Slope Angle	73.1
Profile Curvature	23.3
Plan Curvature	14

**Table A5.42:** Variable importance for the DTC trained using the 20 metre DEM data set for the nine class classification.

	Variable Importance
Slope Angle	100
Aspect	75.376
Estimated Vegetation	63.689
Sodicity Meter	60.20
Geology	54.797

**Table A5.43:** Variable importance for the DTC trained using the field acquired data set for the nine class classification.

et dealther version of esti-	Variable Importance
Flow Accumulation	100
Slope Angle	95.079
Estimated Vegetation	72.411
Aspect	53.27
Geology	42.678
Profile Curvature	41.033
Plan Curvature	31.096
Flow Length	29.936
Sodicity Meter	29.756

**Table A5.44:** Variable importance for the DTC trained using the field acquired data set and the 10 metre DEM variables for the nine class classification.

	Variable Importance
Slope Angle	100
Flow Length	86.86
Aspect	72.12
Estimated Vegetation	67.13
Flow Accumulation	63.80
Sodicity Meter	58.74
Geology	54.92
Profile Curvature	9.73
Plan Curvature	6.43

**Table A5.45:** Variable importance for the DTC trained using the field acquired data set and the 20 metre DEM variables for the nine class classification.

	Variable Importance
Slope Angle	100.00
Classified Vegetation	69.22
Sodicity Meter	65.19
Geology	61.69
Aspect	56.12

**Table A5.46:** Variable importance for the DTC trained using the field acquired data set and classified vegetation for the nine class classification.

	Variable Importance
Flow Accumulation	100
Slope Angle	95.079
Sodicity Meter	68.993
Aspect	53.270
Classified Vegetation	50.519
Geology	42.678
Profile Curvature	41.033
Plan Curvature	31.096
Flow Length	29.936

**Table A5.47:** Variable importance for the DTC trained using the field acquired data set, classified vegetation and 10 metre DEM data set for the nine class classification.

	Variable Importance
Slope Angle	100
Sodicity Meter	67.721
Flow Length	65.845
Geology	65.782
Classified Vegetation	60.593
Aspect	45.909
Flow Accumulation	41.127
Profile Curvature	16.035
Plan Curvature	0

**Table A5.48:** Variable importance for the DTC trained using the field acquired data set, classified vegetation and 20 metre DEM data set for the nine class classification.



# Appendix 6 – Soil Erosion and Risk Maps

Figure A6.1: Classified erosion map derived from the ANN trained using 10 metre DEM variables for a two class classification.



**Figure A6.2:** Classified erosion map drape derived from the ANN trained using 10 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.3: Classified erosion map derived from the ANN trained using 20 metre DEM variables for a two class classification.



**Figure A6.4:** Classified erosion map drape derived from the ANN trained using 20 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.5: Classified erosion map derived from the ANN trained using 10 metre DEM variables for a three class classification.



**Figure A6.6:** Classified erosion map drape derived from the ANN trained using 10 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.7: Classified erosion map derived from the ANN trained using 20 metre DEM variables for a three class classification.



**Figure A6.8:** Classified erosion map drape derived from the ANN trained using 20 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.9: Classified erosion map derived from the ANN trained using 20 metre DEM variables for a nine class classification.



**Figure A6.10:** Classified erosion map drape derived from the ANN trained using 20 metre DEM variables for a nine class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.11: Classified erosion map derived from the DTC trained using 10 metre DEM variables for a two class classification.
Erosion grid1 Value No Appreciable Erosion Erosion 3,000 Meters 1,500 750

**Figure A6.12:** Classified erosion map drape derived from the DTC trained using 10 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.13: Classified erosion map derived from the DTC trained using 20 metre DEM variables for a two class classification.



**Figure A6.14:** Classified erosion map drape derived from the DTC trained using 20 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.15: Classified erosion map derived from the DTC trained using 10 metre DEM variables for a three class classification.



**Figure A6.16:** Classified erosion map drape derived from the DTC trained using 10 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.17: Classified erosion map derived from the DTC trained using 20 metre DEM variables for a three class classification.



**Figure A6.18:** Classified erosion map drape derived from the DTC trained using 20 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.19: Classified erosion map derived from the DTC trained using 10 metre DEM variables for a nine class classification.



**Figure A6.20:** Classified erosion map drape derived from the DTC trained using 10 metre DEM variables for a nine class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.21: Classified erosion map derived from the DTC trained using 20 metre DEM variables for a nine class classification.



**Figure A6.22:** Classified erosion map drape derived from the DTC trained using 20 metre DEM variables for a nine class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.23: Classified erosion map derived from the DA trained using 10 metre DEM variables for a two class classification.



**Figure A6.24:** Classified erosion map drape derived from the DA trained using 10 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.25: Classified erosion map derived from the DA trained using 20 metre DEM variables for a two class classification.



**Figure A6.26:** Classified erosion map drape derived from the DA trained using 20 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.27: Classified erosion map derived from the DA trained using 10 metre DEM variables for a three class classification.

Appendices



**Figure A6.28:** Classified erosion map drape derived from the DA trained using 10 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.29: Classified erosion map derived from the DA trained using 20 metre DEM variables for a three class classification.



**Figure A6.30:** Classified erosion map drape derived from the DA trained using 20 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.31: Classified erosion map derived from the DA trained using 10 metre DEM variables for a nine class classification.



**Figure A6.32:** Classified erosion map drape derived from the DA trained using 10 metre DEM variables for a nine class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.33: Classified erosion map derived from the DA trained using 20 metre DEM variables for a nine class classification.



**Figure A6.34:** Classified erosion map drape derived from the DA trained using 20 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.35: Erosion probability map produced from the DTC trained with the 10 metre DEM variables for a two class classification.



**Figure A6.36:** Erosion probability map drape produced from the DTC trained with the 10 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.37: Erosion probability map produced from the DTC trained with the 10 metre DEM variables for a three class classification.



**Figure A6.38:** Erosion probability map drape produced from the DTC trained with the 10 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.39: Erosion probability map produced from the DTC trained with the 20 metre DEM variables for a three class classification.



**Figure A6.40:** Erosion probability map drape produced from the DTC trained with the 20 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.41: Erosion probability map produced from the DTC trained with the 20 metre DEM variables for a nine class classification.



Figure A6.42: Risk by association map produced from the DTC trained with the 10 metre DEM variables for a two class classification.



**Figure A6.42:** Risk by association map drape produced from the DTC trained with the 10 metre DEM variables for a two class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.43: Risk by association map produced from the DTC trained with the 10 metre DEM variables for a three class classification.



**Figure A6.44:** Risk by association map drape produced from the DTC trained with the 10 metre DEM variables for a three class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).



Figure A6.45: Risk by association map of subsurface erosion produced from the DTC trained with the 10 metre DEM variables for a nine class classification.



**Figure A6.46:** Risk by association map drape of subsurface erosion produced from the DTC trained with the 10 metre DEM variables for a nine class classification (Topographic map reproduced from the 1:25000 National Geographical Institute of Spain).