# CONTENTS

Example: California Data Set

Summary

Exercises

References

Transformations to Achieve Linearity

Box–Cox Transformations

# PREFACE

1	DIMENSION REDUCTION METHODS		1
日本のないのないのないのです。	Need for Dimension Reduction in Data Mining Principal Components Analysis Applying Principal Components Analysis to the <i>Houses</i> Data Set How Many Components Should We Extract? Profiling the Principal Components Communalities Validation of the Principal Components Factor Analysis Applying Factor Analysis to the <i>Adult</i> Data Set Factor Rotation User-Defined Composites Example of a User-Defined Composite Summary References Exercises		1 2 5 9 13 15 17 18 18 20 23 24 25 28
	Autopoint Charles and		2
2	REGRESSION MODELING	Using the Principal Compo	33
	Example of Simple Linear Regression Least-Squares Estimates Coefficient of Determination Standard Error of the Estimate Correlation Coefficient		<ul> <li>34</li> <li>36</li> <li>39</li> <li>43</li> <li>45</li> <li>46</li> </ul>
	Outliers, High Leverage Points, and Influential Observation Regression Model Inference in Regression t-Test for the Relationship Between $x$ and $yConfidence Interval for the Slope of the Regression LinConfidence Interval for the Mean Value of y Given x$	ons de solaris a significante de solaris a solaris a significante de solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris de solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris de solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solaris a solar	48 55 57 58 60 60
	Prediction Interval for a Randomly Chosen Value of y	Jiven x	61
	Example: <i>Baseball</i> Data Set		63 68

vii

74

79

83

84

86

86

xi

viii contents

3	MULTIPLE REGRESSION AND MODEL BUILDING	93
	Example of Multiple Regression	93
	Multiple Regression Model	99
	Inference in Multiple Regression	100
	<i>t</i> -Test for the Relationship Between y and $x_i$	101
	F-Test for the Significance of the Overall Regression Model	102
	Confidence Interval for a Particular Coefficient	104
	Confidence Interval for the Mean Value of y Given $x_1, x_2, \ldots, x_m$	105
	Prediction Interval for a Randomly Chosen Value of y Given $x_1, x_2, \ldots, x_m$	105
	Regression with Categorical Predictors	105
	Adjusting $R^2$ : Penalizing Models for Including Predictors That Are	
	Not Useful	113
	Sequential Sums of Squares	115
	Multicollinearity	116
	Variable Selection Methods	123
	Partial F-Test	123
	Forward Selection Procedure	125
	Backward Elimination Procedure	125
	Stepwise Procedure	126
	Best Subsets Procedure	126
	All-Possible-Subsets Procedure	126
	Application of the Variable Selection Methods	127
	Forward Selection Procedure Applied to the Cereals Data Set	127
	Backward Elimination Procedure Applied to the Cereals Data Set	129
	Stepwise Selection Procedure Applied to the Cereals Data Set	131
	Best Subsets Procedure Applied to the Cereals Data Set	131
	Mallows' $C_p$ Statistic	131
	Variable Selection Criteria	135
	Using the Principal Components as Predictors	142
	Summary	147
	References	149
	Exercises	149
4	LOGISTIC REGRESSION	155
	Simple Example of Logistic Regression	156
	Maximum Likelihood Estimation	158
	Interpreting Logistic Regression Output	159
	Inference: Are the Predictors Significant?	160
	Interpreting a Logistic Regression Model	162
	Interpreting a Model for a Dichotomous Predictor	163
	Interpreting a Model for a Polychotomous Predictor	166
	Interpreting a Model for a Continuous Predictor	170
	Assumption of Linearity	174
	Zero-Cell Problem	177
	Multiple Logistic Regression	179
	Introducing Higher-Order Terms to Handle Nonlinearity	183
	Validating the Logistic Regression Model	189
	WEKA: Hands-on Analysis Using Logistic Regression	194
	Summary	197

	References	
	Exercises	199
5	NAIVE BAYES ESTIMATION AND BAYESIAN NETWORKS	204
- 102	Bayesian Approach	204
	Maximum a Posteriori Classification	206
	Posterior Odds Ratio	210
	Balancing the Data	212
	Naive Bayes Classification	215
	Numeric Predictors	219
	WEKA: Hands-on Analysis Using Naive Bayes	223
	Bayesian Belief Networks	227
	Clothing Purchase Example	227
	Using the Bayesian Network to Find Probabilities	229
	WEKA: Hands-On Analysis Using the Bayes Net Classifier	232
	Summary	234
	References	236
	Exercises	237
6	GENETIC ALGORITHMS	240
	Introduction to Genetic Algorithms	240
	Basic Framework of a Genetic Algorithm	241
	Simple Example of a Genetic Algorithm at Work	243
	Modifications and Enhancements: Selection	245
	Modifications and Enhancements: Crossover	247
	Multipoint Crossover	247
	Uniform Crossover	247
	Genetic Algorithms for Real-Valued Variables	248
	Single Arithmetic Crossover	248
	Simple Arithmetic Crossover	248
	Whole Arithmetic Crossover	249
	Discrete Crossover	249
	Normally Distributed Mutation	249
	Using Genetic Algorithms to Train a Neural Network	249
	WEKA: Hands-on Analysis Using Genetic Algorithms	252
	Summary	261
	References	262
	Exercises	263
7	CASE STUDY MODELING RESPONSE TO DIRECT MAIL MARKETING	265
-	Cross-Industry Standard Process for Data Mining	265
	Business Understanding Phase	205
	Direct Mail Marketing Response Problem	267
	Building the Cost/Benefit Table	267
	Data Understanding and Data Preparation Phases	270
	Clothing Store Data Set	270
	Transformations to Achieve Normality or Symmetry	272
	Standardization and Flag Variables	276
		2.0

## X CONTENTS

Deriving New Variables	277
Exploring the Relationships Between the Predictors and	the Response 278
Investigating the Correlation Structure Among the Predi	ictors 286
Modeling and Evaluation Phases	289
Principal Components Analysis	292
Cluster Analysis: BIRCH Clustering Algorithm	294
Balancing the Training Data Set	298
Establishing the Baseline Model Performance	299
Model Collection A: Using the Principal Components	300
Overbalancing as a Surrogate for Misclassification Cost	s 302
Combining Models: Voting	304
Model Collection B: Non-PCA Models	306
Combining Models Using the Mean Response Probabili	ities 308
Summary	312
References	316

### **INDEX**

317