

CONTENTS

1. EXPLORATORY DATA ANALYSIS

1.1 EDA FOR INDIVIDUAL VARIABLES

- 1.1.1 SUMMARIZATION: MEASURES OF CENTRAL TENDANCY, MEASURES OF DISPERSION
- 1.1.2 DATA VISUALIZATION: HISTOGRAM/BAR CHART, BOX PLOT, STEM AND LEAF DISPLAY
- 1.1.3 OUTLIER DETECTION, SYMMETRY CHECKING
- 1.1.4 MISSING VALUE IMPUTATION
- 1.1.5 TESTING FOR NORMALITY: HISTOGRAM, KUTOSIS SE, QQ PLOT, KS TEST AND SW TEST

1.2 EDA FOR MULTIPLE VARIABLES

- 1.2.1 PAIRWISE SCATTER PLOTS
- 1.2.2 CORRELATION ANALYSIS: PEARSON'S AND SPEARMAN'S CORRELATION AND THEIR SIGNIFICANCE

2. REGRESSION ANALYSIS

2.1 SIMPLE MODEL BUILDING

- 2.1.1 FITTING A LINEAR REGRESSION MODEL
- 2.1.2 TESTING THE SIGNIFICANCE OF INDIVIDUAL REGRESSORS AND OVERALL REGRESSION
- 2.1.3 R SQUARE AND ADJUSTED R SQUARE

2.2 MULTICOLLINEARITY

- 2.2.1 PROBLEM AND ITS CONSEQUENCES
- 2.2.2 DETECTION AND REMOVAL OF MULTICOLLINEARITY USING CORRELATION ANALYSIS
- 2.2.3 DETECTION AND REMOVAL OF MULTICOLLINEARITY USING VARIANCE INFLATION FACTORS (VIFs)

2.3 PARSIMONIOUS MODELLING OR MODEL SELECTION

- 2.3.1 FORWARD SELECTION
- 2.3.2 BACKWARD ELIMINATION
- 2.3.3 STEPWISE SELECTION

2.4 VALIDATION OF ASSUMPTIONS AND RESIDUAL ANALYSIS

- 2.4.1 LINEARITY OF REGRESSION
- 2.4.2 AUTOCORRELATION
- 2.4.3 HETEROSCEDASTICITY
- 2.4.4 NORMALITY OF ERRORS
- 2.4.5 INFLUENTIAL OBSERVATIONS, LEVERAGE AND OUTLIERS

3. CLASSIFICATION PROBLEM

3.1 BINARY CLASSIFICATION

3.1.1 EMAIL SPAM FILTERING

3.1.1.1 EMAIL SPAM FILTERING USING NAÏVE BAYES' CLASSIFIER

3.1.1.2 EMAIL SPAM FILTERING USING LOGISTIC REGRESSION

3.1.2 PREDICTION OF CANCER FROM SMOKING USING LOGISTIC REGRESSION

3.1.3 SKULL TYPE PREDICTION USING LOGISTIC REGRESSION

3.1.4 SENTIMENT ANALYSIS USING LOGISTIC REGRESSION – WHAT MAKES A US PRESIDENTIAL CANDIDATE WIN?

3.1.5 SKULL TYPE PREDICTION USING DISCRIMINANT ANALYSIS

3.1.6 COMPARATIVE STUDY OF SKULL TYPE PREDICTION USING LOGISTIC REGRESSION AND DISCRIMINANT ANALYSIS

3.2 MULTICLASS CLASSIFICATION

3.2.1 MULTICLASS CLASSIFICATION BY DECOMPOSING A MULTICLASS PROBLEM INTO SEVERAL BINARY CLASSIFICATION PROBLEMS – FLOWER SPECIES PREDICTION

3.2.2 MULTICLASS CLASSIFICATION USING MULTINOMIAL LOGISTIC REGRESSION – FLOWER SPECIES PREDICTION

3.2.3 COMPARATIVE STUDY OF THE TWO APPROACHES OF MULTICLASS CLASSIFICATION FOR THE FLOWER SPECIES PREDICTION PROBLEM

4. NON – PARAMETRIC INFERENCE

4.1 FRANK WILCOXON SIGN TEST

4.1.1 ONE SAMPLE SIGN TEST (BINOMIAL TEST)

4.1.1.1 TESTING THE HYPOTHETICAL VALUE OF POPULATION MEDIAN OF WEIGHT AND HEIGHT OF FEMALES

4.1.1.2 TESTING OF EQUALITY OF WIN/LOSS CHANCES OF BASKETBALL MATCHES

4.1.2 TWO SAMPLE SIGN TEST (SIGN TEST) – TESTING FOR SIMILARITY IN GRADES GIVEN BY TWO DIFFERENT PROFESSORS TO SAME SET OF STUDENTS.

4.2 WALD – WOLFOWITZ RUN TEST

4.2.1 TESTING FOR EQUALITY OF PRICE OF A COMMODITY IN TWO DIFFERENT CITIES

4.2.2 TESTING THE RANDOMNESS OF A GIVEN SAMPLE

4.3 MANN – WHITNEY WILCOXON U – TEST

4.3.1 TESTING THE EFFECTIVENESS OF A PROGRAMMED WORK BOOK TEACHING METHOD

4.4 KOLMOGOROV – SMIRNOV (KS) TEST

4.4.1 ONE SAMPLE KS TEST

4.4.1.1 TESTING IF THE SAMPLE IS DRAWN FROM SOME NORMAL DISTRIBUTION

4.4.1.2 TESTING IF THE SAMPLE IS DRAWN FROM SOME EXPONENTIAL DISTRIBUTION

4.4.1.3 TESTING IF THE SAMPLE IS DRAWN FROM SOME POISSON DISTRIBUTION

4.4.1.4 TESTING IF THE SAMPLE IS DRAWN FROM SOME UNIFORM DISTRIBUTION

4.4.2 TWO SAMPLE SIGN TEST (SIGN TEST) – TESTING IF TWO SAMPLES ARE DRAWN FROM SAME POPULATION

4.5 CHI – SQUARE TESTS

4.5.1 KARL PEARSON’S GOODNESS OF FIT – TESTING FOR EQUALITY OF CHANCES OF MALE AND FEMALE BIRTH

4.5.2 INDEPENDENCE OF ATTRIBUTES – TESTING THE SEX DISCRIMINATION IN EMPLOYMENT

4.5.3 MC – NEMAR’S TEST – TESTING THE EFFECTIVENESS OF A DRUG (MARGINAL HOMOGENEITY OF BEFORE-AND-AFTER DRUG MEASUREMENTS)

4.5.4 COCHRAN – MANTEL – HEENSZEL TEST – TESTING FOR EQUALITY OF PRESENCE OF LAP ALLELES IN MARINE AND ESTUARINE HABITAT AFTER CONTROLLING FOR THE AREA/REGION.

5. REFERENCES