

MATH 424: Final Homework

Due on Wednesday, December 20, 2017

Kafai 11:10am

Jonathan Dombrowski

Contents

Q1. Consider the data in q1data.txt. All subjects are asthmatics. For the model with Forced Expiratory Volume (FEV) as the response and Height, Weight, and Age as the predictors,

- a) Examine a plot of the studentized or jackknife residuals versus the predicted values. Are any regression assumption violations apparent? If so, suggest possible remedies.
- b) Examine numerical descriptive statistics, histograms, box-and-whisker plots, and normal probability plots of jackknife residuals. Is the normality assumption violated? If so, suggest possible remedies.
- c) Examine outlier diagnostics, including Cook's distance, leverage statistics, and jackknife residuals, and identify any potential outliers. What course of action, if any, should be taken when outliers are identified?
- d) Examine variance inflation factors, condition indices (unadjusted and adjusted for the intercept), and variance proportions. Are there any important collinearity problems? If so, suggest possible remedies.

Q2. A random sample of data was collected on residential sales in a large city. The data in q2data.txt shows the selling price (Y, in \$1,000s), area (x1, in hundreds of square feet), number of bedrooms (X2), total number of rooms (X3), house age (X4, in years), and location (Z = 0 for in-town and inner suburbs, Z=1 for outer suburbs). In parts a through c, use variables X1, X2, X3, X4, and Z as the predictor variables.

- a) Use the all possible regressions procedure to suggest a best model.
- b) Use the stepwise regression algorithm to suggest a best model.
- c) Use the backward elimination algorithm to suggest a best model.
- d) Which of the models selected in a, b, and c seems to be the best model, and why?

Q3. The data listed in q3data.txt relate to a study by Reiter and others concerning the effects of injecting triethyl-tin (TET) into rats once at age 5 days. The animals were injected with 0, 3, or 6 mg per kilogram of body weight. The response was the log of the activity count, log (ac), for 1 hour, recorded at 21 days of age. The rat was left to move about freely in a figure 8 maze. Analysis of other studies with this type of activity count confirms that log counts should yield Gaussian errors if the model is correct.

- a) Conduct a two-way ANOVA with SEX and DOSAGE as factors.
- b) Using $\alpha = .05$, report your conclusions based on the ANOVA.
- c) Which, if any, families of means should be followed up with multiple-comparison tests? What type of comparisons would you recommend?

Q4. The data in q4data.txt is the record of coronary artery disease (ca, 0=no, 1=yes), age, ECG (0, 1, and 2 based on the reading of ST segment depression), and sex (0=male, 1=female). Based on this model

- a) What is the estimated logistic regression model for the relationship between ca and age, ECG, sex?
- b) What is a 30-year-old male, ECG=2 predicted probability of having coronary artery disease?
- c) Estimate the odds ratio comparing a 30-year-old male, ECG=2 to a 31-year-old male, ECG=2. Interpret this estimated odds ratio.
- d) Find a 95% confidence interval for the population odds ratio being estimated in part (c).