

## STA 2101/442 Assignment Nine<sup>1</sup>

These questions are just practice for the final exam, and are not to be handed in.

1. In a study comparing the effectiveness of different exercise programmes, volunteers were randomly assigned to one of three exercise programmes ( $A$ ,  $B$ ,  $C$ ) or put on a waiting list and told to work out on their own. Aerobic capacity is the body's ability to process oxygen. Aerobic capacity was measured before and after 6 months of participation in the program (or 6 months of being on the waiting list). The response variable was improvement in aerobic capacity. The explanatory variables were age (a covariate) and treatment group. *Treatment group includes the waiting list control condition.*
  - (a) First consider a regression model with an intercept, and no interaction between age and treatment group.
    - i. Make a table showing how you would set up indicator dummy variables for treatment group. Make Waiting List the reference category
    - ii. Write the regression model. Please use  $x$  for age, and make its regression coefficient  $\beta_1$ .
    - iii. In terms of  $\beta$  values, what null hypothesis would you test to find out whether, allowing for age, the three exercise programmes differ in their effectiveness?
    - iv. Write the null hypothesis for the preceding question as  $H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{0}$ . Just give the  $\mathbf{L}$  matrix.
    - v. In terms of  $\beta$  values, what null hypothesis would you test to find out whether Programme  $B$  was better than the waiting list?
    - vi. In terms of  $\beta$  values, what null hypothesis would you test to find out whether Programmes  $A$  and  $B$  differ in their effectiveness?
    - vii. Suppose you wanted to estimate the difference in average benefit between programmes  $A$  and  $C$  for a 27 year old participant. Give your answer in terms of  $\beta$  values.
    - viii. Is it safe to assume that age is independent of the other explanatory variables? Answer Yes or No and briefly explain.
  - (b) Now consider a regression model with an intercept and the interaction (actually a set of interactions) between age and treatment.
    - i. Write the regression model. Make it an extension of your earlier model.
    - ii. Suppose you wanted to know whether the slopes of the 4 regression lines were equal. In terms of  $\beta$  values, what null hypothesis would you test?
    - iii. Suppose you wanted to know whether any differences among mean improvement in the four treatment conditions depends on the participant's age. In terms of  $\beta$  values, what null hypothesis would you test?
    - iv. Write the null hypothesis for the preceding question as  $H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{0}$ . Just give the  $\mathbf{L}$  matrix. It is  $r \times p$ . What is  $r$ ? What is  $p$ ?

---

<sup>1</sup>Copyright information is at the end of the last page.

- v. Suppose you wanted to know whether the difference in effectiveness between Programme *A* and the Waiting List depends on the participant's age. In terms of  $\beta$  values, what null hypothesis would you test?
  - vi. Suppose you wanted to *estimate* the difference in average benefit between programmes *A* and *C* for a 27 year old participant. Give your answer in terms of  $\hat{\beta}$  values.
2. Telephone sales representatives use computer software to help them locate potential customers, answer questions, take credit card information and place orders. Twelve sales representatives were randomly assigned to each of three new software packages the company was thinking of purchasing. The data for each sales representative include the software package (1, 2 or 3), sales last quarter with the old software, and sales this quarter with one of the new software packages. Sales are in number of units sold.

The data are in the file

<http://www.utstat.toronto.edu/~brunner/data/legal/sales.data.txt>.

The explanatory and response variables are what you would think.

- (a) Fit a full model in which the slopes and intercepts of the regression lines relating sales last quarter to sales this quarter might depend on the kind of software the sales representatives are using.
- (b) Carry out an ordinary *F*-test to determine whether the effect of software type on sales depends on the representative's performance last quarter. Be able to state your conclusion in plain, non-statistical language.
- (c) Estimate the slopes of the three regression lines. Base the estimates on numbers from your printout. I don't see how you can do this without making a table.
- (d) Carry out tests to answer these questions. If they are already on the output of `summary`, use that.
  - i. Are the slopes for Software 1 and 2 different?
  - ii. Are the slopes for Software 1 and 3 different?
  - iii. Are the slopes for Software 2 and 3 different?

Protecting the three tests with a Bonferroni correction at the joint 0.05 significance level, what do you conclude? Plain language is not necessary, but you should say what happened.

- (e) The average (sample mean) performance last quarter was 76.56 (please use exactly this number). We are interested in whether the three software packages differ in their effectiveness for sales representatives with average performance last quarter.
  - i. Estimate expected performance this quarter for sales representatives with average performance last quarter. Calculate the estimates with R.
  - ii. State the null hypothesis in symbols.
  - iii. Carry out the *F*-test.
  - iv. In plain language, what do you conclude?

3. Arsenic is a powerful poison, which is why it has been used on farms for many years to kill insects. Even in very small amounts, arsenic can cause cancer in humans, and recently it has been found that rice and foods made from rice tend to be very high in arsenic. Brown rice is worse, by the way.

In a controlled experiment, pots of rice were prepared by either washing the rice first or not, and by cooking the rice in either a low, a medium or a high amount of water. The response variable is amount of arsenic in the cooked rice.

- (a) Use a regression model with *cell means coding*. That's the model with no intercept, and one indicator dummy variable for each treatment combination. You don't have to say how the dummy variables are defined. That will become clear in the next part. Just give the regression equation.
- (b) Write the expected amounts of arsenic in the table below, in terms of the  $\beta$  parameters of your model.

	Amount of Water		
	Low	Medium	High
Washed			
Unwashed			

- (c) If you wanted to test whether the effect of washing the rice depended on how much water you cook it in, what is the null hypothesis? Give your answer in terms of the  $\beta$  values in your model.
- (d) If you wanted to test whether washing the rice before cooking has any effect if the rice is cooked in a lot of water, what is the null hypothesis? Give your answer in terms of  $\beta$  values.
- (e) Suppose you want to test whether the amount of water used to cook the rice makes any difference if the rice has been washed. What is the null hypothesis? Give your answer in terms of  $\beta$  values.
- (f) Averaging across different amounts of water used to cook the rice, does pre-washing affect the amount of arsenic in the rice. What null hypothesis would you test to answer this question? Give your answer in terms of  $\beta$  values.
- (g) If you wanted to test whether the effect of the amount of water used to cook the rice depends on whether you wash it first, what is the null hypothesis? Give your answer in terms of  $\beta$  values.

4. Consider a two-factor analysis of variance in which each factor has two levels. Use this regression model for the problem:

$$Y_i = \beta_0 + \beta_1 d_{i,1} + \beta_2 d_{i,2} + \beta_3 d_{i,1} d_{i,2} + \epsilon_i,$$

where  $d_{i,1}$  and  $d_{i,2}$  are dummy variables.

- (a) Make a two-by-two table showing the four treatment means in terms of  $\beta$  values. Use *effect coding*. In terms of the  $\beta$  values, state the null hypothesis you would use to test for
- Main effect of the first factor
  - Main effect of the second factor
  - Interaction
- (b) Make a two-by-two table showing the four treatment means in terms of  $\beta$  values. Use *indicator dummy variables* (zeros and ones). In terms of the  $\beta$  values, state the null hypothesis you would use to test for
- Main effect of the first factor
  - Main effect of the second factor
  - Interaction
- (c) Which dummy variable scheme do you like more?
5. Effect coding is not the most convenient choice for every purpose. Please consider the two-factor rotten potato example from lecture set 20, the one entitled “Interactions and Factorial ANOVA.” Lecture slide 39 has the expected response for each treatment combination under effect coding. Suppose you wanted to test whether Bacteria Type has an effect just for cool temperatures.
- Write the null hypothesis in scalar form. Simplify as much as possible.
  - Give the  $\mathbf{L}$  matrix in  $H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{h}$ .
6. In that last question, your answer may be different from mine, but we both could be right. Let  $\mathbf{A}$  be an  $r \times r$  matrix with an inverse, and suppose the null hypotheses is written  $H_0 : \mathbf{A}\mathbf{L}\boldsymbol{\beta} = \mathbf{A}\mathbf{h}$  instead of  $H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{h}$ . Show that the  $F^*$  statistic for the general linear test is unaffected.

7. I know this is pretty gruesome, but the data are real. An experiment in dentistry seeks to test the effectiveness of a drug (HEBP) that is supposed to help dental implants become more firmly attached to the jaw bone. This is an initial test on animals. False teeth were implanted into the leg bones of rabbits, and the rabbits were randomly assigned to receive either the drug or a saline solution (placebo). Technicians administering the drug were blind to experimental condition.

Rabbits were also randomly assigned to be "sacrificed" after either 3, 6, 9 or 12 days. At that time, the implants were pulled out of the bone by a machine that measures force in newtons and stiffness in newtons/mm. For both of these measurements, higher values indicate more healing. A measure of "pre-load stiffness" in newtons/mm is also available for each animal. This may be another indicator of how firmly the false tooth was implanted into the bone, but it might even be a covariate. Nobody can seem to remember what "preload" means, so we'll ignore this variable for now.

The explanatory variables are Time and Drug. The response variable is Force required to pull out the tooth. There is more than one reasonable way to do this analysis, but just to keep us together please treat Time as a categorical variable.

The data are available in the file [bunnies.data.txt](#). The variables are

- Identification code
  - Time (3,6,9,12 days of healing)
  - Drug (1=HEBP, 0=saline solution)
  - Stiffness in newtons/mm
  - Force in newtons
  - Preload stiffness in newtons/mm
- (a) Use `table` to find out how many rabbits are in each experimental condition.
- (b) Carry out the standard tests of main effects and interactions. Be prepared to answer the following questions about each test.
- i. What is the value of the test statistic? The answer is a number from your printout.
  - ii. What is the  $p$ -value? The answer is a number from your printout.
  - iii. Do you reject the null hypothesis at the 0.05 level? Yes or No.
  - iv. What, if anything, do you conclude? This is not the place for statistical jargon. "What do you conclude" means say something about the drug, healing, time – something like that.
- (c) I know this is a bit redundant with the preceding question, but *averaging across time, did the drug help the teeth become more firmly attached to the bone?* If the results justify an answer, then answer Yes or No.
- (d) Make a table with a row for each treatment combination. Make columns showing the dummy variables for effect coding.
- (e) Give  $E[Y|\mathbf{X} = \mathbf{x}]$  for a regression model with both main effects and the interaction. Use your variable names from the preceding question.

- (f) In terms of the  $\beta$  values of your regression model, give the null hypothesis you would test in order to answer each of the following questions.
- i. Averaging across time periods, is there a difference between the drug and placebo in mean force required to extract the tooth?
  - ii. Averaging across drug and placebo, does elapsed time affect the mean force required to extract the tooth?
  - iii. Does the effect of the drug depend upon elapsed time?
- (g) Now please return to R. Doing it the easiest way you can, conduct tests to answer the following questions. Just do regular one-at-a-time (custom) tests. Don't bother with any Bonferroni correction this time. Just consider one response variable: Force. As usual, we are guided by the  $\alpha = 0.05$  significance level.
- i. Are the marginal means different at 3 and 6 days?
  - ii. Are the marginal means different at 6 and 9 days?
  - iii. Are the marginal means different at 9 and 12 days?
  - iv. Is there a difference between Drug and Placebo just at 3 days?
  - v. Is there a difference between Drug and Placebo just at 6 days?
  - vi. Is there a difference between Drug and Placebo just at 9 days?
  - vii. Is there a difference between Drug and Placebo just at 12 days?
  - viii. Be able to answer questions like these for each test:
    - A. What is the value of the test statistic? The answer is a number from your printout.
    - B. What is the  $p$ -value? The answer is a number from your printout.
    - C. Do you reject the null hypothesis at the 0.05 level? Yes or No.
    - D. What, if anything, do you conclude? This is not the place for statistical jargon. "What do you conclude" means say something about the drug, healing, time – something like that.

---

This assignment was prepared by [Jerry Brunner](#), Department of Statistics, University of Toronto. It is licensed under a [Creative Commons Attribution - ShareAlike 3.0 Unported License](#). Use any part of it as you like and share the result freely. The L<sup>A</sup>T<sub>E</sub>X source code is available from the course website: <http://www.utstat.toronto.edu/~brunner/oldclass/appliedf18>