# UNIVERSITY OF TORONTO
## Faculty of Arts and Science

December 2011 Final Examination
### STA442H1F/2101HF
Methods of Applied Statistics

Jerry Brunner

Duration - 3 hours

Aids: Calculator Model(s): Any calculator is okay

**Last/Surname** (Print): _____

**First/Given Name** (Print): _____

**Student Number:** _____

**Signature:** _____

| Qn. # | Value | Score |
|:-----:|:-----:|:-----:|
| 1 | 20 | |
| 2 | 20 | |
| 3 | 10 | |
| 4 | 15 | |
| 5 | 5 | |
| 6 | 12 | |
| 7 | 8 | |
| 8 | 10 | |
| Total = 100 Points | | |

20 points

1. Independently for $i = 1, \ldots, n$, let

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i,$$

where $E(X_i) = \mu$, $Var(X_i) = \sigma_x^2$, $E(\epsilon_i) = 0$, $Var(\epsilon_i) = \sigma_\epsilon^2$, and $\epsilon_i$ is independent of $X_i$. Let

$$\widehat{\beta}_1 = \frac{\sum_{i=1}^{n} X_i Y_i}{\sum_{i=1}^{n} X_i^2}.$$

(a) Is $\widehat{\beta}_1$ a consistent estimator of $\beta_1$? Answer Yes or No and **Circle Yes or No**. Prove your answer.

(b) For some special cases we have $\widehat{\beta}_1 \overset{a.s.}{\rightarrow} \beta_1$. When does this happen?

*20 points*

2. Let $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, where $\mathbf{X}$ is an $n \times p$ matrix of known constants, $\boldsymbol{\beta}$ is a $p \times 1$ vector of unknown constants, and $\boldsymbol{\epsilon}$ is multivariate normal with mean zero and covariance matrix $\sigma^2 \mathbf{I}_n$, where $\sigma^2 > 0$ is an unknown constant. You may use without proof the fact that if $\mathbf{T} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then $\mathbf{AT} \sim N(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}')$.

   (a) In the regression model, what is the distribution of $\mathbf{Y}$? No proof is needed.

   (b) The maximum likelihood estimate of $\boldsymbol{\beta}$ is $\widehat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ (no proof needed). What is the distribution of $\widehat{\boldsymbol{\beta}}$? Show the calculations.

*10 points*  3. In a study of math education in elementary school, equal numbers of boys and girls were randomly assigned to one of three training programmes designed to improve spatial reasoning. After five school days of training, the students were given a standardized test of spatial reasoning. Score on the spatial reasoning test is the response variable. You will define a regression model for this factorial analysis of variance. Don't write the model yet.

(a) In the table below, show how your dummy variables are defined. *Use effect coding.* Write the name of each dummy variable at the head of its column.

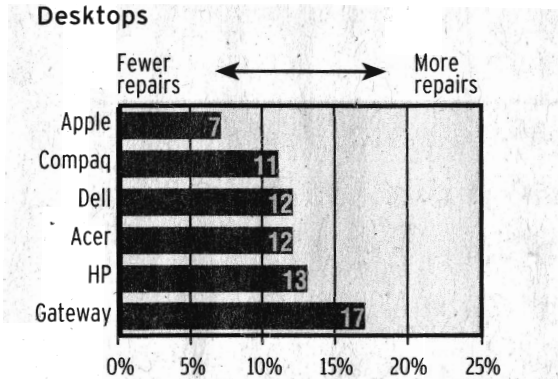| | | | |
|---|---|---|---|
| Girls, Programme 1 | | | |
| Girls, Programme 2 | | | |
| Girls, Programme 3 | | | |
| Boys, Programme 1 | | | |
| Boys, Programme 2 | | | |
| Boys, Programme 3 | | | |

(b) Give $E[Y_i|\mathbf{X}_i = \mathbf{x}_i]$ for the full model. Include the interaction terms.

(c) Suppose you want to test whether, averaging across training programmes, there is a difference between girls and boys in their average performance on the spatial reasoning test. State the null hypothesis in terms of the $\beta$ values in your model.

(d) Suppose you want to test whether the average sex difference in performance depends on which training programme the children are in. State the null hypothesis in terms of the $\beta$ values in your model.

*15 points*

4. The following was scanned from the 2012 *Consumer Reports Buying Guide.* It shows repair history for several desktop computer brands, and seems to be based on a fairly sophisticated statistical analysis.

**Desktops**

| | Fewer repairs ← → More repairs |
|---|---|
| Apple | 7 |
| Compaq | 11 |
| Dell | 12 |
| Acer | 12 |
| HP | 13 |
| Gateway | 17 |

0%   5%   10%   15%   20%   25%

Gateway has been among the more repair-prone brands of desktop computers and Apple has been among the least. That's what we found when we asked almost 20,000 readers who bought a desktop computer between 2008 and 2011 about their experiences. The graph shows the percentage of models for each brand that were repaired or had a serious problem. Differences of fewer than 5 points aren't meaningful, and we've adjusted the data to eliminate differences linked solely to the age of the desktop computer. Models within a brand can vary, and design or manufacture changes might affect future reliability.

(a) What statistical method do you guess they employed?

(b) What words above suggest that they did some kind of hypothesis tests?

(c) What was the covariate?

(d) How would you set up dummy variables for the Brand of Computer? *Make a table in the space beside the scanned material.*

(e) Assuming there is no interaction, write an expression (a function of the $\widehat{\beta}$ values) that would give you the number 17 for Gateway computers. Denote the covariate by $x$.

*5 points*

5. On the `Computer Printout`, the output for the `Beta data` is based on a random sample of size $n = 50$ from a beta distribution.

   (a) What is the maximum likelihood estimate $\widehat{\alpha}$? The answer is a number from the printout.

   (b) Carry out a test of $H_0 : \alpha = 5$ versus $H_1 : \alpha \neq 5$. Calculate the test statistic. The answer is a number. Show a little work. **Circle the number.** Do you reject $H_0$ at $\alpha = 0.05$? **Answer Yes or No.** You have more room than you need.

*12 points*

6. On the `Computer Printout`, the output for the `Cars data` is based on the same metric cars data you analyzed for homework. Potentially there are three regression lines relating weight of car to fuel efficiency.

(a) We wish to know whether country differences in fuel efficiency depend on the weight of the car. Fill in the table below.

| $F$ Statistic | Degrees of Freedom (2 numbers) | $p$-value | Reject $H_0$ at $\alpha = 0.05$? (Yes or No) |
|---|---|---|---|
|  |  |  |  |

(b) Do country differences in fuel efficiency depend on the weight of the car? Answer Yes or No.

(c) Are the three regression lines parallel *in the population*? Answer Yes or No.

(d) What is the estimated expected fuel efficiency for a U.S. car of average weight (meaning average for the entire sample)? The answer is a single number.

(e) What is the estimated slope of the regression line for U.S. cars? The answer is a single number.

(f) To show which slopes are different from one another, make a table whose $i, j$ element is the Bonferroni-adjusted $p$-value for the tests of difference between the slope for country $i$ and country $j$. Just fill in the upper triangular part of the table. Use your calculator to convert $p$-values on the printout to Bonferroni-adjusted $p$-values.

(g) Based on the multiple comparisons, which slopes are really different? Don't just say they're different; say which one is steeper (going down faster).

*8 points*

7. On the `Computer Printout`, the output for the `Birth Weight Data` is based on the same data discussed in lecture.

   (a) The estimated odds of a low birth weight baby are _____ times as great for a mother with a history of premature labour. The answer is a single number; write it on the line.

   (b) You want to know whether *any* of the variables in the model are related to the chances that a baby will have low birth weight. Fill in the table below.

| Wald $\chi^2$ Statistic | Degrees of Freedom | $p$-value | Reject $H_0$ at $\alpha = 0.05$? (Yes or No) |
|---|---|---|---|
|  |  |  |  |

   (c) Does the preceding test indicate that at least one of the explanatory variables is related to the response variable? Answer Yes or No.

   (d) Give an approximate 95% confidence interval for the regression coefficient corresponding to mother's weight. Your answer is a pair of numbers, and there is more than one way to calculate them from the numbers on the printout. You have more room than you need.

*10 points*

8. On the `Computer Printout`, the output for the `Dichotic Listening Study` comes from data given in the Final Assignment. Focus on this question. Does mode of presentation (Left *versus* Right *versus* Both) influence performance?

   (a) There is a single test for the question of interest. Record the numbers below.

   | $F$ Statistic | Degrees of Freedom (2 numbers) | $p$-value | Reject $H_0$ at $\alpha = 0.05$? (Yes or No) |
   |---|---|---|---|
   |  |  |  |  |

   (b) Only one of the pairwise comparisons of marginal means is statistically significant using a multiple comparison method. Give the Bonferroni-adjusted $p$-value. The answer is a single number.

   (c) Describe the difference in simple, non-statistical language. Be sure you say which mean is bigger.