

Student Name: \_\_\_\_\_

ID: \_\_\_\_\_

McGill University  
Faculty of Science  
Department of Mathematics and Statistics

Statistics Part A Comprehensive Exam  
Methodology Paper

Date: Friday, May 16, 2014

Time: 13:00 17:00

Instructions

Answer only two questions out of Section L. If you answer more than two questions, then only the ~~the~~ FIRST TWO questions will be marked.

Answer only two questions out of Section G. If you answer more than two questions, then only the ~~the~~ FIRST TWO questions will be marked.

| Questions | Marks |
|-----------|-------|
| L1        |       |
| L2        |       |
| L3        |       |
| G1        |       |
| G2        |       |
| G3        |       |

This exam comprises the cover page and 16 pages of questions.

Section L (Linear Regression Models)  
Answer only two questions out of L1-L3

L1.

(a) (The use of indicator variables in linear regression analysis)

Consider a simple linear regression model that is used to study the relationship between a response variable  $y$  and a covariate  $x$ . The analysis is based on  $n$  independent observations

$$(x_i, y_i); i = 1, 2, \dots, n$$

Suppose these  $n$  observations are divided into  $M$  groups each having  $n_m$  observations, such that  $\sum_{m=1}^M n_m = n$ . Assume that the groups are independent from each other. The most general simple linear regression model that can be used to analyze such data is

$$y = \beta_0 + \beta_1 x + \epsilon; m = 1, 2, \dots, M$$

where  $\epsilon \sim N(0, \sigma^2)$ . That is, one can fit  $M$  separate (different) simple linear regression models to the data, corresponding to  $M$  independent groups of observations.

Using indicator variables,

| Grade Point Averages |              |             |
|----------------------|--------------|-------------|
| Lower class          | Middle class | Upper class |
| 2.87                 | 3.23         | 2.25        |
| 2.16                 | 3.45         | 3.13        |
| 3.14                 | 2.78         | 2.44        |
| 2.51                 | 3.77         | 2.54        |

Denote

$\mu_0$ : grand mean of the grade point average for a college freshman.

$\mu_i$ : the effect of the  $i$ th socioeconomic class on the grade point average, for  $i = 1, 2, 3$ .

$\mu_{i0}$ : mean of the grade point average for a college freshman in the  $i$ th socioeconomic class, for  $i = 1, 2, 3$ .

In the above notation  $i = 1, 2, 3$ , are corresponding to the three socioeconomic classes mentioned in the table, respectively.

(i) Write down one-way analysis of variance

L2.

An inverter is an electrical device that converts direct current (DC) to alternating current (AC). The data in this question are on measurement of the transient points of an electronic inverter. A portion of the data is given in the following Table. There are 24 observations.

Table 1:

|    | $x_1$ | $x_2$ | $x_3$ | $x_4$ | Y    |
|----|-------|-------|-------|-------|------|
| 1  | 3.00  | 3.00  | 3.00  | 3.00  | 0.79 |
| 2  | 3.00  | 6.00  | 6.00  | 6.00  | 1.71 |
| ⋮  | ⋮     | ⋮     | ⋮     | ⋮     | ⋮    |
| 23 | 2.00  | 3.00  | 8.00  | 6.00  | 1.51 |
| 24 | 3.00  | 3.00  | 8.00  | 8.00  | 0.75 |

The variables of interest are:

Y: Transient point (volts) of PMOS-NMOS inverters.

$x_1$ : Width of the NMOS device.

$x_2$ : Length of the NMOS device.

$x_3$ : Width of the PMOS device.

$x_4$ : Length of the PMOS device

Question L2 is continued on the next page.



L3.

Consider the multiple linear regression model

$$Y = X\beta + \epsilon$$

where  $\beta = (\beta_0; \beta_1; \beta_2; \dots; \beta_k)$  is the unknown vector of regression parameters,  $Y$  is the  $n \times 1$  dimensional vector of observations of the response variable,  $X$  is the  $n \times (k+1)$  dimensional design matrix,  $\epsilon$  is then  $n \times 1$  dimensional vector of errors, and  $\epsilon \sim N(0; \sigma^2 I_n)$  where  $I_n$  is the  $n \times n$  identity matrix.

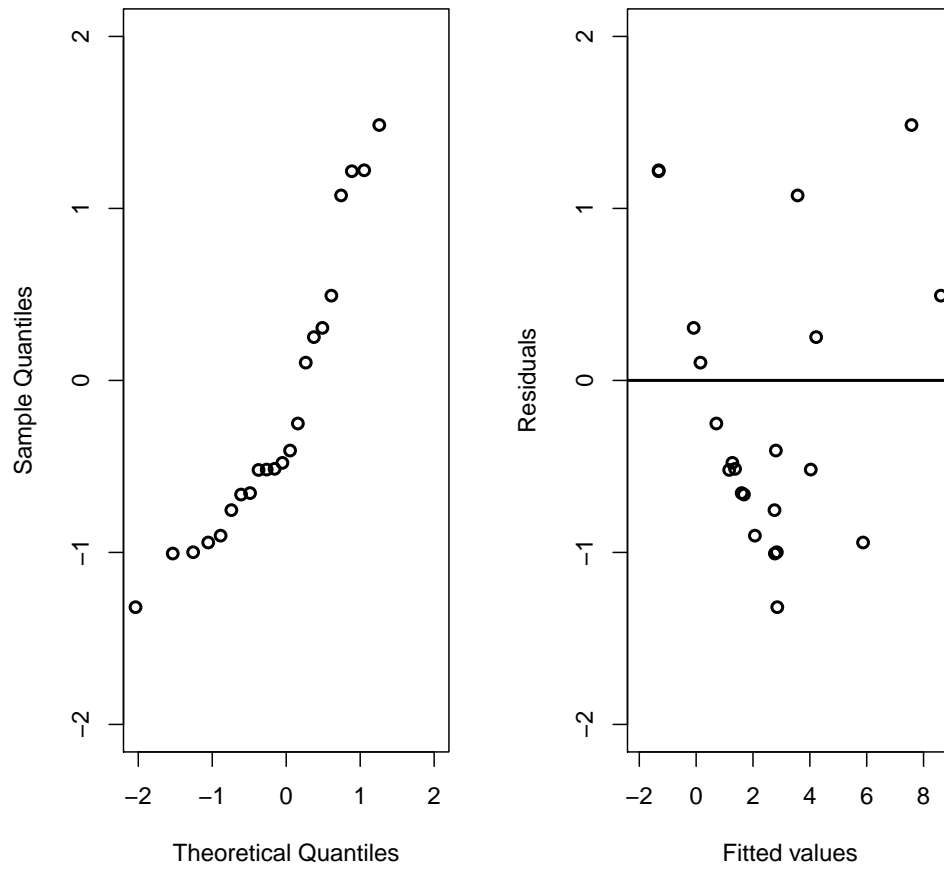
## R-Code for Question L2.

```

1 > D1<-read.table("Q2_data.txt", header=T)
2 > attach(D1)
3 > fit1<-lm(y ~., data=D1)
4 > summary(fit1)
5
6 Call:
7 lm(formula = y ~., data = D1)
8
9 Residuals:
10      Min       1Q   Median       3Q      Max
11 -1.4894 -0.9324 -0.6098  0.7224  3.3659
12
13 Coefficients:
14              Estimate Std. Error t value Pr(>|t|)
15 (Intercept)  1.52430     1.02317   1.490 0.152692
16 x1          -0.30606     0.07736  -3.956 0.000847 ***
17 x2           0.37439     0.05820   6.433 3.63e-06 ***
18 x3           0.44957     0.12354   3.639 0.001746 **
19 x4          -0.46557     0.13750  -3.386 0.003102 **
20 ---
21 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
22
23 Residual standard error: 1.445 on 19 degrees of freedom
24 Multiple R-squared:  0.8044,    Adjusted R-squared:  0.7632
25 F-statistic: 19.53 on 4 and 19 DF,  p-value: 1.604e-06
26
27 > anova(fit1)
28 Analysis of Variance Table
29
30 Response: y
31      Df Sum Sq Mean Sq F value    Pr(>F)
32 x1     1  11.989  11.989   5.7385 0.027056 *
33 x2     1 111.745 111.745  53.4885 6.185e-07 ***
34 x3     1  15.523  15.523   7.4304 0.013418 *
35 x4     1  23.951  23.951  11.4645 0.003102 **
36 Residuals 19  39.694   2.089
37 ---
38 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Figure 1: Residuals plots foModel 1.





```

42 > fit2<-lm(new_y ~., data=D1)
43 > summary(fit2)
44
45 Call:
46 lm(formula = new_y ~., data = D1)
47
48 Residuals:
49      Min       1Q   Median       3Q      Max
50 -0.36567 -0.13258 -0.04697  0.10125  0.40247
51
52 Coefficients:
53             Estimate Std. Error t value Pr(>|t|)
54 (Intercept) -0.07485    0.21386   -0.35    0.73
55 new_x1      -1.24866    0.07939  -15.73 2.38e-12 ***
56 new_x2       1.58623    0.07811   20.31 2.41e-14 ***
57 new_x3       0.93365    0.10515    8.88 3.44e-08 ***
58 new_x4      -1.34004    0.11749  -11.41 6.08e-10 ***
59 ---
60 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
61
62 Residual standard error: 0.2086 on 19 degrees of freedom
63 Multiple R-squared:  0.9772, Adjusted R-squared:  0.9724
64 F-statistic: 203.7 on 4 and 19 DF, p-value: 2.57e-15
65
66 > anova(fit2)
67 Analysis of Variance Table
68
69 Response: new_y
70      Df Sum Sq Mean Sq F value    Pr(>F)
71 new_x1  1  3.2553  3.2553  74.832 5.149e-08 ***
72 new_x2  1 26.0375 26.0375 598.539 7.940e-16 ***
73 new_x3  1  0.4960  0.4960  11.402 0.003167 **
74 new_x4  1  5.6591  5.6591 130.089 6.080e-10 ***
75 Residuals 19  0.8265  0.0435
76 ---
77 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Note: The response and covariates in Model 2 are called new in the output.



Section G (Generalized Linear Models)  
Answer only two questions out of G1 G3

G1. The table below summarizes some data collected on traffic accidents and seat belt usage collected in the state of Florida in 1988:

| Safety Equipment<br>In Use (S) | Whether<br>Ejected (E) | Injury (Y) |       |
|--------------------------------|------------------------|------------|-------|
|                                |                        | Nonfatal   | Fatal |
| Seat belt                      | Yes                    | 1105       | 14    |
|                                | No                     | 411,111    | 483   |
| None                           | Yes                    | 4624       | 497   |
|                                | No                     | 157,342    | 1008  |

Table 2: Source Florida Department of Highway Safety and Motor Vehicles.

- (a) Describe briefly why one cannot test for goodness of fit of the three-way interaction model,  $S:E:I$ , in 3 sentences or fewer (2 marks)
- (b) Consider the homogeneous association or no three-way interaction Poisson log-linear model,  $S:E \quad E:I \quad S:I$ . Describe, in words, what this model assumes about the conditional independence structure of the data. (2 marks)
- (c) Identify the model (model 1a, model 1b, or model 1c) in the R output on pages 14 and 15 which fits the model of homogeneous association to the Florida seat belt data. Assess the goodness of fit of the model you've identified using a significance level  $\alpha = 0.01$ . (8 marks)
- (d) Identify the model (model 1a, model 1b, or model 1c) in the R output on pages 14 and 15 which fits the following model  $S:E \quad E:I$ . Assess the goodness of fit of the model you've identified using a significance level  $\alpha = 0.01$ . (8 marks)

Question G1 is continued on the next page.

- (e) Consider now the two models that you selected in parts (c) and (d) on the previous page. Which is the more appropriate model? Explain your reasoning. (6 marks)
- (f) Another public safety researcher modeled the data in Table 2 using the following logistic regression model:

$$\log \frac{\text{Pr}(\text{Nonfatal} | \mathbf{S}; \boldsymbol{\beta})}{\text{Pr}(\text{Fatal} | \mathbf{S}; \boldsymbol{\beta})} = \beta_0 + \beta_1 \text{Yes}$$

Determine which of the three Poisson loglinear models on pages 14 and 15 yields equivalent inference to this logistic regression model and prove that the models are equivalent. (8 marks)

- (g) Compute the maximum likelihood parameter estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$  for the logistic regression model in part (f) using the equivalent Poisson loglinear model that you identified in part (f). (6 marks)

G2.

- (a) Suppose in an independent, two-sample problem that  $Y_i$  is Poisson with  $\mu_i = \exp(\eta_i)$  for  $i = 1, \dots, n$ , where  $x_i = 1$  for  $i = 1, \dots, n_A$  for the observations from group A and  $x_i = 0$  for  $i = n_A + 1, \dots, n$  for the observations from group B. Show that for any link function  $g$ , the score equations imply that the fitted means  $\hat{\mu}_A$  and  $\hat{\mu}_B$  equal the sample means. (10 marks)
- (b) In a GLM, suppose that  $\text{Var}(Y_i) = \mu_i^2$  for  $E(Y_i) = \mu_i$ . Show that the link function satisfying  $g(\mu) = \ln \mu$  has the same weight matrix  $W$  at each cycle of the Fisher scoring algorithm. Show that this link function for a Poisson GLM is  $g(\mu) = \ln \mu$ . (10 marks)
- (c) Assume that  $Y$  is an ordinal response variable and  $x$  is a continuous predictor. Show that for the cumulative logit model,

$$\text{logit}(\text{Pr}(Y \leq j | x)) = \alpha_j - \beta x$$

cumulative probabilities may be misordered for some  $x$  values. (10 marks)

- (d) Assume that  $X$  and  $Y$  are both ordinal categorical variables and

G3. An often used data set in statistics courses contains observations of female horseshoe crabs. For this problem  $Y_i$  is the binary response variable of interest is whether the horseshoe crab has at least one other male living outside of her nest (a satellite) or not. The only covariate of interest  $X_i$  for this problem is the weight (in kg) of the female crab. The goal of the analysis is to characterize the relationship between weight of the female horseshoe crabs and the presence of at least one male satellite crab.

- (a) A marine biologist collaborator tells you that their past experience indicates that the probability of the presence of at least one satellite crab varies linearly with the weight of the female crab. Your collaborator suggests that one should use simple linear regression with  $Y$  as the response and  $X$  as the sole covariate. Give at least two reasons, without referring to the output, why a simple linear regression model would not necessarily be appropriate here. 4 marks
- (b) You suggest that a GLM with a logistic link would be a better choice. Your collaborator is concerned that this model would not respect the assumption that the probability of a satellite was approximately linear in the weight. Again without referring to the R output, explain why using a generalized linear model with a logistic link would still be reasonable if the association were truly linear over the support of  $X$ . Hint: Consult Section 10.1 of the textbook. 4 marks

s k q r f o f . s 4 marks(



```
99 > summary(gmodel 1b)
100
101 Coefficients:
102             Estimate Std. Error z value Pr(>|z|)
103 (Intercept)    6.026472   0.025986 231.916  <2e-16 ***
104
```





# Table of the Chi-squared distribution

Entries in table are  $\chi^2_{p,q}$  the tail quantile of Chi-squared distribution given in columns,  $q$  given in rows.

|     | Left-tail |          |          |          | Right-tail |           |           |           |           |
|-----|-----------|----------|----------|----------|------------|-----------|-----------|-----------|-----------|
|     | 0.99500   | 0.99000  | 0.97500  | 0.95000  | 0.90000    | 0.05000   | 0.02500   | 0.01000   | 0.00500   |
| 1   | 0.00004   | 0.00016  | 0.00098  | 0.00393  | 0.02577    | 3.84146   | 5.02389   | 6.63490   | 7.87944   |
| 2   | 0.01003   | 0.02010  | 0.05064  | 0.10259  | 0.21475    | 5.99146   | 7.37776   | 9.21034   | 10.59663  |
| 3   | 0.07172   | 0.11483  | 0.21580  | 0.35185  | 0.58424    | 7.81473   | 9.34840   | 11.34487  | 12.83816  |
| 4   | 0.20699   | 0.29711  | 0.48442  | 0.71072  | 1.06387    | 9.48773   | 11.14329  | 13.27670  | 14.86026  |
| 5   | 0.41174   | 0.55430  | 0.83121  | 1.14548  | 1.61074    | 11.07050  | 12.83250  | 15.08627  | 16.74960  |
| 6   | 0.67573   | 0.87209  | 1.23734  | 1.63538  | 2.20413    | 12.59159  | 14.44938  | 16.81189  | 18.54758  |
| 7   | 0.98926   | 1.23904  | 1.68987  | 2.16735  | 2.83307    | 14.06714  | 16.01276  | 18.47531  | 20.27774  |
| 8   | 1.34441   | 1.64650  | 2.17973  | 2.73264  | 3.48954    | 15.50731  | 17.53455  | 20.09024  | 21.95495  |
| 9   | 1.73493   | 2.08790  | 2.70039  | 3.32511  | 4.16846    | 16.91898  | 19.02277  | 21.66599  | 23.58935  |
| 10  | 2.15586   | 2.55821  | 3.24697  | 3.94030  | 4.86515    | 18.30704  | 20.48318  | 23.20925  | 25.18818  |
| 11  | 2.60322   | 3.05348  | 3.81575  | 4.57481  | 5.57778    | 19.67514  | 21.92005  | 24.72497  | 26.75685  |
| 12  | 3.07382   | 3.57057  | 4.40379  | 5.22603  | 6.30249    | 21.02607  | 23.33666  | 26.21697  | 28.29952  |
| 13  | 3.56503   | 4.10692  | 5.00875  | 5.89186  | 7.04180    | 22.36203  | 24.73560  | 27.68825  | 29.81947  |
| 14  | 4.07467   | 4.66043  | 5.62873  | 6.57063  | 7.78954    | 23.68479  | 26.11895  | 29.14124  | 31.31935  |
| 15  | 4.60092   | 5.22935  | 6.26214  | 7.26094  | 8.54786    | 24.99579  | 27.48839  | 30.57791  | 32.80132  |
| 16  | 5.14221   | 5.81221  | 6.90766  | 7.96165  | 9.32641    | 26.29623  | 28.84535  | 31.99993  | 34.26719  |
| 17  | 5.69722   | 6.40776  | 7.56419  | 8.67176  | 10.12716   | 27.58711  | 30.19101  | 33.40866  | 35.71847  |
| 18  | 6.26480   | 7.01491  | 8.23075  | 9.39046  | 10.94984   | 28.86930  | 31.52638  | 34.80531  | 37.15645  |
| 19  | 6.84397   | 7.63273  | 8.90652  | 10.11701 | 11.79203   | 30.14353  | 32.85233  | 36.19087  | 38.58226  |
| 20  | 7.43384   | 8.26040  | 9.59078  | 10.85081 | 12.65119   | 31.41043  | 34.16961  | 37.56623  | 39.99685  |
| 21  | 8.03365   | 8.89720  | 10.28290 | 11.59131 | 13.52650   | 32.67057  | 35.47888  | 38.93217  | 41.40106  |
| 22  | 8.64272   | 9.54249  | 10.98232 | 12.33801 | 14.41813   | 33.92444  | 36.78071  | 40.28936  | 42.79565  |
| 23  | 9.26042   | 10.19572 | 11.68855 | 13.09051 | 15.32706   | 35.17246  | 38.07563  | 41.63840  | 44.18128  |
| 24  | 9.88623   | 10.85636 | 12.40115 | 13.84843 | 16.25356   | 36.41503  | 39.36408  | 42.97982  | 45.55851  |
| 25  | 10.51965  | 11.52398 | 13.11972 | 14.61141 | 17.19781   | 37.65248  | 40.64647  | 44.31410  | 46.92789  |
| 26  | 11.16024  | 12.19815 | 13.84390 | 15.37916 | 18.15983   | 38.88514  | 41.92317  | 45.64168  | 48.28988  |
| 27  | 11.80759  | 12.87850 | 14.57338 | 16.15140 | 19.13971   | 40.11327  | 43.19451  | 46.96294  | 49.64492  |
| 28  | 12.46134  | 13.56471 | 15.30786 | 16.92788 | 20.13672   | 41.33714  | 44.46079  | 48.27824  | 50.99338  |
| 29  | 13.12115  | 14.25645 | 16.04707 | 17.70837 | 21.15108   | 42.55697  | 45.72229  | 49.58788  | 52.33562  |
| 30  | 13.78672  | 14.95346 | 16.79077 | 18.49266 | 22.18152   | 43.77297  | 46.97924  | 50.89218  | 53.67196  |
| 40  | 20.70654  | 22.16426 | 24.43304 | 26.50930 | 29.65552   | 55.75848  | 59.34171  | 63.69074  | 66.76596  |
| 50  | 27.99075  | 29.70668 | 32.35736 | 34.76425 | 37.68912   | 67.50481  | 71.42020  | 76.15389  | 79.48998  |
| 60  | 35.53449  | 37.48485 | 40.48175 | 43.18796 | 45.65737   | 79.08194  | 83.29767  | 88.37942  | 91.95170  |
| 70  | 43.27518  | 45.44172 | 48.75756 | 51.73928 | 53.67204   | 90.53123  | 95.02318  | 100.42518 | 104.21490 |
| 80  | 51.17193  | 53.54008 | 57.15317 | 60.39148 | 61.65832   | 101.87947 | 106.62857 | 112.32879 | 116.32106 |
| 90  | 59.19630  | 61.75408 | 65.64662 | 69.12603 | 70.07890   | 113.14527 | 118.13589 | 124.11632 | 128.29894 |
| 100 | 67.32756  | 70.06489 | 74.22193 | 77.92947 | 78.78178   | 124.34211 | 129.56120 | 135.80672 | 140.16949 |



