Stats 413 Fall 2019

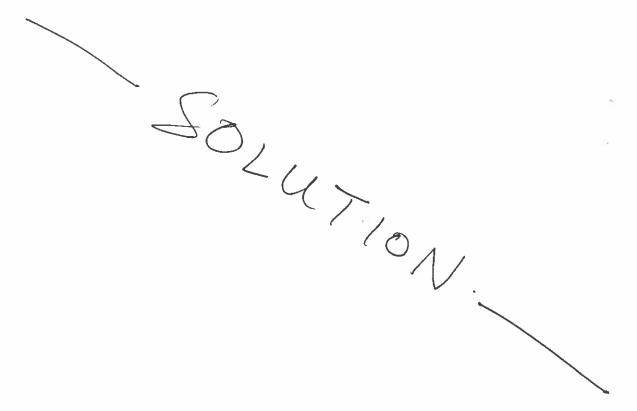
Exam I

Total Points: 50

2:30 pm - 3:50pm. Wednesday, Oct 9

Name: _	Pele			
ID:		-		

Instruction: Answer each question in the space provided. You may continue the answer on the back of the sheet, but please do not continue the answer for any question on the sheet for a different question. If you need more space, please use extra blank sheets provided and label them clearly.



1. [22 pts] Consider a multiple linear regression model $Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \beta_3 X_{i,3} + e_i, i = 1, \dots, 100$. We get the following R-output (values are rounded to two digits).

Call:

 $lm(formula = y \sim x1 + x2 + x3)$

Residuals:

Min 1Q Median 3Q Max -4.50 -1.35 0.05 1.18 4.00

Coefficients:

(Intercept) x1 x2	0.68 -0.38 1,40	Std.	0.20 0.20	t	3.5 -1.9	8e-04 0.06	***
	1.40		0.19		7.2	1e-10	***
x3	0.98		0.21		4.6	1e-05	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.9 on 96 degrees of freedom Multiple R-squared: 0.44, Adjusted R-squared: 0.42 F-statistic: 25 on 3 and 96 DF, p-value: 3.9e-12

- (a) [4 pts] For the considered linear regression model, write out the corresponding mean function. In addition, specify what assumptions are made for the statistical errors in this model.
- · E[41x] = Bo+ Bix1 + B2x2 + B3X3.
- . assumptions on errors:

E[ei]
$$X = 0$$
 eis are independent.
 $Var(ei)X) = o^{2}$ ($Cov(e) = oTn$)

- (b) [5 pts] What is the ordinary least square estimate of β₃, i.e., β̂₃? How to interpret this value? What is the value of the standard error of β̂₃? How to interpret this value? Give a 95% confidence interval of β₃ (suppose the 97.5% quantile of the used t-distribution is 2).
- · B3 = 0.98. expected change in 4 for a 1 unit increase in X3; holding X1 & Xe constant
- Se($\vec{\beta}_3$) = 0.21, estimate of standard deviation of the estimata $\vec{\beta}_3$.
- . 95% CI for B3: [0.98-2.0.21, 0.98+2-021]=[0.56,1.4]
 - (c) [1 pt] Calculate the predicted value of the response of a new observation with $X_1^* = 1, X_2^* = 1, X_3^* = 0$.

$$\hat{Y}^{+} = \hat{\beta}_{0} + \hat{\beta}_{1} \times \hat{1} + \hat{\beta}_{2} \times \hat{2} + \hat{\beta}_{3} \times \hat{3}$$

$$= 0.48 + (-0.38) \cdot 1 + (-4.1) + 0$$

$$= 1.7$$

(d) [3 pts] What are the values of R^2 and adjusted R^2 for the above regression? How shall we interpret the value of R^2 ?

R2 is the proportion of variance explained by the linear relationship w/ x1, x2, x3.

- (e) [3 pts] For the F-statistic in the last line of the R-output, write out the corresponding mean functions under the null and alternative hypotheses. What distribution does the corresponding test statistic follow under the null hypothesis? Give your conclusion for this test.
- · Ho: ETYIXT = Bo, Ha: ETYIXT = Bo+BIXI+ B2X2+ B2X3.
- . Under the, the F-statistic follows on F3,96 distribution
- · P= 3.9×10 × 0.05, so we reject to.
 - (f) [6 pts] We have the following ANOVA (type II) table for the regression $Y \sim X_1 + X_2 + X_3$.

Anova Table (Type II tests)

Response: y
Sum Sq Df F value Pr(>F)
x1 13.37 1 3.5426 0.06284
x2 198.27 1 52.5220 1.073e-10
x3 80.89 1 21.4266 1.152e-05
Residuals 362.40 96

For each F-test in this ANOVA table, specify the mean functions under the null and alternative hypotheses. What distribution does the corresponding test statistic follow under each null hypothesis? Based on the p-values, what are your conclusions about these tests?

row XI: Ho: ETYIM = BO+ B2X2+B2X3

Ma: ETYIM = BO+ B2X2+ B3 X3+ B1X1.

FNF,96. P>0.05, do not reject flo.

, row X2: Ho: ETYIM = Bo+ Bix, + B3X3

Ha: ETYIM = BO+ BIXIT BOX3+ B2X2

FNF1,96. PROIOS, reject to.

· rowks: Ho! ETYIXT = Bo+ BIKI+ P2K2

Ma: ETYLX] = Bo+ Bixi+ Bzxz+ Bix3,

FNF1,96

Proces, reject the

- 2. [18 pts] Consider a multiple linear regression with n observations and p predictors. $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$, where $\mathbf{Y} = (Y_1, \dots, Y_n)^{\top}$ is the $n \times 1$ vector of responses: \mathbf{X} is the $n \times (p+1)$ matrix of predictors, including a column of 1's for the intercept: and $\mathbf{e} = (e_1, \dots, e_n)^{\top}$ is the $n \times 1$ vector of statistical errors. Let $\bar{Y} = n^{-1} \sum_{i=1}^{n} Y_i$.
 - (a) [5 pts] Derive the formula of the ordinary least square (OLS) estimator $\hat{\beta}$ (please show details).
 - (b) [4 pts] Show that $\hat{\beta}$ is an unbiased estimator of β .
 - (c) [4 pts] Let $\widehat{\mathbf{Y}} = (\widehat{Y}_1, \dots, \widehat{Y}_n)^\mathsf{T}$ be the fitted values from the OLS estimation. Prove that the sample mean of $\{\widehat{Y}_1, \dots, \widehat{Y}_n\}$ is the same to the sample mean of $\{Y_1, \dots, Y_n\}$, i.e., $n^{-1} \sum_{i=1}^n \widehat{Y}_i = \widehat{Y}$.
 - (d) [5 pts] Prove that TSS = RSS + SS_{reg}, where TSS is the total sum of squares (i.e. $\sum_{i=1}^{n} (Y_i \bar{Y})^2$). RSS is the residual sum of squares (i.e. $\sum_{i=1}^{n} (Y_i \bar{Y})^2$), and SS_{reg} is the regression sum of squares (i.e. $\sum_{i=1}^{n} (\widehat{Y}_i \bar{Y})^2$).

(a)
$$\beta$$
 minimizes $RSS = \sum_{i=1}^{N} (y_i - x_i^T \beta) = (y_i - x_i \beta)^T (y_i - x_i \beta) = \cdots$

$$(\beta) = Tr((4-x\beta)^{T}(4-x\beta))$$

$$= Tr((4^{T}x)^{T} - 2Tr((4^{T}x\beta) + Tr(\beta^{T}x^{T}x\beta))$$

$$= Tr((4^{T}x\beta) + Tr(\beta^{T}x^{T}x\beta)$$

$$\nabla_{\beta}J(\beta) = -2(Y^{T}x)^{T} + 2X^{T}x\beta$$

$$0 = -2x^{T}y + 2x^{T}x^{S} \implies \hat{\beta} = cx^{T}x^{T}x^{T}y \quad (if cx^{T}x^{T}exisk)$$

(b)
$$\hat{\beta} = (x^T x)^T x^T y \qquad y = x\beta + e$$
.

$$= (x^T x)^T x^T (x\beta + e)$$

$$= \beta + (x^T x)^T x^T e$$
.

$$E[\hat{\beta}] = E[\beta + (x^{T}x)^{T}x^{T}e] = \beta + (x^{T}x)^{T}x^{T}E[e] = 0.$$

$$= \beta.$$

thu, Bis unbiasul for B.

equivalently show
$$\frac{S}{i\pi} \dot{y}_i = \frac{S}{i\pi} \dot{y}_i$$

LHS = $\frac{S}{i\pi} \dot{y}_i = \frac{S}{i\pi} (\dot{y}_i + \dot{e}_i) = \frac{S}{i\pi} \dot{y}_i + \frac{S}{i\pi} \dot{e}_i$

= $\frac{S}{i\pi} \dot{y}_i$ from womal equations.

normal equation: X7(y-xp3)=0, so from first column of X.

(d).
$$TSS = \sum_{i=1}^{n} (y_i - \overline{y})^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i + \hat{y}_i - \overline{y})^2 = \dots$$

$$= \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \sum_{i=1}^{n} (\hat{y}_i - \hat{y}_i)^2 + 2\sum_{i=1}^{n} (y_i - \hat{y}_i)(\hat{y}_i - \hat{y}_i)$$

RSS SSreg. Since
$$yi-\hat{y}i=\hat{e}i$$
 $(\uparrow x)=\sum_{i=1}^{\infty}\hat{e}_{i}\cdot x_{i}^{T}\hat{\beta}$
 $\hat{y}_{i}=x_{i}^{T}\hat{\beta}$ $=(\sum_{i=1}^{\infty}\hat{e}_{i}\cdot x_{i}^{T})\hat{\beta}$

inns TSS = RSS + SSreg.

=o by normal equ.

3. [10 pts] Consider a simple linear regression model. Give the definition form of \mathbb{R}^2 . Prove that \mathbb{R}^2 is equal to the square of the correlation between the predictor and the response.

$$R^{2} = 1 - \frac{\hat{S}_{SS}}{T_{SS}} = 1 - \frac{\hat{S}_{S}(y_{1} - \hat{y}_{1})^{2}}{\hat{S}_{S}(y_{1} - \hat{y}_{1})^{2}} = \frac{\hat{S}_{S}(\hat{y}_{1} - \hat{y}_{2})^{2}}{\hat{S}_{S}(y_{1} - \hat{y}_{2})^{2}}.$$

$$\hat{S}_{SS}(y_{1} - \hat{y}_{2})^{2}.$$

$$\rho^{2} = \left(\frac{\hat{\Sigma}(x_{i}-\bar{x})(y_{i}-\bar{y}_{i})^{2}}{\hat{\Sigma}(x_{i}-\bar{x}_{i})^{2}\hat{\Sigma}(y_{i}-\bar{y}_{i})^{2}}\right)^{2}$$

$$R^{2} = \frac{\hat{\chi}(\hat{\beta}_{0} + \hat{\beta}_{1}\hat{\chi}_{1} - \hat{\gamma})^{2}}{\hat{\chi}(\hat{\gamma}_{1} - \hat{\gamma}_{1})^{2}}$$

$$\hat{\chi}(\hat{\gamma}_{1} - \hat{\gamma}_{1})^{2}$$

$$\hat{\chi}(\hat{\gamma}_{1} - \hat{\gamma}_{1})^{2}$$

$$=\frac{\sum_{i=1}^{n}(\beta_{i}(x_{i}-\overline{x}))^{2}}{\sum_{i=1}^{n}(\gamma_{i}-\overline{y})^{2}}$$

$$=\frac{\sum_{i=1}^{n}(\beta_{i}(x_{i}-\overline{x}))^{2}}{\sum_{i=1}^{n}(\gamma_{i}-\overline{y})^{2}}$$

$$=\frac{\sum_{i=1}^{n}(\beta_{i}(x_{i}-\overline{x}))^{2}}{\sum_{i=1}^{n}(\gamma_{i}-\overline{y})^{2}}$$

$$\beta_{i} = \frac{\sum_{i=1}^{n} (x_{i} - \overline{x}) (y_{i} - \overline{y})}{\sum_{i=1}^{n} (x_{i} - \overline{x})^{2}}$$

$$= \left(\frac{\hat{S}(x_i - \bar{X})(y_i - \bar{Y})}{\hat{S}(x_i - \bar{X})^2}\right)^2 \frac{\hat{S}(x_i - \bar{X})^2}{\hat{S}(x_i - \bar{X})^2} = \frac{\hat{S}(x_i - \bar{X})(y_i - \bar{Y})}{\hat{S}(x_i - \bar{X})}$$