

MAT 3375, Regression Analysis

Assignment 1 (Total = 50)

Professor: Termeh Kousha

Deadline: Friday, Sep 23, 2016 SHARP 3pm (Math department (585 KED) drop boxes)

Late assignments will NOT be accepted; nor will unstapled assignments. Professors in the math department will not lend you a stapler.

You should complete ALL the questions in the assignments. It is possible, however, that not all the questions will be marked. In that case, the same questions will be marked in all assignments. You will not be informed beforehand which questions will be marked.

Student Name _____

Student Number _____

By signing below, you declare that this work was your own and that you have not copied from any other individual or other source.

Signature _____

8 points

1. Question 2.25

part (b):

4 points

$$\text{Cov}(\bar{Y}, \hat{\beta}_1) = \text{Cov}\left(\frac{\sum Y_i}{n}, \frac{S_{XY}}{S_{XX}}\right) = \frac{1}{n S_{XX}} \text{Cov}\left(\sum Y_i, \sum (x_i - \bar{x}) Y_i\right)$$

$$= \frac{1}{n S_{XX}} \sum \text{Cov}(Y_i, (x_i - \bar{x}) Y_i) =$$

$$= \frac{1}{n S_{XX}} \sum (x_i - \bar{x}) \underbrace{\text{Cov}(Y_i, Y_i)}_{\sigma^2} = \frac{\sigma^2}{n S_{XX}} \sum (x_i - \bar{x}) = 0$$

4 points

$$\text{Part (a)} \quad \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = \text{Cov}(\bar{Y} - \hat{\beta}_1 \bar{x}, \hat{\beta}_1) =$$

$$\underbrace{\text{Cov}(\bar{Y}, \hat{\beta}_1)}_{=0 \text{ by part (b)}} - \text{Cov}(\hat{\beta}_1 \bar{x}, \hat{\beta}_1) = -\bar{x} \text{Cov}(\hat{\beta}_1, \hat{\beta}_1)$$

$$= -\bar{x} \text{Var}(\hat{\beta}_1)$$

$$= -\frac{\bar{x} \sigma^2}{S_{XX}}$$

2. Consider the simple linear regression model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, \quad i = 1, \dots, n,$$

where $\epsilon_1, \dots, \epsilon_n$ are independent, $E[\epsilon_i] = 0$ and $Var[\epsilon_i] = \sigma^2$. Show that

$$Var[\hat{\beta}_0] = \sigma^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{S_{XX}} \right].$$

Hint: Use question 1. $(Cov(\hat{Y}_i, \hat{\beta}_0) = 0)$.

$$\begin{aligned} Var(\hat{\beta}_0) &= Var(\bar{Y} - \hat{\beta}_1 \bar{X}) = Var(\bar{Y}) + \bar{X}^2 Var(\hat{\beta}_1) \\ &\quad + 2 Cov(\bar{Y}, \hat{\beta}_1) \\ &= \frac{\sigma^2}{n} + \overset{0}{Var(\hat{\beta}_1)} = \frac{\sigma^2}{n} + \frac{\bar{X}^2}{S_{XX}} \sigma^2 \\ &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{S_{XX}} \right). \end{aligned}$$

3. Consider the following sum of cross-products:

$$S_{XY} = \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}).$$

Show that S_{XY} can be written as $S_{XY} = \sum_{i=1}^n (X_i - \bar{X})Y_i$ and $S_{XY} = (\sum_{i=1}^n X_i Y_i) - n\bar{X}\bar{Y}$.
Consequently, by plugging-in $Y_i = X_i$ we get

$$S_{XX} = \sum_{i=1}^n (X_i - \bar{X})^2 = \left(\sum_{i=1}^n X_i^2 \right) - n\bar{X}\bar{X}.$$

$$S_{XX} = \sum_{i=1}^n (X_i - \bar{X})X_i.$$

$$\begin{aligned} S_{XY} &= \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum (x_i - \bar{x})y_i - \underbrace{\sum (x_i - \bar{x})\bar{y}}_{\bar{y} \sum (x_i - \bar{x})} \\ &= \sum (x_i - \bar{x})y_i = \sum x_i y_i - \bar{x} \sum y_i \\ &= \sum x_i y_i - n\bar{x}\bar{y} \end{aligned}$$

20 points

4. Simple linear regression method was used to analyze the data from a study investigating the relationship between roadway surface temperature in Fahrenheit (X) and pavement deflection (Y). Summary statistics were $n = 20$,

$$\sum_{i=1}^n Y_i = 12.575, \quad \sum_{i=1}^n Y_i^2 = 8.86, \quad \sum_{i=1}^n X_i = 1478, \quad \sum_{i=1}^n X_i^2 = 143,215.8, \quad \sum_{i=1}^n X_i Y_i = 1083.67.$$

Note: Alternatively, SSCP table can be given.

- Verify that $S_{XY} = 154.3775$, $S_{XX} = 33991.6$, $S_{YY} = 0.9534688$. (Hint: Use Question 3)
- Obtain the least squares estimates of the slope and intercept.
- Give an estimate of the mean pavement deflection when the surface temperature is 85F.
- Give a point estimate of the pavement deflection when the surface temperature is 90F.
- Compute the coefficient of correlation.
- Compute the coefficient of determination and interpret it.
- Estimate σ^2 . Hint: The alternative formula for SSE is $SSE = S_{YY} - \hat{\beta}_1 S_{XY}$.
- Use t -test to check for significance of regression using $\alpha = 0.05$. Find the p -value for this test. What conclusion can you draw?
- Use F -test to check for significance of regression using $\alpha = 0.05$
- Find a 99% confidence interval for slope. Hint: $t_{0.01/2,18} = 2.878$.
- Find a 99% confidence interval for intercept.
- Find a 99% confidence interval for the mean deflection when temperature is 85F.
- Find a 99% prediction interval on pavement deflection when the temperature is 90F.

(a) verify.

← b) $\hat{\beta}_1 = \frac{S_{XY}}{S_{XX}}$ $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$

so $\hat{\beta}_1 = 0.004542$ $\hat{\beta}_0 = 0.293123$

$\hat{y}_i = 0.293123 + 0.004542(x_i)$

c) $x_i = 85$ $\hat{y}_i = 0.68$ $x_i = 90 \rightarrow \hat{y}_i = 0.7$

e) $r = \frac{S_{XY}}{\sqrt{S_{XX} S_{YY}}} = 0.857$

f) $R^2 = r^2 = 0.735^4 \rightarrow$ That means 73.5% \rightarrow

Continue

← of the variation in pavement deflection is explained by the linear relationship with the surface temperature.

$$\textcircled{2} \textcircled{g) } \hat{\sigma}^2 = \text{MSE} = \frac{\text{SSE}}{n-2} = \frac{\text{SYY} - \hat{\beta}_1 \text{SXY}}{n-2} = 0.01402$$

$$\textcircled{2} \textcircled{h) } H_0: \beta_1 = 0, H_1: \beta_1 \neq 0$$

$$t = \frac{\hat{\beta}_1 - 0}{\sqrt{\frac{\hat{\sigma}^2}{\text{Sxx}}}} = 7.072.$$

$$t_{20-2} = t_{18} \quad 2(P(t_{18} > 7.072)) < 0.001$$

Reject the Null hypo. → there is a linear relationship between x and y .

$$\textcircled{3} \textcircled{i) } F^* = \frac{\text{SSR}}{\text{MSE}} = \frac{\text{SST} - \text{SSE}}{0.01402} = \frac{0.9534688 - 18(0.01402)}{0.01402}$$

$$= \frac{0.7011}{0.01402} = 50.007$$

$F^* > \underbrace{F_{0.05, 1, 18}}_{4.41} \rightarrow$ we can reject.

$$\textcircled{2} \textcircled{j) } \text{SD}[\hat{\beta}_1] = \sqrt{\frac{\hat{\sigma}^2}{\text{Sxx}}} \quad \hat{\beta}_1 \in (0.0027, 0.0064)$$

$$\textcircled{2} \textcircled{k) } \text{SD}[\hat{\beta}_0] = \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\text{Sxx}} \right]}$$
$$\hat{\beta}_0 \in (0.137, 0.45)$$

2) $x_n = 85$ $E(Y_n) = \beta_0 + \beta_1 x_n$

$$\hat{Y}_n \pm t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{(x_n - \bar{x})^2}{S_{xx}} \right]}$$

$$\hat{Y}_n = \hat{\beta}_0 + \hat{\beta}_1 x_n = 0.68$$

$$\hat{Y}_n \in (0.601, 0.761).$$

(m) $x_{n(\text{new})} = 90.$

2) $\hat{Y}_n(\text{new}) = \hat{\beta}_0 + \hat{\beta}_1(90) = 0.7$

$$\hat{Y}_n(\text{new}) \pm t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[1 + \frac{1}{n} + \frac{(x_n - \bar{x})^2}{S_{xx}} \right]}$$

$$Y_n(\text{new}) \in (0.3513, 1.0524).$$

12 points.

5. Consider the simple regression model with $\beta_0 = 0$ (regression passing through the origin). Thus, the regression model is

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

where $\epsilon_1, \dots, \epsilon_n$ are independent, $E[\epsilon_i] = 0$ and $Var[\epsilon_i] = \sigma^2$.

- 4 P₅ (a) Find the least-squares estimator of β_1 . Is it unbiased?

- 3 P₃ (b) Compute $Var[\hat{\beta}_1]$.

- 5 P₃ (c) Let $\hat{Y}_i = \hat{\beta}_1 X_i$. Show that

$$\hat{\sigma}^2 := \frac{1}{n-1} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

is an unbiased estimator of σ^2 .

(a). We want to minimize

$$Q(\beta_1) = \sum_{i=1}^n (Y_i - \beta_1 X_i)^2$$

$$\frac{dQ}{d\beta_1} = -2 \sum (Y_i - \beta_1 X_i) X_i = 0$$

$$\sum Y_i X_i - \beta_1 \sum X_i^2 = 0 \quad \hat{\beta}_1 = \frac{\sum Y_i X_i}{\sum X_i^2} *$$

$$E(\hat{\beta}_1) = E\left(\frac{\sum Y_i X_i}{\sum X_i^2}\right) = \frac{\sum X_i E(Y_i)}{\sum X_i^2} = \frac{\sum X_i \beta_1 X_i}{\sum X_i^2} = \beta_1 \frac{\sum X_i^2}{\sum X_i^2} = \beta_1$$

2 points.

$\hat{\beta}_1$ is an unbiased estimator for β_1 .

$$(b) \text{ var}[\hat{\beta}_1] = \text{var}\left[\frac{\sum Y_i X_i}{\sum X_i^2}\right] = \frac{1}{(\sum X_i^2)^2} \sum_{i=1}^n X_i^2 \text{var}(Y_i)$$

$$= \frac{1}{(\sum X_i^2)^2} \sigma^2 \sum_{i=1}^n X_i^2 = \frac{\sigma^2}{\sum X_i^2}$$

$$(c) \text{ we have } \sum (Y_i - \hat{Y}_i)^2 = \sum (Y_i - \hat{\beta}_1 X_i)^2 = \sum Y_i^2 + \hat{\beta}_1^2 \sum X_i^2 - 2 \hat{\beta}_1 \sum X_i Y_i$$
$$= \sum Y_i^2 - \hat{\beta}_1^2 \sum X_i^2$$

→

Note that

$$E(Y_i^2) = \text{var}(Y_i) + E(Y_i)^2 = \sigma^2 + \beta_1^2 x_i^2$$

and

$$E(\hat{\beta}_1^2) = \text{var}(\hat{\beta}_1) + E(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum x_i^2} + \beta_1^2$$

Thus

$$\begin{aligned} E\left(\sum (Y_i - \hat{Y}_i)^2\right) &= E\left(\sum Y_i^2 - \hat{\beta}_1^2 \sum x_i^2\right) \\ &= \sum E(Y_i^2) - E(\hat{\beta}_1^2) \sum x_i^2 \\ &= \sum (\sigma^2 + \beta_1^2 x_i^2) - \left(\frac{\sigma^2}{\sum x_i^2} + \beta_1^2\right) \sum x_i^2 \\ &= n\sigma^2 + \sum \beta_1^2 x_i^2 - \sigma^2 - \beta_1^2 \sum x_i^2 = \\ &= n\sigma^2 - \sigma^2 = (n-1)\sigma^2. \end{aligned}$$

Therefore :

$$E(\hat{\sigma}^2) = E\left(\frac{\sum (Y_i - \hat{Y}_i)^2}{n-1}\right) = \sigma^2$$

which shows that $\hat{\sigma}^2$ is an unbiased estimator for σ^2 .

5 points.

6. Consider the simple linear regression model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, \quad i = 1, \dots, n,$$

where $\epsilon_1, \dots, \epsilon_n$ are independent, $E[\epsilon_i] = 0$ and $Var[\epsilon_i] = \sigma^2$. Consider the test

$$H_0 : \beta_1 = 0, \quad H_1 : \beta_1 \neq 0.$$

Show that t -test and F -test are equivalent.

Hint: Use the following result: if $Z \sim t(\nu)$, then $Z^2 \sim F(1, \nu)$.

The t -test :
$$t^* = \frac{\hat{\beta}_1}{\frac{s}{\sqrt{S_{XX}}}} \sim t_{(n-2)}.$$

$T \sim t_{n-2}$ P-value is $2P(T > |t^*|)$

for the F -test
$$F^* = \frac{MSR}{MSE} = \frac{\hat{\beta}_1^2 S_{XX}}{\hat{\sigma}^2} =$$

$S^2 = MSE$
$$\frac{\hat{\beta}_1^2}{\frac{s^2}{S_{XX}}} = (t^*)^2 \sim F_{1, n-2}$$

The p -value is

$$2P(F_{1, n-2} > F^*) =$$

$$2P((T)^2 > F^* = t^{*2}) =$$

$T \sim t_{n-2}$ $2P(|T| > |t^*|) = 2P(T > t^*)$
 two test are equivalent.
 Always positive

5 points.

7. Let r_{XY} be the sample correlation coefficient between X and Y :

$$r_{XY} = \frac{S_{XY}}{\sqrt{S_{XX}S_{YY}}}$$

Let $Z = aX + b$, $a > 0$. Show that $r_{ZY} = r_{XY}$.

We first show $S_{ZY} = aS_{XY}$

$$\begin{aligned} S_{ZY} &= \sum z_i y_i - n \bar{z} \bar{y} = \sum (ax_i + b) y_i - n(a\bar{x} + b) \bar{y} \\ &= a \sum x_i y_i + b \sum y_i - na \bar{x} \bar{y} - nb \bar{y} = \\ &= a \sum x_i y_i - na \bar{x} \bar{y} = a (\sum x_i y_i - n \bar{x} \bar{y}). \end{aligned}$$

Also you can show similarly $S_{ZZ} = a^2 S_{XX}$.

Therefore:

$$r_{ZY} = \frac{S_{ZY}}{\sqrt{S_{ZZ} S_{YY}}} = \frac{a S_{XY}}{\sqrt{a^2 S_{XX} S_{YY}}} =$$

$$\frac{S_{XY}}{\sqrt{S_{XX} S_{YY}}}$$