

# STAT 3008 Applied Regression Analysis

## Tutorial 3.

XU yongze

Feb 23 & 24 & 26, 2013

### Contents

|   |          |
|---|----------|
| <b>1 Review Questions</b>                                   | <b>2</b> |
| <b>2 ANOVA Test(Analysis of Variance)</b>                   | <b>4</b> |
| <b>3 Hypothesis Testing of Parameter Estimates</b>          | <b>5</b> |
| <b>4 Confidence Interval of Parameter Estimates</b>         | <b>6</b> |
| <b>5 Exercises</b>  | <b>8</b> |
| 5.1 Exercise 3 Q1 . . . . .                                 | 8        |
| 5.2 Exercise 3 Q2 . . . . .                                 | 8        |
| 5.3 Example from 12-13 Midterm(similar to Ex3 Q3) . . . . . | 8        |
| 5.4 Exercise 3 Q5 2.12.4 . . . . .                          | 9        |
| 5.5 Example from 11-12 Midterm . . . . .                    | 9        |

# 1 Review Questions

For model  $y_i = \beta_0 + \beta_1 x_i + e_i$ ,  $e_i \sim i.i.d.N(0, \sigma^2)$ , derive the followings:

- a.  $\sum y_i = \sum \hat{y}_i$
- b.  $\sum x_i \hat{e}_i = 0$
- c.  $\sum \hat{y}_i \hat{e}_i = 0$
- d.  $\rho(\hat{\beta}_0, \hat{\beta}_1) = \frac{\bar{x}}{\sqrt{\frac{S_{XX}}{n} + \bar{x}^2}}$ .

Hint: In matrix notation,

$$\mathbf{y} = X\beta + \mathbf{e},$$

$H = X(X'X)^{-1}X'$  is symmetric and idempotent,

$$\hat{\mathbf{y}} = H\mathbf{y},$$

For a)

$$\sum y_i = \sum \hat{y}_i \Leftrightarrow \mathbf{1}'\mathbf{y} = \mathbf{1}'H\mathbf{y}, \mathbf{1}' = (1 \cdots 1),$$

$$HX = X, X = (\mathbf{1} \ *), H\mathbf{1} = \mathbf{1}.$$

For b)

$$X'\hat{\mathbf{e}} = \mathbf{0} \Rightarrow \sum x_i \hat{e}_i = 0,$$

$$X'\hat{\mathbf{e}} = X'(I - H)\mathbf{y} = (X' - X')\mathbf{y} = \mathbf{0}.$$

For c)

$$\hat{\mathbf{y}}'\hat{\mathbf{e}} = (H\mathbf{y})'(I - H)\mathbf{y} = \mathbf{y}'(H - H)\mathbf{y} = 0$$

Not in matrix notation,

a)

$$\begin{aligned} & \sum \hat{y}_i \\ &= \sum (\hat{\beta}_0 + \hat{\beta}_1 x_i) \\ &= \sum ((\bar{y} - \hat{\beta}_1 \bar{x}) + \hat{\beta}_1 x_i) \\ &= n\bar{y} - n\hat{\beta}_1 \bar{x} + \sum \hat{\beta}_1 x_i = n\bar{y} \\ &= \sum y_i \end{aligned}$$

b)

$$\begin{aligned}
& \sum x_i \hat{e}_i \\
&= \sum x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \\
&= \sum x_i y_i - \hat{\beta}_0 \sum x_i - \hat{\beta}_1 \sum x_i^2 \\
&= \sum x_i y_i - (\bar{y} - \hat{\beta}_1 \bar{x})(n \bar{x}) - \hat{\beta}_1 \sum x_i^2 \\
&= \sum x_i y_i - n \bar{x} \bar{y} + \hat{\beta}_1 (n \bar{x}^2 - \sum x_i^2) \\
&= \sum x_i y_i - n \bar{x} \bar{y} - \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} (\sum x_i^2 - n \bar{x}^2) \\
&= 0
\end{aligned}$$

c)

$$\begin{aligned}
& \sum \hat{y}_i \hat{e}_i \\
&= \sum (\hat{\beta}_0 + \hat{\beta}_1 x_i) \hat{e}_i \\
&= \hat{\beta}_0 \sum \hat{e}_i + \hat{\beta}_1 \sum x_i \hat{e}_i \\
&= 0
\end{aligned}$$

d)

$$\rho(\hat{\beta}_0, \hat{\beta}_1)$$

$$\begin{aligned}
&= \frac{Cov(\hat{\beta}_0, \hat{\beta}_1)}{\sqrt{Var(\hat{\beta}_0)Var(\hat{\beta}_1)}} \\
&= \frac{-\frac{\bar{x}}{SXX} \sigma^2}{\sqrt{(\frac{1}{n} + \frac{\bar{x}^2}{SXX}) \sigma^2 (\frac{1}{SXX}) \sigma^2}} \\
&= -\frac{\bar{x}}{SXX} * \frac{1}{\sqrt{(\frac{SXX + \bar{x}^2}{nSXX}) \frac{1}{SXX}}} \\
&= \frac{\bar{x}}{\sqrt{\frac{SXX + \bar{x}^2}{n}}}
\end{aligned}$$

## 2 ANOVA Test(Analysis of Variance)

### ♣ What is ANOVA test?

A test to compare different models, such as  $E(Y|X) = \beta_0$  and  $E(Y|X) = \beta_0 + \beta_1 x$ .

### ♣ Comparing Models.

$H_0 : E(Y|X) = \beta_0$  v.s.  $H_1 : E(Y|X) = \beta_0 + \beta_1 x$

Under  $H_0$ ,  $\beta_1 = 0$ , which means x and y are independent.

We have two methods to finish the comparison:

(1) Test  $\beta_1 = 0$

(2) ANOVA Test:

Under 0,  $\hat{\beta}_0 = \bar{y}$ ,  $RSS_0 = \sum(y_i - \bar{y})^2 = SYY$ .

Under 1,  $RSS_1 = \sum[y_i - (\hat{\beta}_0 + \hat{\beta}_1 x)]^2 = SYY - \frac{SXY^2}{SXX}$

Obviously,  $RSS_0 > RSS_1$ , but it is not necessarily model under  $H_1$  is better.

$SS_{reg}$  = Difference sum of squares due to regression.(as we add another parameter  $\beta_1$  to the model).

$$SS_{reg} = \Delta RSS = RSS_0 - RSS_1 = \frac{SXY^2}{SXX}$$

If  $SS_{reg}$  is large, which means the addition of the parameter  $\beta_1$  is useful to explain the variation of the data.

Q: How large is large?

Notice:

$$\frac{SS_{reg}}{\sigma^2} \sim \chi_1^2 \quad (2.1)$$

$$\frac{RSS_1}{\sigma^2} \sim \chi_{n-2}^2 \quad (2.2)$$

We can derive:

$$T.S. = \frac{SS_{reg}}{\frac{RSS_1}{n-2}} \sim F_{1,n-2}. \quad (2.3)$$

If obs.T.S.  $> F_{(\alpha,1,n-2)}$ ,  $H_0$  is rejected at  $\alpha$  level.

The following is the very important ANOVA Table.

| Source           | DF  | SS  | MS   | F                 |
|------------------|-----|---|--|-------------------|
| Regression       | 1   | $SS_{reg} = \Delta RSS = \frac{SXY^2}{SXX}$ | $MS_{reg} = \frac{SS_{reg}}{1} = SS_{reg}$ | $\frac{MSR}{MSE}$ |
| Residuals(Error) | n-2 | $SSE = RSS_1 = SYY - \frac{SXY^2}{SXX}$     | $MSE = \frac{SSE}{n-2}$                    |                   |
| Total            | n-1 | $SST = RSS_0 = SYY$                         |  |                   |

### 3 Hypothesis Testing of Parameter Estimates

1)  $\beta_0$  (Intercept):

$$H_0 : \beta_0 = \beta_0^* \text{ v.s. } H_1 : \beta_0 \neq \beta_0^*$$

$$T.S. = \frac{\hat{\beta}_0 - \beta_0^*}{\sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SXX} \right)}} \sim t_{n-2} \quad (3.1)$$

If  $|obs.T.S.| > t_{(\frac{\alpha}{2}, n-2)}$ ,  $H_0$  is rejected at  $\alpha$  level.

2)  $\beta_1$  (Slope):

$$H_0 : \beta_1 = \beta_1^* \text{ v.s. } H_1 : \beta_1 \neq \beta_1^*$$

$$T.S. = \frac{\hat{\beta}_1 - \beta_1^*}{\sqrt{\hat{\sigma}^2 \frac{1}{SXX}}} \sim t_{n-2} \quad (3.2)$$

If  $|obs.T.S.| > t_{(\frac{\alpha}{2}, n-2)}$ ,  $H_0$  is rejected at  $\alpha$  level.

3)  $\hat{y}_* | x_*$  (Prediction):

$$H_0 : y_* | x_* = y_* \text{ v.s. } H_1 : y_* | x_* \neq y_*$$

$$T.S. = \frac{\hat{y}_* | x_* - y_*}{\sqrt{\hat{\sigma}^2 \left( 1 + \frac{1}{n} + \frac{(x_* - \bar{x})^2}{SXX} \right)}} \sim t_{n-2} \quad (3.3)$$

If  $|obs.T.S.| > t_{(\frac{\alpha}{2}, n-2)}$ ,  $H_0$  is rejected at  $\alpha$  level.

4)  $\hat{y} | x_*$  (Fitted Value): (Just for your reference)

$$H_0 : E(y | x_*) = y_* \text{ v.s. } H_1 : E(y | x_*) \neq y_*$$

$$T.S. = \frac{(\hat{y} | x_* - y_*)^2}{2\hat{\sigma}^2 \left( \frac{1}{n} + \frac{(x_* - \bar{x})^2}{SXX} \right)} \sim F_{(2, n-2)} \quad (3.4)$$

If obs.T.S.  $> F_{(\alpha, 2, n-2)}$ ,  $H_0$  is rejected at  $\alpha$  level.

5)  $\sigma^2$ (Variance):

$$H_0 : \sigma^2 = \alpha \text{ v.s. } H_1 : \sigma^2 > \alpha$$

$$T.S. = \frac{(n-2)\hat{\sigma}^2}{a} \sim \chi_{n-2}^2 \quad (3.5)$$

If obs.T.S.  $> \chi_{(\alpha, n-2)}^2$ ,  $H_0$  is rejected at  $\alpha$  level.

## 4 Confidence Interval of Parameter Estimates

Construction of Confidence Interval

1)  $\beta_0$ (Intercept):

$$\hat{\beta}_0 \sim N(\beta_0, \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SXX} \right)) \quad (4.1)$$

$$(n-2) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2 \quad (4.2)$$

We can derive:

$$\frac{\hat{\beta}_0 - \beta_0}{\sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SXX} \right)}} \sim t_{n-2} \quad (4.3)$$

Then, the  $(1-\alpha)100\%$ confidence interval of  $\beta_0$  is:

$$\left[ \hat{\beta}_0 - t_{(n-2, \frac{\alpha}{2})} \sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SXX} \right)}, \hat{\beta}_0 + t_{(n-2, \frac{\alpha}{2})} \sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{SXX} \right)} \right] \quad (4.4)$$

2)  $\beta_1$ (Slope):

$$\hat{\beta}_1 \sim N(\beta_1, \sigma^2 \frac{1}{SXX}) \quad (4.5)$$

$$(n-2) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2 \quad (4.6)$$

We can derive:

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\hat{\sigma}^2 \frac{1}{SXX}}} \sim t_{n-2} \quad (4.7)$$

Then, the  $(1-\alpha)100\%$ confidence interval of  $\beta_1$  is:

$$\left[ \hat{\beta}_1 - t_{(n-2, \frac{\alpha}{2})} \sqrt{\hat{\sigma}^2 \frac{1}{SXX}}, \hat{\beta}_1 + t_{(n-2, \frac{\alpha}{2})} \sqrt{\hat{\sigma}^2 \frac{1}{SXX}} \right] \quad (4.8)$$

3) $\hat{y}_* | x_*$  (**Prediction**):

Prediction  $\hat{y}_*|x_*$  is the value of y given that value of  $x_*$  is **observed**.

$$\hat{y}_*|x_* = \hat{\beta}_0 + \hat{\beta}_1 x_* \quad (4.9)$$

Thus, variance of prediction is the deviation of prediction value and the mean value of y.

$$Var(\hat{y}_*|x_*) = \sigma^2 \left( 1 + \frac{1}{n} + \frac{(x_* - \bar{x})^2}{SXX} \right) \quad (4.10)$$

$$\widehat{Var}(\hat{y}_*|x_*) = \hat{\sigma}^2 \left( 1 + \frac{1}{n} + \frac{(x_* - \bar{x})^2}{SXX} \right) \quad (4.11)$$

Therefore, the prediction interval (confidence interval) of  $\hat{y}_*|x_*$  is:

$$\left[ \hat{y}_*|x_* - t_{(n-2, \frac{\alpha}{2})} \sqrt{\hat{\sigma}^2 \left( 1 + \frac{1}{n} + \frac{(x_* - \bar{x})^2}{SXX} \right)}, \hat{y}_*|x_* + t_{(n-2, \frac{\alpha}{2})} \sqrt{\hat{\sigma}^2 \left( 1 + \frac{1}{n} + \frac{(x_* - \bar{x})^2}{SXX} \right)} \right] \quad (4.12)$$

4) $\hat{y}|x_*$  (**Fitted Value**):

Fitted value is:

$$E(Y|X = x_*) = \hat{\beta}_0 + \hat{\beta}_1 x_* \quad (4.13)$$

Here we estimate the **mean of the Y population** associated with  $x_*$ .

$$Var(\hat{y}|x_*) = \sigma^2 \left( \frac{1}{n} + \frac{(x_* - \bar{x})^2}{SXX} \right) \quad (4.14)$$

$$\widehat{Var}(\hat{y}|x_*) = \hat{\sigma}^2 \left( \frac{1}{n} + \frac{(x_* - \bar{x})^2}{SXX} \right) \quad (4.15)$$

The confidence interval of the fitted value is:

$$\hat{y}|x_* \pm t_{(n-2, \frac{\alpha}{2})} \sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{(x_* - \bar{x})^2}{SXX} \right)} \quad (4.16)$$

The simultaneous confidence band of  $\hat{y}|x_*$  is:

$$\left[ \hat{y}|x_* - [2F_{(a,2,n-2)}]^{\frac{1}{2}} \sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{(x_* - \bar{x})^2}{SXX} \right)}, \hat{y}|x_* + [2F_{(a,2,n-2)}]^{\frac{1}{2}} \sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{(x_* - \bar{x})^2}{SXX} \right)} \right] \quad (4.17)$$

(Please refer to "Statistical Inference" by G. Casella & R. Berger ).

5)  $\sigma^2$  (Variance):

$$(n - 2) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2 \quad (4.18)$$

The  $(1-\alpha)100\%$  confidence interval of  $\sigma^2$  is:

$$\left[ \frac{(n-2)\hat{\sigma}^2}{\chi_{(1-\frac{\alpha}{2}, n-2)}^2}, \frac{(n-2)\hat{\sigma}^2}{\chi_{(\frac{\alpha}{2}, n-2)}^2} \right] \quad (4.19)$$

## 5 Exercises

### 5.1 Exercise 3 Q1

### 5.2 Exercise 3 Q2

### 5.3 Example from 12-13 Midterm(similar to Ex3 Q3)

In simple linear regression, if the value of the predictor  $X$  is replaced by  $cX$ , where  $c$  is some non-zero constant, which of the following will be affected? (Circle the answer(s))

(a)  $\hat{\beta}_0$ , (b)  $\hat{\beta}_1$ , (c)  $\hat{\sigma}^2$ , (d)  $R^2$ , (e) t-test statistic of the null hypothesis  $H_0 : \beta_1 = 0$ .

Hint:

$$\hat{\beta}_{1(new)} = \frac{c \sum (x_i - \bar{x})(y_i - \bar{y})}{c^2 \sum (x_i - \bar{x})^2} = \frac{1}{c} \hat{\beta}_1.$$

$$\hat{\beta}_{0(new)} = \bar{y} - (\frac{1}{c} \hat{\beta}_1)(c \bar{x}) = \bar{y} - \hat{\beta}_1 \bar{x} = \hat{\beta}_0.$$

$$\hat{e}_{i(new)} = y_i - \hat{\beta}_0 - \hat{\beta}_{1(new)} x_{i(new)} = y_i - \hat{\beta}_0 - (\frac{1}{c} \hat{\beta}_1)(c x_i) = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i = \hat{e}_i,$$

$$\therefore \hat{\sigma}_{(new)}^2 = \frac{\sum \hat{e}_i^2}{n-2} = \hat{\sigma}^2.$$

$$\therefore \hat{e}_{i(new)} = \hat{e}_i, RSS_{(new)} = RSS.$$

$\therefore SYY = \sum (y_i - \bar{y})^2$  is unaffected,

$$\therefore R_{(new)}^2 = 1 - \frac{RSS}{SYY} = R^2.$$

$$t_{(new)} = \frac{\hat{\beta}_0}{se(\hat{\beta}_0)_{(new)}}, \text{ where } se(\hat{\beta}_0)_{(new)} = \hat{\sigma} \left( \frac{1}{n} + \frac{\bar{x}_{(new)}^2}{SXX_{(new)}} \right)^{1/2}.$$

$$\therefore \frac{\bar{x}_{(new)}^2}{SXX_{(new)}} = \frac{c^2 \bar{x}^2}{c^2 \sum (x_i - \bar{x})} = \frac{\bar{x}^2}{\sum (x_i - \bar{x})} = \frac{\bar{x}^2}{SXX},$$

$$\therefore t_{(new)} = t.$$

## 5.4 Exercise 3 Q5 2.12.4

## 5.5 Example from 11-12 Midterm

1. (30) Fill in the missing values in the following tables of regression output .

| ANOVA Table |                |       |             |              |         |
|-------------|----------------|-------|-------------|--------------|---------|
| Source      | Sum of Squares | d.f.  | Mean Square | F-statistics | p-value |
| Regression  | _____          | _____ | _____       | _____        | 9.5e-09 |
| Residuals   | _____          | _____ | _____       | _____        | _____   |
| Total       | _____          | _____ | _____       | _____        | _____   |

| Coefficient Table |                        |               |              |         |
|-------------------|------------------------|---------------|--------------|---------|
| Variable          | Coefficient            | s.e.          | t-statistics | p-value |
| Constant          | _____                  | _____         | _____        | 0.00322 |
| X                 | -2.04245               | _____         | _____        | _____   |
| $n =$ _____       | $\hat{\sigma} =$ _____ | $R^2 =$ _____ | _____        | _____   |

p.s. It is found in R that  $qf(1 - 9.5e^{-9}, 1, 6) = 1917.3$ . Also  $\bar{x} = 5.125$ ,  $\bar{y} = -9.1974$ ,  $SXX = 54.875$ .

**Solution:**

| ANOVA Table |          |    |          |        |         |
|-------------|----------|----|----------|--------|---------|
| Source      | SS       | Df | MS       | F      | p-value |
| Regression  | 228.9167 | 1  | 228.9167 | 1917.3 | 9.5e-09 |
| Residuals   | 0.7164   | 6  | 0.1194   |        |         |
| Total       | 229.6331 |    |          |        |         |

| Coefficient Table |                         |                |         |         |
|-------------------|-------------------------|----------------|---------|---------|
| Variable          | Coefficient             | s.e.           | t       | p-value |
| Constant          | 1.2702                  | 0.2685         | 4.7313  | 0.00322 |
| X                 | -2.04245                | 0.04665        | 43.7861 | 9.5e-09 |
| $n = 8$           | $\hat{\sigma} = 0.3455$ | $R^2 = 0.9969$ |         |         |

$$\therefore qf(1 - 9.5 \times 10^{-9}, 1, 6) = 1917.3$$

$$\therefore F = 1917.3$$

$$t_{\beta_1}^2 = F, \therefore |t_{\beta_1}| = \sqrt{1917.3} = 43.7861$$

the corresponding p - value =  $9.5^{-9}$

$$df_{reg} = 1, df_{residuals} = 6, n = df_{residuals} + 2 = 8.$$

$$s.e.(\hat{\beta}_1) = |\hat{\beta}_1| / |t_{\beta_1}| = 0.04665$$

$$\hat{\beta}_0 = \bar{y} - \bar{x}\hat{\beta}_1 = 1.2702$$

$$SXY = \hat{\beta}_1 SXX = -112.0794$$

$$SS_{reg} = \frac{SXY^2}{SXX} = 228.9167$$

$$MS_{reg} = SS_{reg}/1 = 228.9167$$

$$MS_{residuals} = MS_{reg}/F = 0.1194$$

$$SS_{residuals} = 6 \times MS_{residuals} = 0.7164$$

$$SS_{total} = SS_{residuals} + SS_{reg} = 229.6331$$

$$\hat{\sigma}^2 = MS_{residuals} = 0.1194, \therefore \hat{\sigma} = 0.3455$$

$$R^2 = \frac{SS_{reg}}{SS_{total}} = 0.9969$$

$$s.e.(\hat{\beta}_0) = \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{SXX}} = 0.2685$$

$$t_{\beta_0} = \frac{\hat{\beta}_0}{s.e.(\hat{\beta}_0)} = 4.7313$$