

LINEAR REGRESSION MODELS: Homework #2

Due on September 25, 2017

Professor Jingchen Liu

Fan Yang
UNI: fy2232

Problem 1

a.

My answer is yes. There is a linear association between Y and X. Because in the diagram above, we can find 95% confidence interval. And the estimated slope value is 0.755048, which lies in the 95% confidence interval [0.452886, 1.05721]. What's more, it is now obvious that the implied level of significance is 95%.

b.

The plausible reason for this question is 'dollar sales cannot be negative even if the population in a district is zero.' But actually population in a district can not be 0. What's more, $x=0$ is not included in our sample space. So there is no sense considering the value of sales when $x=0$.

Problem 2

a.

The summary of model in Problem 1.19 lists as below:

```
> summary(lm.pr19)

Call:
lm(formula = da$V1 ~ da$V2)

Residuals:
    Min       1Q   Median       3Q      Max
-2.74004 -0.33827  0.04062  0.44064  1.22737

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   2.11405     0.32089   6.588 1.3e-09 ***
da$V2          0.03883     0.01277   3.040 0.00292 **
---
Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
```

As we already know the confidence interval of β_1 : $[\hat{\beta}_1 - \xi_{0.005}SD(\beta_1), \hat{\beta}_1 + \xi_{0.005}SD(\beta_1)]$

Since $\xi_{0.005} = 2.58$, and we know $SD(\beta_1) = 0.01277$, therefore, the CI is:

$[0.03883 - 2.58 \times 0.01277, 0.03883 + 2.58 \times 0.01277]$, which is
 $[0.0058834, 0.0717766]$

this interval does not include 0.

if the interval includes 0, which means GPA(Y) can not be predicted by ACT test score(X) since they shows little asso ciation. And this conclusion could reach a confidence level of 99%.

b.

two alternatives are:

$$\begin{aligned} H_0 : \beta_1 &= 0 \\ H_1 : \beta_1 &\neq 0 \end{aligned}$$

test statistic is

$$t^* = \frac{\hat{\beta}_1}{sd(\beta)}$$

when $t^* > t(1 - \alpha/2, n - 2)$, we reject H_0 ; otherwise, conclude H_0 .

$$\text{Calculate } t^* \text{ first: } t^* = \frac{\hat{\beta}_1}{sd(\beta)} = \frac{0.03883}{0.01277} = 3.04072$$

And we can get from table that $t(1 - \alpha/2, n - 2) = t(1 - 0.005, 118) \approx 2.58 < t^* = 3.04072$.

Therefore, we can conclude that $\beta_1 \neq 0$.

Which means linear association exists between student's ACT score (X) and GPA.

c.

The two-sided P-value for the sample outcome is obtained by first finding the onesided P-value, $Pr(|t(118)| > t^* = 3.04072)$.

We use R to compute that this probability is about 0.003.

Since the P-value is less 0.01, which is the level of significance. Therefore we can conclude H_1 directly.

Problem 3

a. we use R to get the summary of the model:

```
Call:
lm(formula = d22$V1 ~ d22$V2)

Residuals:
    Min       1Q   Median       3Q      Max
-5.1500 -2.2188  0.1625  2.6875  5.5750

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  168.60000     2.65702   63.45  < 2e-16 ***
d22$V2        2.03438     0.09039   22.51 2.16e-12 ***
---
Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
```

The regression function is $Y_i = 168.6 + 2.03438X_i$

We already know that $\frac{\hat{\beta}_1 - \beta_1}{sd(\hat{\beta}_1)}$ follows $t(n - 2)$

Thus, $[\hat{\beta}_1 - t(0.005, 14)sd(\hat{\beta}_1), \hat{\beta}_1 + t(0.005, 14)sd(\hat{\beta}_1)]$ is the confidence interval.

$$\hat{\beta}_1 = 2.03438 \text{ and } t(0.005, 14) = 2.98 \text{ and } sd(\hat{\beta}_1) = 0.09039$$

So the confidence interval is $[2.03438 - 0.09039 \times 2.98, 2.03438 + 0.09039 \times 2.98]$

which is $[1.765018, 2.303742]$

This is also interval estimate for change in the mean hardness.

b. two alternatives are:

$$\begin{aligned} H_0 : \beta_1 &= 2 \\ H_1 : \beta_1 &\neq 2 \end{aligned}$$

test statistic is

$$t^* = \frac{\hat{\beta}_1 - \beta_1}{sd(\hat{\beta}_1)}$$

when $t^* > t(1 - \alpha/2, n - 2)$, we reject H_0 ; otherwise, conclude H_0 .

$$\text{Calculate } t^* \text{ first: } t^* = \frac{\hat{\beta}_1 - \beta_1}{sd(\hat{\beta}_1)} = \frac{2.03438 - 2}{0.09039} = 0.3803518$$

And we can get from table that $t(1 - \alpha/2, n - 2) = t(1 - 0.005, 14) \approx 2.98 > t^* = 0.3803518$.

Therefore, we can not reject H_0 . So we conclude that $\beta_1 = 2$.

The two-sided P-value for the sample outcome is obtained by first finding the onesided P-value, $Pr(|t(14)| > t^* = 0.3803518)$.

We use R to compute that this probability is about 0.7093927. So P-value is 0.7093927.

c.

From textbook equation (2.27) we know that

$$\begin{aligned} \delta &= \frac{|\beta_1 - \beta_{10}|}{\sigma(\hat{\beta}_1)} \\ &= \frac{0.3}{\sigma(0.1)} = \frac{0.3}{\sigma(0.1)} = 3 \end{aligned}$$

And $\alpha = 0.01$ $\delta = 3$, we can find the power using table: Power=0.53.

Problem 4

Let's derive $\sigma^2\{\hat{Y}_h\}$ first:

$$\sigma^2\{\hat{Y}_h\} = \sigma^2 \left[\frac{1}{n} + \frac{(X_h - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right] = \sigma^2 \left[\frac{1}{n} + \frac{\frac{1}{n}(X_h - \bar{X})^2}{\frac{1}{n}\sum (X_i - \bar{X})^2} \right]$$

as n becomes large, $\frac{1}{n} \rightarrow 0$

According to law of large number, when n is large enough

$$\frac{1}{n} \sum (X_i - \bar{X})^2 \rightarrow E((X_i - \bar{X})^2) ; \text{ While } E(X_i - \bar{X}) = 0$$

$$\text{Thus } \frac{1}{n} \sum (X_i - \bar{X})^2 \rightarrow \text{var}(X)$$

$$\text{Therefore } \sigma^2\{\hat{Y}_h\} \rightarrow \sigma^2 \left[\frac{1}{n} + \frac{(X_h - \bar{X})^2}{n \times \text{var}(X)} \right] \rightarrow \sigma^2 [0 + 0] = 0$$

Now we will consider $\sigma^2\{pred\}$:

$$\sigma^2\{pred\} = \sigma^2 + \sigma^2\{\hat{Y}_h\}$$

Although $\sigma^2\{\hat{Y}_h\} \rightarrow 0$ while n becomes large,

σ^2 is variance of the distribution of Y at $X = X_h$ and will always larger than 0

So $\sigma^2\{pred\}$ can not be close to 0

We can conclude that with sample size becomes larger, we can get accurate mean of Y ;
but the prediction mean could not be accurately got no matter how large n is.

Problem 5

a.

The $1 - \alpha$ confidence limits are:

$$\hat{Y}_h \pm t(1 - \alpha/2; n - 2)s\{\hat{Y}_h\}$$

where

$$\hat{Y}_h = 2.11405 + 0.03883 * 28 = 3.20129$$

$$s^2\{\hat{Y}_h\} = \sigma^2 \left[\frac{1}{n} + \frac{(X_h - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right]$$

$$= \sigma^2 \left[\frac{1}{n} + \frac{(X_h - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right]$$

$$\text{While } \hat{\sigma}^2 = MSE = 0.3882848$$

$$\frac{(X_h - \bar{X})^2}{\sum (X_i - \bar{X})^2} = \frac{(28 - 24.725)^2}{\sum (X_i - 24.725)^2} = 0.004506707$$

$$\text{So } s\{\hat{Y}_h\} = \sqrt{0.3882848 * \left(\frac{1}{118} + 0.004506707 \right)} = 0.07099602$$

And we know $t(0.975, 118) = 1.980272$, therefore the confidence limits is

$$3.20129 \pm 1.980272 * 0.07099602 \text{ which is } [3.060699, 3.341881]$$

b.

The $1 - \alpha$ confidence interval is:

$$\hat{Y}_h \pm t(1 - \alpha/2; n - 2)s\{pred\}$$

We already know that:

$$\hat{Y}_h = 3.20129 \quad MSE = 0.3882848 \quad s\{\hat{Y}_h\} = 0.07099602 \text{ and } t(0.975, 118) = 1.980272$$

Then we need to get $s\{pred\}$

$$s^2\{pred\} = MSE + s^2\{\hat{Y}_h\} = 0.3882848 + 0.07099602^2 = 0.3933252$$

$$s\{pred\} = 0.6271564$$

Therefore the confidence limits is

$$3.20129 \pm 1.980272 * 0.6271564 \text{ which is } [1.95935, 4.44323]$$

c.

The prediction interval in part (b) is wider than the confidence interval in part (a).

Because the prediction for new observation involves new error randomness. In order to get the same confidence level, it must be wider.

d.

The two boundary values at $1 - \alpha$ confidence level is:

$$\hat{Y}_h \pm W s\{\hat{Y}_h\}$$

where

$$W^2 = 2F(1 - \alpha; 2, n - 2) = 2F(0.95; 2, 118) = 6.146181$$

So $W = 2.479149$

We already know that:

$$\hat{Y}_h = 3.20129 \quad MSE = 0.3882848 \quad s\{\hat{Y}_h\} = 0.07099602$$

Therefore the confidence band are

$$3.02528 \leq \beta_0 + \beta_1 X_h \leq 3.3773$$

The confidence band is wider at this point than the confidence interval in part (a)

Because the bands must cover all possible values no matter what observation is, which means the same level of confidence interval estimate of the mean must be included in the bands.

Problem 6

$$\begin{aligned} b_0 &= \bar{Y} - b_1 \bar{X} \\ E(b_0) &= E(\bar{Y} - b_1 \bar{X}) \\ &= E(\beta_0 + \beta_1 \bar{X} + \bar{\varepsilon} - b_1 \bar{X}) \\ &= E(\beta_0 + \bar{X}(\beta_1 - b_1)) \\ &= \beta_0 + \bar{X}(\beta_1 - E(b_1)) \quad (b_1 \text{ is unbiased.}) \\ &= \beta_0 \end{aligned}$$

Thus, b_0 is a unbiased estimator of β_0 .

Problem 7

$$\begin{aligned}
\sigma^2\{b_0\} &= \text{var}(\bar{y} - b_1\bar{X}) \\
&= \text{var}(\bar{y}) + \text{var}(b_1\bar{X}) - 2\text{cov}(\bar{y}, b_1\bar{X}) \\
&= \text{var}(\bar{y}) + \bar{X}^2 \text{var}(b_1) - 2\bar{X}^2 \text{cov}(\bar{y}, b_1) \quad (\text{using (2.31), then :}) \\
&= \text{var}(\bar{y}) + \bar{X}^2 \text{var}(b_1) \\
&= \text{var}\left(\frac{1}{n} \sum y_i\right) + \bar{X}^2 \frac{\sigma^2}{\sum (X_i - \bar{X})^2} \\
&= \frac{1}{n^2} \sum \text{var}(y_i) + \frac{\sigma^2 \bar{X}^2}{\sum (X_i - \bar{X})^2} \\
&= \frac{1}{n^2} \sum \sigma^2 + \frac{\sigma^2 \bar{X}^2}{\sum (X_i - \bar{X})^2} \\
&= \frac{1}{n} \sigma^2 + \frac{\sigma^2 \bar{X}^2}{\sum (X_i - \bar{X})^2} \\
&= \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum (X_i - \bar{X})^2} \right)
\end{aligned}$$

$$\sigma^2\{\hat{Y}_h\} = \sigma^2 \left(\frac{1}{n} + \frac{(X_h - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right)$$

$$\text{while } X_h \text{ is 0 } \sigma^2\{\hat{Y}_h\} = \sigma^2\{b_0 + b_1 * 0\} = \sigma^2\{b_0\} = \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum (X_i - \bar{X})^2} \right)$$

That's to say, variance (2.22b) is a special case of variance (2.29b) when $X_h = 0$

Problem 8

the $1 - \alpha$ confidence limits for β_1 are:

$$b_1 \pm t(1 - \alpha/2; n - 2) \text{sd}\{b_1\}$$

For the **first** region

```

> b_1 <- lm.cdi1$coefficients[2]
> b_1
522.1588
> t <- qt(0.95, lm.cdi1$df.residual)
> sd2 <- sum((lm.cdi1$residuals)^2) / lm.cdi1$df.residual / sum((X - mean(X))^2)
> b_1 - t * sqrt(sd2)
460.5177
> b_1 + t * sqrt(sd2)
583.8

```

interval estimate of β_1 is [460.5177, 583.8]

For the **second** region

```
> b_1 <- lm.cdi2$coefficients[2]
> b_1
238.6694
> t <- qt(0.95,lm.cdi2$df.residual)
> sd2 <- sum((lm.cdi2$residuals)^2)/lm.cdi2$df.residual / sum((X-mean(X))^2)
> b_1-t*sqrt(sd2)
193.4858
> b_1+t*sqrt(sd2)
283.853
```

interval estimate of β_1 is [193.4858, 283.853]

For the **third** region

```
> b_1 <- lm.cdi3$coefficients[2]
> b_1
330.6117
> t <- qt(0.95,lm.cdi3$df.residual)
> sd2 <- sum((lm.cdi3$residuals)^2)/lm.cdi3$df.residual / sum((X-mean(X))^2)
> b_1-t*sqrt(sd2)
285.7076
> b_1+t*sqrt(sd2)
375.5158
```

interval estimate of β_1 is [285.7076, 375.5158]

For the **forth** region

```
> b_1 <- lm.cdi4$coefficients[2]
> b_1
440.3157
> t <- qt(0.95,lm.cdi4$df.residual)
> sd2 <- sum((lm.cdi4$residuals)^2)/lm.cdi4$df.residual / sum((X-mean(X))^2)
> b_1-t*sqrt(sd2)
364.7585
> b_1+t*sqrt(sd2)
515.8729
```

interval estimate of β_1 is [364.7585, 515.8729]

Different regression lines of the 4 regions have different slopes as shown above. Their values of bound of interval vary from each other very much.