

p8130\_final\_jn2761

jiaying Ning

12/11/2020

### Problem 1

(a)

- **Outcome:**
  - Antibodies : level of antibodies at 6-month
- **Predictor:**
  - age : age (years) of patients at baseline
  - Vaccine 1
  - Vaccine 2
  - Placebo
- **equation of the linear regression model**
  - $\text{Antibodies}(i) = \beta_0 + \beta_1 \text{age}(i) + \beta_2 I(\text{Trt}=\text{Vaccine1})(i) + \beta_3 I(\text{Trt}=\text{Vaccine2})(i) + \text{error}$
  - $i=1,2,\dots,30$
  - With placebo as reference category
- **distribution assumption imposed on the error terms**
  - Error terms should be uncorrelated and follow a normal distribution with mean = 0 and a constant variance.  $N(0, \text{var})$

(b)

### Problem 2

(a)

$$\text{salaries}(\text{hat}) = 17.847 + 1.103 * \text{Years}(\text{hat}) + 0.322 * \text{Publications}(\text{hat}) + 1.593 * \text{Gender}(\text{hat}) + 1.289 * \text{Grants}(\text{hat})$$

(b)

$$\text{slopeGender} + (-)t(30, 0.975) * \text{se}(\text{Gender})$$

```
qt(0.975, 30)
```

```
## [1] 2.042272
```

$$\varepsilon \sim N(0, \sigma^2)$$

(b)

$$\begin{bmatrix} \text{Antibodies} \\ y_1 \\ y_2 \\ \vdots \\ y_{30} \end{bmatrix} = \begin{bmatrix} 1 & \text{Age}_1 & 0 & 0 \\ 1 & \text{Age}_2 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & \text{Age}_{30} & 0 & 0 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_{30} \end{bmatrix}$$

↓

$$\text{Antibodies}_i = \beta_0 + \text{Age}_i \beta_1 + I(\text{Vaccine1}) \beta_2 + I(\text{Vaccine2}) \beta_3 + \varepsilon_i$$

with placebo as reference  $i = 1 \dots 30$

Figure 1: Matrix notation

```
gender_upper=1.593+(qt(0.975,30)*0.688)
gender_lower=1.593-(qt(0.975,30)*0.688)
```

- 95% CI for the coefficient associated with 'Gender' is (0.1879166,2.9980834)

(c)

### Hypothesis

- H0: There is no linear association exists between publication and faculty salaries (in thousands of dollars).  $\beta_{\text{publication}} = 0$
- H1: There is linear association exists between publication and faculty salaries (in thousands of dollars).  $\beta_{\text{publication}} \neq 0$

### Calculations

$$(\beta_{\text{publication}} - 0) / \text{se}(\text{publication}) = 0.322 / 0.521 = 0.6180422$$

### Decision

```
qt(0.975,30)
```

```
## [1] 2.042272
```

- **Critical Value:**  $t_{\sim}(30,0.975) = 2.042272$
- **Decision Rule:**
  - Reject  $H_0$ : if  $|t| > t(30,0.975)$
  - Fail to reject  $H_0$ :  $|t| < t(30,0.975)$

## Conclusion

- For the current data, we have test stats of 0.6180422 for slope. Since  $t\_stats:0.6180422 < \text{critical value } 2.042272$ , we fail to reject the null and conclude that there is no significant linear association exists between publication and faculty salaries (in thousands of dollars).

**comment** - Now we know that there is no significant linear association exists between publication and faculty salaries, we can conclude that the the most significant variables in predicting salary is the success in obtaining grant support.

(d)

- We could not generalize our findings to the entire academic community (i.e., US) because even though we yield some significance result from our current data, all of our samples are drawn from the MSPH faculty which might not represent all faculty in the United States. Differences between department and location might make our current conclusion not applicable.

## Problem 3

Among the three models, I would recommend model B.

- **Mallow's Cp Criteria:** Cp compares the predictive ability for each model to the full model, we want to choose Cp values that is  $\leq$  number of parameter. From the three models, Model A's Cp value is larger than the number of parameter. Model B and C both have Cp smaller than the number of parameter, but model B has Cp value that are closer to the number of parameter. Also Model B has smaller number of parameter but similar Cp, which meet the goal of model selection: "get a small number of variables while maintaining the same predictive ability"
- **Adjusted R-squared:** Adjusted  $R^2$  tells us how well the variance of outcome is explained by the relationship with predictors. Therefore, we will want to have models that have larger adjusted R-squared value. Overall, Model B has the largest adjusted R-squared value.
- **MSE:**
  - We first compare the MSE between models and see that all three models have similar MSE. Model B and Model C have smaller MSE than Model A.
  - We then compare MSE on Testing and Training data within each model. Which I calculate the difference between MSE and MSPE within each model. The difference between MSE and MSPE is also smallest in Model B comparing to other models.

0.078-0.045

## [1] 0.033

```
0.072-0.043
```

```
## [1] 0.029
```

```
0.079-0.042
```

```
## [1] 0.037
```

- Therefore, overall, I would recommend model B.

## Problem 4

(a)

```
#Import Data
hospital_df=
  read.csv("./data/Hospital.csv")
```

### fit regression model

fit a simple linear regression with length of stay (LOS) as the outcomes and number of beds (BEDS) as only predictor. Use R to address the following points:

```
linear = lm(LOS ~ BEDS, data=hospital_df)
```

summary

```
summary(linear)
```

```
##
## Call:
## lm(formula = LOS ~ BEDS, data = hospital_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.8291 -1.0028 -0.1302  0.6782  9.6933
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  8.6253643   0.2720589   31.704 < 2e-16 ***
## BEDS         0.0040566   0.0008584    4.726 6.77e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.752 on 111 degrees of freedom
## Multiple R-squared:  0.1675, Adjusted R-squared:  0.16
## F-statistic: 22.33 on 1 and 111 DF, p-value: 6.765e-06
```

95% CI for the true slope

$$\text{slope} + (-)t(111, 0.975) * se(\text{beta})$$

```
confint(linear, level=0.95)
```

```
##                2.5 %        97.5 %  
## (Intercept) 8.086261517 9.164467086  
## BEDS        0.002355649 0.005757623
```

- 95% CI for the true slope is (0.002355649, 0.005757623)

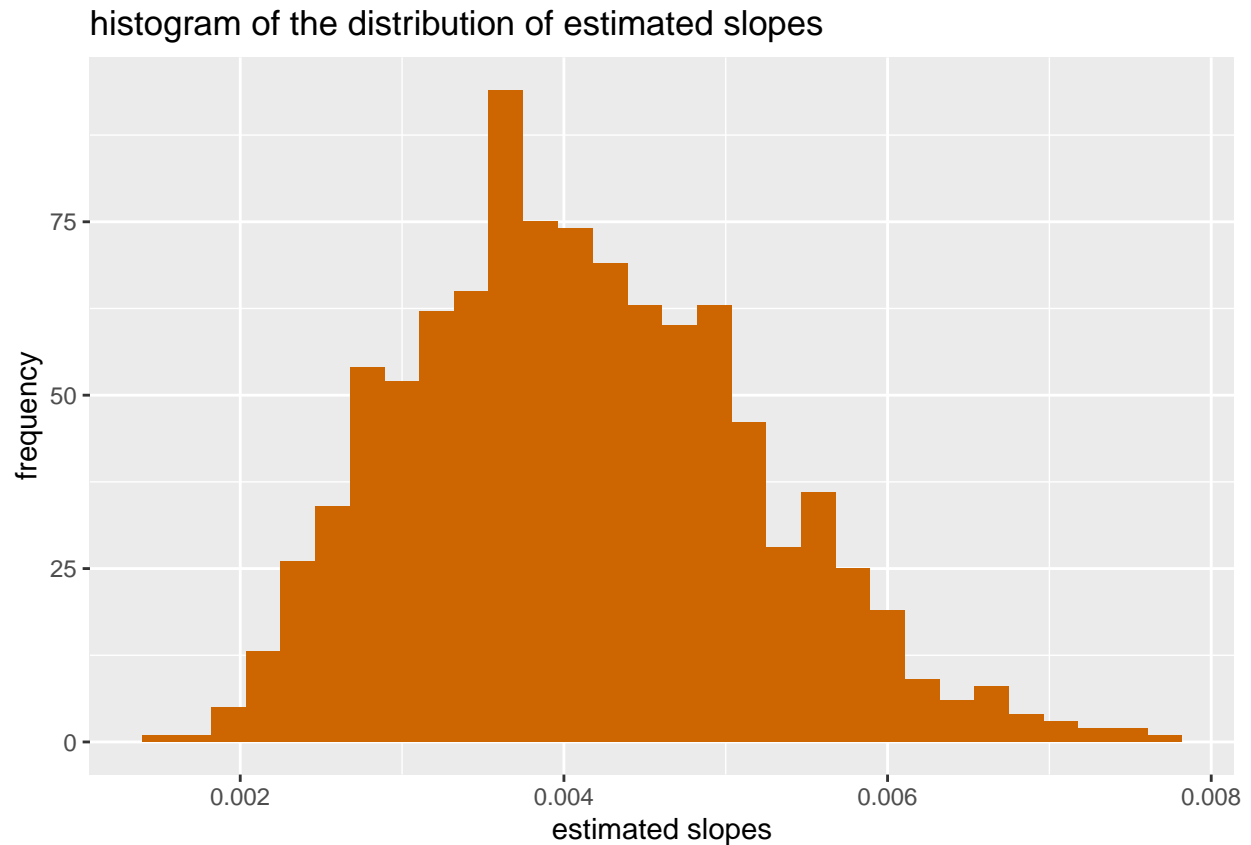
(b)

```
library(infer)  
boot_dist <- hospital_df %>%  
  specify(formula = LOS ~ BEDS) %>%  
  generate(reps = 1000, type = "bootstrap") %>%  
  calculate(stat = "slope")  
# Save the estimated slope values (1,000)  
boot_dist_slope <- boot_dist[[2]]
```

(c)

```
library(ggplot2)  
boot_dist %>%  
  ggplot(aes(x = stat)) + geom_histogram(fill="darkorange3") + labs(  
    titles = "histogram of the distribution of estimated slopes",  
    x = "estimated slopes",  
    y = "frequency"  
  )
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



Comment: In order to make inferences, the sampling distribution of slope should follow normal distribution. From the histogram we see that the estimated slope roughly follow a normal distribution, thus we conclude that after bootstrap, the distribution of slope does approximate the shape of the sampling distribution of  $\beta_1$  which meet the assumption of normal distribution for parameters.

(d)

```
lower=quantile(boot_dist_slope,0.025)
upper=quantile(boot_dist_slope,0.975)
```

- The 95% CI for the true slope is (0.0022891,0.0063338)

(e)

- Before Bootstrap:95% CI for the true slope is (0.002355649,0.005757623)
- After bootstrap:The 95% CI for the true slope is (0.0022891,0.0063338)

Overall, both confidence interval have similar value, meaning the original data does represent the underlying distribution. To be exact, the 95% CI for the true slope after bootstrap has wider interval comparing to the 95% CI before bootstrap.