

# Stat Theory Homework 2.

Bo hao Tang

- 1: (i) Correctly specified means that under true distribution  $P$  there exist  $\beta_0, \beta_1 \in \mathbb{R}$  such that

$$E_p(Y|x) = \beta_0 + \beta_1 x$$

Misspecified means there are no such  $\beta_0, \beta_1$ .  
For example if  $x \in \mathbb{R}$  and  $Y = 1_{\{x < 0\}}$

then  $(X, Y)$  can not be described by this model since  
 $E(Y|x) = 1_{\{x < 0\}} = 0, \text{ or } 1$  but  $\beta_0 + \beta_1 x$   
can only have one value  $\beta_0$  or the whole real line for its range.

Consider  $V = \text{span}\{1, x\}$  then the model constrains parts that orthogonal to  $V$  and not constrain parts in  $V$ .

$$\begin{aligned} \text{(ii)} \quad \hat{\beta} &= \left[ \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_n \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \right]^{-1} \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix} \\ &= \begin{pmatrix} 1 & \bar{x} \\ \bar{x} & \bar{x}^2 \end{pmatrix}^{-1} \begin{pmatrix} \bar{y} \\ \bar{xy} \end{pmatrix} = \frac{1}{\bar{x}^2 - (\bar{x})^2} \begin{pmatrix} \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix} \begin{pmatrix} \bar{y} \\ \bar{xy} \end{pmatrix} \\ &= \frac{1}{\bar{x}^2 - (\bar{x})^2} \begin{pmatrix} \bar{x}^2 \bar{y} - \bar{x} \bar{xy} \\ \bar{xy} - \bar{x} \bar{y} \end{pmatrix} \end{aligned}$$

By WLLN, we know that

$$\hat{\beta} \xrightarrow{P} \frac{1}{E X^2 - (E X)^2} \begin{pmatrix} E X^2 E Y - E X E X Y \\ E X Y - E X E Y \end{pmatrix}$$

$$= \frac{1}{\text{var } X} \begin{pmatrix} \text{var } X E Y - \text{cov}(X, Y) E X \\ \text{cov}(X, Y) \end{pmatrix} = \beta^*$$

For limit distribution of  $\sqrt{n}(\hat{\beta} - \beta^*)$ , we use delta method, we know by CLT that

$$\sqrt{n} \left\{ \begin{pmatrix} \bar{X} \\ \bar{Y} \\ \frac{\bar{XY}}{\bar{X}^2} \end{pmatrix} - \begin{pmatrix} E X \\ E Y \\ E X Y \\ E X^2 \end{pmatrix} \right\} \xrightarrow{d} N(0, \underbrace{\begin{bmatrix} \text{var } X & \text{cov}(X, Y) & \text{cov}(X, XY) & \text{cov}(X, X^2) \\ & \text{var } Y & \text{cov}(Y, X Y) & \text{cov}(Y, X^2) \\ & & \text{var } X Y & \text{cov}(X Y, X^2) \\ & & & \text{var } X^2 \end{bmatrix}}_{D})$$

therefore  $\sqrt{n}(\hat{\beta} - \beta^*) \xrightarrow{d} N(0, V)$

where  $V = W D W^T$

$$\text{where } W = \begin{pmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \frac{\partial f}{\partial x_3} & \frac{\partial f}{\partial x_4} \\ \frac{\partial g}{\partial x_1} & \frac{\partial g}{\partial x_2} & \frac{\partial g}{\partial x_3} & \frac{\partial g}{\partial x_4} \end{pmatrix} \begin{matrix} x_1 = E X, x_2 = E Y, x_3 = E X Y \\ x_4 = E X^2 \end{matrix}$$

$$\text{where } f = \frac{x_4 x_2 - x_1 x_3}{x_4 - x_1^2}, g = \frac{x_3 - x_1 x_2}{x_4 - x_1^2}$$

(iii) if  $X \perp Y$ ,  $\text{cov}(X, Y) = 0 \Rightarrow \beta_1^* = 0$

(iv) see the code part.

(v)  $\hat{\beta}_{\text{ols}} \xrightarrow{P} \beta^*$  in the setting  $(x_i, y_i) \stackrel{\text{i.i.d}}{\sim} p$

We compare it with estimator  $\tilde{\beta}$  where

$$\tilde{\beta}[(x_i, y_i)_{i=1}^n] = \begin{pmatrix} \bar{y} - \tilde{\beta}_1 \bar{x} \\ \tilde{\beta}_1 \end{pmatrix}$$

$$\text{where } \tilde{\beta}_1 = \text{Median} \left\{ \frac{y_{2k} - y_{2k-1}}{x_{2k} - x_{2k-1}} \right\}_{k=1}^{\lfloor \frac{n}{2} \rfloor}$$

Then see the coding part.

2: We need to prove that if there exists two sequences of random vectors  $U_n, V_n$  such that

$$\exists C < +\infty \quad P(\|U_n\| > C) \rightarrow 0$$

$$V_n / \|Y_n\| \xrightarrow{P} \vec{0}$$

$$\text{and } Y_n = U_n + V_n$$

$$\text{Then } \exists C' < +\infty \quad P(\|Y_n\| > C') \rightarrow 0$$

$$\begin{aligned} \text{Proof: Let } C' = 2C \quad \text{then } P(\|Y_n\| > C') &= P(2\|Y_n\| > C' + \|Y_n\|) \\ &= P(\|Y_n\| > \frac{C'}{2} + \frac{1}{2}\|Y_n\|) = P(\|U_n + V_n\| > C + \frac{1}{2}\|Y_n\|) \\ &\leq P(\|U_n\| > C) + P(\|V_n\| > \frac{1}{2}\|Y_n\|) \rightarrow 0 \end{aligned}$$

$$\text{therefore } Y_n = O_p(1)$$

# Coding Part

*Bohao Tang*

*April 1, 2018*

iv

Two distribution is

$$Y = 1 + 2X + \epsilon$$

and

$$Y = 1 + 2X^2 + \epsilon$$

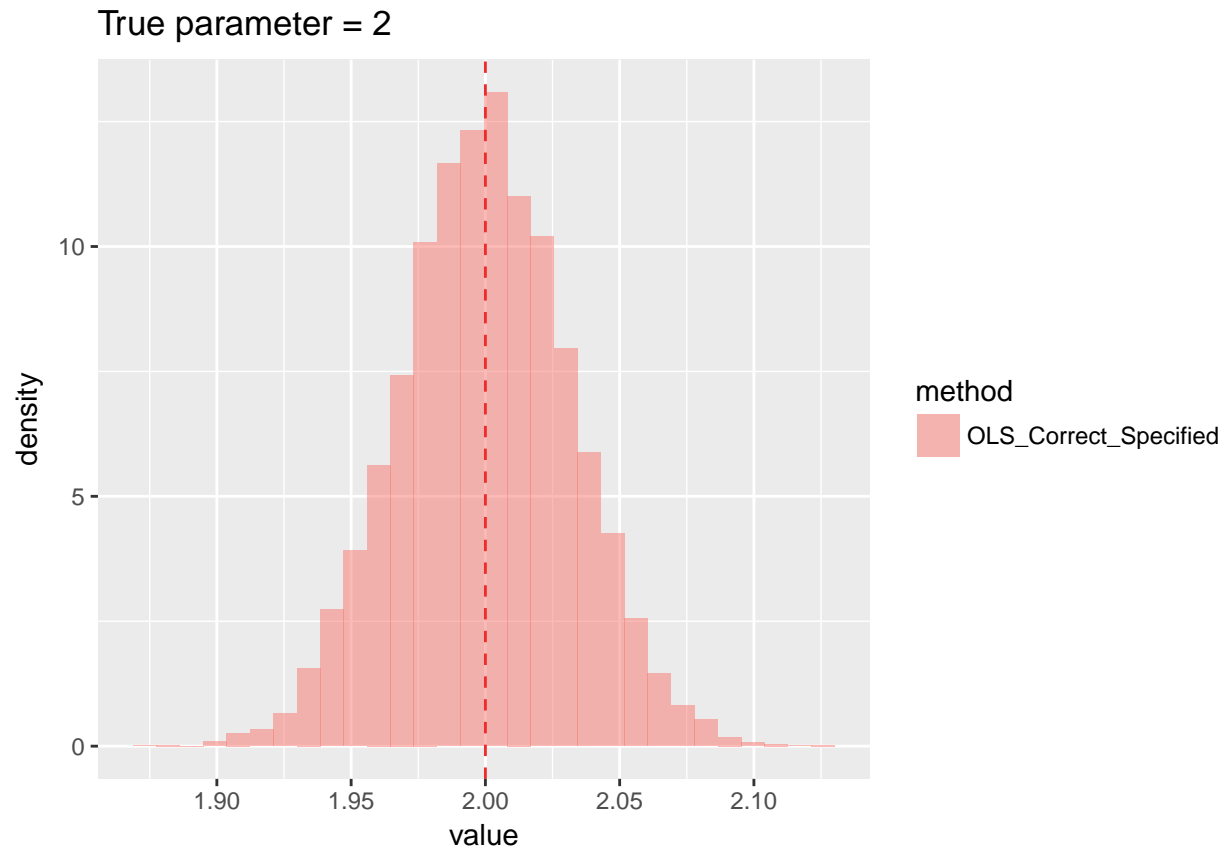
Where  $X, \epsilon$  mutual independent and  $\epsilon, X \sim N(0, 1)$ .

Then the regression  $E(Y|X) = \beta_0 + \beta_1 X$  correctly specifies first distribution and misspecifies the second.

```
modell1 = function(X, eps){  
  1 + 2*X + eps  
}  
modell2 = function(X, eps){  
  1 + 2 * X^2 + eps  
}
```

For the first situation, we run estimation 10000 times, each have 1000 samples and only focus on  $\beta_1$ . The red dashed line is the true  $\beta_1$ :

```
library(ggplot2)  
  
beta1 = c()  
  
n = 1000  
for(i in 1:10000){  
  X = rnorm(n, 0, 1)  
  eps = rnorm(n, 0, 1)  
  
  Y = modell1(X, eps)  
  
  estimator = cov(X, Y) / var(X)  
  
  beta1 = c(beta1, estimator)  
}  
  
data = data.frame(value = beta1)  
data$method = "OLS_Correct_Specified"  
  
ggplot(data, aes(value, fill = method)) +  
  geom_histogram(alpha = 0.5, aes(y = ..density..), position = 'identity') +  
  geom_vline(xintercept = 2, linetype="dashed", color="firebrick2") +  
  ggtitle("True parameter = 2")
```



Then the bias and variance for  $\hat{\beta}_1$  is:

```
bias = mean(beta1) - 2
bias
```

```
## [1] -9.523031e-05
```

```
variance = var(beta1)
variance
```

```
## [1] 0.0009997933
```

In the second situation, now true  $\beta_1$  doesn't exist, we need to compare the result to  $\beta_1^* = \frac{Cov(X,Y)}{Var(X)} = 0$ . We do the same thing as above:

```
betalmis = c()

n = 1000
for(i in 1:10000){
  X = rnorm(n, 0, 1)
  eps = rnorm(n, 0, 1)

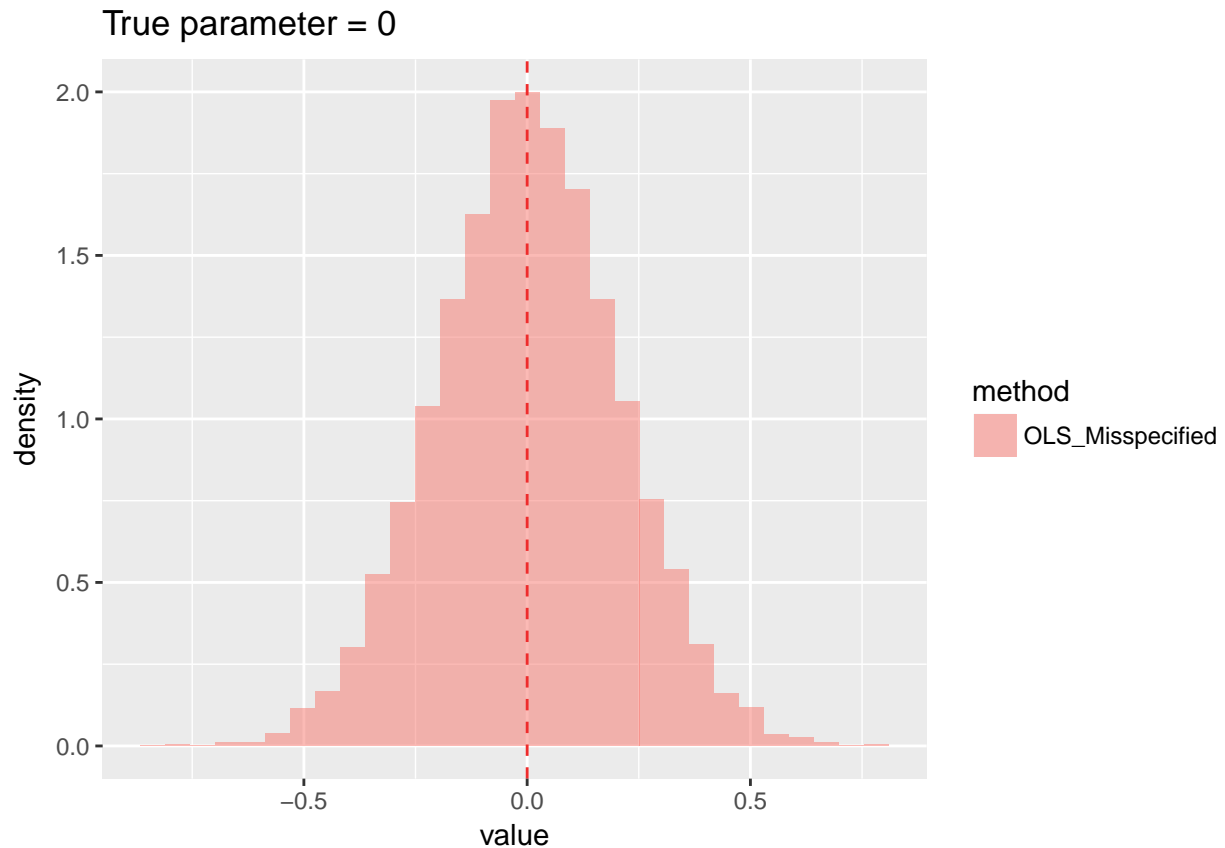
  Y = model2(X, eps)

  estimator = cov(X, Y) / var(X)

  betalmis = c(betalmis, estimator)
}
```

```
datamis = data.frame(value = betalmis)
datamis$method = "OLS_Misspecified"

ggplot(datamis, aes(value, fill = method)) +
  geom_histogram(alpha = 0.5, aes(y = ..density..), position = 'identity') +
  geom_vline(xintercept = 0, linetype="dashed", color="firebrick2") +
  ggtitle("True parameter = 0")
```



Then the bias and variance for  $\hat{\beta}_1^*$  is:

```
bias = mean(betalmis) - 0
bias
```

```
## [1] 0.00111693
```

```
variance = var(betalmis)
variance
```

```
## [1] 0.04135902
```

**v**

We do the model:

$$Y = 1 + 2X + \epsilon$$

where  $X, \epsilon$  mutual independent and  $\epsilon, X \sim N(0, 1)$ .

And the second estimator mentioned in homework solution, which is just a median of slopes. We use data from iv for OLS estimator and for second estimator:

```
slope = c()

n = 1000
index = 1:1000
oddi = index[index%%2 == 1]
eveni = index[index%%2 == 0]
for(i in 1:10000){
  X = rnorm(n, 0, 1)
  eps = rnorm(n, 0, 1)

  Y = model1(X, eps)

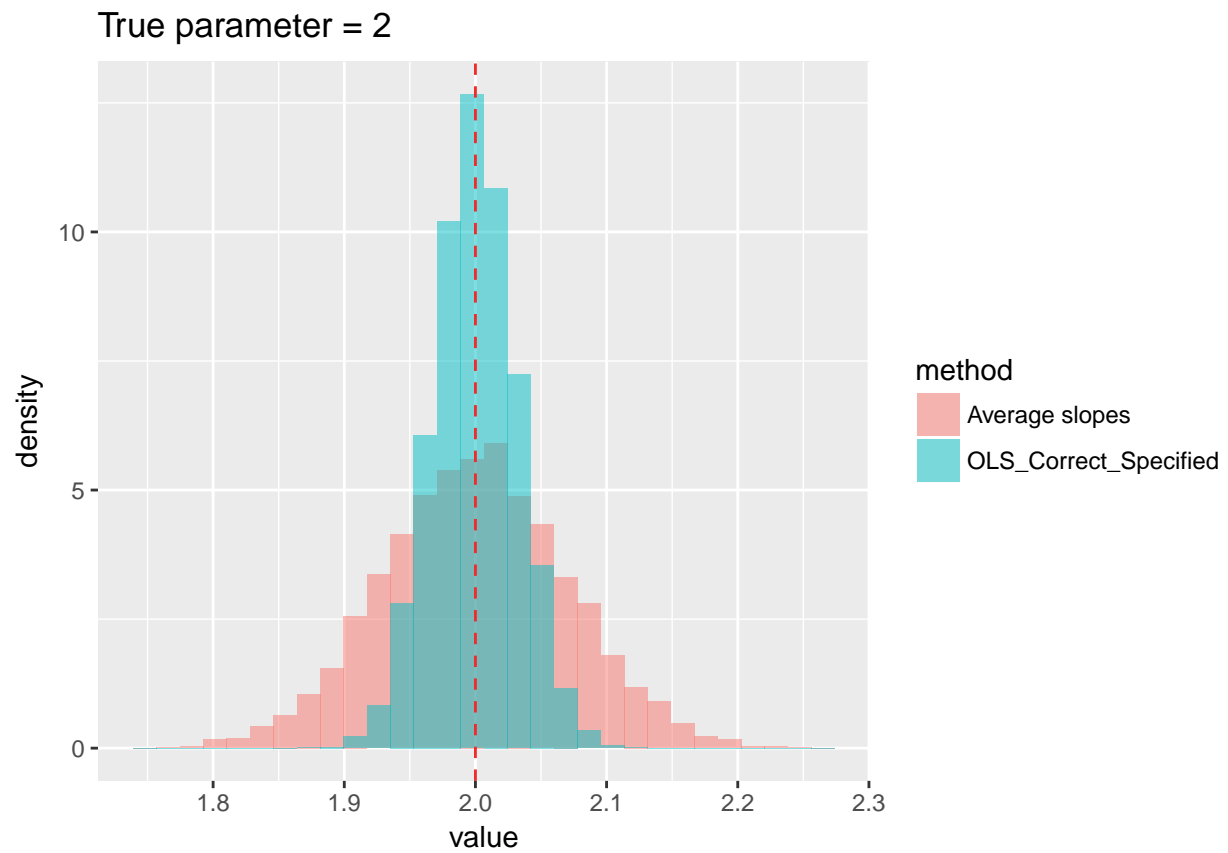
  estimator =
    median( (Y[eveni] - Y[oddi]) / (X[eveni] - X[oddi]) )

  slope = c(slope, estimator)
}

slopes = data.frame(value = slope)
slopes$method = "Average slopes"

estimators = rbind(data, slopes)

ggplot(estimators, aes(value, fill = method)) +
  geom_histogram(alpha = 0.5, aes(y = ..density..), position = 'identity') +
  geom_vline(xintercept = 2, linetype="dashed", color="firebrick2") +
  ggtitle("True parameter = 2")
```



Therefore we can see OLS is a better estimator for  $\beta_1$ .