

# ECON 251 Discussion

## Multiple Linear Regression

Elird Haxhiu

Fall 2022

# Outline

1. Mincer (1972) regression framework
2. Why use multiple (vs simple) linear regression?
3. Omitted variable bias (OVB) formula
4. Multiple linear regression (MLR) model assumptions
5. Ordinary least squares (OLS) estimator
6. Main theorems on bias, consistency, and efficiency
7. Hypothesis tests: t and F
8. Confidence intervals

# Mincer (1972) regression framework

- $\log Y_i \geq 0$  denotes log earnings (outcome)
- $S_i \in \{0,1\}$  is whether  $i$  finished college (treatment)
- $U_i$  is unobserved error term (ex: ability)
- simple linear population regression function (PRF)
- $\beta \approx$  Mincer (1972) returns to college

$$\log Y_i = \alpha + \beta \cdot S_i + U_i$$

- $\log Y_i \geq 0$  denotes log earnings (outcome)
- $S_i \in \{0,1\}$  is whether  $i$  finished college (treatment)
- $U_i$  is unobserved error term (ex: ability)
- simple linear population regression function (PRF)
- $\beta \approx$  Mincer (1972) returns to college

$$\log Y_i = \alpha + \beta \cdot S_i + U_i$$

- Potential outcomes + treatment effects

$$\text{ATE} := E[\log Y_i(1) - \log Y_i(0)]$$

- Need **independence** to *identify* **ATE** with simple **comparison**, which is given by

$$S_i \perp \log Y_i(1), \log Y_i(0) \\ \Leftrightarrow E[U_i|S_i] = E[U_i] = 0$$

$$\hat{\beta}^{OLS}$$

- $\log Y_i \geq 0$  denotes log earnings (outcome)
- $S_i \in \{0,1\}$  is whether  $i$  finished college (treatment)
- $U_i$  is unobserved error term (ex: ability)
- simple linear population regression function (PRF)
- $\beta \approx$  Mincer (1972) returns to college

$$\log Y_i = \alpha + \beta \cdot S_i + U_i$$

- Potential outcomes + treatment effects

$$\text{ATE} := E[\log Y_i(1) - \log Y_i(0)]$$

- Need **independence** to *identify* **ATE** with simple **comparison**, which is given by

$$S_i \perp \log Y_i(1), \log Y_i(0) \\ \Leftrightarrow E[U_i | S_i] = E[U_i] = 0$$

$$\begin{aligned} \hat{\beta}^{OLS} &= \hat{E}[\log Y_i | S_i = 1] - \hat{E}[\log Y_i | S_i = 0] \\ &= \overline{\log Y_1} - \overline{\log Y_0} \\ &= \frac{1}{N_1} \sum_{i|S_i=1} \log Y_i - \frac{1}{N_0} \sum_{i|S_i=0} \log Y_i \end{aligned}$$

# Why use multiple (vs simple) linear regression?

- Multiple regression: we can get closer to satisfying the hypothetical (but necessary, and luckily also sufficient assumption known as) **random assignment/independence** by conditioning on some observable **characteristics  $X_i$**  (provocative example: IQ test score)

# Why use multiple (vs simple) linear regression?

- Multiple regression: we can get closer to satisfying the hypothetical (but necessary, and luckily also sufficient assumption known as) **random assignment/independence** by conditioning on some observable **characteristics**  $X_i$  (provocative example: IQ test score)

$$\log Y_i = \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i$$

$$\begin{aligned} & S_i \perp \log Y_i(1), \log Y_i(0) \mid X_i \\ \Leftrightarrow & E[U_i \mid S_i, X_i] = 0 \end{aligned}$$

- Inclusion of  $X_i$  allows us to “control” for any reasons why there may not be truly random assignment of treatment (in simple PRF)

# Omitted variable bias (OVB)

“True” model

$$\log Y_i = \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i$$

$$\text{Cov}(S_i, U_i) = 0$$

Our model

$$\log Y_i = a + b \cdot S_i + E_i$$

Auxiliary model

$$X_i = c + \gamma \cdot S_i + \eta_i$$



# Omitted variable bias (OVB)

“True” model	$\log Y_i = \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i$	$\text{Cov}(S_i, U_i) = 0$
Our model	$\log Y_i = a + b \cdot S_i + E_i$	
Auxiliary model	$X_i = c + \gamma \cdot S_i + \eta_i$	

By (naively) assuming  $\text{Cov}(S_i, E_i) = 0$  in our model, the population value of slope is

$$b = \frac{\text{Cov}(S_i, \log Y_i)}{\text{Var}(S_i)}$$

$$= \beta + \delta \cdot \gamma$$

# Omitted variable bias (OVB)

“True” model

$$\log Y_i = \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i$$

$$\text{Cov}(S_i, U_i) = 0$$

Our model

$$\log Y_i = a + b \cdot S_i + E_i$$

Auxiliary model

$$X_i = c + \gamma \cdot S_i + \eta_i$$

By (naively) assuming  $\text{Cov}(S_i, E_i) = 0$  in our model, the population value of slope is

$$b = \frac{\text{Cov}(S_i, \log Y_i)}{\text{Var}(S_i)} = \frac{\text{Cov}(S_i, \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i)}{\text{Var}(S_i)}$$

# Omitted variable bias (OVB)

“True” model	$\log Y_i = \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i$	$\text{Cov}(S_i, U_i) = 0$
Our model	$\log Y_i = a + b \cdot S_i + E_i$	
Auxiliary model	$X_i = c + \gamma \cdot S_i + \eta_i$	

By (naively) assuming  $\text{Cov}(S_i, E_i) = 0$  in our model, the population value of slope is

$$\begin{aligned} b &= \frac{\text{Cov}(S_i, \log Y_i)}{\text{Var}(S_i)} = \frac{\text{Cov}(S_i, \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i)}{\text{Var}(S_i)} \\ &= \frac{\text{Cov}(S_i, \alpha) + \text{Cov}(S_i, \beta S_i) + \text{Cov}(S_i, \delta X_i) + \text{Cov}(S_i, U_i)}{\text{Var}(S_i)} \end{aligned}$$

# Omitted variable bias (OVB)

“True” model	$\log Y_i = \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i$	$\text{Cov}(S_i, U_i) = 0$
Our model	$\log Y_i = a + b \cdot S_i + E_i$	
Auxiliary model	$X_i = c + \gamma \cdot S_i + \eta_i$	

By (naively) assuming  $\text{Cov}(S_i, E_i) = 0$  in our model, the population value of slope is

$$\begin{aligned} b &= \frac{\text{Cov}(S_i, \log Y_i)}{\text{Var}(S_i)} = \frac{\text{Cov}(S_i, \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i)}{\text{Var}(S_i)} \\ &= \frac{\text{Cov}(S_i, \alpha) + \text{Cov}(S_i, \beta S_i) + \text{Cov}(S_i, \delta X_i) + \text{Cov}(S_i, U_i)}{\text{Var}(S_i)} \\ &= \frac{\beta \cdot \text{Var}(S_i) + \delta \cdot \text{Cov}(S_i, X_i)}{\text{Var}(S_i)} \end{aligned}$$

# Omitted variable bias (OVB)

“True” model

$$\log Y_i = \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i$$

$$\text{Cov}(S_i, U_i) = 0$$

Our model

$$\log Y_i = a + b \cdot S_i + E_i$$

Auxiliary model

$$X_i = c + \gamma \cdot S_i + \eta_i$$

By (naively) assuming  $\text{Cov}(S_i, E_i) = 0$  in our model, the population value of slope is

$$\begin{aligned} b &= \frac{\text{Cov}(S_i, \log Y_i)}{\text{Var}(S_i)} = \frac{\text{Cov}(S_i, \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i)}{\text{Var}(S_i)} \\ &= \frac{\text{Cov}(S_i, \alpha) + \text{Cov}(S_i, \beta S_i) + \text{Cov}(S_i, \delta X_i) + \text{Cov}(S_i, U_i)}{\text{Var}(S_i)} \\ &= \frac{\beta \cdot \text{Var}(S_i) + \delta \cdot \text{Cov}(S_i, X_i)}{\text{Var}(S_i)} = \beta + \delta \cdot \frac{\text{Cov}(S_i, X_i)}{\text{Var}(S_i)} \end{aligned}$$

# Omitted variable bias (OVB)

“True” model	$\log Y_i = \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i$	$\text{Cov}(S_i, U_i) = 0$
Our model	$\log Y_i = a + b \cdot S_i + E_i$	
Auxiliary model	$X_i = c + \gamma \cdot S_i + \eta_i$	

By (naively) assuming  $\text{Cov}(S_i, E_i) = 0$  in our model, the population value of slope is

$$\begin{aligned} b &= \frac{\text{Cov}(S_i, \log Y_i)}{\text{Var}(S_i)} = \frac{\text{Cov}(S_i, \alpha + \beta \cdot S_i + \delta \cdot X_i + U_i)}{\text{Var}(S_i)} \\ &= \frac{\text{Cov}(S_i, \alpha) + \text{Cov}(S_i, \beta S_i) + \text{Cov}(S_i, \delta X_i) + \text{Cov}(S_i, U_i)}{\text{Var}(S_i)} \\ &= \frac{\beta \cdot \text{Var}(S_i) + \delta \cdot \text{Cov}(S_i, X_i)}{\text{Var}(S_i)} = \beta + \delta \cdot \frac{\text{Cov}(S_i, X_i)}{\text{Var}(S_i)} = \beta + \delta \cdot \gamma \end{aligned}$$

# Assumptions

- MLR1 (linear outcome model)
- MLR2 (random sampling)
- MLR3 (no collinearity)

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$$

$\{Y_i, X_{i1}, \dots, X_{ik}\}_{i=1}^N$  is random draw  
no  $X_{ij}$  linear function of any other  $X_{il}$

# Assumptions

- MLR1 (linear outcome model)
- MLR2 (random sampling)
- MLR3 (no collinearity)
- MLR4 (independence)

$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$   
 $\{Y_i, X_{i1}, \dots, X_{ik}\}_{i=1}^N$  is random draw  
no  $X_{ij}$  linear function of any other  $X_{il}$

$$E[U_i | X_{i1}, \dots, X_{ik}] = 0$$



# Assumptions

- MLR1 (linear outcome model)
- MLR2 (random sampling)
- MLR3 (no collinearity)
- MLR4 (independence)
- MLR5 (homoskedasticity)
- MLR6 (normality)

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$$

$\{Y_i, X_{i1}, \dots, X_{ik}\}_{i=1}^N$  is random draw  
no  $X_{ij}$  linear function of any other  $X_{il}$

$$E[U_i | X_{i1}, \dots, X_{ik}] = 0$$

$$\text{Var}(U_i | X_{i1}, \dots, X_{ik}) = \sigma^2$$

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

$$\Rightarrow Y_i \sim N(\beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik}, \sigma^2)$$

# Ordinary Least Squares (OLS) Estimator

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$$

$$\min_{\{\beta_0, \beta_1, \dots, \beta_k\}} \frac{1}{N} \sum_{i=1}^N (Y_i - \beta_0 - \beta_1 X_{i1} - \cdots - \beta_k X_{ik})^2$$

$$\Rightarrow \hat{\beta}_j^{OLS}$$

# Ordinary Least Squares (OLS) Estimator

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$$

$$\min_{\{\beta_0, \beta_1, \dots, \beta_k\}} \frac{1}{N} \sum_{i=1}^N (Y_i - \beta_0 - \beta_1 X_{i1} - \cdots - \beta_k X_{ik})^2$$

$$\Rightarrow \hat{\beta}_j^{OLS} = \frac{\widehat{\text{Cov}}(\tilde{X}_{ij}, Y_i)}{\widehat{\text{Var}}(\tilde{X}_{ij})} \quad \forall j = \{0, 1, \dots, k\}$$

$$= \frac{\widehat{\text{Cov}}(X_{ij} - \hat{\theta}_1 X_{i1} - \cdots - \hat{\theta}_k X_{ik}, Y_i)}{\widehat{\text{Var}}(X_{ij} - \hat{\theta}_1 X_{i1} - \cdots - \hat{\theta}_k X_{ik})}$$

# OLS Results

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$$

• T1 (unbiased)      MLR1+2+3+4       $\Rightarrow E[\hat{\beta}_j^{OLS}] = \beta_j \quad \forall j = \{0, 1, \dots, k\}$

# OLS Results

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$$

- T1 (unbiased)      MLR1+2+3+4       $\Rightarrow E[\hat{\beta}_j^{OLS}] = \beta_j \quad \forall j = \{0, 1, \dots, k\}$
- T2 (efficient)  
(Gauss-Markov)      MLR1+2+3+4+5       $\Rightarrow E[\hat{\beta}_j^{OLS}] = \beta_j \quad \forall j = \{0, 1, \dots, k\}$   
 $\text{Var}[\hat{\beta}_j^{OLS}] \leq \text{Var}[\hat{\beta}_j^{\text{other linear}}]$

# OLS Results

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$$

- T1 (unbiased)      MLR1+2+3+4       $\Rightarrow E[\hat{\beta}_j^{OLS}] = \beta_j \quad \forall j = \{0, 1, \dots, k\}$
- T2 (efficient)  
(Gauss-Markov)      MLR1+2+3+4+5       $\Rightarrow E[\hat{\beta}_j^{OLS}] = \beta_j \quad \forall j = \{0, 1, \dots, k\}$   
$$\text{Var}[\hat{\beta}_j^{OLS}] \leq \text{Var}[\hat{\beta}_j^{\text{other linear}}]$$
$$\text{Var}[\hat{\beta}_j^{OLS}] = \frac{\sigma^2}{\text{Var}(X_{ij}) \cdot [1 - R_{\text{reg } X_j \text{ on all } X_k}^2]}$$
$$E[\hat{\sigma}^2] = E\left[\frac{1}{N-k-1} \sum_{i=1}^N \hat{U}_i^2\right] = \sigma^2$$
$$\text{se}[\hat{\beta}_j^{OLS}] := \sqrt{\widehat{\text{Var}}[\hat{\beta}_j^{OLS}]}$$

# OLS Results

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$$

- T3 (efficient) MLR1+2+3+4+5+6  $\Rightarrow \hat{\beta}_j^{OLS} \sim N(\beta_j, \text{Var}[\beta_j]) \quad \forall j = \{0, 1, \dots, k\}$   
(Classical)

# OLS Results

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$$

- T3 (efficient) MLR1+2+3+4+5+6  $\Rightarrow \hat{\beta}_j^{OLS} \sim N(\beta_j, \text{Var}[\beta_j]) \quad \forall j = \{0, 1, \dots, k\}$   
(Classical)

$$\frac{\hat{\beta}_j^{OLS} - \beta_j}{\text{sd}[\beta_j]} \sim N(0, 1)$$

$$\frac{\hat{\beta}_j^{OLS} - \beta_j}{\text{se}[\beta_j]} \sim t(N - k - 1)$$

$$\text{se}[\hat{\beta}_j^{OLS}] := \sqrt{\widehat{\text{Var}}[\hat{\beta}_j^{OLS}]}$$



t and F tests

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$$

- Individual hypothesis test about slope parameter (t-test)

$$H_0: \beta_j = 0$$

$$t_{\widehat{\beta}_j^{OLS}} := \frac{\widehat{\beta}_j^{OLS} - 0}{\text{se}[\widehat{\beta}_j^{OLS}]} \sim t(N - k - 1)$$

t and F tests

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_k X_{ik} + U_i$$

- Individual hypothesis test about slope parameter (t-test)

$$H_0: \beta_j = 0$$

$$t_{\hat{\beta}_j^{OLS}} := \frac{\hat{\beta}_j^{OLS} - 0}{\text{se}[\hat{\beta}_j^{OLS}]} \sim t(N - k - 1)$$

- Joint hypothesis test about entire linear model

$$SSR_U := \sum_{i=1}^N \hat{U}_i^2 = \sum_{i=1}^N (Y_i - \hat{Y}_i)^2$$

$$SSR_R := \sum_{i=1}^N (Y_i - \hat{\beta}_0)^2$$

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_k = 0 \quad F := \frac{\frac{SSR_R - SSR_U}{k}}{\frac{SSR_U}{N - k - 1}} \sim F(k, N - k - 1)$$

# Confidence Intervals

$$P \left[ \hat{\beta}_j^{OLS} - c_{\alpha} \cdot \text{se} \left[ \hat{\beta}_j^{OLS} \right] \leq \beta_j \leq \hat{\beta}_j^{OLS} + c_{\alpha} \cdot \text{se} \left[ \hat{\beta}_j^{OLS} \right] \right] = 1 - \alpha$$

significance level (rate we tolerate Type 1 errors)

$$\alpha \in \{0.01, 0.05, 0.1\}$$

critical value associated w/  $\alpha$  in distribution

$$c_{\alpha} \approx 1.96 \text{ if } 5\%$$

estimated standard error

$$\text{se} \left[ \hat{\beta}_j^{OLS} \right] := \sqrt{\widehat{\text{Var}} \left[ \hat{\beta}_j^{OLS} \right]}$$

interpretation = ???

# Confidence Intervals

$$P \left[ \hat{\beta}_j^{OLS} - c_{\alpha} \cdot \text{se} \left[ \hat{\beta}_j^{OLS} \right] \leq \beta_j \leq \hat{\beta}_j^{OLS} + c_{\alpha} \cdot \text{se} \left[ \hat{\beta}_j^{OLS} \right] \right] = 1 - \alpha$$

significance level (rate we tolerate Type 1 errors)

$$\alpha \in \{0.01, 0.05, 0.1\}$$

critical value associated w/  $\alpha$  in distribution

$$c_{\alpha} \approx 1.96 \text{ if } 5\%$$

estimated standard error

$$\text{se} \left[ \hat{\beta}_j^{OLS} \right] := \sqrt{\widehat{\text{Var}} \left[ \hat{\beta}_j^{OLS} \right]}$$

interpretation = this procedure to estimate bounds will cover true  $\beta_j$  parameter  
95% of the time (over many hypothetical repeated samples)