

POLS 6481. Research Design and Quantitative Methods II

Lecture 2. Simple Regression

Readings: Wooldridge, *Introductory Econometrics 5e*, 2.1 + 2.2 + 2.3a

Outline:

Prelude – PED Use in Baseball

1. Review of the Linear Simple Regression Model
2. Summation (Review Appendix A.1, A.2)
3. Definitions (Review Appendix B.3, B.4)

Interlude – PED Use in Baseball

4. Assumptions of the Linear Simple Regression Model

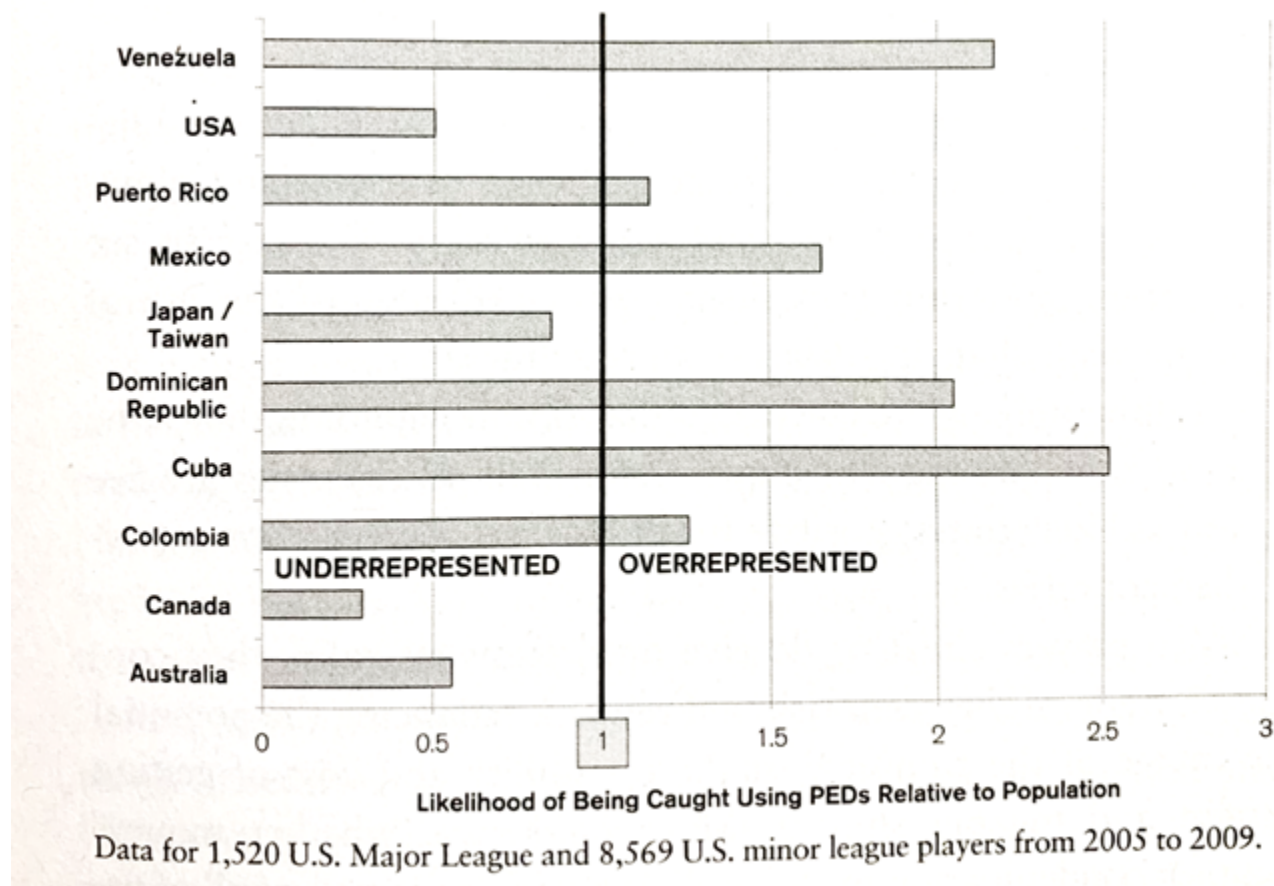
Postlude – PED Use in Baseball

POLS 6481. Research Design and Quantitative Methods II

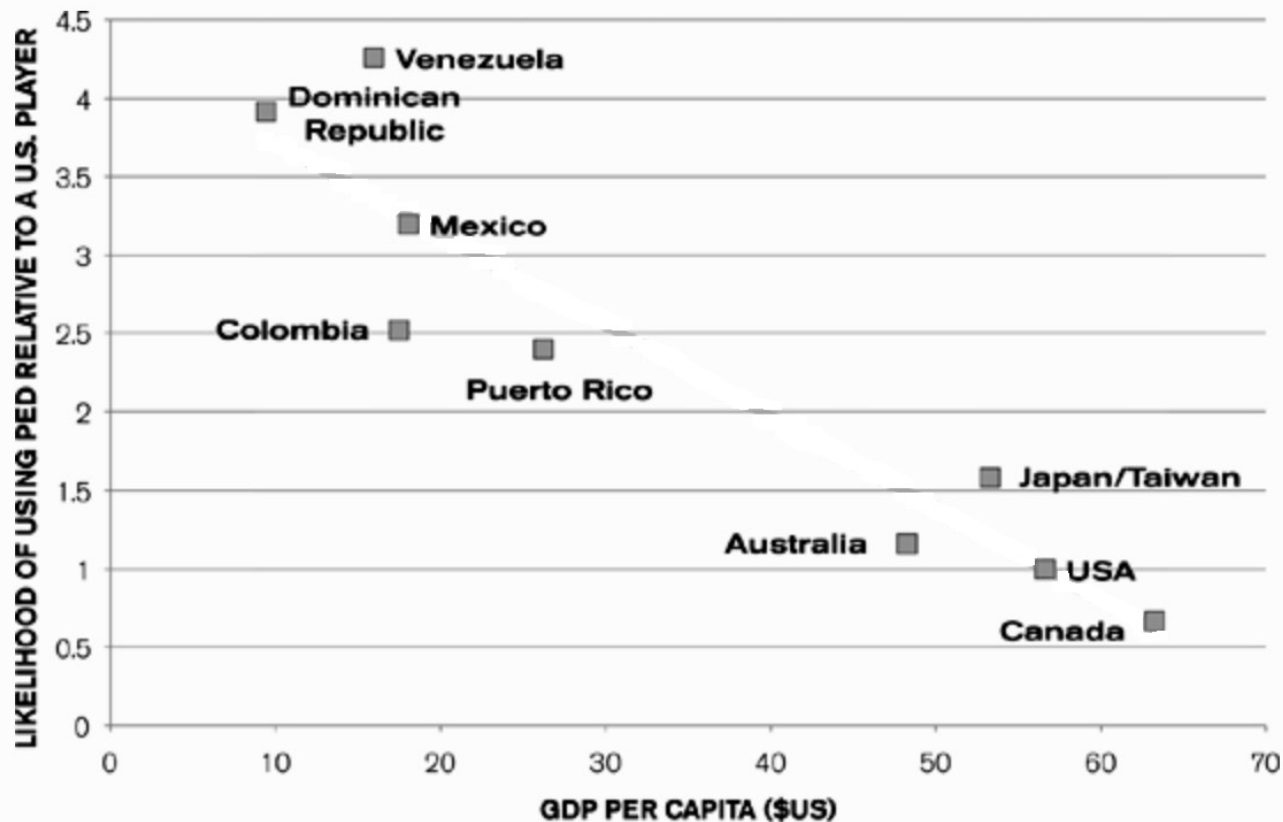
Lecture 2. Simple Regression

Readings: Wooldridge, *Introductory Econometrics 5e*, 2.1 + 2.2 + 2.3a

Prelude – PED Use in Baseball



Prelude – PED Use in Baseball



POLS 6481. Research Design and Quantitative Methods II

Lecture 2. Simple Regression

Readings: Wooldridge, *Introductory Econometrics 5e*, 2.1 + 2.2 + 2.3a

1. Review of the Linear Simple Regression Model

“explained” variable	y	a.k.a. dependent variable, <i>regressand</i>	
“explanatory” variable	x	a.k.a. independent variable, <i>regressor</i>	
intercept parameter	β_0	$E(y x = 0)$	
slope parameter	β_1	$\frac{\Delta E(y x)}{\Delta x}$	
Population regression function	$E(y x) = \beta_0 + \beta_1 x$	<i>systematic component</i>	[2.8]
disturbance	u	<i>stochastic component</i>	
Simple linear regression model	$y = \beta_0 + \beta_1 x + u$		

2. Summation (Review Appendix A.1)

For the following, assume summation is over $i = 1, \dots, n$

$$[A.1] \quad \sum x_i = x_1 + x_2 + \dots$$

$$[A.5] \quad \frac{\sum x_i}{n} = \bar{x}$$

$$[A.6] \quad \sum (x_i - \bar{x}) = 0$$

$$[A.2] \quad \sum c = n \cdot c$$

$$[A.3] \quad \sum c \cdot x_i = c \cdot \sum x_i$$

$$[A.7] \quad \sum x_i \cdot (x_i - \bar{x}) = \sum (x_i - \bar{x})^2 = \sum (x_i)^2 - \frac{(\sum x_i)^2}{n}$$

$$\sum x_i \cdot (y_i - \bar{y}) = \sum (x_i - \bar{x}) \cdot (y_i - \bar{y})$$

3. Definitions (Review Appendix B.3, B.4)

For the following, assume summation is over $i = 1, \dots, n$

Mean $E(x) = \frac{1}{n} \sum x_i$ $E(y) = \frac{1}{n} \sum y_i$

Variance $E[(x - \mu_x)^2] = \frac{1}{n} \sum (x_i - \bar{x})^2$ $E[(y - \mu_y)^2] = \frac{1}{n} \sum (y_i - \bar{y})^2$ [B.23]

Covariance $E[(x - \mu_x) \cdot (y - \mu_y)] = \frac{1}{n} \sum (x_i - \bar{x}) \cdot (y_i - \bar{y})$ [B.26]

Correlation $r_{xy} = \frac{cov(x,y)}{\sqrt{var(x) \cdot var(y)}} = \frac{\sum (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2}}$ [B.29]

Regression slope $\beta_1 = \frac{cov(x,y)}{var(x)} = \frac{\sum (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$ [2.19]

Regression intercept $\beta_0 = \bar{y} - \beta_1 \bar{x}$ [2.17]

3. Definitions, continued

Sample regression function $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 \cdot x$ [2.23]

Fitted values $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 \cdot x_i$ for $i = 1, \dots, n$ [2.20]

Residual
$$\begin{aligned}\hat{u}_i &= y_i - \hat{y}_i \\ &= y_i - \hat{\beta}_0 - \hat{\beta}_1 \cdot x_i\end{aligned}$$
 [2.21]

Sum of squared residuals $\sum \hat{u}_i^2 = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 \cdot x_i)^2$ [2.22]

4. Assumptions of the Linear Simple Regression Model

Assumption 1. Zero mean residual

$$E(\hat{u}) = 0$$

[2.10]

[[used to prove $\hat{\beta}_0 = \bar{y} - \beta_1 \bar{x}$; see whiteboard, p. 28]]

① Assume $E(\hat{u}) = 0$; prove $\hat{\beta}_0 = \bar{y} - \beta_1 \bar{x}$

define $\hat{u}_i = y_i - \bar{y} = y_i - \hat{\beta}_0 - \beta_1 x_i$

then $\frac{1}{n} \sum \hat{u}_i = 0 = \frac{1}{n} \sum (y_i - \hat{\beta}_0 - \beta_1 x_i)$

$$0 = \frac{1}{n} \sum y_i - \frac{1}{n} \sum \hat{\beta}_0 - \frac{1}{n} \sum \beta_1 x_i$$

$$0 = \bar{y} - \frac{1}{n} \cdot n \hat{\beta}_0 - \beta_1 \bar{x}$$

\uparrow def. of \bar{y} \uparrow rule that $\sum_{i=1}^n c = nc$ \uparrow def. of \bar{x} and rule that $\sum c x_i = c \sum x_i$

$$0 = \bar{y} - \hat{\beta}_0 - \beta_1 \bar{x}$$

$$\hat{\beta}_0 = \bar{y} - \beta_1 \bar{x}$$

4. Assumptions of the Linear Simple Regression Model

Assumption 2. Zero correlation of residual and regressor $E(\hat{u} \cdot x) = 0$ [2.11]

[[used to prove $\hat{\beta}_1 = \frac{\text{cov}(x,y)}{\text{var}(x)}$; see whiteboard, p. 28-29]]

② Assume $E(\hat{u} \cdot x) = 0$, prove that $\beta_1^* = \frac{\text{cov}(x,y)}{\text{var}(x)}$ =

$$\frac{1}{n} \sum u_i x_i = 0 = \frac{1}{n} \sum x_i (y_i - \beta_0^* - \beta_1^* x_i)$$

$$\begin{aligned} \text{substitute } \bar{y} - \beta_1^* \bar{x} \text{ for } \beta_0^*: & \frac{1}{n} \sum x_i (y_i - \bar{y} + \beta_1^* \bar{x} - \beta_1^* x_i) \\ &= \frac{1}{n} \sum x_i ((y_i - \bar{y}) - \beta_1^* (x_i - \bar{x})) \\ &= \frac{1}{n} \sum x_i (y_i - \bar{y}) - \beta_1^* \frac{1}{n} \sum x_i (x_i - \bar{x}) \\ &= \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}) - \beta_1^* \cdot \frac{1}{n} \sum (x_i - \bar{x})^2 \\ 0 &= \text{cov}(x, y) - \beta_1^* \text{var}(x) \end{aligned}$$

$$\beta_1^* \text{var}(x) = \text{cov}(x, y)$$

$$\beta_1^* = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

Interlude – PED Use in Baseball

Scorecasting, by T.J. Moskowitz & L.J. Wertheim (2011)

Dataset: *PEDs in baseball.csv*

Sample: National aggregate data ($n = 9$ ish) from 2005 to 2009

Dependent variable: Relative likelihood of PED use (1 = American ballplayer)

Independent variable: Gross Domestic Product per capita (US \$)

```
x <- c(48, 63, 17, 9, 53, 18, 26, 56, 16)
```

```
y <- c(1.2, 0.7, 2.5, 3.9, 1.6, 3.2, 2.4, 1.0, 4.2)
```

```
origin <- c("Australia", "Canada", "Colombia", "Dominican", "JapanTaiwan", "Mexico",  
           "PuertoRico", "USA", "Venezuela")
```

```
plot(x, y, pch=16, xlim = c(0,70), ylim = c(0,4.5))
```

```
text(x, y, labels=origin, pos=1, xpd=TRUE)
```

POLS 6481. Research Design and Quantitative Methods II

Lecture 2. Simple Regression

Readings: Wooldridge, *Introductory Econometrics 5e*, 2.1 + 2.2 + 2.3a

Interlude – PED Use in Baseball

i	x_i	y_i	$(x_i - \bar{x})$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
Australia	48	1.2					
Canada	63	0.7					
Colombia	17	2.5					
Dominican	9	3.9					
JapanTaiwan	53	1.6					
Mexico	18	3.2					
PuertoRico	26	2.4					
USA	56	1.0					
Venezuela	16	4.2					
Σ				$SST_x =$		$SST_y =$	
Σ/n	$\bar{x} =$	$\bar{y} =$		$var(x) =$		$var(y) =$	$cov(x,y) =$

Slope: $\widehat{\beta}_1 = \frac{cov(x,y)}{var(x)} = \text{————} =$

Intercept: $\widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \bar{x} =$

POLS 6481. Research Design and Quantitative Methods II

Lecture 2. Simple Regression

Readings: Wooldridge, *Introductory Econometrics 5e*, 2.1 + 2.2 + 2.3a

Postlude – PED Use in Baseball

i	x_i	y_i	$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$	$\hat{u}_i = (y_i - \hat{y}_i)$	$\hat{u}_i \cdot x_i$	$\hat{u}_i^2 = (y_i - \hat{y}_i)^2$
Australia	48	1.2				
Canada	63	0.7				
Colombia	17	2.5				
Dominican	9	3.9				
JapanTaiwan	53	1.6				
Mexico	18	3.2				
PuertoRico	26	2.4				
USA	56	1.0				
Venezuela	16	4.2				
Σ			$\Sigma \hat{y}_i =$	$\Sigma \hat{u}_i =$	$\Sigma \hat{u}_i \cdot x_i =$	$SSR =$
Σ/n	$\bar{x} =$	$\bar{y} =$	$\hat{\bar{y}} =$			

sigma: $\sqrt{\hat{\sigma}^2} = \sqrt{\frac{SSR}{n-2}} =$

$se(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{\Sigma(x_i - \bar{x})^2}} = -- =$

$se(\hat{\beta}_0) = \frac{\hat{\sigma}}{\sqrt{\Sigma(x_i - \bar{x})^2}} \cdot \sqrt{\frac{\Sigma x_i^2}{n}} =$