

# Precept 6: Regression

Soc 500: Applied Social Statistics

Alex Kindel

Princeton University

October 18, 2018

# Today's plan<sup>1</sup>

- ① Some homework items
  - PS4: Properties of variance
  - PS5: OLS assumptions
- ② Dummy variables
- ③ Practice
  - Interpreting regression output
  - Partialing out
  - Interactions
  - Multicollinearity
- ④ If time...
  - Measurement error
  - `lm()` internals

---

<sup>1</sup>Thanks to Ian Lundberg and Brandon Stewart for some of today's examples and slides.

## Properties of variance

Let  $X_j \sim ?(\mu, \sigma^2)$ .

$$\begin{aligned} \text{Var}[aX_j - bX_k] &= a^2 \text{Var}[X_j] + (-b)^2 \text{Var}[X_k] \\ &= a^2 \sigma^2 + b^2 \sigma^2 \end{aligned} \tag{1}$$

# OLS assumptions

- ① **Linearity in parameters:** We likely have omitted variable bias.

# OLS assumptions

- ① **Linearity in parameters:** We likely have omitted variable bias.
- ② **Homoskedasticity:** The spread of the residuals increases in  $X$ .

# OLS assumptions

- ① **Linearity in parameters:** We likely have omitted variable bias.
- ② **Homoskedasticity:** The spread of the residuals increases in  $X$ .
- ③ **Zero conditional mean:** Residuals are more negative for high values of  $X$ .

# OLS assumptions

- ① **Linearity in parameters:** We likely have omitted variable bias.
- ② **Homoskedasticity:** The spread of the residuals increases in  $X$ .
- ③ **Zero conditional mean:** Residuals are more negative for high values of  $X$ .
- ④ **Random sampling:** FFCWS has a complex sample design; the unweighted observations are not quite representative (what's the population?).

## Dummy variables



Today's example:

Today's example:

Today's example: **Pokémon again!**

Just kidding.

## Today's example: **LaLonde (1986)**

Famous<sup>2</sup> econometric analysis of the effects of a job training program

- re78 is earnings in 1978
- age is age
- educ is education, in years

Let's go through some examples in the R file.

---

<sup>2</sup>At this point in your education, you may have realized that academics should never be trusted when they refer to something as well-known.



## Attenuation bias: When $X$ has random noise

What happens when  $X$  is measured with error?

$$\begin{aligned}\hat{\beta}_1 &= \frac{\text{Cov}(\tilde{X}, Y)}{\text{Var}(\tilde{X})} \\ &= \frac{\text{Cov}(X + u, \beta X + \epsilon)}{\text{Var}(X + u)} \\ &= \frac{\beta \text{Cov}(X, X) + \text{Cov}(X, \epsilon) + \text{Cov}(u, X) + \text{Cov}(u, \epsilon)}{\text{Var}(X) + \text{Var}(u) + 2\text{Cov}(X, u)} \\ &= \beta \frac{\text{Var}(X) + 0 + 0 + 0}{\text{Var}(X) + \text{Var}(u) + 0} \\ &= \beta \frac{\sigma_x^2}{\sigma_x^2 + \sigma_u^2} = \beta \frac{\sigma_x^2}{\sigma_{\tilde{x}}^2}\end{aligned}$$

$\hat{\beta}$  will thus be biased toward 0. We call this attenuation.

## No bias when $Y$ has random noise

What happens when  $Y$  is measured with error? No bias.

$$\begin{aligned}\hat{\beta}_1 &= \frac{\text{Cov}(X, \tilde{Y})}{\text{Var}(X)} \\ &= \frac{\text{Cov}(X, \beta X + u + \epsilon)}{\text{Var}(X)} \\ &= \frac{\beta \text{Cov}(X, X) + \text{Cov}(X, u) + \text{Cov}(X, \epsilon)}{\text{Var}(X)} \\ &= \beta \frac{\text{Var}(X) + 0 + 0}{\text{Var}(X)} \\ &= \beta\end{aligned}$$

Bigger standard error:

$$\widehat{SE}(\hat{\beta}_1) = \frac{\sigma^2}{\text{Var}(X)}$$