

# 4.5 SAMPLING DISTRIBUTION OF THE OLS ESTIMATOR

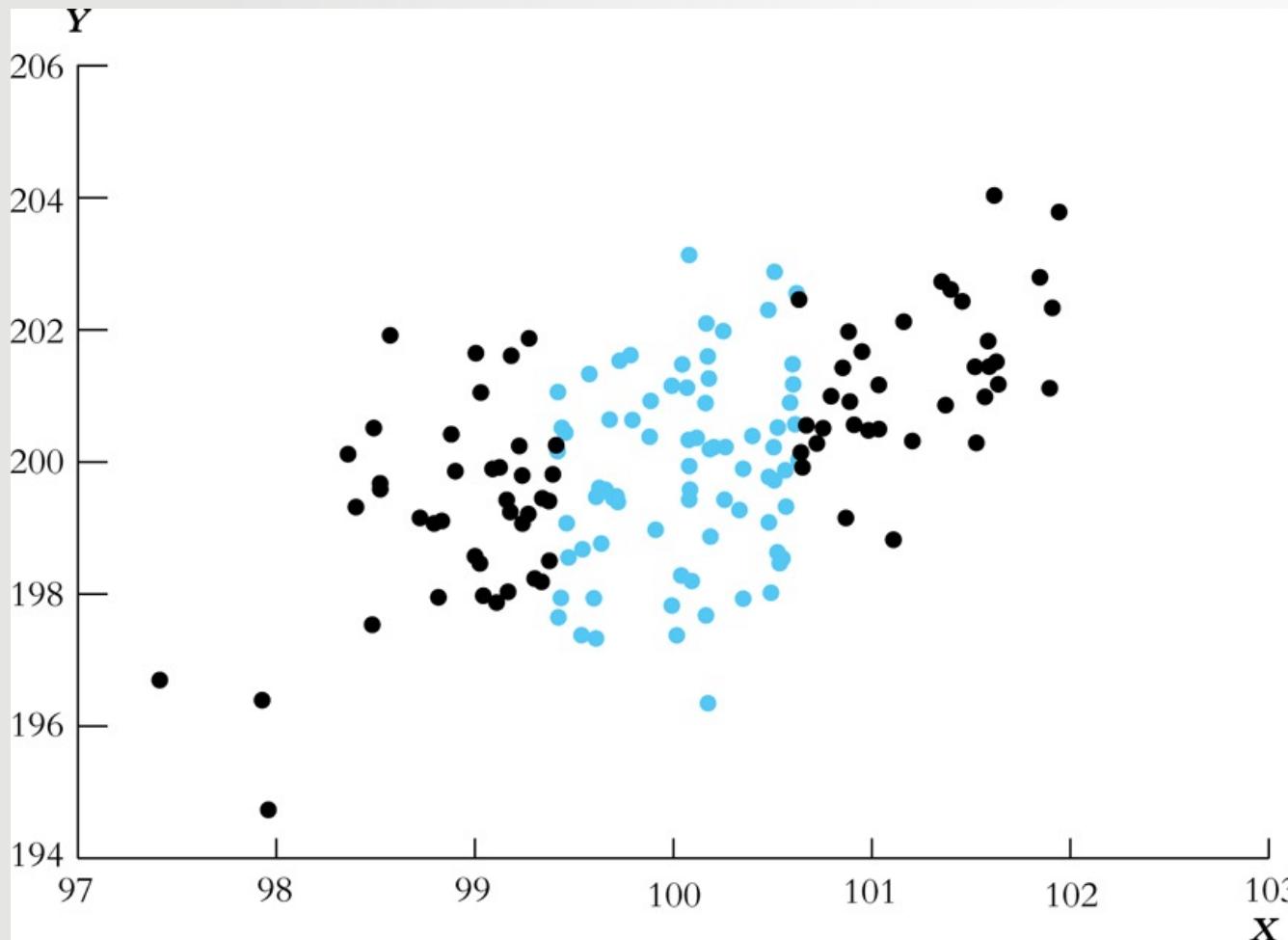
# SAMPLING DISTRIBUTION OF $\hat{\beta}_0$ AND $\hat{\beta}_1$

- $\hat{\beta}_0$  and  $\hat{\beta}_1$  = random variables (*why?*)
- $E(\hat{\beta}_0) = \beta_0$  and  $E(\hat{\beta}_1) = \beta_1$
- $\hat{\beta}_0$  and  $\hat{\beta}_1$  tend to be normally distributed in large samples.
- $\hat{\beta}_0 \sim N(\beta_0, \sigma_{\hat{\beta}_0}^2)$  and  $\hat{\beta}_1 \sim N(\beta_1, \sigma_{\hat{\beta}_1}^2)$
- What determines  $\sigma_{\hat{\beta}_0}^2$  and  $\sigma_{\hat{\beta}_1}^2$ ?

# THE VARIANCE OF THE OLS ESTIMATOR

$$\text{var}(\hat{\beta}_1) = \frac{1}{n} \times \frac{\text{var}[(X_i - \mu_x)u_i]}{(\sigma_X^2)^2}.$$

# THE VARIANCE OF THE OLS ESTIMATOR



- The number of black and blue dots is the same and they come from the same joint distribution.
- Using which would you get a more accurate regression line?
- Increasing the spread of  $X$  decreases  $\text{var}(\hat{\beta}_1)$

# 4.6 HYPOTHESIS TESTS ABOUT REGRESSION COEFFICIENTS

# HYPOTHESIS TESTS

- Null and alternative hypotheses:

$$H_0: \beta_1 = \beta_{1,0} \text{ vs. } H_1: \beta_1 \neq \beta_{1,0}$$

- Three steps for testing  $H_0$  :
  1. Compute  $\hat{\beta}_1$  and  $SE(\hat{\beta}_1)$  using sample data.
  2. Compute the t-statistics
  3. Compute the p-value

# HYPOTHESIS TESTS

- Null and alternative hypotheses:

$$H_0: \beta_1 = \beta_{1,0} \text{ vs. } H_1: \beta_1 \neq \beta_{1,0}$$

- Three steps for testing  $H_0$  :
  1. Compute  $\hat{\beta}_1$  and  $SE(\hat{\beta}_1)$  using sample data.
  2. Compute the t-statistics
  3. Compute the p-value

# THE STANDARD ERROR OF $\hat{\beta}_1$

- $SE(\hat{\beta}_1)$  is an estimator of  $\sigma_{\hat{\beta}_1}$ .
- $SE(\hat{\beta}_1) = \sqrt{\hat{\sigma}_{\hat{\beta}_1}^2} = \frac{1}{n} \times \frac{\frac{1}{n-2} \sum_{i=1}^n (X_i - \bar{X})^2 \hat{u}_i^2}{\left[ \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right]^2}$
- (complicated, but STATA will do it for you)
- Also called *robust* standard error.

```
regress testscr str, robust
```

Regression with robust standard errors

Number of obs = 420  
F( 1, 418) = 19.26  
Prob > F = 0.0000  
R-squared = 0.0512  
Root MSE = 18.581

testscr	Robust					
	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
str	-2.279808	.5194892	-4.39	0.000	-3.300945	-1.258671
_cons	698.933	10.36436	67.44	0.000	678.5602	719.3057

- $\hat{\beta}_1 = -2.28$  and  $SE(\hat{\beta}_1) = 0.52$
- $\hat{\beta}_0 = 698.9$  and  $SE(\hat{\beta}_1) = 10.36$

# HYPOTHESIS TESTS

- Null and alternative hypotheses:

$$H_0: \beta_1 = \beta_{1,0} \text{ vs. } H_1: \beta_1 \neq \beta_{1,0}$$

- Three steps for testing  $H_0$  :
  1. Compute  $\hat{\beta}_1$  and  $SE(\hat{\beta}_1)$  using sample data.
  2. Compute the t-statistics
  3. Compute the p-value

# T-STATISTICS FOR OLS PARAMETERS

$$t = \frac{\text{estimator} - \text{hypothesized value}}{\text{standard error of estimator}}$$

$$t = \frac{\hat{\beta}_1 - \hat{\beta}_{1,0}}{SE(\hat{\beta}_1)}$$

- t has a *standard normal distribution* in large samples
- $t \sim N(0,1)$

```
regress testscr str, robust
```

Regression with robust standard errors

Number of obs = 420  
F( 1, 418) = 19.26  
Prob > F = 0.0000  
R-squared = 0.0512  
Root MSE = 18.581

testscr	Robust					
	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
str	-2.279808	.5194892	-4.39	0.000	-3.300945	-1.258671
_cons	698.933	10.36436	67.44	0.000	678.5602	719.3057

- $\hat{\beta}_1 = -2.28$  and  $SE(\hat{\beta}_1) = 0.52$  and  $t = -4.39$
- $\hat{\beta}_0 = 698.9$  and  $SE(\hat{\beta}_1) = 10.36$  and  $t = 67.44$

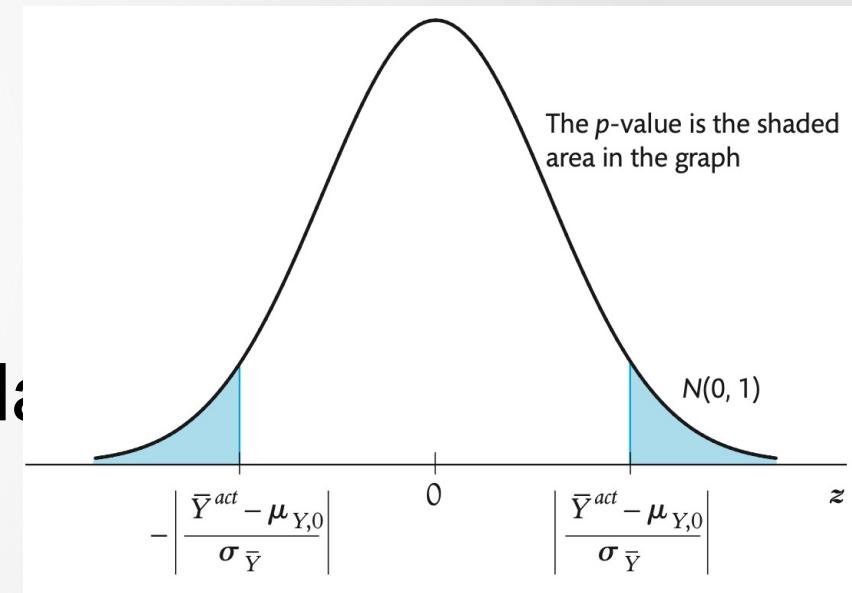
# HYPOTHESIS TESTS

- Null and alternative hypotheses:

$$H_0: \beta_1 = \beta_{1,0} \text{ vs. } H_1: \beta_1 \neq$$

- Three steps for testing  $H_0$ :

1. Compute  $\hat{\beta}_1$  and  $SE(\hat{\beta}_1)$  using sample data
2. Compute the t-statistics
3. Compute the p-value



# COMPUTING THE P-VALUE

- p-value  $= Pr_{H_0} [|\hat{\beta}_1 - \beta_{1,0}| > |\hat{\beta}_1^{act} - \hat{\beta}_{1,0}|]$  $= Pr_{H_0} \left[ \left| \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)} \right| > \left| \frac{\hat{\beta}_1^{act} - \hat{\beta}_{1,0}}{SE(\hat{\beta}_1)} \right| \right]$  $= Pr_{H_0} (|t| > |t^{act}|)$  $= Pr_{H_0} (|Z| > |t^{act}|)$  $= 2\phi(-|t^{act}|)$

```
regress testscr str, robust
```

Regression with robust standard errors

Number of obs = 420  
F( 1, 418) = 19.26  
Prob > F = 0.0000  
R-squared = 0.0512  
Root MSE = 18.581

testscr	Robust					
	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
str	-2.279808	.5194892	-4.39	0.000	-3.300945	-1.258671
_cons	698.933	10.36436	67.44	0.000	678.5602	719.3057

- $\hat{\beta}_1 = -2.28$  and  $SE(\hat{\beta}_1) = 0.52$  and  $t = -4.39$  and  $p < 0.001$
- $\hat{\beta}_0 = 698.9$  and  $SE(\hat{\beta}_1) = 10.36$  and  $t = 67.44$  and  $p < 0.001$

regress testscr str, robust					
Regression with robust standard errors					
					Number of obs = 420
					F( 1, 418) = 19.26
					Prob > F = 0.0000
					R-squared = 0.0512
					Root MSE = 18.581
<hr/>					
		Robust			
testscr	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
<hr/>					
str	-2.279808	.5194892	-4.39	0.000	-3.300945 -1.258671
_cons	698.933	10.36436	67.44	0.000	678.5602 719.3057
<hr/>					

1. Compute  $SE(\hat{\beta}_1) = 0.52$

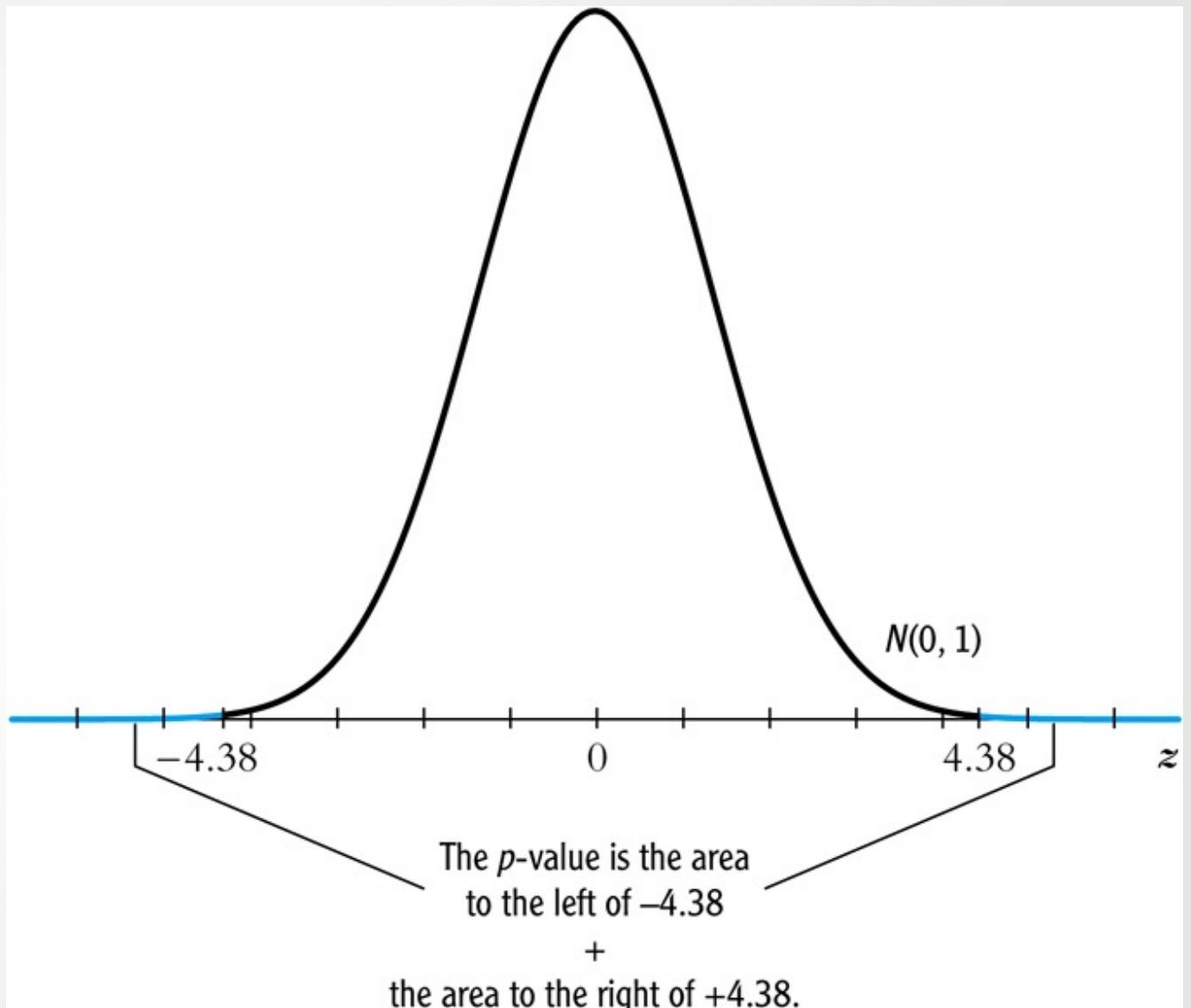
2. Compute the t-statistic:

$$t = \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)} = \frac{-2.28 - 0}{0.52} = -4.39$$

3. Compute the p-value:  $2\Phi(-|t|) = 2\Phi(-4.39) = 0.00001$

**$p$ -value =0.00001 (or  
 $10^{-5}$ )**

We can reject the null hypothesis: smaller classes do have higher test scores.



# 4.7 CONFIDENCE INTERVALS FOR REGRESSION COEFFICIENTS

# CONFIDENCE INTERVAL FOR $\beta_1$

- **95% confidence interval:** a range of values that is 95% likely to include the “true” population coefficient  $\beta_1$ .
- Includes all  $\beta_1$  values that we *cannot* reject at the 5% significance level.
- 95% confidence interval for  $\beta_1$ :

$$\hat{\beta}_1 - 1.96 * SE(\hat{\beta}_1) \leq \beta_1 \leq \hat{\beta}_1 + 1.96 * SE(\hat{\beta}_1)$$

# CONFIDENCE INTERVAL FOR $\beta_1$

- 95% confidence interval for the effect of test scores.
- We estimated  $\hat{\beta}_1 = -2.28$  and  $SE(\hat{\beta}_1) = 0.52$
- So the true  $\beta_1$  is 95% likely to be between:
  - $-2.28 - (1.96 \times 0.52) = -3.30$
  - $-2.28 + (1.96 \times 0.52) = -1.26$

```

regress testscr str, robust
Regression with robust standard errors
                                                Number of obs =      420
                                                F(  1,    418) =    19.26
                                                Prob > F        = 0.0000
                                                R-squared       = 0.0512
                                                Root MSE        = 18.581

```

	Robust					
testscr	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
str	-2.279808	.5194892	-4.39	0.000	<b>-3.300945</b>	<b>-1.258671</b>
_cons	698.933	10.36436	67.44	0.000	678.5602	719.3057

- Confidence interval for  $\beta_1$ :  $[-3.30 \leq \beta_1 \leq -1.26]$

Confidence interval for  $\beta_1$  (coefficient of STR):

$$[-3.30 \leq \beta_1 \leq -1.26]$$

**YOUR TURN:**

Can you compute a confidence interval for the average effect of a 3.5 increase in STR?

Confidence interval for  $\beta_1$  (coefficient of STR):

$$[-3.30 \leq \beta_1 \leq -1.26]$$

Lower bound:  $-3.30 * 3.5 = -11.55$

Upper bound:  $-1.26 * 3.5 = -4.41$

Increasing STR by 3 students will decrease test scores by between 4.41 and 11.55 points.

# CONFIDENCE INTERVAL FOR PREDICTED EFFECTS

- Confidence interval for the effect of a  $\Delta x$  change in X:

$$[ (\hat{\beta}_1 \text{ lower bound}) \times \Delta x ; (\hat{\beta}_1 \text{ lower bound}) \times \Delta x ]$$

$$[ (\hat{\beta}_1 - 1.96 * SE(\hat{\beta}_1)) \times \Delta x ; \hat{\beta}_1 + 1.96 * SE(\hat{\beta}_1) \times \Delta x ]$$

# 4.8 REGRESSION WHEN X IS A BINARY VARIABLE

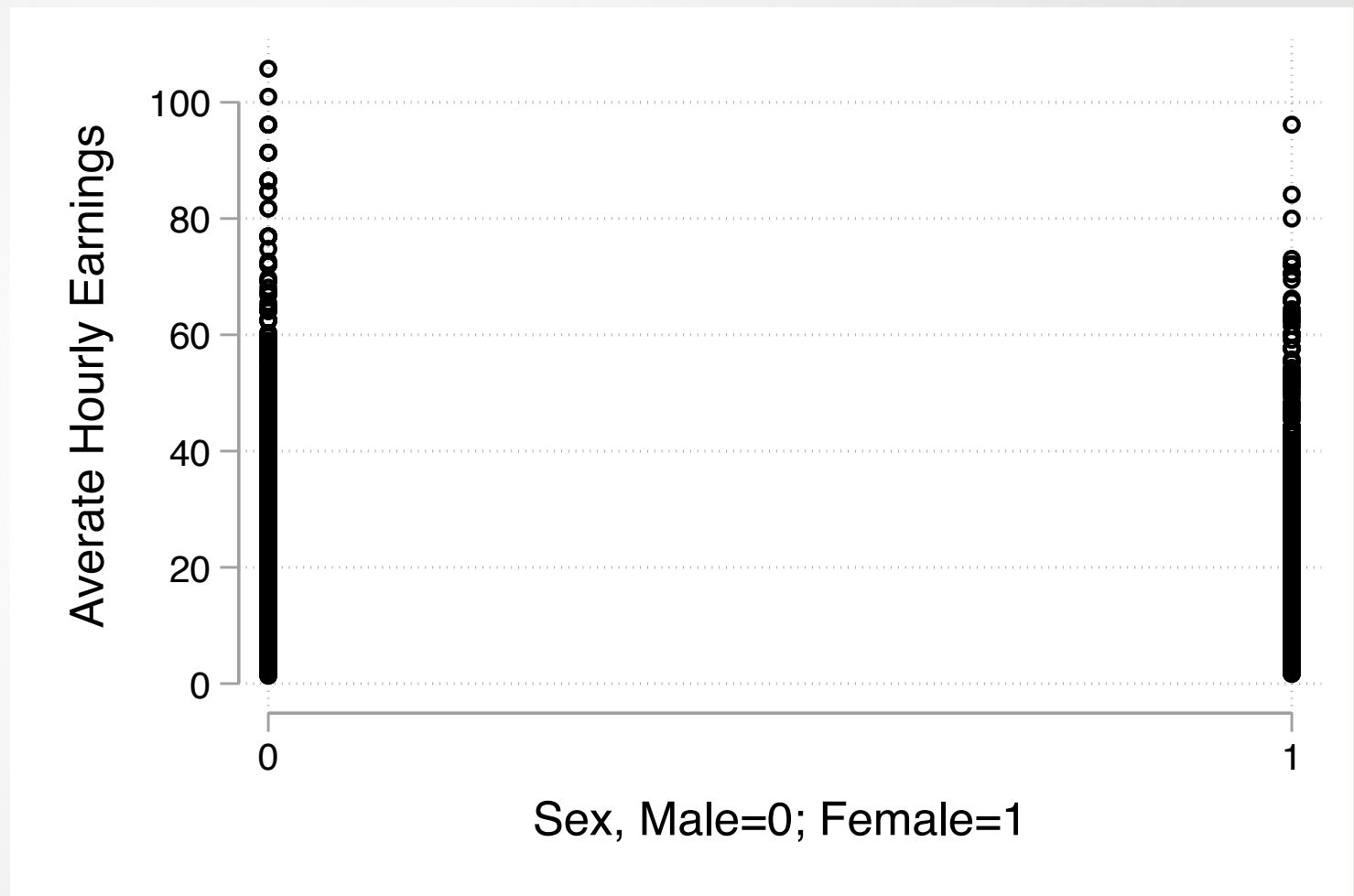
# REGRESSION WHEN X IS BINARY

- **Binary (or *indicator* or *dummy*) variables**
  - Sex at birth (1 = female; 0 = male)
  - Urban or rural (1 = urban; 0 = rural)
  - Treatment or placebo  
(1 = treatment; 0 = placebo)
  - ....



# CPS 2015 data

*scatter ahe female*



# REGRESSION WHEN X IS BINARY

$$E(Y|D) = \beta_0 + \beta_1 D$$

$$Y_i = \beta_0 + \beta_1 D_i + u_i$$

When  $D_i = 0$ :

$$E(Y|D = 0) = \beta_0 + \beta_1 \times 0 = \beta_0$$

When  $D_i = 1$ :

$$E(Y|D = 1) = \beta_0 + \beta_1 \times 1 = \beta_0 + \beta_1$$



$$\beta_1 = E(Y|D = 1) - E(Y|D = 0)$$

# EXAMPLE: GENDER GAP IN EARNINGS

```
. reg ahe female, robust
```

Linear regression

Number of obs = 13,201  
F(1, 13199) = 184.93  
Prob > F = 0.0000  
R-squared = 0.0131  
Root MSE = 10.695

ahe	Robust					
	Coefficient	std. err.	t	P> t	[95% conf. interval]	
female	-2.495648	.1835205	-13.60	0.000	-2.855375	-2.135922
_cons	18.32845	.1300679	140.91	0.000	18.0735	18.5834

- AHE for men:  
 $\beta_0 = 18.33$
- Difference between women and men:  
 $\beta_1 = -2.50$
- AHE for women:  
$$\begin{aligned}\beta_0 + \beta_1 &= \\ &= 18.32 - 2.50 \\ &= 15.83\end{aligned}$$

# EXAMPLE: GENDER GAP IN EARNINGS

```
. ttest ahe, by(female) unequal unpaired
```

Two-sample t test with unequal variances

Group	Obs	Mean	Std. err.	Std. dev.	[95% conf. interval]	
0	7,694	18.32845	.1300665	11.40884	18.07348	18.58341
1	5,507	15.8328	.1294705	9.6079	15.57899	16.08661
Combined	13,201	17.28735	.0936909	10.76467	17.1037	17.471
diff		2.495648	.1835209		2.13592	2.855377

diff = mean(0) - mean(1) t = 13.5987  
H0: diff = 0 Satterthwaite's degrees of freedom = 12855.9

- T-test for difference in means (see “review of statistics”).
- Regression does exactly the same thing!

# REGRESSION WHEN X IS BINARY: SUMMING UP

$$Y_i = \beta_0 + \beta_1 D_i + u_i$$

- $\beta_0$  = mean of Y when  $X=0$
- $\beta_0 + \beta_1$  = mean of Y when  $X=1$
- $\beta_1$  = difference in group means:  $E(Y|X = 1) - E(Y|X = 0)$
- T-stats, p-value, confidence intervals calculated as usual.
- Will give the same result as a t-test for difference in means.

# 4.9 HETEROSKEDASTICITY AND HOMOSKEDASTICITY

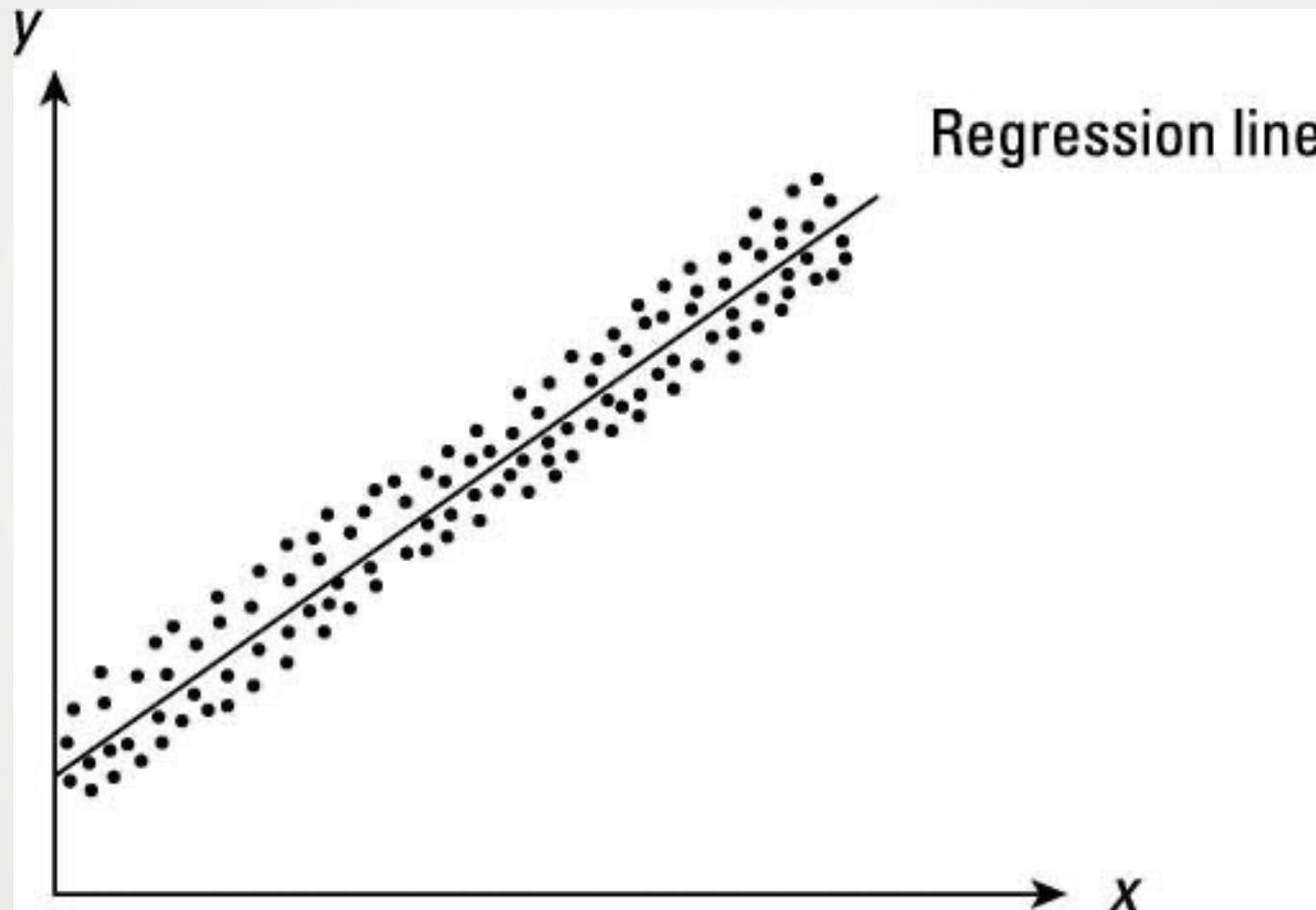
# HETEROSKEDASTICITY & HOMOSKEDASTICITY

- Concerns the *conditional variance* of the error term:

$$\text{var}(u|X = x)$$

- Homoskedasticity:  $\text{var}(u|X = x)$  is constant (does not depend on X).
- Heteroskedasticity:  $\text{var}(u|X = x)$  varies with X.

# *HOMOSKEDASTICITY IN A PICTURE:*



# VARIANCE OF $\hat{\beta}_1$ UNDER HOMOSKEDASTICITY

- In general the variance of  $\hat{\beta}_1$  is

$$var(\hat{\beta}_1) = \frac{1}{n} \times \frac{var[(X_i - \mu_x)]u_i]}{[var(X_i)]^2}$$

- If  $u_i$  is homoscedastic we get a simpler formula:

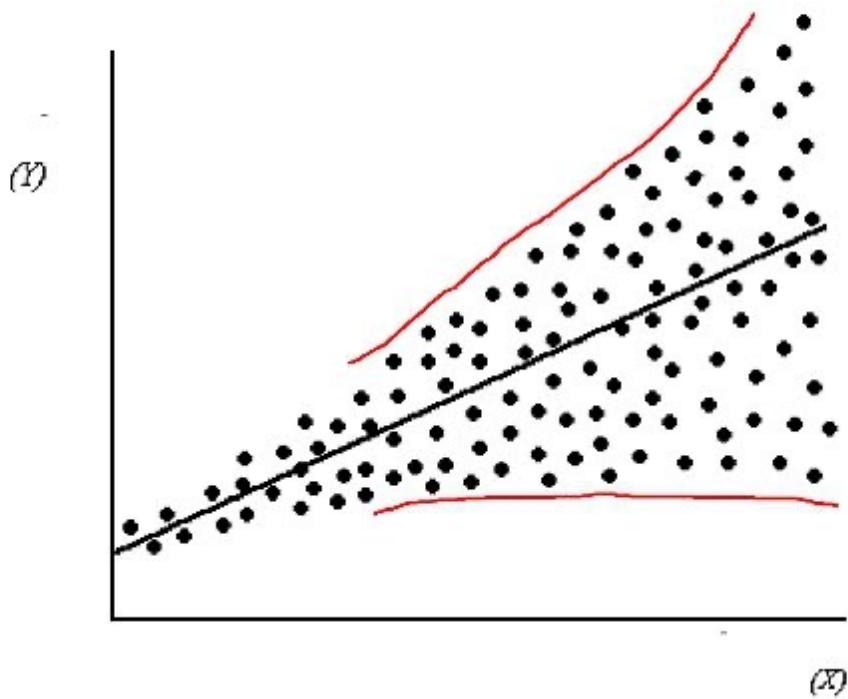
$$var(\hat{\beta}_1) = \frac{1}{n} \times \frac{var[u_i]}{[var(X_i)]^2}$$

- So the homoskedasticity-only standard error of  $\hat{\beta}_1$  is also simpler:

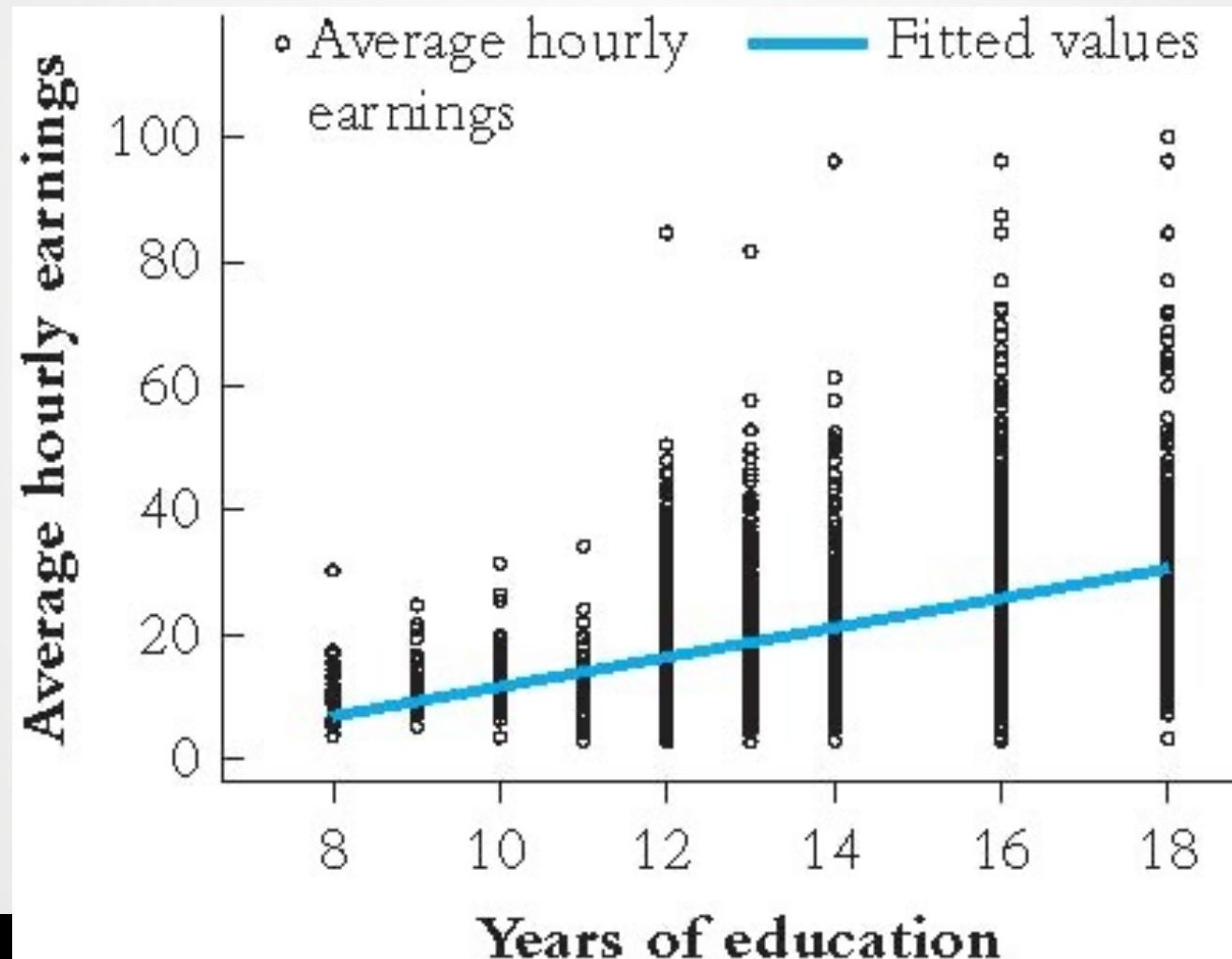
$$SE(\hat{\beta}_1) = \sqrt{\hat{\sigma}_{\hat{\beta}_1}^2} = \sqrt{\frac{\sigma_u^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}$$

# HETEROSKEDASTICITY IN A PICTURE:

Heteroskedasticity



## REAL-DATA EXAMPLE OF HETROSKEDEASTICITY (CPS DATA)



# SO, WHEN SHOULD YOU USE HOMOSKEDASTICITY-ONLY STANDARD ERRORS?

- Never.
- Our usual (heteroskedasticity-robust) SEs are always fine.
- Always use the *robust* option in STATA!