

EEP/IAS 118 - Introductory Applied  
Econometrics, Lecture 5

July 2017

# Lecture Outline

This Lecture:

- Omitted Variable Bias
- Stat Review II: Confidence Intervals

Assignments:

- Problem Set 2 due next Wednesday
- Quiz 2 Next Wednesday

# Omitted Variable Bias

Recall the end of last lecture we showed how the coefficient on *profocc* changed dramatically when we removed *educ*

$$\begin{array}{l} \widehat{\log(\text{wage})} = 1.17 + .106 \text{ educ} + .011 \text{ exp} - .26 \text{ female} + .012 \text{ profocc} \quad R^2 = .28 \\ \quad \quad \quad (.08) \quad (.005) \quad (.0009) \quad (.02) \quad (.03) \quad n = 2000 \\ \\ \widehat{\log(\text{wage})} = 2.57 + \quad \quad \quad + .011 \text{ exp} - .26 \text{ female} + .358 \text{ profocc} \quad R^2 = .16 \\ \quad \quad \quad (.03) \quad \quad \quad (.0009) \quad (.02) \quad (.03) \quad n = 2000 \end{array}$$

We will investigate why this happened and what this implies about our confidence in  $\hat{\beta}$

# Omitted Variable Bias: Motivation

Let's examine another real world example:

- School lunch programs were developed to fight undernutrition and boost learning. We have a sample of 408 schools (unit of observation is the school) with data on % of kids in school lunches and % of kids who passed a math test.
- To investigate we regress the math scores on the percentage of kids taking advantage of the free lunch:

$$\mathit{Math} = 61.4 - 0.45\mathit{lunch}\%$$

- Interpret  $\Rightarrow$  a one percentage point increase in free lunches decreases predicted math scores by 0.45 points. Implies free lunch *harms* school performance. Does this makes sense?

# Omitted Variable Bias: Motivation

- What's wrong / missing: in the areas with lots of school lunch programs there are a higher poverty levels, which contributes to lower test scores
- The model we should have had is:

$$Math = \beta_0 + \beta_1 lunch\% + \beta_2 povertyrate + u$$

If the effect of lunches are positive but the effect of poverty is highly negative, then we will falsely attribute the negative poverty effect to the school lunch program by not including an indicator for poverty rate.

- **We will get downward bias in our estimate of the effect of school lunch on math scores**

# Omitted Variable Bias

Let's examine more technically what is happening. Recall:

- **Assumption MLR3:**  $\mathbb{E}[u|x_1 \cdots x_k] = 0$ 
  - This is necessary to obtain an unbiased estimate of  $\beta$
  - $\mathbb{E}[\hat{\beta}_1] = \beta_1$
- **As we've seen this assumption can fail**
  - One way it can fail is if we fail to include a relevant variable (i.e. that explains  $y$ ) that is also correlated with the included  $x$ .
- **Consequence:** Biased estimates
  - $\mathbb{E}[\hat{\beta}_1] \neq \beta_1$
  - Commonly referred to as Omitted variable bias (OVB)

*Let's draw a graph using previous example to build intuition why failure of  $\mathbb{E}[u|x_1 \cdots x_k] = 0$  implies  $\mathbb{E}[\hat{\beta}_1] \neq \beta_1$*

# Omitted Variable Bias: Math

Let's work through the math of why  $\mathbb{E}[\hat{\beta}_1] \neq \beta_1$  when this assumption fails

- The true population regression is given by:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

- But we choose an underspecified model:

$$y = \tilde{\beta}_0 + \tilde{\beta}_1 x_1 + \tilde{u}$$

- We want to know how does  $\tilde{\beta}_1$  relate to  $\beta_1$ ?

# Omitted Variable Bias: Math

We have already seen that bias occurs when the two  $x$  are correlated. We can express this relationship with:

$$x_2 = a + \rho x_1 + v$$

Then the true model can be expressed as:

$$\begin{aligned} y &= \beta_0 + \beta_1 x_1 + \beta_2 (a + \rho x_1 + v) + u \\ &= \underbrace{\beta_0 + \beta_2 a}_{\tilde{\beta}_0} + \underbrace{(\beta_1 + \beta_2 \rho)}_{\tilde{\beta}_1} x_1 + \underbrace{u + \beta_2 v}_{\tilde{u}} \end{aligned}$$

Then:

$$\underbrace{\tilde{\beta}_1}_{\text{underspecified}} = \underbrace{\beta_1}_{\text{true effect of lunch}} + \rho \underbrace{\beta_2}_{\text{effect of poverty}}$$



# Omitted Variable Bias

- We have found that that omitting a relevant variable that is correlated with  $x$  leads to the following expression:

$$E[\tilde{\beta}_1] = \beta_1 + \beta_2\rho_{12}$$

Summary of these terms:

- $\tilde{\beta}_1$  = coefficient on  $X$  from the “biased” regression
- $\beta_1$  = coefficient on  $X$  from the “unbiased” regression  
→  $E[\tilde{\beta}_1] - \beta_1 = \text{bias}$
- $\beta_2$  = coefficient on  $X_{\text{omitted}}$  from the “unbiased” regression  
→ sign of the relationship between  $X_{\text{omitted}}$  and  $Y$
- $\rho$  = correlation between  $X$  and  $X_{\text{omitted}}$   
→ sign of the relationship between  $X$  and  $X_{\text{omitted}}$

# Omitted Variable Bias

- Looking at the expression for bias reveals another important fact:

$$E[\tilde{\beta}_1] = \beta_1 + \beta_2\rho_{12}$$

- In order for there to be bias we need both

$$\rho_{12} \neq 0, \text{ and } \beta_2 \neq 0$$

- Say if  $x_2$  is unrelated to  $y$  ( $\beta_2 = 0$ ), then leaving  $x_2$  out does not induce bias
- In this case, including  $x_2$  simply increases the variance of  $\hat{\beta}_1$  (makes the estimate less precise)

# Omitted Variable Bias

This chart summarizes resulting bias on our included variable  $x$  when we omit the variable  $x_{ov}$  depending on the covariance between  $x_{ov}$  and  $y$  and  $x_{ov}$  and the included  $x$ :

	$Cov(x, x_{ov}) > 0$	$Cov(x, x_{ov}) < 0$
$Cov(y, x_{ov}) > 0$	Upward bias	Downward bias
$Cov(y, x_{ov}) < 0$	Downward bias	Upward bias

You can also use this table to determine the covariance of  $x$  and  $x_{ov}$  if you know the sign of the bias and the relationship between  $x_{ov}$  and  $y$

## Omitted Variable Bias: Practice example

Let's think about another example: a model of homicide as a function of police financing (assume more police financing in reality leads to fewer killings):

$$\text{homicide} = \beta_0 + \beta_1 \text{police finance} + u$$

- What is missing from this model?
- Many things, but let's focus on one variable: the level of gang presence in an area

## Omitted Variable Bias: Practice example

Why does leaving out gang violence lead to a biased estimate of  $\beta_1$ ?

$$\text{homicide} = \beta_0 + \beta_1 \text{policefinance} + u$$

- Police financing and the level of gang violence are themselves likely to be correlated. In areas with lots of gangs, the police likely receive more money  $\Rightarrow \rho_{12} = \text{cov}(\text{police}, \text{gangs}) > 0$
- At the same time gangs also lead to more homicides in an area  $\Rightarrow \beta_2 = \text{cov}(\text{homicide}, \text{gangs}) > 0$
- Therefore, if we don't account for gang activity, it might seem like more police financing actually causes *more* homicides. But really we are just picking up the effect of gang activity!
- This implies that our estimator  $\hat{\beta}_1$  will be **upward** biased.

## OVB: Question Types

In an over simplification, there are two core things you will be asked to do in an OVB question:

- 1 Given the sign of two out of three of  $\beta_1$ ,  $\beta_2$ ,  $\rho_{12}$  and then asked to find the sign of the third
- 2 Given a model, asked to think about plausible omitted variables and then asked to sign the bias

# OVB Example - Question Type 1

We ran the following two regressions:

$$\widehat{\ln(wage)} = 1.19 + 0.101educ + 0.011exp$$

$$\widehat{\ln(wage)} = 1.06 + 0.117educ + 0.011exp - 0.25female$$

- 1 Interpret the coefficients on *educ* in both regressions
- 2 In what direction was the coefficient on *educ* biased due to the exclusion of *female* from the regression?
- 3 Discuss the coefficient on *female*
- 4 Based on 1) the direction of bias and 2) the coefficient on *female*, what does this imply about the covariance between *female* and *educ*?

## OVB Example - Question Type 1

$$\widehat{\ln(wage)} = 1.19 + 0.101educ + 0.011exp$$

$$\widehat{\ln(wage)} = 1.06 + 0.117educ + 0.011exp - 0.25female$$

- 1 Interpret the coefficients on *educ* in both regressions  
*A one year increase in educ leads to a predicted 10.1% (11.7%) increase in wages*
- 2 In what direction was the coefficient on *educ* biased?  
*0.101 - 0.117 = -0.016, so we have downward bias*
- 3 What does this imply about the covariance between *female* and *educ*?

*We see that*

- $cov(female, wage) < 0$
- Downward bias
- $\Rightarrow (-) = cov(fem, educ) * (-) \Rightarrow cov(fem, educ) > 0$



## OVB Example - Question Type 2

Anderson (2008) examine whether state "primary" seat belt laws (e.g., cops can pull you over just for not wearing your seat belt) reduces traffic fatalities. Suppose we run this regression on population, and the presence of the law:

$$\widehat{fatalities} = \hat{\beta}_0 + \hat{\beta}_1 pop + \hat{\beta}_2 primary$$

$$\widehat{fatalities} = 156.002 + 0.1232pop + 17.258primary$$

- 1 If we were naive (i.e., weren't concerned about OVB), how would we interpret this regression?
- 2 Identify a possible important omitted variable
- 3 Sign the bias this omission would cause on  $\hat{\beta}_{primary}$

## OVB Example - Question Type 2

$$\widehat{fatalities} = 156.002 + 0.1232pop + 17.258primary$$

- 1 If we were naive (i.e. weren't concerned about OVB), how would we interpret this regression?

*Having a primary seatbelt law actually **increases** traffic fatalities! This surprising result should tip us off that OVB is a possible problem*

- 2 Identify a possible important omitted variable

*State speed limit is an important omitted variable. States with high speed limits are more likely to pass a primary seatbelt law*

- 3 Sign the bias this omission would cause on  $\hat{\beta}_{primary}$

*We know:*

- $cov(speed, primary) > 0$
- $cov(speed, fatalities) > 0$
- $\Rightarrow$  upward bias

# OVB - Review

As a check we can go back to the three examples from last lecture and answer these more rigorously

What happens when we:

- 1 Omitting an irrelevant  $x$  variable? Why?
- 2 Omitting an important variable that is not correlated with the other independent variable? Why?
- 3 Omitting an important variable that is correlated with the other independent variable? Why?

# OVB - Review

What happens when:

- 1 Omitting an irrelevant  $x_2$  variable? Why?
  - $\beta_2 = 0$  therefore there will be no bias from leaving out  $x_2$ .  
Including  $x_2$  could also reduce the precision of  $\hat{\beta}_1$  because it increases  $R_{x_1}^2$  in the variance formula
- 2 Omitting an important variable that is not correlated with the other independent variable? Why?
  - $\rho_{12} = 0$  therefore there will be no bias from leaving out  $x_2$ .  
However, we could increase the precision of  $\hat{\beta}_1$  if  $\beta_2$  is large
- 3 Omitting an important variable that is correlated with the other independent variable? Why?
  - Both  $\rho_{12} \neq 0$  and  $\beta_2 \neq 0$ , therefore leaving out  $x_2$  will introduce bias into  $\beta_1$

# OVB: Final Thoughts

We've seen how omitting important variables can lead to biased estimates of  $\beta$

- This is a *very* common problem - we almost never have enough data so that we haven't omitted anything important
  - How do we get data on things like ability? Commitment? Family connections?
- Then why do we bother?
  - Regression can still be useful for *predictions* of  $y$  even though  $\hat{\beta}$ s are biased
  - Econometrics is primarily concerned with developing techniques and finding conditions under which we are more confident that we satisfy  $\mathbb{E}[u|x_1 \cdots x_k] = 0$
  - Much of the second part of the course will be covering some of these techniques

# Switching gears: Confidence Intervals & Hypothesis Testing

- We are going to be transitioning from point estimates (and bias) to hypothesis testing
- Why? Because now we want to use statistics to tell us how much uncertainty we should have in our estimates
- While OLS gives us the best possible fit of our model to our sample data, we would like to know how close the estimate is likely to be from the true **population** parameter.
- Today, introduce the notion of confidence intervals

## Statistics Reminders

We have a random variable  $Y$ , and we know in the *population* that  $E[y] = \mu$  and  $Var(Y) = \sigma^2$ . If we then get a sample  $Y_1, \dots, Y_n$ , we can build an estimator for  $\mu$ :

$$\bar{Y} = \frac{1}{n} \sum_i Y_i$$

This estimator  $\bar{Y}$  is itself a random variable and therefore has an expected value and a variance:

$$E(\bar{Y}) = \frac{1}{n} \sum_i E(Y_i) = \mu$$

$$Var(\bar{Y}) = Var\left(\frac{1}{n} \sum_i Y_i\right) = \frac{1}{n^2} \left(\sum_i Var(Y_i)\right) = \frac{\sigma^2}{n}$$

# Statistics Reminders

Using these facts, we can write the distribution of  $\bar{Y}$  as

$$\bar{Y} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

- the first term ( $\mu$ ) is the mean of the distribution
- the second term ( $\frac{\sigma^2}{n}$ ) is the variance of the distribution
- the  $N(\cdot)$  indicates the Normal distribution - but why a normal distribution?



# Important Theorem in Statistics

The central limit theorem (CLT) states that the average from a random sample for any population (with finite variance), when standardized, has an asymptotic standard normal **distribution**.

Consider a random sample  $X_1, \dots, X_n$  from a population with mean  $\mu$  and variance  $\sigma^2$ , then

$$Z_n = \frac{\bar{X}_n - \mu}{\sigma_X / \sqrt{n}} \xrightarrow{d} N(0, 1)$$

⇒ If we take many samples and calculate the sample means ( $\bar{X}_n$ ), these will be normally distributed. If we then subtract the true population mean and divide by the true population variance, the distribution of this new random variable  $Z_n$  has PDF that is a standard normal (mean zero, standard deviation one)

# Confidence Intervals

*Why is this useful?*

- We have that (by application of the CLT)

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

Which means that

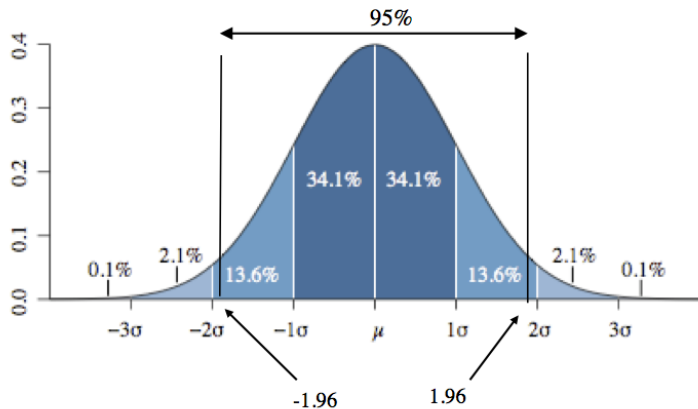
$$\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$$

- Which says that this transformation of  $\bar{X}$  is distributed as a *standard normal*. This is an easy PDF to work with

# Standard Normal Distribution

- 72% of observations lie within one standard deviation
- 95% of observations lie within approximately two standard deviations
- More specifically, we know that for any standard normal variable  $v$ ,  $Pr(-1.96 < v < 1.96) = 95\%$ 
  - Also  $Pr(-1.65 < v < 1.56) = 90\%$
  - And  $Pr(-2.56 < v < 2.56) = 99\%$
- **Note:** we call this value (1.65, 1.96, 2.56) the “critical value”  
- these values correspond to know points in a particular distribution such that 90%, 95%, or 99% of the probability falls between their positive and negative values

# Standard Normal Distribution



# Confidence Intervals

We can take these facts to write:

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \rightarrow \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

Then from our knowledge of the standard normal:

$$\Pr\left(-1.96 < \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} < 1.96\right) = 0.95$$

$$\Pr\left(-1.96 \frac{\sigma}{\sqrt{n}} < \bar{X} - \mu < 1.96 \frac{\sigma}{\sqrt{n}}\right) = 0.95$$

$$\Pr\left(-\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}} < -\mu < -\bar{X} + 1.96 \frac{\sigma}{\sqrt{n}}\right) = 0.95$$

$$\Pr\left(\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}}\right) = 0.95$$

# Confidence Intervals

$$\Pr(\bar{X} - 1.96 \frac{\sigma_X}{\sqrt{n}} < \mu < \bar{X} + 1.96 \frac{\sigma_X}{\sqrt{n}}) = 0.95$$

*This* now is very powerful, we have quantified the uncertainty of our estimator

- This equation indicates that the random range defined by  $[\bar{X} - 1.96 \frac{\sigma_X}{\sqrt{n}}, \bar{X} + 1.96 \frac{\sigma_X}{\sqrt{n}}]$  contains the true  $\mu$  with 95% probability
- **Note:** the wording here is *very* specific. We do not want to imply that  $\mu$  is a random variable. The “randomness” comes from  $\bar{X}$  depending on the sample
  - **DON'T** say “there is a 95% chance that  $\mu$  is in the confidence interval” - you will lose points

# Confidence Intervals

So, we have defined a 95% confidence interval:

$$\left[ \bar{X} - 1.96 \frac{\sigma_X}{\sqrt{n}}, \bar{X} + 1.96 \frac{\sigma_X}{\sqrt{n}} \right]$$

- Remember though that we don't ever observe the true  $\sigma_X$
- As before, we have to estimate it with  $s_X$

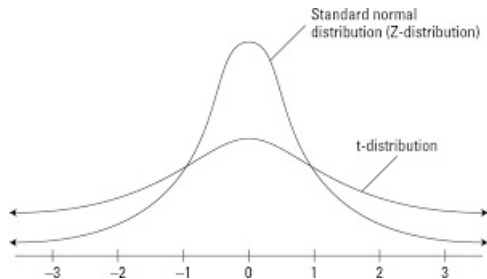
$$s^2 = \frac{(x_i - \bar{X})^2}{n - 1}$$

- This costs us something: **we lose the normality of the resulting distribution!**

# Confidence Intervals

- Instead, we have to use the (student) t-distribution, which will widen our confidence interval

$$\frac{\bar{X} - \mu}{s_x / \sqrt{n}} \sim t_{n-1}$$





# Confidence Intervals

$$\frac{\bar{X} - \mu}{s_x / \sqrt{n}} \sim t_{n-1}$$

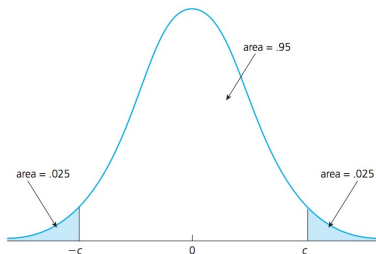
Fortunately, the t-distribution is also easy (ish) to work with:

- the  $n - 1$  is called the “degrees of freedom” - this affects how wide the t-distribution is
- The specific values for which 95% of observations fall between is now NOT 1.96. Instead the number we use in the 95% confidence will generally be higher than 1.96 (and depends on the sample size)
- **Note:** When  $n$  is large, the t-distribution is indistinguishable from the normal distribution
  - Roughly, once  $n$  is larger than 200 use a standard normal

# Constructing Confidence Intervals: Five Steps

We take a random sample of 121 UCB students' heights in inches. Now, to construct a confidence interval for the average height of UCB students:

- 1 Determine the confidence level - standard is 95%, but 99% and 90% are also used.
- 2 Compute  $\bar{X}$  and  $s_X$ . Let's say  $\bar{X} = 65$  and  $s_X^2 = 4$
- 3 Find critical value,  $c$ , from the t-table.  $c$  will depend on sample size ( $n$ ) and the confidence level:



## t-table

TABLE B: *t*-DISTRIBUTION CRITICAL VALUES

df	Tail probability <i>p</i>											
	.25	.20	.15	.10	.05	.025	.02	.01	.005	.0025	.001	.0005
1	1.000	1.376	1.963	3.078	6.314	12.71	15.89	31.82	63.66	127.3	318.3	636.6
2	.816	1.061	1.386	1.886	2.920	4.303	4.849	6.965	9.925	14.09	22.33	31.60
3	.765	.978	1.250	1.638	2.353	3.182	3.482	4.541	5.841	7.453	10.21	12.92
4	.741	.941	1.190	1.533	2.132	2.776	2.999	3.747	4.604	5.598	7.173	8.610
5	.727	.920	1.156	1.476	2.015	2.571	2.757	3.365	4.032	4.773	5.893	6.869
6	.718	.906	1.134	1.440	1.943	2.447	2.612	3.143	3.707	4.317	5.208	5.959
7	.711	.896	1.119	1.415	1.895	2.365	2.517	2.998	3.499	4.029	4.785	5.408
8	.706	.889	1.108	1.397	1.860	2.306	2.449	2.896	3.355	3.833	4.501	5.041
9	.703	.883	1.100	1.383	1.833	2.262	2.398	2.821	3.250	3.690	4.297	4.781
10	.700	.879	1.093	1.372	1.812	2.228	2.359	2.764	3.169	3.581	4.144	4.587
11	.697	.876	1.088	1.363	1.796	2.201	2.328	2.718	3.106	3.497	4.025	4.437
12	.695	.873	1.083	1.356	1.782	2.179	2.303	2.681	3.055	3.428	3.930	4.318
13	.694	.870	1.079	1.350	1.771	2.160	2.282	2.650	3.012	3.372	3.852	4.221
14	.692	.868	1.076	1.345	1.761	2.145	2.264	2.624	2.977	3.326	3.787	4.140
15	.691	.866	1.074	1.341	1.753	2.131	2.249	2.602	2.947	3.286	3.733	4.073
16	.690	.865	1.071	1.337	1.746	2.120	2.235	2.583	2.921	3.252	3.686	4.015
17	.689	.863	1.069	1.333	1.740	2.110	2.224	2.567	2.898	3.222	3.646	3.965
18	.688	.862	1.067	1.330	1.734	2.101	2.214	2.552	2.878	3.197	3.611	3.922
19	.688	.861	1.066	1.328	1.729	2.093	2.205	2.539	2.861	3.174	3.579	3.883
20	.687	.860	1.064	1.325	1.725	2.086	2.197	2.528	2.845	3.153	3.552	3.850
21	.686	.859	1.063	1.323	1.721	2.080	2.189	2.518	2.831	3.135	3.527	3.819
22	.686	.858	1.061	1.321	1.717	2.074	2.183	2.508	2.819	3.119	3.505	3.792
23	.685	.858	1.060	1.319	1.714	2.069	2.177	2.500	2.807	3.104	3.485	3.768
24	.685	.857	1.059	1.318	1.711	2.064	2.172	2.492	2.797	3.091	3.467	3.745
25	.684	.856	1.058	1.316	1.708	2.060	2.167	2.485	2.787	3.078	3.450	3.725
26	.684	.856	1.058	1.315	1.706	2.056	2.162	2.479	2.779	3.067	3.435	3.707
27	.684	.855	1.057	1.314	1.703	2.052	2.158	2.473	2.771	3.057	3.421	3.690
28	.683	.855	1.056	1.313	1.701	2.048	2.154	2.467	2.763	3.047	3.408	3.674
29	.683	.854	1.055	1.311	1.699	2.045	2.150	2.462	2.756	3.038	3.396	3.659
30	.683	.854	1.055	1.310	1.697	2.042	2.147	2.457	2.750	3.030	3.385	3.646
40	.681	.851	1.050	1.303	1.684	2.021	2.123	2.423	2.704	2.971	3.307	3.551
50	.679	.849	1.047	1.299	1.676	2.009	2.109	2.403	2.678	2.937	3.261	3.496
60	.679	.848	1.045	1.296	1.671	2.000	2.099	2.390	2.660	2.915	3.232	3.460
80	.678	.846	1.043	1.292	1.664	1.990	2.088	2.374	2.639	2.887	3.195	3.416
100	.677	.845	1.042	1.290	1.660	1.984	2.081	2.364	2.626	2.871	3.174	3.390
1000	.675	.842	1.037	1.282	1.646	1.962	2.056	2.330	2.581	2.813	3.098	3.300
∞	.674	.841	1.036	1.282	1.645	1.960	2.054	2.326	2.576	2.802	3.091	3.291

# Constructing Confidence Intervals: Five Steps

- 4 Plug everything into the confidence interval formula:

$$CI = \left[ \bar{X} - c \cdot \underbrace{\left( \frac{s_x}{\sqrt{n}} \right)}_{se(\bar{X})}, \bar{X} + c \left( \frac{s_x}{\sqrt{n}} \right) \right]$$

- Remember,  $c$  is found by looking at the t-table for  $n - 1$  degrees of freedom for the desired confidence level
  - $\bar{X}$ ,  $s_x$ , and  $n$  we can calculate from the sample
- 5 Interpret: There is a 95% probability that this interval covers our true value.

# Constructing Confidence Intervals: Five Steps

From our example:

- $c = 1.98$  (found in t-table for 120 ( $n-1$ ) degrees of freedom)
- $\bar{X} = 65$
- $s_X = 2$
- $n = 121$

plugging everything in yields:

$$CI = \left[ 65 - 1.98 \left( \frac{2}{\sqrt{121}} \right), 65 + 1.98 \left( \frac{2}{\sqrt{121}} \right) \right]$$

Doing the math, the 95% confidence interval is [64.64, 65.36].

## Special Case: Binary Variables

Let's say  $x$  can only take on a zero or one value. Then  $p$  is the true (but unknown) proportion of 1 in the population. Recall that for one observation  $x$ :

$$E(x) = p$$

$$\text{Var}(x) = p(1 - p)$$

For a sample  $x_1, \dots, x_n$ , if we find the sample average  $\bar{X}$ :

$$E(\bar{X}) = p$$

$$\text{Var}(\bar{X}) = \frac{p(1 - p)}{n}$$

$$\text{Std}(\bar{X}) = \sqrt{\frac{p(1 - p)}{n}}$$

What does this mean about how we calculate a confidence interval for  $p$  from an estimator  $\hat{p}$ ?

## Special Case: Binary Variables

From the CI formula:

$$\left[ \hat{p} - 1.96\sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}, \hat{p} + 1.96\sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \right]$$

- With a continuous variable, we had to estimate  $\sigma_X$  in the standard deviation formula
- For a binary variable, we don't have to do this! We have already estimated everything we need to the standard deviation in  $\hat{p}$
- Implies we don't lose our normal distribution! Therefore, with binary variables you always choose the critical values from the standard normal (z) distribution

## Special Case: Binary Variable Example

Let's say we there is a poll that asks support for presidential candidate A. The poll asks 130 registered voters and 45% support candidate A. What is the 95% confidence interval for candidate A's support?

- 1 Confidence level is given: 95%
- 2 Compute  $\hat{p}$  and the standard error

$$\hat{p} = .45$$
$$se = \sqrt{\frac{0.45 * 0.55}{130}} = 0.043$$

- 3 Find the critical value from the z-table (it is 1.96)



## Special Case: Binary Variable Example

- 4 Plug everything into the formula

$$[0.45 - 1.96 * 0.043, 0.45 + 1.96 * 0.043]$$

$$[0.365, 0.534]$$

- 5 Interpret: There is a 95% probability that the true  $p$  is contained in the random interval  $[0.365, 0.534]$