

## 6. Statistical Inference: Estimation

LPO.8800: Statistical Methods in Education Research

Sean P. Corcoran

### Last time

- The *sampling distribution* of  $\bar{x}$
- *Standard error*: a measure of variability in a sampling distribution
- Central Limit Theorem: tells us the distribution of  $\bar{x}$  when  $x$  is normally distributed or otherwise
- $\bar{x}$  is *unbiased*
- Efficiency of estimators
- Sampling distribution for proportions ( $\hat{\pi}$ )
- Simulating sampling distributions in Stata

## Estimating population parameters

The objective of inferential statistics is to make inferences about **population parameters** (e.g.,  $\mu$ ,  $\sigma^2$ ) using **sample statistics** (e.g.,  $\bar{x}$ ,  $s^2$ ).

- Estimators vs. estimates
  - ▶ The term **estimator** refers to the method, or statistic (e.g.  $\bar{x}$ ), used to estimate a population parameter (e.g.  $\mu$ ). Often represented by a “hat” over the parameter being estimated, such as  $\hat{\mu}$  or  $\hat{\sigma}$ .
  - ▶ When used as a noun, an **estimate** refers to a specific realization of a sample statistic. An *estimator* gives you an *estimate*.
- Point vs. interval estimation
  - ▶ **Point estimation** is the process of estimating a specific parameter.
  - ▶ **Interval estimation** is the process of estimating a *range* of likely values for the parameter. Takes into account variability in the sampling distribution of the estimator.

## Point estimates

Lecture 5 was devoted to the sampling distribution of  $\bar{x}$ , a point estimator of the population mean  $\mu$ . We saw that:

- The sampling distribution of  $\bar{x}$  has a mean of  $\mu$  and standard error
$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$
- According to the Central Limit Theorem:
  - ▶ If  $x$  is distributed *normal*: the sampling distribution of  $\bar{x}$  is normal.
  - ▶ If  $x$  is *not* distributed normal: the sampling distribution of  $\bar{x}$  is *approximately* normal if  $n$  is sufficiently large.

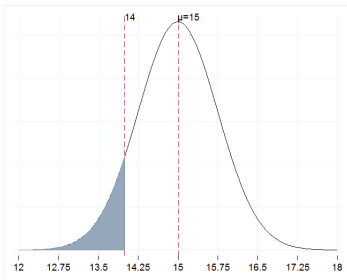
## Example From Lecture 5

Suppose we plan to draw a random sample of 16 from the population of  $x$  and compute the sample mean  $\bar{x}$ . We know  $x$  is normally distributed with  $\mu=15$  and  $\sigma=3$ . What will the sampling distribution of  $\bar{x}$  look like? From the CLT we know:

- $\bar{x}$  will have a normal distribution
- $\bar{x}$  will have a mean of 15
- $\bar{x}$  will have a standard error  $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{3}{\sqrt{16}} = 0.75$

## Example From Lecture 5

Knowing its sampling distribution, we can make statements about how likely particular realized values of  $\bar{x}$  are. For example, in the above scenario, what is the probability we will draw a random sample with an  $\bar{x}$  of 14 or lower?

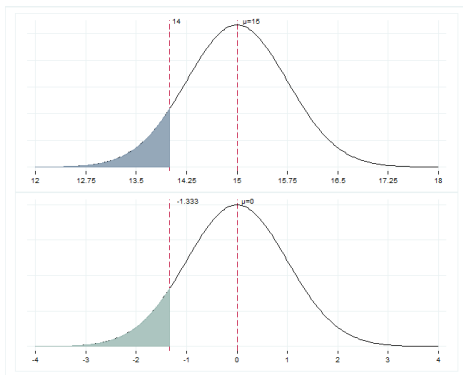


## Example From Lecture 5

Putting in z-score (standard normal) terms:

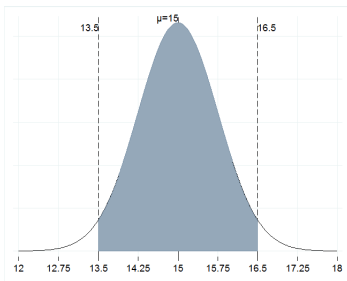
- The value  $\bar{x} = 14$  is  $z = \frac{14-15}{0.75} = -1.33$  standard errors below the population mean.
- Thus:  $\Pr(\bar{x} \leq 14) = \Pr(z \leq -1.33)$
- Stata: `display normal(-1.33) = 0.092`

## Example From Lecture 5



## Example From Lecture 5

What is the probability we will draw a random sample with an  $\bar{x}$  **between** 13.5 and 16.5?

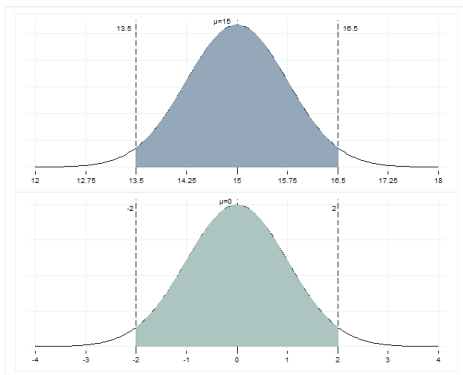


## Example From Lecture 5

Putting in z-score (standard normal) terms:

- The value  $\bar{x} = 13.5$  is  $z = \frac{13.5-15}{0.75} = -2$  standard errors below the population mean.
- The value  $\bar{x} = 16.5$  is  $z = \frac{16.5-15}{0.75} = +2$  standard errors above the population mean.
- Thus  $\Pr(13.5 \leq \bar{x} \leq 16.5) = \Pr(-2 \leq z \leq +2)$
- Stata: `display normal(2) - normal(-2) = 0.954` (about 95%)

## Review example, cont.



## Confidence intervals

The goal of interval estimation is to construct an interval that will—under repeated samples—contain the true population parameter  $(1 - \alpha)\%$  of the time.

- $(1 - \alpha)\%$  is the **confidence level**, the percentage of times in repeated samples that the interval will contain the population parameter:
  - ▶ 95% confidence level ( $\alpha = 0.05$ )
  - ▶ 99% confidence level ( $\alpha = 0.01$ )
  - ▶ 90% confidence level ( $\alpha = 0.10$ )
- $\alpha$  is the **error probability**, a measure of our tolerance for error (more on this later). We typically choose this.
- An interval estimate with a specified confidence level is a **confidence interval**.

## Confidence intervals

All else equal, the cost of narrower (more **precise**) confidence intervals is a greater probability that the interval will *not* contain the true population parameter.

- A *low confidence level* is associated with narrower (more *precise*) confidence intervals, but a higher likelihood of error—i.e., computing an interval that does not contain the population parameter.
- A *high confidence level* is associated with wider (less *precise*) confidence intervals, but a lower likelihood of error.
- This terminology can be confusing, because “precise” does not mean more *accurate*, but rather a narrower range of values.

Ideally, you'd like both narrow CIs *and* a high confidence level.

## Confidence intervals

Another approach...

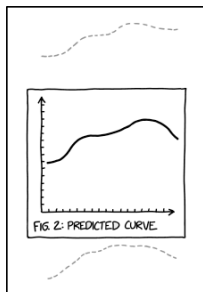


FIG. 2: PREDICTED CURVE

SCIENCE TIP: IF YOUR MODEL IS BAD ENOUGH, THE CONFIDENCE INTERVALS WILL FALL OUTSIDE THE PRINTABLE AREA.

## Confidence interval for $\mu$

Under certain assumptions, we saw that  $\bar{x}$  will fall within two standard errors of the true population mean roughly 95% of the time. More accurately,  $\bar{x}$  falls within **1.96** standard errors of the true population mean 95% of the time. If this is the case, then the interval:

$$\bar{x} - 1.96 \left( \frac{\sigma}{\sqrt{n}} \right), \bar{x} + 1.96 \left( \frac{\sigma}{\sqrt{n}} \right)$$

will include the population mean ( $\mu$ ) 95% of the time. This is known as a **95% confidence interval for  $\mu$** . By the same logic, this interval will *not* contain the population mean 5% of the time ( $\alpha = 0.05$ ).

## Confidence interval for $\mu$

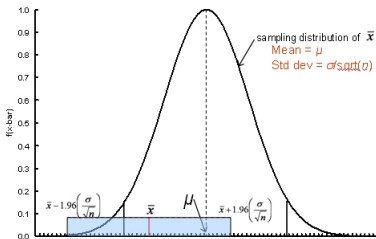


Figure: Confidence interval around  $\bar{x}$



## Confidence interval for $\mu$

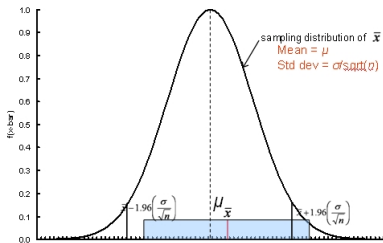


Figure: Confidence interval around  $\bar{x}$

## Confidence interval for $\mu$

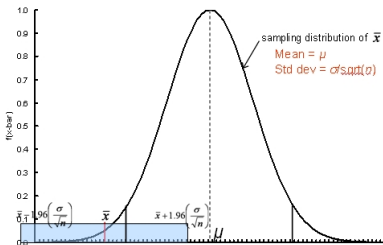


Figure: Confidence interval around  $\bar{x}$

## Example 1a

A sample of 121 women is taken ( $n=121$ ) where the mean height is estimated to be  $\bar{x} = 64$  inches. It is known that  $\sigma = 3$ . Construct a 95% confidence interval for the true population mean height ( $\mu$ ).

$$64 \pm 1.96 \left( \frac{3}{\sqrt{121}} \right) = (63.47, 64.53)$$

This is our interval estimate of the range of likely values for  $\mu$ . For later reference, note the width of this interval is 1.06 inches.

## Confidence interval for $\mu$

Changing the confidence level is only a matter of changing the  $z$  value used in the interval estimator. For a confidence level of  $(1 - \alpha)\%$ , the appropriate  $z$  value is the one for which there is a probability  $\alpha/2$  of exceeding. Then the  $(1 - \alpha)\%$  confidence interval is:

$$\bar{x} \pm z_{\alpha/2} \left( \frac{\sigma}{\sqrt{n}} \right)$$

## Confidence interval for $\mu$

Choosing the appropriate z-values:

- There is a 0.005 probability that z falls above 2.576 and a 0.005 probability that z falls below -2.576. Thus **2.576** is the z value used for a 99% confidence interval ( $\alpha = 0.01$ ).
- There is a 0.05 probability that z falls above 1.645 and a 0.05 probability that z falls below -1.645. Thus **1.645** is the z value used for a 90% confidence interval ( $\alpha = 0.10$ ).

`display (-1)*invnormal(0.005)`

`display (-1)*invnormal(0.05)`

### Example 1b

A sample of 121 women is taken ( $n=121$ ) where the mean height is estimated to be  $\bar{x} = 64$  inches. It is known that  $\sigma = 3$ . Construct a **99%** confidence interval for the true population mean height ( $\mu$ ).

$$64 \pm \mathbf{2.576} \left( \frac{3}{\sqrt{121}} \right) = (63.29, 64.71)$$

Note this is a width of 1.42 inches, a wider range of values than the 95% confidence interval.

## Example 1c

A sample of 121 women is taken ( $n=121$ ) where the mean height is estimated to be  $\bar{x} = 64$  inches. It is known that  $\sigma = 3$ . Construct a **90%** confidence interval for the true population mean height ( $\mu$ ).

$$64 \pm \mathbf{1.645} \left( \frac{3}{\sqrt{121}} \right) = (63.55, 64.45)$$

Note this has a width of 0.9 inches, a narrower range of values than the 95% confidence interval.

## Confidence intervals in Stata

Confidence intervals for the mean are straightforward in Stata with the `mean` command. (In this example we are using  $s$  as an estimator for  $\sigma$ , which was assumed to be known in the earlier examples. More on this later). For example, using the NELS data the 95% and 99% confidence intervals for mean 8th grade math achievement can be found using:

- `mean achmat08`
- `mean achmat08, level(99)`

A 95% confidence interval is the default in Stata, but the confidence level can be changed with the `level` option.

## Example 2

From a sample of 3,650 days, the average daily high temperature in Richmond, VA is estimated to be 69 degrees. It is known that  $\sigma = 17.5$ . Construct a **95%** confidence interval for the true population mean temperature in Richmond ( $\mu$ ).

$$69 \pm 1.96 \left( \frac{17.5}{\sqrt{3650}} \right) = (68.43, 69.57)$$

## Example 3

IQ scores are scaled to have a mean of 100 and a standard deviation of 16. Suppose in a sample of 42 people,  $\bar{x} = 103$ . Construct a **95%** confidence interval for mean IQ.

$$103 \pm 1.96 \left( \frac{16}{\sqrt{42}} \right) = (98.17, 107.83)$$

## Confidence interval for a proportion

- Thus far we've described a confidence interval for the mean of a random variable  $x$ . A special case is a dichotomous variable that can only take on the values 1 or 0 (e.g., will vote for Biden or not).
- Recall that the proportion of  $x$  that equal one in the population is denoted  $\pi$ . The proportion of  $x$  equal to zero is  $1 - \pi$ .
- We estimate  $\pi$  with the *sample* proportion  $\hat{\pi}$ , the mean of the dichotomous variable  $x$ . (E.g., in a sample of 900 likely voters,  $\hat{\pi} = 0.65$ ).

## Confidence interval for a proportion

- From Lectures 4-5, the population standard deviation of Bernoulli  $x$  is  $\sqrt{\pi(1 - \pi)}$ .
- From Lecture 5: the Central Limit Theorem tells us our sample proportion  $\hat{\pi}$  will be normally distributed with a mean of  $\pi$  and a standard error of  $\sqrt{\pi(1 - \pi)/n}$ .
- Knowing this, we can construct confidence intervals for  $\pi$  in the same manner as for  $\mu$ . For example, the following interval will contain  $\pi$  in 95% of samples:

$$\hat{\pi} \pm 1.96 \sqrt{\frac{\pi(1 - \pi)}{n}}$$

## Confidence interval for a proportion

When we do not know  $\pi$  (which is almost always true, otherwise we would not need to estimate it), we can use our estimate of it,  $\hat{\pi}$ , in place of  $\pi$ :

$$\hat{\pi} \pm 1.96 \sqrt{\frac{\hat{\pi}(1 - \hat{\pi})}{n}}$$

In a sample of 900 likely voters in which 65% of voters say they will vote for Joe Biden, a 95% confidence interval for the population proportion is:

$$0.65 \pm 1.96 \sqrt{\frac{0.65(1 - 0.65)}{900}}$$

$$0.65 \pm 1.96 * (0.0159) = (0.6188, 0.6812)$$

With a sample of only 900 likely voters, we are able to produce a confidence interval with a relatively small range of likely values.

## Margin of error

The product of  $z$  and the standard error is often called the **margin of error**. For example, the margin of error in the mayoral poll above is:

$$1.96 \sqrt{\frac{0.65(1 - 0.65)}{900}} = 0.031$$

or 3.1 percentage points. The margin of error is what is added/subtracted from  $\hat{\pi}$  to get the confidence interval.

## Confidence interval assumptions

- Our ability to say a confidence interval will contain the population parameter in (say) 95% of samples rests on a key assumption about the estimator's sampling distribution: *normality*.
- When the sampling distribution of  $\bar{x}$  departs from normality, our approach to estimating confidence intervals will be less valid.
- Another condition in which the above confidence interval formula may not apply is when  $\sigma$  is *unknown*, which is almost always the case.

## Confidence interval when $\sigma$ is unknown

- When  $\sigma$  is *unknown*, we can estimate it with the sample standard deviation  $s$ .
- In doing so, we introduce additional variability that would not be present if we knew  $\sigma$ .
- The standardized version of  $\bar{x}$  no longer follows a standard normal ( $z$ ) distribution, but rather a **Student's  $t$  distribution** (see next slide).
- The confidence interval for  $\mu$  then becomes:

$$\bar{x} \pm t_{\alpha/2} \left( \frac{s}{\sqrt{n}} \right)$$



# The $t$ distribution

The  $t$ -distribution is similar to the standard normal ( $z$ ):

- It is symmetric and bell-shaped.
- Its mean and median are zero.
- The standard deviation depends on its *degrees of freedom*, which in this case is  $df = n - 1$ .
- For small  $n$ , the tails of the  $t$ -distribution are thicker than those of the standard normal (i.e., greater variance).
- For large  $n$ , the  $t$ -distribution looks approximately like the standard normal.

## The $t$ distribution

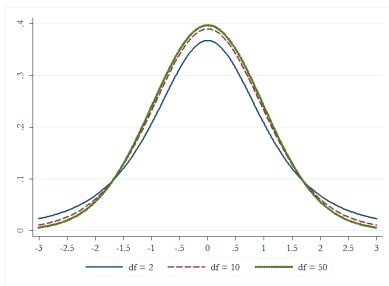


Figure:  $t$ -distribution with 2, 10, and 350  $df$

## The $t$ distribution

What are the implications of thicker tails? The probability of falling within 1.96 standard deviations of the mean (0) is lower for  $t$  than for  $z$ . Also the likelihood of more extreme values (the tails) is somewhat higher.

- Using the  $df$ , one can easily use Stata, statistical tables, or an online calculator to find values of  $t$  associated with particular probabilities. (Similar to what we did for  $z$ ).

## The $t$ distribution—using Stata

Using Stata to find probabilities and values from the  $t$  distribution:

- `display ttail( $df$ ,  $t$ )` finds the probability of exceeding a given  $t$  value. Comparable to `display normal( $z$ )`, except that `ttail()` gives you the probability *above* a certain  $t$ , while `normal()` gives you the probability *below* a certain  $z$ .
- `display invttail( $df$ ,  $p$ )` finds the value of  $t$  for which there is a probability  $p$  of exceeding. Comparable to `display invnormal( $p$ )`, except that `invttail()` gives you the  $t$  for which the probability *above*  $t$  is  $p$ , while `invnormal()` gives you the  $z$  for which the probability *below*  $z$  is  $p$ .
- When finding the  $t$  to use in a confidence interval,  $\alpha/2$  is your  $p$  in these functions!

# The $t$ distribution—using Stata

- For example, suppose  $n=21$ . To find the  $t$  value for a 95% confidence interval, use display `invttail(20,0.025)`.
- $t = 2.086$ , which you should note is *larger* than the 1.96 value used in 95% confidence intervals when  $\sigma$  was known.

## The $t$ distribution

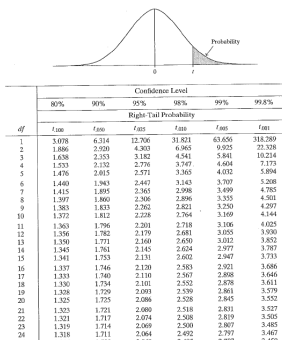


Figure:  $t$ -distribution lookup table

## Example 1, revisited

This time assume a sample of 41 women is taken ( $n=41$ ) where the mean height is estimated to be  $\bar{x} = 64$  inches. We estimate  $s = 3$ . Construct a 95% confidence interval for the true population mean height ( $\mu$ ).

$$\bar{x} \pm t_{0.025} \left( \frac{s}{\sqrt{n}} \right)$$
$$64 \pm 2.021 \left( \frac{3}{\sqrt{41}} \right) = (63.05, 64.95)$$

The  $t$  value of 2.021 comes from the  $t$ -distribution table with  $df=40$  and  $\alpha/2 = 0.025$ . In Stata:

```
display invttail(40,0.025)
```

## Example 4

From a sample of 20, suppose the mean test score is 76.1 with  $s = 15.2$ . Construct a **90%** confidence interval for the true population mean test score ( $\mu$ ).

$$\bar{x} \pm t_{0.05} \left( \frac{s}{\sqrt{n}} \right)$$
$$76.1 \pm 1.729 \left( \frac{15.2}{\sqrt{20}} \right) = (70.22, 81.98)$$

The  $t$  value of 1.729 comes from the  $t$ -distribution table with  $df=19$  and  $\alpha/2 = 0.05$ . In Stata:

```
display invttail(19,0.05)
```

## Example 5

Nabisco, makers of Chips Ahoy, claims that there are 1000 chocolate chips in each 18 oz bag of cookies. The company offered \$25,000 in scholarships to schools if students could verify the claim. The Air Force Academy collected a sample of 42 bags.  $\bar{x} = 1261$  and  $s = 117.6$ . Construct a 95% confidence interval for the true mean number of chips ( $\mu$ ).

$$\bar{x} \pm t_{0.025} \left( \frac{s}{\sqrt{n}} \right)$$
$$1261 \pm 2.02 \left( \frac{117.6}{\sqrt{42}} \right) = (1222.4, 1297.6)$$

The  $t$  value of 2.02 comes from the  $t$ -distribution table with  $df=41$  and  $\alpha/2 = 0.025$ . In Stata:

```
display invttail(41,0.025)
```

## Confidence intervals in Stata

- Because  $\sigma$  is always estimated, Stata's `mean` command uses  $t$  values in constructing confidence intervals. These values are close to  $z$  whenever  $n$  is large.

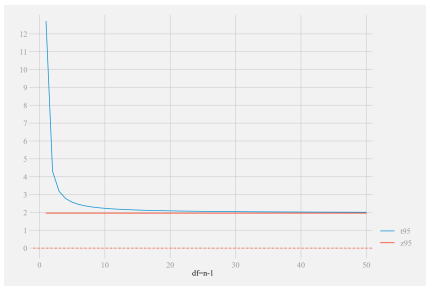


Figure:  $z$  versus  $t$  for a 95% confidence interval

## Margins of error and sample size

In all of the above examples, confidence intervals have taken the form of **point estimate  $\pm$  margin of error**. The margin of error was determined by three things: the desired confidence level, the underlying population standard deviation, and the sample size. For example:

$$\bar{x} \pm z_{\alpha/2} \left( \frac{\sigma}{\sqrt{n}} \right)$$

$$\bar{x} \pm t_{\alpha/2} \left( \frac{s}{\sqrt{n}} \right)$$

$$\hat{\pi} \pm z_{\alpha/2} \sqrt{\frac{\pi(1-\pi)}{n}}$$

## Margins of error and sample size

We can use the margin of error to determine how large a sample needs to be in order for point estimates to be within some distance of the true population mean  $(1 - \alpha)\%$  of the time. Example:

- We would like  $\bar{x}$  to be within 4 of  $\mu$  in 95% of samples
- Put another way, our margin of error will need to be 4
- Assume  $\sigma = 25$
- How large will  $n$  need to be?

## Margins of error and sample size

$$\begin{aligned}1.96 \left( \frac{25}{\sqrt{n}} \right) &= 4 \\ \frac{25}{\sqrt{n}} &= 2.0408 \\ \frac{25}{2.0408} &= \sqrt{n} \\ 12.25 &= \sqrt{n} \\ 150 &= n\end{aligned}$$

## Margins of error and sample size

Notice the minimum required sample size depends on both the significance level (which determines the z-score, in this case 1.96) and the population standard deviation. For a given desired margin of error  $m$  we can write the minimum sample size required as:

$$n = \left( \frac{z\sigma}{m} \right)^2$$

Of course,  $\sigma$  is not known and can't be estimated with  $s$  before collecting data. But one may be able to use other information about the distribution of  $x$  to approximate this value.

## Margins of error and sample size: proportions

The procedure for estimating minimum required sample sizes is similar for proportions. For a given desired margin of error  $m$  (between 0 and 1), we can write the minimum sample size required as:

$$n = \left(\frac{z}{m}\right)^2 \pi(1 - \pi)$$

Again,  $\pi$  is not known before collecting data, but may be approximated given other information. (For example, in an opinion poll, the researcher could use previous polls on the same question to approximate  $\pi$ ). The most conservative estimate for  $\pi$  is 0.50.

## Margins of error and sample size: proportions

**Table:** Minimum sample size required for margin of error  $m$  and  $\pi$

m:	0.05	0.04	0.03	0.02	0.01
$\pi$	Minimum sample size required				
0.30	323	504	896	2,017	8,067
0.35	350	546	971	2,185	8,740
0.40	369	576	1,024	2,305	9,220
0.45	380	594	1,056	2,377	9,508
0.50	384	600	1,067	2,401	9,604
0.55	380	594	1,056	2,377	9,508
0.60	369	576	1,024	2,305	9,220
0.65	350	546	971	2,185	8,740
0.70	323	504	896	2,017	8,067



## Margins of error and sample size: proportions

For a given  $n$ , sampling variability is greatest when  $\pi = 0.5$

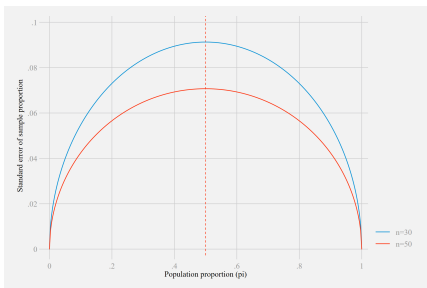


Figure: Standard error of sample proportion as a function of  $\pi$

## Simulations

Let's try simulating  $(1 - \alpha)\%$  confidence intervals in Stata:

- Draw a sample of size  $n$  from a population distribution
- Calculate  $\bar{x}$  and calculate a  $(1 - \alpha)\%$  confidence interval, first assuming we know  $\sigma$ , and then assuming we don't know  $\sigma$  and have to estimate it using  $s$ .
- In what percent of simulations does the confidence interval contain the true population mean in each case?