

## Sampling Distributions

Quiz

1

### Parameters and statistics

The average fuel tank capacity of all cars made by Ford is 14.7 gallons. This value represents a

- a) Parameter because it is an average from all possible cars.
- b) Parameter because it is an average from all Ford cars.
- c) Statistic because it is an average from a sample of all cars.
- d) Statistic because it is an average from a sample of American cars.

### Parameters and statistics

The fraction of all American adults who received at least one speeding ticket last year can be represented by

- a)  $\mu$ .
- b)  $\bar{x}$ .
- c)  $p$ .
- d)  $\hat{p}$ .

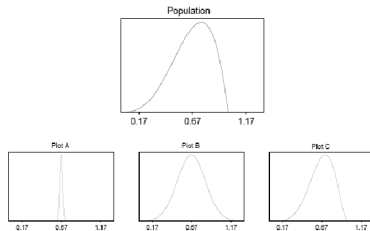
### Sample size

We wish to estimate the mean price,  $\mu$ , of all hotel rooms in Las Vegas. The Convention Bureau of Las Vegas did this in 1999 and used a sample of  $n = 112$  rooms. In order to get a better estimate of  $\mu$  than the 1999 survey, we should

- a) Take a larger sample because the sample mean will be closer to  $\mu$ .
- b) Take a smaller sample since we will be less likely to get outliers.
- c) Take a different sample of the same size since it does not matter what  $n$  is.

### Sampling distribution

If the first graph shows the population, which plot could be the sampling distribution of  $\bar{x}$  if all samples of size  $n = 50$  are drawn?



- a) Plot A
- b) Plot B
- c) Plot C

### Sampling distribution

Which of the following is true?

- a) The shape of the sampling distribution of  $\bar{x}$  is always bell-shaped.
- b) The shape of the sampling distribution of  $\bar{x}$  gets closer to the shape of the population distribution as  $n$  gets large.
- c) The shape of the sampling distribution of  $\bar{x}$  gets approximately normal as  $n$  gets large.
- d) The mean of the sampling distribution of  $\bar{x}$  gets closer to  $\mu$  as  $n$  gets large.

### Sampling distribution

True or false: The shape of the sampling distribution of  $\bar{x}$  becomes more normal the larger your sample size is.

- a) True
- b) False

### Sampling distribution

True or false: The standard deviation of the sampling distribution of  $\bar{x}$  always less than the standard deviation of the population when the sample size is at least 2.

- a) True
- b) False

### Sampling distribution

The theoretical sampling distribution  $\bar{x}$

- a) Gives the values of  $\bar{x}$  from all possible samples of size  $n$  from the same population.
- b) Provides information about the shape, center, and spread of the values in a single sample.
- c) Can only be constructed from the results of a single random sample of size  $n$ .
- d) Is another name for the histogram of values in a random sample of size  $n$ .

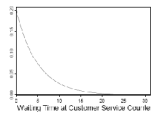
### Standard error

What does  $\frac{\sigma}{\sqrt{n}}$  measure?

- a) The spread of the population.
- b) The spread of the  $\bar{x}$  's.
- c) Different values  $\bar{x}$  could be.

### Sampling distributions

The following density curve represents waiting times at a customer service counter at a national department store. The mean waiting time is 5 minutes with standard deviation 5 minutes. If we took all possible samples of size  $n = 100$ , how would you describe the sampling distribution of the  $\bar{x}$  's?



- a) Shape = right skewed, center = 5, spread = 5
- b) Shape = less right skewed, center = 5, spread = 0.5
- c) Shape = approx. normal, center = 5, spread = 5
- d) Shape = approx. normal, center = 5, spread = 0.5

### Central Limit Theorem

Which is a true statement about the Central Limit Theorem?

- a) We need to take repeated samples in order to estimate  $\mu$ .
- b) It only applies to populations that are Normally distributed.
- c) It says that the distribution of  $\bar{x}$  's will have the same shape as the population.
- d) It requires the condition that the sample size,  $n$ , is large and that the samples were drawn randomly.

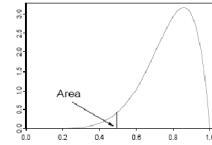
## Central Limit Theorem

True or false: The Central Limit Theorem allows us to compute probabilities on  $\bar{x}$  when the conditions are met.

- a) True
- b) False

## Sampling distributions

The following density curve shows a sampling distribution of  $\bar{x}$  created by taking all possible samples of size  $n = 6$  from a population that was very left-skewed. Which of the following would result in a decrease of the area to the left of 0.5 (denoted by the vertical line)?



- a) Increasing the number of samples taken.
- b) Taking a different sample.
- c) Decreasing  $n$ .
- d) Increasing the number of observations in each sample.

## Sampling distributions

What effect does increasing the sample size,  $n$ , have on the center of the sampling distribution of  $\bar{x}$ ?

- a) The mean of the sampling distribution gets closer to the mean of the population.
- b) The mean of the sampling distribution gets closer to 0.
- c) The variability of the population mean is decreased.
- d) It has no effect. The mean of the sampling distribution always equals the mean of the population.

## Sampling distributions

What effect does increasing the sample size,  $n$ , have on the spread of the sampling distribution of  $\bar{x}$ ?

- a) The spread of the sampling distribution gets closer to the spread of the population.
- b) The spread of the sampling distribution gets larger.
- c) The spread of the sampling distribution gets smaller.
- d) It has no effect. The spread of the sampling distribution always equals the spread of the population.

## Sampling distributions

What effect does increasing the sample size,  $n$ , have on the shape of the sampling distribution of  $\bar{x}$ ?

- a) The shape of the sampling distribution gets closer to the shape of the population.
- b) The shape of the sampling distribution gets more bell-shaped.
- c) It has no effect. The shape of the sampling distribution always equals the shape of the population.

## Sampling distributions

Which of the following would result in a decrease in the spread of the approximate sampling distribution of  $\bar{x}$ ?

- a) Increasing the sample size.
- b) Increasing the number of samples taken.
- c) Increasing the population standard deviation
- d) Decreasing the value of the population mean.

## Sampling distributions

Time spent working out at a local gym is normally distributed with mean  $\mu = 43$  minutes and standard deviation  $\sigma = 6$  minutes.

The gym took a sample of size  $n = 24$  from its patrons. What is the distribution of  $\bar{x}$ ?

- a) Normal with mean  $\mu = 43$  minutes and standard deviation  $\sigma = 6$  minutes.
- b) Normal with mean  $\mu = 43$  minutes and standard deviation  $\sigma = \frac{6}{\sqrt{24}}$  minutes.
- c) Cannot be determined because the sample size is too small.

## Confidence Intervals

Chapter 7

20

## Inference

- We are in the fourth and final part of the course - **statistical inference**, where we draw conclusions about the population based on the data obtained from a sample chosen from it.

21

## Confidence Intervals (CI)

- **The goal:** to give a range of plausible values for the estimate of the unknown population parameter (the population mean,  $\mu$ , or the population proportion,  $p$ )
- We start with our best guess: the sample statistic (the sample mean  $\bar{x}$ , or the sample proportion  $\hat{p}$ )
- Sample statistic = point estimate

22

## Point estimate

- Use a single statistic based on sample data to estimate a population parameter
- Simplest approach
- But not always very precise due to variation in the sampling distribution

23

## Confidence Intervals (CI)

**CI = point estimate  $\pm$  margin of error**



24

## Margin of error

- Shows how accurate we believe our estimate is
- The smaller the margin of error, the more precise our estimate of the true parameter
- Formula:

$$E = \left( \begin{array}{c} \text{critical} \\ \text{value} \end{array} \right) \cdot \left( \begin{array}{c} \text{standard deviation} \\ \text{of the statistic} \end{array} \right)$$

## Confidence Intervals (CI) for a Mean

- Suppose a random sample of size  $n$  is taken from a normal population of values for a quantitative variable whose mean  $\mu$  is unknown, when the population's standard deviation  $\sigma$  is known.
- A confidence interval (CI) for  $\mu$  is:

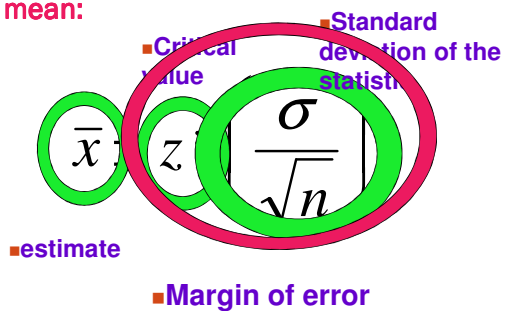
$$\text{CI} = \text{point estimate} \pm \text{margin of error}$$

$$\bar{x} \pm z^* \cdot \frac{\sigma}{\sqrt{n}}$$

Point estimate      Margin of error (m or E)

26

## Confidence interval for a population mean:



## Confidence level

- A confidence interval is associated with a confidence level. The confidence level gives us the success rate of the procedure used to construct the CI. We will say: “the 95% confidence interval for the population mean is ...”
- The most common choices for a confidence level are 90% ( $z^* = 1.645$ ), 95% ( $z^* = 1.96$ ), and 99% ( $z^* = 2.576$ ).

28

## Statement: (memorize!!)

We are \_\_\_\_\_% confident that the true mean context lies within the interval \_\_\_\_\_ and \_\_\_\_\_.

## Using the calculator

- Calculator: STAT → TESTS → 7:Interval...

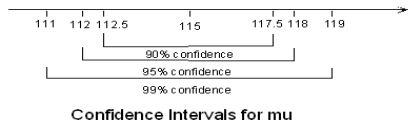
- Inpt: Data    Stats

Use this when you have data in one of your lists

Use this when you know  $\bar{x}$  and  $\sigma$

30

## The Trade-off



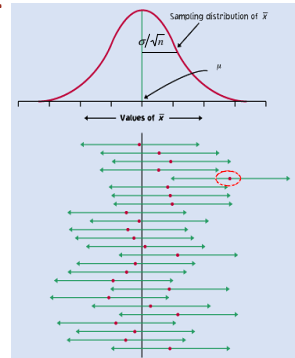
- **There is a trade-off between the level of confidence and precision in which the parameter is estimated.**
- higher level of confidence – wider confidence interval
- lower level of confidence – narrower confidence interval

31

## 95% confident means:

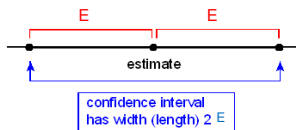
In 95% of all possible samples of this size  $n$ ,  $\mu$  will indeed fall in our confidence interval.

In only 5% of samples would miss  $\mu$ .



## The Margin of Error

- The width (or length) of the CI is exactly twice the margin of error (E):



- The margin of error is therefore "in charge" of the width of the confidence interval.

33

## Comment

- The margin of error (E) is

$$E = z^* \cdot \frac{\sigma}{\sqrt{n}}$$

and since  $n$ , the sample size, appears in the denominator, **increasing  $n$  will reduce the margin of error** for a fixed  $z^*$ .

34

## How can you make the margin of error smaller?

- $z^*$  smaller  
(lower confidence level)
- $\sigma$  smaller  
(less variation in the population)
- $n$  larger  
(to cut the margin of error to be 4 times as big)

Really cannot change!

## Margin of Error and the Sample Size

- In situations where a researcher has some flexibility as to the sample size, the **researcher can calculate in advance what the sample size is that he/she needs in order to be able to report a confidence interval with a certain level of confidence and a certain margin of error.**

36

## Calculating the Sample Size

$$E = z^* \cdot \frac{\sigma}{\sqrt{n}}$$

$$n = \left( z^* \cdot \frac{\sigma}{E} \right)^2$$

Clearly, the sample size  $n$  must be an integer. Calculation may give us a non-integer result. In these cases, we should always **round up to the next highest integer.**

37

## Example

- IQ scores are known to vary normally with standard deviation 15. How many students should be sampled if we want to estimate population mean IQ at 99% confidence with a margin of error equal to 2?

$$n = \left( z^* \frac{\sigma}{E} \right)^2 = \left( 2.576 \frac{15}{2} \right)^2 = 373.26 \Rightarrow n = 374$$

**They should take a sample of 374 students.**

38

## Assumptions for the validity of

- The sample must be random
- The standard deviation,  $\sigma$ , is known
- and either
  - the sample size must be large ( $n \geq 30$ ) or
  - for smaller sample the variable of interest must be normally distributed in the population.

$$\bar{x} \pm z^* \cdot \frac{\sigma}{\sqrt{n}}$$

39

## Steps to follow

1. **Check conditions:** SRS,  $\sigma$  is known, and either  $n \geq 30$  or the population distribution is normal
2. **Calculate** the CI for the given confidence level
3. **Interpret** the CI

40

## Example 1

- A college admissions director wishes to estimate the mean age of all students currently enrolled. In a random sample of 20 students, the mean age is found to be 22.9 years. From past studies, the standard deviation is known to be 1.5 years and the population is normally distributed. Construct a 90% confidence interval of the population mean age.

41

## Step 1: Check conditions

- A college admissions director wishes to estimate the mean age of all students currently enrolled. In a **random sample** of 20 students, the mean age is found to be 22.9 years. From past studies, the **standard deviation is known** to be 1.5 years and the population is **normally distributed**.
- SRS ✓
- $\sigma$  is known ✓
- The population is normally distributed ✓

42

Step 2: Calculate the 90% CI using the formula

$$\bar{x} = 22.9$$

$$\sigma = 1.5$$

$$n = 20$$

$$z^* = 1.645$$

$$\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}} = 22.9 \pm 1.645 \frac{1.5}{\sqrt{20}} = 22.9 \pm 0.6 = (22.3, 23.5)$$

43

Step 2: Calculate the 90% CI using the calculator

• **Calculator:** STAT→TESTS→7:ZInterval...

- Inpt: Data
- $\sigma = 1.5$
- $\bar{x} = 22.9$
- $n = 20$
- C-Level: .90
- Calculate

**ZInterval : (22.3, 23.5)**

44

Step 3: Interpretation

- **We are 90% confident that the mean age of all students at that college is between 22.3 and 23.5 years.**

45

Example 1

- How many students should he ask if he wants the margin of error to be no more than 0.5 years with 99% confidence?

$$n = \left( z^* \cdot \frac{\sigma}{E} \right)^2 = \left( 2.576 \frac{1.5}{0.5} \right)^2 = 59.72$$

- **Thus, he needs to have at least 60 students in his sample.**

46

Example 2



**A scientist wants to know the density of bacteria in a certain solution.** He makes measurements of 10 randomly selected sample:

24, 31, 29, 25, 27, 27, 32, 25, 26, 29 \*10<sup>6</sup> bacteria/ml.

From past studies the scientist knows that the distribution of bacteria level is normally distributed and the population standard deviation is 2\*10<sup>6</sup> bacteria/ml.

a. What is the point estimate of  $\mu$ ?

$$\bar{x} = 27.5 * 10^6 \text{ bacteria/ml.}$$

Example 2

- b. Find the 95% confidence interval for the mean level of bacteria in the solution.

• **Step 1:** check conditions: SRS, normal distribution,  $\sigma$  is known. All satisfied.

• **Step 2:** CI:

$$\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}} = 27.5 \pm 1.96 \frac{2}{\sqrt{10}} = 27.5 \pm 1.24 = (26.26, 28.74)$$

- **Step 3:** Interpret: we are 95% confident that the mean bacteria level in the whole solution is between 26.26 and 28.74 \*10<sup>6</sup> bacteria/ml.

48



## Example 2

Using the calculator:

- Enter the number into one of the lists, say L1
- STAT → TESTS → 7: ZInterval
  - Inpt: Data
  - $\sigma$ : 2
  - List: L1
  - Freq: 1 (it's always 1)
  - C-Level: .95
  - Calculate
- (26.26, 28.74)

49

## Example 2

- c. What is the margin of error?

From part b:

$$\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}} = 27.5 \pm 1.96 \frac{2}{\sqrt{10}} = 27.5 \pm 1.24 = (26.26, 28.74)$$

Thus, the margin of error is

$$E = 1.24 * 10^6 \text{ bacteria/ml.}$$

50

## Example 2

- d. How many measurements should he make to obtain a margin of error of at most  $0.5 * 10^6$  bacteria/ml with a confidence level of 95%?

$$n = \left( z^* \cdot \frac{\sigma}{E} \right)^2 = \left( 1.96 \frac{2 \times 10^6}{0.5 \times 10^6} \right)^2 = 61.4656$$

- Thus, he needs to take 62 measurements.

51

## Assumptions for the validity of $\bar{x} \pm z^* \cdot \frac{\sigma}{\sqrt{n}}$

- The sample must be random
- The standard deviation,  $\sigma$ , is known and **either**
  - The sample size must be large ( $n \geq 30$ ) or
  - For smaller sample the variable of interest must be normally distributed in the population.
- The only situation when we cannot use this confidence interval, then, is when the sample size is small and the variable of interest is not known to have a normal distribution. In that case, other methods called non-parametric methods need to be used.

52

## Example 3



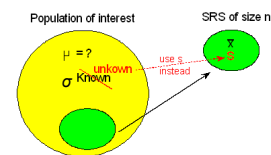
In a randomized comparative

experiment on the effects of calcium on blood pressure, researchers divided 54 healthy, white males at random into two groups, takes calcium or placebo. The paper reports a mean seated systolic blood pressure of 114.9 with standard deviation of 9.3 for the placebo group. Assume systolic blood pressure is normally distributed.

Can you find a z-interval for this problem? Why or why not?

## BUT what if $n < 30$ and $\sigma$ is unknown?

- Well, there is some good news and some bad news!



The good news is that we can easily replace the population standard deviation,  $\sigma$ , with the *sample standard deviation*  $s$ .

54

## And the bad news is...

- that once  $\sigma$  has been replaced by  $s$ , we lose the Central Limit Theorem together with the normality of  $\bar{X}$  and therefore the **confidence multipliers  $z^*$**  for the different levels of confidence are (generally) not accurate any more.
- The new multipliers come from a different distribution called the "***t* distribution**" and are therefore denoted by  **$t^*$**  (instead of  $z^*$ ).

55

## CI for the population mean when $n < 30$ and $\sigma$ is unknown

- The confidence interval for the population mean  $\mu$  when  $n < 30$   **$\sigma$  is unknown** is therefore:

$$\bar{x} \pm t^* \cdot \frac{s}{\sqrt{n}}$$

56

## $z^*$ vs. $t^*$

- There is an important difference between the confidence multipliers we have used so far ( $z^*$ ) and those needed for the case when  $\sigma$  is unknown ( $t^*$ ).
  - $z^*$ , depends only on the level of confidence,
  - $t^*$  depend on both the level of confidence and on the sample size (for example: the  $t^*$  used in a 95% confidence when  $n=10$  is different from the  $t^*$  used when  $n=40$ ).

57

## $t$ -distribution

- There is a different  $t$  distribution for each sample size. We specify a particular  $t$  distribution by giving its degrees of freedom. The degrees of freedom for the one-sample  $t$  statistic come from the sample standard error  $s$  in the denominator of  $t$ . Since  $s$  has  $n-1$  degrees of freedom, the  $t$ -distribution has  $n-1$  degrees of freedom.

58

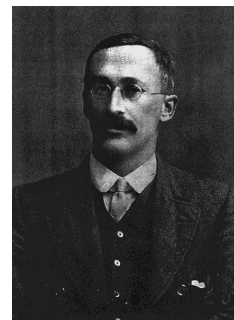
## $t$ -distribution

- The  $t$ -distribution is bell shaped and symmetric about the mean.
- The total area under the  $t$ -curve is 1
- The mean, median, and mode of the  $t$ -distribution are equal to zero.
- The tails in the  $t$ -distribution are "thicker" than those in the standard normal distribution.
- As the df (sample size) increases, the  $t$ -distribution approaches the normal distribution. After 29 df the  $t$ -distribution is very close to the standard normal  $z$ -distribution.

59

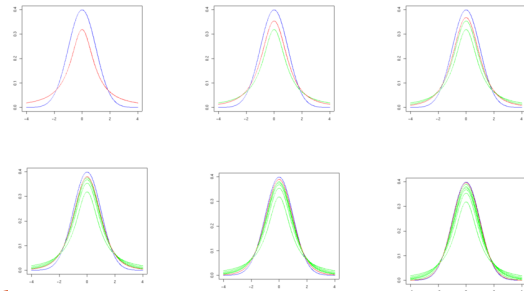
## Historical Reference

- William Gosset (1876-1937) developed the  $t$ -distribution while employed by the Guinness Brewing Company in Dublin, Ireland. Gosset published his findings using the name "Student". The  $t$ -distribution is, therefore, sometimes referred to as "Student's  $t$ -distribution".



60

Density of the t-distribution (red and green) for 1, 2, 3, 5, 10, and 30 df compared to normal distribution (blue)



61

## Calculator

- **Calculator:**
- STAT → TESTS → 8:TInterval...
- Inpt: Data Stats

Use this when  
you have data  
in one of your lists

Use this when  
you know  $\bar{x}$  and  $s$

62

## Example

- To study the metabolism of insects, researchers fed cockroaches measured amounts of a sugar solution. After 2, 5, and 10 hours, they dissected some of the cockroaches and measured the amount of sugar in various tissues. Five roaches fed the sugar solution and dissected after 10 hours had the following amounts of sugar in their hindguts:



63

## Example

- 55.95, 68.24, 52.73, 21.50, 23.78
- Find the 95% CI for the mean amount of sugar in cockroach hindguts:

$$\bar{x} = 44.44 \quad s = 20.741$$

The degrees of freedom,  $df = n - 1 = 4$ , and from the table we find that for the 95% confidence,  $t^* = 2.776$ .

- Then

$$\bar{x} \pm t^* \cdot \frac{s}{\sqrt{n}} = 44.44 \pm 2.776 \cdot \frac{20.741}{\sqrt{5}} = (18.69, 70.19)$$

64

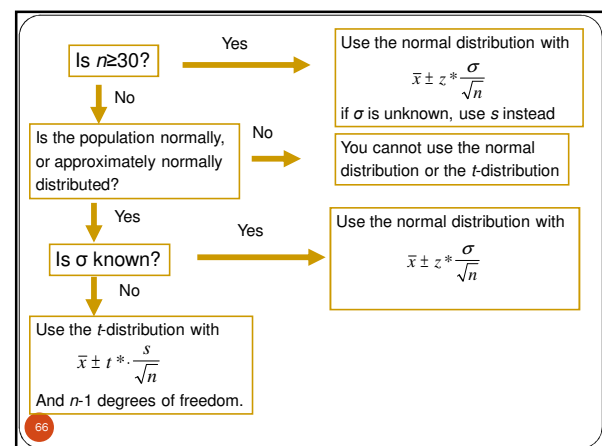
## Example

- The large margin of error is due to the small sample size and the rather large variation among the cockroaches.

### Calculator:

- Put the data in  $L_1$ .
- STAT → TESTS → 8:TInterval...
  - Inpt: Data Stats
  - List:  $L_1$
  - Freq: 1
  - C-level: .95

65



66

### Examples: You take:

- 24 samples, the data are normally distributed,  $\sigma$  is known
  - normal distribution with  $\sigma$   $\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}}$
- 14 samples, the data are normally distributed,  $\sigma$  is unknown
  - t-distribution with  $s$   $\bar{x} \pm t^* \frac{s}{\sqrt{n}}$
- 34 samples, the data are not normally distributed,  $\sigma$  is unknown
  - normal distribution with  $s$   $\bar{x} \pm z^* \frac{s}{\sqrt{n}}$
- 12 samples; the data are not normally distributed,  $\sigma$  is unknown
  - cannot use the normal distribution or the t-distribution

67

### Some Cautions:

- The data **MUST** be a SRS from the population
- The formula is not correct for more complex sampling designs, i.e., stratified, etc.
- No way to correct for bias in data
- Outliers can have a large effect on confidence interval
- Must know  $\sigma$  to do a z-interval – which is unrealistic in practice

### Estimating a Population Proportion

When the variable of interest is **categorical**, the population parameter that we will infer about is a **population proportion ( $p$ )** associated with that variable.

- For example, if we are interested in studying opinions about the death penalty among U.S. adults, and thus our variable of interest is "death penalty (in favor/against)," we'll choose a sample of U.S. adults and use the collected data to make inference about  $p$  - the proportion of US adults who support the death penalty.



69

### Example 2

- Suppose that we are interested in the opinions of U.S. adults regarding legalizing the use of marijuana. In particular, we are interested in the parameter  $p$ , the proportion of U.S. adults who believe marijuana should be legalized.
- Suppose a poll of 1000 U.S. adults finds that 560 of them believe marijuana should be legalized.



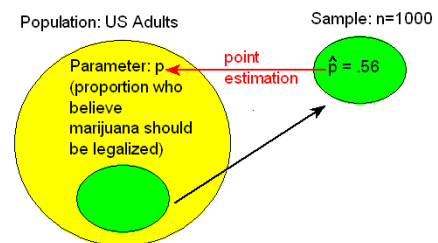
70

### Example 2

- If we wanted to estimate  $p$ , the population proportion by a single number based on the sample, it would make intuitive sense to use the corresponding quantity in the sample, the sample proportion  $\hat{p} = 560/1000 = 0.56$ . We say in this case that 0.56 is the **point estimate** for  $p$ , and that in general, we'll always use  $\hat{p}$  as the **point estimator** for  $p$ .
- Note, again, that when we talk about the *specific value* (.56), we use the term estimate, and when we talk in general about the *statistic* we use the term *estimator*. Here is a visual summary of this example:

71

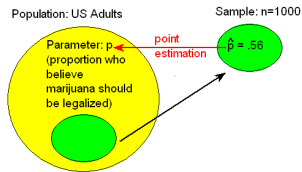
### Example 2



72

## Back to Example 2

- Suppose a poll of 1000 U.S. adults finds that 560 of them believe marijuana should be legalized.



73

## The CI for $p$

- Thus, the **confidence interval for  $p$**  is

$$\hat{p} \pm E = \hat{p} \pm z^* \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

- For a 95% CI use  $z^*=1.96$
- For a 90% CI use  $z^*=1.645$
- For a 99% CI use  $z^*=2.576$

74

## Calculator:

- STAT → TESTS → A: 1-PropZInt...

- $x$  is the number of successes:  $x = n\hat{p}$

75

## Conditions

- The CI is reasonably accurate when three conditions are met:

- The sample was a simple random sample (SRS) from a binomial population
- Both  $n\hat{p} \geq 10$  and  $n(1-\hat{p}) \geq 10$
- The size of the population is at least 10 times the size of the sample

76

## Example

- Suppose you have a random sample of 40 buses from a large city and find that 24 buses have a safety violation. Find the 90% CI for the proportion of all buses that have a safety violation.

- Conditions:

- SRS

- both ✓

$$n\hat{p} = 40\left(\frac{24}{40}\right) = 24 \geq 10 \quad \text{and}$$

$$n(1-\hat{p}) = 40\left(1-\frac{24}{40}\right) = 16 \geq 10 \quad \checkmark$$

- The size of the population (all the buses) is at least 10 times the size of the sample (40) ✓

77

## 90% CI

$$\hat{p} = \frac{24}{40} = 0.6$$

- For 90% CI  $z^*=1.645$

$$\hat{p} \pm E = \hat{p} \pm z^* \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = 0.6 \pm 1.645 \cdot \sqrt{\frac{0.6(1-0.6)}{40}} = 0.6 \pm 0.13 = (0.47, 0.73)$$

78

## Interpretation

- 1. **What is it that you are 90% sure is in the confidence interval?**  
The proportion of all of the buses in this population that have safety violations if we could check them all.
- 2. **What is the meaning (or interpretation) of the confidence interval of 0.47 to 0.73?**  
We are 90% confident that if we could check all of the buses in this population, between 47% and 73% of them would have safety violations.
- 3. **What is the meaning of 90% confidence?**  
If we took 100 random samples of buses from this population and computed the 90% confidence interval from each sample, then we would expect that 90 of these intervals would contain the proportion of all buses in this population that have safety violations. In other words, we are using a method that captures the true population proportion 90% of the time.

79

## Margin of Error and Sample Size

- When we have some level of flexibility in determining the sample size, we can set a desired margin of error for estimating the population proportion and find the sample size that will achieve that.
- For example, a final poll on the day before an election would want the margin of error to be quite small (with a high level of confidence) in order to be able to predict the election results with the most precision. This is particularly relevant when it is a close race between the candidates. The polling company needs to figure out how many eligible voters it needs to include in their sample in order to achieve that.
- Let's see how we do that.

80

## Margin of Error and Sample Size

- The confidence interval for  $p$  is

$$\hat{p} \pm E = \hat{p} \pm z^* \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \quad E = z^* \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

- Thus, the margin of error is

Using some algebra we have

$$n = \left(\frac{z^*}{E}\right)^2 \hat{p}(1-\hat{p})$$

81

$$n = \left(\frac{z^*}{E}\right)^2 \hat{p}(1-\hat{p})$$

- If you have a good estimate  $\hat{p}$  of  $p$ , use it in this formula, otherwise take the conservative approach by setting  $\hat{p} = \frac{1}{2}$ .
- You have to decide on a level of confidence so you know what value of  $z^*$  to use (most common one is the 95% level).
- Also, obviously, you have to set the margin of error (the most common one is 3%).

82

What sample size should we use for a survey if we want a margin of error to be at most 3%?

Let's use the 95% confidence here, so  $z^*=1.96$ . Also, since we don't have an estimate of  $p$ , we will use  $\hat{p} = 0.5$ . Then

$$n = \left(\frac{z^*}{E}\right)^2 \hat{p}(1-\hat{p}) = \left(\frac{1.96}{0.03}\right)^2 (0.5)(1-0.5) = 1067.111$$

Because you must have a sample size of at least 1067.111, **round up** to 1068. So  $n$  should be at least 1068.

83

## Summary: CI for ...

|                                                                          |                                                                                    |
|--------------------------------------------------------------------------|------------------------------------------------------------------------------------|
| a population proportion                                                  | $\hat{p} \pm z^* \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$                              |
| a population mean, $n \geq 30$<br>$\sigma$ is known/ $\sigma$ is unknown | $\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}} \quad \bar{x} \pm z^* \frac{s}{\sqrt{n}}$ |
| a population mean, $\sigma$ is unknown<br>and $n < 30$                   | $\bar{x} \pm t^* \frac{s}{\sqrt{n}}$                                               |

84