

# Math 140

## Introductory Statistics

Professor Silvia Fernández

Chapter 7

Based on the book *Statistics in Action*  
by A. Watkins, R. Scheaffer, and G. Cobb.

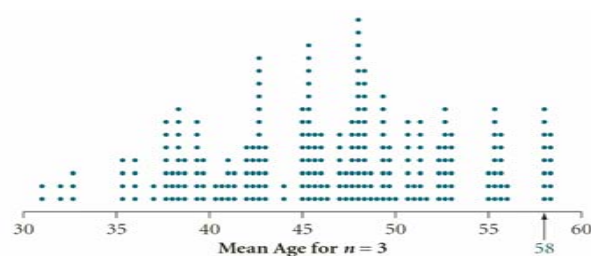
## 7.1 Generating Sampling Distributions

- **Sampling Distributions.** Distribution of summary statistics obtained from taking repeated random samples.
- Steps for generating a sampling distribution:
  - I. Take a random sample of a fixed size  $n$  from a population.
  - II. Compute a Summary Statistic for this sample.
  - III. Repeat steps I and II many times.
  - IV. Display the distribution of the Summary Statistic.

**Note:** A way to remember these steps is, Random Sample, Summary Statistic, Repetition, Distribution.

## Example

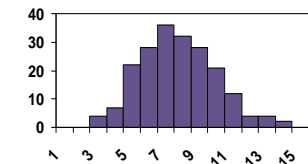
- Westvaco case.  
Randomly select three workers from the group of 10 with ages 25, 33, 35, 38, 48, 55, 55, 55, 56, and 64, and calculate the mean age of the three selected.



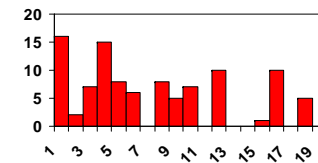
**Display 7.1** A simulated sampling distribution of the mean age from random samples of three people who could have been laid off at Westvaco.

## Shape, Center, and Spread

- A good description of a sampling distribution is the trio shape, center, and spread.
- Recall the rectangles activity 4.2. (See displays 5.7 and 5.8)

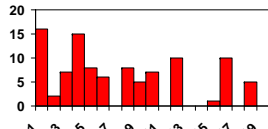


■ Sample mean of the areas of 5 rectangles



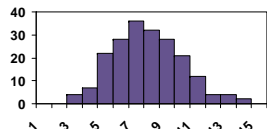
■ Population of Rectangle Areas

# Shape, Center, and Spread



Population of Rectangle Areas

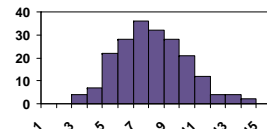
- Shape: Irregular
- Center= Mean =  $\mu = 7.41$
- Spread = Standard Deviation =  $\sigma = 5.23$



Sample mean of the areas of 5 rectangles

- Shape: Normal with a hint of skew to the right.
- Center= Mean=  $\bar{x} = 7.377$
- Spread = Standard Deviation =  $SE = 2.23$

# Shape, Center, and Spread



Sample mean of the areas of 5 rectangles

- Shape: Normal with a hint of skew to the right.
- Center= Mean=  $\bar{x} = 7.377$
- Spread = Standard Deviation =  $SE = 2.23$

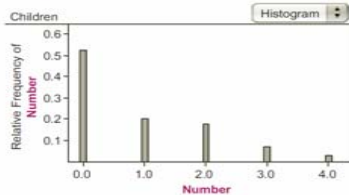
### Notes.

- The standard deviation of the sampling distribution is often called the Standard Error (SE)
- Most sample distributions are nearly normal, we'll see more about this later.
- Values that are in the middle 95% of a random distribution are called **Reasonably Likely**.
- Values that are in the outer 5% of a random distribution are called **Rare Events**.

# 7.2 Sampling distribution of the sample mean.

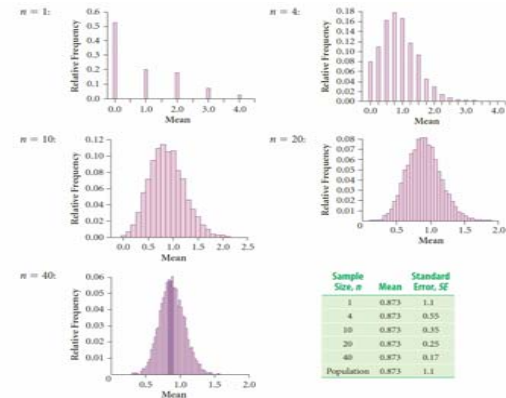
Example: Population

Number of Children	Proportion of Families
0	0.524
1	0.201
2	0.179
3	0.070
4 (or more)	0.026



Display 7.24 The number of children per family in the United States. [Source: U.S. Census Bureau, Statistical Abstract of the United States, 2004-2005, www.census.gov.]

# Sampling distribution of the sample mean for sample sizes 1, 4, 10, 20, and 40.



Display 7.25 Sampling distributions of the sample mean for samples of size 1, 4, 10, 20, and 40.

## Notation

	Population	Sample	Sampling Distribution
Mean	$\mu$	$\bar{x}$	$\mu_{\bar{x}}$
Standard Deviation	$\sigma$	$s$	$\sigma_{\bar{x}}$ or $SE$
Size	$N$	$n$	

## Properties of The Sampling Distribution of The Sample Mean

- The mean  $\mu_{\bar{x}}$  of the sampling distribution of  $\bar{x}$  equals the mean of the population  $\mu$ .

$$\mu_{\bar{x}} = \mu$$

- \*The standard deviation  $\sigma_{\bar{x}}$  of the sampling distribution of  $\bar{x}$ , also called the **standard error** of the mean, equals the standard deviation of the population  $\sigma$  divided by the square root of the sample size  $n$ :

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

- The Shape of the sampling distribution will be approximately normal if the population is approximately normal; for other populations, the sampling distribution becomes more normal as  $n$  increases. This property is called the **Central Limit Theorem**.

\*This holds as long as you sample with replacement or your sample size is less than 10% of the population size. (See exercise E30.)

## Example 1

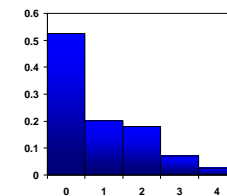
- Problems usually involve a combination of the three properties of the Sampling Distribution of the Sample Mean, together with what we learned about the normal distribution.
- Example: **Average Number of Children**  
What is the probability that a random sample of 20 families in the United States will have an average of 1.5 children or fewer?

## Example 1

- Example: **Average Number of Children**  
What is the probability that a random sample of 20 families in the United States will have an average of 1.5 children or fewer?

Number of Children (per family), $x$	Proportion of families, $P(x)$
0	0.524
1	0.201
2	0.179
3	0.070
4 or more	0.026

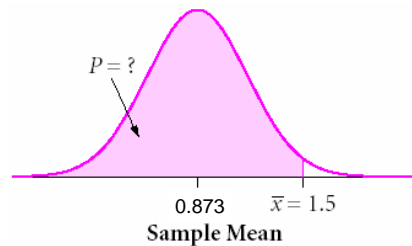
- Mean (of population)  
 $\mu = 0.873$
- Standard Deviation  
 $\sigma = 1.095$



## Example 1

$$\mu_{\bar{x}} = \mu = 0.873$$

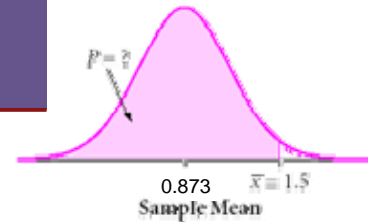
$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{1.095}{\sqrt{20}} = 0.2448$$



## Example 1

$$\mu_{\bar{x}} = \mu = 0.873$$

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{1.095}{\sqrt{20}} = 0.2448$$



- Find z-score of the value 1.5

$$z = \frac{\bar{x} - \text{mean}}{SD} = \frac{\bar{x} - \mu_{\bar{x}}}{\sigma_{\bar{x}}}$$

$$= \frac{1.5 - 0.873}{0.2448} \approx 2.56$$

$$\text{normalcdf}(-99999, 2.56) \approx .9948$$

- So in a random sample of 20 families there is a 99.47% probability that the mean number of children per family will be less than 1.5

OR  $\text{normalcdf}(-99999, 1.5, 0.873, 0.2448) \approx .9948$

## Example 2

- Example: **Reasonably Likely Averages**

What average numbers of children are reasonably likely in a random sample of 20 families?

- Recall that the values that are in the middle 95% of a random distribution are called **Reasonably Likely**.

Note that by calculating the z-scores of 2.5% and 97.5% we find that the **Reasonably Likely** values are those values within 1.96 standard deviations from the mean.

That is, between  $\mu - 1.96\sigma$  and  $\mu + 1.96\sigma$

## Finding Probabilities for Sample Totals

- Sometimes situations are stated in terms of the total number in the sample rather than the average number: "What is the probability that there are 30 or fewer children in a random sample of 20 families in the United States?" You have the choice of two equivalent ways to do this problem.
- Method I:** Find the equivalent average number of children,  $\bar{x}$ , by dividing the total number of children, 30, by the sample size, 20:

$$\bar{x} = \frac{30}{20} = 1.5$$

Then you can use the same formulas and procedure as in the previous examples.

- Method II:** Convert the formulas from the previous examples to equivalent formulas for the sum, then proceed as in the next example.

## Sampling Distribution of the Sum of a Sample

- If a random sample of size  $n$  is selected with mean  $\mu$  and standard deviation  $\sigma$ , then
  - the mean of the sampling distribution of the sum is

$$\mu_{sum} = n\mu$$

- the standard error of the sampling distribution of the sum is

$$\sigma_{sum} = \sqrt{n} \cdot \sigma$$

- the shape of the sampling distribution will be approximately normal if the population is approximately normal; for other populations, the sampling distribution becomes more normal as  $n$  increases.

Note: To get the "sum" formulas just multiply by  $n$

## Examples 3 and 4

- **Ex3: The Probability of 25 or fewer Children**  
What is the probability that a random sample of 20 families in the United States will have a total of 25 children or fewer?
- **Ex4: Reasonably Likely Totals**  
In a random sample of 20 families, what total numbers of children are reasonably likely?

## Sample Size vs. Population Size

- As long as the sample size is a small percentage (around 10% or less) of the population size, it doesn't matter much if you sample with or without replacement, and, in fact, the population size will have little effect on the statistical analysis.
- If the sample size is more than 10% of the population size then a more complex formula needs to be used for the standard error, we will not do this here since it rarely happens in practice. (See Exercise E30.)

## 7.3 Sampling Distribution of the Sample Proportion

- You often hear reports of percentages or proportions: About 60% of automobile drivers in Mississippi use seat belts. (The national average is about 82%.)
- To make intelligent decisions based on data that is reported this way, you must understand the behavior of proportions that arise from random samples.
- The properties of sample proportions are similar to the properties of sample means.

## The Sample Proportion p-hat

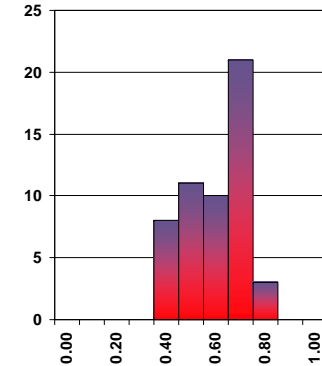
- In a certain population we say that  $p$  is the proportion of the population having a certain property. We say that  $p$  is the proportion of “success” according to our property. (e.g. using a seat belt)
- Note that  $p$  is always a number between 0 and 1.
- When we select a sample of size  $n$ , we calculate the proportion of successes in our sample by dividing the number of successes by the sample size. We call this **the sample proportion** and we denote it by p-hat.

$$\hat{p} = \frac{\text{number of successes}}{\text{sample size}} = \frac{\text{number of successes}}{n}$$

## Simulation (Activity 7.3a)

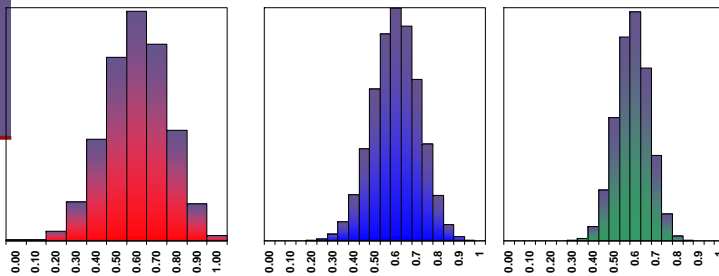
- Choose 10 random numbers between 1 and 10. (use your calculator or a random row in the Table D on page 828)
- Count the number of successes the following way: Numbers between 1 and 6 are successes, 7 to 10 (or 0) are not.
- Calculate

$$\hat{p} = \frac{\text{number of successes}}{10}$$



## Computer Simulations

- The following diagram shows the exact sampling distributions of the sample proportion for samples of size 10, 20, and 40; when  $p = 0.6$



## Center and Spread for Sample Proportions

- We can assign successes as follows:
  - Non-user of seatbelt 0
  - User of seat-belt 1
- We then have the following relative frequency table:

Use Seat Belts	Relative Frequency	In general
0	0.4	$1 - p$
1	0.6	$p$

- And by adding the ones we get

$$\hat{p} = \frac{\text{sum of values}}{\text{sample size}} = \bar{x}$$

- Then we can calculate the mean of the population as follows

$$\mu = \sum x \cdot P(x) = 0(.4) + 1(0.6)$$

- And in general

$$\mu = \sum x \cdot P(x) = 0(1 - p) + 1(p) = p$$

- On the other hand we know that the mean of the sampling distribution of the sample mean is equal to the mean of the population, that is

$$\mu_{\hat{p}} = \mu_{\bar{x}} = \mu = p$$

## Center and Spread for Sample Proportions

Use Seat Belts	Relative Frequency	In general
0	0.4	$1-p$
1	0.6	$p$

- Similarly we can calculate the standard deviation of the population as follows

$$\begin{aligned}\sigma &= \sqrt{\sum (x - \mu)^2 P(x)} = \\ &= \sqrt{(0 - 0.6)^2 0.4 + (1 - 0.6)^2 0.6} \\ &= \sqrt{(0.6)^2 0.4 + (0.4)^2 0.6} \\ &= \sqrt{(0.6)(0.4)(0.6 + 0.4)} = \sqrt{(0.6)(0.4)}\end{aligned}$$

and in general  $\sigma = \sqrt{\sum (x - \mu)^2 P(x)} =$

$$\begin{aligned}&= \sqrt{(0 - p)^2 (1 - p) + (1 - p)^2 p} \\ &= \sqrt{p(1 - p)}\end{aligned}$$

## Center and Spread for Sample Proportions

- On the other hand we know that the standard deviation of the sampling distribution of the sample mean is equal to the SD of the population divided by the square root of the sample size, that is

$$\begin{aligned}\sigma_{\hat{p}} &= \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{\sqrt{p(1-p)}}{\sqrt{n}} \\ \sigma_{\hat{p}} &= \sqrt{\frac{p(1-p)}{n}}\end{aligned}$$

## Properties of The Sampling Distribution of The Sample Proportion

- The mean  $\mu_{\hat{p}}$  of the sampling distribution of  $\hat{p}$  equals the proportion of successes  $p$ :

$$\mu_{\hat{p}} = p$$

- The standard deviation  $\sigma_{\hat{p}}$  of the sampling distribution of  $\hat{p}$ , equals the standard deviation of the population divided by the square root of the sample size  $n$ :

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$$

- As the sample size gets larger, the shape of the sampling distribution gets more normal and will be approximately normal if  $n$  is large enough.
- As a guideline, if both  $np$  and  $n(1-p)$  are at least 10, then using the normal distribution as an approximation for the shape of the sampling distribution will give reasonably accurate results.

## Example 5

- Drivers in the Northeast and Mid-Atlantic states had the highest failure rate, 20%, on the GMAC Insurance National Driver's Test. (They also were the drivers most likely to speed.) [Source: Insurance Journal, www.insurancejournal.com.]
- Describe the shape, center, and spread of the sampling distribution of the proportion of drivers who would fail the test in a random sample of 60 drivers from these states.
- What are the reasonably likely proportions of drivers who would fail the test?

## More Examples

- **Example 6.** Calculate  $\sigma_{\hat{p}}$  with  $p = 0.8$  and  $n = 100, 200, 400$ .
- **Example 7. Using the Properties to Find Probabilities**  
About 60% of Mississippians use seat belts. Suppose your class conducts a survey of 40 randomly selected Mississippians.
  - a. What is the chance that 75% or more of those selected wear seat belts?
  - b. Would it be quite unusual to find fewer than 25% of the Mississippians selected wear seat belts?

## Finding Probabilities for the Number of Successes

- Same as with the sample mean, sometimes problems are stated in terms of the number of successes rather than the proportion of successes. In that case we can again use either of two methods. If we use the second method we need to know about the distribution of the sum of the successes.
- If a random sample of size  $n$  is selected from a population with proportion of success  $p$ , then the sampling distribution satisfies:

$$\mu_{sum} = n\mu_{\hat{p}} = np$$

$$\sigma_{sum} = n \cdot \sigma_{\hat{p}} = \sqrt{np(1-p)}$$

- the shape of the sampling distribution will be approximately normal if if both  $np$  and  $n(1-p)$  are at least 10

Note: To get the "sum" formulas just multiply by  $n$

## Example 8

- **Probability of 30 or More Wearing Seat Belts**  
In a random sample of 40 Mississippians, what is the probability that 25 or more use seat belts?

## E35

- The ethnicity of about 92% of the population of China is Han Chinese. Suppose you take a random sample of 1000 Chinese. [Source: CIA World Factbook.]
  - a. Make an accurate sketch, with a scale on the horizontal axis, of the sampling distribution of the proportion of Han Chinese in your sample.
  - b. Make an accurate sketch, with a scale on the horizontal axis, of the sampling distribution of the number of Han Chinese in your sample.
  - c. What is the probability of getting 90% or fewer Han Chinese in your sample?
  - d. What is the probability of getting 925 or more Han Chinese?
  - e. What numbers of Han Chinese would be rare events? What proportions?

## Summary: Sampling distributions

Population

$\mu$  = mean  
 $\sigma$  = standard  
deviation  
 $N$  = Size

Sample

$\bar{x}$  = mean  
 $s$  = standard  
deviation  
 $n$  = size  
 $p$  = probability  
of success

Sample mean

$$\mu_{\bar{x}} = \mu$$
$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

Sample sum (total)

$$\mu_{sum} = n\mu$$
$$\sigma_{sum} = \sqrt{n} \cdot \sigma$$

Sample proportion

$$\mu_p = p$$
$$\sigma_p = \sqrt{\frac{p(1-p)}{n}}$$

Sample number  
of successes

$$\mu_{sum} = np$$
$$\sigma_{sum} = \sqrt{np(1-p)}$$