

Today: Finish Chapter 9 (Sections 9.6 to 9.8 and 9.9 Lesson 3)

**ANNOUNCEMENTS:**

- Quiz #7 begins after class today, ends Monday at 3pm.
- Quiz #8 will begin next Friday and end at **10am** Monday (day of final). There will be clicker questions in all lectures next week.
- The last homework assignment (#8) will be from the lectures on Mon and Wed and will be due next Friday. Problems to be assigned next Friday are already on the web, with solutions, and are to help you review that material (not to hand in).
- Review for the final exam is posted and will be covered in discussion sections next Friday (*not* in class). Two files:
  - Material since 2<sup>nd</sup> midterm
  - Concepts from the quarter that need extra review

**HOMEWORK:** (Due Monday, March 11)

Chapter 9: #68, 72, 146

**Recall, general format for all sampling distributions in Ch. 9:**

Assuming sample size conditions are met, the sampling distribution of the **sample statistic**:

- Is approximately normal
- Mean = **population parameter** ( $p$ ,  $p_1 - p_2$ ,  $\mu$ , etc.)
- Standard deviation = *standard deviation of* \_\_\_\_\_; the blank is filled in with the **statistic** ( $\hat{p}$ ,  $\hat{p}_1 - \hat{p}_2$ ,  $\bar{x}$  etc.)
- Often the standard deviation must be estimated, and then it is called the *standard error of* \_\_\_\_\_.

See summary table on page 353 for all details!

Update on the five situations we will cover for the rest of this quarter:

Parameter name and description	Population parameter	Sample statistic
<b>For Categorical Variables: [Done!]</b>		
One population proportion (or probability)	$p$	$\hat{p}$
Difference in two population proportions	$p_1 - p_2$	$\hat{p}_1 - \hat{p}_2$
<b>For Quantitative Variables: [Today, M, W]</b>		
One population mean	$\mu$	$\bar{x}$
Population mean of paired differences (dependent samples, paired)	$\mu_d$	$\bar{d}$
Difference in two population means (independent samples)	$\mu_1 - \mu_2$	$\bar{x}_1 - \bar{x}_2$

For each situation will we:

- Learn about the **sampling distribution** for the sample statistic
- Learn how to find a **confidence interval** for the true value of the parameter
- **Test hypotheses** about the true value of the parameter

Today: Sampling distributions for means of quantitative data:

- one mean
- mean difference for paired data
- difference between means for independent samples

Remember, two samples are called **independent samples** when the measurements in one sample are not related to the measurements in the other sample. Could come from:

- Separate samples
- One sample, divided into two groups by a categorical variable (such as male or female)
- Randomization into two groups where each unit goes into only one group

**Paired data** occur when two measurements are taken on the same individuals, or individuals are paired in some way.

## Sampling Distribution for a Sample Mean (Section 9.6)

Suppose we take a random sample of size  $n$  from a population and measure a quantitative variable.

Notation for **Population** (uses Greek letters):

$\mu$  = mean of the **population** of measurements.

$\sigma$  = standard deviation of the **population** of measurements.

Notation for **Sample**:

$\bar{x}$  = **sample** mean of a random **sample** of  $n$  individuals.

$s$  = **sample** standard deviation of the random **sample**

The sampling distribution of the **sample mean**  $\bar{x}$  is:

- approximately normal

- Mean = **population parameter** =  $\mu$

- Standard deviation = *standard deviation of  $\bar{x}$*  =

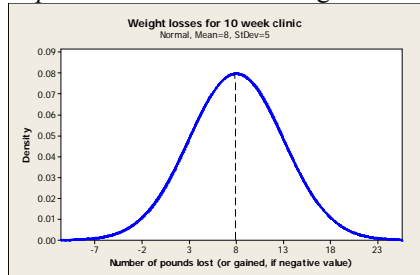
$$s.d.(\bar{x}) = \frac{\sigma}{\sqrt{n}}$$

- Often the **standard deviation** of  $\bar{x}$  must be estimated, and then it is called the **standard error of  $\bar{x}$** . Replace **population  $\sigma$**  with **sample standard deviation  $s$** , so

$$s.e.(\bar{x}) = \frac{s}{\sqrt{n}}$$

Consider the **mean weight loss** for the **population** of people who attend weight loss clinics for 10 weeks. *Suppose* the **population of individual** weight losses is approximately normal,  $\mu = 8$  pounds,  $\sigma = 5$  pounds. (Empirical rule: see picture)

Population of individual weight losses



- We plan to take a random sample of 25 people from this population and record weight loss for each person, then find **sample mean  $\bar{x}$** .

- We know the value of the **sample mean** will vary for different samples of  $n = 25$ . How much will they vary? Where is the center of the distribution of possibilities?

Results for four possible random samples of 25 people, with the corresponding **sample mean  $\bar{x}$**  and **sample standard deviation  $s$** :

**Sample 1:**  $\bar{x} = 8.32$  pounds,  $s = 4.74$  pounds.

**Sample 2:**  $\bar{x} = 6.76$  pounds,  $s = 4.73$  pounds.

**Sample 3:**  $\bar{x} = 8.48$  pounds,  $s = 5.27$  pounds.

**Sample 4:**  $\bar{x} = 7.16$  pounds,  $s = 5.93$  pounds.

**Note:**

- Each sample had a different **sample mean**, which did not always match the **population mean of 8 pounds**.
- Although we cannot determine whether one sample mean will accurately reflect the **population mean**, statisticians have determined what to expect for *all possible sample means*.

$\mu$  = mean of population of interest = 8 pounds

$\sigma$  = standard deviation of population of interest = 5 pounds.

$\bar{x}$  = **sample mean** of a random sample of  $n$  individuals.

Then the *sampling distribution of  $\bar{x}$*  is approximately normal, with

- Mean =  $\mu$
- Standard deviation =  $s.d.(\bar{x}) = \frac{\sigma}{\sqrt{n}}$

Example: **Mean of 25 weight losses**, the distribution of possible values is approximately normal with:

- mean = 8 pounds
- standard deviation =  $\frac{\sigma}{\sqrt{n}} = \frac{5}{\sqrt{25}} = 1$  pound
- What if  $n = 100$  instead of 25?

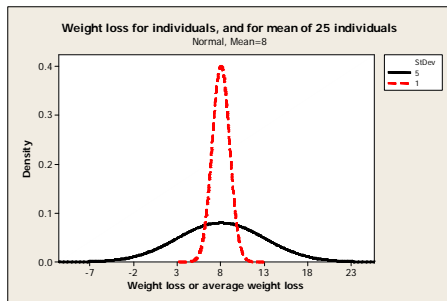
Compare: *individual* weight loss,  $\bar{x}$  for  $n = 25$ ,  $\bar{x}$  for  $n = 100$

	Individuals (wt loss)	Mean of 25	Mean of 100
Mean	8 pounds	8 pounds	8 pounds
St. Dev.	5 pounds	1 pound	½ pound

**Conditions for sampling distribution of  $\bar{x}$  to be approximately normal:**

- Population (individual values) are approx. bell-shaped OR
- Sample size is *large* (at least 30, more if outliers)

Comparing original population with sampling distribution of  $\bar{x}$  :

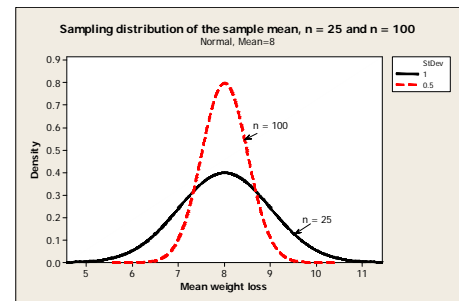


From the empirical rule:

	68%	95%	99.7%
Individuals	3 to 13 lbs	-2 to 18 lbs	-7 to 23 lbs
Mean of $n = 25$	7 to 9 lbs	6 to 10 lbs	5 to 11 lbs

Note that *larger* sample size will result in *smaller* s.d. ( $\bar{x}$ )

Compare sampling distribution of  $\bar{x}$  for  $n = 25$  and  $n = 100$ :



In other words, for *larger* samples, the sample mean  $\bar{x}$  will be closer to  $\mu$  in general, and thus will be a better estimate for  $\mu$ .

**Example where the original population is *not* bell-shaped:**

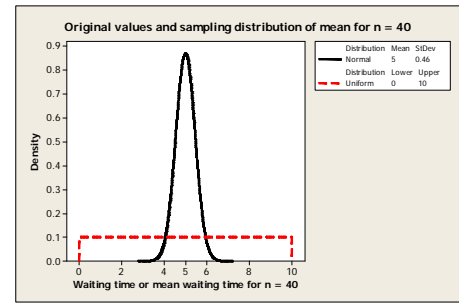
A bus runs every 10 minutes. When you show up at the bus stop, it could come immediately, or any time up to 10 minutes. So the time you wait for it is *uniform*, from 0 to 10 minutes, and independent from day to day.

Population mean =  $\mu = 5$  minutes, population s.d. =  $\sigma = \sqrt{\frac{10^2}{12}} = 2.9$

What is the sampling distribution of  $\bar{x}$  for  $n = 40$  days?

Even though the *original times* are *uniform* (flat shape), the possible values of the sample mean  $\bar{x}$  are:

- Approximately normal
- Mean of  $\bar{x} = 5$  minutes
- Standard deviation of  $\bar{x} = \frac{\sigma}{\sqrt{n}} = \frac{2.9}{\sqrt{40}} = 0.46$  minutes



**Examples of possible samples:**

6.9	6.2	8.8	8.3	7.1	6.5	7.3	9.5	3.4	9.9	5.8	1.4	3	4.1	4.4	0.6	$\bar{x} = 5.29$
2.7	1.2	7.4	0.7	6.8	7.7	6.2	6.1	3.3	1	5.3	9.4	1	0.8	1	9.4	
8.1	3.9	7.2	8.6	1.1	0.4	9.9	9.2									
0.6	3.2	0.8	0.2	8.5	1.4	4.7	0.5	9.7	8.9	6.3	3.3	0.8	4.1	2.6	3.7	$\bar{x} = 4.3$
5.7	3.2	8.9	2.3	1.1	9.9	3	0.8	7.9	0.8	5.9	2.5	7.9	7.6	4	2.2	
0.6	0.1	6.1	6.9	8.1	2.6	9.6	5.3									

Sections 9.7 and 9.8: Sampling distributions for *mean of paired differences*, and for *differences in means for independent samples*.  
Need to learn to distinguish between these two situations.

Notation for **paired differences**:

- $d_i$  = difference in the two measurements for individual  $i = 1, 2, \dots, n$
- $\mu_d$  = mean of the *population of differences*, if all possible pairs were to be measured
- $\sigma_d$  = the *standard deviation of the population of differences*
- $\bar{d}$  = the *mean of the sample of differences*
- $s_d$  = the *standard deviation of the sample of differences*

Example: IQ measured after listening to Mozart and to silence  
 $d_i$  = difference in IQ for student  $i$  for the two conditions  
 $\mu_d$  = *population mean difference*, if all students measured (unknown)  
 $\bar{d}$  = the *mean of the sample of differences* = 9 IQ points  
Based on *sample*, we want to *estimate mean population difference*

**Notation for difference in means for independent samples:**

$\mu_1$  = *population mean* of the first population  
 $\mu_2$  = *population mean* of the second population  
Parameter of interest is  $\mu_1 - \mu_2$  = the *difference in population means*

$\bar{x}_1$  = *sample mean* of the sample from the first population  
 $\bar{x}_2$  = *sample mean* of the sample from the second population  
The *sample statistic* is  $\bar{x}_1 - \bar{x}_2$  = the difference in *sample means*

$\sigma_1$  = *population standard deviation* of the first population  
 $\sigma_2$  = *population standard deviation* of the second population

$s_1$  = *sample standard deviation* of the sample from the 1<sup>st</sup> population  
 $s_2$  = *sample standard deviation* of the sample from the 2<sup>nd</sup> population

$n_1$  = size of the sample from the 1<sup>st</sup> population  
 $n_2$  = size of the sample from the 2<sup>nd</sup> population

Examples where **paired data** might be used:

- Estimate average difference in income for husbands and wives
- Compare SAT scores before and after a training program
- Difference between what you earned in 2012 and what you hope to have as your starting salary when you graduate.

Note that **paired differences** are similar to the “one mean” situation, except special notation tells us that the means are of the differences.

Examples where **independent samples** might be used:

- Compare hours of study for men and women in our class.
- Compare number of sick days off from work for people who had a flu shot and people who didn't
- Compare change in blood pressure for people randomly assigned to a meditation program or to an exercise program for 3 months.

Conditions for the sampling distributions for these two situations are the same as for a single mean, with a slight twist:

- For **paired differences**, population of *differences* must be bell-shaped OR sample must be large.
- For **difference in means for independent samples**, *both* populations must be bell-shaped OR *both* sample sizes must be large.

In both cases, the sampling distribution of the sample statistic is approximately normal, with mean = **population parameter** of interest.

For **paired differences**:

$$\text{mean} = \mu_d, \quad \text{s.d.}(\bar{d}) = \frac{\sigma_d}{\sqrt{n}} \quad (\text{same as one mean, but with } d\text{'s})$$

For **difference in two means**:

$$\text{mean} = \mu_1 - \mu_2, \quad \text{s.d.}(\bar{x}_1 - \bar{x}_2) = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

### Standardized Statistics:

For all 5 cases in Chapter 9, as long as the conditions are satisfied for the sampling distribution to be approximately normal, the **standardized statistic** for a sample statistic is:

$$z = \frac{\text{sample statistic} - \text{population parameter}}{\text{s.d.}(\text{sample statistic})}$$

Note that the denominator has s.d., *not* s.e.

**For one mean:**

$$z = \frac{\bar{x} - \mu}{\text{s.d.}(\bar{x})} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{\sqrt{n}(\bar{x} - \mu)}{\sigma}$$

### Ex: Stat 7, Winter 2011, Hours of study per week for the class

Speculation over the long run is that students study an average of about 5 hours per week for Statistics 7. Have data from 264 students in Winter 2011, with **sample mean** of 5.36 hours. *Suppose* for the population of all possible Statistics 7 students (not just in Winter 2011), **population mean** =  $\mu = 5$  hours a week, and **population standard deviation**  $\sigma = 4$  hours a week (they are *not* bell-shaped – definitely skewed to the right).

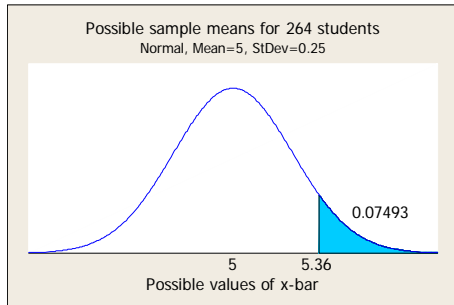
For survey,  $n = 264$ . What are possible values of  $\bar{x}$  (sampling distribution of  $\bar{x}$ )?

- Approximately normal
- Mean = 5 hours
- $\text{s.d.}(\bar{x}) = \frac{4}{\sqrt{264}} = 0.25$

In our **sample**,  $n = 264$  and  $\bar{x} = 5.36$

Standardized statistic for 5.36 hours is  $z = \frac{5.36 - 5}{0.25} = \frac{.36}{.25} = 1.44$

If **population mean** is really 5 hours, with  $\sigma = 4$  hours, how unlikely is a **sample mean** of 5.36 hours or more for  $n = 264$ ?



**How to compute this answer if given  $\mu$ ,  $\sigma$ ,  $n$ , and  $\bar{x}$ :**

Sample means for  $n = 264$  are:

- approximately normal
- with mean of 5 hours and
- s.d. of 4 hours.
- So, the standardized score for 5.36 is:
- 

$$z = \frac{\sqrt{n}(\bar{x} - \mu)}{\sigma} = \frac{\sqrt{264}(5.36 - 5)}{4} = 1.44$$

Area above z-score of 1.44 is about .075. So it is *feasible* that the true population mean for all Stat 7 students (not just Winter 2011) is indeed 5 hours.

### Unknown Population Standard Deviation

When  $\sigma$  is *not* known, we must use the **sample** standard deviation  $s$  instead.

Standard deviation of  $\bar{x}$  :

$$s.d.(\bar{x}) = \frac{\sigma}{\sqrt{n}}$$

Standard error of  $\bar{x}$  :

$$s.e.(\bar{x}) = \frac{s}{\sqrt{n}}$$

### Major consequence:

When using standard error in situations involving means, *standardized statistic* has a *t-distribution* instead of a *z-distribution*; also called **Student's *t* distribution**.

### Student's *t* distribution



In 1908 William Sealy Gossett figured out the formula for the *t* distribution. Called Student's *t* because... explained in class!

## Standardized Statistic Using Standard Error

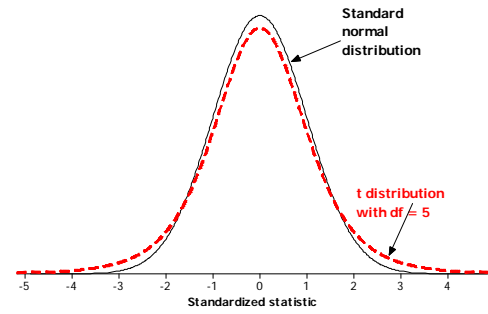
Usually we don't know  $\sigma$  (population standard deviation), so we need to use  $s$  (sample standard deviation). In that case, the standardized statistic for  $\bar{x}$  is

$$t = \frac{\bar{x} - \mu}{s.e.(\bar{x})} = \frac{\bar{x} - \mu}{s / \sqrt{n}} = \frac{\sqrt{n}(\bar{x} - \mu)}{s}$$

This has a *Student's t distribution* with degrees of freedom =  $n - 1$

- It looks almost exactly like the normal distribution
- It is completely specified by knowing the “df”
- It gets closer and closer to the normal distribution, and when degrees of freedom = infinity, it is exactly the normal distribution.

## Comparison of $t$ distribution with $df = 5$ and standard normal distribution



For example, middle 95% for  $t$  with  $df = 5$  is  $-2.57$  to  $+2.57$

For standard normal, it is about  $-2$  to  $+2$

In Chapter 11 (Monday) we will learn how to find probabilities.

Summary of sampling distributions for the 5 parameters (p. 353):

- The **statistic** has a sampling distribution.
- It is *approximately normal* if the sample(s) is (are) large enough.
- The *mean of the sampling distribution* = the **parameter**.
- The **standard deviation** of the sampling distribution is in the table below, in the column “standard deviation of the statistic.”
- Sometimes it needs to be estimated, then “**standard error**” is used.

	Parameter	Statistic	Standard Deviation of the Statistic	Standard Error of the Statistic	Standardized Statistic with s.e.
One proportion	$p$	$\hat{p}$	$\sqrt{\frac{p(1-p)}{n}}$	$\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$	$z$
Difference Between Proportions	$p_1 - p_2$	$\hat{p}_1 - \hat{p}_2$	$\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$	$\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$	$z$
One Mean	$\mu$	$\bar{x}$	$\frac{\sigma}{\sqrt{n}}$	$\frac{s}{\sqrt{n}}$	$t$
Mean Difference, Paired Data	$\mu_d$	$\bar{d}$	$\frac{\sigma_d}{\sqrt{n}}$	$\frac{s_d}{\sqrt{n}}$	$t$
Difference Between Means	$\mu_1 - \mu_2$	$\bar{x}_1 - \bar{x}_2$	$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$	$\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$	$t$

The three situations involving means:

	Parameter	Statistic	Standard Deviation of the Statistic	Standard Error of the Statistic	z or t? (with s.e.)
One Mean	$\mu$	$\bar{x}$	$\frac{\sigma}{\sqrt{n}}$	$\frac{s}{\sqrt{n}}$	$t$
Mean Difference, Paired Data	$\mu_d$	$\bar{d}$	$\frac{\sigma_d}{\sqrt{n}}$	$\frac{s_d}{\sqrt{n}}$	$t$
Difference Between Means	$\mu_1 - \mu_2$	$\bar{x}_1 - \bar{x}_2$	$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$	$\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$	$t$