



#### Inference With Variance Unknown...

Previously, we looked at estimating and testing the population mean when the population standard deviation ( $\sigma$ ) was known or given:  $\overline{x} = u$ 

$$z = \frac{\overline{x} - \mu}{\sigma / \sqrt{n}}$$

But how often do we know the actual population variance?

Instead, we use the *Student t-statistic*, given by:

$$t = \frac{\overline{x} - \mu}{s / \sqrt{n}}$$

12.3

Inference With Variance Unknown... When  $\sigma$  is unknown, we use its point estimator s $z = \overline{x} - \mu$   $(t = \overline{x} - \mu)$  $\sqrt{t} = \sqrt{\sqrt{n}}$ and the z-statistic is replaced by the the t-statistic, where the number of "degrees of freedom" v, is n–1.

## Testing $\mu$ when $\sigma$ is unknown...

When the population standard deviation is unknown and the population is normal, the test statistic for testing hypotheses about  $\mu$  is:

$$t = \frac{\overline{x} - \mu}{s / \sqrt{n}}$$

which is Student *t* distributed with  $\mathcal{V} = n-1$  degrees of freedom. The confidence interval estimator of  $\mu$  is given by:

$$\overline{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}$$

Example 12.1

It is likely that in the near future nations will have to do more to save the environment.

Possible actions include reducing energy use and recycling.

Currently (2007) most products manufactured from recycled material are considerably more expensive than those manufactured from material found in the earth.

12.6

## Example 12.1

Newspapers are an exception.

It can be profitable to recycle newspaper.

A major expense is the collection from homes. In recent years a number of companies have gone into the business of collecting used newspapers from households and recycling them.

A financial analyst for one such company has recently computed that the firm would make a profit if the mean weekly newspaper collection from each household exceeded 2.0 pounds.

12.7

#### Example 12.1

In a study to determine the feasibility of a recycling plant, a random sample of 148 households was drawn from a large community, and the weekly weight of newspapers discarded for recycling for each household was recorded. Xm12-01\*

Do these data provide sufficient evidence to allow the analyst to conclude that a recycling plant would be profitable?





Example 12.1  
From the data we determine :  

$$\sum x_i = 322.7$$
 and  $\sum x_i^2 = 845.1$   
 $\overline{x} = \frac{\sum x_i}{n} = \frac{322.7}{148} = 2.18$   
 $s^2 = \frac{\sum x_i^2 - \frac{(\sum x_i)^2}{n-1}}{n-1} = \frac{845.1 - \frac{(322.7)^2}{148}}{148-1} = .962$   
 $s = \sqrt{s^2} = \sqrt{.962} = .981$ 



Example 12.1	COMPUTE
t-Test: Mean	
Mean	Newspaper 2 1804
Standard Deviation	0.9812
Hypothesized Mean	2
df	147
t Stat	2.2369
P(T<=t) one-tail	0.0134
t Critical one-tail	2.352
P(T<=t) two-tail	0.0268
t Critical two-tail	2.6097
	12.13



## Example 12.2

In 2004 (the latest year reported) 130,134,000 tax returns were files in the United States.

The Internal Revenue Service (IRS) examined 0.77% or 1,008,000 of them to determine if they were correctly done.

To determine how well the auditors are performing a random sample of these returns was drawn and the additional tax was reported.  $\underline{Xm12-02}$ 

Estimate with 95% confidence the mean additional income tax collected from the 1,008,000 files audited.







Exa	Example 12.2 COMPUTE				<b>MPUTE</b>
	A	В	С	D	
1	t-Estimate	: Mean			
2					
3				Taxes	
4	Mean			6001	
5	Standard D	Deviation		2864	
6	LCL			5611	
7	UCL			6392	
					12.19



# Check Required Conditions

The Student t distribution is *robust*, which means that if the population is nonnormal, the results of the t-test and confidence interval estimate are still valid provided that the population is "not *extremely* nonnormal".

To check this requirement, *draw a histogram* of the data and see how "bell shaped" the resulting figure is. If a histogram is extremely skewed (say in the case of an exponential distribution), that could be considered "extremely nonnormal" and hence t-statistics would be not be valid in this case.





## **Estimating Totals of Finite Populations**

The inferential techniques introduced thus far were derived by assuming infinitely large populations. In practice however, most populations are finite.

When the population is small we must adjust the test statistic and interval estimator using the finite population correction factor introduced in Chapter 9.

However, in populations that are large relative to the sample size we can ignore the correction factor. Large populations are defined as populations that are at least 20 times the sample size.

# **Estimating Totals of Finite Populations**

Finite populations allow us to use the confidence interval estimator of a mean to produce a confidence interval estimator of the population total.

To estimate the total we multiply the lower and upper confidence limits of the estimate of the mean by the population size.

Thus, the confidence interval estimator of the total is

$$N\left[\overline{x}\pm t_{\alpha/2}\,\frac{s}{\sqrt{n}}\right]$$

#### Estimating Totals of Finite Population

For example, a sample of 500 households (in a city of 1 million households) reveals a 95% confidence interval estimate that the *household mean* spent on Halloween candy lies between \$20 & \$30.

We can estimate the *total* amount spent in the city by multiplying these lower and upper confidence limits by the total population:

$$N\left[\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}\right] = 1,000,000[\$20...\$30]$$

Thus we estimate that the *total* amount spent on Halloween in the city lies between \$20 million and \$30 million.

12.26

Developing an Understanding of Statistical Concepts

The t-statistic like the z-statistic measures the difference between the sample mean and the hypothesized value of in terms of the number of standard errors.

However, when the population standard deviation is unknown we estimate the standard error as

 $s/\sqrt{n}$ 

Developing an Understanding of Statistical Concepts

When we introduced the Student t distribution in Section 8.4 we pointed out that it is more widely spread out than the standard normal.

This circumstance is logical.

The only variable in the z-statistic is the sample mean, which will vary from sample to sample.

12.28

Developing an Understanding of Statistical Concepts

The t-statistic has two variables, the sample mean and the sample standard deviation s, both of which will vary from sample to sample.

Because of this feature the t-statistic will display greater variability.





## **Testing & Estimating Population Variance**

The test statistic used to test hypotheses about  $\sigma^2$  is:

$$\chi^2 = \frac{(n-1)s^2}{\sigma^2}$$

which is chi-squared with v = n-1 degrees of freedom.





## Example 12.3

To examine the veracity of the claim, a random sample of 25 l-liter fills was taken and the results (cubic centimeters) recorded.  $\underline{Xm12-03}$ 

Do these data allow the president to make this claim at the 5% significance level?







	Α	В	С	D
	Chi Squared Test: Variance			
				Fills
	Sample Variance			0.6333
5	Hypothesized Variance			1
6	df			24
7	chi-squared Stat			15.20
8	P (CHI<=chi) one-tail			0.0852
9	chi-squared Critical one tail	Left-tail		13.85
0	· · · ·	Right-tail		36.42
1	P (CHI<=chi) two-tail			0.1705
2	chi-squared Critical two tail	Left-tail		12.40
3		Right-tail		39.36



Example 12.4  

$$x_{\alpha/2,n-1}^2 = x_{.005,24}^2 = 45.56$$
  
 $x_{1-\alpha/2,n-1}^2 = x_{.995,24}^2 = 9.89$   
 $LCL = \frac{(n-1)s^2}{x_{\alpha/2}^2} = \frac{15.20}{45.56} = .3336$   
 $UCL = \frac{(n-1)s^2}{x_{1-\alpha/2}^2} = \frac{15.20}{9.89} = 1.537$ 

Exa	mple 12.4				
Estin	Estimate with 99% confidence the variance of fills in				
Exan	nple 12.3. <u>Xm12-03</u>				
	A	В			
1	Chi Squared Estimate: Variance				
2					
3		Fills			
4	Sample Variance	0.6333			
5	df	24			
6	LCL	0.3336			
7	UCL	1.5375			
		12.42			

Example 12.4	INTERPRET
In Example 12.3, we saw that there was evidence to infer that the population va	s not sufficient riance is less than 1.
Here we see that is estimated to lie betw 1.5375.	ween .3336 and
Part of this interval is above 1, which to variance may be larger than 1, confirm reached in Example 12.3.	ells us that the ing the conclusion we
	12.43





## Inference: Population Proportion...

When data are nominal, we count the number of occurrences of each value and calculate proportions. Thus, the parameter of interest in describing a population of nominal data is the population proportion p.

This parameter was based on the binomial experiment.

Recall the use of this statistic:  $\hat{p} = \frac{x}{n}$ 

where p-hat  $(\hat{p})$  is the sample proportion: **x** successes in a sample size of **n** items.

## Inference: Population Proportion...

When np and n(1-p) are both greater than 5, the sampling distribution of  $\hat{p}$  is approximately normal with



## Inference: Population Proportion

Test statistic for *p*:

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

The confidence interval estimator for  $\boldsymbol{p}$  is given by:

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

(both of which require that np > 5 and n(1-p) > 5)

# Example 12.5

After the polls close on election day networks compete to be the first to predict which candidate will win.

The predictions are based on counts in certain precincts and on exit polls.

Exit polls are conducted by asking random samples of voters who have just exited from the polling booth (hence the name) for which candidate they voted.

12.49

## Example 12.5

In American presidential elections the candidate who receives the most votes in a state receives the state's entire Electoral College vote.

In practice, this means that either the Democrat or the Republican candidate will win.

Suppose that the results of an exit poll in one state were recorded where 1 = Democrat and 2 = Republican.

<u>Xm12-05</u>\*



Example 12.5	IDENTIFY
The problem objective is to describe to in the state. The data are nominal beca "Democrat" and "Republican." Thus tested is the proportion of votes in the the Republican candidate. Because we whether the network can declare the F winner at 8:01 P.M., the alternative hy	the population of votes ause the values are the parameter to be e entire state that are for e want to determine Republican to be the ypothesis is
$H_1: p > .50$	
And hence our null hypothesis becom $H_0: p = .50$	nes:
	12.52



Example 12.5  
The rejection region is  

$$z > z_{\alpha} = z_{.05} = 1.645$$
  
 $\hat{p} = \frac{x}{n} = \frac{407}{765} = .532$   
 $z = \frac{\hat{p} - p}{\sqrt{p(1 - p)/n}} = \frac{.532 - .5}{\sqrt{.5(1 - .5)/765}} = 1.77$   
 $p - value = p(z > 1.77) = 1 - p(z < 1.77) = 1 - .9616 = .0384$ 

z-Test: $Proportion$ VotesSample $Proportion$ 0.532Observations765Hypothesized Proportion0.5z Stat1.77 $P(Z<=z)$ one-tail0.0382z Critical one-tail1.6449 $P(Z<=z)$ two-tail0.0764		А	В	С	D	
Sample ProportionVotesSample Proportion $0.532$ Observations765Hypothesized Proportion $0.5$ z Stat $1.77$ P(Z<=z) one-tail	1	z-Test: Pro	oportion			
Sample Proportion         Votes           Sample Proportion         0.532           Observations         765           Hypothesized Proportion         0.5           z Stat         1.77           P(Z<=z) one-tail	2					
Sample Proportion       0.532         Observations       765         Hypothesized Proportion       0.5         z Stat       1.77         P(Z<=z) one-tail	3				Votes	
Observations         765           Hypothesized Proportion         0.5           z Stat         1.77           P(Z<=z) one-tail	4	Sample Pro	oportion		0.532	
Hypothesized Proportion       0.5         z Stat       1.77         P(Z<=z) one-tail	5	Observatio	ns		765	
z Stat       1.77         P(Z<=z) one-tail	6	Hypothesiz	ed Proportion	on	0.5	
P(Z<=z) one-tail         0.0382           z Critical one-tail         1.6449           ) P(Z<=z) two-tail	7	z Stat			1.77	
z Critical one-tail1.6449) P(Z<=z) two-tail	8	P(Z<=z) or	ne-tail		0.0382	
) P(Z<=z) two-tail 0.0764	9	z Critical or	ne-tail		1.6449	
	10	P(Z<=z) tw	o-tail		0.0764	
1 z Critical two-tail 1.96	11	z Critical tw	vo-tail		1.96	



Example 12.5	INTERPRET
Such an error would mean that a m 8:01 P.M. that the Republican has evening would have to admit to a m	etwork would announce at won and then later in the mistake.
If a particular network were the or it would cast doubt on their integra number of viewers.	nly one that made this error ity and possibly affect the





## **Estimating Totals for Large Populations**

In much the same way as we saw earlier, when a population is *large* and *finite* we can estimate the total number of successes in the population by taking the product of the size of the population (**N**) and the confidence interval estimator:

$$N\left[\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right]$$

The Nielsen Ratings (used to measure TV audiences) uses this technique. Results from a small sample audience (5,000 viewers) is extrapolated to the total number of TV households (110 million).

# **Nielsen Ratings**

Statistical techniques play a vital role in helping advertisers determine how many viewers watch the shows that they sponsor.

Although several companies sample television viewers to determine what shows they watch, the best known is the A. C. Nielsen firm.

The Nielsen ratings are based on a random sample of approximately 5,000 of the 110 million households in the United States with at least one television (in 2007).

12.61

## **Nielsen Ratings**

A meter attached to the televisions in the selected households keeps track of when the televisions are turned on and what channels they are tuned to.

The data are sent to the Nielsen's computer every night from which Nielsen computes the rating and sponsors can determine the number of viewers and the potential value of any commercials.

# Nielsen Ratings

The results from Sunday, April 1, 2007 for the time slot 9:00 to 9:30 P.M. have been recorded using the following codes:

Network	Show	Code
ABC	Desperate Housewives	1
CBS	The Amazing Race 11	2
NBC	Deal or No Deal	3
Fox	Family Guy	4
Television turn	ed off or watched another channel	5

12.63

12.64

 Nielsen Ratings
 IDENTIFY

 The problem objective is to describe the population of television shows watched by viewers across the country.

 The data are nominal.

 The combination of problem objective and data type make the parameter to be estimated the proportion of the entire population that watched the "Deal or No Deal."

 The confidence interval estimator of the proportion is

  $\hat{\rho}(1-\hat{p})$ 

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

Niels	Nielsen Ratings Example <b>Сомр</b> ит				
	A	В			
1	z-Estimate: Proportion				
2		Program			
3	Sample Proportion	0.0836			
4	Observations	5000			
5	LCL	0.0759			
6	UCL	0.0913			
			12.65		





## Selecting the Sample Size

When we introduced the sample size selection method to estimate a mean in Section 10.3, we pointed out that the sample size depends on the confidence level and the bound on the error of estimation that the statistics practitioner is willing to tolerate.

When the parameter to be estimated is a proportion the bound on the error of estimation is

$$\mathbf{B} = \mathbf{z}_{\alpha/2} \sqrt{\frac{\hat{\mathbf{p}}(1-\hat{\mathbf{p}})}{n}}$$

# Selecting the Sample Size

Solving for n we produce the required sample size to estimate p and where B is the bound on the error of Estimation

$$n = \left(\frac{z_{\alpha/2}\sqrt{\hat{p}(1-\hat{p})}}{B}\right)^2$$

Unfortunately we do not know the value of  $\hat{p}$ 

## Selecting the Sample Size

Two methods – in each case we choose a value for  $\hat{p}$  then solve the equation for n.

Method 1 : no knowledge of even a rough value of  $\hat{p}$  This is a 'worst case scenario' so we substitute  $\hat{p} = .50$ 

Method 2 : we have some idea about the value of  $\hat{p}$  This is a better scenario and we substitute in our estimated  $\hat{p}$  value.

12.70

#### Selecting the Sample Size

Method 1 :: no knowledge of value of  $\hat{p}$ , use 50%:

$$n = \left(\frac{1.96\sqrt{.50(1 - .50)}}{.03}\right)^2 = 1,068$$

Method 2 :: some idea about a possible  $\hat{p}$  value, say 20%:  $(106\sqrt{20(1-20)})^2$ 

$$n = \left(\frac{1.96\sqrt{.20(1 - .20)}}{.03}\right) = 683$$

Thus, we can sample fewer people if we already have a reasonable estimate of the population proportion before starting.





*Mass marketing* refers to the mass production and marketing by a company of a single product for the entire market.

Mass marketing is especially effective for commodity goods such as gasoline, which are very difficult to differentiate from the competition,

It has given way to target marketing, which focuses on satisfying the demands of a particular segment of the entire market.

12.73

#### APPLICATIONS IN MARKETING: Market Segmentation

Because there is no single way to segment a market, managers must consider several different variables (or characteristics) that could be used to identify segments.

Surveys of customers are used to gather data about various aspects of the market, and statistical techniques are applied to define the segments.

Market segmentation separates consumers of a product into different groups in such a way that members of each group are similar to each other and there are differences between groups.

There are many ways to segment a market.

Table 12.1 lists several different segmentation variables and their market segments.

Table 12.1 Marke	t Segmentation	
Segmentation variable Geographic	Segments	
Countries	Brazil, Canada, China, France, United States	
Country regions	Midwest, Northeast, Southwest, Southeast	
Demographic		
Age	Under 5, 5- 12, 13-19, 20-29, 30-50, over 50	
Education	Some high school, high school graduate, some college, college or university graduate	
Income	Under \$30,000, \$30,000-49,999, \$50,000- 79,999, over \$80,000	
Marital status	Single, married, divorced, widowed	
		12.76

It is important for marketing managers to know the size of the segment because the size (among other parameters) determines its profitability.

Not all segments are worth pursuing. In some instances the size of the segment is too small or the costs of satisfying it may be too high.

#### APPLICATIONS IN MARKETING: Market Segmentation

The size can be determined in several ways.

The census provides useful information.

For example, we can determine the number of Americans in various age categories or the size of geographic residences.

For other segments we may need to survey members of a general population and use the inferential techniques introduced in the previous section where we showed how to estimate the total number of successes.

12.78

We can survey large populations to estimate the proportion of the population that fall into each segment.

From these estimates we can estimate the size of markets using the confidence interval estimator

$$N\left(\hat{p}\pm z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right)$$

Example 12.6

In segmenting the breakfast cereal market a food manufacturer uses health and diet consciousness as the segmentation variable. Four segments are developed:

- 1. Concerned about eating healthy foods
- 2. Concerned primarily about weight
- 3. Concerned about health because of illness
- 4. Unconcerned

12.80

#### Example 12.6

To distinguish between groups surveys are conducted. On the basis of a questionnaire people are categorized as belonging to one of these groups.

A recent survey asked a random sample of 1250 American adults (20 and over) to complete the questionnaire. The categories were recorded using the codes.  $\underline{Xm12-06}$ 

The most recent census reveals that there are 207,347,000 Americans who are 20 and over.

Estimate with 95% confidence the number of American adults who are concerned about eating healthy foods.



E	kam	ple 12.6	COM	PUTE
		A	В	
	1	z-Estimate: Proportion		
	2		Group	
	3	Sample Proportion	0.2152	
	4	Observations	1250	
	5	LCL	0.1924	
	6	UCL	0.2380	
				12.83



LAAIIIPIE IZ.0	Examp	le	12.	6
----------------	-------	----	-----	---

INTERPRET

We estimate that the size of this market segment lies between 39,893,563 and 49,348,586.