# Chapter 8: Confidence Intervals

Yu Yang

School of Statistics
University of Minnesota

October 7, 2022

# Difference Between $\mu$ and $p$

From here on, it is important to identify the difference between $\mu$ and $p$. The goal is to identify whether a problem is a $p$ problem or a $\mu$ problem!

$\mu$: when the variable we are dealing with is quantitative (numerical)
$p$: when the outcome is binary (only 2 outcomes are possible).

- true average height of North American females $\mu$
- true # of hot dogs eaten in a day by students at the U $\mu$
- true proportion of people who will vote for the Republican party in 2016. $p$
- true proportion of smokers at the U $p$

# Point Estimation

We have already done point estimation: we used $\bar{x}$ to estimate $\mu$ and $\hat{p}$ to estimate $p$

For any parameter, there can be many point estimates. Example: we can estimate $\mu$ with $\bar{x}$ or $M$ (median), or some other estimator of the mean.

Which estimator do we chose?

# Assessing an Estimator

1. How close can we expect our estimate (a statistic) to be to the truth (a parameter)?
2. How much does the estimate vary from sample to sample?

# Bias and Standard Error

### Bias
The bias of a statistic is the difference between the mean of its sampling distribution and the true parameter value.

A statistic is **unbiased** if this difference is zero. An unbiased estimator doesn't systematically over- or underestimate the parameter, so we want an unbiased estimator.

### Standard Error (se)
A **standard error** is the estimated standard deviation of a sampling distribution.

Note: Standard errors are commonly used in practice, because TRUE standard deviations are usually unknown.

For example, we plug in $\hat{p}$ to estimate the standard deviation of the sampling distribution and get $\sqrt{\hat{p}(1 - \hat{p})/n}$.

Oct 7 Lecture Stopped Here

# Assessing an Estimator

A good estimator should have
1. a small bias (ideally **no bias**)
2. a small standard error

## Unbiased Estimator

Is $\bar{x}$ an unbiased estimator of $\mu$? Lets look at the Sampling distribution: (if $n \geq 30$)

$$\bar{x} \mathrel{\dot\sim} N(\mu, \sigma/\sqrt{n})$$

Bias = Mean of sampling distribution - true parameter value

$= \mu - \mu = 0$

$\bar{x}$ is unbiased!

What about $\hat{p}$?

# Unbiased Estimator

Is $\bar{x}$ an unbiased estimator of $\mu$? Lets look at the Sampling distribution: (if $n \geq 30$)

$$\bar{x} \stackrel{.}{\sim} N(\mu, \sigma/\sqrt{n})$$

Bias = Mean of sampling distribution - true parameter value
$\quad = \mu - \mu = 0$
$\bar{x}$ is unbiased!

What about $\hat{p}$?

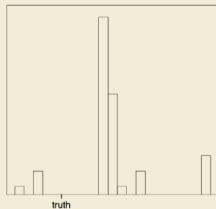$$\hat{p} \stackrel{.}{\sim} N\left(p, \sqrt{\frac{p(1-p)}{n}}\right)$$

Bias = $p - p = 0$
$\hat{p}$ is unbiased!

# Example 8.1

Here's a scenario:

- We want to choose from 4 different statistics
- The value of each statistic will vary from sample to sample
- We drew multiple samples and noted the value of each statistic
- We plotted histograms and have approximate sampling distributions. The truth is marked on the histograms
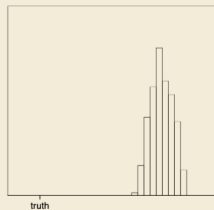
# Example 8.1



biased, large standard error          unbiased, small standard error

unbiased, large standard error        biased, small standard error

## Interval Estimation

Point estimates are our best guess for the parameter. But, it might be too high, too low, or just about right!

Therefore, we also need a measure of its reliability, i.e., we want to be able to say how *confident* we are in the estimate.

**Interval Estimate:** An interval estimate is an interval of numbers within which the parameter value is believed to fall.

# Example 8.2

According to the 2012 Resident Satisfaction Survey, 32% of Minneapolis residents felt that public safety is among the biggest challenges of the city in the next five years.

Considering the reported margin of error of plus or minus 3 percentage points, find an interval estimate for $p$, the true proportion of Minneapolis residents who felt that public safety is among the biggest challenges of the city in the next five years.

$$\hat{p} \pm \textbf{moe} = .32 \pm .03 = (.29, .35).$$

We are fairly confident that the true proportion of Minneapolis residents who felt that public safety is among the biggest challenges of the city in the next five years is between 29% and 35%.

# Confidence Interval

## Margin of Error (moe)

The margin of error measures how accurate a point estimate is likely to be in estimating a parameter.

## Confidence Interval (CI)

A confidence interval is an interval containing the most believable values for a parameter. The probability that this method produces an interval that contains the parameter is called the **confidence level**.

## General Form of a CI

point $\pm$ moe = (point estimate - moe, point estimate + moe)

The larger the moe, the wider the interval and the less accurate our estimate is.

# Example 8.3

We want to estimate the number of contracts in a city that are awarded to minority owned firms. We take a random sample of 389 contracts from the city and find that 58 were awarded to minority-owned firms.

Our goal is to use the sample information to construct a confidence interval with a confidence level of .95 for $p$, the true proportion of contracts that are assigned to minority-owned firms.

# Example 8.3

(a) What is the sampling distribution of $\hat{p}$, the sample proportion of contracts that went to minority-owned firms?

$$\hat{p} \overset{.}{\sim} N\left(p, \sqrt{\frac{p(1-p)}{n}}\right)$$

(b) Between what 2 values does the middle 95% of this distribution fall? $z = \pm 1.96$ corresponds to 95% of the distribution. Now we unstandardize

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

In 95% of all possible samples, $\hat{p}$ will be within $1.96\sqrt{p(1-p)/n}$ of $p$. The other 5% will be "unlucky" and $\hat{p}$ will be further than $1.96\sqrt{p(1-p)/n}$ from $p$.

# Example 8.3

(c) Using the estimated standard deviation (i.e., standard error) of $\hat{p}$, construct a 95% confidence interval for $p$.

Part (b) $\Rightarrow$

$$\left(\hat{p} - 1.96\sqrt{\frac{p(1-p)}{n}}, \hat{p} + 1.96\sqrt{\frac{p(1-p)}{n}}\right)$$

Estimate $p$ by $\hat{p}$. $\hat{p} = 58/389 = .149$. A 95% CI for $p$ is

$$.149 \pm 1.96\sqrt{\frac{.149(1-.149)}{n}} = (.149 - .035, .149 + .035) = (.114, .184)$$

Interpretation: On average, approximately 95 out of every 100 samples drawn from the population would produce CIs that contain the true population proportion $p$.

# Example 8.3

Suppose an equal opportunity law requires that least 20% of all city contracts go to minority-owned firms. What information does the CI give us in this context?

We are 95% confident that the true proportion of contracts that go to minority-owned firms is between .114 and .184.

Therefore it is unlikely that the true proportion of contracts that go to the minority-owned firms is 20% or higher.

Oct 10 Lecture Stopped Here

# Confidence Intervals

NOTE: We can calculate confidence intervals for any confidence level, not just 0.95.

**Notation:**
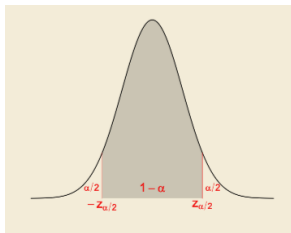$\alpha$ = error probability of the confidence interval ($0 < \alpha < 1$)
$1 - \alpha$ = confidence level of the confidence interval

# Large Sample Confidence Interval for $p$

We take a random sample size $n$, and number of successes and failures are both greater than 15. Then the large sample CI for $p$ with confidence level $1 - \alpha$ is

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

where $z_{\alpha/2}$ depends on the confidence level. Specifically, $\pm z_{\alpha/2}$ mark the **middle** $1 - \alpha$ proportion of $N(0,1)$ distribution

# Finding $z_{\alpha/2}$

Though we can find $z_{\alpha/2}$ for any specified confidence level, 90%, 95% and 99% confidence intervals are the most common. We provide $z_{\alpha/2}$ for these confidence levels below:

| Confidence Level | $\alpha$ | $\alpha/2$ | $z_{\alpha/2}$ |
|:---:|:---:|:---:|:---:|
| 90% | .10 | .050 | 1.645 |
| 95% | .05 | .025 | 1.960 |
| 99% | .01 | .005 | 2.575 |

# Interpreting Confidence Intervals

1. *In the long run*, 95% of 95% CIs will contain the true parameter. The other 5% that do not are based on the "unlucky" samples that produce unusually high or low values of the statistic.

2. We are 95% confident that the true value of the parameter is between the upper and the lower bounds of the 95% CI.

3. We are 95% confident the interval will include the true value of the population parameter.

# Example 8.4

When the 2010 General Social Survey asked subjects if they would be willing to accept cuts in their standard of living to protect the environment, 486 of 1374 said yes.

(a) Estimate the population proportion $p$, who would answer yes.

# Example 8.4

When the 2010 General Social Survey asked subjects if they would be willing to accept cuts in their standard of living to protect the environment, 486 of 1374 said yes.

(a) Estimate the population proportion $p$, who would answer yes.
$\hat{p} = 486/1374 = 0.354$.

(b) Can we use this sample to calculate a valid confidence interval?

# Example 8.4

When the 2010 General Social Survey asked subjects if they would be willing to accept cuts in their standard of living to protect the environment, 486 of 1374 said yes.

(a) Estimate the population proportion $p$, who would answer yes.
$\hat{p} = 486/1374 = 0.354$.

(b) Can we use this sample to calculate a valid confidence interval?
Yes, we have more than 15 successes and 15 failures, so we can apply the CLT and construct a confidence interval based on that.

# Example 8.4

(c) Calculate a 90% CI for $p$.

# Example 8.4

(c) Calculate a 90% CI for $p$.

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = .354 \pm 1.645 \sqrt{\frac{.354(1-.354)}{1374}}$$
$$= .354 \pm .021 = (.333, .375)$$

We are 90% confident that the true proportion of people willing to lower their standard of living is between .333 and .375.

(d) Now calculate a 98% CI for $p$.

# Example 8.4

(c) Calculate a 90% CI for $p$.

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = .354 \pm 1.645\sqrt{\frac{.354(1-.354)}{1374}}$$
$$= .354 \pm .021 = (.333, .375)$$

We are 90% confident that the true proportion of people willing to lower their standard of living is between .333 and .375.

(d) Now calculate a 98% CI for $p$.

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = .354 \pm 2.33\sqrt{\frac{.354(1-.354)}{1374}}$$
$$= .354 \pm .030 = (.324, .384)$$

(e) Similarly, a 99% CI for $p$ can be shown to equal $(.321, .388)$. What pattern do you see as we increase the confidence level?

# Example 8.4

(c) Calculate a 90% CI for $p$.

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = .354 \pm 1.645 \sqrt{\frac{.354(1-.354)}{1374}}$$
$$= .354 \pm .021 = (.333, .375)$$

We are 90% confident that the true proportion of people willing to lower their standard of living is between .333 and .375.

(d) Now calculate a 98% CI for $p$.

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = .354 \pm 2.33 \sqrt{\frac{.354(1-.354)}{1374}}$$
$$= .354 \pm .030 = (.324, .384)$$

(e) Similarly, a 99% CI for $p$ can be shown to equal $(.321, .388)$. What pattern do you see as we increase the confidence level?
As the confidence level increases, the CI becomes wider.

# Example 8.5: Confidence Intervals in R

A woman who smokes during pregnancy increases health risks to her infant. Let $p$ = proportion of women smokers that quit during pregnancy .

(a) Suppose that in a random sample of 300 women who smoked prior to pregnancy, 51 quit smoking during pregnancy. Use this sample information to calculate a 98% CI for $p$ in R.

```
> # 'x' = number of successes, 'n' = sample size
> # prop.test builds confidence intervals for proportions.
> prop.test(x = 51, n = 300, conf.level = .98,
+            alternative = "two.sided")

        1-sample proportions test with continuity correction

data:  51 out of 300, null probability 0.5
X-squared = 129.3633, df = 1, p-value < 2.2e-16
alternative hypothesis: true p is not equal to 0.5
98 percent confidence interval:
0.1240578 0.2280177
sample estimates:
   p
0.17
```

We are 98% confident that the true proportion of women smokers quit during pregnancy is between .124 and .228.

# Example 8.5

(b) Now, suppose that we have a random sample of 1000 women smokers in which 170 quit smoking. Use this sample information to calculate a 98% CI for $p$ in R.

```
> prop.test(x = 170, n = 1000, conf.level = .98,
+            alternative = "two.sided")

        1-sample proportions test with continuity correction

data:  170 out of 1000, null probability 0.5
X-squared = 434.281, df = 1, p-value < 2.2e-16
alternative hypothesis: true p is not equal to 0.5
98 percent confidence interval:
 0.1436947 0.1999219
sample estimates:
   p
0.17
```

We are 98% confident that the true proportion of women smoker that quit during pregnancy is between .144 and .200.

# Example 8.5

(c) What pattern do you see?

# Example 8.5

(c) What pattern do you see?

The samples produce the same point estimate ($\hat{p} = 0.17$). However, at the same level of confidence, a larger sample size results in a smaller moe.

Oct 12 Lecture Stopped Here

# Small Sample Confidence Interval for $p$

When $n$ is small, the normal approximation to the sampling distribution of $\hat{p}$ can be awful, thus large sample confidence interval for $p$ can be misleading.

See the paragraph "Exact Intervals for Small Numbers of Failures and/or Small Sample Sizes" in Section 7.2.4.1 of *NIST/SEMATECH e-Handbook of Statistical Methods* at `https://www.itl.nist.gov/div898/handbook/prc/section2/prc241.htm` for details. We will not talk about this case any further.

# Properties of a Confidence Interval

1. moe increases as confidence level increases
2. for fixed $\hat{p}$, moe decreases as sample size increases.

# Quick Recap

- Point estimates are good, but we want to say how confident we are in the estimate.
- A confidence interval is an interval of possible values for the parameter
- For proportion $p$, the $1 - \alpha$ confidence interval is:

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

- As moe increases, confidence level increases.
- We want high levels of confidence and small intervals.
- For small intervals, we need a large sample size.

# Sample Size Calculations

Goal: Ideally, we would like **both**

1. High level of confidence in our interval estimate; and
2. Small margin of error

If the sample size is large enough, we can achieve **both** of these goals simultaneously. If this is the case, then why don't we always pick really large samples?

# Sample Size Calculations

Goal: Ideally, we would like **both**

1. High level of confidence in our interval estimate; and
2. Small margin of error

If the sample size is large enough, we can achieve **both** of these goals simultaneously. If this is the case, then why don't we always pick really large samples?

Sample size is limited by money and other resources. Instead, our goal is to pick $n$ that is "big enough".

# Example 8.2 Continued

Recall that in the 2012 Resident Satisfaction Survey, the City of Minneapolis found out that 32% of the residents felt that public safety is among the biggest challenges of the city in the next five years. However, suppose we want to obtain a more recent estimate of $p$, the true proportion of Minneapolis residents who felt that public safety is among the biggest challenges of the city in the next five years.

How many people do we need to poll in order to estimate $p$ within 3 percentage points with 95% confidence?

# Example 8.2 Continued

Recall that in the 2012 Resident Satisfaction Survey, the City of Minneapolis found out that 32% of the residents felt that public safety is among the biggest challenges of the city in the next five years. However, suppose we want to obtain a more recent estimate of $p$, the true proportion of Minneapolis residents who felt that public safety is among the biggest challenges of the city in the next five years.

How many people do we need to poll in order to estimate $p$ within 3 percentage points with 95% confidence?

$$\text{moe} = z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq .03 \Rightarrow \frac{\hat{p}(1-\hat{p})}{n} \leq \left(\frac{.03}{1.96}\right)^2$$

$$\Rightarrow n \geq \hat{p}(1-\hat{p})\left(\frac{1.96}{.03}\right)^2$$

# Example 8.2 Continued

Problem: Our lower bound for $n$ depends on $\hat{p}$. However, we haven't collected any data yet so we don't have an estimate $\hat{p}$.

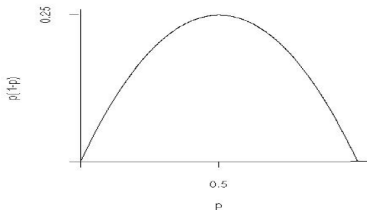Solution: substitute an educated guess of $p$ for $\hat{p}$

Back to the problem...

Old survey $\Rightarrow \hat{p} = .32$.
$n \geq .32(1 - .32)(\frac{1.96}{.03})^2 = 928.81 \Rightarrow n = 929$.

Always round up!

# What if we had no prior guess for $p$?

Use the worst case scenario. moe is maximum when $\sqrt{\hat{p}(1-\hat{p})/n}$ is maximum (i.e., when $\hat{p}(1-\hat{p})$ is maximum):



This happens at $\hat{p} = 1/2$.

Therefore, using $\hat{p} = 1/2$ will give a larger value for $n$ than any other $\hat{p}$. However, this will sometimes lead us to collect a lot more data than is necessary if $p$ is either close to 0 or close to 1.

# Sample Size for a Desired MOE

The sample size required to estimate $p$ within margin of error of (at most) $m$ with a confidence level of $1 - \alpha$ is

$$n = p^*(1 - p^*)\left(\frac{z_{\alpha/2}}{m}\right)^2,$$

where $z^*$ depends on the desired confidence level $(1 - \alpha)$, and

1. $p^*$ is based on other information or our past experience or
2. $p^* = .5$ when we don't have any reasonable guess for $p$.

NOTE: **Always round up for sample size calculations!**

## Example 8.6

A campus group is interested in estimating the proportion of U students that feel that their binge drinking has negatively impacted their academic performance. How many students should they poll in order to estimate the proportion to within $\pm.05$ with 90% confidence if

(a) We have no prior information?

# Example 8.6

A campus group is interested in estimating the proportion of U students that feel that their binge drinking has negatively impacted their academic performance. How many students should they poll in order to estimate the proportion to within $\pm.05$ with 90% confidence if

(a) We have no prior information?

$$n = p^*(1 - p^*)\left(\frac{z_{\alpha/2}}{m}\right)^2 = .5(1 - .5)\left(\frac{1.645}{.05}\right)^2 = 270.6 \approx 271$$

(b) Based on UMADD statistics, we expect the proportion to be around 0.25?

# Example 8.6

A campus group is interested in estimating the proportion of U students that feel that their binge drinking has negatively impacted their academic performance. How many students should they poll in order to estimate the proportion to within $\pm.05$ with 90% confidence if

(a) We have no prior information?

$$n = p^*(1 - p^*)\left(\frac{z_{\alpha/2}}{m}\right)^2 = .5(1 - .5)\left(\frac{1.645}{.05}\right)^2 = 270.6 \approx 271$$

(b) Based on UMADD statistics, we expect the proportion to be around 0.25?

$$n = p^*(1 - p^*)\left(\frac{z_{\alpha/2}}{m}\right)^2 = .25(1 - .25)\left(\frac{1.645}{.05}\right)^2 = 202.95 \approx 203$$

(don't need as many observations as in the worst case scenario!)

Oct 14 Lecture Stopped Here

# Confidence Intervals for Population Mean $\mu$

The procedure here is very similar to what we did for population proportion $p$. The general form of a confidence interval is:

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

We will use the sampling distribution of $\bar{X}$ to construct confidence interval for $\mu$. Recall that if $n \geq 30$, then

$$\bar{X} \stackrel{.}{\sim} N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

# Confidence Intervals for Population Mean $\mu$

If we standardize $\bar{X}$, we have:

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

**The Problem:** $\sigma$ is unknown

**The Solution:** Estimate $\sigma$ using the sample standard deviation $s$

$$t = \frac{\bar{X} - \mu}{s/\sqrt{n}}$$

NOTE: $t$ does not have a normal distribution! Since we replace $\sigma$ with $s$, we introduce some "error". Therefore, the distribution for $t$ changes due to increased variability.
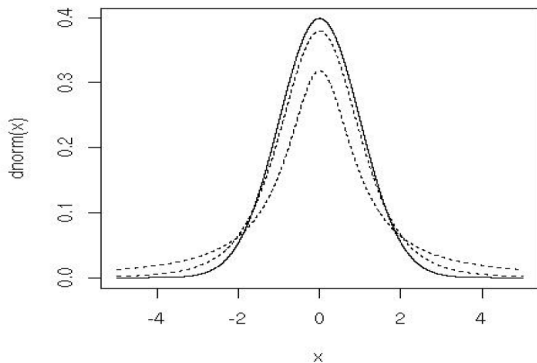
# The $t$-distribution

If $X_1, X_2, \ldots, X_n$ is a random sample from a normal population with unknown mean $\mu$ and unknown standard deviation $\sigma$, then

$$t = \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t_{n-1}$$

where $t_{n-1}$ denotes the "$t$-distribution with $n-1$ degrees of freedom".

NOTE: The above equation is also *approximately* true when $X_1, X_2, \ldots, X_n$ is a random sample from *any* distribution with mean $\mu$ and standard deviation $\sigma$ is $n$ is large ($\geq 30$)
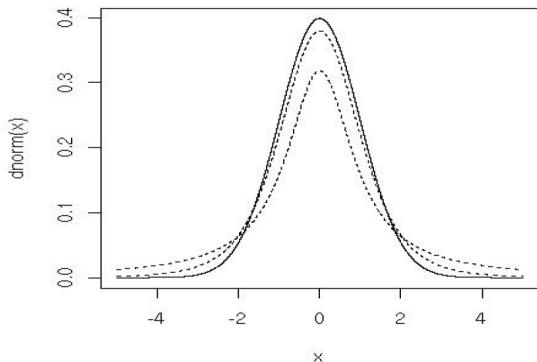
# Properties of $t$-distribution



**Properties:**

- bell-shaped
- symmetric about 0
- has fatter tails than $N(0, 1)$ (i.e., more spread out)

# $t$-distribution and $N(0, 1)$



The exact shape and spread of the $t$ distribution is determined by degrees of freedom(df)

As df increases, $t$-distribution gets closer to $N(0, 1)$. This is because $s$ estimates $\sigma$ better when $n$ increases.

# $t$-distribution Calculations

Use R to find the probabilities and quantiles for $t$-distributions with varying $df$.

1. If df $= 11$, find the value of $t$ for which $P(t_{11} \geq t) = .005$

```
> qt(1 - 0.005, df = 11)
[1] 3.105807
> qt(0.005, df = 11, lower = FALSE)
[1] 3.105807
```

2. If df $= 21$, find the value of $t$ for which $P(t_{21} \leq t) = .99$)

```
> qt(0.99, df = 21)
[1] 2.517648
```

3. If df $= 14$, What is $P(t_{14} \geq 2)$?

```
> 1 - pt(2, df = 14)
[1] 0.03264398
> pt(2, df = 14, lower = FALSE)
[1] 0.03264398
```

# Confidence Interval for $\mu$

**Assumptions**:

1. randomness: $x_1, x_2, \ldots, x_n$ are a random sample from some population with unknown mean $\mu$ and unknown standard deviation $\sigma$.
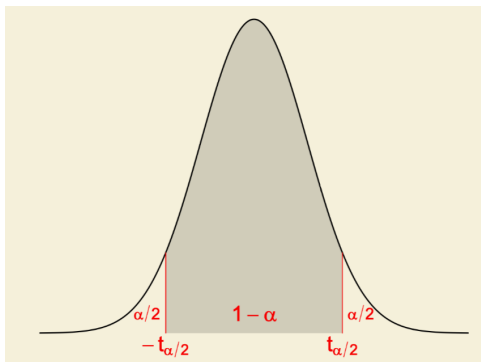
2. normality: the population distribution is approximately normal

If assumptions above are satisfied, the confidence interval with level $1 - \alpha$ for $\mu$ is

$$\bar{x} \pm t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}$$

$$(df = n - 1)$$

# Confidence Interval for $\mu$

Here $t_{\alpha/2, n-1}$ depends on the confidence level. Specifically $\pm t_{\alpha/2, n-1}$ marks the middle $1 - \alpha$ proportion of the $t_{n-1}$ distribution.



1. moe increases as the confidence level increases
2. moe decreases as sample size increases
3. the smaller the standard deviation, the smaller the moe

# Example 8.8

A survey of 51 current adult smokers in the U.S. asked "On average, how many cigarettes do you smoke per day?"

(a) Based on the following stem-and -leaf plot of the raw data, does it appear that the underlying assumptions of the confidence interval are satisfied?

```
> cigs <- c(1,3,3,8,9,9,9,10,11,11,12,13,14,14,15,15,15,16,
+ 16,16,16,17,17,17,17,18,19,19,20,20,20,20,20,22,22,
+ 23,23,24,25,25,25,28,30,30,30,30,32,32,35,38,40)
> stem(cigs)
0 | 133
0 | 8999
1 | 0112344
1 | 55566667777899
2 | 0000022334
2 | 5558
3 | 000022
3 | 58
4 | 0
```

Looks approximately normal

# Example 8.8

(b) Calculate and interpret a 99% confidence interval for $\mu$, the true mean number of cigarettes smoker per day by U.S. smokers.

```
> mean(cigs)
[1] 19.09804
> sd(cigs)
[1] 8.775545
> length(cigs)
[1] 51
> qt(1 - 0.005, df = 50)
[1] 2.677793
> (moe <- qt(1 - 0.005, df = 50) * sd(cigs) / sqrt(length(cigs)))
[1] 3.290532
> mean(cigs) + c(-1, 1) * moe
[1] 15.80751 22.38857
```

$$\bar{x} \pm t_{\alpha/2, n-1} \frac{s}{\sqrt{n}} = 19.1 \pm t_{.005, 50} \frac{8.78}{51}$$

$$= 19.1 \pm (2.68)(1.23)$$

$$= 19.1 \pm 3.3 = (15.8, 22.4)$$

Interpretation: We are 99% confident that the average number of cigarettes that a smoker has each day is between 15.8 and 22.4

## Example 8.8

3. Repeat the analysis using the t.test function in R.

```
> t.test(x = cigs, conf.level = 0.99,
+         alternative = "two.sided")

One Sample t-test

data:  cigs
t = 15.542, df = 50, p-value < 2.2e-16
alternative hypothesis: true mean is not equal to 0
99 percent confidence interval:
 15.80751 22.38857
sample estimates:
mean of x
 19.09804
```

# Robustness of $t$-Procedures

When we use the $t$ distribution to construct a confidence interval for $\mu$, we assume that the sample is **randomly** drawn from a population that is **approximately** normal. However, in practice these assumptions are rarely satisfied.

How does this affect the reliability of interval estimates?

# Robustness

**Robustness:** A statistical method is robust with respect to a particular assumption if it performs adequately even when that assumption is violated.

NOTES: The $t$ procedure for calculating a confidence interval

1. is robust to non-normality (if no outliers present)
2. does not perform well with extreme outliers
3. **not** robust to violations of the assumption of a random sample.

# Sample Size Calculations

Goal: Determine how large of a sample size is needed to obtain **both** a small margin of error and high confidence in our interval estimate.
For instance, how large of a sample do we need if we want a confidence interval for $\mu$ with a $1 - \alpha$ confidence level and moe that is no more than $m$?

$$t_{\alpha/2, n-1} \frac{s}{\sqrt{n}} \leq m \Rightarrow \sqrt{n} \geq t_{\alpha/2, n-1} \frac{s}{m}$$

$$\Rightarrow n \geq \left( \frac{t_{\alpha/2, n-1} s}{m} \right)^2$$

The Problem: Since no sample yet, $s$ and df for $t$ are unknown!

# Sample Size Calculations

The solutions:

- substitute an educated guess of $\sigma$ for $s$

- use the normal distribution ($z_{\alpha/2}$) instead of $t_{\alpha/2,n-1}$. (This will be fairly close for large $n$)

# Sample Size for a Desired MOE

The sample size required to estimate $\mu$ within moe $m$ with a confidence level of $1 - \alpha$ is approximately,

$$n = \left( \frac{z_{\alpha/2} s^*}{m} \right)^2$$

where $s^*$ is an estimate of $\sigma$ based on other information or past experience.

**Reminder: Always round up for sample size**

## Example 8.10

Suppose the Minnesota wildlife service wishes to estimate the mean number of days spent hunting, per hunter, for all licensed hunters in Minnesota. How many hunters must they survey in order to be 95% confident that their estimate is within 1 day of the true mean? (Use the fact that for a previous study conducted in 2000 the sample standard deviation was s = 10.)

$$n = \left( \frac{z_{\alpha/2} s^*}{m} \right)^2$$
$$= \left( \frac{(1.96)(10)}{1} \right)^2$$
$$= 384.16$$

So we need to survey 385 hunters.