

# Statistics

# What is Statistics?

- **Statistics** is the application of what we know to what we want to know.

*Is the S&P 500 a good model  
of the entire U.S. economy?*

*Does the population of Texas  
reflect the entire U.S. population?*

# Population vs. Sample

- These terms come up again and again
- **Population** is every member of a group we want to study
- **Sample** is a small set of (hopefully) random members of the population

# Parameter vs. Statistic

- A **parameter** is a characteristic of a population. Often we want to understand parameters.
- A **statistic** is a characteristic of a sample. Often we apply **statistical inferences** to the sample in an attempt to describe the population.

# Variable

- A **variable** is a characteristic that describes a member of the sample.
- Variables can be **discrete**, or **continuous**  
age salary  
gender birthplace

# Sampling

# Sampling

- One of the great benefits of statistical models is that a reasonably sized ( $>30$ ) random sample will almost always reflect the population.
- The challenge becomes, how do we select members randomly, and avoid bias?

# Sampling Bias

- There are several forms of bias:

## Selection Bias

Perhaps the most common, this type of bias favors those members of a population who are more inclined and able to answer polls.



# Sampling Bias

## Selection Bias

**Undercoverage Bias:** making too few observations or omitting entire segments of a population

# Sampling Bias

## Selection Bias

**Self-selection Bias:** people who volunteer may differ significantly from those in the population who don't

# Sampling Bias

## Selection Bias

**Healthy-user Bias:** the sample may come from a healthier segment of the overall population – people who walk/jog, work outside, follow healthier behaviors, etc.

# Undercoverage Bias

- A hospital survey of employees conducted during daytime hours
- Neglects to poll people who work the night shift.



# Self-Selection Bias

- An online survey about a sports team
- Only people who feel strongly about the team will answer the survey.



# Healthy-User Bias

- Polling customers at a fruit stand to study a connection between diet and health.
- Those polled likely do *other* things that have greater impact on their health.



# Sampling Bias

## Survivorship Bias

If a population improves over time, it may be due to lesser members leaving the population due to death, expulsion, relocation, etc.

# A Classic Puzzle

- At the start of World War I, British soldiers wore cloth caps.
- The war office became alarmed at the high number of head injuries, so they issued metal helmets to all soldiers.





# A Classic Puzzle

- They were surprised to find that the number of head injuries *increased* with the use of metal helmets.
- If the intensity of fighting was the same before and after the change, why should the number of head injuries increase?

# A Classic Puzzle

- Answer: You have to consider *all* of the data
- Before the switch, many things that gave head injuries to soldiers wearing metal helmets would have caused fatalities for those wearing cloth caps!



# Another Survivorship Example

- In World War II, statistician Abraham Wald worked for America's Statistical Research Group (SRG)



Adapted from [https://en.wikipedia.org/wiki/Abraham\\_Wald](https://en.wikipedia.org/wiki/Abraham_Wald)

# Another Survivorship Example

- One problem the SRG worked on was to examine the distribution of damage to aircraft by enemy fire and to advise the best placement of additional armor.



# Another Survivorship Example

- Common logic was to provide greater protection to parts that received more damage.



# Another Survivorship Example

- Wald saw it differently – he felt that damage must be more uniformly distributed and that aircraft that could return had been hit in less vulnerable parts.



# Another Survivorship Example

- Wald proposed that the Navy reinforce the areas where returning aircraft were undamaged, since those were areas that, if hit, would cause the plane to be lost!



# Types of Sampling

- Random
- Stratified Random
- Cluster



# Random Sampling

- As its name suggests, **random sampling** means every member of a population has an equal chance of being selected.
- However, since samples are usually much smaller than populations, there's a chance that entire demographics might be missed.

# Stratified Random Sampling

- **Stratified random sampling** ensures that groups within a population are adequately represented.
- First, divide the population into segments based on some characteristic.
- Members cannot belong to two groups at once.

# Stratified Random Sampling

- Next, take random samples from each group
- The size of each sample is based on the size of the group relative to the population.

# Stratified Random Sampling Example

- A company wants to conduct a survey of customer satisfaction
- They can only survey 10% of their customers
- They want to ensure that every age group is fairly represented

# Stratified Random Sampling Example

- The customer breakdown by age group is as follows:

20-29	30-39	40-49	50+	TOTAL
1400	4450	3200	950	10,000

  
stratum

  
strata

# Stratified Random Sampling Example

- To obtain a 10% sample,  
take 10% from each group:

20-29	30-39	40-49	50+	TOTAL
1400	4450	3200	950	10,000
140	445	320	95	1,000

# Clustering

- A third – and often less precise – method of sampling is **clustering**
- The idea is to break the population down into groups and sample a random selection of groups, or *clusters*.
- Usually this is done to reduce costs.

# Clustering Examples

- A marketing firm sends pollsters to a handful of neighborhoods  
(instead of canvassing an entire city)
- A researcher samples fishing boats that are in port on a particular day  
(also known as **convenience sampling**)



# Central Limit Theorem

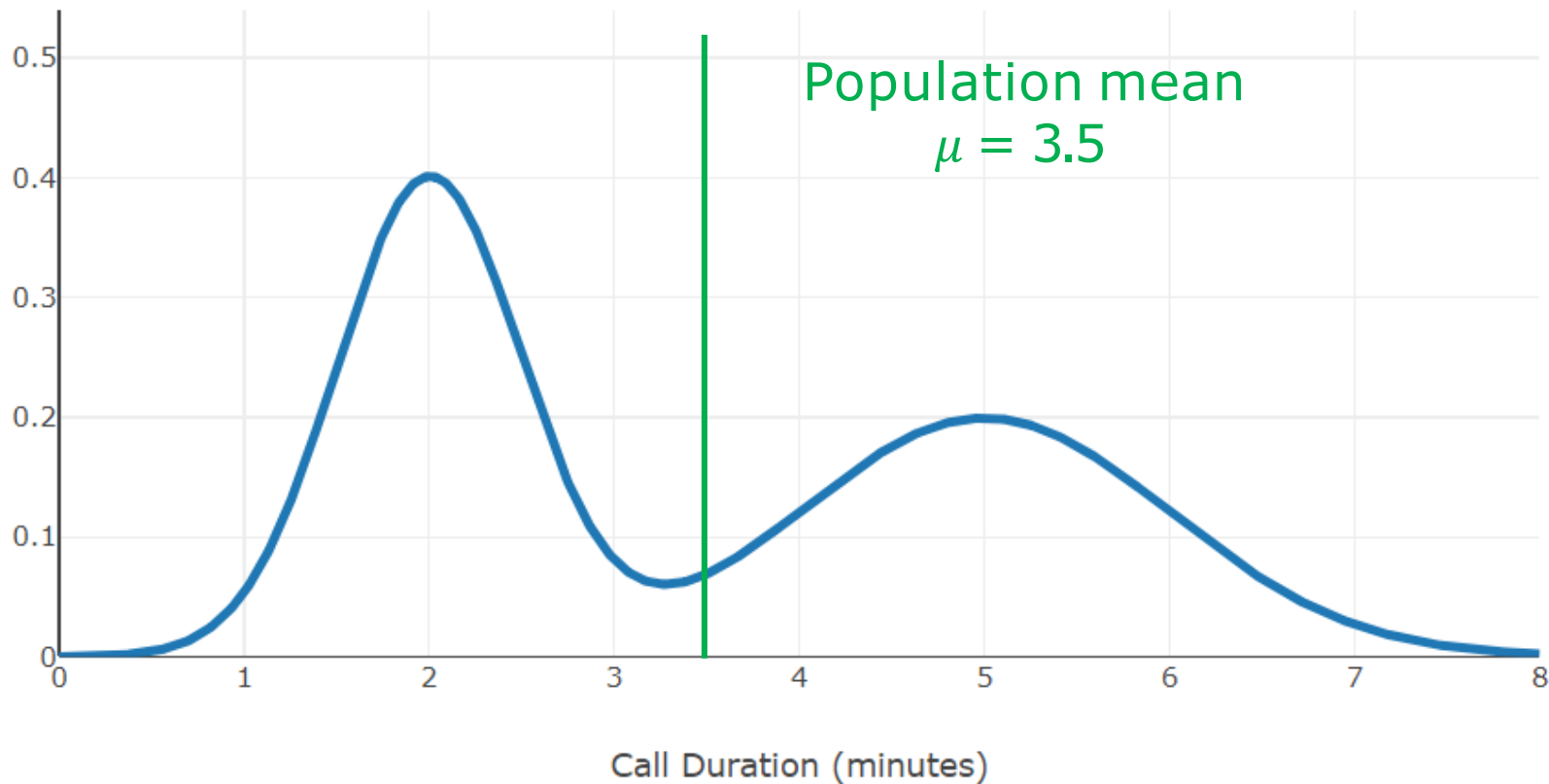
# Central Limit Theorem

- What makes sampling such a good statistical tool is the **Central Limit Theorem**
- Recall that a sample mean often varies from the population mean.
- The CLT considers a large number of random sample tests.

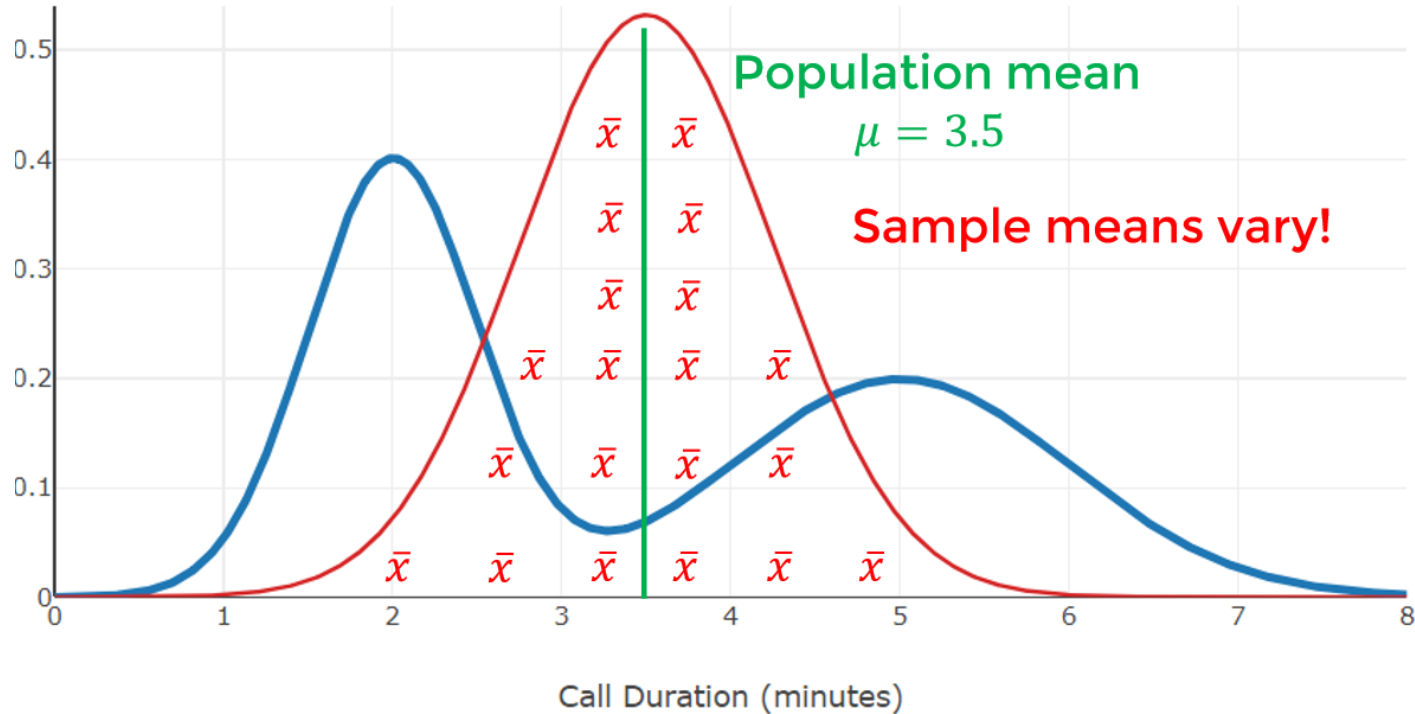
# Central Limit Theorem

- The CLT states that the mean values from a group of samples will be *normally distributed* about the population mean, even if the population itself is not normally distributed.
- That is, 95% of all sample means should fall within  $2\sigma$  of the population mean

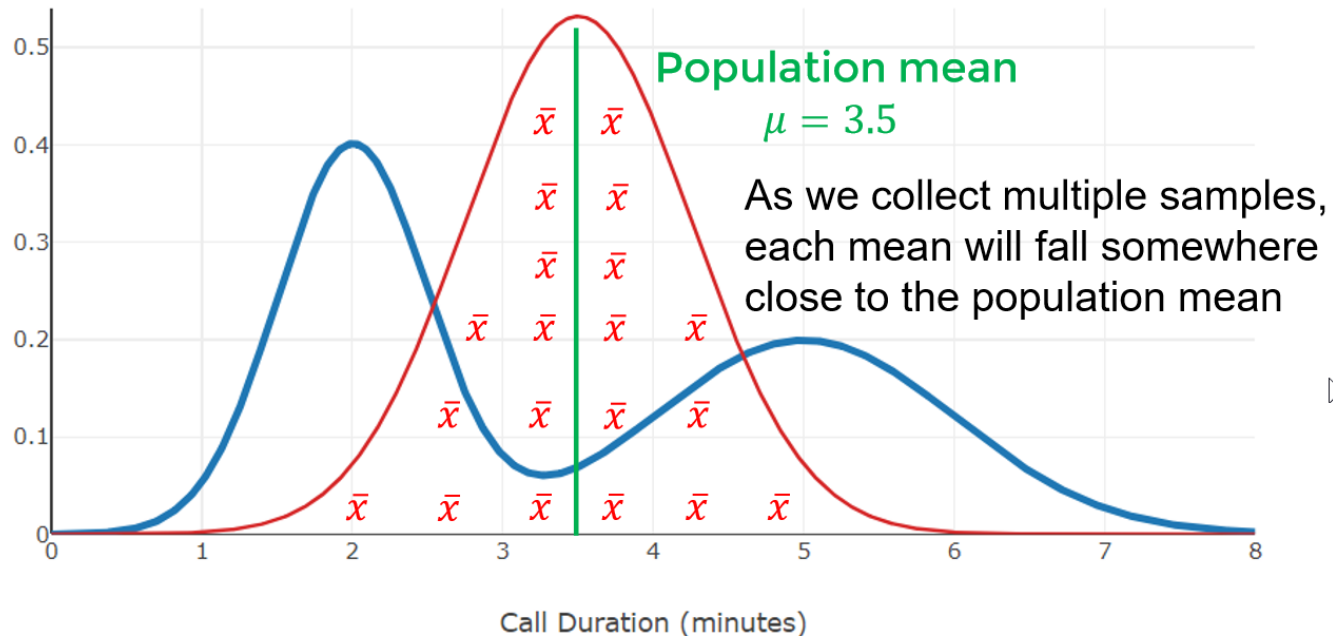
# Central Limit Theorem



# Central Limit Theorem



# Central Limit Theorem



# Proof of CLT Available on Wikipedia

For those who are curious, the full proof of the Central Limit Theorem is available at

[https://en.wikipedia.org/wiki/Central\\_limit\\_theorem](https://en.wikipedia.org/wiki/Central_limit_theorem)



W Central limit theorem - W x

Secure | [https://en.wikipedia.org/wiki/Central\\_limit\\_theorem](https://en.wikipedia.org/wiki/Central_limit_theorem)

Remarks [edit]

**Proof of classical CLT** [edit]

For a theorem of such fundamental importance to statistics and applied probability, the central limit theorem has a remarkably simple proof using characteristic functions.<sup>[16]</sup> It is similar to the proof of the (weak) law of large numbers.

As stated above, suppose  $\{X_1, \dots, X_n\}$  are independent and identically distributed random variables, each with mean  $\mu$  and finite variance  $\sigma^2$ . The sum  $X_1 + \dots + X_n$  has mean  $n\mu$  and variance  $n\sigma^2$ . Consider the random variable

$$Z_n = \frac{X_1 + \dots + X_n - n\mu}{\sqrt{n\sigma^2}} = \sum_{i=1}^n \frac{X_i - \mu}{\sqrt{n\sigma^2}} = \sum_{i=1}^n \frac{1}{\sqrt{n}} Y_i,$$

where in the last step we defined the new random variables  $Y_i = \frac{X_i - \mu}{\sigma}$ , each with zero mean and unit variance ( $\text{var}(Y) = 1$ ). The characteristic function of  $Z_n$  is given by

$$\varphi_{Z_n}(t) = \varphi_{\sum_{i=1}^n \frac{1}{\sqrt{n}} Y_i}(t) = \varphi_{Y_1}\left(\frac{t}{\sqrt{n}}\right) \varphi_{Y_2}\left(\frac{t}{\sqrt{n}}\right) \cdots \varphi_{Y_n}\left(\frac{t}{\sqrt{n}}\right) = \left[\varphi_{Y_1}\left(\frac{t}{\sqrt{n}}\right)\right]^n,$$

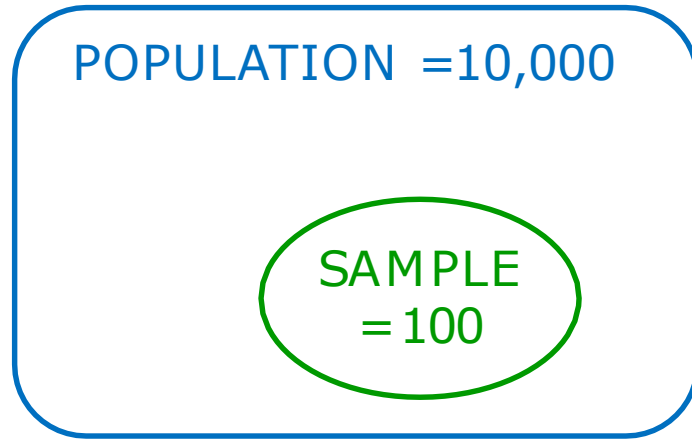
# Standard Error



# Standard Error

- Let's quickly review terminology
- Let's say we have a **population** of voters
- It is unrealistic to poll the entire population, so we poll a **sample**
- We calculate a **statistic** from that sample that lets us estimate a **parameter** of the population

# Standard Error



$N$  = # population members

$P$  = population parameter

$\sigma$  = pop. standard deviation

$n$  = # sample members

$\hat{p}$  = sample statistic

$SE_{\hat{p}}$  = standard error of the  
sample

# Standard Error

- If for the population of Australia, the mean height is 5'9", and for our 100-person sample the mean height is 5'10", Then

$$P = 5'9"$$

$$\hat{p} = 5'10"$$

$$SE_{\hat{p}} = \text{Standard Error of the Mean}$$

POPULATION = 10,000

SAMPLE  
= 100



# Standard Error of the Mean

- Where the population standard deviation describes how wide individual values stray from the population mean, the Standard Error of the Mean describes how far a sample mean may stray from the population mean.

# Standard Error of the Mean

- If the population standard deviation  $\sigma$  is known, then the sample standard error of the mean can be calculated as:

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

# Standard Error Exercise

- An IQ Test is designed to have a mean score of 100 with a standard deviation of 15 points.
- If a sample of 10 scores has a mean of 104, can we assume they come from the general population?



# Standard Error Exercise

- Sample of 10 IQ Test scores:

$$n = 10 \quad \bar{x} = 104 \quad \sigma = 15$$

$$SE_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{15}{\sqrt{10}} = 4.743$$

- 68% of 10-item sample means are expected to fall between 95.257 and 104.743

# Confidence Intervals

POPULATION = 10,000

SAMPLE  
= 100

"We can say with a 95% **confidence level** that the population parameter lies within a **confidence interval** of plus-or-minus two standard errors of the sample statistic"

$N$  = # population members

$P$  = population parameter

$\sigma$  = pop. standard deviation

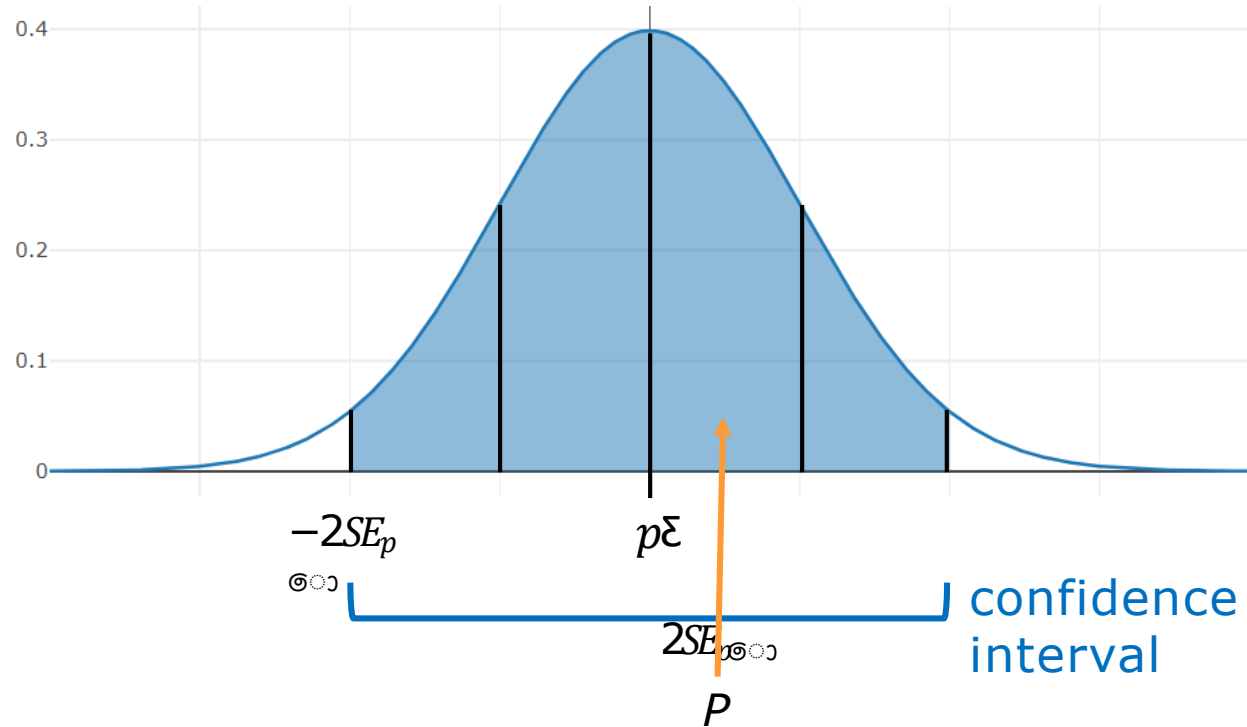
$n$  = # sample members

$\hat{p}$  = sample statistic

$SE_{\hat{p}}$  = standard error of the sample



# Confidence Intervals



# Point Estimators

- In the above example, the sample statistic  $\hat{p}$  is a **point estimator** of the population parameter  $P$ .

# Hypothesis Testing

# Hypothesis Testing

- **Hypothesis Testing** is the application of statistical methods to real-world questions.
- We start with an assumption, called the **null hypothesis**
- We run an experiment to test this null hypothesis

# Hypothesis Testing

- Based on the results of the experiment, we either **reject** or **fail to reject** the null hypothesis
- If the null hypothesis is rejected, then we say the data supports another, mutually exclusive **alternate hypothesis**
- We never “PROVE” a hypothesis!

# Framing the Hypothesis

- How do we frame the question that forms our null hypothesis?
- At the start of the experiment, the null hypothesis is assumed to be true.
- If the data fails to support the null hypothesis, only then can we look to an alternative hypothesis

# Framing the Hypothesis

If testing something assumed to be true,  
the null hypothesis can reflect the assumption:

Claim: *"Our product has an average  
shipping weight of 3.5kg."*

Null hypothesis: average weight = 3.5kg

Alternate hypothesis: average weight  $\neq$  3.5kg

# Framing the Hypothesis

If testing a claim we *want* to be true,  
but can't assume, we test its opposite:

Claim: *"This prep course improves  
test scores."*

Null hypothesis:            old scores  $\geq$  new scores

Alternate hypothesis: old scores  $<$  new scores



# Framing the Hypothesis

The null hypothesis should contain an equality ( $=, \leq, \geq$ ):

average shipping weight  $= 3.5\text{kg}$   $H_0: \mu = 3.5$

The alternate hypothesis should not have an equality ( $\neq, <, >$ ):

average shipping weight  $\neq 3.5\text{kg}$   $H_1: \mu \neq 3.5$

# Framing the Hypothesis

The null hypothesis should contain an equality ( $=, \leq, \geq$ ):

old scores  $\geq$  new scores

$$H_0: \mu_0 \geq \mu_1$$

The alternate hypothesis should not have an equality ( $\neq, <, >$ ):

old scores  $<$  new scores

$$H_1: \mu_0 < \mu_1$$

# Hypothesis Testing

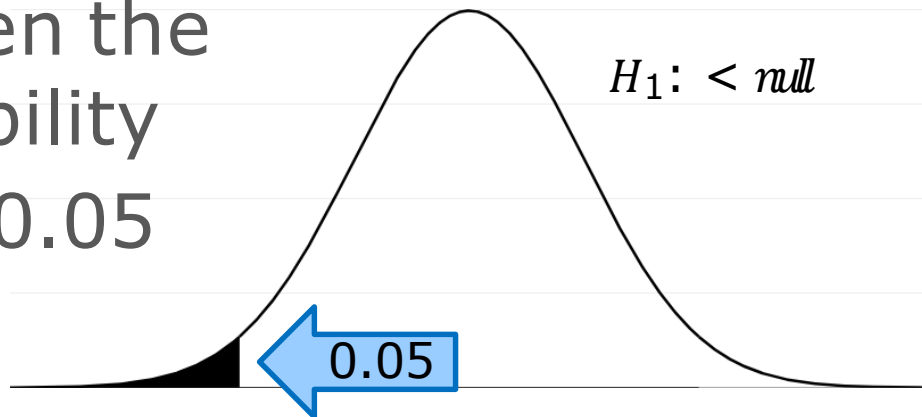
- So what lets us reject or fail to reject the null hypothesis?

# Hypothesis Testing

- We run an experiment and record the result.
- Assuming our null hypothesis is valid, if the probability of observing these results is very small (inside of 0.05) then we reject the null hypothesis.
- Here 0.05 is our level of significance  
 $\alpha = 0.05$

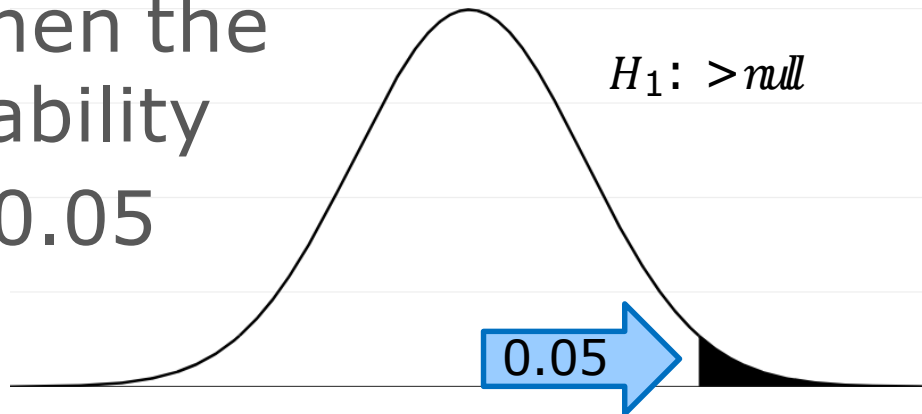
# Hypothesis Testing -Tails

- The level of significance  $\alpha$  is the area inside the *tail(s)* of our null hypothesis.
- If  $\alpha = 0.05$  and the alternative hypothesis is *less than* the null, then the left-tail of our probability curve has an area of 0.05



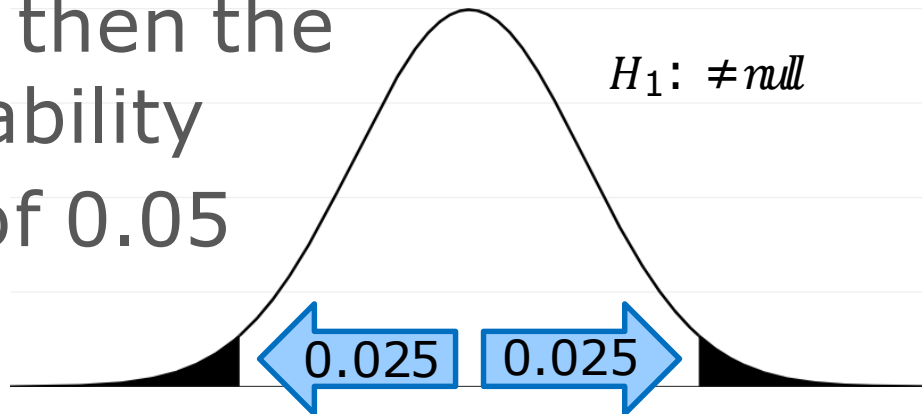
# Hypothesis Testing -Tails

- The level of significance  $\alpha$  is the area inside the *tail(s)* of our null hypothesis.
- If  $\alpha = 0.05$  and the alternative hypothesis is *more than* the null, then the right-tail of our probability curve has an area of 0.05



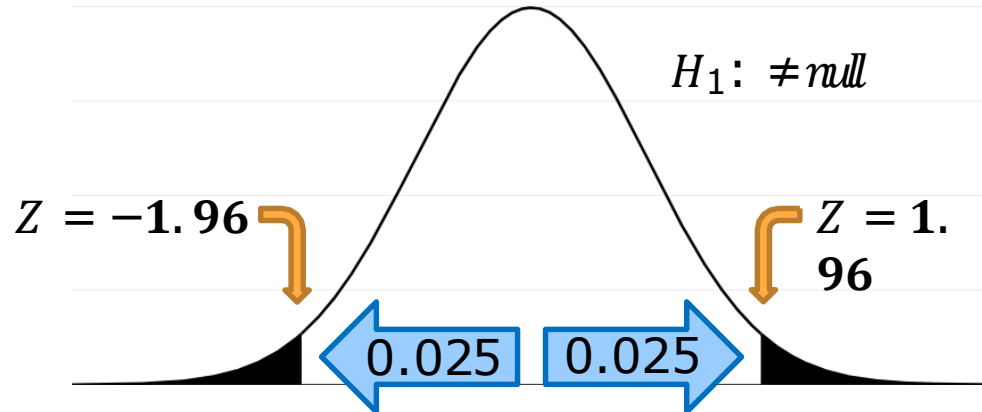
# Hypothesis Testing -Tails

- The level of significance  $\alpha$  is the area inside the *tail(s)* of our null hypothesis.
- If  $\alpha = 0.05$  and the alternative hypothesis is *not equal to* the null, then the two tails of our probability curve *share* an area of 0.05



# Hypothesis Testing - Tails

- These areas establish our **critical values** or Z-scores:





# Tests of Mean vs. Proportion

- There are two main types of tests:
- Test of Means
- Test of Proportions

# Tests of Mean vs. Proportion

- Each of these two types of tests has their own test statistic to calculate.
- Let's review the situation for each test before we work through some examples in the upcoming lectures.

# Tests of Mean vs. Proportion

- Mean

when we look to find an **average**, or specific value in a population we are dealing with means

- Proportion

whenever we say something like "**35%**" or "**most**" we are dealing with proportions

# Test Statistics

- When working with means:

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

← assumes we know  
the population  
standard deviation

- When working with proportions:

$$Z = \frac{\hat{p} - p}{\sqrt{\frac{p \cdot q}{n}}} = \frac{\hat{p} - p}{\sqrt{\frac{p \cdot (1 - p)}{n}}}$$

# Hypothesis Testing – P-value Test

In a **traditional test**:

- take the level of significance  $\alpha$
- use it to determine the critical value
- compare the test statistic to the critical value

In a **P-value test**:

- take the test statistic
- use it to determine the P-value
- compare the P-value to the level of significance  $\alpha$

# Hypothesis Testing – P-value Test

“If the P-value is low,  
the null must go!”

reject  $H_0$

“If the P-value is high,  
the null must fly!”

fail to reject  
 $H_0$

# Testing Example

## Exercise #1

# Testing Exercise #1- Mean

- For this next example we'll work in the left-hand side of the probability distribution, with negative  $z$ -scores
- We'll show how to run the hypothesis test using the traditional method, and then with the P-value method



# Testing Exercise #1- Mean

- A company is looking to improve their website performance.

$$\mu = 3.125$$
$$\sigma = 0.700$$

- Currently pages have a mean load time of 3.125 seconds, with a standard deviation of 0.700 seconds.
- They hire a consulting firm to improve load times.

# Testing Exercise #1- Mean

- Management wants a 99% confidence level
- A sample run of 40 of the new pages has a mean load time of 2.875 seconds.
- Are these results statistically faster than before?

$$\mu = 3.125$$

$$\sigma = 0.700$$

$$\alpha = 0.01$$

$$n = 40$$

$$\bar{x} = 2.875$$

# Testing Solution #1- Mean

1. State the null hypothesis:

$$H_0: \mu \geq 3.125$$

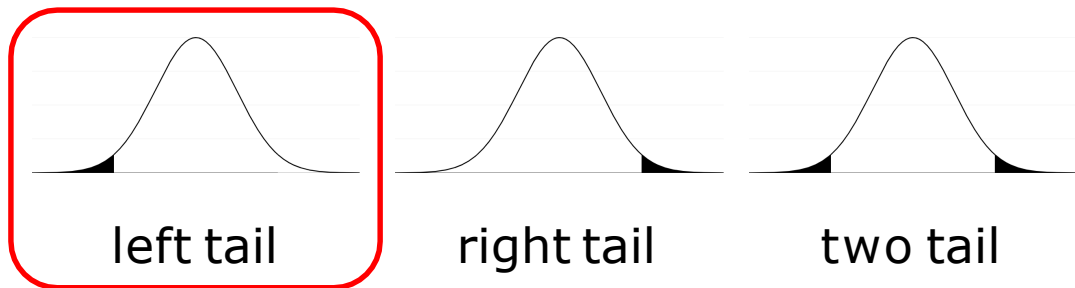
2. State the alternative hypothesis:

$$H_1: \mu < 3.125$$

3. Set a level of significance:

$$\alpha = 0.01$$

4. Determine the test type:



# Testing Solution #1- Mean

TRADITIONAL METHOD:

5. Test Statistic:

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{2.875 - 3.125}{0.7/\sqrt{40}} = -2.259$$

6. Critical Value:

*z-table lookup on 0.01*     $z = -2.325$

$$\mu = 3.125$$

$$\sigma = 0.700$$

$$\alpha = 0.01$$

$$n = 40$$

$$\bar{x} = 2.875$$

$$Z = -2.259$$

$$z = -2.325$$

z	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
-3.4	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0002
-3.3	.0005	.0005	.0005	.0004	.0004	.0004	.0004	.0004	.0004	.0003
-3.2	.0007	.0007	.0006	.0006	.0006	.0006	.0006	.0005	.0005	.0005
-3.1	.0010	.0009	.0009	.0009	.0008	.0008	.0008	.0008	.0007	.0007
-3.0	.0013	.0013	.0013	.0012	.0012	.0011	.0011	.0011	.0010	.0010
-2.9	.0019	.0018	.0018	.0017	.0016	.0016	.0015	.0015	.0014	.0014
-2.8	.0026	.0025	.0024	.0023	.0023	.0022	.0021	.0021	.0020	.0019
-2.7	.0035	.0034	.0033	.0032	.0031	.0030	.0029	.0028	.0027	.0026
-2.6	.0047	.0045	.0044	.0043	.0041	.0040	.0039	.0038	.0037	.0036
-2.5	.0062	.0060	.0059	.0057	.0055	.0054	.0052	.0051	.0049	.0048
-2.4	.0082	.0080	.0078	.0075	.0073	.0071	.0069	.0068	.0066	.0064
-2.3	.0107	.0104	.0102	.0099	.0096	.0094	.0091	.0089	.0087	.0084
-2.2	.0139	.0136	.0132	.0129	.0125	.0122	.0119	.0116	.0113	.0110
-2.1	.0179	.0174	.0170	.0166	.0162	.0158	.0154	.0150	.0146	.0143
-2.0	.0228	.0222	.0217	.0212	.0207	.0202	.0197	.0192	.0188	.0183
-1.9	.0287	.0281	.0274	.0268	.0262	.0256	.0250	.0244	.0239	.0233
-1.8	.0359	.0351	.0344	.0336	.0329	.0322	.0314	.0307	.0301	.0294
-1.7	.0446	.0436	.0427	.0418	.0409	.0401	.0392	.0384	.0375	.0367
-1.6	.0548	.0537	.0526	.0516	.0505	.0495	.0485	.0475	.0465	.0455
-1.5	.0668	.0655	.0643	.0630	.0618	.0606	.0594	.0582	.0571	.0559
-1.4	.0808	.0793	.0778	.0764	.0749	.0735	.0721	.0708	.0694	.0681
-1.3	.0968	.0951	.0934	.0918	.0901	.0885	.0869	.0853	.0838	.0823
-1.2	.1151	.1131	.1112	.1093	.1075	.1056	.1038	.1020	.1003	.0985
-1.1	.1357	.1335	.1314	.1292	.1271	.1251	.1230	.1210	.1190	.1170
-1.0	.1587	.1562	.1539	.1515	.1492	.1469	.1446	.1423	.1401	.1379
-0.9	.1841	.1814	.1788	.1762	.1736	.1711	.1685	.1660	.1635	.1611

# Testing Solution #1- Mean

TRADITIONAL METHOD:

7. Fail to Reject the Null Hypothesis

Since  $-2.259 > -2.325$ , the  
test statistic falls outside  
the rejection region

We can't say that the new web  
pages are statistically faster.

$$\mu = 3.125$$

$$\sigma = 0.700$$

$$\alpha = 0.01$$

$$n = 40$$

$$\bar{x} = 2.875$$

$$Z = -2.259$$

$$z = -2.325$$

# Testing Solution #1- Mean

P-VALUE METHOD:

5. Test Statistic:

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{2.875 - 3.125}{0.7/\sqrt{40}} = -2.259$$

6. P-Value:

*z-table lookup on -2.26*  $P = 0.0119$

$$\mu = 3.125$$

$$\sigma = 0.700$$

$$\alpha = 0.01$$

$$n = 40$$

$$\bar{x} = 2.875$$

$$Z = -2.259$$

$$P = 0.0119$$

# Testing Solution #1- Mean

P-VALUE METHOD:

7. Fail to Reject the Null Hypothesis

Since  $0.0119 > 0.01$ , the  
P-value is greater than the  
level of significance  $\alpha$

We can't say that the new web  
pages are statistically faster.

$$\mu = 3.125$$

$$\sigma = 0.700$$

$$\alpha = 0.01$$

$$n = 40$$

$$\bar{x} = 2.875$$

$$Z = -2.259$$

$$P = 0.0119$$



# Testing Example

## Exercise #2

## Testing Exercise #2 - Proportion

- A video game company surveys 400 of their customers and finds that 58% of the sample are teenagers.
- Is it fair to say that most of the company's customers are teenagers?

# Testing Solution #2 - Proportion

1. Set the null hypothesis:  $H_0: P \leq 0.50$
2. Set the alternative hypothesis:  $H_1: P > 0.50$
3. Calculate the test statistic:

$$Z = \frac{\hat{p} - p}{\sqrt{\frac{p \cdot q}{n}}} = \frac{0.58 - 0.50}{\sqrt{\frac{0.50(1 - 0.50)}{400}}} = \frac{0.08}{0.025} = 3.2$$

# Testing Solution #2 - Proportion

4. Set a significance level:

$$\alpha =$$

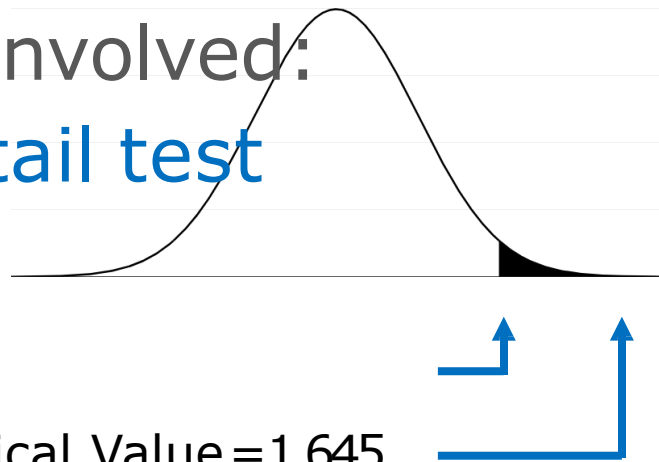
0.05

5. Decide what type of tail is involved:

$H_1: P > 0.50$  means a right-tail test

6. Look up the critical value:

$$Z = 1.645$$

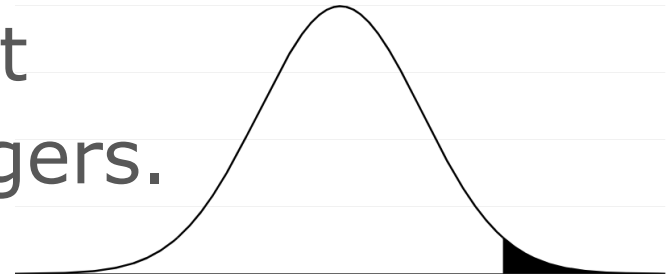


Critical Value = 1.645

Test Statistic = 3.2

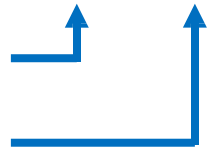
# Testing Solution #2 - Proportion

7. Based on the sample,  
we reject the null hypothesis,  
and support the claim that  
most customers are teenagers.



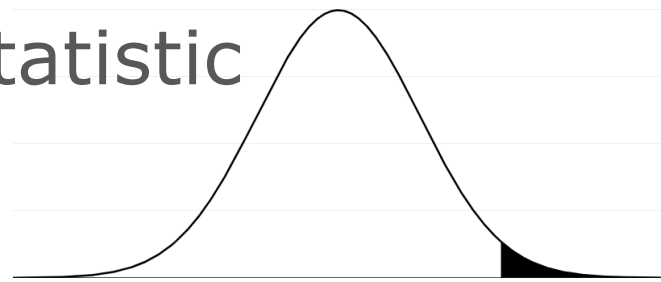
Critical Value = 1.645

Test Statistic = 3.2



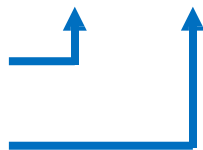
# Testing Solution #2 - Proportion

NOTE: The size of the sample matters!  
If we had started with a sample size of 40 instead of 400, our test statistic would have been only 1.01, and we would fail to reject the null hypothesis.



Critical Value = 1.645

Test Statistic = 3.2



# Type 1 and Type 2 Errors

# Type I and Type II Errors

- Often in medical fields (and other scientific fields) hypothesis testing is used to test against results where the "truth" is already known.
- For example, testing a new diagnostic test for cancer for patients you have already successfully diagnosed by other means.



# Type I and Type II Errors

- In this situation, you already know if the Null Hypothesis is True or False.
- In these situations where you already know the "truth", then you would know its possible to commit an error with your results .

# Type I and Type II Errors

- This type of analysis is common enough that these errors already have specific names:
- Type I Error
- Type II Error

# Type I and Type II Errors

- If we reject a null hypothesis that should have been supported, we've committed a **Type I Error**

$H_0$ : *There is no fire*

Pull the fire alarm,  
only to find out there  
really was no fire.



# Type I and Type II Errors

- If we fail to reject a null hypothesis that should have been rejected we've committed a **Type II Error**

$H_0$ : *There is no fire*

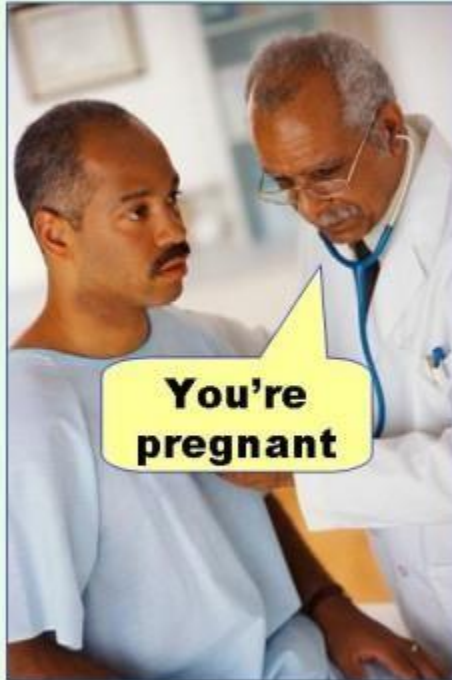
Don't pull the fire alarm, only to find there really is a fire.



$H_0$ : Not pregnant

$H_1$ : Are pregnant

**Type I error**  
(false positive)



**Type II error**  
(false negative)



# Student's T-Distribution

# Student's T-Distribution

- Developed by William Sealy Gossett while he was working at Guinness Brewery
- Published under the pseudonym "Student" as Guinness wouldn't let him use his name.
- Goal was to select the best barley from small samples, when the population standard deviation was unknown!

# Purpose of a t-test

- Using the t-table, the Student's t-test determines if there is a significant difference between two sets of data
- Due to variance and outliers, it's not enough just to compare mean values
- A t-test also considers sample variances



# Types of Student's t-test

- One-sample t-test

Tests the null hypothesis that the population mean is equal to a specified value  $\mu$  based on a sample mean  $\bar{x}$

# Types of Student's t-test

- Independent two-sample t-test  
Tests the null hypothesis that two sample means  $\mu_1$  and  $\mu_2$  are equal

# Types of Student's t-test

- **Dependent, paired-sample t-test**

Used when the samples are dependent:

- one sample has been tested twice (repeated measurements)
- two samples have been matched or "paired"

# One-Sample Student's t-test

- Calculate the t-statistic

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

$\bar{x}$  = sample mean

$\mu$  = population mean

$s$  = sample standard error

$n$  = sample size

# One-Sample Student's t-test

- Compare to a t-score

$$t \leq t_{n-1,\alpha}$$

$t$  = t-statistic

$t_{n-1,\alpha}$  = t-critical

$n - 1$  = degrees of freedom

$\alpha$  = significance level

# Independent Two-Sample t-test

The calculation of the t-statistic differs slightly for the following scenarios:

- equal sample sizes, equal variance
- unequal sample sizes, equal variance
- equal or unequal sample sizes, unequal variance

# Independent Two-Sample t-test

- Calculate the t-statistic

$$t = \frac{\text{signal}}{\text{noise}} = \frac{\text{difference in means}}{\text{sample variability}} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$\bar{x}_1, \bar{x}_2$  = sample means

$s_1^2, s_2^2$  = sample variances

$n_1, n_2$  = sample sizes

# Independent Two-Sample t-test

- Compare to a t-score

$$t \leq t_{df,\alpha}$$

$t$  = t-statistic

$t_{df,\alpha}$  = t-critical

$df$  = degrees of freedom

$\alpha$  = significance level

Since we have two, potentially unequal-sized samples with different variances, determining the degrees of freedom is a little more complicated.



# Degrees of Freedom

- The Satterthwaite Formula:

$$df = \frac{\left( \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2}{\frac{1}{n_1 - 1} \left( \frac{s_1^2}{n_1} \right)^2 + \frac{1}{n_2 - 1} \left( \frac{s_2^2}{n_2} \right)^2}$$

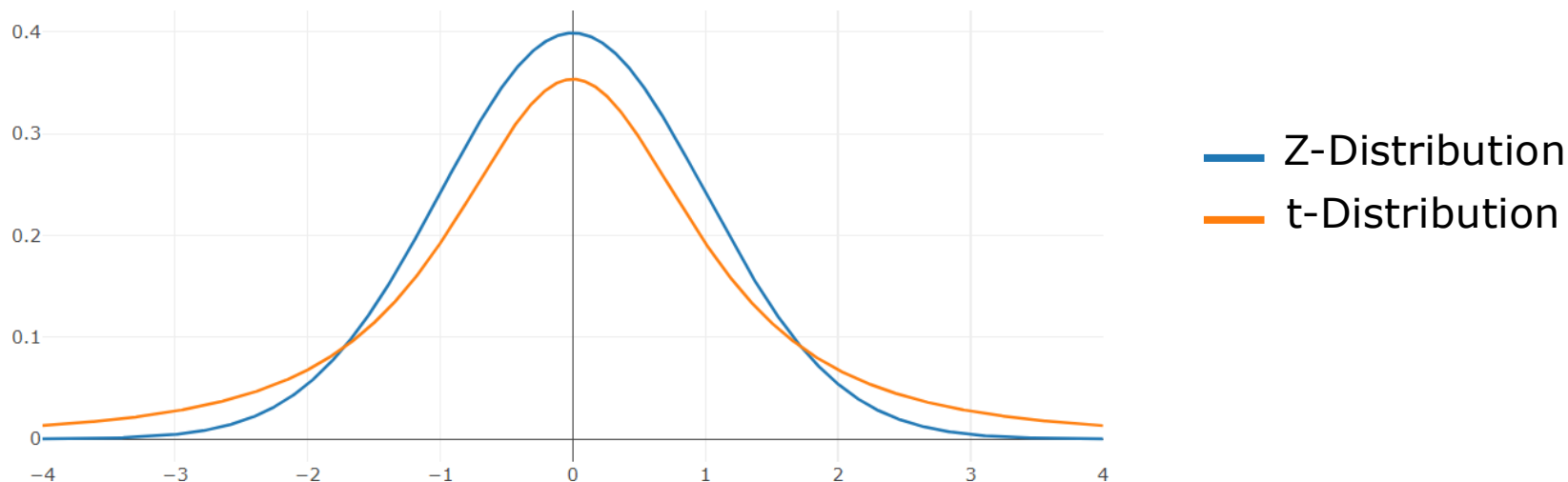
# Degrees of Freedom

- The General Formula:

$$df = n_1 + n_2 - 2$$

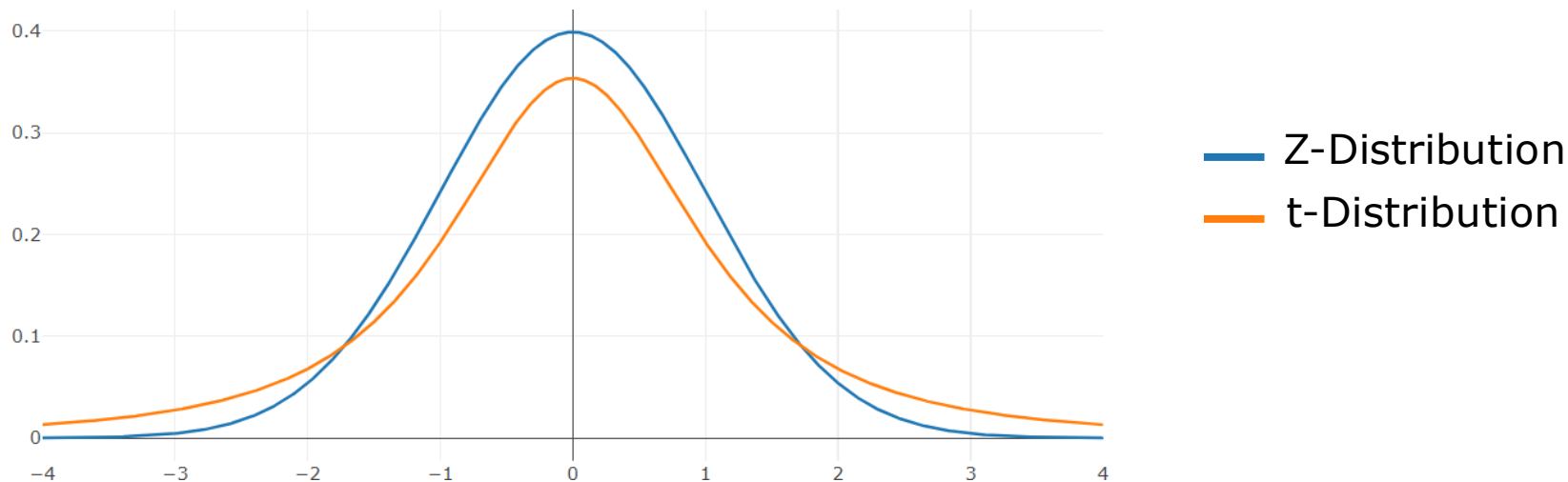
# Student's t-Distribution

- t-Distributions have fatter tails than normal Z-Distributions



# Student's t-Distribution

- They approach a normal distribution as the degrees of freedom increase.



# Student's T-Distribution Example Exercise

# Student's t-test Example

An auto manufacturer has two plants that produce the same car.



# Student's t-test Example

They are forced to close one of the plants.



# Student's t-test Example

The company wants to know if there's a significant difference in production between the two plants.





# Student's t-test Example

Daily production over the same 10 days is as follows:



Plant A	Plant B
1184	1136
1203	1178
1219	1212
1238	1193
1243	1226
1204	1154
1269	1230
1256	1222
1156	1161
1248	1148

# Student's t-test Example

First compare sample means

$$\bar{x}_A - \bar{x}_B = 1222 - 1186 = 36$$

From this sample, it looks like Plant A produces 36 more car per day than Plant B

	Plant A	Plant B
	1184	1136
	1203	1178
	1219	1212
	1238	1193
	1243	1226
	1204	1154
	1269	1230
	1256	1222
	1156	1161
	1248	1148
	$\bar{x}_A$	$\bar{x}_B$
Mean	1222	1186

# Student's t-test Example

Is 36 more cars enough to say that the plants are different?

$$H_0: X_A \leq X_B$$

$$H_1: X_A > X_B$$

one-tailed test

$$(10 + 10 - 2) = 18 \text{ degrees of freedom}$$

	Plant A	Plant B
	1184	1136
	1203	1178
	1219	1212
	1238	1193
	1243	1226
	1204	1154
	1269	1230
	1256	1222
	1156	1161
	1248	1148
	$\bar{x}_A$	$\bar{x}_B$
Mean	1222	1186

# Student's t-test Example

Compute the variance

A	(x-1222)	(x-1222) <sup>2</sup>
1184	-38	1444
1203	-19	361
1219	-3	9
1238	16	256
1243	21	441
1204	-18	324
1269	47	2209
1256	34	1156
1156	-66	4356
1248	26	676
		<b>11232</b>

$$s^2 = \frac{\Sigma(x - \bar{x})^2}{n - 1}$$

$\Sigma(x-1222)^2$	11232
$\frac{\Sigma(x-1222)^2}{9}$	<b>1248</b>

Plant A	Plant B
1184	1136
1203	1178
1219	1212
1238	1193
1243	1226
1204	1154
1269	1230
1256	1222
1156	1161
1248	1148
	$\bar{x}_A$
Mean	1222
Variance	1248

# Student's t-test Example

Compute the t-value

$$\begin{aligned}
 &= \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \\
 &= \frac{36}{\sqrt{\frac{1248}{10} + \frac{1246}{10}}} = \frac{36}{15.79} \\
 &= 2.28
 \end{aligned}$$

	Plant A	Plant B
	1184	1136
	1203	1178
	1219	1212
	1238	1193
	1243	1226
	1204	1154
	1269	1230
	1256	1222
	1156	1161
	1248	1148
	$\bar{x}_A$	$\bar{x}_B$
Mean	1222	1186
Variance	1248	1246

# Student's t-test Example

Look up our critical value from a t-table

a one-tailed test

95% confidence

18 degrees of  
freedom

critical value = 1.734

cum. prob	$t_{.90}$	$t_{.95}$	$t_{.975}$	$t_{.99}$	$t_{.995}$
one-tail	0.10	0.05	0.025	0.01	0.005
two-tails	0.20	0.10	0.05	0.02	0.01
17	1.333	1.740	2.110	2.567	2.898
18	1.330	1.734	2.101	2.552	2.878
19	1.328	1.729	2.093	2.539	2.861

# Student's t-test Example

Compare our t-value (2.28) to the critical value (1.734):

$$2.28 > 1.734$$

since our computed t-value is *greater* than the critical value, we reject the null hypothesis.

Plant A	Plant B
1184	1136
1203	1178
1219	1212
1238	1193
1243	1226
1204	1154
1269	1230
1256	1222
1156	1161
1248	1148

# Student's t-test Example

We believe with 95% confidence that Plant A produces more cars per day than Plant B.

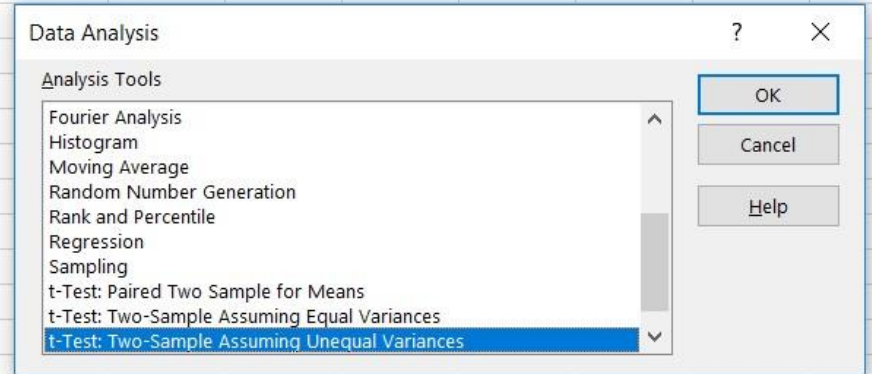
We decide to close Plant B.





# Student's t-Test with Excel

	A	B	C	D	E	F	G	H	I	J	K
1	t-Test: Two-Sample Assuming Unequal Variances										
2											
3		Variable 1	Variable 2								
4	Mean	1186	1222								
5	Variance	1246	1248								
6	Observations	10	10								
7	Hypothesized Mean Difference	0									
8	df	18									
9	t Stat	-2.279577051									
10	P(T<=t) one-tail	0.017522528									
11	t Critical one-tail	1.734063607									
12	P(T<=t) two-tail	0.035045056									
13	t Critical two-tail	2.10092204									
14											



# Student's t-Test with Python

```
>>> from scipy.stats import ttest_ind
>>> a = [1184, 1203, 1219, ... 1248]
>>> b = [1136, 1178, 1212, ... 1148]
>>> ttest_ind(a,b).statistic
2.2795770510504845
>>> ttest_ind(a,b).pvalue/2
0.017522528133638322
```