



Data Science

Statistics



Agenda

01

Introduction to Statistics

02

Categories of Statistics

03

Measures of Central
Tendency & Spread

04

Probability

05

Covariance & Correlation

06

Hypothesis Testing

Introduction to Statistics

Statistics is a very broad subject, with applications in a vast number of different fields. Generally, one can say that statistics is the methodology for collecting, analyzing, interpreting and drawing conclusions from information.

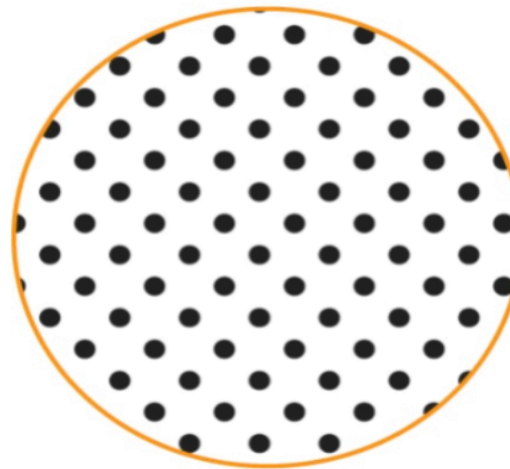


Statistical Terminology

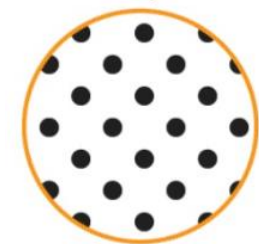
Statistical Terminology

What is a
population and
what is a
sample?

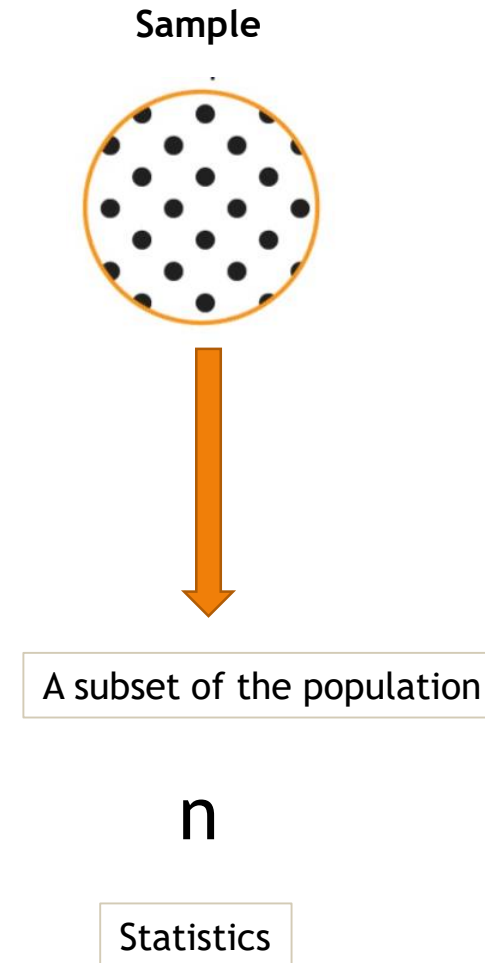
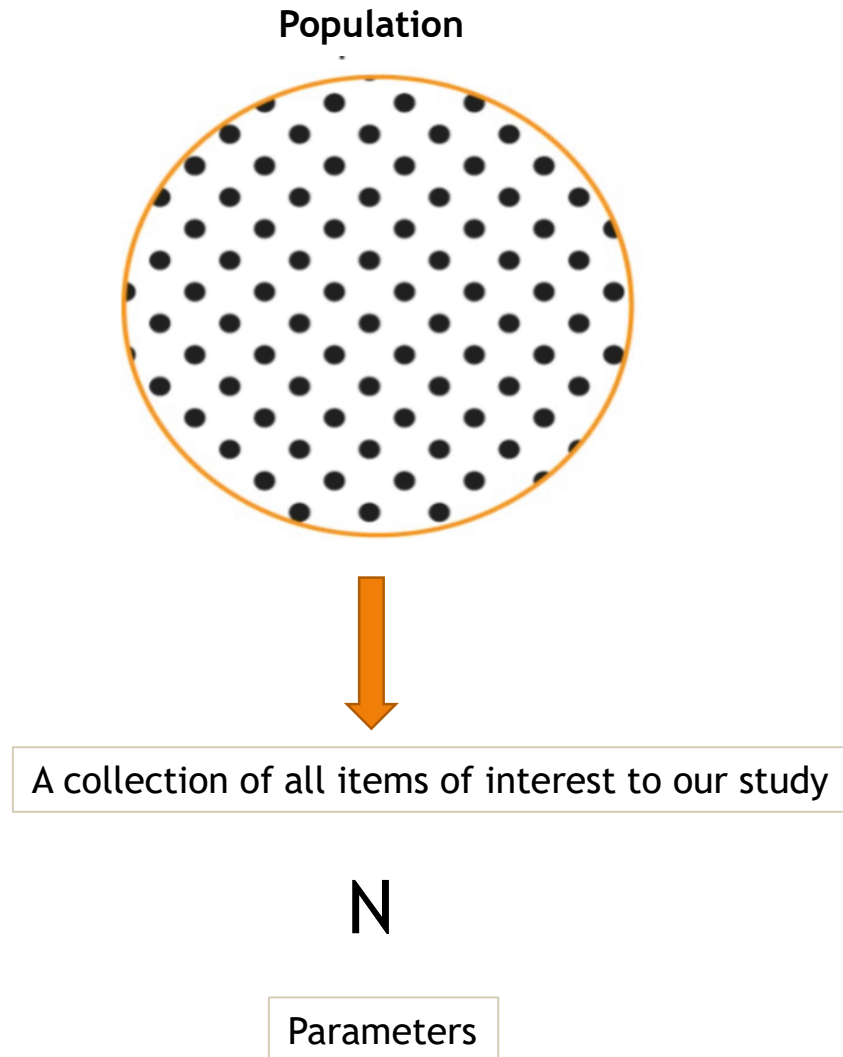
A **population** includes
all of the elements
from a set of data.



A **sample** consists one or more
observations drawn from the
population.



Statistical Terminology



Statistical Terminology

What is a census and what is a survey?



Census

Gathering data from the whole population of interest



Survey

Gathering data from a sample in order to make conclusions about the population

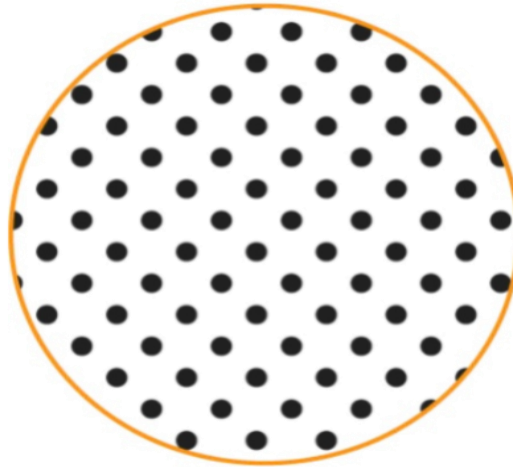


Statistical Terminology

What is a
parameter and
what is a
statistic?



Population

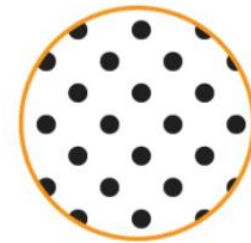


Parameter

A descriptive measure of the population

Example: Population mean, population variance, population standard deviation, etc.

Sample

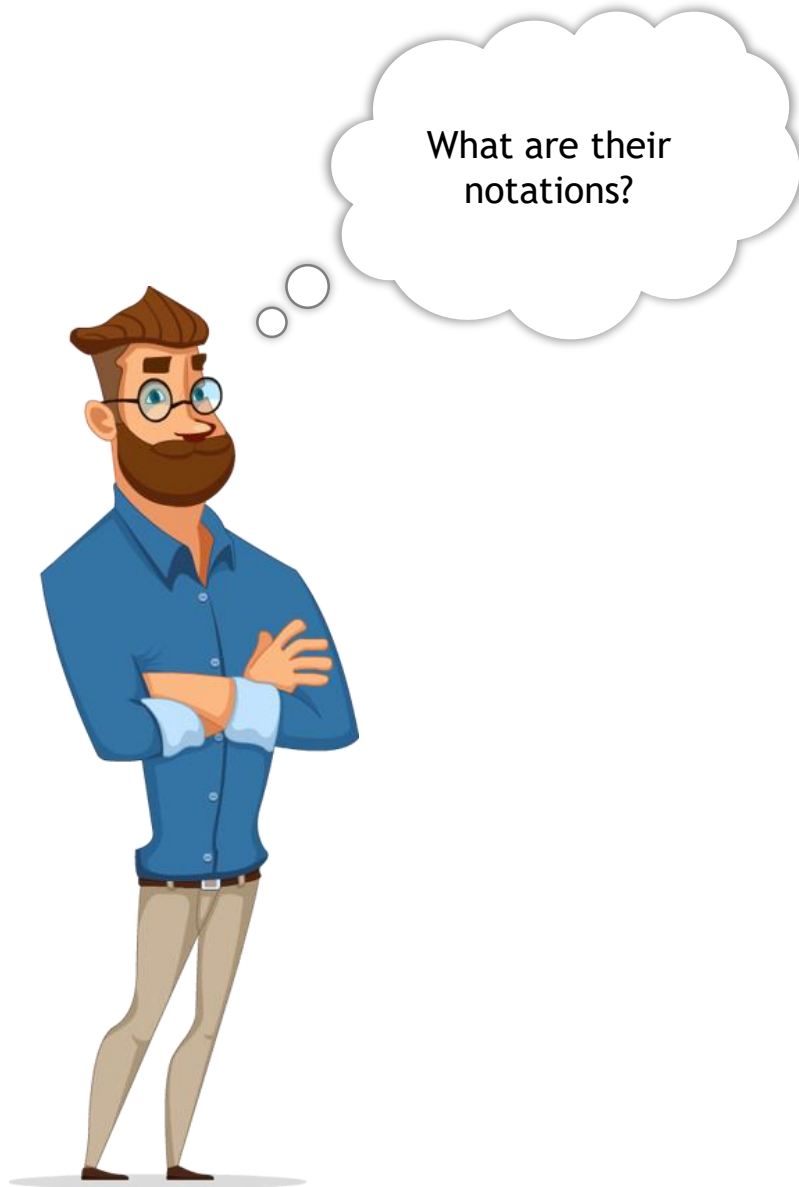


Statistic

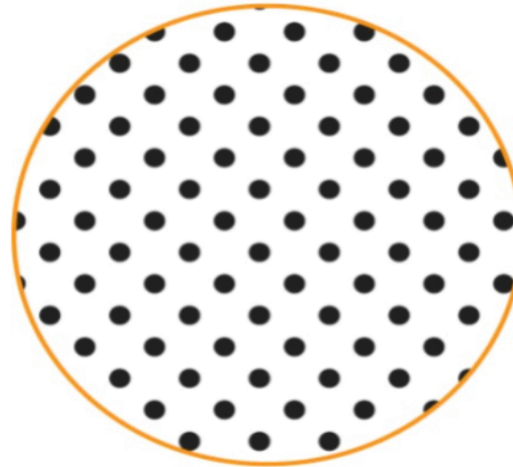
A descriptive measure of the sample

Example: Sample mean, sample variance, sample standard deviation, etc.

Statistical Terminology



Population



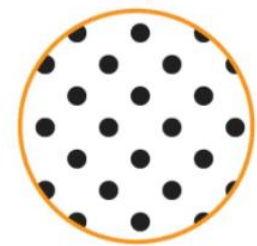
Greek - Population Parameter

Mean - μ

Variance - σ^2

Standard Deviation - σ

Sample



Roman - Sample Statistic

Mean - \bar{x}

Variance - s^2

Standard Deviation - s

Sampling

Sampling

What is
Sampling?

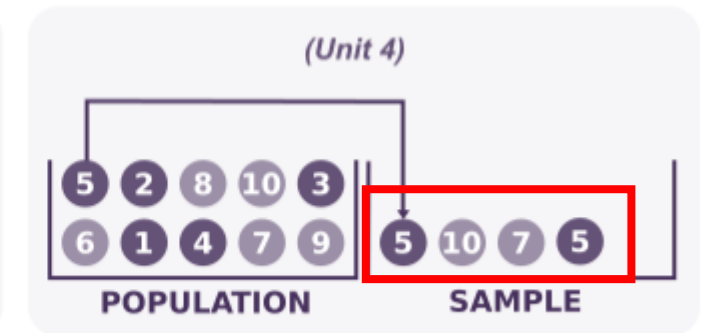
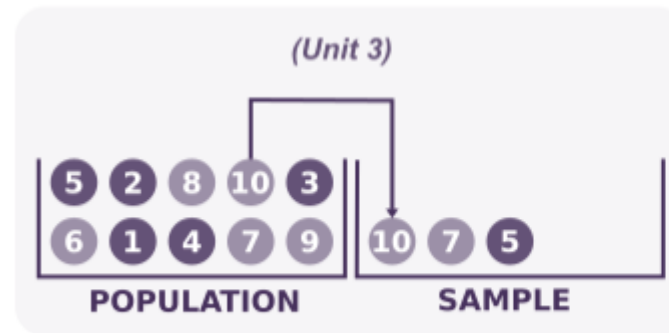
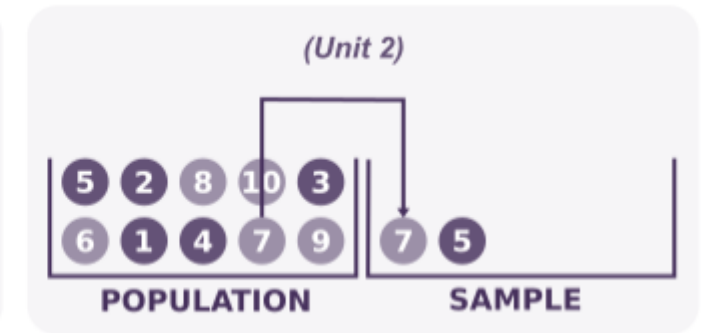
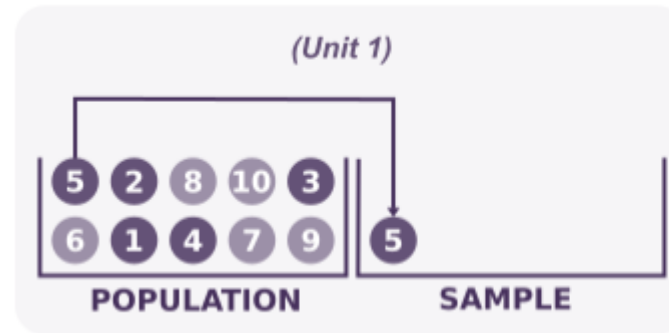


Sampling Techniques

Simple Random Sampling

What is Simple Random Sampling?

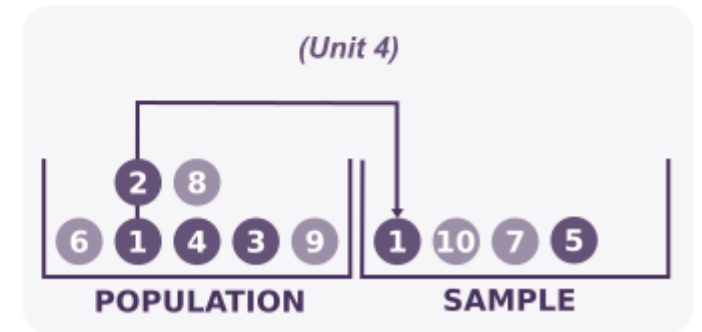
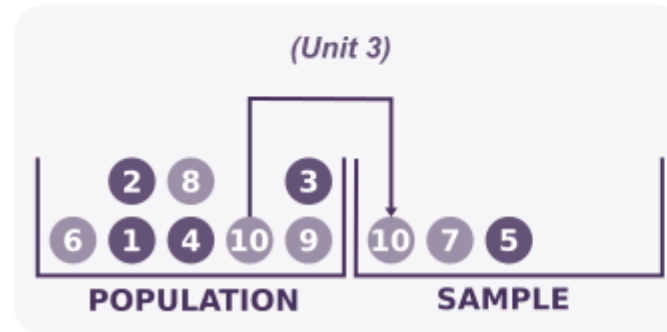
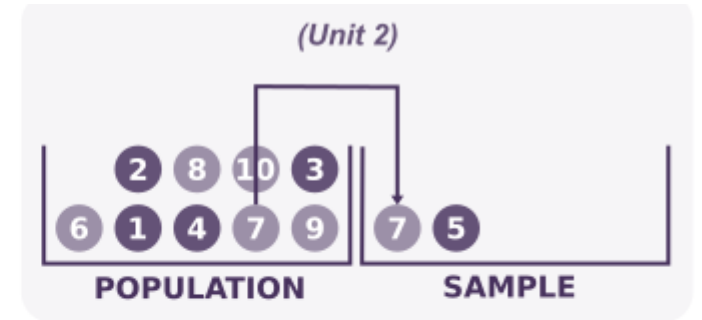
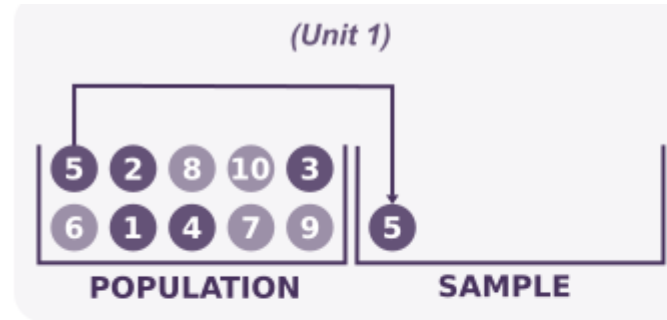
Every subject in the population has an equal chance of being selected.
Randomly picking subjects from the whole population



Simple Random Sampling with Replacement

Simple Random Sampling

Simple Random
Sampling without
Replacement

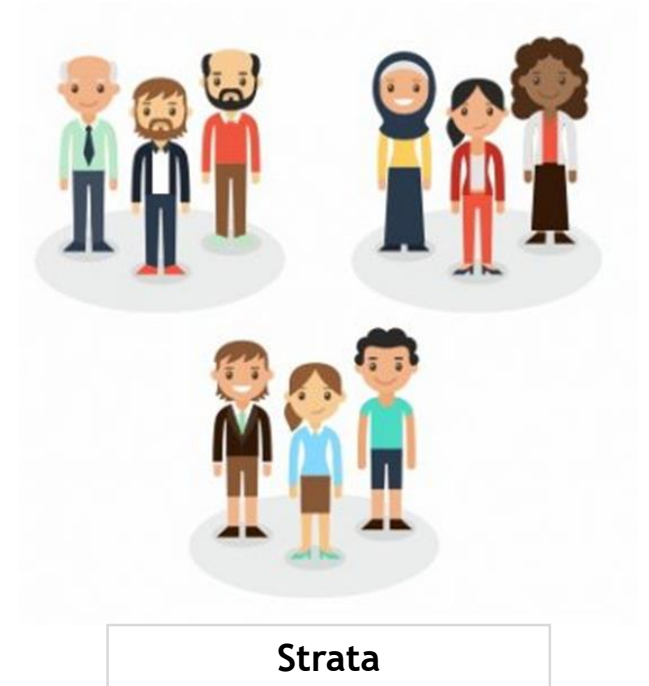
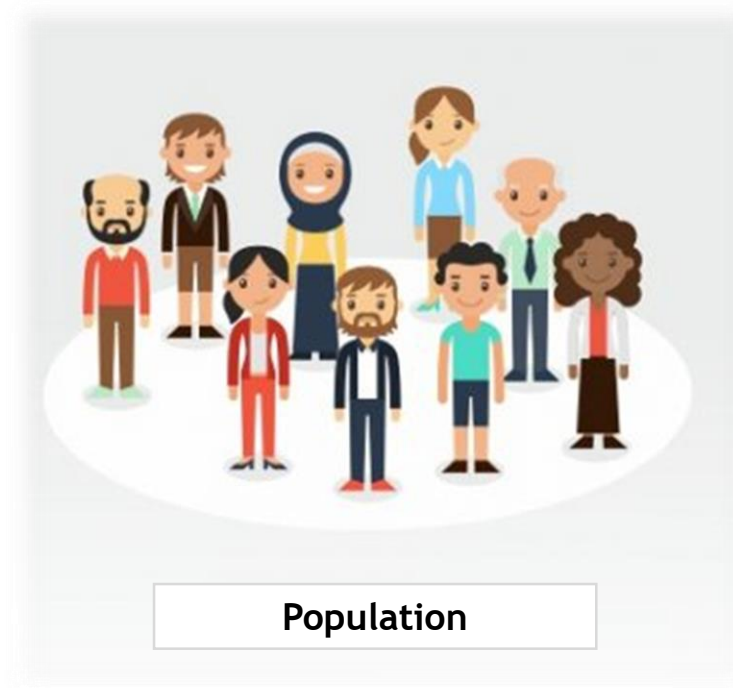


Simple Random Sampling without Replacement

Stratified Sampling

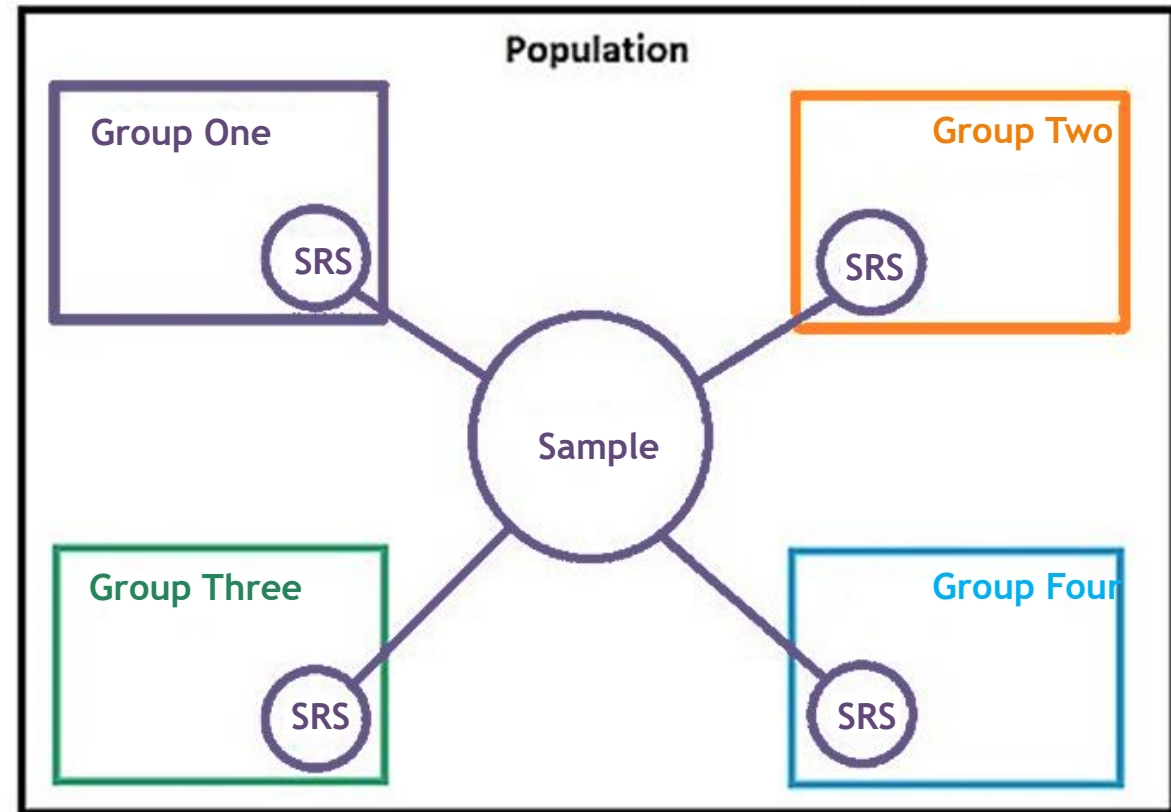
What is Stratified Sampling?

The population is partitioned into non-overlapping groups called strata, and a sample is selected by some design within each stratum.



Stratified Sampling

What is
Stratified
Sampling?

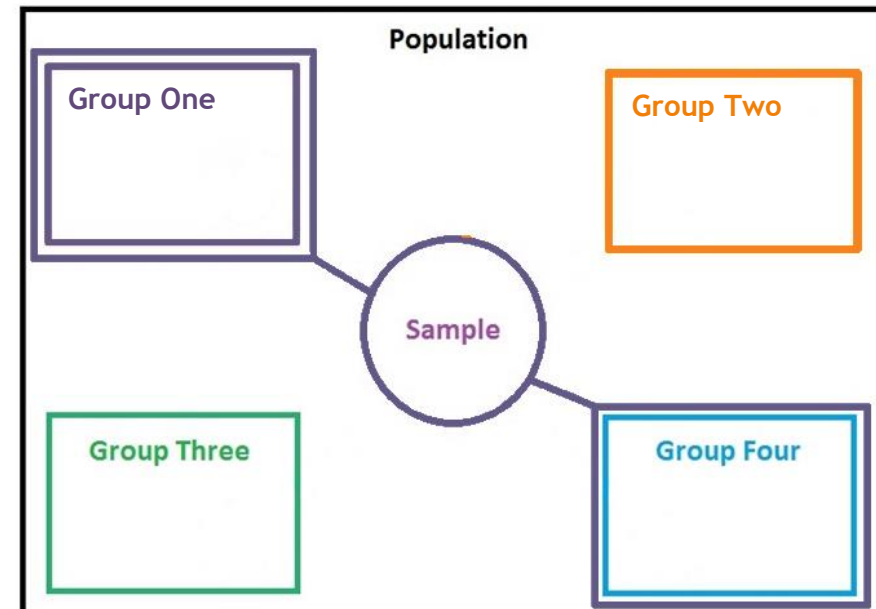


Cluster Sampling

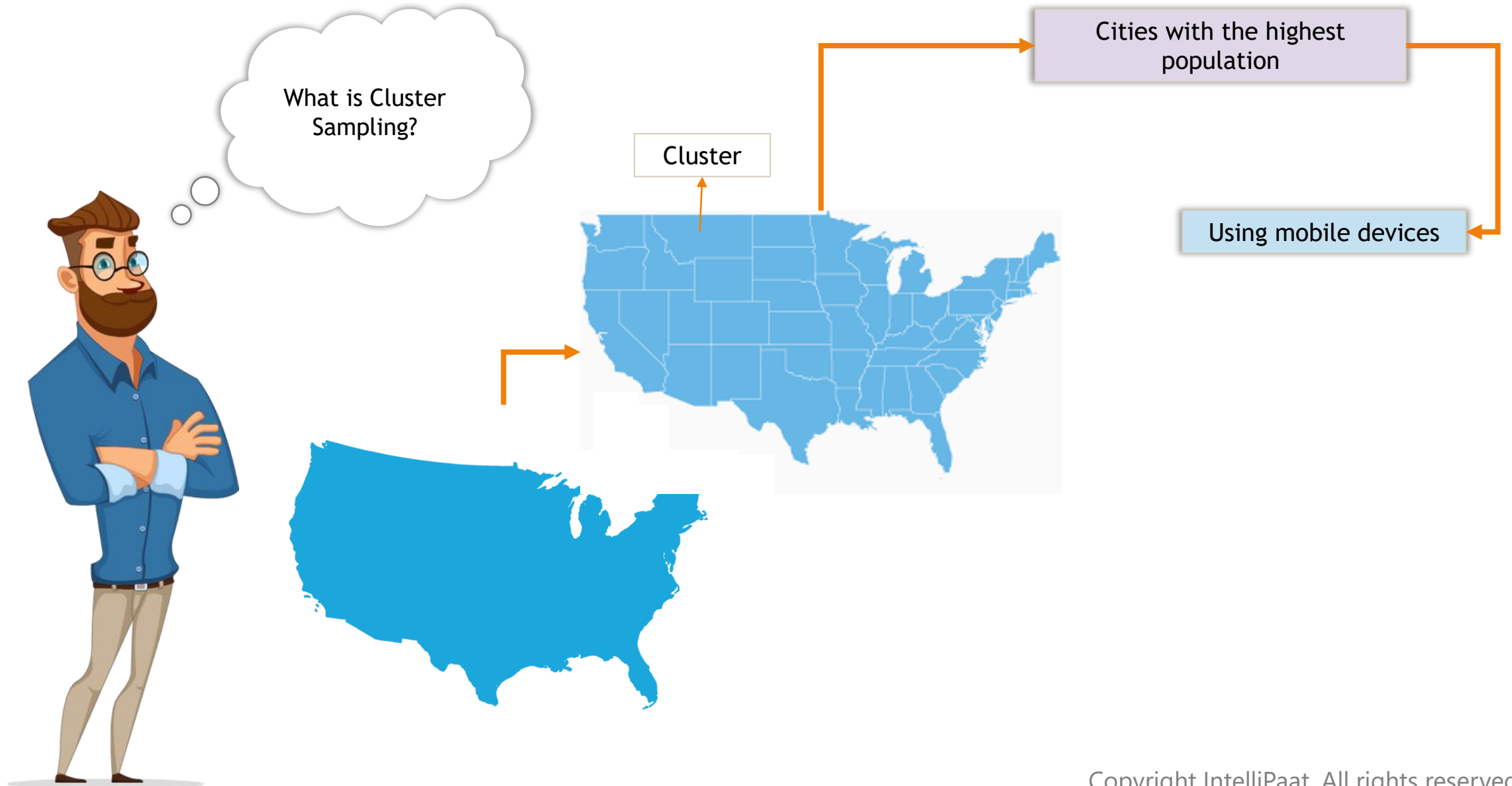
What is Cluster Sampling?



- Divide the population into groups (clusters)
- Obtain a simple random sample of clusters from all possible clusters
- Obtain data on every sampling unit in each of the randomly selected clusters

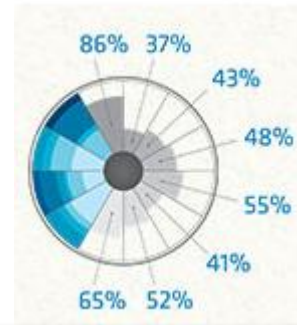
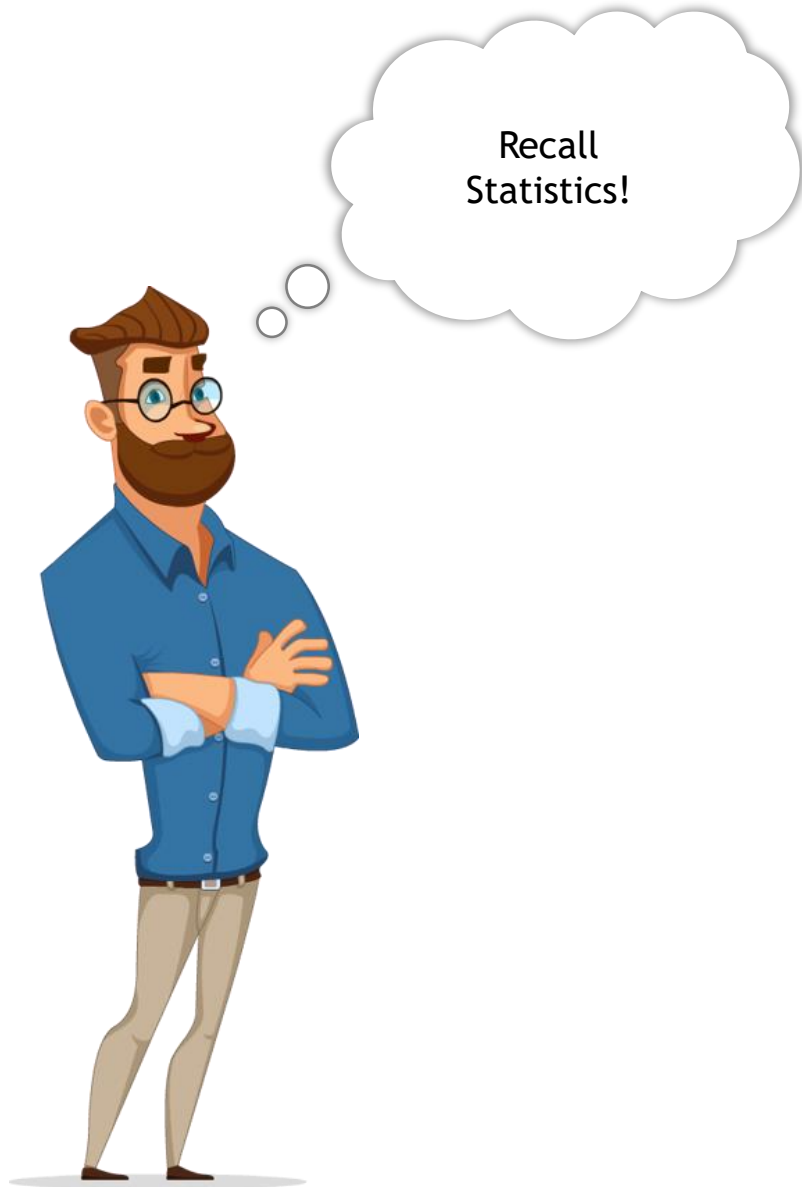


Cluster Sampling

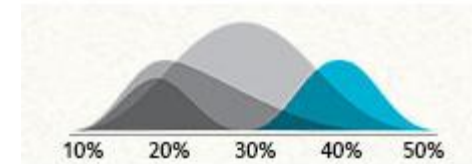


Categories of Statistics

Categories of Statistics



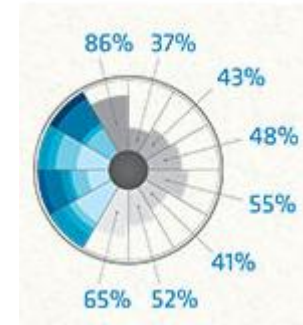
Descriptive Statistics



Inferential Statistics

Categories of Statistics

What is
Descriptive
Statistics?

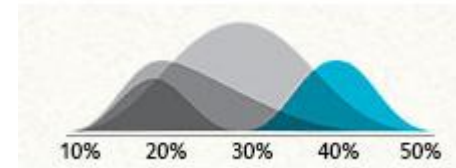


The analysis of data that helps describe, show or summarize the data in a meaningful way such that, for example, patterns might emerge from the data.

Just describes and summarizes data

Categories of Statistics

What is
Inferential
Statistics ?



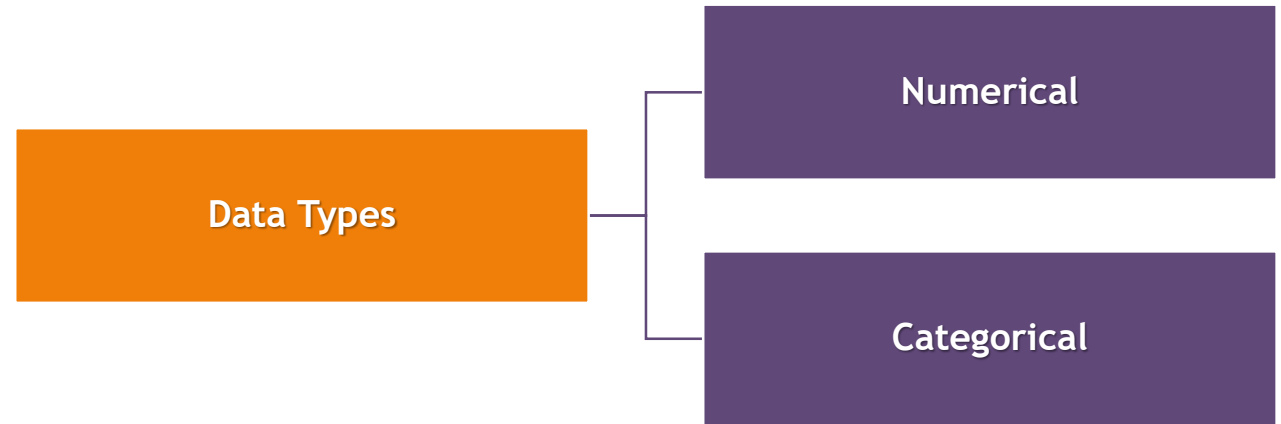
Inferential Statistics takes data from a sample and makes inferences about the larger population from which the sample was drawn.

Studies a sample of the same data

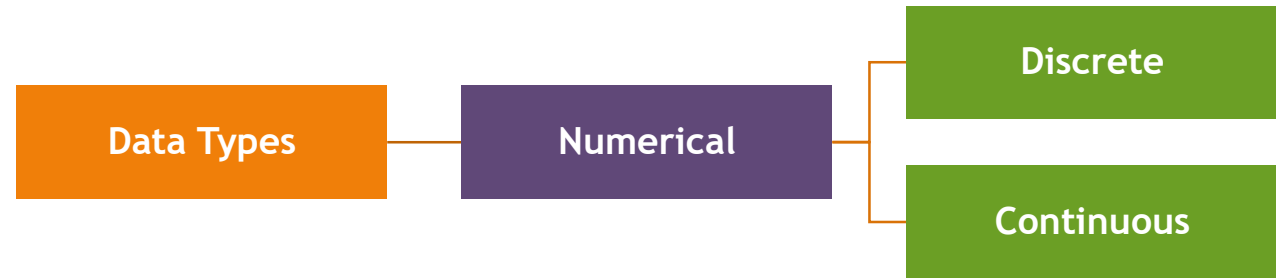
Types of Data

Types of Data

What are the
Data Types in
Statistics?



Types of Data - Numerical Data



- Deals with numbers
- Data which can be measured
- Length, height, area, volume, weight, speed, time, temperature, humidity, sound, etc.

Types of Data - Numerical Data


A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed. A thought bubble above his head contains the text "What is Discrete Data?".

What is
Discrete Data?

Values or observations that is counted as distinct and separate and can only take particular **values**

Counted

Types of Data - Numerical Data

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble above him contains the text "What is Continuous Data?".

What is
Continuous
Data?

You can measure continuous data. **Values** or observations may take on any **value** within a finite or infinite interval.

Measured


Temperature?

- ☐ - 10
- ☐ -5
- ☐ 0
- ☐ + 5
- ☐ + 10
- ☐ + 15

Types of Data - Categorical Data



Types of Data - Categorical Data

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed. A thought bubble above his head contains the text "What is Ordinal Data?".

What is Ordinal Data?

Values or observations can be ranked (put in order) or have a rating scale attached. You can count and order but not measure ordinal data.

Logical Ordering


How do you feel today?

- ☒ 1 - Very Unhappy
- ☐ 2 - Unhappy
- ☐ 3 - OK
- ☐ 4 - Happy
- ☐ 5 - Very Happy

How satisfied are you with our service?

- ☒ 1 - Very Unsatisfied
- ☐ 2 - Somewhat Unsatisfied
- ☐ 3 - Neutral
- ☐ 4 - Somewhat Satisfied
- ☐ 5 - Very Satisfied

Types of Data - Categorical Data

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed and looking thoughtful. A thought bubble above him contains the text "What is Nominal Data?".

What is
Nominal Data?

Values or observations can be assigned a code in the form of a number where the numbers are simply labels. You can count but not order or measure nominal data.

Name or Label

What is your gender?

- ☒ M – Male
- ☐ F – Female

What is your hair color?

- ☒ 1 – Brown
- ☐ 2 – Black
- ☐ 3 – Blonde
- ☐ 4 – Gray
- ☐ 5 – Other

Where do you live?

- ☒ A – North of the equator
- ☐ B – South of the equator
- ☐ C – Neither: In the international space station

Describing Data through Statistics

Describing Data through Statistics




Describing Data through Statistics

Measures of Central Tendency

Measures of Spread

Describing Data through Statistics

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble is above his head.

What are the
measures of
Central
Tendency?

These are ways of describing the central position of a frequency distribution for a group of data.

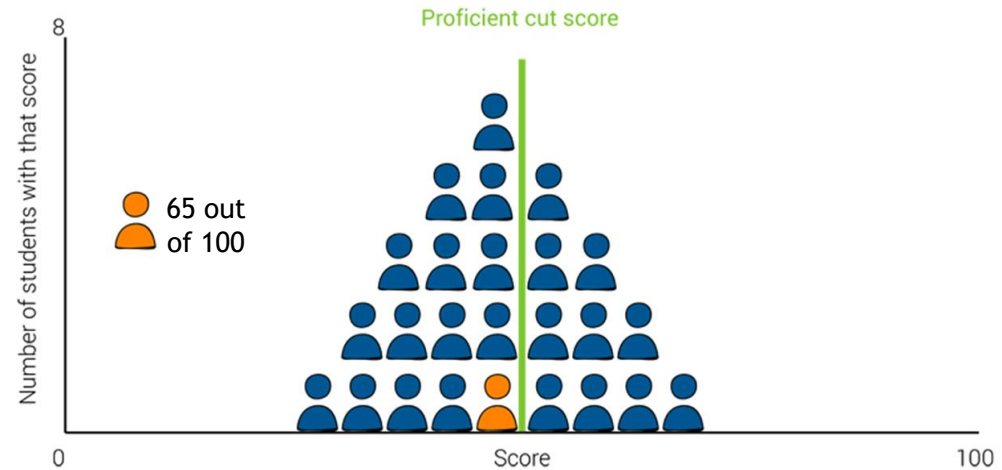
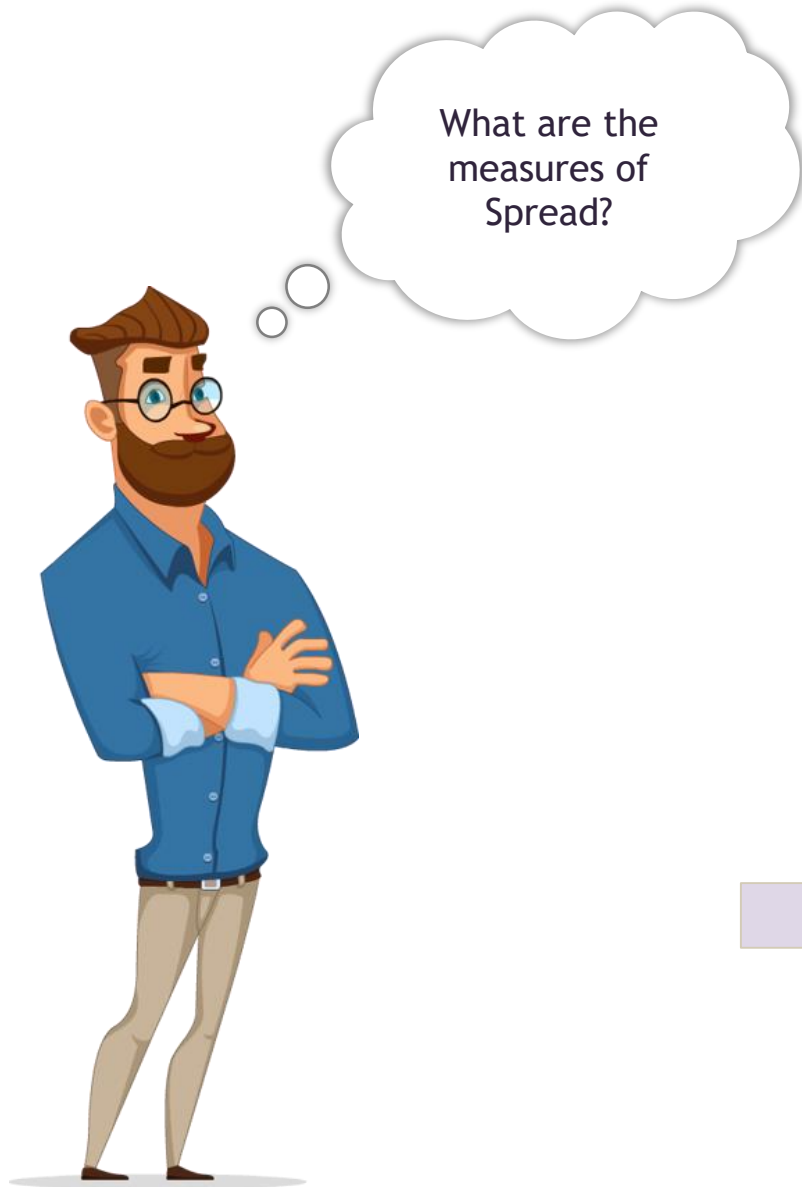
Mean

Median

Mode

Summarize a bunch of scores with a single number

Describing Data through Statistics



These are ways of summarizing a group of data by describing how spread out the scores are.

Range

Quartiles

Variance

Standard Deviation

Know the spread of scores within a bunch of scores

The Central Tendencies

Mean



The mean is the most common measure of center. It is what most people think of when they hear the word "average". However, the mean is affected by extreme values so it may not be the best measure of center to use in a skewed distribution.

Used with both discrete and continuous data

$$\text{Mean} = \bar{x} = \frac{\sum_i^n x_i}{n} = \frac{x_1, x_2, x_3, \dots x_i}{i}$$

$$\text{Mean} = \frac{(1 + 2 + 3 + 4 + 5)}{5} \longrightarrow \frac{15}{5} \longrightarrow 3$$

Trimmed Mean



What is
Trimmed Mean?

A variation of the mean, which you calculate by dropping a fixed number of sorted values at each end and then taking average of the remaining values.

Trimmed mean with p smallest and largest values omitted

Marks of 5 student = 67, 15, 75, 72, 85

$$\text{Trimmed Mean} = \bar{x} = \frac{\sum_{i=p+1}^{n-p} x_i}{n-2p}$$

$$\text{Trimmed Mean} = \frac{(67+75+72)}{3} \rightarrow \frac{214}{3} \rightarrow 71.3$$

Implementation in R

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed. A thought bubble above him contains the text "Let's implement mean".

Let's
implement
mean

customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService
7590-VHVEG	Female	0	Yes	No	1	No
5575-GNVDE	Male	0	No	No	34	Yes
3668-QPYBK	Male	0	No	No	2	Yes
7795-CFOCW	Male	0	No	No	45	No
9237-HQITU	Female	0	No	No	2	Yes
9305-CDSKC	Female	0	No	No	8	Yes
1452-KIOVK	Male	0	No	Yes	22	Yes
6713-OKOMC	Female	0	No	No	10	No
7892-POOKP	Female	0	Yes	No	28	Yes
6388-TABGU	Male	0	No	Yes	62	Yes
9763-GRSKD	Male	0	Yes	Yes	13	Yes

```
mean(customer_churn$MonthlyCharges)
```

Median

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed. A thought bubble above him contains the text "What is Median?".

What is
Median?

The median is the middle number in a sorted list of the data.


$$\text{Median} = \frac{n+1}{2}$$

$$\text{If } n=7, \text{ Median} = \frac{7+1}{2} = 4^{\text{th}} \text{ value}$$

$$\text{Median} = \frac{n}{2}, \left(\frac{n}{2}\right) + 1$$

$$\text{If } n=6, \text{ Median} = \frac{6}{2}, \left(\frac{6}{2}\right) + 1 = 3^{\text{rd}} \text{ and } 4^{\text{th}} \text{ value}$$

Median

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble is above his head.

What are the advantages median has over mean?

If we have more outliers in our datasets, then it is best to use median than mean because mean will give poor central tendency.

An outlier is a rare chance of occurrence within a given dataset.

Implementation in R

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed. A thought bubble above him contains the text "Let's implement median".

Let's
implement
median

customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService
7590-VHVEG	Female	0	Yes	No	1	No
5575-GNVDE	Male	0	No	No	34	Yes
3668-QPYBK	Male	0	No	No	2	Yes
7795-CFOCW	Male	0	No	No	45	No
9237-HQITU	Female	0	No	No	2	Yes
9305-CDSKC	Female	0	No	No	8	Yes
1452-KIOVK	Male	0	No	Yes	22	Yes
6713-OKOMC	Female	0	No	No	10	No
7892-POOKP	Female	0	Yes	No	28	Yes
6388-TABGU	Male	0	No	Yes	62	Yes
9763-GRSKD	Male	0	Yes	Yes	13	Yes

```
median(customer_churn$MonthlyCharges)
```

Mode

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble above him contains the text "What is Mode?".

What is Mode?

The most frequently occurring data point

Mean and median need not be in a dataset, but mode has to be in it.

Variability and Spread

Range

What is Range?

$$\text{Range} = \text{Max} - \text{Min}$$

Player 1 = $13 - 7$

Points scored per game	7	8	9	10	11	12	13
Frequency, f	1	1	2	2	2	1	1

Player 2 = $13 - 7$

Points scored per game	7	9	10	11	13
Frequency, f	1	2	4	2	1

Player 3 = $30 - 3$

Points scored per game	3	6	7	10	11	13	30
Frequency, f	2	1	2	3	1	1	1



Implementation in R


A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed. A thought bubble above him contains the text "Let's implement range".

Let's
implement
range

customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService
7590-VHVEG	Female	0	Yes	No	1	No
5575-GNVDE	Male	0	No	No	34	Yes
3668-QPYBK	Male	0	No	No	2	Yes
7795-CFOCW	Male	0	No	No	45	No
9237-HQITU	Female	0	No	No	2	Yes
9305-CDSKC	Female	0	No	No	8	Yes
1452-KIOVK	Male	0	No	Yes	22	Yes
6713-OKOMC	Female	0	No	No	10	No
7892-POOKP	Female	0	Yes	No	28	Yes
6388-TABGU	Male	0	No	Yes	62	Yes
9763-GRSKD	Male	0	Yes	Yes	13	Yes

```
range(customer_churn$MonthlyCharges)
```

Quartiles

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble above him contains the text "What are Quartiles?".

What are Quartiles?

Quartiles tell us about the spread of a dataset by breaking the dataset into quarters, just like the median breaks it in half.

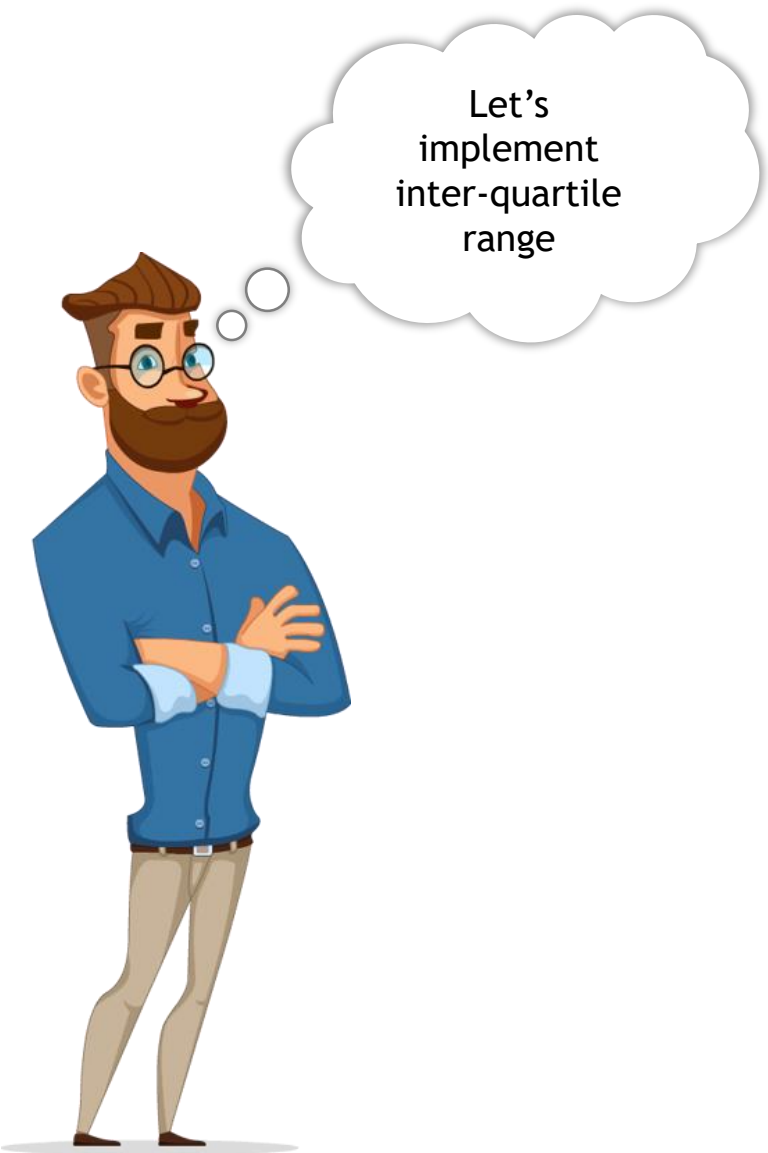
Lower quartile (25th percentile, Q1) = $\frac{(n+1)}{4}$ th

Middle quartile = Median = $\frac{2*(n+1)}{4}$ th

Upper quartile (75th percentile, Q3) = $\frac{3*(n+1)}{4}$ th

Interquartile range, IQR = Q3-Q1 (central 50% of data)

Implementation in R


A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed. A thought bubble above him contains the text "Let's implement inter-quartile range".

Let's
implement
inter-quartile
range

customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService
7590-VHVEG	Female	0	Yes	No	1	No
5575-GNVDE	Male	0	No	No	34	Yes
3668-QPYBK	Male	0	No	No	2	Yes
7795-CFOCW	Male	0	No	No	45	No
9237-HQITU	Female	0	No	No	2	Yes
9305-CDSKC	Female	0	No	No	8	Yes
1452-KIOVK	Male	0	No	Yes	22	Yes
6713-OKOMC	Female	0	No	No	10	No
7892-POOKP	Female	0	Yes	No	28	Yes
6388-TABGU	Male	0	No	Yes	62	Yes
9763-GRSKD	Male	0	Yes	Yes	13	Yes

```
IQR(customer_churn$MonthlyCharges)
```

Variance & Standard Deviation

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble is above his head.

What is Variance
and what is
Standard
Deviation?

The variance of the data is the average squared distance between the mean and each data value.

$$\text{Variance} = (\text{Unit})^2$$

$$\text{Variance of length} = (\text{m})^2$$

Standard deviation measures how spread out the values in a dataset are around the mean. Square root of the variance

Variance & Standard Deviation


What is Variance
and what is
Standard
Deviation?



$$\text{Variance} = S^2 = \frac{\sum (x - \bar{x})^2}{n - 1}$$

$$\text{Standard Deviation} = S = \sqrt{\text{variance}}$$

Variance & Standard Deviation

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed and looking thoughtful. A thought bubble is above his head.

What is Variance
and what is
Standard
Deviation?


Standard deviation is the measure of spread most commonly used in statistical practice when the mean is used to calculate central tendency.

Standard deviation is also influenced by **outliers**. **One** value could contribute largely to the results of the standard deviation.

Standard deviation is also useful when comparing the spread of two separate datasets that have approximately the same mean.

The more widely spread the values are the larger the standard deviation is.

Implementation in R

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed and a thoughtful expression. A thought bubble above him contains the text "Let's implement variance and standard deviation".

Let's implement
variance and
standard
deviation

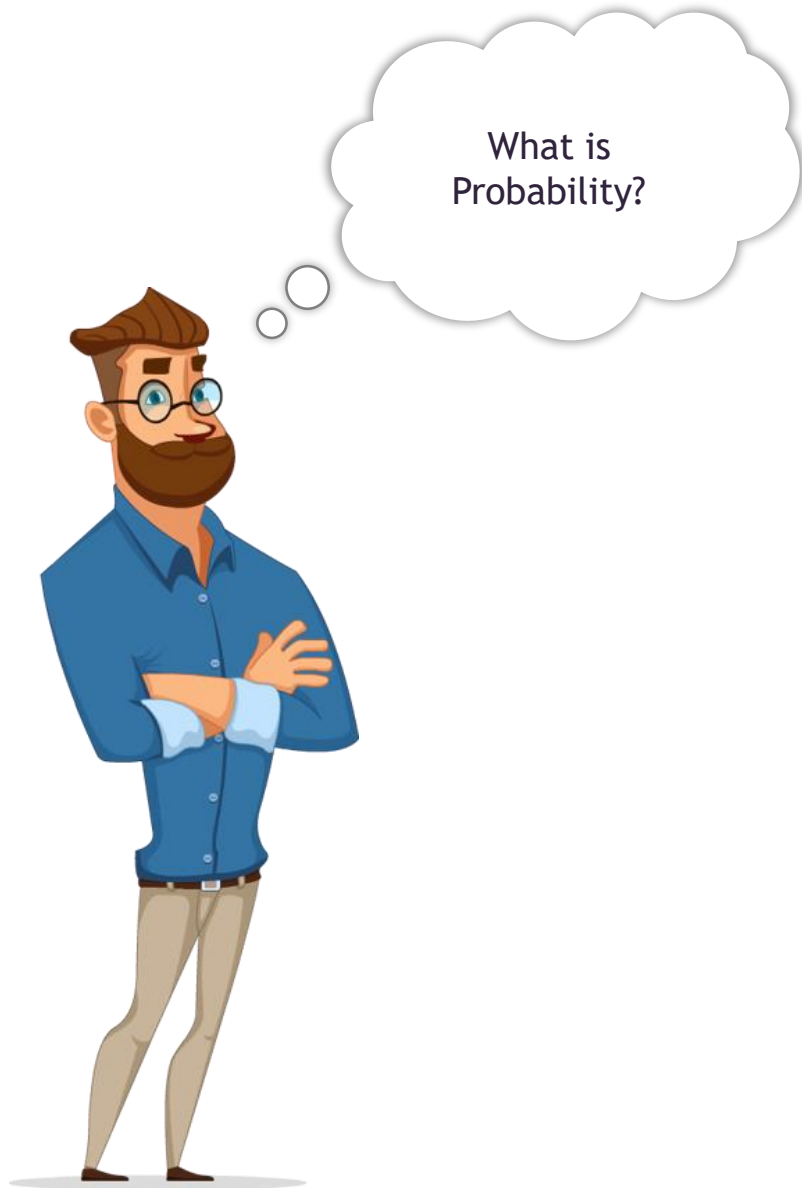
customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService
7590-VHVEG	Female	0	Yes	No	1	No
5575-GNVDE	Male	0	No	No	34	Yes
3668-QPYBK	Male	0	No	No	2	Yes
7795-CFOCW	Male	0	No	No	45	No
9237-HQITU	Female	0	No	No	2	Yes
9305-CDSKC	Female	0	No	No	8	Yes
1452-KIOVK	Male	0	No	Yes	22	Yes
6713-OKOMC	Female	0	No	No	10	No
7892-POOKP	Female	0	Yes	No	28	Yes
6388-TABGU	Male	0	No	Yes	62	Yes
9763-GRSKD	Male	0	Yes	Yes	13	Yes

```
var(customer_churn$MonthlyCharges)
```

```
sd(customer_churn$MonthlyCharges)
```

Probability

Probability



Probability can be determined post conducting a thought experiment.

$$P(E) = \frac{\text{\textit{\# of times an event occurred}}}{\text{\textit{total \# of opportunities for the event to have occurred}}}$$

Probability

Probability
Example



$$\text{Blue fish} = 14$$

$$\text{Green fish} = 5$$

$$\text{Orange fish} = 5$$

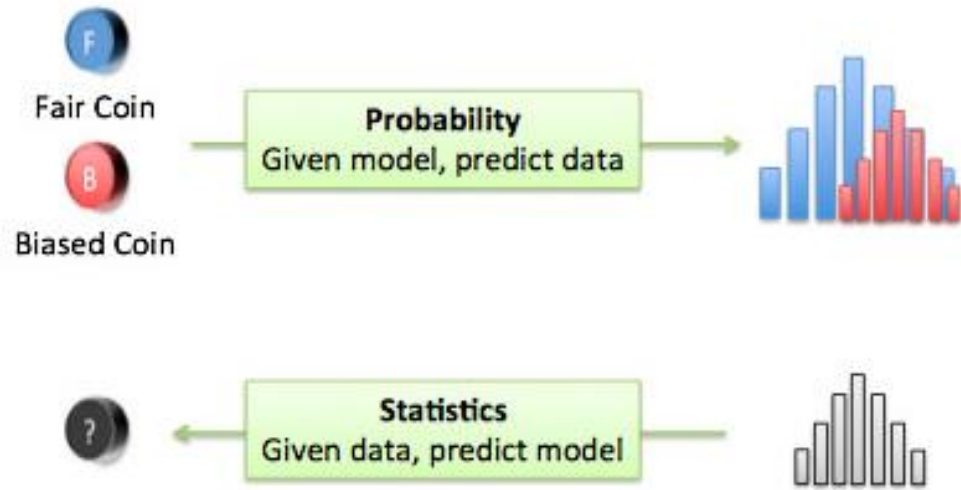
$$\text{Purple fish} = 6$$

$$\text{Total} = 30$$

$$\text{Probability of a purple fish} = \frac{6}{30} = \frac{1}{5}$$

Probability Vs. Statistics

Probability Vs. Statistics



Probability	<ul style="list-style-type: none">• Predict the likelihood of a future event• What will happen in a given ideal world?
Statistics	<ul style="list-style-type: none">• Analyze the past events• How ideal is the world?

Probability Terminology

Probability Terminology

Probability Key
Terms

Experiment

A test to see what will happen in case you do something.



Outcome

Refers to a single (one) result of an experiment.



Event

The set of a group of different outcomes of an experiment.



Sample Space

The total number of all different possible outcomes of a given experiment.

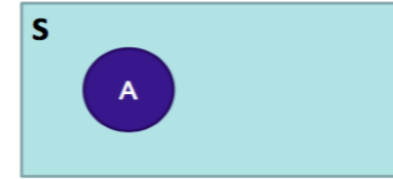
	(H) heads	(T) tails
heads (H)	(H) (H)	(T) (H)
tails (T)	(H) (T)	(T) (T)

Probability Rules

Probability Rules



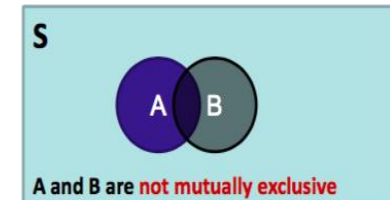
$$P(S) = 1$$



$$0 \leq P(A) \leq 1$$



$$\begin{aligned} P(A \text{ or } B) \\ = P(A) + P(B) \end{aligned}$$



$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

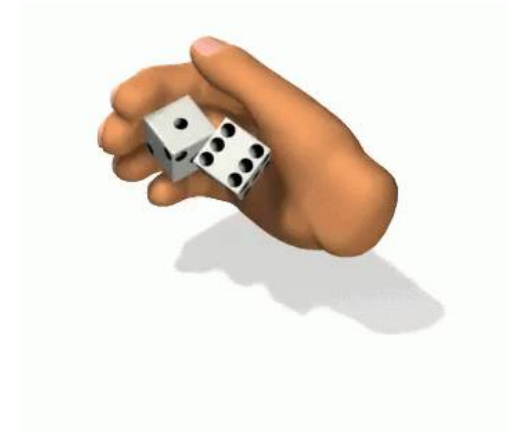
Independent Events

What are
Independent
Events?



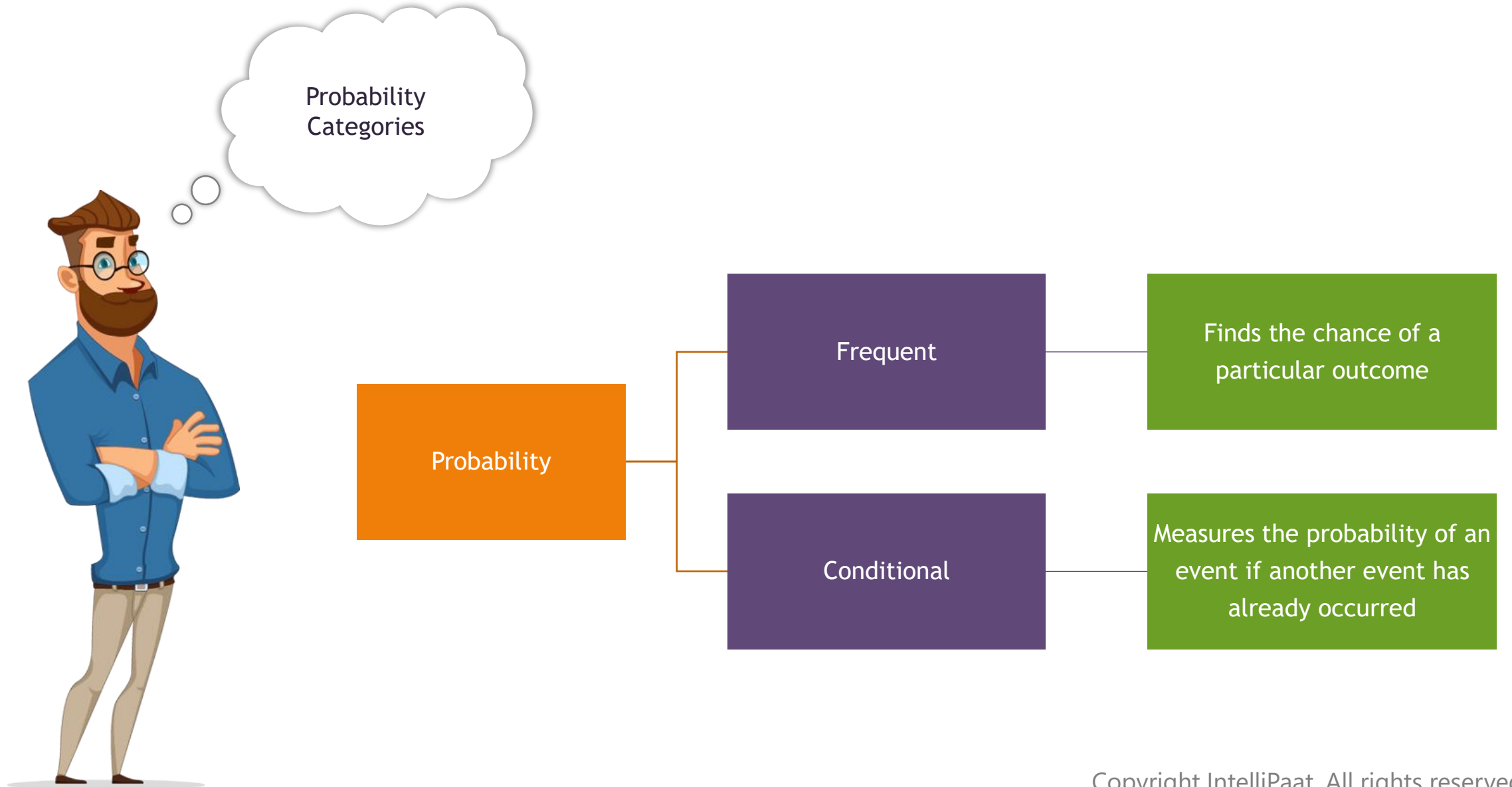
The outcome of the event B is not dependent on the outcome of the event A

$$P(A \text{ and } B) = P(A) * P(B)$$

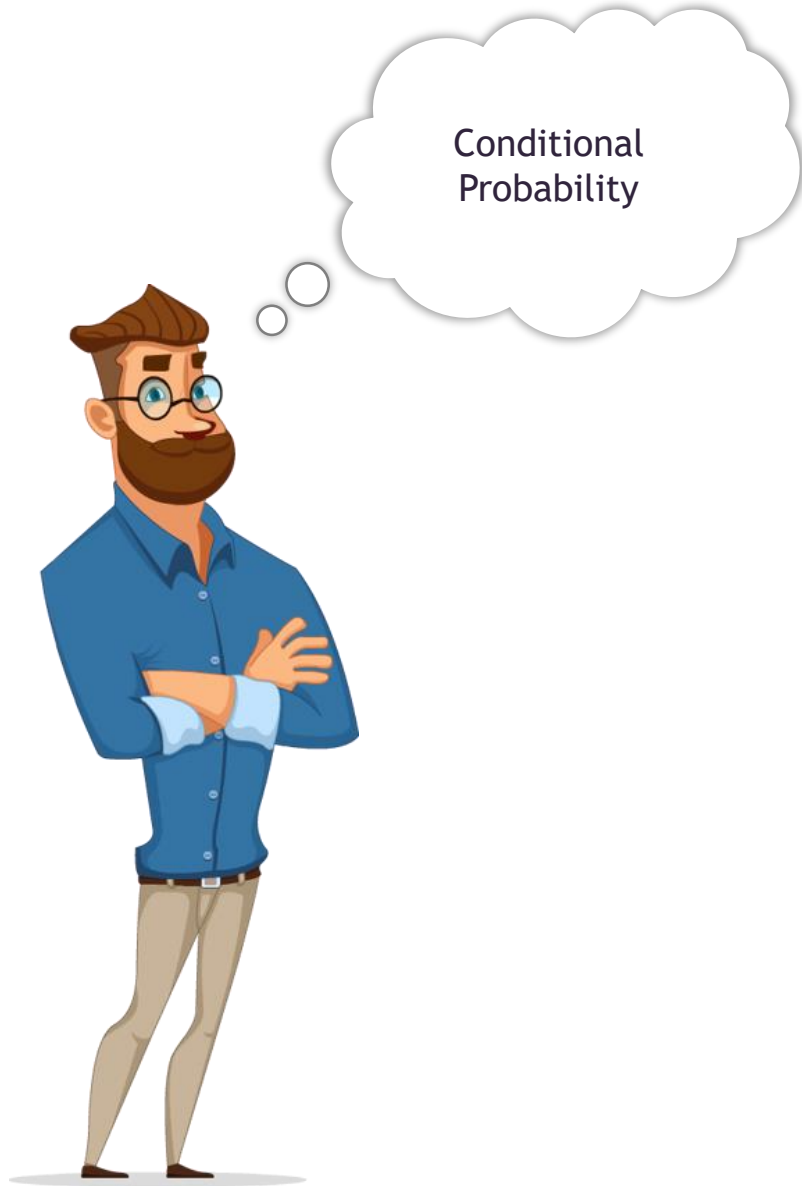


Probability Categories

Probability Categories



Conditional Probability



$$P(B | A) = \frac{P(A \cap B)}{P(A)}$$

$$\text{or } P(A \cap B) = P(A)P(B | A)$$

Conditional Probability

Formula for
Conditional
Probability



$$P(A|B) = \frac{P(A \text{ and } B)}{P(B)} \Rightarrow P(A \text{ and } B) = P(B) * P(A|B)$$

Similarly

$$P(B|A) = \frac{P(A \text{ and } B)}{P(A)} \Rightarrow P(A \text{ and } B) = P(A) * P(B|A)$$

Equating, we get

$$P(A|B) * P(B) = P(A) * P(B|A)$$

$$\therefore P(A|B) = \frac{P(A) * P(B|A)}{P(B)}$$

Probability Distribution

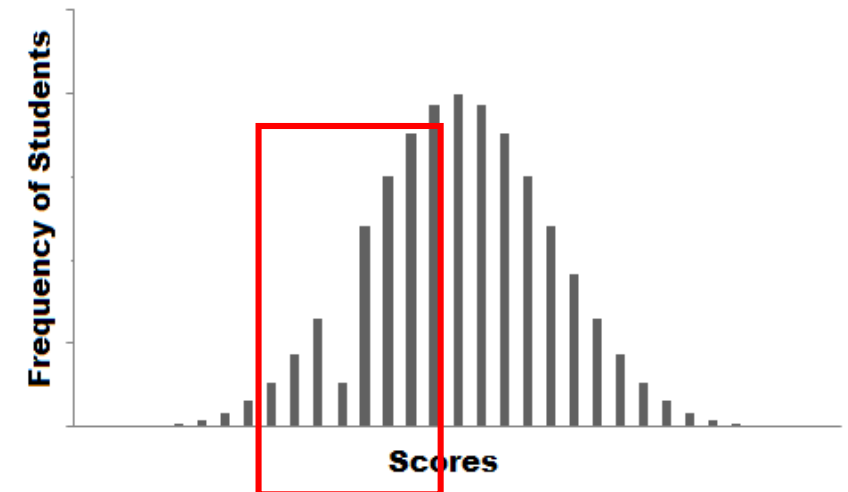
Probability Distribution

What is
Probability
Distribution?

The probability function for a discrete random variable is the probability mass function.

It shows the exact probabilities for a particular value of the random variable.

S. No.	Scores
1	25
2	27
3	38
4	42
5	16
6	35
7	46
8	48
9	31
10	31

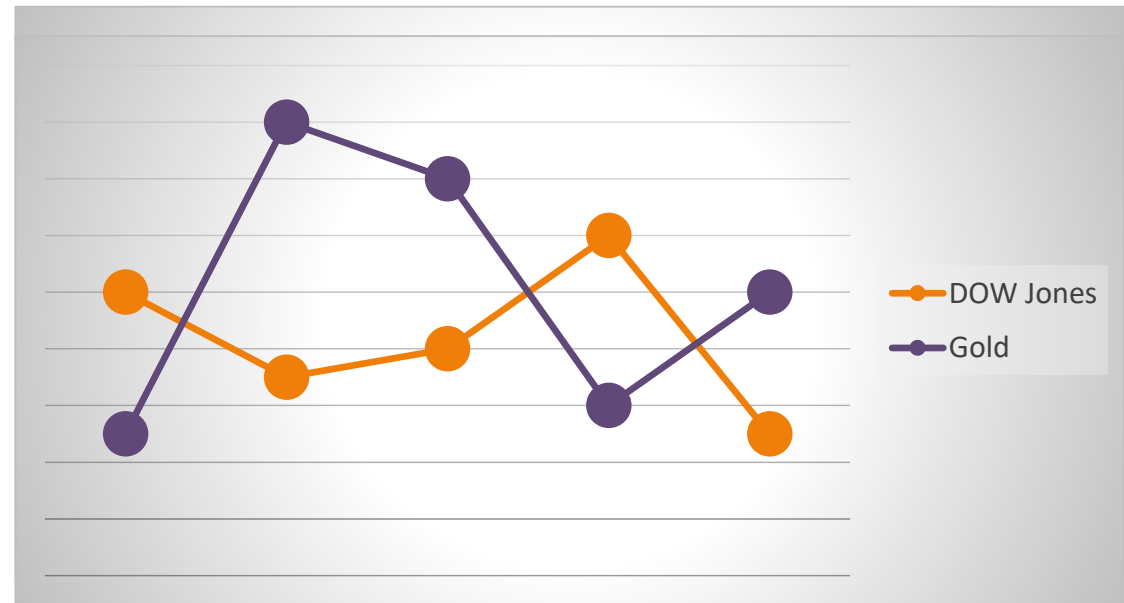


Covariance & Correlation

Correlation


What is
Correlation ?

When you say that two items correlate, you are saying that a change in one item effects a change in another item. You will always talk about correlation as a range between -1 and 1 .



Correlation values are dependent on units of measure of “X” and “Y”.

Covariance

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble is above his head.

What is
Covariance ?

If you say that two items tend to vary together, then you are talking about the covariance between the two items which can be positive or negative covariance.

Positive covariance indicates that higher than average values of one variable tend to be paired with higher than average values of the other variable.

Negative covariance indicates that higher than average values of one variable tend to be paired with lower than average values of the other variable.

Covariance is not standardized, unlike the correlation coefficient. Therefore, covariance values can range from negative infinity to positive infinity.

Covariance & Correlation



Covariance

Positive Covariance

Negative Covariance

Correlation

Positive Correlation

Negative Correlation

Zero Correlation

Covariance & Correlation

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed. A thought bubble above his head contains the text "Covariance and Correlation".

Covariance and
Correlation

Correlation is dimensionless; i.e., it is a unit-free measure of the relationship between variables, unlike covariance where the value is obtained by the product of the units of the two variables.

Covariance & Correlation



$$COV(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$

$$Correlation = \rho = \frac{cov(X, Y)}{\sigma_X \sigma_Y}$$

Covariance & Correlation

Key differences
between
Covariance and
Correlation



Basis for comparison	Covariance	Correlation
Meaning	A measure indicating the extent to which two random variables change in tandem	A measure indicating that how strongly two random variables are related
What is it?	Measure of correlation	Scaled version of covariance
Values	Lies between $-\infty$ and $+\infty$	Lies between -1 and $+1$
Change in scale	Affects covariance	Does not affect correlation


Covariance & Correlation

Are sunshine
and concert
attendance
correlated?

Sunshine (hours)	1.9	2.5	3.2	3.8	4.7	5.5	5.9	7.2
Concert attendance (100s)	22	33	30	42	38	49	42	55



Covariance & Correlation

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed and looking thoughtful. A thought bubble is above his head.

Let's see how
concert
attendance
varies with
sunshine

Independent variable (explanatory): Sunshine (Plotted on x-axis)

Dependent variable (response): Concert attendance (Plotted on y-axis)

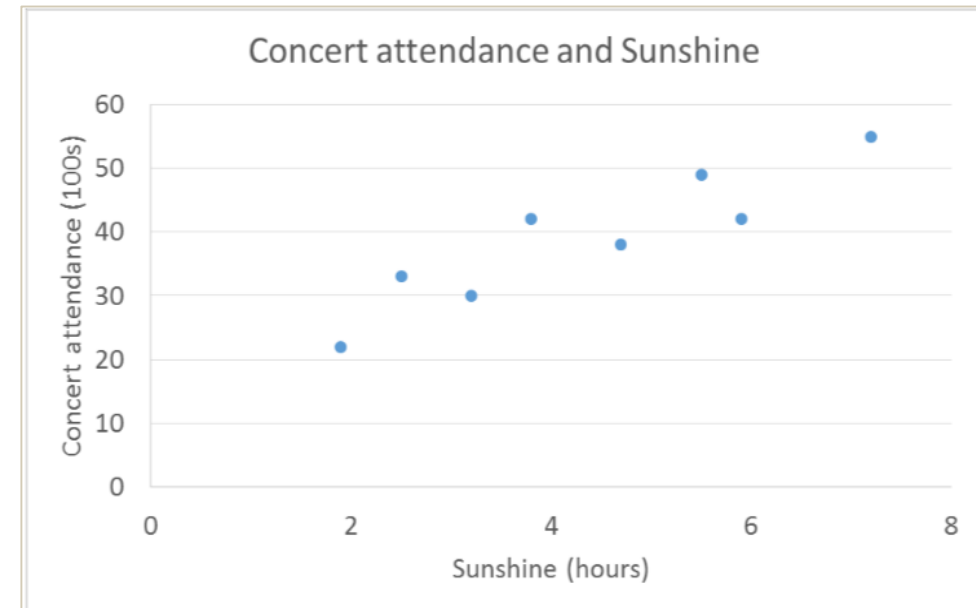
Sunshine (hours)	1.9	2.5	3.2	3.8	4.7	5.5	5.9	7.2
Concert attendance (100s)	22	33	30	42	38	49	42	55

Covariance & Correlation

This is an
interesting
correlation!



The number of attendees increase with increase in sunshine



Covariance & Correlation

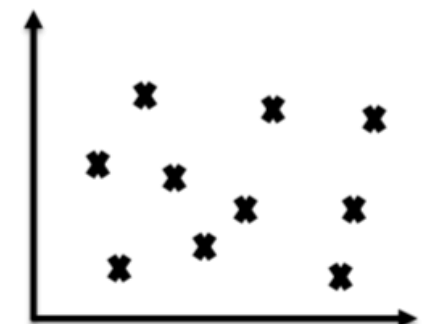
These are the
different
correlation
lines.



Positive Linear
Correlation




Negative Linear
Correlation



No Correlation

Correlation Coefficient

Correlation Coefficient

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed. A thought bubble above his head contains the text "What is Correlation Coefficient?".

What is
Correlation
Coefficient?

Correlation coefficient is used to find how strong a relationship is between data.

It returns a value between -1 and 1 , where:

1

indicates a
strong positive
relationship.

-1

indicates a
strong negative
relationship.

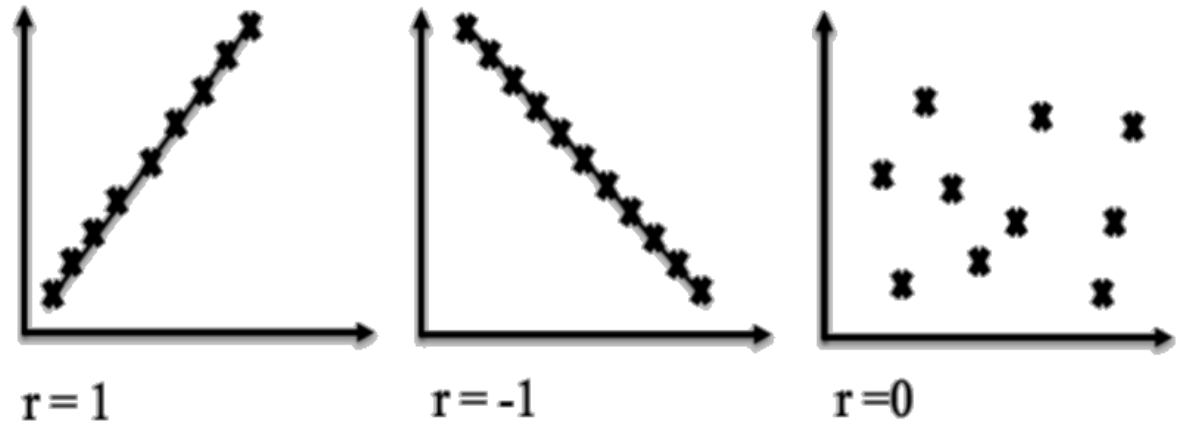
0

indicates no
relationship at
all.

Correlation Coefficient

What is
Correlation
Coefficient?


$$r = r_{xy} = \frac{\text{Cov}(x, y)}{S_x \times S_y}$$



Correlation Is Not Causation!

Because two things correlating does not necessarily mean that one causes the other

Correlation Does Not Cause Causation

A cartoon illustration of a man with a brown beard and glasses, wearing a blue button-down shirt and tan pants. He is standing with his arms crossed, looking thoughtful.

Correlation is
not causation!

Correlation describes the size and direction of a relationship between two or more variables.


Causation indicates that one event is the result of the occurrence of the other event.

Smoking → Correlates → Alcoholism

Smoking → Causes → The risk of developing Lung Cancer

Standardization & Normalization

Standardization

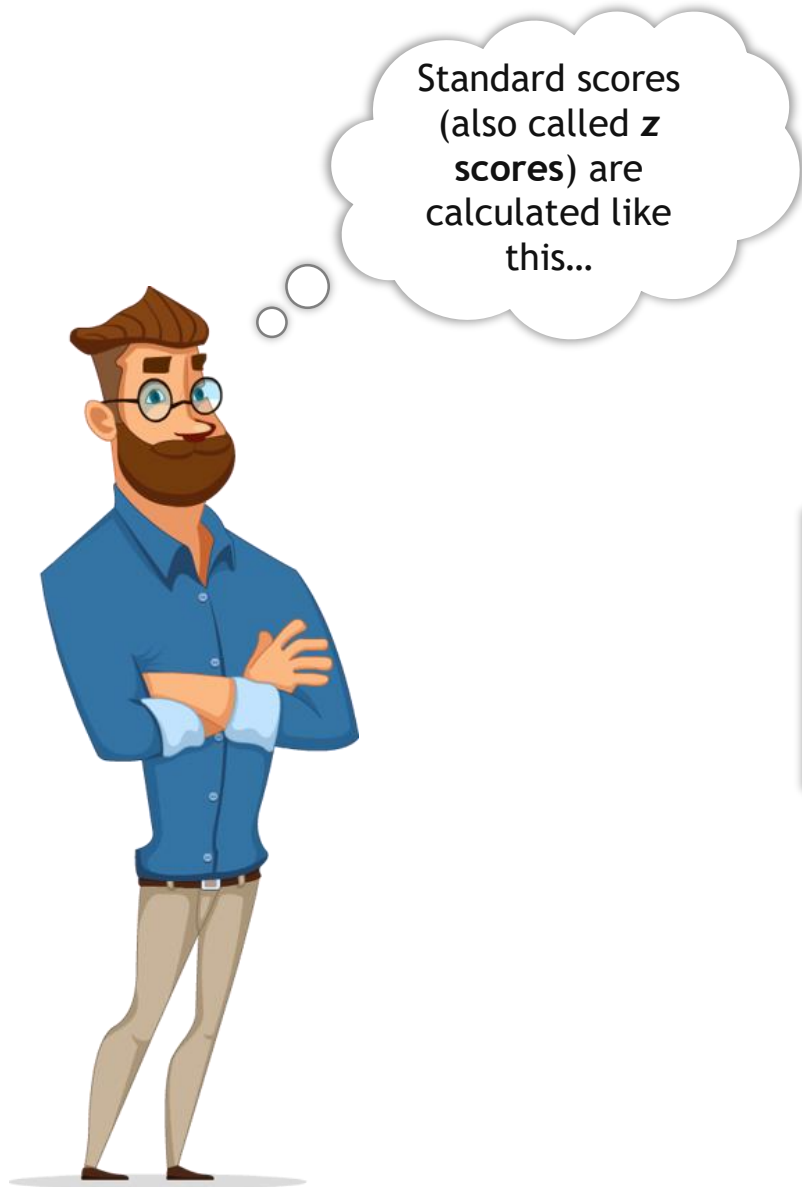
A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble above him contains the text "What is Standardization?".

What is
Standardization?

With standardization (or Z-score normalization), the features can be rescaled so that they'll have the properties of a standard normal distribution with $\mu = 0$ and $\sigma = 1$

Where, μ is the mean and σ is the standard deviation from the mean


Standardization



$$z = \frac{x - \mu}{\sigma}$$

Standardizing the features so that they are centered around 0 with a standard deviation of 1 is not only important if we are comparing measurements that have different units, but it is also a general requirement for many machine learning algorithms

Normalization

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble is above his head.

Alternative
approach to z-
score
normalization

An alternative approach to Z-score normalization (or standardization) is the so-called **Min-Max scaling**.

Often, it is also simply called “Normalization”, a common cause for ambiguities.

The data is scaled to a fixed range: 0 to 1.

Normalization

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble is above his head.

Formula for
Normalization

The cost of having this bounded range—in contrast to standardization—is that we will end up with smaller standard deviations, which can suppress the effect of outliers

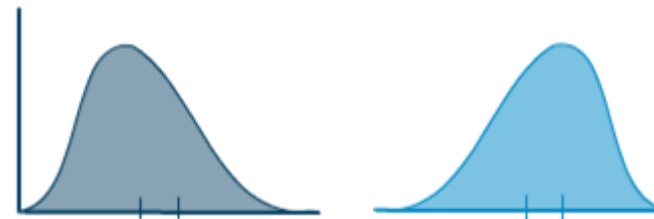
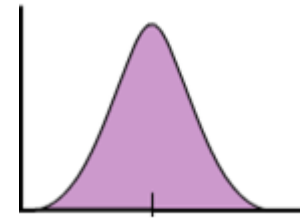
$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

Skewness

Skewness

What is
Skewness?

Skewness represents an imbalance and asymmetry from the mean of a data distribution

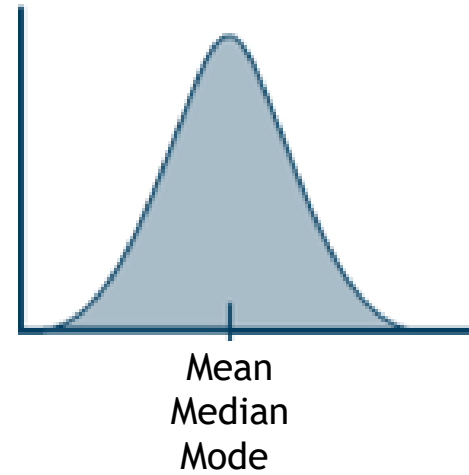


Skewness

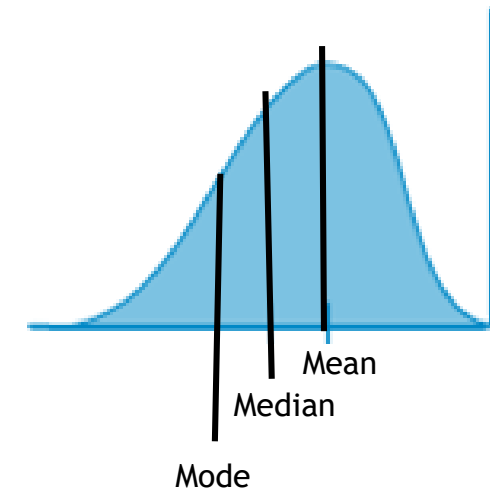
Properties of
Skewed Bell
Curves



Symmetric Distribution

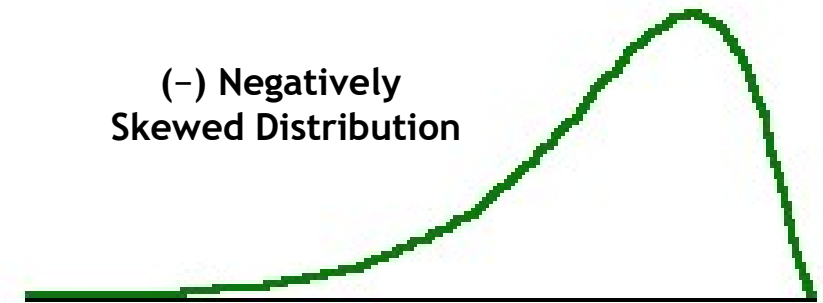
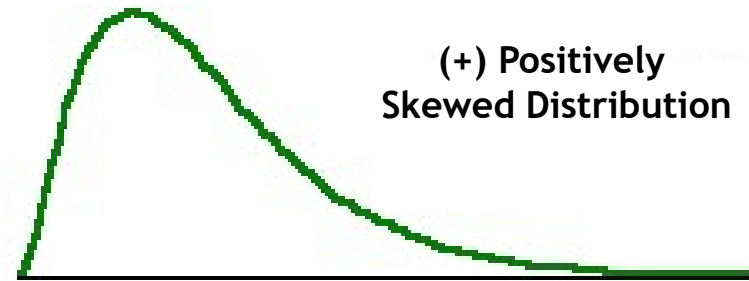


Skewed Distribution



Skewness

Properties of
Skewed Bell
Curves



Binomial Distribution

Binomial Distribution

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble is above his head.

What is Binomial Distribution?

A binomial distribution is a specific probability distribution. It is used to model the probability of obtaining one of two outcomes, a certain number of times (k), out of fixed number of trials (N) of a discrete random event.

Expected Outcome
Successful Outcome

p

Other Outcome
Failure

$1 - p$

Binomial Distribution

Criteria for
using Binomial
Distributions



Rule #1

There are only two mutually exclusive outcomes for a discrete random variable.

Rule #2

There is a fixed number of repeated trials.

Rule #3

Each trial is an independent event.

Rule #4

The probability of success for each trial is fixed.

Binomial Distribution

What is the
Binomial
Distribution
formula?



The probability of getting x successes in n Bernoulli trials

This starts the count of
number of ways an event can
occur.

$n!$

This is the
probability of
success for x trials.

$$P(x) = \frac{n!}{(n-x)! x!} p^x q^{n-x}$$

This is the
probability of
failure for x trials.

$(n-x)!$

This ends the count of
number of ways an event can
occur.

This deletes
duplications.

Geometric Distribution

Geometric Distribution

What is
Geometric
Distribution?



The geometric distribution represents the number of failures before you get a success in a series of Bernoulli trials.

$$f(x) = (1 - p)^{x-1}p$$

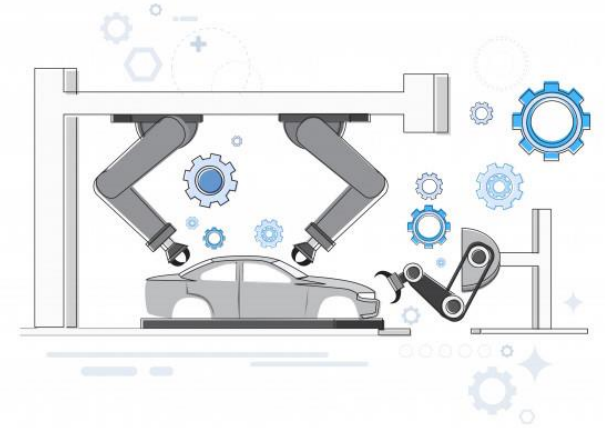


Normal Distribution

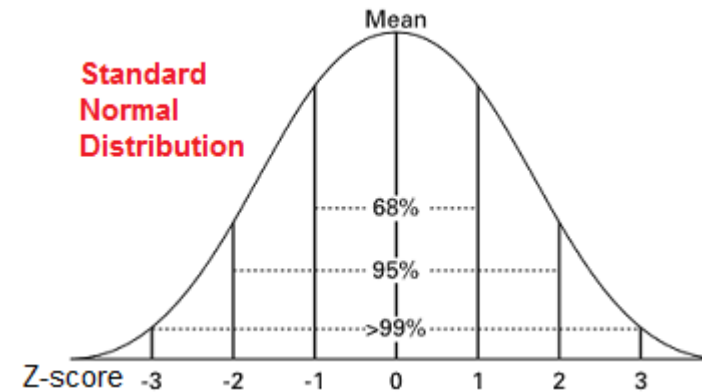
Normal Distribution



84th percentile or above



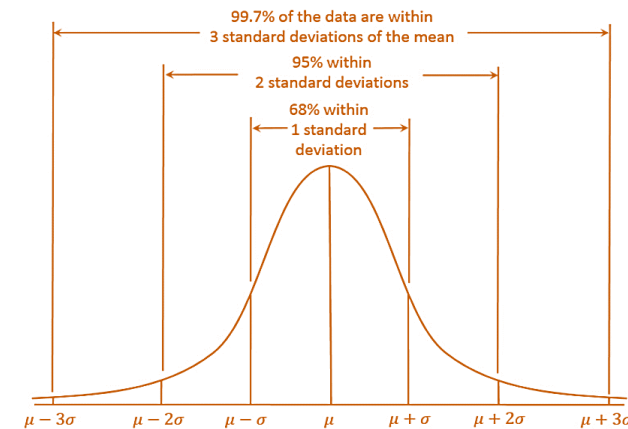
95% of the lithium ion batteries



Normal Distribution

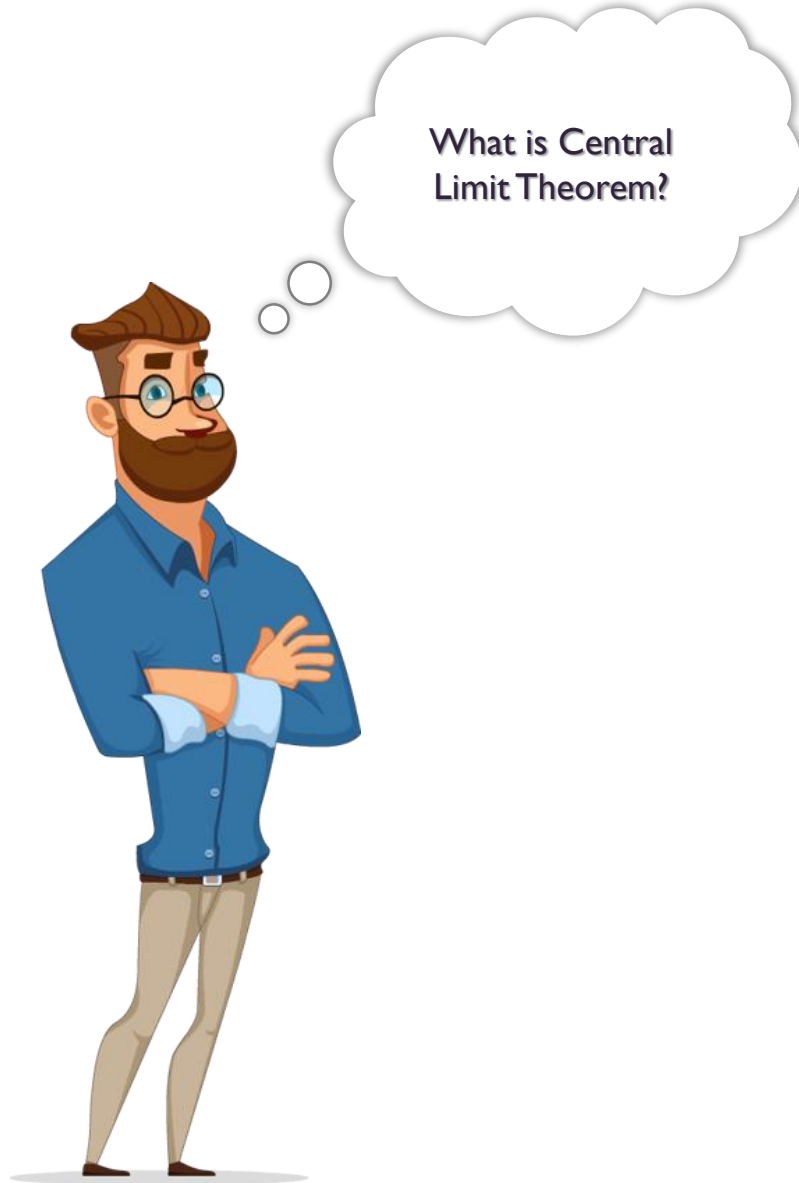
Features of Normal Distributions

- Symmetric bell shape
- Mean and median are equal
- The area under the normal curve is equal to 1.0
- Normal distributions are denser in the center and less dense in the tails
- 68% of the area of a normal distribution is within one standard deviation of the mean
- Approximately 95% of the area of a normal distribution is within two standard deviations of the mean



Central Limit Theorem

Central Limit Theorem

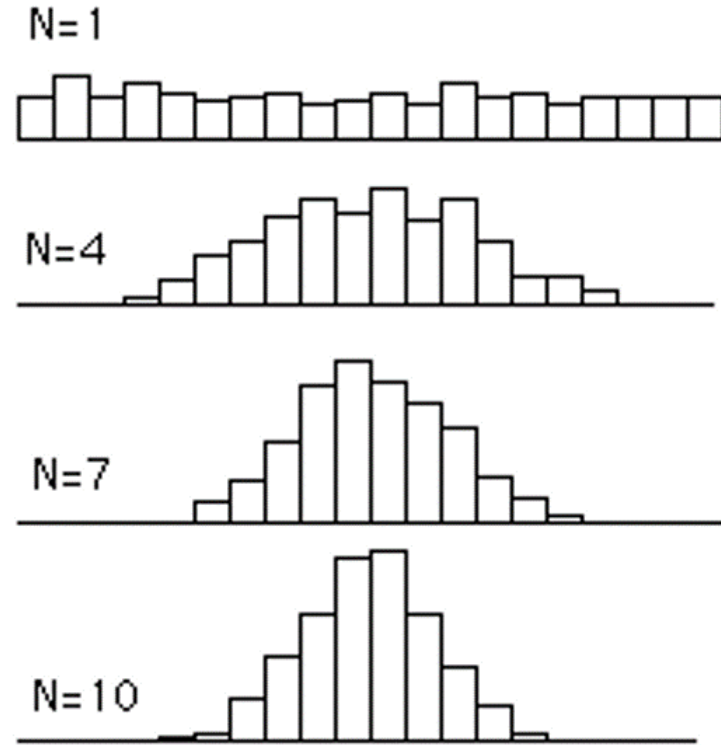


The Central Limit Theorem states that the sampling distribution of the sample means approaches a normal distribution as the sample size gets larger, no matter what the shape of the population distribution is.

This fact holds especially true for sample sizes over 30.

Central Limit Theorem

What is Central Limit Theorem?




When N increases:

1. the distributions become more and more normal.
2. the spread of the distributions decreases.

Hypothesis Testing

Hypothesis Testing

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed and looking thoughtful. A thought bubble is above his head.

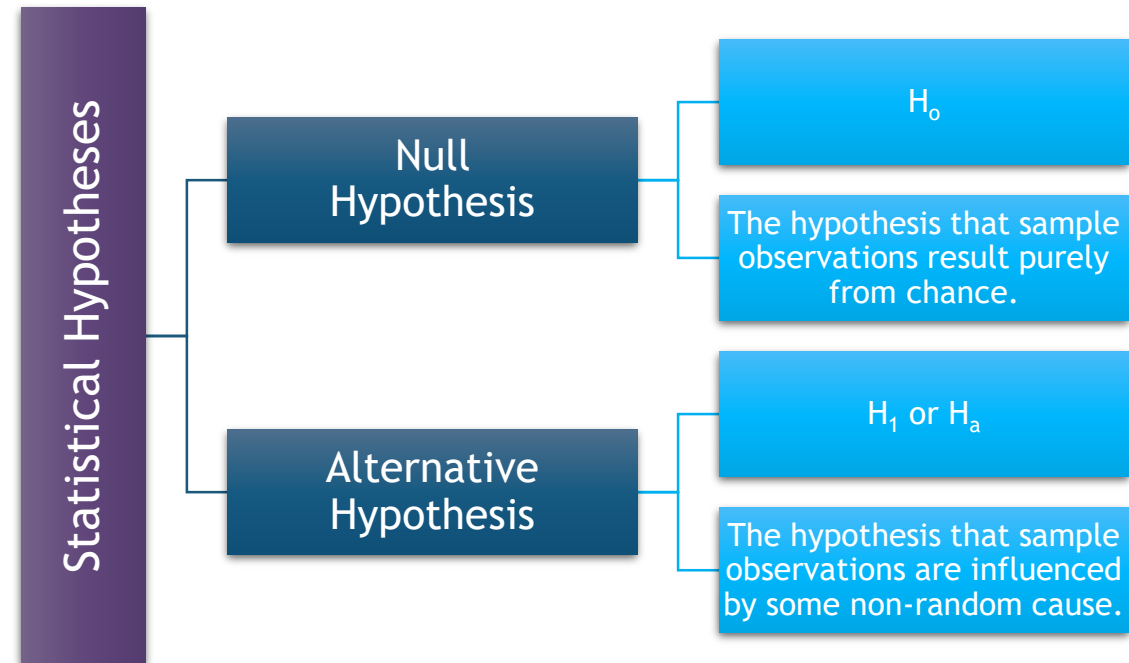
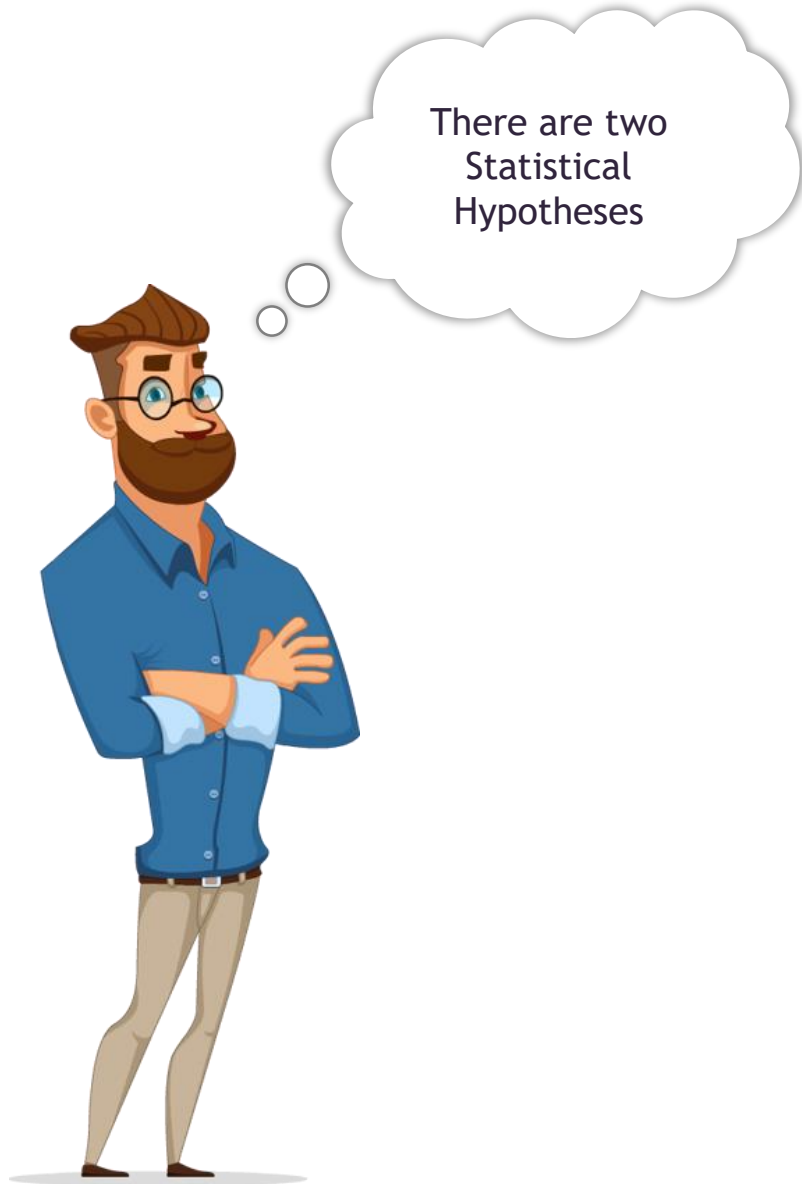
What is
Hypothesis
Testing?

An assumption about a population parameter

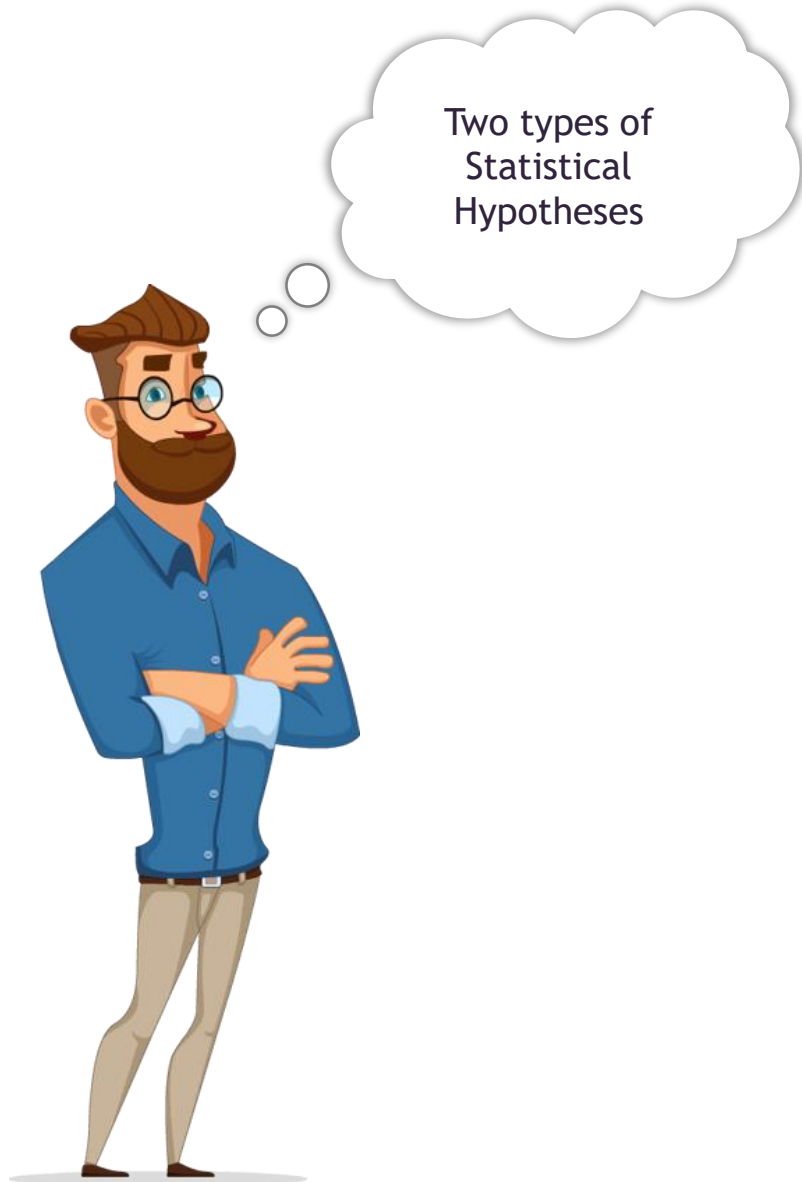
This assumption may or may not be true.

Formal procedures used by statisticians to accept or reject statistical hypotheses

Hypothesis Testing



Hypothesis Testing



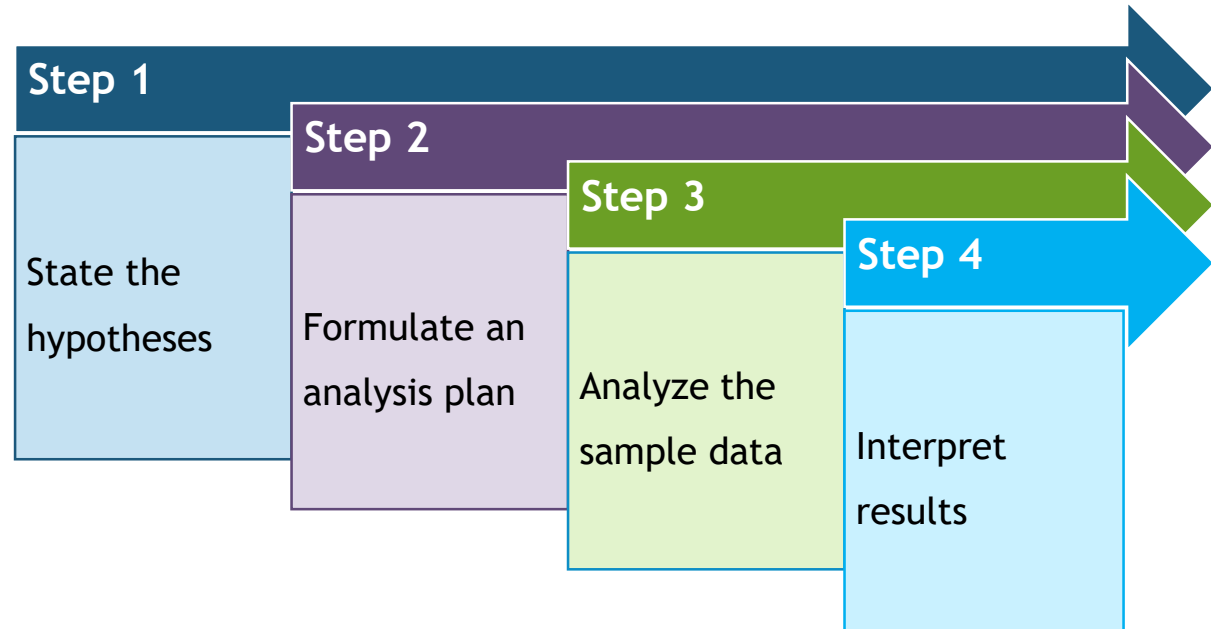
$$H_0: P = 0.5$$
$$H_1: P \neq 0.5$$

Coin Flip = 50 times → 40 Heads and 10 Tails

Reject the Hypothesis

Hypothesis Testing

The Structure
of
Hypothesis
Tests



Decision Error

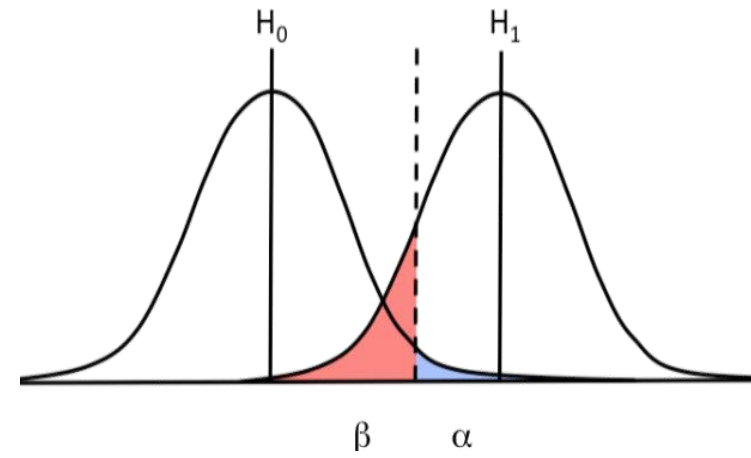
Decision Error

Types of Error



Type II Error

Type I Error

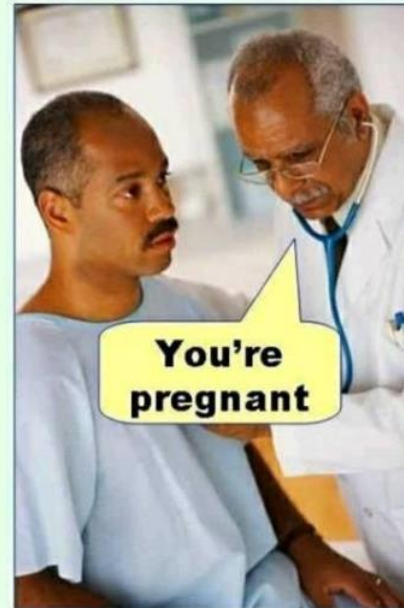


Decision Error

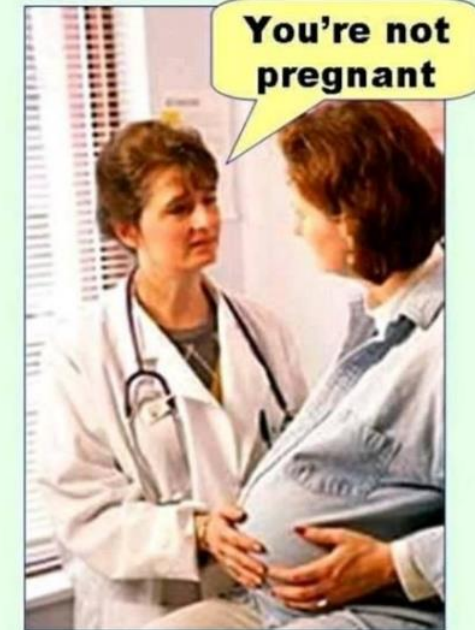
Types of Error



Type I error
(false positive)



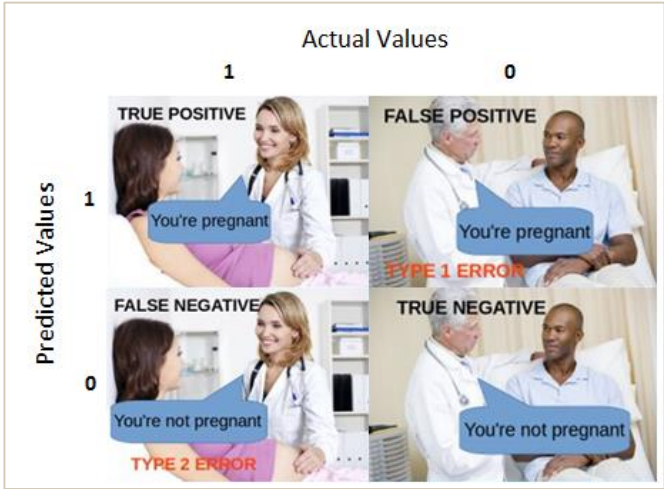
Type II error
(false negative)



Decision Error


Types of Error -
Metrix

		Condition	
		Condition Positive	Condition Negative
Test Outcome	Test Outcome Positive	True Positive	False Positive (Type I error)
	Test Outcome Negative	False Negative (Type II error)	True Negative



Decision Rules

Decision Rules

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and khaki pants, standing with his arms crossed and looking thoughtful. A thought bubble above him contains the text "What is Z-value?".

What is Z-value?

Z-value is a measure of standard deviation, i.e., how many standard deviations away from the mean is the observed value.

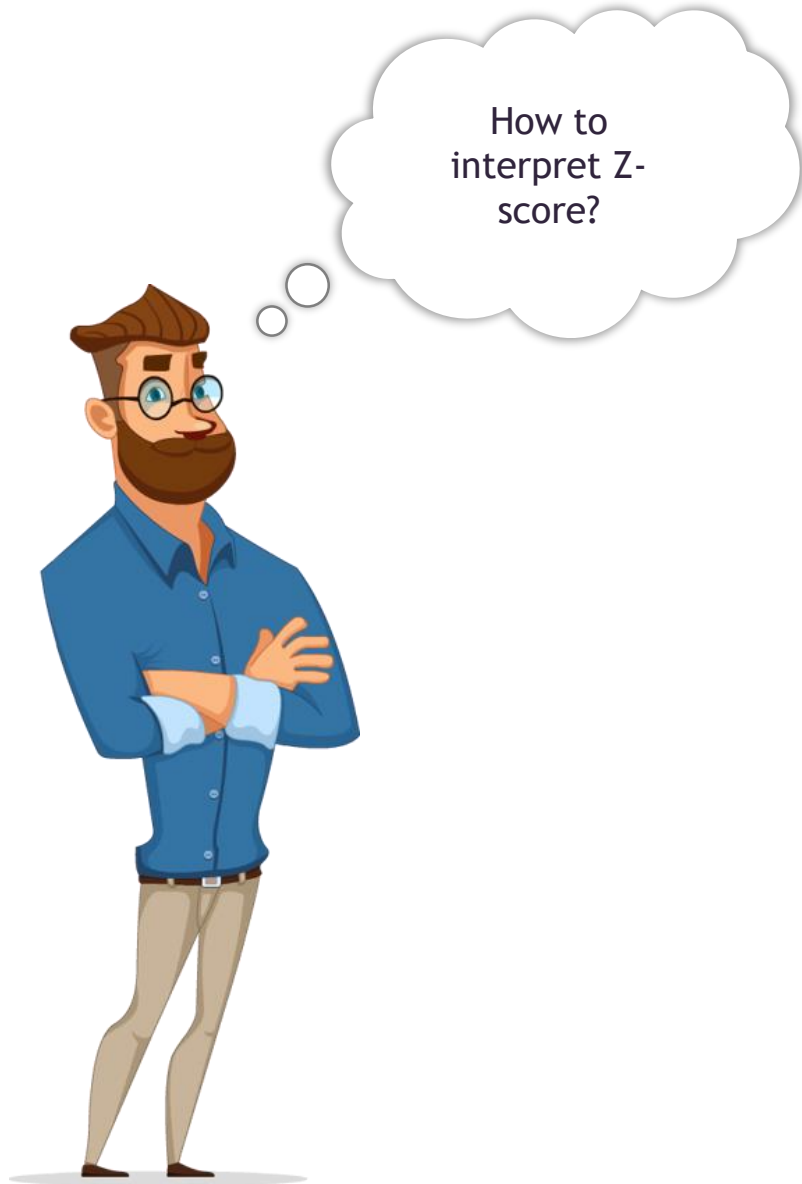
For example, the value of z-value = +1.8 which can be interpreted as the observed value is +1.8 standard deviations away from the mean.

P-values are probabilities.

$$z = \frac{(X - \mu)}{\sigma}$$

$$z = \frac{(190 - 150)}{25} = 1.6$$

Decision Rules



How to
interpret Z-
score?

$Z\text{-score} < 0$, an element less than the mean.

$Z\text{-score} > 0$, an element greater than the mean.

$Z\text{-score} = 0$, an element equal to the mean.

$Z\text{-score} = 1$, an element that is 1 standard deviation greater than the mean; $z\text{-score} = 2$, an element 2 standard deviations greater than the mean; etc.

$Z\text{-score} = -1$, an element that is 1 standard deviation less than the mean; $z\text{-score} = -2$, an element 2 standard deviations less than the mean; etc.

If the number of elements in the set is large, about 68% of the elements have a z-score between -1 and 1 ; about 95% have a z-score between -2 and 2 ; and about 99% have a z-score between -3 and 3 . $Z\text{-score} > 3$, an element is an outlier.

Decision Rules

A cartoon illustration of a man with a beard and glasses, wearing a blue shirt and tan pants, standing with his arms crossed and looking thoughtful. A thought bubble above him contains the text "What is P-value?".

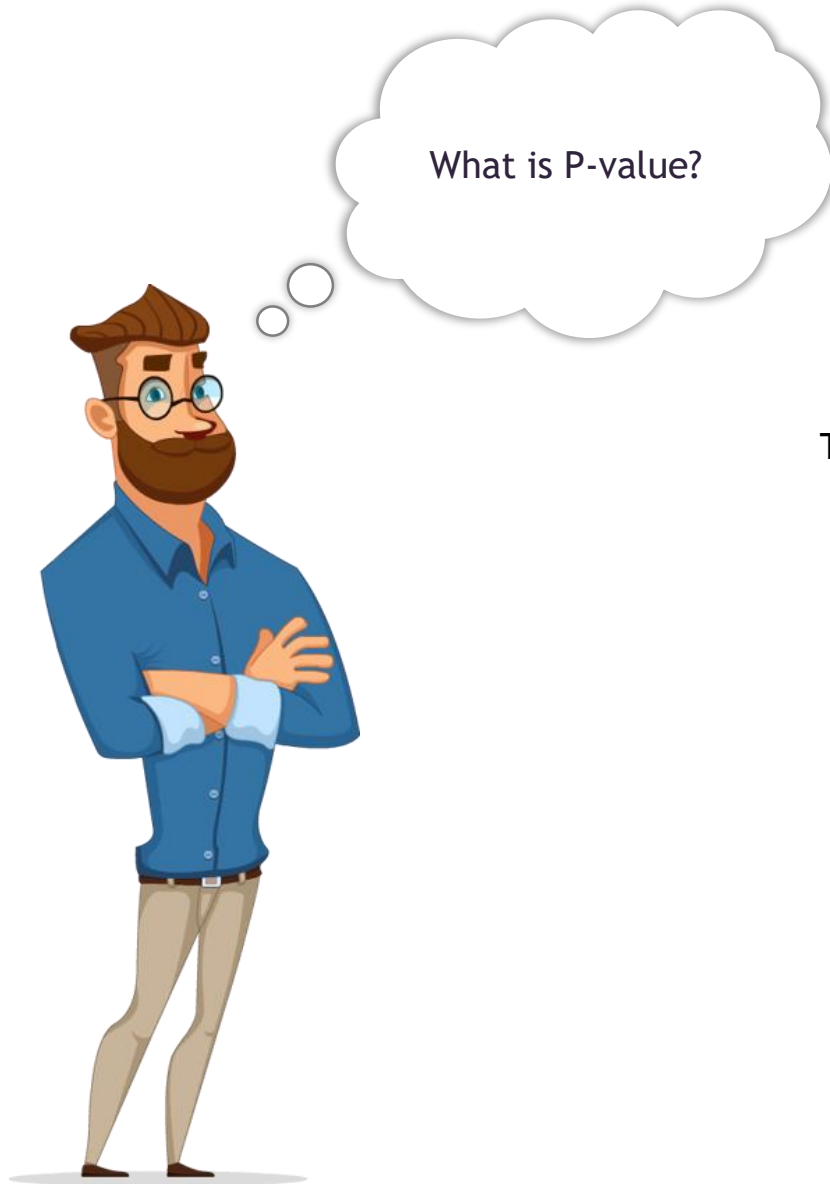
What is P-value?

When you perform a hypothesis test in statistics, a p-value helps you determine the significance of your results.

The p-value, or calculated probability, is the probability of finding the observed, or more extreme, results when the null hypothesis (H_0) of a study question is true.

The term “significance level” (alpha) is used to refer to a pre-chosen probability and the term “p-value” is used to indicate a probability that you calculate after a given study.

Decision Rules



The p-value is a number between 0 and 1 and interpreted in the following way:

- A small p-value (typically ≤ 0.05) indicates strong evidence against the null hypothesis, so you reject the null hypothesis.
- A large p-value (> 0.05) indicates weak evidence against the null hypothesis, so you fail to reject the null hypothesis.
- p-values very close to the cutoff (0.05) are considered to be marginal (could go either way). Always report the p-value so that your readers can draw their own conclusions.

Decision Rules



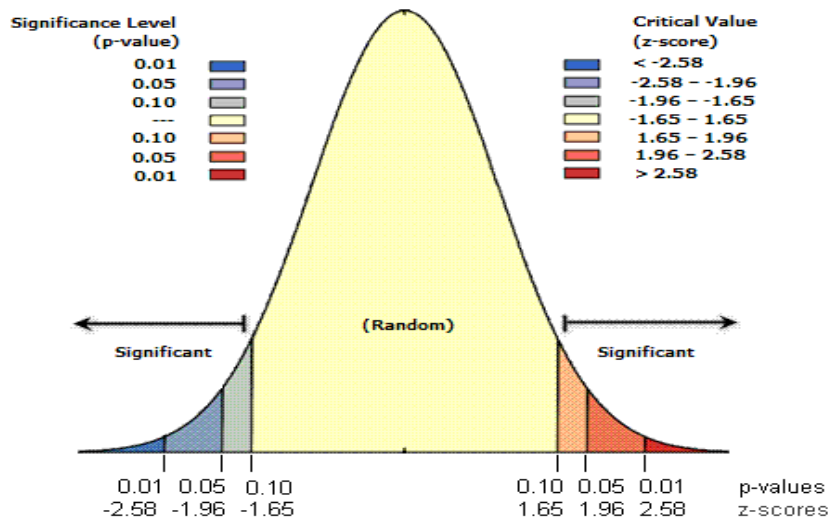
30 MINUTES OR LESS ON AVERAGE

Mean Delivery Time is 30 minutes max, is incorrect

p -value turns out to be 0.001, which is much less than 0.05

Decision Rules

z-scores and
p-value



z-score (Standard Deviations)	p-value (Probability)	Confidence level
< -1.65 or > +1.65	< 0.10	90%
< -1.96 or > +1.96	< 0.05	95%
< -2.58 or > +2.58	< 0.01	99%

When the absolute value of the z-score is large and the probabilities are small

Quiz

Which of the following statements describes descriptive statistics?

- a. Descriptive statistics involves organizing, displaying and describing data
- b. Descriptive statistics uses a random sample of data taken from a population to describe and make inferences about the population
- c. All of the above
- d. None of the above

Which of the following statements describes discrete data?

- a. Discrete data are continuous
- b. Discrete data are whole numbers and are usually a count of objects
- c. Discrete data have labels
- d. All of the above

Which of the following are four measurement scales?

- a. arrange, mutate, summary and ordinal
- b. mutate, nominal, interval and arrange
- c. nominal, ordinal, interval and ratio
- d. All of the above

Which of the following are the properties of normal distribution?

- a. Mean = median = mode
- b. Symmetry about the center
- c. The total area under the curve is 1
- d. All of the above

Which of the following is the formula for z-score?

- a. $z = x - \mu / \sigma$
- b. $z = x - \mu / s$
- c. $z = x / \sigma$
- d. None of the above

Thank You



India : +91-7847955955

US : 1-800-216-8930 (TOLL FREE)



sales@intellipaate.com



24/7 Chat with Our Course Advisor