

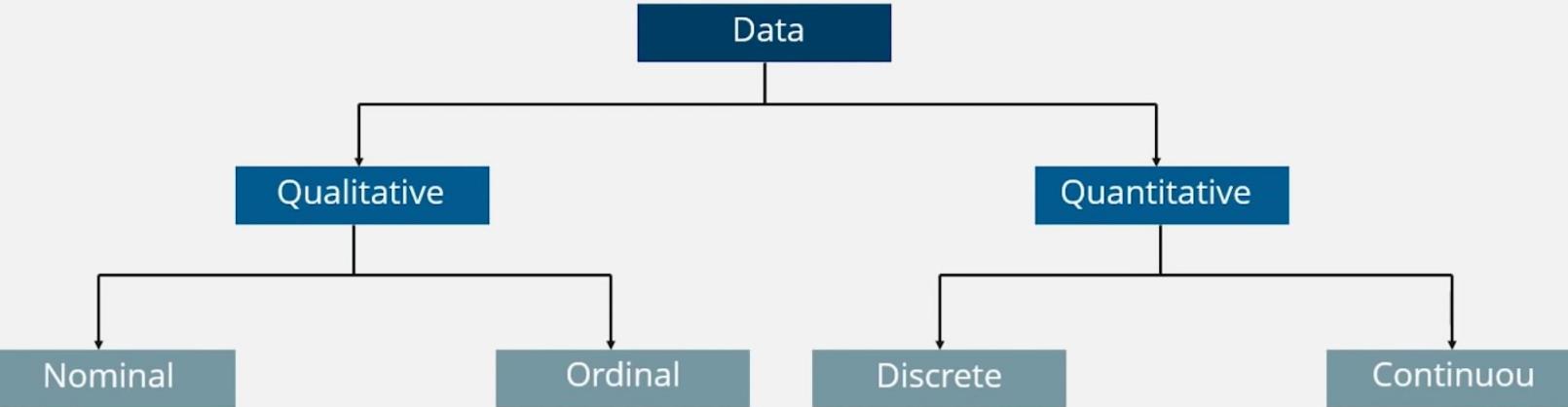
Statistic and Probability



What is Data?

WHAT IS DATA?

Data refers to facts and statistics collected together for reference or analysis.



Types Of Data

QUALITATIVE DATA

Qualitative data deals with characteristics and descriptors that can't be easily measured, but can be observed subjectively.

Nominal Data

Data with no inherent order or ranking such as gender or race, such kind of data is called Nominal data



Gender
Male
Female
Male
Male

Ordinal Data

Data with an ordered series, such as shown in the table, such kind of data is called Ordinal data

Customer ID	Rating
001	Good
002	Average
003	Average
004	Bad

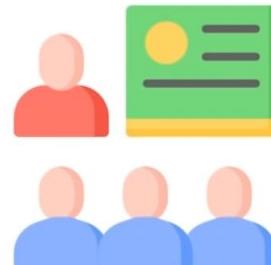
QUANTITATIVE DATA

Quantitative data deals with numbers and things you can measure objectively.

Discrete Data

Also known as categorical data, it can hold finite number of possible values.

Example: Number of students in a class



Continuous Data

Data that can hold infinite number of possible values.

Example: Weight of a person



WHAT IS STATISTICS?

Statistics is an area of applied mathematics concerned with the data collection, analysis, interpretation and presentation.



WHAT IS STATISTICS?

Statistics is an area of applied mathematics concerned with the data collection, analysis, interpretation and presentation.



Your company has created a new drug that may cure cancer. How would you conduct a test to confirm the drug's effectiveness?

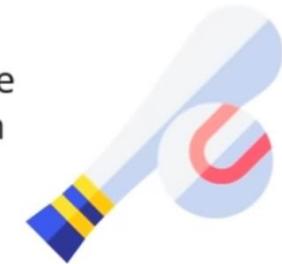


WHAT IS STATISTICS?

Statistics is an area of applied mathematics concerned with the data collection, analysis, interpretation and presentation.



You and a friend are at a baseball game, and out of the blue he offers you a bet that neither team will hit a home run in that game. Should you take the bet?



WHAT IS STATISTICS?

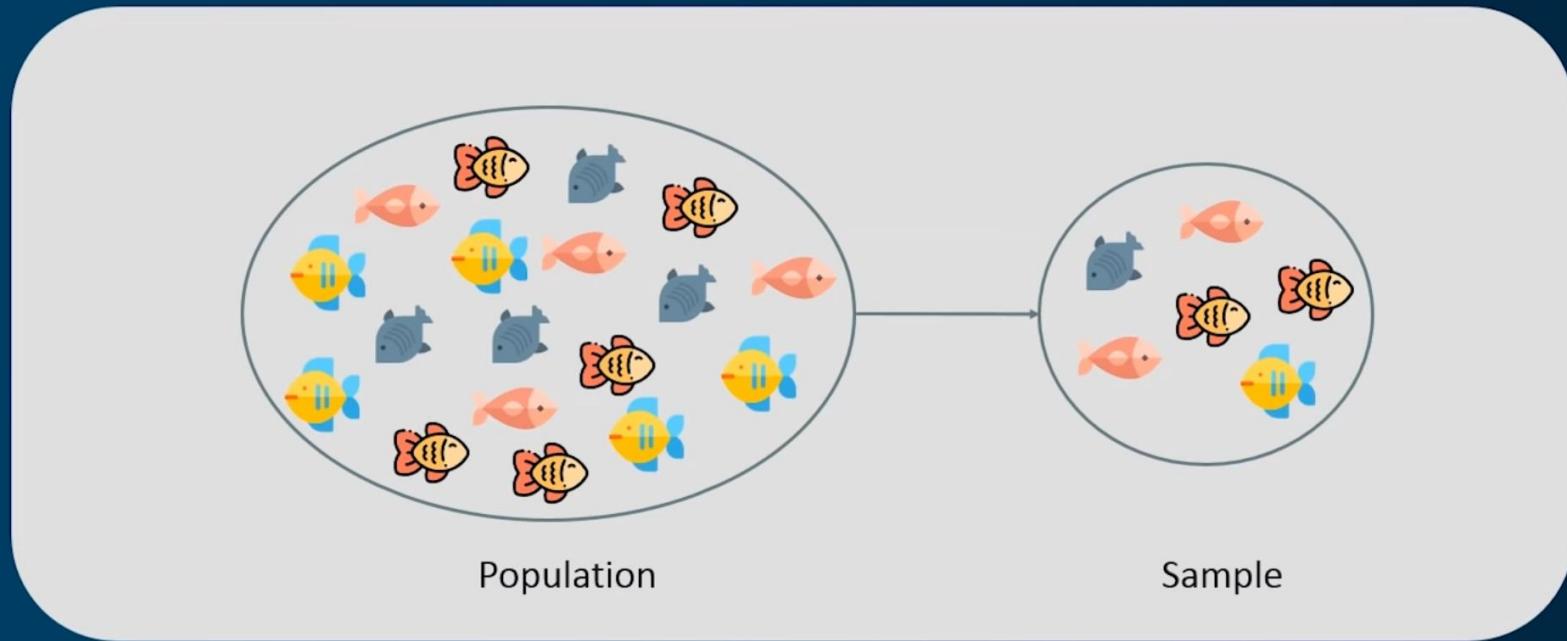
Statistics is an area of applied mathematics concerned with the data collection, analysis, interpretation and presentation.



The latest sales data have just come in, and your boss wants you to prepare a report for management on places where the company could improve its business. What should you look for? What should you not look for?



Basic Terminology in Statistics

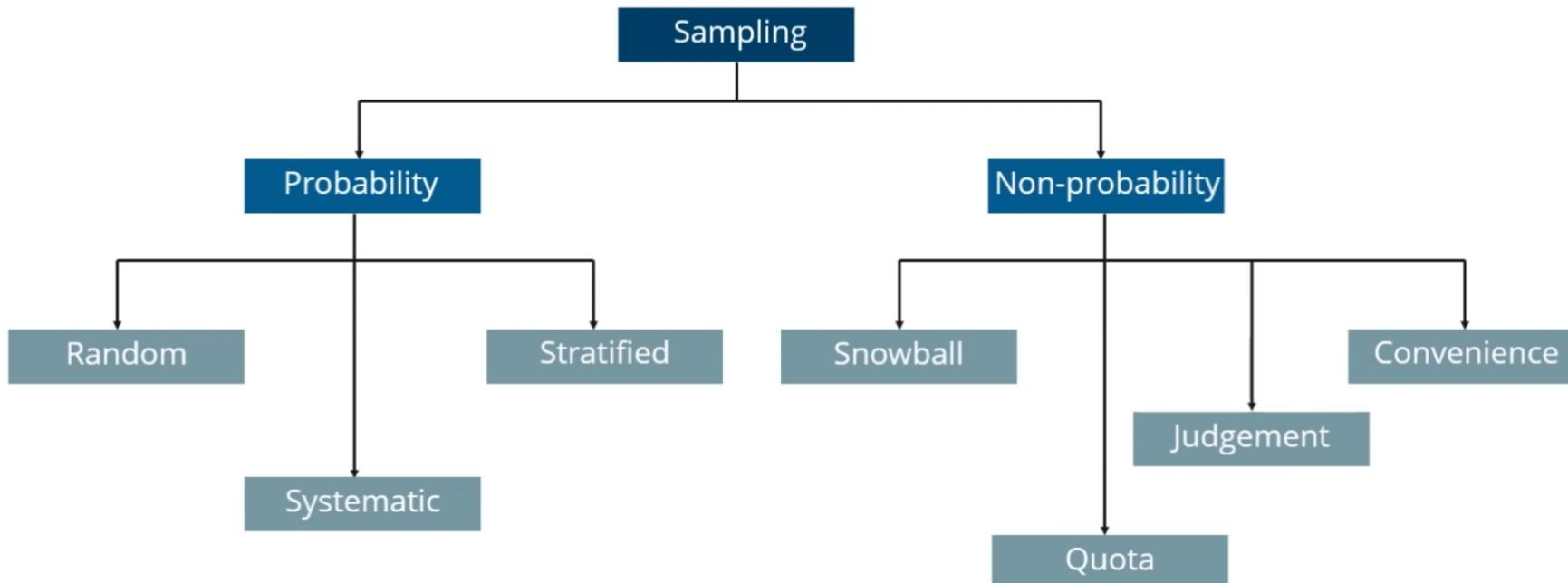


Statistics Terminologies

Population: A collection or set of individuals or objects or events whose properties are to be analyzed.

Sample: A subset of population is called 'Sample'. A well chosen sample will contain most of the information about a particular population parameter

SAMPLING TECHNIQUES



RANDOM SAMPLING

Random Sampling

Systematic Sampling

Stratified Sampling



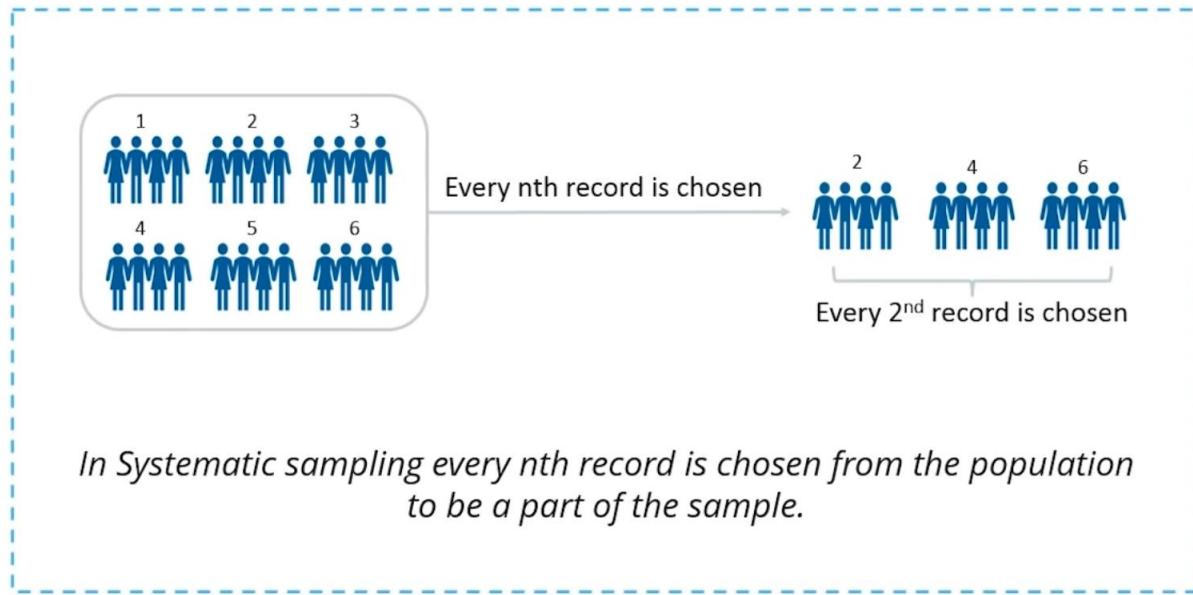
Each member of the population has equal chance of being selected in the sample.

SYSTEMATIC SAMPLING

Random Sampling

Systematic Sampling

Stratified Sampling

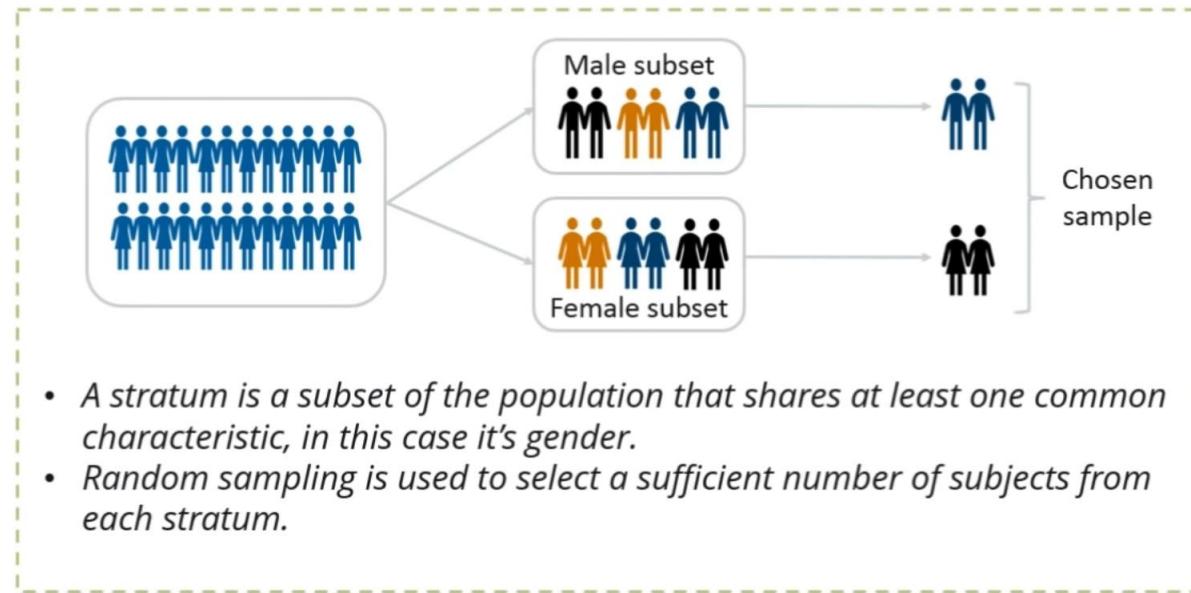


STRATIFIED SAMPLING

Random Sampling

Systematic Sampling

Stratified Sampling



Descriptive Statistics

DESCRIPTIVE STATISTICS

Descriptive statistics uses the data to provide descriptions of the population, either through numerical calculations or graphs or tables.



Maximum

Average

Minimum



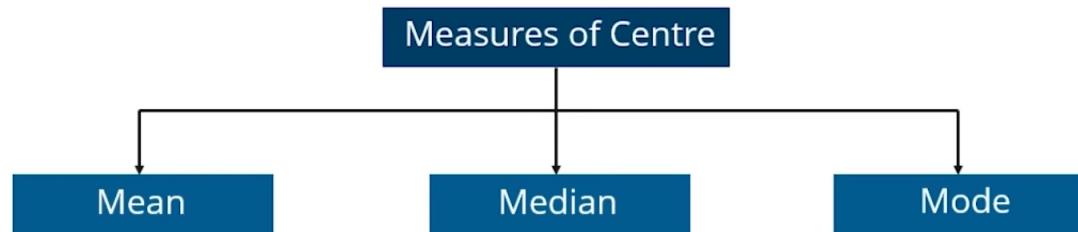
Descriptive Statistics is mainly focused upon the main characteristics of data. It provides graphical summary of the data.

DESCRIPTIVE STATISTICS

Descriptive statistics is a method used to describe and understand the features of a specific data set by giving short summaries about the sample and measures of the data.

Descriptive statistics are broken down into two categories:

- **Measures of Central tendency**
- Measures of Variability (spread)



Here is a sample dataset of cars containing the variables:

- Cars,
- Mileage per Gallon(mpg)
- Cylinder Type (cyl)
- Displacement (disp)
- Horse Power(hp)
- Real Axle Ratio(drat)

Cars	mpg	cyl	disp	hp	drat
MazdaRX4	21	6	160	110	3.9
MazdaRX4_WAG	21	6	160	110	3.9
Datsun_710	22.8	4	108	93	3.85
Alto	21.3	6	108	96	3
WagonR	23	4	150	90	4
Toyota_11	23	6	108	110	3.9
Honda_12	23	4	160	110	3.9
Ford_11	23	6	160	110	3.9

Mean

Measure of average of all the values in a sample is called Mean.

Here is a sample dataset of cars containing the variables:

- Cars,
- Mileage per Gallon(mpg)
- Cylinder Type (cyl)
- Displacement (disp)
- Horse Power(hp)
- Real Axle Ratio(drat)

Cars	mpg	cyl	disp	hp	drat
MazdaRX4	21	6	160	110	3.9
MazdaRX4_WAG	21	6	160	110	3.9
Datsun_710	22.8	4	108	93	3.85
Alto	21.3	6	108	96	3
WagonR	23	4	150	90	4
Toyota_11	23	6	108	110	3.9
Honda_12	23	4	160	110	3.9
Ford_11	23	6	160	110	3.9

Mean

To find out the average horsepower of the cars among the population of cars, we will check and calculate the average of all values:

$$\frac{110 + 110 + 93 + 96 + 90 + 110 + 110 + 110}{8} = 103.625$$

Here is a sample dataset of cars containing the variables:

- Cars,
- Mileage per Gallon(mpg)
- Cylinder Type (cyl)
- Displacement (disp)
- Horse Power(hp)
- Real Axle Ratio(drat)

Cars	mpg	cyl	disp	hp	drat
MazdaRX4	21	6	160	110	3.9
MazdaRX4_WAG	21	6	160	110	3.9
Datsun_710	22.8	4	108	93	3.85
Alto	21.3	6	108	96	3
WagonR	23	4	150	90	4
Toyota_11	23	6	108	110	3.9
Honda_12	23	4	160	110	3.9
Ford_11	23	6	160	110	3.9

Median

Measure of the central value of the sample set is called **Median**.

Here is a sample dataset of cars containing the variables:

- Cars,
- Mileage per Gallon(mpg)
- Cylinder Type (cyl)
- Displacement (disp)
- Horse Power(hp)
- Real Axle Ratio(drat)

Cars	mpg	cyl	disp	hp	drat
MazdaRX4	21	6	160	110	3.9
MazdaRX4_WAG	21	6	160	110	3.9
Datsun_710	22.8	4	108	93	3.85
Alto	21.3	6	108	96	3
WagonR	23	4	150	90	4
Toyota_11	23	6	108	110	3.9
Honda_12	23	4	160	110	3.9
Ford_11	23	6	160	110	3.9

Median

To find out the center value of mpg among the population of cars, arrange records in *Ascending order*, i.e., **21, 21, 21.3, 22.8, 23, 23, 23, 23**

In case of even entries, take average of the two middle values, i.e. $(22.8+23)/2 = 22.9$

Here is a sample dataset of cars containing the variables:

- Cars,
- Mileage per Gallon(mpg)
- Cylinder Type (cyl)
- Displacement (disp)
- Horse Power(hp)
- Real Axle Ratio(drat)

Cars	mpg	cyl	disp	hp	drat
MazdaRX4	21	6	160	110	3.9
MazdaRX4_WAG	21	6	160	110	3.9
Datsun_710	22.8	4	108	93	3.85
Alto	21.3	6	108	96	3
WagonR	23	4	150	90	4
Toyota_11	23	6	108	110	3.9
Honda_12	23	4	160	110	3.9
Ford_11	23	6	160	110	3.9

Mode

The value most recurrent in the sample set is known as Mode.

Here is a sample dataset of cars containing the variables:

- Cars,
- Mileage per Gallon(mpg)
- Cylinder Type (cyl)
- Displacement (disp)
- Horse Power(hp)
- Real Axle Ratio(drat)

Cars	mpg	cyl	disp	hp	drat
MazdaRX4	21	6	160	110	3.9
MazdaRX4_WAG	21	6	160	110	3.9
Datsun_710	22.8	4	108	93	3.85
Alto	21.3	6	108	96	3
WagonR	23	4	150	90	4
Toyota_11	23	6	108	110	3.9
Honda_12	23	4	160	110	3.9
Ford_11	23	6	160	110	3.9

Mode

To find the most common type of cylinder among the population of cars, check the value which is repeated most number of times, i.e., cylinder type 6

MEASURES OF SPREAD

A measure of spread, sometimes also called a measure of dispersion, is used to describe the variability in a sample or population.

Range

Inter Quartile Range

Variance

Standard Deviation

Range is the given measure of how spread apart the values in a dataset are.

$$\text{Range} = \text{Max}(x_i) - \text{Min}(x_i)$$

MEASURES OF SPREAD

A measure of spread, sometimes also called a measure of dispersion, is used to describe the variability in a sample or population.

Range

Inter Quartile Range

Variance

Standard Deviation

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

Q1 Q2 Q3

1	2	3	4	5	6	7	8
---	---	---	---	---	---	---	---

MEASURES OF SPREAD

Consider the marks of the 100 students below, ordered from the lowest to the highest scores

The first quartile (Q1) lies between the 25th and 26th.
 $Q1 = (45 + 45) \div 2 = 45$

Order	Score								
1st	35	21st	42	41st	53	61st	64	81st	74
2nd	37	22nd	42	42nd	53	62nd	64	82nd	74
3rd	37	23rd	44	43rd	54	63rd	65	83rd	74
4th	38	24th	44	44th	55	64th	66	84th	75
5th	39	25th	45	45th	55	65th	67	85th	75
6th	39	26th	45	46th	56	66th	67	86th	76
7th	39	27th	45	47th	57	67th	67	87th	77
8th	39	28th	45	48th	57	68th	67	88th	77
9th	39	29th	47	49th	58	69th	68	89th	79
10th	40	30th	48	50th	58	70th	69	90th	80
11th	40	31st	49	51st	59	71st	69	91st	81
12th	40	32nd	49	52nd	60	72nd	69	92nd	81
13th	40	33rd	49	53rd	61	73rd	70	93rd	81
14th	40	34th	49	54th	62	74th	70	94th	81
15th	40	35th	51	55th	62	75th	71	95th	81
16th	41	36th	51	56th	62	76th	71	96th	81
17th	41	37th	51	57th	63	77th	71	97th	83
18th	42	38th	51	58th	63	78th	72	98th	84
19th	42	39th	52	59th	64	79th	74	99th	84
20th	42	40th	52	60th	64	80th	74	100th	85

The second quartile (Q2) between the 50th and 51st.
 $Q2 = (58 + 59) \div 2 = 58.5$

The third quartile (Q3) between the 75th and 76th.
 $Q3 = (71 + 71) \div 2 = 71$

MEASURES OF SPREAD

A measure of spread, sometimes also called a measure of dispersion, is used to describe the variability in a sample or population.

Range

Inter Quartile Range

Variance

Standard Deviation

Inter Quartile Range(IQR) is the measure of variability, based on dividing a dataset into quartiles.

- *Quartiles divide a rank-ordered data set into four equal parts, denoted by Q1, Q2, and Q3, respectively*
- *The interquartile range is equal to Q3 minus Q1, i.e.. $IQR = Q3 - Q1$*

MEASURES OF SPREAD

A measure of spread, sometimes also called a measure of dispersion, is used to describe the variability in a sample or population.

Range

Inter Quartile Range

Variance

Standard Deviation

Standard Deviation is the measure of the dispersion of a set of data from its mean.

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

www.edureka.co/data-science

SUBS
CRIBE
e! ▶

STANDARD DEVIATION

Standard Deviation Use Case: Daenerys has 20 Dragons. They have the numbers 9, 2, 5, 4, 12, 7, 8, 11, 9, 3, 7, 4, 12, 5, 4, 10, 9, 6, 9, 4. Work out the Standard Deviation.

STEP 1

Find out the mean for your sample set.

The Mean is:

$$\frac{9+2+5+4+12+7+8+11+9+3+7+4+12+5+4+10+9+6+9+4}{20}$$

$$\mu=7$$

STANDARD DEVIATION

Standard Deviation Use Case: Daenerys has 20 Dragons. They have the numbers 9, 2, 5, 4, 12, 7, 8, 11, 9, 3, 7, 4, 12, 5, 4, 10, 9, 6, 9, 4. Work out the Standard Deviation.

STEP 2

Then for each number, subtract the Mean and square the result.

$$(x_i - \mu)^2$$

$$(9-7)^2 = 2^2 = 4$$

$$(2-7)^2 = (-5)^2 = 25$$

$$(5-7)^2 = (-2)^2 = 4$$

And so on...

□ We get the following results:

4, 25, 4, 9, 25, 0, 1, 16, 4, 16, 0, 9, 25, 4, 9, 9, 4, 1, 4, 9

STANDARD DEVIATION

Standard Deviation Use Case: Daenerys has 20 Dragons. They have the numbers 9, 2, 5, 4, 12, 7, 8, 11, 9, 3, 7, 4, 12, 5, 4, 10, 9, 6, 9, 4. Work out the Standard Deviation.

STEP 3

Then work out the mean of those squared differences.

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$
$$\frac{4+25+4+9+25+0+1+16+4+16+0+9+25+4+9+9+4+1+4+9}{20}$$

$\square \sigma^2 = 8.9$

Probability

WHAT IS PROBABILITY?

Probability is the measure of how likely an event will occur.

- Probability is the ratio of desired outcomes to total outcomes:
(desired outcomes) / (total outcomes)
 - Probabilities of all outcomes always sums to 1
- Example:
- On rolling a dice, you get 6 possible outcomes
 - Each possibility only has one outcome, so each has a probability of 1/6
 - For example, the probability of getting a number '2' on the dice is 1/6



TERMINOLOGIES IN PROBABILITY

Random Experiment

An experiment or a process for which the outcome cannot be predicted with certainty



Sample Space

The entire possible set of outcomes of a random experiment is the sample space (S) of that experiment

Event

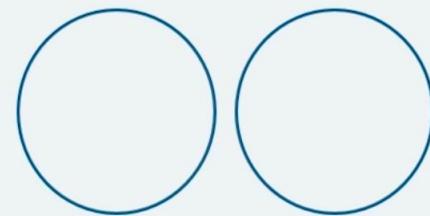
One or more outcomes of an experiment. It is a subset of sample space(S)



TYPES OF EVENTS

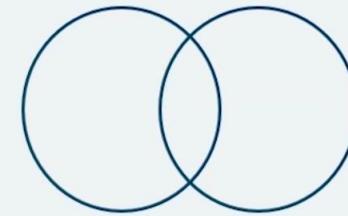
Disjoint Events do not have any common outcomes.

- The outcome of a ball delivered cannot be a sixer and a wicket
- A single card drawn from a deck cannot be a king and a queen
- A man cannot be dead and alive



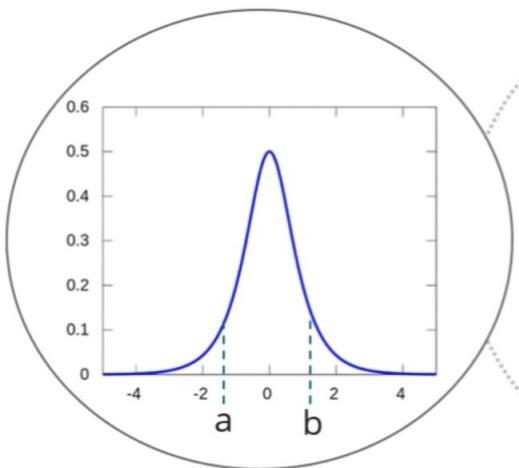
Non-Disjoint Events can have common outcomes

- A student can get 100 marks in statistics and 100 marks in probability
- The outcome of a ball delivered can be a no ball and a six



PROBABILITY DENSITY FUNCTION

The equation describing a continuous probability distribution is called a Probability Density Function



Property 01



Graph of a PDF will be continuous over a range



Property 02



Area bounded by the curve of density function and the x-axis is equal to 1



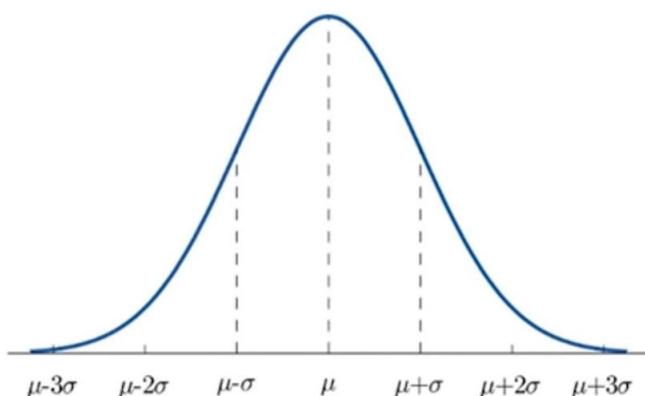
Property 03



Probability that a random variable assumes a value between a & b is equal to the area under the PDF bounded by a & b

NORMAL DISTRIBUTION

The Normal Distribution is a probability distribution that associates the normal random variable X with a cumulative probability



$$Y = [1/\sigma * \sqrt{2\pi}] * e^{-(x - \mu)^2/2\sigma^2}$$

Where,

- X is a normal random variable
- μ is the mean and
- σ is the standard deviation

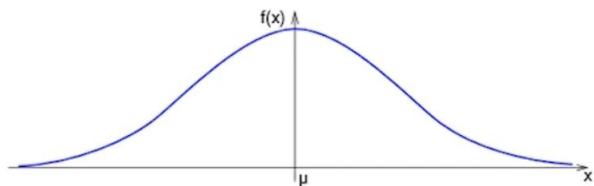


Note: Normal Random variable is variable with mean at 0 and variance equal to 1

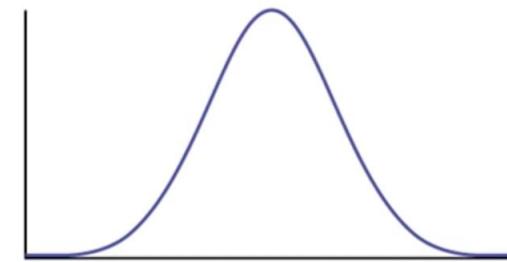
STANDARD DEVIATION & CURVE

The graph of the Normal Distribution depends on two factors: the *Mean* and the *Standard Deviation*

- **Mean:** Determines the location of center of the graph
- **Standard Deviation:** Determines the height of the graph



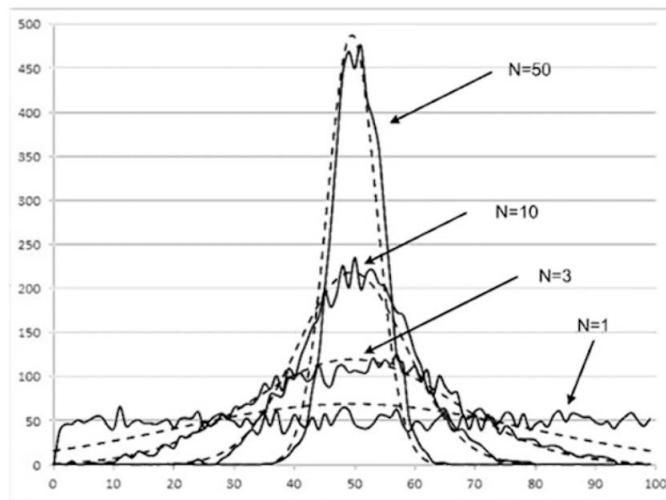
If the standard deviation is large,
the curve is short and wide.



If the standard deviation is small,
the curve is tall and narrow.

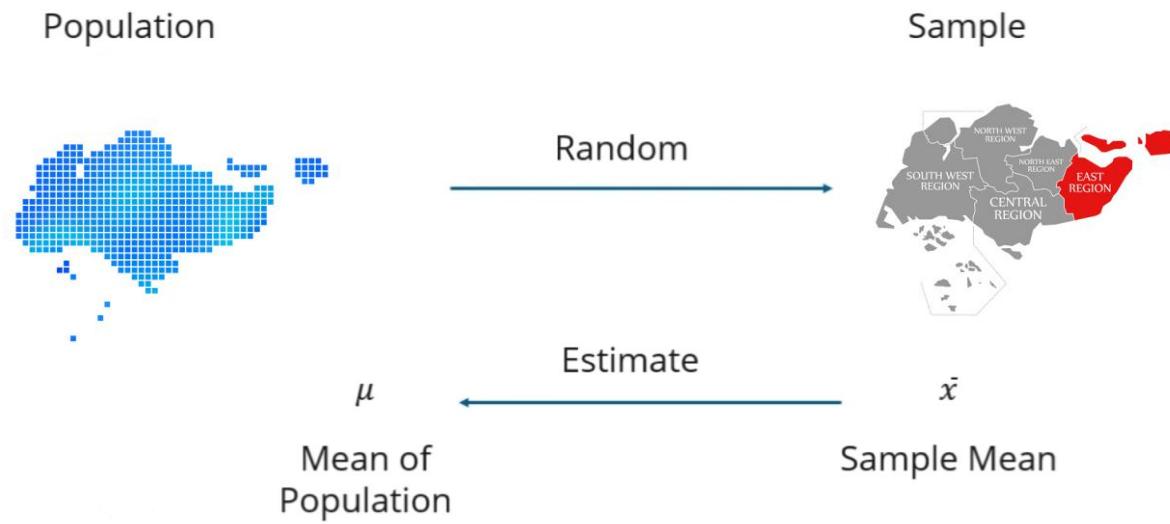
CENTRAL LIMIT THEOREM

The **Central Limit Theorem** states that the sampling distribution of the mean of any independent, random variable will be normal or nearly normal, if the sample size is large enough



POINT ESTIMATION

Point Estimation is concerned with the use of the sample data to measure a single value which serves as an approximate value or the best estimate of an unknown population parameter.



FINDING THE ESTIMATES



Method of Moments

Estimates are found out by equating the first k sample moments to the corresponding k population moments

Maximum of Likelihood

Uses a model and the values in the model to maximize a likelihood function. This results in the most likely parameter for the inputs selected

Bayes' Estimators

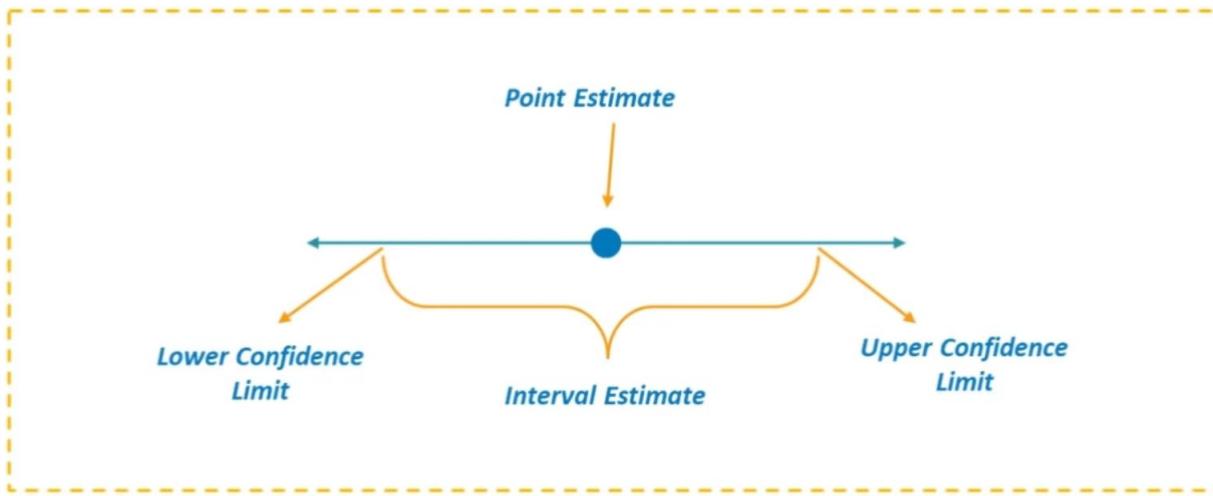
Minimizes the average risk (an expectation of random variables)

Best Unbiased Estimators

Several unbiased estimators can be used to approximate a parameter (which one is “best” depends on what parameter you are trying to find)

INTERVAL ESTIMATE

An Interval, or range of values, used to estimate a population parameter is called Interval Estimate.



CONFIDENCE INTERVAL

01

Confidence Interval is the measure of your confidence, that the interval estimate contains the population mean, μ

Statisticians use a confidence interval to describe the amount of uncertainty associated with a sample estimate of a population parameter

02

03

Technically, a range of values so constructed that there is a specified probability of including the true value of a parameter within it

MARGIN OF ERROR

- Difference between the point estimate and the actual population parameter value is called the **Sampling Error**
- When μ is estimated, the sampling error is the difference $\mu - \bar{x}$

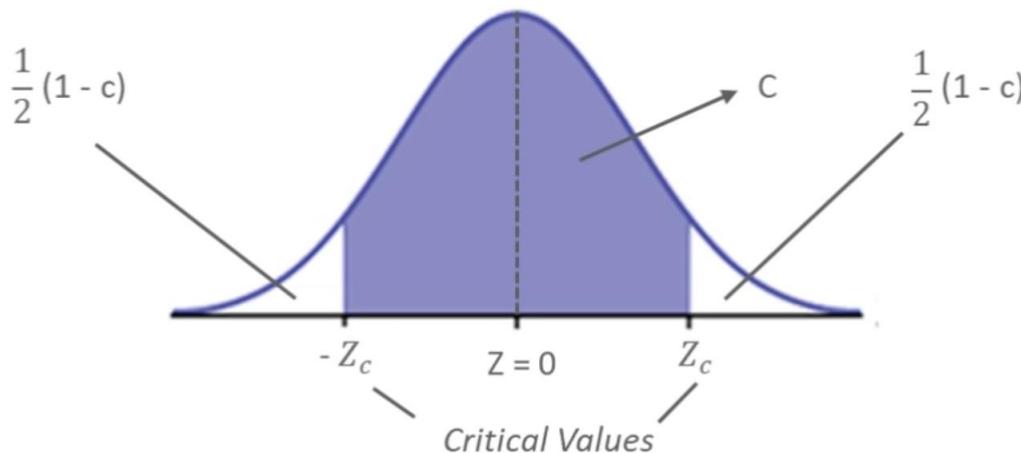
Margin of Error E, for a given level of confidence is the greatest possible distance between the point estimate and the value of the parameter it is estimating



$$E = Z_c \frac{\sigma}{\sqrt{n}}$$

ESTIMATING LEVEL OF CONFIDENCE

The level of confidence c , is the probability that the interval estimate contains the population parameter.



C is the area beneath the normal curve between the critical values
Corresponding Z score can be calculated using the standard normal table

MARGIN OF ERROR - USE CASE

A random sample of 32 textbook prices is taken from a local college bookstore. The mean of the sample is $\bar{x} = 74.22$, and the sample standard deviation is $S = 23.44$. Use a 95% confidence level and find the margin of error for the mean price of all textbooks in the bookstore

You know by formula,

$$E = Z_c \frac{\sigma}{\sqrt{n}}$$

$$E = 1.96 * (23.44/\sqrt{32}) \approx 8.12$$

HYPOTHESIS TESTING

Statisticians use hypothesis testing to formally check whether the hypothesis is accepted or rejected.

Hypothesis testing is conducted in the following manner:

- ❖ **State the Hypotheses** – This stage involves stating the null and alternative hypotheses.
- ❖ **Formulate an Analysis Plan** – This stage involves the construction of an analysis plan.
- ❖ **Analyse Sample Data** – This stage involves the calculation and interpretation of the test statistic as described in the analysis plan.
- ❖ **Interpret Results** – This stage involves the application of the decision rule described in the analysis plan.

HYPOTHESIS TESTING EXAMPLE



Nick



John



Bob



Harry



$$P(\text{John not picked for a day}) = \frac{3}{4}$$

$$P(\text{John not picked for 3 days}) = \frac{3}{4} \times \frac{3}{4} \times \frac{3}{4} = 0.42 \text{ (approx)}$$

$$P(\text{John not picked for 12 days}) = \left(\frac{3}{4}\right)^{12} = \mathbf{0.032} < 0.05$$

HYPOTHESIS TESTING EXAMPLE



Nick



John



Bob



Harry



Null Hypothesis (H_0) : Result is no different from assumption.

Alternate Hypothesis (H_a) : Result disproves the assumption.

Probability of Event < 0.05 (5%)