

[Connect, Message, Like, Follow & Share, 100% Free Counselling → Thank You](#)



Statistics for Data Science



Contents

1• Importance Of Statistics

2• Type Of Analytics

3• Probability

4• Properties Of Statistics

5• Central Tendency

6• Variability

7• Relationship Between Variables

8• Probability Distribution

9• Hypothesis Testing And
Statistical Significance

10• Regression

Importance of Statistics

- 1) Using various statistical tests, determine the relevance of features.
- 2) To avoid the risk of duplicate features, find the relationship between features.
- 3) Putting the features into the proper format.
- 4) Data normalization and scaling This step also entails determining the distribution of data as well as the nature of data.
- 5) Taking the data for further processing and making the necessary modifications.
- 6) Determine the best mathematical approach/model after processing the data.
- 7) After the data are acquired, they are checked against the various accuracy measuring scales.

Acknowledge the Different Types of Analytics in Statistics



1. Descriptive Analytics – What happened?

It tells us what happened in the past and helps businesses understand how they are performing by providing context to help stakeholders interpret data.

Descriptive analytics should serve as a starting point for all organizations. This type of analytics is used to answer the fundamental question “what happened?” by analyzing data, which is often historical.

It examines past events and attempts to identify specific patterns within the data. When people talk about traditional business intelligence, they’re usually referring to Descriptive Analytics.

Pie charts, bar charts, tables, and line graphs are common visualizations for Description Analytics.

This is the level at which you should begin your analytics journey because it serves as the foundation for the other three tiers. To move forward with your analytics, you must first determine what happened.

Consider some sales use cases to gain a better understanding of this. For instance, how many sales occurred in the previous quarter? Was it an increase or a decrease?



2. Diagnostic Analytics – Why did it happen?

It goes beyond descriptive data to assist you in comprehending why something occurred in the past.

Diagnostic analytics is the next step in analytics, a sort of advanced analytics that examines data or content to answer the question, “Why did it happen?”

Drill-down, data discovery, data processing, and correlations are several techniques used.

This is the second step because you want to first understand what occurred to work out why it occurred. Typically, once an organisation has achieved descriptive insights, diagnostics will be applied with a bit more effort.



Diagnostic analytics diagnoses issues based on data relationships, identifying patterns, and discovering anomalies in the data to help users answer questions.

3. Predictive Analytics – What is likely to happen?

It forecasts what is likely to happen in the future and provides businesses with data-driven actionable insights.

Once an organisation encompasses a firm grasp on what happened and why it happened, it can advance to the subsequent level of analytics, Predictive.

Predictive Analytics is another style of advanced analytics that seeks to answer the question “What is probably going to happen?” using data and knowledge.

The transition from Predictive Analytics to Diagnostics Analytics is critical. multivariate analysis, forecasting, multivariate statistics, pattern matching, predictive modelling, and forecasting are all a part of predictive analytics.

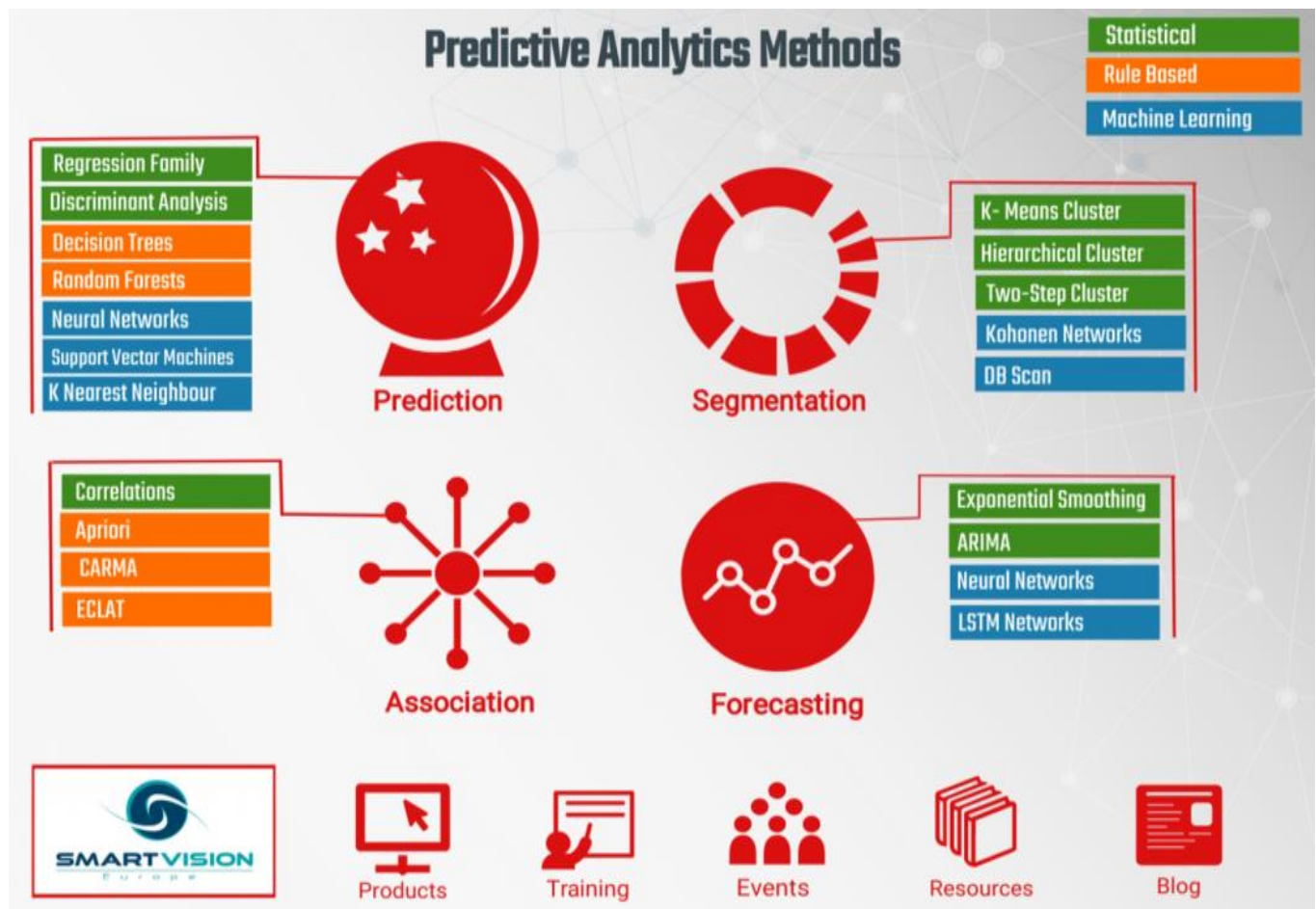
These techniques are more difficult for organisations to implement because they necessitate large amounts of high-quality data. Furthermore, these techniques necessitate a thorough understanding of statistics as well as programming languages such as R and Python.

Many organisations may lack the internal expertise required to effectively implement a predictive model.

So, why should any organisation bother with it? Although it can be difficult to achieve, the value that Predictive Analytics can provide is enormous.

A Predictive Model, for example, will use historical data to predict the impact of the next marketing campaign on customer engagement.

If a company can accurately identify which action resulted in a specific outcome, it can predict which actions will result in the desired outcome. These types of insights are useful in the next stage of analytics.



4. Prescriptive Analytics – What should be done?

It makes recommendations for actions that will capitalise on the predictions and guide the potential actions toward a solution.

Prescriptive analytics is the final and most advanced level of analytics.

Prescriptive Analytics is an analytics method that analyses data to answer the question “What should be done?”

Techniques used in this type of analytics include graph analysis, simulation, complex event processing, neural networks, recommendation engines, heuristics, and machine learning.

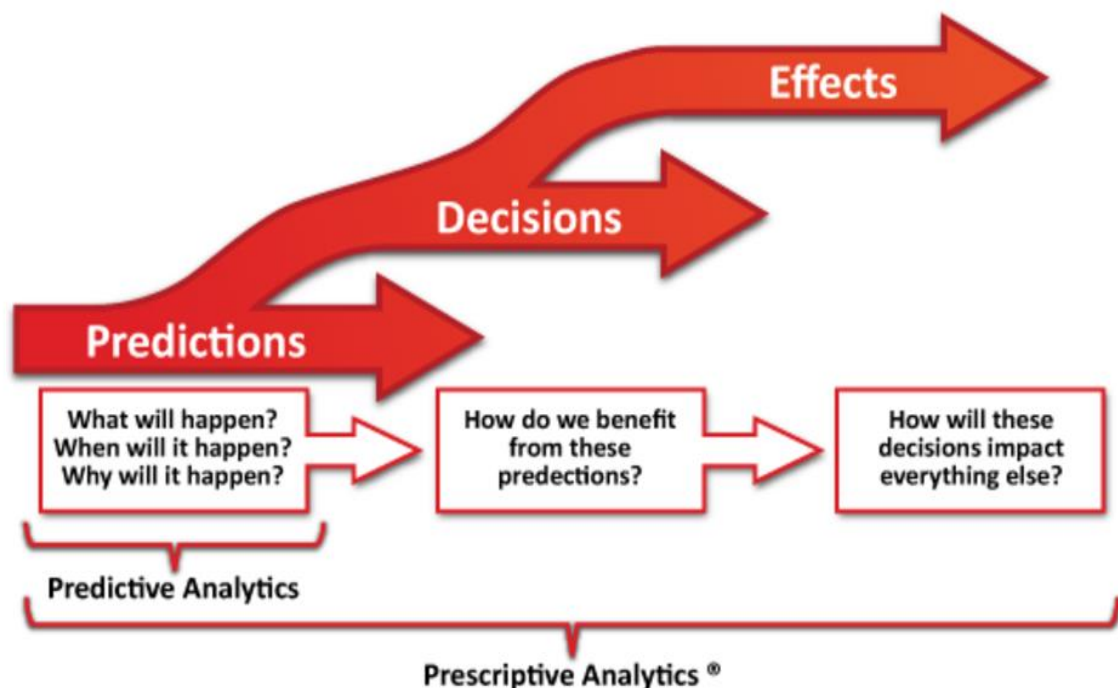
This is the toughest level to reach. The accuracy of the three levels of the analytics below has a significant impact on the dependability of Prescriptive

Analytics. The techniques required to obtain an effective response from a prescriptive analysis are determined by how well an organisation has completed each level of analytics.

Considering the quality of data required, the appropriate data architecture to facilitate it, and the expertise required to implement this architecture, this is not an easy task.

Its value is that it allows an organisation to make decisions based on highly analysed facts rather than instinct. That is, they are more likely to achieve the desired outcome, such as increased revenue.

Once again, a use case for this type of analytics in marketing would be to assist marketers in determining the best mix of channel engagement. For instance, which segment is best reached via email?



PROBABILITY

In a Random Experiment, the probability is a measure of the likelihood that an event will occur. The number of favorable outcomes in an experiment with n outcomes is denoted by x . The following is the formula for calculating the probability of an event.

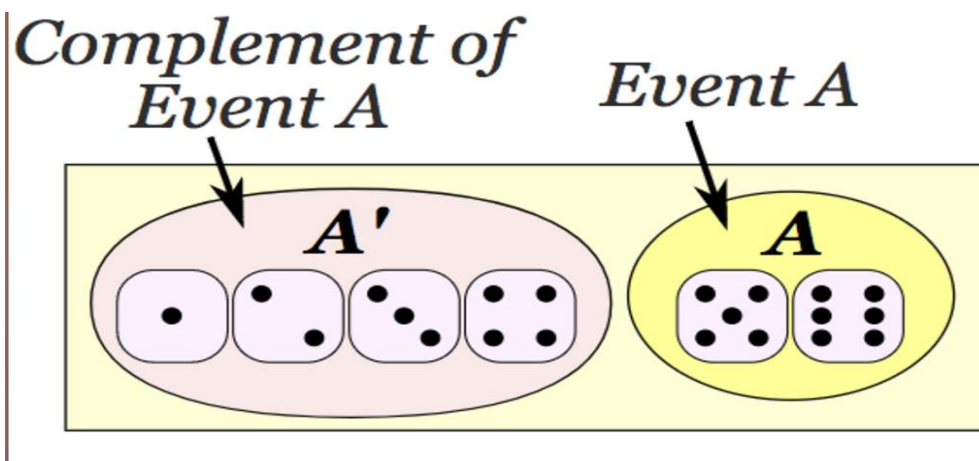
$$\text{Probability (Event)} = \text{Favourable Outcomes} / \text{Total Outcomes} = x/n$$

Let's look at a simple application to better understand probability. If we need to know if it's raining or not. There are two possible answers to this question: "Yes" or "No." It is possible that it will rain or not rain. In this case, we can make use of probability. The concept of probability is used to forecast the outcomes of coin tosses, dice rolls, and card draws from a deck of playing cards.

Properties of Statistics

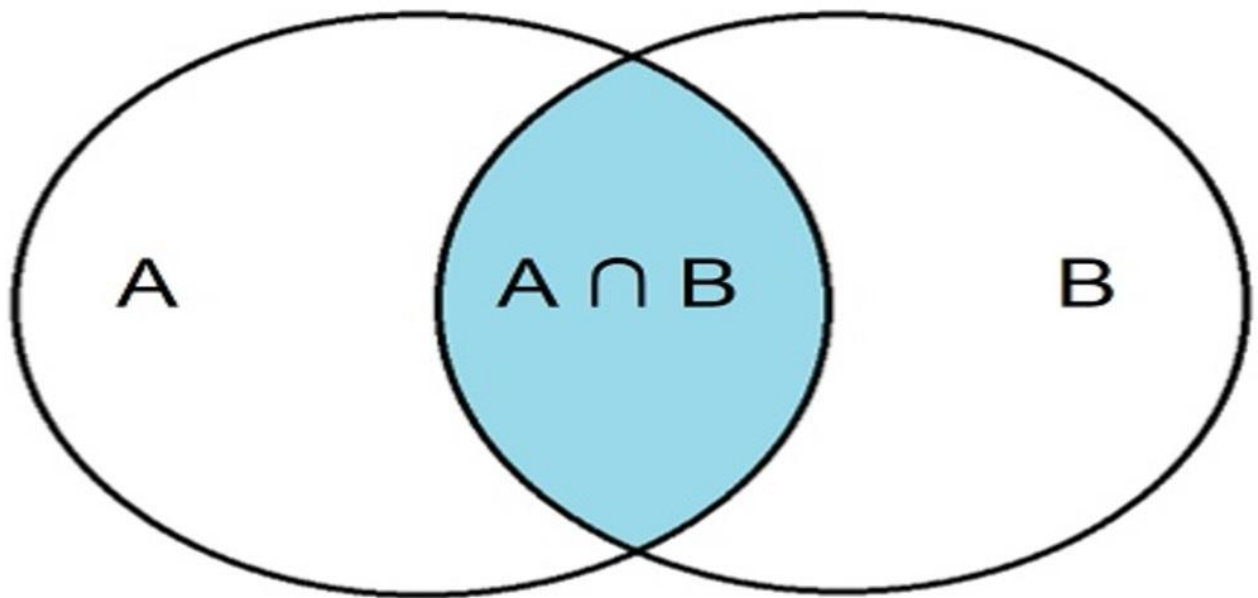
· **Complement:** A' , the complement of an event A in a sample space S , is the collection of all outcomes in S that are not members of set A . It is equivalent to rejecting any verbal description of event A .

$$P(A) + P(A') = 1$$



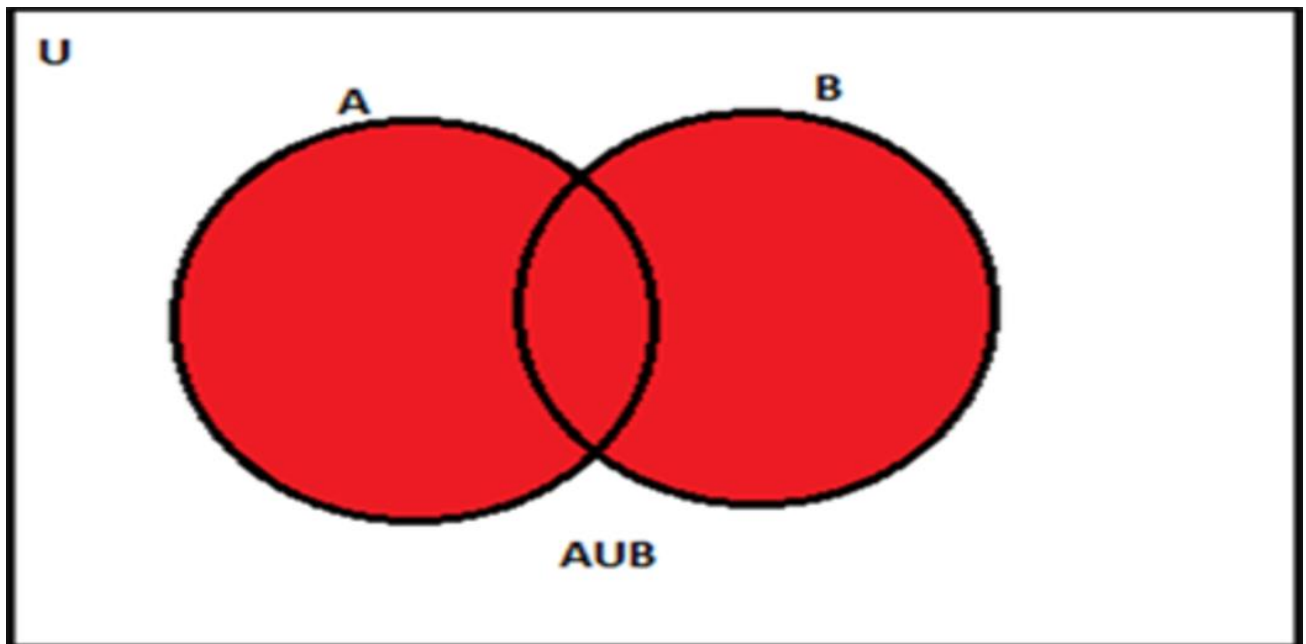
- **Intersection:** The intersection of events is a collection of all outcomes that are components of both sets A and B. It is equivalent to combining descriptions of the two events with the word “and.”

$$P(A \cap B) = P(A)P(B)$$



- **Union:** The union of events is the collection of all outcomes that are members of one or both sets A and B. It is equivalent to combining descriptions of the two events with the word “or.”

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



- **Conditional Probability:** $P(A|B)$ is a measure of the likelihood of one event happening in relation to one or more other events. When $P(B) > 0$,

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

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

Probability of
A and B
Probability of
A given B
Probability of B

Conditional Probability Formula

- **Independent Events:** Two events are considered independent if the occurrence of one has no effect on the likelihood of the occurrence of the other.
 $P(A|B) = P(A)P(B)$, where $P(A) \neq 0$ and $P(B) \neq 0$, $P(A|B) = P(A)$, $P(B|A) = P(B)$,
 $P(A|B) = P(A)$, $P(A|B) = P(A)$, $P(B|A) = P(B)$, $P(B|A) = P(B)$, $P(B|A) = P(B)$, $P(B|A) = P(B)$

Independent Events

If A and B are *independent events*, then

$$\begin{aligned}P(B | A) &= P(B) \\P(A | B) &= P(A)\end{aligned}$$

Test for Independent Events

Events A and B are independent events if and only if

$$P(A \cap B) = P(A)P(B)$$

Note: this generalizes to more than two independent events.

· **Mutually Exclusive Events:** If events A and B share no elements, they are mutually exclusive. Because A and B have no outcomes in common, it is impossible for both A and B to occur on a single trial of the random experiment. This results in the following rule

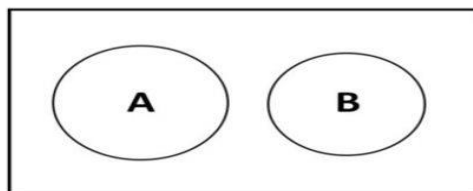
$$P(A \cap B) = 0$$

Any event A and its complement A^c are mutually exclusive if and only if A and B are mutually exclusive, but A and B can be mutually exclusive without being complements.

Mutually Exclusive Events

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

$$P(A \text{ and } B) = 0$$



Note: No overlap between A and B

· **Bayes' Theorem:** it is a method for calculating conditional probability. The probability of an event occurring if it is related to one or more other events is known as conditional probability. For example, your chances of finding a parking space are affected by the time of day you park, where you park, and what conventions are taking place at any given time.

The diagram illustrates Bayes' Theorem with the following components:

- LIKELIHOOD**
The probability of "B" being True, given "A" is True
- PRIOR**
The probability "A" being True. This is the knowledge.
- POSTERIOR**
The probability of "A" being True, given "B" is True
- MARGINALIZATION**
The probability "B" being True.

The equation is shown as:

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

Arrows indicate the mapping: LIKELIHOOD points to $P(B|A)$, PRIOR points to $P(A)$, POSTERIOR points to $P(A|B)$, and MARGINALIZATION points to $P(B)$.

Central Tendency in Statistics

1) Mean: The mean (or average) is that the most generally used and well-known measure of central tendency. It will be used with both discrete and continuous data, though it's most typically used with continuous data (see our styles of Variable guide for data types). The mean is adequate the sum of all the values within the data set divided by the number of values within the data set. So, if we have n values in a data set and they have values x_1, x_2, \dots, x_n , the sample mean, usually denoted by "**x bar**", is:

Population Mean Formula

$$\text{Population Mean} = \frac{\text{Sum of All the Items}}{\text{Number of Items}}$$

$$\text{Sample Mean} = \frac{\text{Sum of All the Items in Sample}}{(\text{Number of Items in Sample} - 1)}$$

2) Median: The median value of a dataset is the value in the middle of the dataset when it is arranged in ascending or descending order. When the dataset has an even number of values, the median value can be calculated by taking the mean of the middle two values.

The following image gives an example for finding the median for odd and even numbers of samples in the dataset.

1, 3, 3, **6**, 7, 8, 9

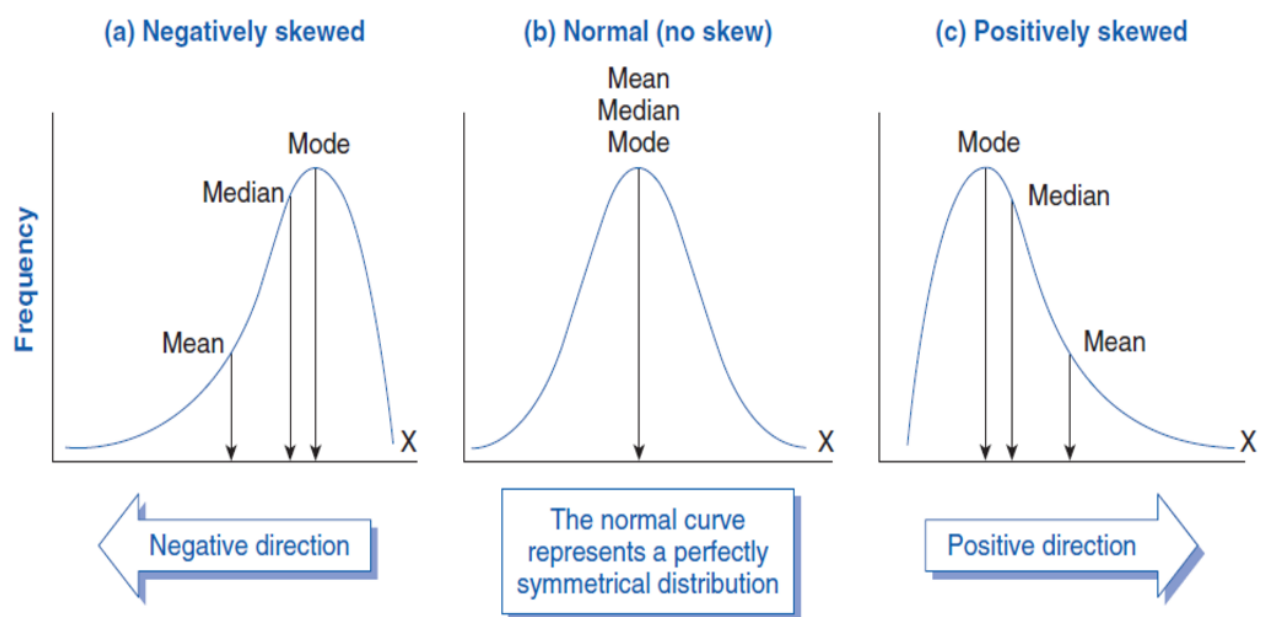
Median = **6**

1, 2, 3, **4**, **5**, 6, 8, 9

Median = $(4 + 5) \div 2$
= **4.5**

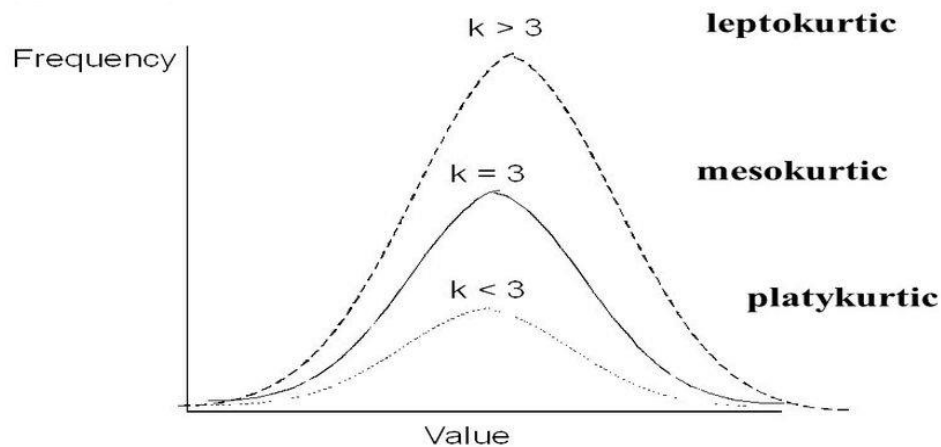
3) Mode: The mode is the value that appears the most frequently in your data set. The mode is the highest bar in a bar chart. A multimodal distribution exists when the data contains multiple values that are tied for the most frequently occurring. If no value repeats, the data does not have a mode.

4) Skewness: Skewness is a metric for symmetry, or more specifically, the lack of it. If a distribution, or data collection, looks the same to the left and right of the centre point, it is said to be symmetric.



5) Kurtosis: Kurtosis is a measure of how heavy-tailed or light-tailed the data are in comparison to a normal distribution. Data sets having a high kurtosis are more likely to contain heavy tails or outliers. Light tails or a lack of outliers are common in data sets with low kurtosis.

Kurtosis



Variability in Statistics

Range: In statistics, the range is the smallest of all dispersion measures. It is the difference between the distribution's two extreme conclusions. In other words, the range is the difference between the distribution's maximum and minimum observations.

$$\text{Range} = X_{\max} - X_{\min}$$

Where X_{\max} represents the largest observation and X_{\min} represents the smallest observation of the variable values.

Percentiles, Quartiles and Interquartile Range (IQR)

· **Percentiles** — It is a statistician's unit of measurement that indicates the value below which a given percentage of observations in a group of observations fall.

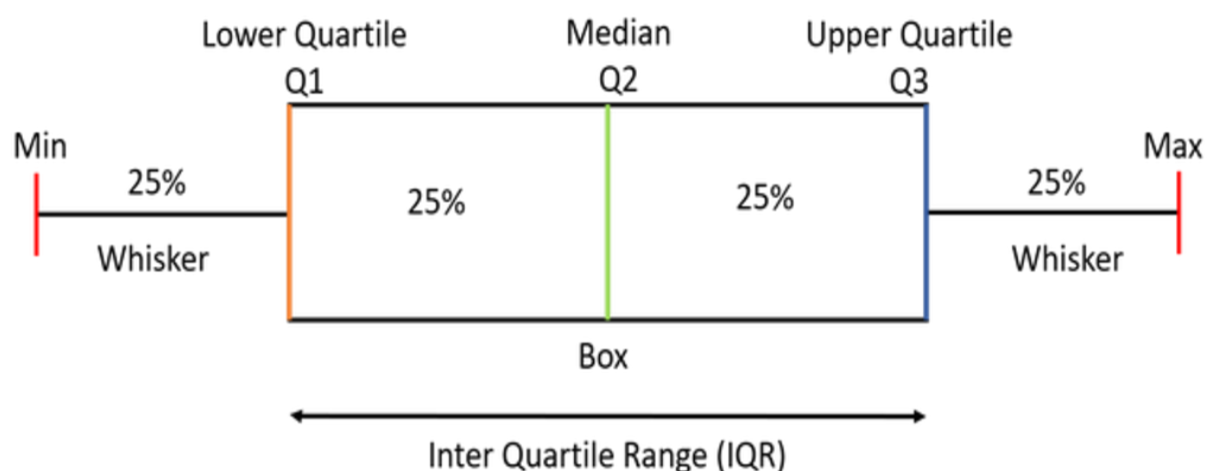
For instance, the value Q_X represents the 40th percentile of XX (0.40)

· **Quantiles**— Values that divide the number of data points into four more or less equal parts, or quarters. Quantiles are the 0th, 25th, 50th, 75th, and 100th percentile values or the 0th, 25th, 50th, 75th, and 100th percentile values.

· **Interquartile Range (IQR)**— The difference between the third and first quartiles is defined by the interquartile range. The partitioned values that divide the entire series into four equal parts are known as quartiles. So, there are three quartiles. The first quartile, known as the lower quartile, is denoted by Q1, the second quartile by Q2, and the third quartile by Q3, known as the upper quartile. As a result, the interquartile range equals the upper quartile minus the lower quartile.

IQR = Upper Quartile – Lower Quartile

= Q3 – Q1



- **Variance:** The dispersion of a data collection is measured by variance. It is defined technically as the average of squared deviations from the mean.

Population Variance	Sample Variance
$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$ <p> σ^2 = population variance x_i = value of i^{th} element μ = population mean N = population size </p>	$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$ <p> s^2 = sample variance x_i = value of i^{th} element \bar{x} = sample mean n = sample size </p>

- **Standard Deviation:** The standard deviation is a measure of data dispersion WITHIN a single sample selected from the study population. The square root of the variance is used to compute it. It simply indicates how distant the individual values in a sample are from the mean. To put it another way, how dispersed is the data from the sample? As a result, it is a sample statistic.

Standard Deviation Formula



Population	Sample
$\sigma = \sqrt{\frac{\sum (X - \mu)^2}{N}}$ <p> X - The Value in the data distribution μ - The population Mean N - Total Number of Observations </p>	$s = \sqrt{\frac{\sum (X - \bar{x})^2}{n-1}}$ <p> X - The Value in the data distribution \bar{x} - The Sample Mean n - Total Number of Observations </p>

· **Standard Error (SE):** The standard error indicates how close the mean of any given sample from that population is to the true population mean. When the standard error rises, implying that the means are more dispersed, it becomes more likely that any given mean is an inaccurate representation of the true population mean. When the sample size is increased, the standard error decreases – as the sample size approaches the true population size, the sample means cluster more and more around the true population mean.

Standard Error Formula



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



Relationship Between Variables

· **Causality:** The term “causation” refers to a relationship between two events in which one is influenced by the other. There is causality in statistics when the value of one event, or variable, grows or decreases as a result of other events.

Each of the events we just observed may be thought of as a variable, and as the number of hours worked grows, so does the amount of money earned. On the other hand, if you work fewer hours, you will earn less money.

- **Covariance:** Covariance is a measure of the relationship between two random variables in mathematics and statistics. The statistic assesses how much – and how far – the variables change in tandem. To put it another way, it's a measure of the variance between two variables. The metric, on the other hand, does not consider the interdependence of factors. Any positive or negative value can be used for the variance.

The following is how the values are interpreted:

- Positive covariance: When two variables move in the same direction, this is called positive covariance.
- Negative covariance indicates that two variables are moving in opposite directions.

The diagram illustrates the formula for covariance, $Cov(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$, with several annotations in orange:

- total count of sample values:** An arrow points to the summation symbol \sum .
- single observed value of dependent variable:** An arrow points to x_i .
- mean of all values of independent variable:** An arrow points to \bar{x} .
- single observed value of independent variable:** An arrow points to y_i .
- mean of all values of independent variable:** An arrow points to \bar{y} .
- population count minus one (Bessel's Correction):** An arrow points to the denominator $n - 1$.

• **Correlation:** Correlation is a statistical method for determining whether or not two quantitative or categorical variables are related. To put it another way, it's a measure of how things are connected. Correlation analysis is the study of how variables are connected.

Ø Here are a few examples of data with a high correlation:

- 1) Your calorie consumption and weight.
- 2) Your eye colour and the eye colours of your relatives.
- 3) The amount of time you spend studying and your grade point average

Ø Here are some examples of data with poor (or no) correlation:

- 1) Your sexual preference and the cereal you eat are two factors to consider.
- 2) The name of a dog and the type of dog biscuit that they prefer.
- 3) The expense of vehicle washes and the time it takes to get a Coke at the station.

Correlations are useful because they allow you to forecast future behaviour by determining what relationship variables exist. In the social sciences, such as government and healthcare, knowing what the future holds is critical. Budgets and company plans are also based on these facts.

$$\text{Correlation} = \frac{\text{Cov}(x, y)}{\sigma x * \sigma y}$$

Probability Distributions

-

Probability Distribution Functions

1) Probability Mass Function (PMF): The probability distribution of a discrete random variable is described by the PMF, which is a statistical term.

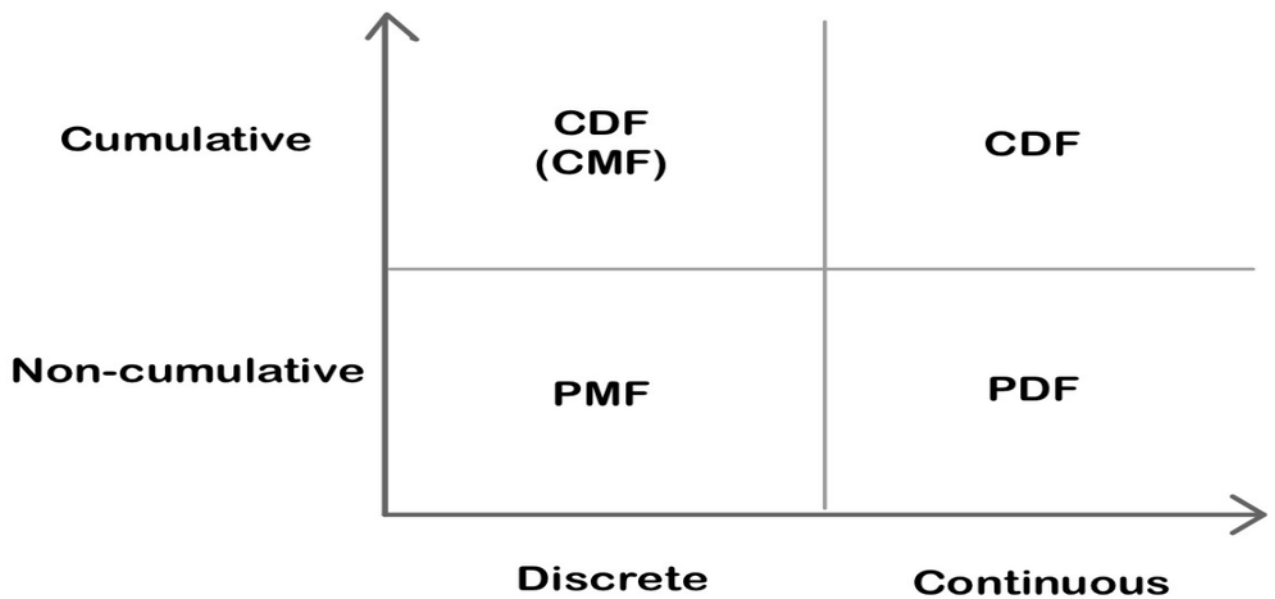
The terms PDF and PMF are frequently misunderstood. The PDF is for continuous random variables, whereas the PMF is for discrete random variables. Throwing a dice, for example (you can only choose from 1 to 6 numbers (countable))

2) Probability Density Function (PDF): The probability distribution of a continuous random variable is described by the word PDF, which is a statistical term.

The Gaussian Distribution is the most common distribution used in PDF. If the features / random variables are Gaussian distributed, then the PDF will be as well. Because the single point represents a line that does not span the area under the curve, the probability of a single outcome is always 0 on a PDF graph.

3) Cumulative Density Function (CDF): The cumulative distribution function can be used to describe the continuous or discrete distribution of random variables.

If X is the height of a person chosen at random, then $F(x)$ is the probability of the individual being shorter than x . If $F(180 \text{ cm})=0.8$, then an individual chosen at random has an 80% chance of being shorter than 180 cm (equivalently, a 20 per cent chance that they will be taller than 180cm).

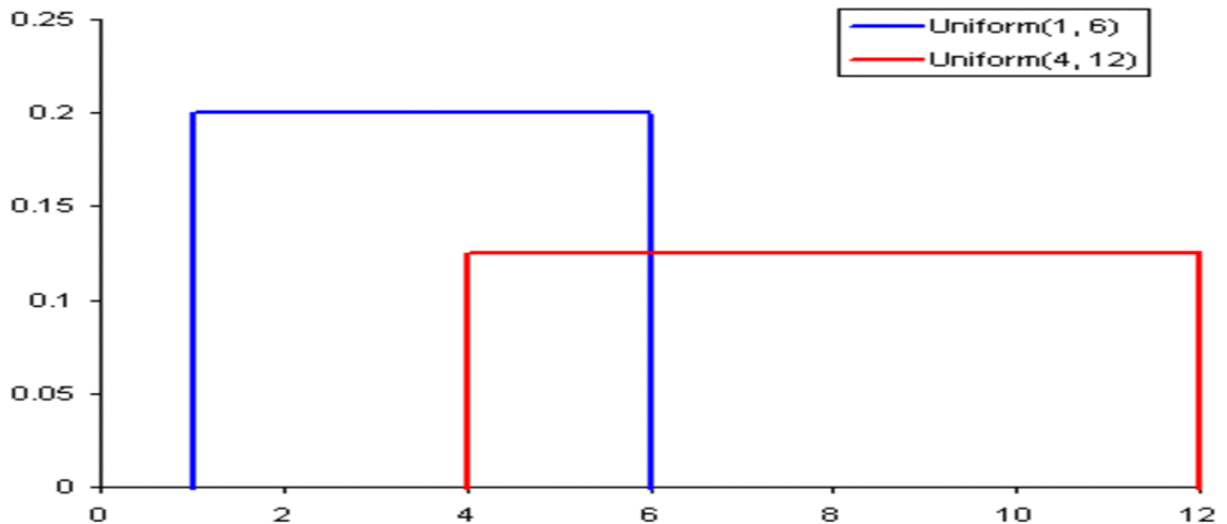


-

Continuous Probability Distribution

1) Uniform Distribution: Uniform distribution is a sort of probability distribution in statistics in which all events are equally likely. Because the chances of drawing a heart, a club, a diamond, or a spade are equal, a deck of cards contains uniform distributions. Because the likelihood of receiving heads or tails in a coin toss is the same, a coin has a uniform distribution.

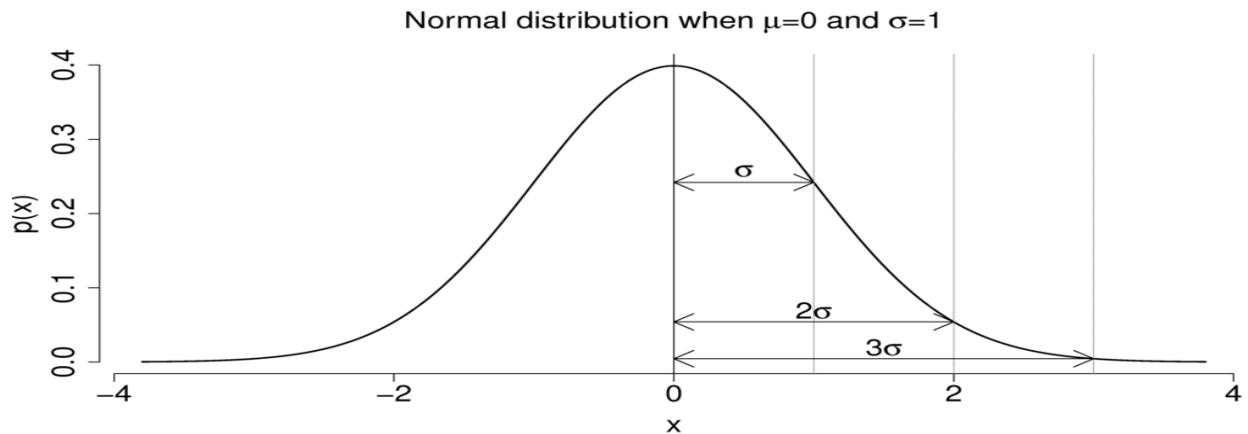
A coin flip that returns a head or tail has a probability of $p = 0.50$ and would be represented by a line from the y-axis at 0.50.



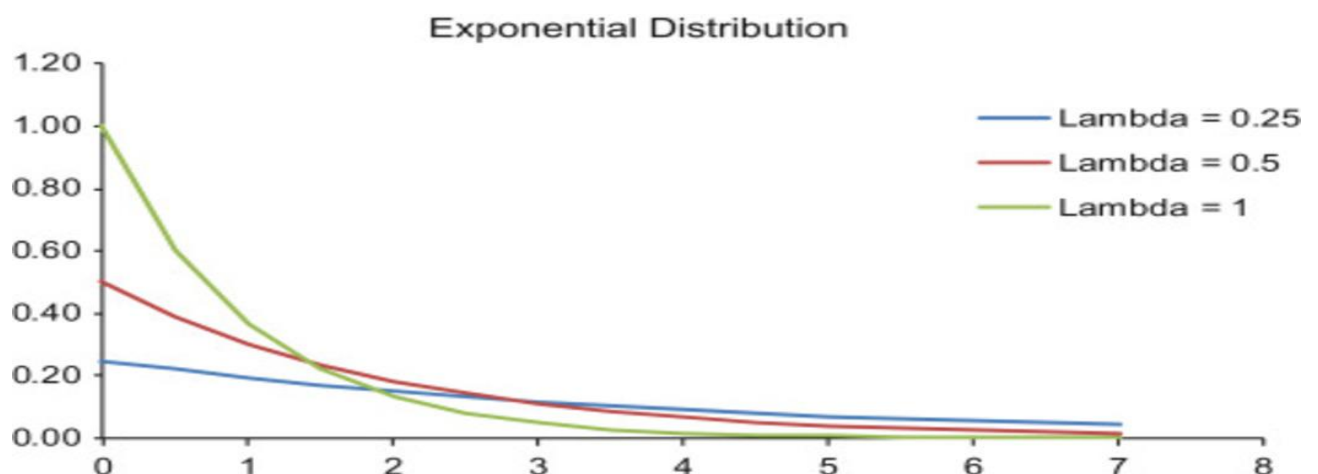
2) Normal/Gaussian Distribution: The normal distribution, also known as the Gaussian distribution, is a symmetric probability distribution centred on the mean, indicating that data around the mean occur more frequently than data far from it. The normal distribution will show as a bell curve on a graph.

Points to remember: –

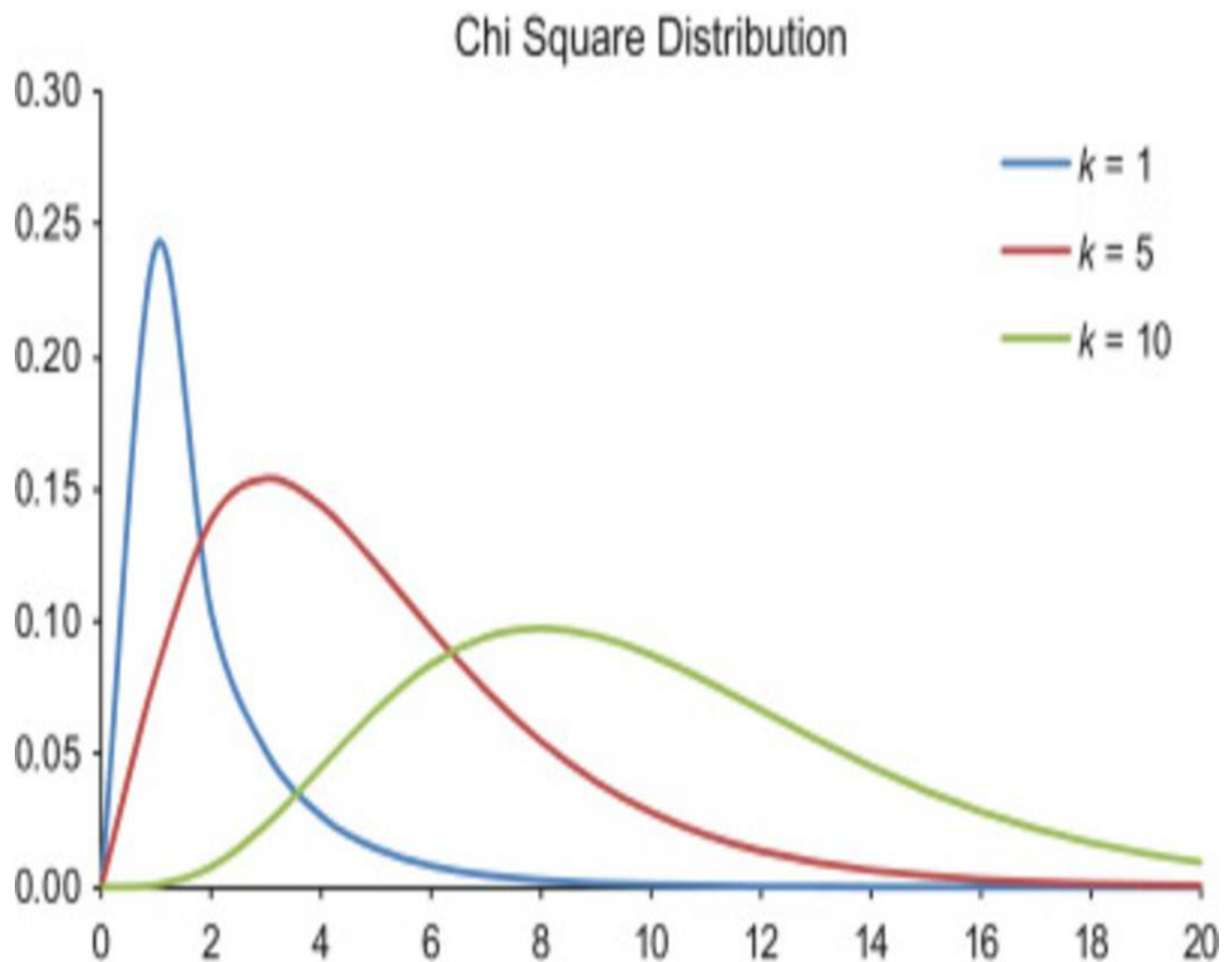
- A probability bell curve is referred to as a normal distribution.
- The mean of a normal distribution is 0 and the standard deviation is 1. It has a kurtosis of 3 and zero skew.
- Although all symmetrical distributions are normal, not all normal distributions are symmetrical.
- Most pricing distributions aren't totally typical.



3) Exponential Distribution: The exponential distribution is a continuous distribution used to estimate the time it will take for an event to occur. For example, in physics, it is frequently used to calculate radioactive decay, in engineering, it is frequently used to calculate the time required to receive a defective part on an assembly line, and in finance, it is frequently used to calculate the likelihood of a portfolio of financial assets defaulting. It can also be used to estimate the likelihood of a certain number of defaults occurring within a certain time frame.



4) Chi-Square Distribution: A continuous distribution with degrees of freedom is called a chi-square distribution. It's used to describe a sum of squared random variable's distribution. It's also used to determine whether a data distribution's goodness of fit is good, whether data series are independent, and to estimate confidence intervals around variance and standard deviation for a random variable from a normal distribution. Furthermore, the chi-square distribution is a subset of the gamma distribution.



Discrete Probability Distribution

1) Bernoulli Distribution: A Bernoulli distribution is a discrete probability distribution for a Bernoulli trial, which is a random experiment with just two

outcomes (named “Success” or “Failure” in most cases). When flipping a coin, the likelihood of getting ahead (a “success”) is 0.5. “Failure” has a chance of $1 - P$. (where p is the probability of success, which also equals 0.5 for a coin toss). For $n = 1$, it is a particular case of the binomial distribution. In other words, it’s a single-trial binomial distribution (e.g. a single coin toss).

BERNOULLI DISTRIBUTION

- A Bernoulli trial is an experiment with only two outcomes. An r.v. X has Bernoulli(p) distribution if

$$X = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1 - p \end{cases}; 0 \leq p \leq 1$$

$$P(X = x) = p^x (1 - p)^{1-x} \text{ for } x = 0, 1; \text{ and } 0 < p < 1$$

11

2) Binomial Distribution: A discrete distribution is a binomial distribution. It’s a well-known probability distribution. The model is then used to depict a variety of discrete phenomena seen in business, social science, natural science, and medical research.

Because of its relationship with a binomial distribution, the binomial distribution is commonly employed. For binomial distribution to be used, the following conditions must be met:

1. There are n identical trials in the experiment, with n being a limited number.
2. Each trial has only two possible outcomes, i.e., each trial is a Bernoulli’s trial.

3. One outcome is denoted by the letter S (for success) and the other by the letter F (for failure) (for failure).
4. From trial to trial, the chance of S remains the same. The chance of success is represented by p, and the likelihood of failure is represented by q (where $p+q=1$).
5. Each trial is conducted independently.
6. The number of successful trials in n trials is the binomial random variable x.

If X reflects the number of successful trials in n trials under the preceding conditions, then x is said to follow a binomial distribution with parameters n and p.

Binomial Distribution Formula

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

where

n = the number of trials (or the number being sampled)

x = the number of successes desired

p = probability of getting a success in one trial

$q = 1 - p$ = the probability of getting a failure in one trial

3) Poisson Distribution: A Poisson distribution is a probability distribution used in statistics to show how many times an event is expected to happen over a certain amount of time. To put it another way, it's a count distribution. Poisson distributions are frequently accustomed comprehend independent events that occur at a gradual rate during a selected timeframe.

The Poisson distribution is a discrete function, which means the variable can only take values from a (possibly endless) list of possibilities. To put it another way, the variable can't take all of the possible values in any continuous range. The variable can only take the values 0, 1, 2, 3, etc., with no fractions or decimals, in the Poisson distribution (a discrete distribution).

Poisson Distribution Formula

$$P(X = x) = \frac{\lambda^x e^{-\lambda}}{x!}$$

where

$x = 0, 1, 2, 3, \dots$

λ = mean number of occurrences in the interval

e = Euler's constant ≈ 2.71828

Hypothesis Testing and Statistical Significance in Statistics

Hypothesis testing may be a method within which an analyst verifies a hypothesis a couple of population parameters. The analyst's approach is set by

the kind of the info and also the purpose of the study. the utilization of sample data to assess the plausibility of a hypothesis is thought of as hypothesis testing.

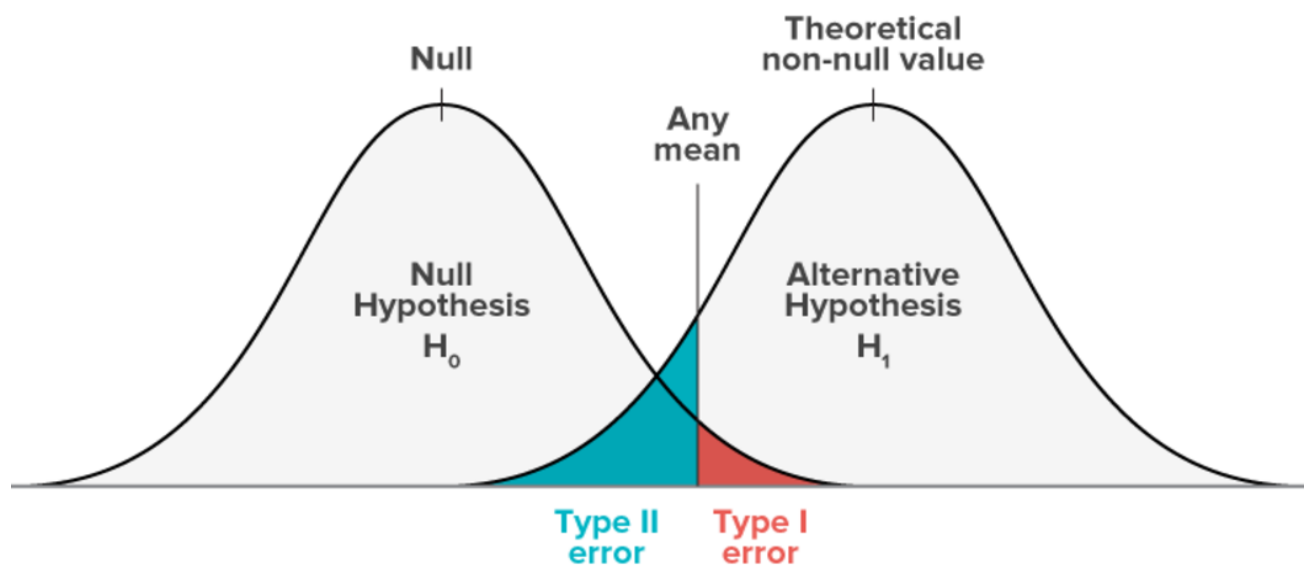
Null and Alternative Hypothesis

Null Hypothesis (H_0)

A population parameter (such as the mean, standard deviation, and so on) is equal to a hypothesised value, according to the null hypothesis. The null hypothesis is a claim that is frequently made based on previous research or specialised expertise.

Alternative hypothesis (H_1)

The alternative hypothesis says that a population parameter is less, more, or different than the null hypothesis's hypothesised value. The alternative hypothesis is what you believe or want to prove to be correct.



Type 1 and Type 2 error

Type 1 error:

A type 1 error, often referred to as a false positive, happens when a researcher rejects a real null hypothesis incorrectly. This suggests you're claiming your findings are noteworthy after they actually happened by coincidence.

Your alpha level (α), which is that the p-value below which you reject the null hypothesis, represents the likelihood of constructing a sort I error. When rejecting the null hypothesis, a p-value of 0.05 suggests that you simply are willing to tolerate a 5% probability of being mistaken.

By setting p to a lesser value, you'll lessen your chances of constructing a kind I error.

Type 2 error:

A type II error commonly said as a false negative happens when a researcher fails to reject a null hypothesis that's actually true. During this case, a researcher finds that there's no significant influence when, in fact, there is.

Beta (β) is that the probability of creating a sort II error, and it's proportional to the statistical test's power ($\text{power} = 1 - \beta$). By ensuring that your test has enough power, you'll reduce your chances of constructing a sort II error.

This can be accomplished by ensuring that your sample size is sufficient to spot a practical difference when one exists.

	Reject H_0	Fail to Reject H_0
H_0 Is True	Type I Error α (FP)	Correct $1 - \alpha$ (TN)
H_0 Is False	Correct $1 - \beta$ ("Statistic Power") (TP)	Type II Error β (FN)

-

Interpretation

P-value: The p-value in statistics is that the likelihood of getting outcomes a minimum of as extreme because the observed results of a statistical hypothesis test, given the null hypothesis is valid. The p-value, instead of rejection points, is employed to work out the smallest amount level of significance at which the null hypothesis is rejected. A lower p-value indicates that the choice hypothesis has more evidence supporting it.

Critical Value: it is a point on the test distribution that is compared to the test statistic to see if the null hypothesis should be rejected. You can declare statistical significance and reject the null hypothesis if the absolute value of your test statistic is larger than the crucial value.

Significance Level and Rejection Region: The probability that an event (such as a statistical test) occurred by chance is the significance level of the occurrence. We call an occurrence significant if the level is very low, i.e., the possibility of it happening by chance is very minimal. The rejection region

depends on the importance level. the importance level is denoted by α and is that the probability of rejecting the null hypothesis if it's true.

Z-Test: The z-test may be a hypothesis test within which the z-statistic is distributed normally. The z-test is best utilized for samples with quite 30 because, in line with the central limit theorem, samples with over 30 samples are assumed to be approximately regularly distributed.

The null and alternative hypotheses, also because the alpha and z-score, should all be reported when doing a z-test. The test statistic should next be calculated, followed by the results and conclusion. A z-statistic, also called a z-score, could be a number that indicates what number of standard deviations a score produced from a z-test is above or below the mean population.

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

\bar{x} = sample mean

μ = population mean

σ = population standard deviation

n = sample size

T-Test: A t-test is an inferential statistic that's won't see if there's a major difference within the means of two groups that are related in how. It's most ordinarily employed when data sets, like those obtained by flipping a coin 100 times, are expected to follow a traditional distribution and have unknown

variances. A t-test could be a hypothesis-testing technique that will be accustomed to assess an assumption that's applicable to a population.

t-Test Formula



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



$$t = \frac{(\bar{X}_1 - \bar{X}_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

ANOVA (Analysis of Variance): ANOVA is the way to find out if experimental results are significant. **One-way ANOVA** compares two means from two independent groups using only one independent variable. **Two-way ANOVA** is the extension of one-way ANOVA using two independent variables to calculate the main effect and interaction effect.

Chi-Square Test: it is a test that assesses how well a model matches actual data. A chi-square statistic requires data that is random, raw, mutually exclusive, collected from independent variables, and drawn from a large enough sample. The outcomes of a fair coin flip, for example, meet these conditions.

In hypothesis testing, chi-square tests are frequently utilised. Given the size of the sample and the number of variables in the relationship, the chi-square

statistic examines the size of any disparities between the expected and actual results.

$$\chi^2 = \sum \frac{(\text{Observed value} - \text{Expected value})^2}{\text{Expected value}}$$

Image Source:

Image 1 –

https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcSFIQDc_CIXRtqjzRHbCEpiaoEiCiQKIm5M7g&usqp=CAU

Image 2 – <https://pbs.twimg.com/media/EKSH8PIXUAAfR2w.jpg>

Image 3 – <https://www.altexsoft.com/media/2020/12/word-image-6.png>

Image 4 – <https://i.pinimg.com/originals/3c/3f/94/3c3f9482e1ea366cdcf3002cb6bc3620.png>

Image 5 –

<https://www.sv-europe.com/wp-content/uploads/predictive-analytics-methods-infographics-1024x724.png>

Image 6 – https://upload.wikimedia.org/wikipedia/commons/c/c7/Three_Phases_of_Analytics.png

Image 7 – <https://magoosh.com/statistics/files/2018/01/Screen-Shot-2018-01-12-at-1.24.59-PM.png>

Image 8 – <https://cetking.com/wp-content/uploads/2019/03/IntersectionAB.jpg>

Image 9 – <https://cdn1.byjus.com/wp-content/uploads/2018/10/n..png>

Image 10 – https://miro.medium.com/max/424/1*ubXttqa7n5A-ILtZKwO1zA.jpeg

Image 11 –

<https://slideplayer.com/slide/8995780/27/images/5/Test+for+Independent+Events.jpg>

Image 12 –

<https://andymath.com/wp-content/uploads/2019/06/Mutually-Exclusive-1024×676.jpg>

Image 13 – https://miro.medium.com/max/1400/1*CnoTGGO7XeUpUMeXDrIfvA.png

Image 14 –

<https://cdn.educba.com/academy/wp-content/uploads/2019/04/Population-Mean-Formula.jpg>

Image 15 –

https://upload.wikimedia.org/wikipedia/commons/thumb/c/cf/Finding_the_median.png/1200px-Finding_the_median.png

Image 16 –

https://www.biologyforlife.com/uploads/2/2/3/9/22392738/c101b0da6ea1a0dab31f80d9963b0368_orig.png

Image 17 – https://images.slideplayer.com/12/3490002/slides/slide_12.jpg

Image 18 –

<https://media.geeksforgeeks.org/wp-content/uploads/20201127012952/boxplot-660×233.png>

Image 19 – <https://www.onlinemathlearning.com/image-files/population-variance.png>

Image 20 –

<https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcQBCuWgGCiwOxxXRiSDdEj-IT5QxltsvQeH8g&usqp=CAU>

Image 21 –

<https://cdn.wallstreetmojo.com/wp-content/uploads/2020/01/Standard-Error-Formula.jpg>

Image 22 –

<https://www.alpharithms.com/wp-content/uploads/1156/covariance-equation-diagram.jpg>

Image 23 –

<https://zerodha.com/varsity/wp-content/uploads/2017/04/Correlation-Formula.png>

Image 24 – https://miro.medium.com/max/1400/1*ktIttLCFRAqdUILE180v9g.png

Image 25 – <https://www.vosesoftware.com/riskwiki/images/imagec247.gif>

Image 26 – https://learnche.org/pid/_images/normal-distribution-standardized.png

Image 27 –

https://www.statisticshowto.com/wp-content/uploads/2014/06/360px-Exponential_pdf.svg_.png

Image 28 –

<https://ars.els-cdn.com/content/image/3-s2.0-B9780128051634000049-u04-08-9780128051634.jpg>

Image 29 –

https://slidetodoc.com/presentation_image_h/dab38137c14122e391230d25b9bef01a/image-11.jpg

Image 30 – https://miro.medium.com/max/518/0*nY0id5Q-Vj_EgerW.png

Image 31 – <https://www.onlinemathlearning.com/image-files/poisson-distribution-formula.png>

Image 32 – https://miro.medium.com/max/740/1*BVoAV9v4RMi6vAcinhvJKA.png

Image 33 – https://miro.medium.com/max/1400/1*uF5aBZ63BZ8IDEyF8wo-CA.png

Image 34 –

<https://www.researchgate.net/publication/297600508/figure/fig4/AS:670388308680713@1536844436654/Formula-for-the-Z-test.ppm>

Image 35 – <https://cdn.educba.com/academy/wp-content/uploads/2019/12/t-Test-Formula.jpg>

Image 36 – <https://cdn1.byjus.com/wp-content/uploads/2020/10/Chi-Square-Test.png>