

Statistics And Probability For Data Scientist And Data Analyst

Chapter 3: Descriptive Statistics: -

A descriptive statistic is a summary statistic that quantitatively describes or summarizes features from a collection of information from sample or entire population.

Characteristics of Data:

1. **Centre:** A representative or average value that indicates where the middle of the data set is located.
2. **Variation:** A measure of the amount that the data values vary.
3. **Distribution:** The nature or shape of the spread of the data over the range of values (Such as bell-shaped, uniform, or skewed).
4. **Outliers:** An observation of data that does not fit the rest of the data. Extreme high or extreme low.

❖ Measure Of Central Tendency: -

A measure of central tendency is a single value that attempts to describe a set of data by identifying the central position within that set of data.

In statistics, the mean, median, and mode are the three most common measures of central tendency.

Mean: - The **mean** (or average) of a set of data values is the sum of all of the data values divided by the number of data values.

$$\text{Mean} = \frac{\sum x}{n}$$

$\sum x$ = denotes the sum of a set of data values

x = is the variable usually used to represent the individual data values.

n = represents the number of data values in a sample.

N = represents the number of data values in a population.

$$\bar{x} = \frac{\sum x}{n} \quad \text{is the mean of a set of } \textit{sample} \text{ values.}$$

$$\mu = \frac{\sum x}{N} \quad \text{is the mean of all values in a } \textit{population}.$$

Example: - Mount Rival hosts a soccer tournament each year. This season, in 10 games, the lead scorer for the home team scored 7, 5, 0, 7, 8, 5, 5, 4, 1 and 5 goals. What is the mean score of this player?

The sum of all values is 47 and there are 10 values. Therefore, the mean is $47 \div 10 = 4.7$ goals per game.

Statistics And Probability For Data Scientist And Data Analyst

Median: - The median is the middle value of the ordered data. The median is often denoted by \tilde{x} (pronounced “x-tilde”).

The most important step in finding the median is to first order the data from smallest to largest.

- 1] Arrange the data in increasing order, i.e. smallest to largest.
- 2] Find the location of the median in the ordered data by $\frac{n+1}{2}$, where n is the sample size.

Example: - First, order the data. 69, 76, 76, 78, 80, 82, 86, 88, 91, 95

With $n = 10$, the median position is found by $(10 + 1) / 2 = 5.5$. Thus, the median is the average of the fifth (80) and sixth (82) ordered value and the median = 81.

Note on Odd or Even Sample Sizes

If the sample size is an odd number then the location point will produce a median that is an observed value. If the sample size is an even number, then the location will require one to take the mean of two numbers to calculate the median. The result may or may not be an observed value as the example below illustrates.

Mode: - The mode is the value that occurs most often in the data. It is important to note that there may be more than one mode in the dataset.

Example: - First, order the data. 69, 76, 76, 78, 80, 82, 86, 88, 91, 95.

The most frequent value in this data set is 76. Therefore, the mode is 76.

Mean From a Frequency Distribution: -

mean from frequency distribution:

First multiply each frequency and class midpoint, then add the products.

$$\bar{x} = \frac{\sum(f \cdot x)}{\sum f}$$

sum of frequencies

Statistics And Probability For Data Scientist And Data Analyst

Table 3-1 Finding the Mean from a Frequency Distribution

Word Counts from Men	Frequency f	Class Midpoint x	$f \cdot x$
0-9,999	46	4,999.5	229,977.0
10,000-19,999	90	14,999.5	1,349,955.0
20,000-29,999	40	24,999.5	999,980.0
30,000-39,999	7	34,999.5	244,996.5
40,000-49,999	3	44,999.5	134,998.5
Totals:	$\Sigma f = 186$		$\Sigma(f \cdot x) = 2,959,907$
			$\bar{x} = \frac{\Sigma(f \cdot x)}{\Sigma f} = \frac{2,959,907}{186} = 15,913.5$

❖ **Measures of Dispersion:** - How the Data is Spread around a centre of distribution.

1]Range: - The range is the difference in the maximum and minimum values of a data set. The range is easy to calculate but it is very much affected by extreme values.

$$\text{Range} = \text{maximum} - \text{minimum}$$

2]Standard Deviation: - A standard deviation (**or** σ) is a measure of how dispersed the data is in relation to the mean. Low standard deviation means data are clustered around the mean, and high standard deviation indicates data are more spread out.

Or

Measures the average distance your data values are from the mean.

1. Never negative and Never Zero unless all entries are the same.
2. Greatly affected by outliers.

Sample standard deviation is denoted by 's'.

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

Statistics And Probability For Data Scientist And Data Analyst

Problem: India has 1 satellite used for military and intelligence purposes, Japan has 3, and Russia has 14. Find the Standard Deviation of the sample values of 1, 3, and 14.

X	$x - \bar{x}$	$(x - \bar{x})^2$
1	1 - 6 = -5	25
3	3 - 6 = -3	9
14	14 - 6 = 8	64
N = 3		$\sum(x - \bar{x})^2 = 98$

Mean (\bar{x}) = $1 + 3 + 14 / 3 = 6$

Other Formula

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

$$s = \sqrt{\frac{98}{3 - 1}} = \sqrt{\frac{98}{2}} = \sqrt{49} = 7$$

x	x^2
1	1
3	9
14	196
N = 3 $\sum x = 18$	$\sum x^2 = 206$

$$s = \sqrt{\frac{n \sum (x^2) - (\sum x)^2}{n(n - 1)}} = \sqrt{\frac{3 * 206 - (18)^2}{3(3 - 1)}} = \sqrt{\frac{618 - 324}{6}} = \sqrt{\frac{324}{6}} = \sqrt{49} = 7$$

SOLUTION

Shown below is the computation of the standard deviation of 1 satellite, 3 satellites, and 14 satellites using Formula 3-5.

$n = 3$ (because there are 3 values in the sample)

$\sum x = 18$ (found by adding the sample values: $1 + 3 + 14 = 18$)

$\sum x^2 = 206$ (found by adding the squares of the sample values, as in $1^2 + 3^2 + 14^2 = 206$)

Using Formula 3-5, we get

$$s = \sqrt{\frac{n(\sum x^2) - (\sum x)^2}{n(n - 1)}} = \sqrt{\frac{3(206) - (18)^2}{3(3 - 1)}} = \sqrt{\frac{294}{6}} = 7.0 \text{ satellites}$$

Statistics And Probability For Data Scientist And Data Analyst

Standard Deviation of a Population: -

population standard deviation $\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}}$

3]Variance: -The variance of a set of values is a measure of variation equal to the square of the standard deviation.

Population variance

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

where μ is the population mean and the summation is over all possible values of the population and N is the population size.

σ^2 is often estimated by using the sample variance.

Sample Variance

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

Where n is the sample size and \bar{x} is the sample mean.

Coefficient of Variance: -When comparing variation in two different sets of data, the standard deviations should be compared only if the two sets of data use the same scale and units and they have approximately the same mean. If the means are substantially different, or if the samples use different scales or measurement units, we can use the coefficient of variation, defined as follows.

The coefficient of variation (or CV) for a set of the nonnegative sample or population data, expressed as a per cent describes the standard deviation relative to the mean, and is given by the following:

Sample

$$CV = \frac{s}{\bar{x}} \cdot 100\%$$

Population

$$CV = \frac{\sigma}{\mu} \cdot 100\%$$

Statistics And Probability For Data Scientist And Data Analyst

Problem: - Heights and Weights of Men: Compare the variation in heights of men to the variation in weights of men, using these sample results obtained from Data Set 1 in Appendix B: for men, the heights yield $\bar{x} = 68.34$ in. and $s = 3.02$ in; the weights yield $\bar{x} = 172.55$ lb and $s = 26.33$ lb. Note that we want to compare variation among heights to variation among weights.

Solution: - We can compare the standard deviations if the same scales and units are used and the two means are approximately equal, but here we have different scales (heights and weights) and different units of measurement (inches and pounds), so we use the coefficients of variation:

$$\text{heights: } CV = \frac{s}{\bar{x}} \cdot 100\% = \frac{3.02 \text{ in.}}{68.34 \text{ in.}} \cdot 100\% = 4.42\%$$

$$\text{weights: } CV = \frac{s}{\bar{x}} \cdot 100\% = \frac{26.33 \text{ lb}}{172.55 \text{ lb}} \cdot 100\% = 15.26\%$$

Although the standard deviation of 3.02 in. cannot be compared to the standard deviation of 26.33 lb, we can compare the coefficients of variation, which have no units. We can see that heights (with $CV = 4.42\%$) have considerably less variation than weights (with $CV = 15.26\%$).

❖ **Probability:** - The probability of an event is a number that indicates how likelihood the event is to occur.

$P = \text{Number of Favorable Outcomes} / \text{Total Number of Outcomes.}$

Notation for Probabilities:

- P denotes a probability.
- A, B, and C denote specific events.
- $P(A)$ denotes the probability of event A occurring.

Vocabulary –

- **Event:** - An event is any collection of results or outcomes of a procedure.
- **Simple Event:** - A single outcome (A simple event is an outcome or an event that cannot be further broken down into simpler components.)
- **Sample Space:** - The sample space for a procedure consists of all possible simple events. That is, the sample space consists of all outcomes that cannot be broken down any further.

Example: -

Procedure	Event	Sample space
Flipping a coin Flip one time	Head	Head/Tail {H, T}

Statistics And Probability For Data Scientist And Data Analyst

Procedure	Event	Sample space
Flipping a coin Flip a coin 3 times	1 Head 2 Tails	<p>{H H H / H H T / H T T / *H T H / T T T / *T T H / *T H T / T H H}</p> <p>* = 3 Ways we can achieve the event.</p> <p>All single outcomes are separately called a simple event. If we put everything in flower brackets it's called sample space.</p>

Procedure	Event	Sample space
Flipping a coin Flip a coin 3 times	1 Tail 2 Heads 3 Heads 3 Tails	<p>{**H H H / *H H T / H T T / *H T H / ++T T T / T T H / T H T / *T H H}</p> <p>* = First event can be achieved 3 times</p> <p>** = Second event can be achieved 1 time</p> <p>++ = Third event can be achieved one time.</p> <p>All single outcomes are separately called a simple event. If we put everything in flower brackets it's called sample space.</p>

Interpretations of Probability: -

- **Classical Interpretation of Probability:** - The probability that event E occurs is denoted by P(E). When all outcomes are equally likely, then:

$$P(E) = \frac{\text{number of outcomes in } E}{\text{number of possible outcomes}}$$

- **Subjective Probability:** - **Subjective probability** reflects personal belief which involves personal judgment, information, intuition, etc.
- **Relative Frequency concept of Probability (Empirical Approach):** - If a particular outcome happens over a large number of events, then the percentage of that outcome is close to the true probability.

$$P(E) \approx \frac{\text{number of outcomes in } E}{\text{number of possible outcomes}}$$

Statistics And Probability For Data Scientist And Data Analyst

Properties of Probability: -

- The probability of an impossible event is 0.
- The probability of an event that is certain to occur is 1.
- For any event A, the probability of A is between 0 and 1 inclusive.
That is, $0 \leq P(A) \leq 1$
- The more of a procedure is repeated, the closer the observed probability will get to classical probability.
- The probability of an event plus the probability of the complement must be equal to 1. $P(A) + P(A') = 1$
- If A and B are mutually exclusive, then $A \cap B = \emptyset$. Therefore, $P(A \cap B) = 0$. This is important when we consider mutually exclusive (or disjoint) events.

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

Rule of Multiplication: - The rule of multiplication applies to the situation when we want to know the probability of the intersection of two events; that is, we want to know the probability that two events (Event A and Event B) both occur.

Rule of Multiplication The probability that Events A and B both occur is equal to the probability that Event A occurs times the probability that Event B occurs, given that A has occurred.

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

Rule of Addition: - The probability that Event A or Event B occurs is equal to the probability that Event A occurs plus the probability that Event B occurs minus the probability that both Events A and B occur.

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

Note: Invoking the fact that $P(A \cap B) = P(A)P(B|A)$, the Addition Rule can also be expressed as:

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

Conditional Probability: - The Probability of an event occurring given that some other event has already occurred.

$P(B|A)$ represents the probability of event B occurring after it is assumed that event A has already occurred. (We can read $(B|A)$ as “B given A” or as “event B occurring after event A has already occurred.”).

Independent Events: The occurrence of one event does not depend on the occurrence of another event or subsequent event. (Non-Independent events are Dependent).

If A & B are independent $P(B|A) = P(B)$

Dependent Events: The occurrence of one event depends on the occurrence of another event or subsequent event. (Dependent events are dependent).

Computing Conditional Probability

The Probability of A given B:

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

The Probability of B given A:

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

Statistics And Probability For Data Scientist And Data Analyst

The Probability of the Intersection of Dependent Events

The probability of dependent events A and B derived from the formulas for conditional probability:

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

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

Note ! usually $P(A|B) \neq P(B|A)$

Permutation: - Permutation relates to the act of arranging all the members of a set into some sequence or order.

A permutation is the choice of r things from a set of n things without replacement and where the order matters.

$${}^n P_r = (n!) / (n-r)!$$

Eg: - A bet on an exact in a race is won by correctly selecting the horses that finish first and second, and you must select those two horses in the correct order. The 132nd running of the Kentucky Derby had a field of 20 horses. If a better randomly selects two of those horses for an exact bet, what is the probability of winning?

We have n = 20 horses available, and we must select r = 2 of them without replacement. The number of different sequences of arrangements is found as shown:

$${}_n P_r = \frac{n!}{(n-r)!} = \frac{20!}{(20-2)!} = 380$$

Combination: - Combination is a way of selecting items from a collection, such that (unlike permutations) the order of selection does not matter.

A combination is the choice of r things from a set of n things without replacement and where order does not matter.

$${}_n C_r = \frac{n!}{(n-r)! r!}$$

Baye's Theorem: - Bayes' Theorem states that the conditional probability of an event, based on the occurrence of another event, is equal to the likelihood of the second event given the first event multiplied by the probability of the first event.

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \text{ where } P(B) \neq 0$$

Statistics And Probability For Data Scientist And Data Analyst

❖ Probability Distribution: -

Random Variable: - A **random variable** is a variable that takes on different values determined by chance.

Probability Distribution: - A probability distribution is a description that gives the probability for each value of the random variable. It is often expressed in the format of a graph, table, or formula.

Eg: Probability distribution for rolling a die.

X	P(X)
1	1/6
2	1/6
3	1/6
4	1/6
5	1/6
6	1/6

Discrete Random Variable: - When the random variable can assume only a countable, sometimes infinite, number of values.

Eg: Number of people in the class, Number of eggs.

Continuous Random Variable: - When the random variable can assume an uncountable number of values in a line interval.

Eg: Height, Weight, Temperature.

Probability Function: - A probability function is a mathematical function that provides probabilities for the possible outcomes of the random variable, **X**. It is typically denoted as $f(x)$.

There are two classes of probability functions: Probability Mass Functions and Probability Density Functions.

- **Probability Mass Function (PMF):** - If the random variable is a **discrete random variable**, the probability function is usually called the **probability mass function (PMF)**. If **X** is discrete, then $f(x) = P(X = x)$. In other words, the PMF for a constant, x , is the probability that the random variable **X** is equal to x . The PMF can be in the form of an equation or it can be in the form of a table.

Properties of probability mass functions:

1. $f(x) > 0$, for x in the sample space and 0 otherwise.
2. $\sum_x f(x) = 1$. In other words, the sum of all the probabilities of all the possible outcomes of an experiment is equal to 1.

Statistics And Probability For Data Scientist And Data Analyst

- **Probability Density Function (PDF):** - If the random variable is a **continuous random variable**, the probability function is usually called the **probability density function (PDF)**. Contrary to the discrete case, $f(x) \neq P(X = x)$.

Properties of a probability density function:

1. $f(x) > 0$, for x in the sample space and 0 otherwise.
2. The area under the curve is equal to 1.

- The probability of a random variable being less than or equal to a given value is calculated using another probability function called the **cumulative distribution function**.
- **Cumulative Distribution Function (CDF):** - A **cumulative distribution function (CDF)**, usually denoted $F(x)$, is a function that gives the probability that the random variable, X , is less than or equal to the value x .

$$F(x) = P(X \leq x)$$

Note! The definition of the cumulative distribution function is the same for a discrete random variable or a continuous random variable. For a continuous random variable, however, $P(X = x) = 0$. Therefore, the CDF, $F(x) = P(X \leq x) = P(X < x)$, for the continuous case.

Expected mean of a Discrete Random Variable: - For a discrete random variable, the expected value, usually denoted as μ or $E(X)$, is calculated using:

$$\mu = E(X) = \sum x_i f(x_i)$$

The formula means that we multiply each value, x , in the support by its respective probability, $f(x)$, and then add them all together. It can be seen as an average value but weighted by the likelihood of the value.

Variance of Discrete Random Variable: -

$$\sigma^2 = \text{Var}(X) = \sum (x_i - \mu)^2 f(x_i)$$

The formula means that we take each value of x , subtract the expected value, square that value and multiply that value by its probability. Then sum all of those values.

There is an easier form of this formula we can use.

$$\sigma^2 = \text{Var}(X) = \sum x_i^2 f(x_i) - E(X)^2 = \sum x_i^2 f(x_i) - \mu^2$$

The formula means that first, we sum the square of each value times its probability then subtract the square of the mean. We will use this form of the formula in all of our examples.

Statistics And Probability For Data Scientist And Data Analyst

Standard Deviation of a Discrete Random Variable: - The standard deviation of a random variable, X , is the square root of the variance.

$$\sigma = SD(X) = \sqrt{\text{Var}(X)} = \sqrt{\sigma^2}$$

Binominal Distribution: - The binomial distribution is a special discrete distribution where there are two distinct complementary outcomes, a “success” and a “failure”. We have a binomial experiment if ALL of the following four conditions are satisfied:

1. The experiment consists of n identical trials.
2. Each trial results in one of the two outcomes, called success and failure.
3. The probability of success, denoted p , remains the same from trial to trial.
4. The n trials are independent. That is, the outcome of any trial does not affect the outcome of the others.

If the four conditions are satisfied, then the random variable X =number of successes in n trials, is a **binomial random variable** with

$$\begin{aligned}\mu &= E(X) = np && \text{(Mean)} \\ \text{Var}(X) &= np(1-p) && \text{(Variance)} \\ SD(X) &= \sqrt{np(1-p)}, \text{ where } p \text{ is the probability of the “success.”} && \text{(Standard Deviation)}\end{aligned}$$

Notation for Binomial Probability Distributions

S and F (success and failure) denote the two possible categories of all outcomes.

$$P(S) = p \quad (p = \text{probability of a success})$$

$$P(F) = 1 - p = q \quad (q = \text{probability of a failure})$$

n denotes the fixed number of trials.

x denotes a specific number of successes in n trials, so x can be any whole number between 0 and n , inclusive.

p denotes the probability of *success* in *one* of the n trials.

q denotes the probability of *failure* in *one* of the n trials.

$P(x)$ denotes the probability of getting exactly x successes among the n trials.

The Binomial Formula

For a binomial random variable with probability of success, p , and n trials...

$$f(x) = P(X = x) = \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x} \text{ for } x = 0, 1, 2, \dots, n$$

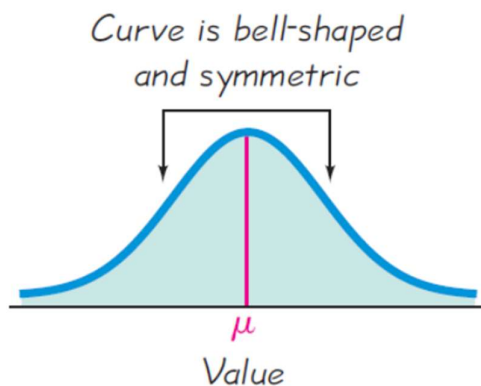
Statistics And Probability For Data Scientist And Data Analyst

Poisson Distribution: - A Poisson distribution is defined as a discrete frequency distribution that gives the probability of the number of independent events that occur in the fixed time.

- The number of trials “n” tends to infinity
- Probability of success “p” tends to zero
- $np = 1$ is finite.

Continuous Probability Distribution: -

If a continuous random variable has a distribution with a graph that is symmetric and bell-shaped and it can be described by the equation given as given below, we say that it has a normal distribution.



$$y = \frac{e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}}{\sigma\sqrt{2\pi}}$$

Expected mean of a Continuous Random variable: - The expected value (or mean) of a continuous random variable is denoted by $\mu = E(Y)$.

Variance of Continuous Random Variable: - The variance of a continuous random variable is denoted by $\sigma^2 = \text{Var}(Y)$.

Standard Deviation of Continuous Random Variable: - The standard deviation of a continuous random variable is denoted by $\sigma = \sqrt{\text{Var}(Y)}$.

Normal Distribution: - The **Normal Distribution**, also called the **Gaussian Distribution**, The Normal Distribution is a family of continuous distributions that can model many histograms of real-life data which are mound-shaped (bell-shaped) and symmetric.

Mean μ (center of the curve)

standard deviation σ (spread about the center) (..and variance σ^2)

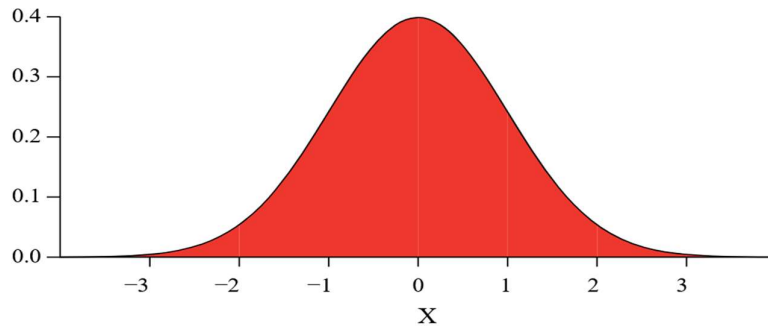
The range can also extend to $-\infty$ to $+\infty$ and still we can find a smooth curve.

The mean can be any real number and the standard deviation is greater than zero.

Statistics And Probability For Data Scientist And Data Analyst

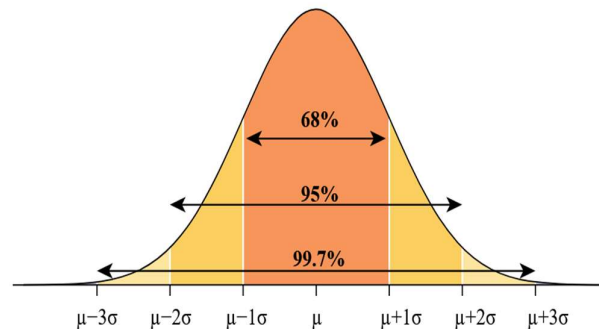
Standard Normal Distribution: -A standard normal distribution has a mean of 0 and variance of 1. This is also known as a **z distribution**. You may see the notation $N(\mu, \sigma^2)$ where N signifies that the distribution is normal, μ is the mean, and σ^2 is the variance. A Z distribution may be described as $N(0, 1)$. Note that since the standard deviation is the square root of the variance then the standard deviation of the standard normal distribution is 1.

Standard Normal Distribution, $N(0,1)$



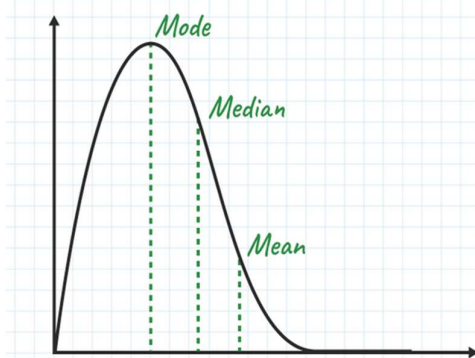
Empirical Rule: -

- 68% of the observations lie within one standard deviation to either side of the mean.
- 95% of the observations lie within two standard deviations to either side of the mean.
- 99.7% of the observations lie within three standard deviations to either side of the mean.



Skewness: - Skewness is a measure of the asymmetry of a distribution.

Positive Skew



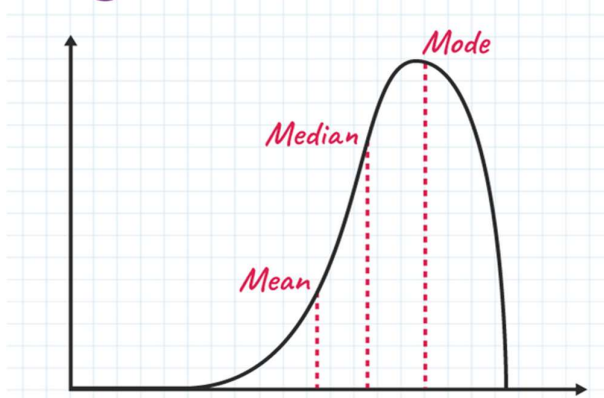
• **Positive Skewness:** - Positive Skewness means the tail on the right side of the distribution is longer. The mean and median will be greater than the mode.

Condition for positive skewness = **Mean > Median > Mode**

Statistics And Probability For Data Scientist And Data Analyst

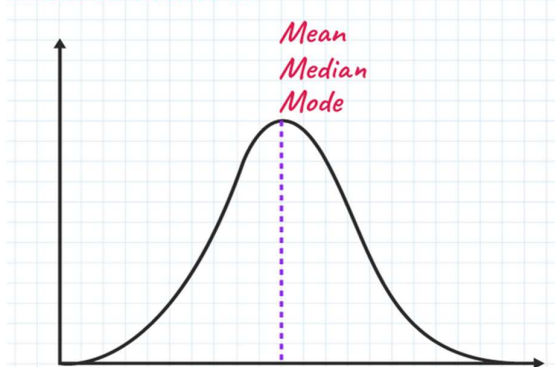
- **Negative Skewness:** - Negative Skewness means when the tail of the left side of the distribution is longer than the tail on the right side. **The mean and median will be less than the mode.** Condition for negative skewness is **Mode > Median > Mean.**

Negative Skew

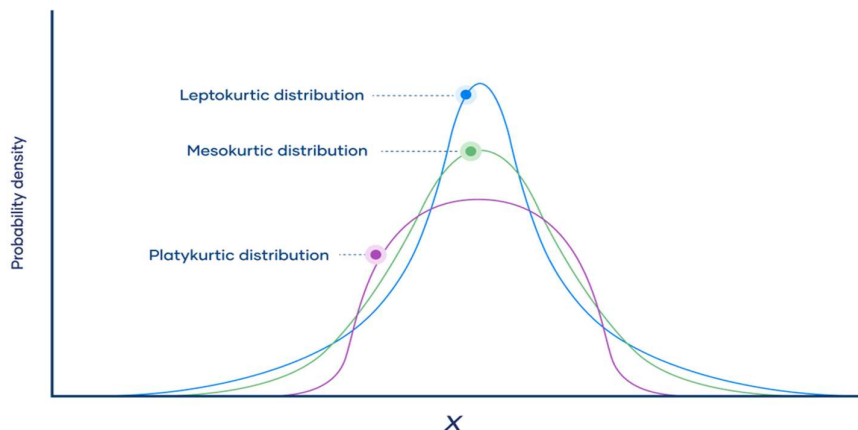


- **Zero Skewness:** - It is also known as a “symmetric distribution”. It signifies that distribution of data is evenly distributed around the mean. Condition for zero skewness is **Mean = Mode = Median.**

Zero Skew



Kurtosis: - **Kurtosis** is a measure of the Tailedness of a distribution.
Excess kurtosis is the tailedness of a distribution relative to a normal distribution.



There are 3 Types of Kurtosis:

Statistics And Probability For Data Scientist And Data Analyst

- **Mesokurtic Distribution:** - A **mesokurtic** distribution is medium-tailed, so outliers are neither highly frequent, nor highly infrequent.
Kurtosis is measured in comparison to normal distributions.
Normal distributions have a kurtosis of 3, so any distribution with a kurtosis of approximately 3 is mesokurtic.
Normal distributions have an excess kurtosis of 0, so any distribution with an excess kurtosis of approximately 0 is mesokurtic.
- **Platykurtic Distribution:** - A **platykurtic** distribution is thin-tailed, meaning that outliers are infrequent. **Platykurtosis is sometimes called** negative kurtosis.
A kurtosis of less than 3.
An excess kurtosis of less than 0.
- **Leptokurtic Distribution:** - A **leptokurtic** distribution is fat-tailed, meaning that there are a lot of outliers. **Leptokurtosis is sometimes called** positive kurtosis,
A kurtosis of more than 3
An excess kurtosis of more than 0.