

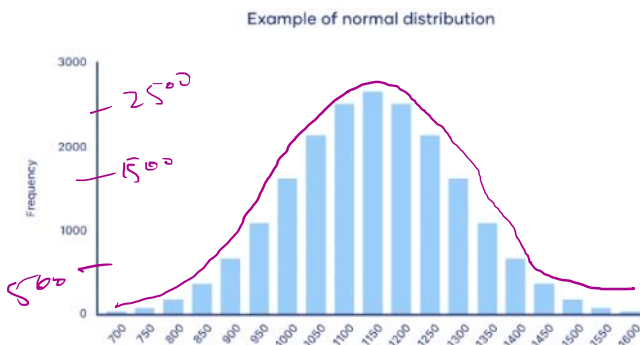
Agenda

- Normal or Gaussian Distribution
- Properties of Normal Distribution
- Empirical Rule in Normal Distribution
- Central Limit Theorem
- Covariance
- Pearson Coefficient Correlation

• Normal or Gaussian Distribution

In a normal distribution, data is symmetrically distributed with no skew. When plotted on a graph, the data follows a bell shape, with most values clustering around a central region and tapering off as they go further away from the center.

Normal distributions are also called Gaussian distributions or bell curves because of their shape.

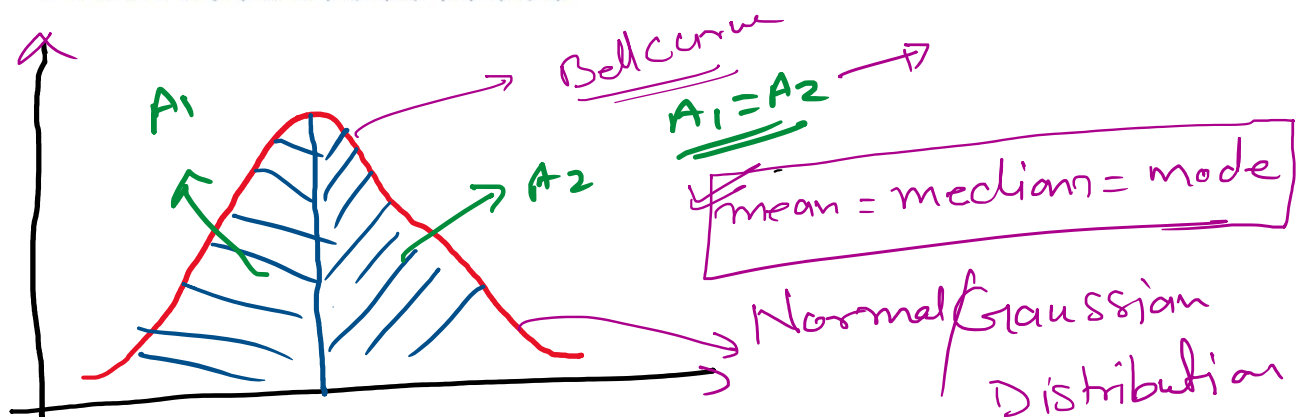


eg = Height

170 -
180 -
165
163
173
:
180

} = 100

mean = median = mode



* It is a type of continuous Probability distribution for a real-valued random variable.

Notation

$N(\mu, \sigma^2)$

mean

std

Parameter $\rightarrow \mu \in \mathbb{R} = \underline{\text{mean}}$

$\sigma^2 \in \mathbb{R}^+ = \text{variance.}$

$x \in \mathbb{R}$

$$\text{PDF} = \frac{1}{\sigma \sqrt{2\pi}} \times e^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2}$$

↓
Probability density function

Mean of Normal Distribution

mean(μ) \Rightarrow Average value

variance & std

var = σ^2

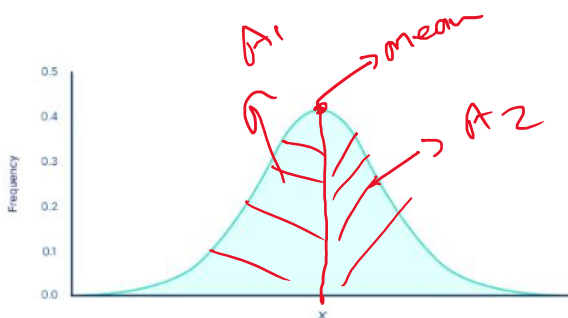
std = $\sqrt{\text{var}}$

What are the properties of normal distributions?

Normal distributions have key characteristics that are easy to spot in graphs:

- The mean, median and mode are exactly the same.
- The distribution is symmetric about the mean—half the values fall below the mean and half above the mean.
- The distribution can be described by two values: the mean and the standard deviation.

μ & std



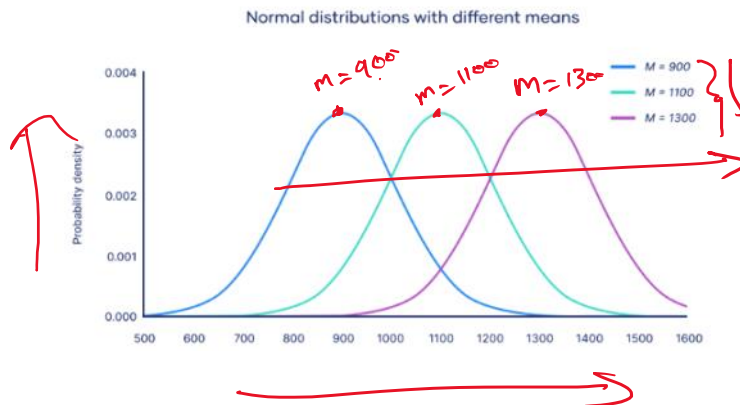
$A_1 = A_2 = A_{\text{area}}$

μ, std

Note 2

The mean is the location parameter while the standard deviation is the scale parameter.

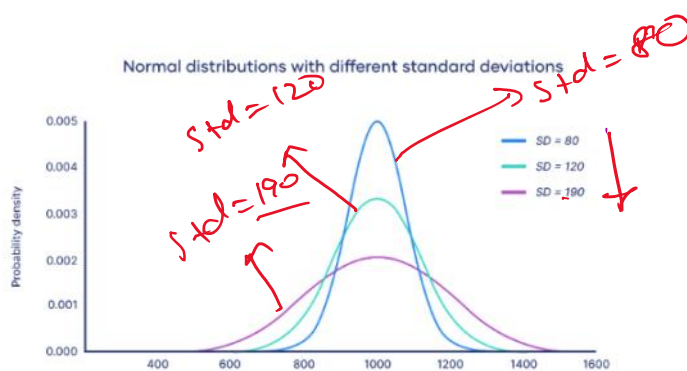
1. The mean determines where the peak of the curve is centered. Increasing the mean moves the curve right, while decreasing it moves the curve left.



$\mu \uparrow \& \text{ location}$

$\leftarrow \rightarrow$

2. The standard deviation stretches or squeezes the curve. A small standard deviation results in a narrow curve, while a large standard deviation leads to a wide curve.



Empirical rule

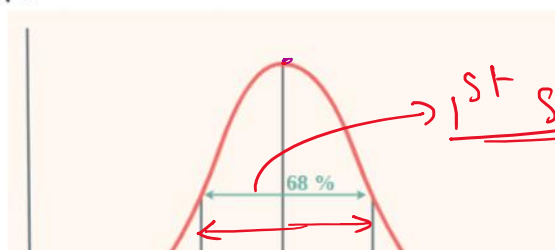
The **empirical rule**, or the 68-95-99.7 rule, tells you where most of your values lie in a normal distribution:

- Around 68% of values are within 1 standard deviation from the mean.
- Around 95% of values are within 2 standard deviations from the mean.
- Around 99.7% of values are within 3 standard deviations from the mean.

68%, 95%, & 99.7%

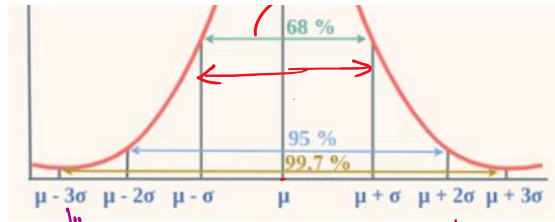
1st std
2nd std
3rd std

Normal Distribution Graph

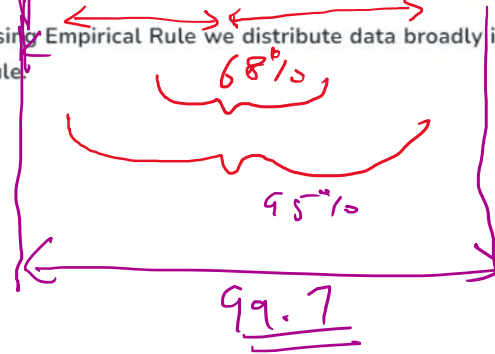


1st std = 68%

$\sigma^2 = \sigma$



Studying the graph it is clear that using Empirical Rule we distribute data broadly in three parts. And thus, empirical rule is also called "68 – 95 – 99.7" rule.



Probability

- ① $P(\mu - \sigma \leq x \leq \mu + \sigma) \approx 68\%$
- ② $P(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 95\%$
- ③ $P(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 99.7\%$

Central Limit Theorem

The **central limit theorem** states that if you take sufficiently large samples from a population, the samples' means will be **normally distributed**, even if the population isn't normally distributed.

The central limit theorem relies on the concept of a **sampling distribution**, which is the **probability distribution** of a **statistic** for a large number of **samples** taken from a population.

Imagining an experiment may help you to understand sampling distributions:

- Suppose that you draw a **random sample** from a population and calculate a **statistic** for the sample, such as the mean.
- Now you draw another random sample of the same size, and again calculate the **mean**.
- You repeat this process many times, and end up with a large number of means, one for each sample.

The distribution of the sample means is an example of a **sampling distribution**.

The central limit theorem says that the sampling distribution of the mean will always be **normally distributed**, as long as the sample size is large enough. Regardless of whether the population has a normal, Poisson, binomial, or any other distribution, the sampling distribution of the mean will be normal.

A normal distribution is a symmetrical, bell-shaped distribution, with increasingly fewer observations the further from the center of the distribution.

What is Z-Score?

Z-score, also known as the **standard score**, tells us the **deviation of a data point from the mean** by expressing it in terms of standard deviations above or below the mean. It gives us an idea of how far a data point is from the mean. Hence, the Z-Score is measured in terms of standard deviation from the mean. For example, a **Z-score of 2** indicates the value is **2 standard deviations away from the mean**. To use a z-score, we need to know the population mean (μ) and also the population standard deviation (σ).

Z-score is a statistical measure that describes a value's position relative to the mean of a group of values. It is

✓ Z-score is a statistical measure that describes a value's position relative to the mean of a group of values. It is expressed in terms of standard deviations from the mean. The Z-score indicates how many standard deviations an element is from the mean. stat u

-3 to +3



Z-Score Formula

To calculate the z- score for any given data we need the value of the element along with the mean and standard deviation.

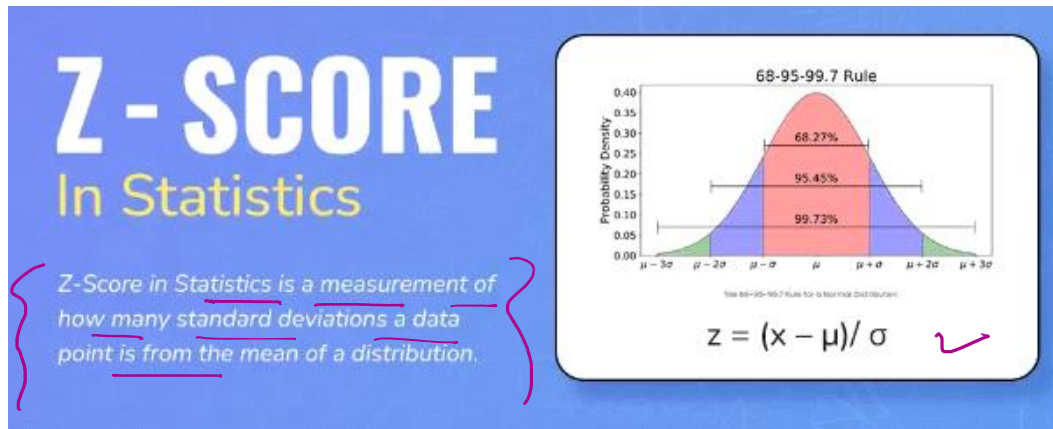
A z-score can be calculated using the following Z- score formula.

$$z = (X - \mu) / \sigma$$

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

where,

- z = Z-Score
- X = Value of Element
- μ = Population Mean
- σ = Population Standard Deviation



Ex

	A	B	C
	Factor (x)	Mean (μ)	St. Dev. (σ)
1	3	12.17	6.4
2	13	12.17	6.4
3	8	12.17	6.4
4	21	12.17	6.4
5	17	12.17	6.4
6	11	12.17	6.4



Z score

=> same distribution

	A	B	C	D
	Factor (x)	Mean (μ)	St. Dev. (σ)	Z-Score
1	3	12.17	6.4	-1.43
2	13	12.17	6.4	1.43

-3 to +3

	Factor (x)	Mean (μ)	St. Dev. (σ)	Z-Score
1				
2	3	12.17	6.4	-1.43
3	13	12.17	6.4	0.13
4	8	12.17	6.4	-0.65
5	21	12.17	6.4	1.38
6	17	12.17	6.4	0.75
7	11	12.17	6.4	-0.18

Covariance & correlation

Q: what is relationship b/w x & y

x	y
2	3
4	5
6	7
8	9

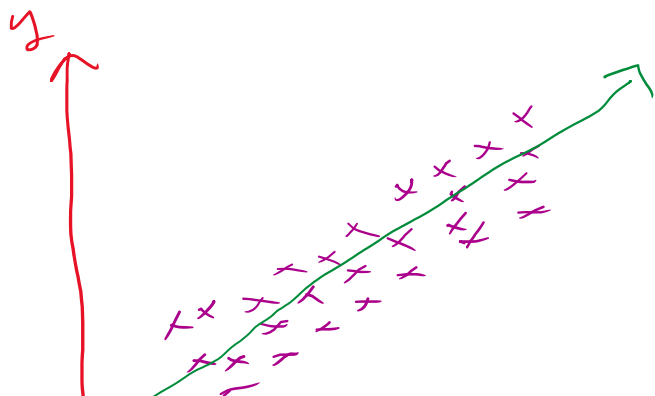
- (a) $x \uparrow$ $y \uparrow$
 (b) $x \downarrow$ $y \downarrow$
 (c) $x \uparrow$ $y \downarrow$
 (d) $x \downarrow$ $y \uparrow$

①

$x \uparrow$ $y \uparrow$

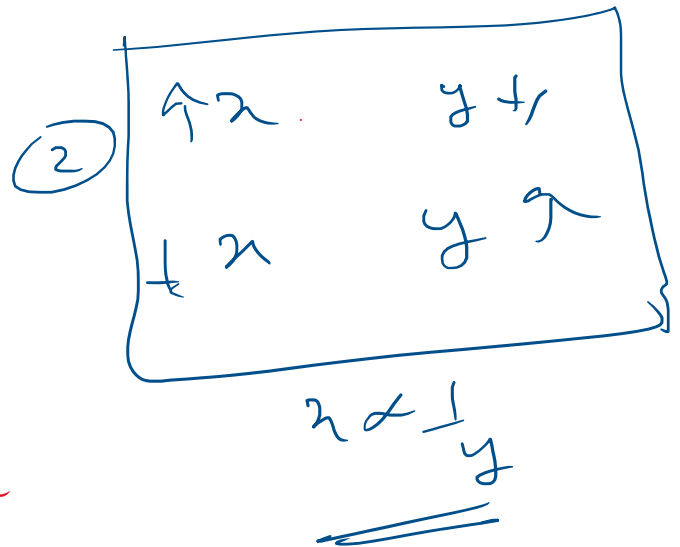
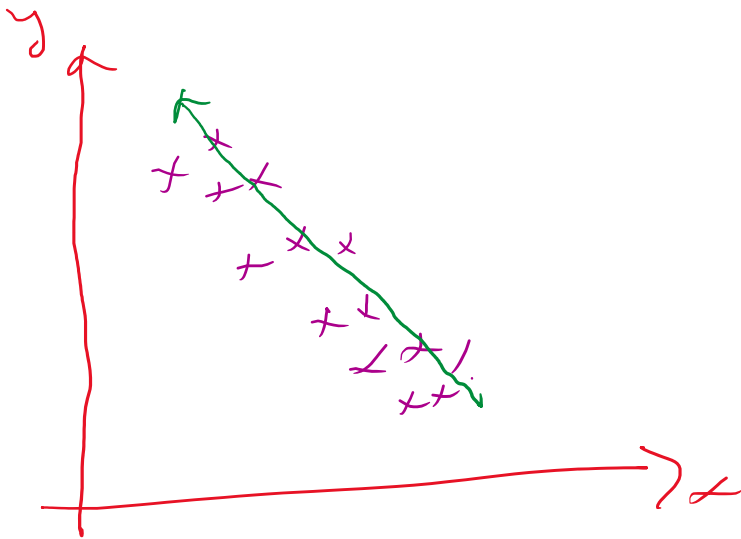
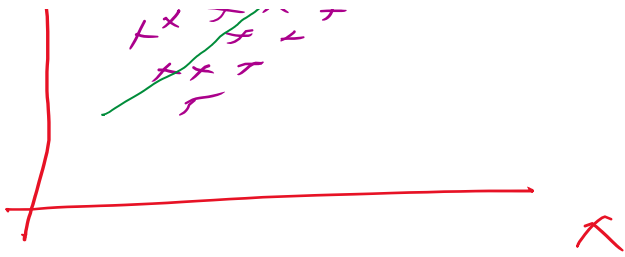
②

$x \downarrow$ $y \downarrow$



①

$x \uparrow$ $y \uparrow$
 $x \downarrow$ $y \downarrow$



Covariance

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$x_i \rightarrow$ Data point of x

$\bar{x} \rightarrow$ Sample mean

$y_i \rightarrow$ Data points of y

$\bar{y} \rightarrow$ Sample mean of y

(1) +ve Covariance

$\text{cov}(x, y) \rightarrow \underline{\underline{+ve}}$

$x \uparrow$	$y \uparrow$
$x \downarrow$	$y \downarrow$

$x \propto y$

(2) -ve Covariance

$x \uparrow$	$y \downarrow$
$x \downarrow$	$y \uparrow$

$x \propto \frac{1}{y}$

Ex

$\frac{x}{2}$

$\frac{y}{3}$

4

5

6

7

$\bar{x} = 4$

$\bar{y} = 5$

$$\text{cov}(x, y) = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$\Rightarrow \frac{[(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)]}{n-1}$$

A hand-drawn red coordinate system on a white background. It features two perpendicular axes: a vertical y-axis pointing upwards and a horizontal x-axis pointing to the right. Both axes have arrowheads at their ends. A third red line, representing a vector, originates from a point in the first quadrant and points towards the origin. The axes are labeled with 'y' at the top and 'x' at the right end.

(x & y having +ve covariance)

Motu

① Relationship $x \Delta y$
(+ve or -ve)

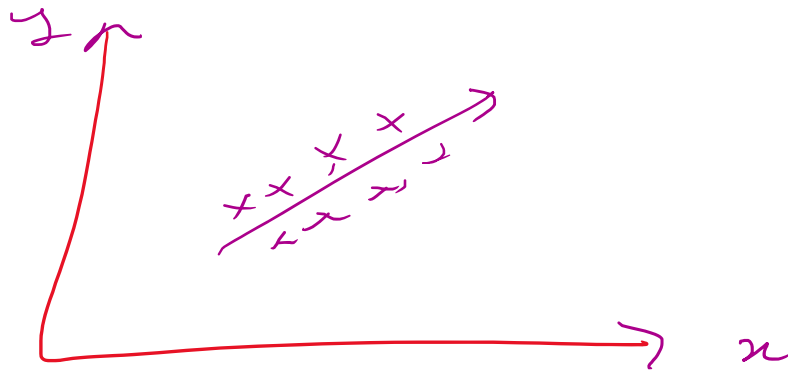
(11) Covariance does not have a specific limit value

Pearson Correlation Coefficient

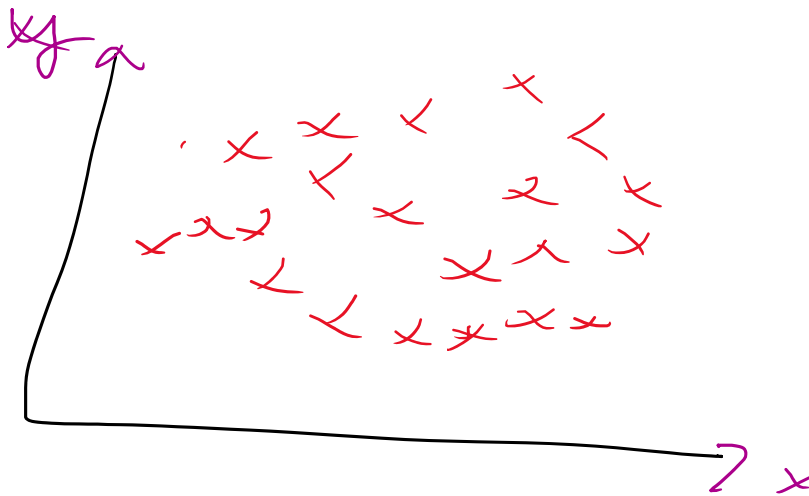
~~$[-1 \quad t \quad 1)$~~

$$\rho_{[x,y]} = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y}$$

① Between 0 & 1 \rightarrow positive correlation



② 0 \Rightarrow No correlation
 \Downarrow
No relationship



③ Between 0 & -1
Negative correlation



