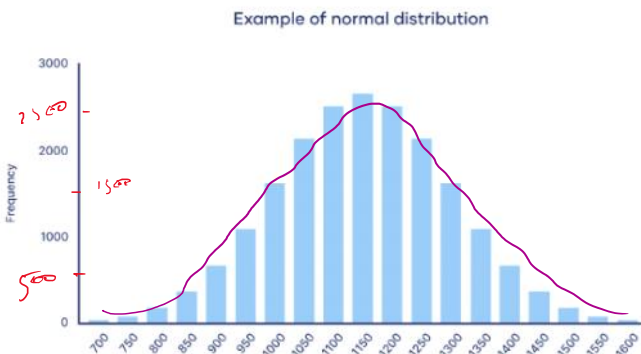# Session-16 Statistics

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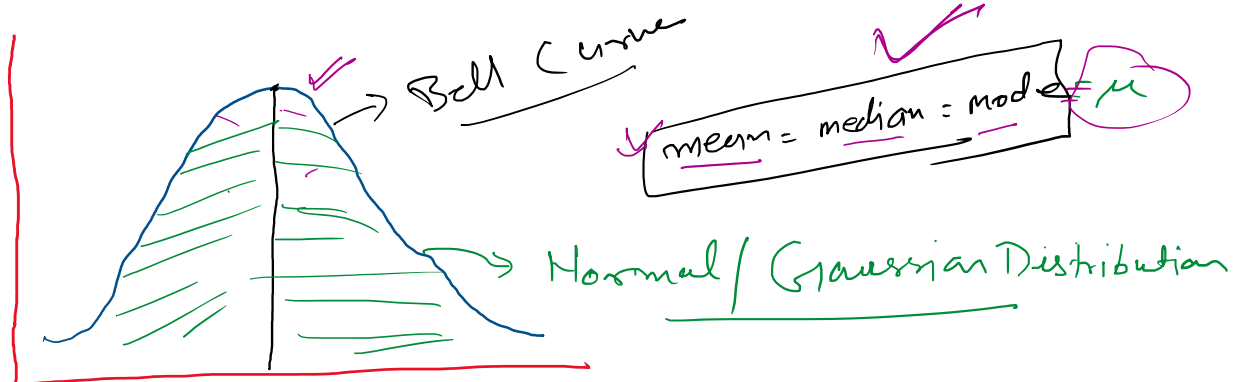
- **Normal or Gaussian Distribution**

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 }
=> Weight }

-) Age
-) IQ }

mean ≠ median ≠ mode

**Example of normal distribution**

Bell Curve

mean = median = mode = $\mu$

Normal / Gaussian Distribution

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

Parameters : $\mu \in \mathbb{R}$ = mean

$\sigma^2 \in \mathbb{R} > 0$ = variance

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

mean $(\mu)$ = Average Value

variance & std

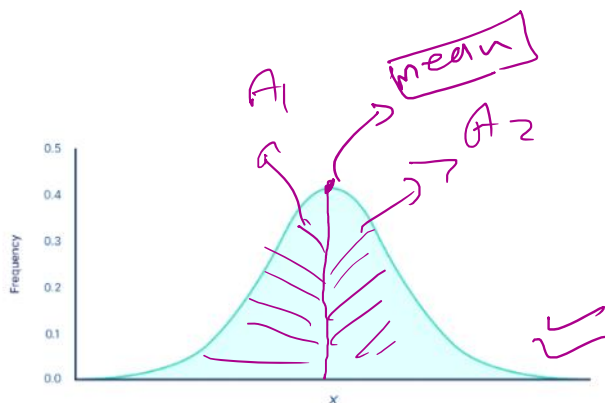var $= \sigma^2$

std $= \sqrt{var}$

## 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.
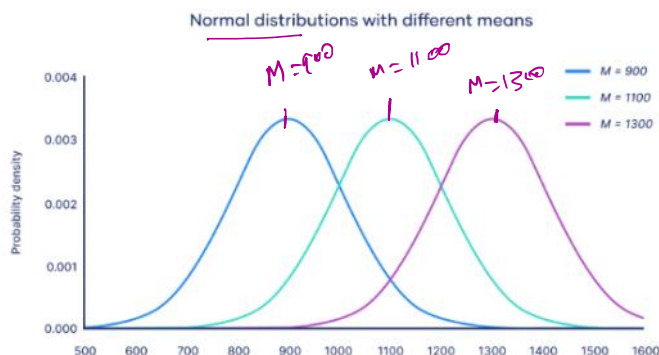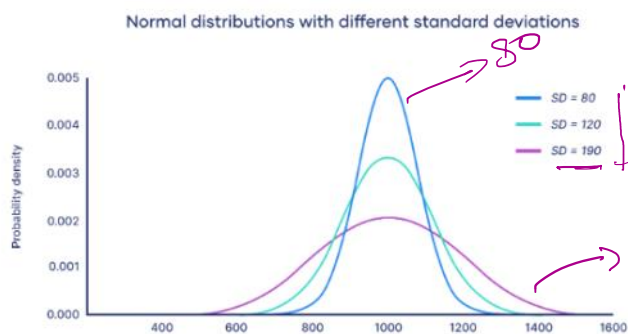


$A_1 = A_2 = $ Area

$(\mu, std)$

**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.

*m*

### Normal distributions with different means

*M = 900    M = 1100    M = 1300*

| | M = 900 |
| | M = 1100 |
| | M = 1300 |

*M↑*

*M↓*

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.
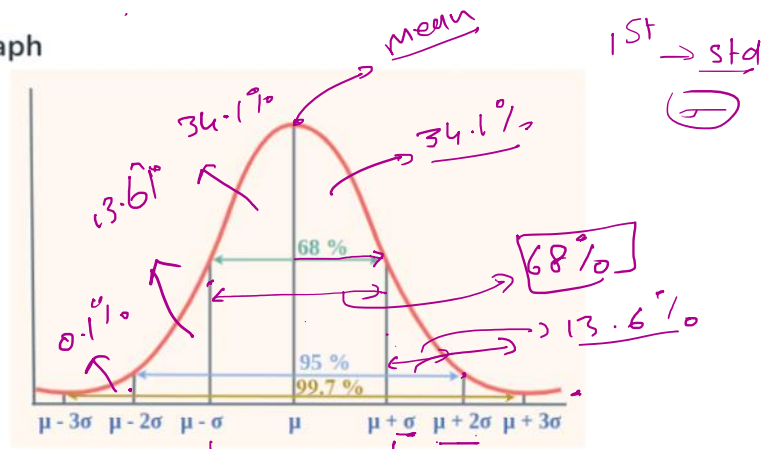
### Normal distributions with different standard deviations

*80*

| | SD = 80 |
| | SD = 120 |
| | SD = 190 |

*190*

*↑ std ⇒ wide curve*
*↓ std ⇒ narrow*

## 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.

## Normal Distribution Graph



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.**

$$Pr\left(\mu - \sigma \leq x \leq \mu + \sigma\right) \approx 68\%$$

$$Pr\left(\mu - 2\sigma \leq x \leq \mu + 2\sigma\right) \approx 95\%$$

$$Pr\left(\mu - 3\sigma \leq x \leq \mu + 3\sigma\right) \approx 99.75\%$$

et IRIS $\rightarrow$ (Sepal width, Sepal length) $\rightarrow$

## 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.
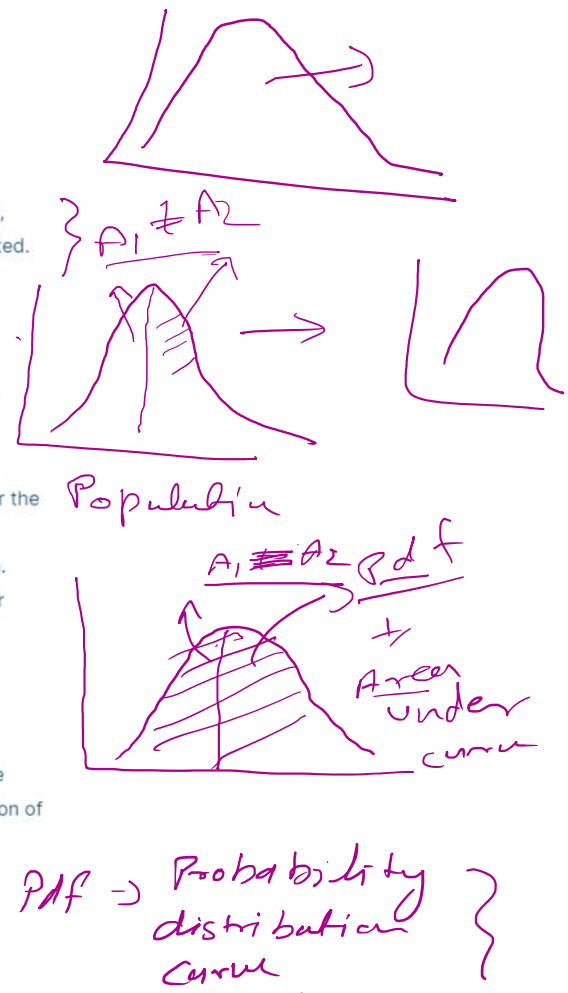
Imagining an experiment may help you to understand sampling distributions:

- 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.
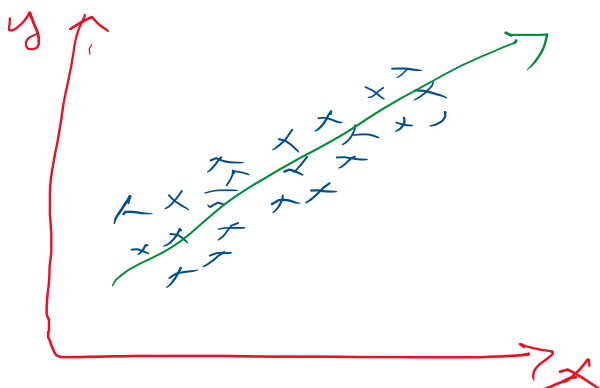
Population

$A_1 \neq A_2$ pdf + Area under curve

Pdf $\rightarrow$ Probability distribution curve

# Covariance & Correlation

| x | y |
|---|---|
| 2 | 3 |
| 4 | 5 |
| 6 | 7 |
| 8 | 9 |

Q what is relationship, $x$ & $y$

(a) $x \uparrow \quad y \uparrow$

(b) $x \downarrow \quad y \uparrow$

(c) $x \uparrow \quad y \downarrow$

(d) $x \downarrow \quad y \downarrow$

$$\boxed{\begin{array}{l} \uparrow x \propto y \uparrow \downarrow, \\ \uparrow \uparrow x \propto \dfrac{1}{\downarrow y \uparrow} \end{array}}$$



a) $\boxed{\begin{array}{ll} x \uparrow & y \uparrow \\ x \downarrow & y \downarrow \end{array}}$

$$x \propto y$$



b) $\boxed{\begin{array}{ll} x \uparrow & y \downarrow \\ x \downarrow & y \uparrow \end{array}}$

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

## Covariance

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

$x_i = $ Data Points of $x$

$\bar{x} \rightarrow$ Sample mean

$y_i \rightarrow$ DataPoint of $y$

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

$\text{cov}(x, y)$

$x \uparrow \quad y \uparrow$
$x \downarrow \quad y \downarrow$ $\rightarrow$ +ve Covariance

$x \uparrow \quad y \downarrow$
$x \downarrow \quad y \uparrow$ $\rightarrow$ -ve Covariance

eg

| $x$ | $y$ |
|-----|-----|
| 2 | 3 |
| 4 | 5 |
| 6 | 7 |
| $\bar{x} = 4$ | $\bar{y} = 5$ |

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

$$= \left[ \frac{(2-4)(3-5) + (4-4)(5-5)}{} + (6-4)(7-5) \right] \Big/ 2$$

$$= \frac{4 + 0 + 4}{2} = \frac{8}{2} = 4$$

$\{ x \Delta y$ having +ve covariance $)$

+ve covariance

( +ve covariance )

Note ⇒ i) Relationship $x$ & $y$
( +ve or -ve value )

ii) covariance does not have a specific
limit value.

## Pearson Correlation Coefficient (r)

$$[\, -1 \text{ to } 1 \,]$$

$$P_{(x,y)} = \frac{cov(x,y)}{\sigma_x \, \sigma_y}$$

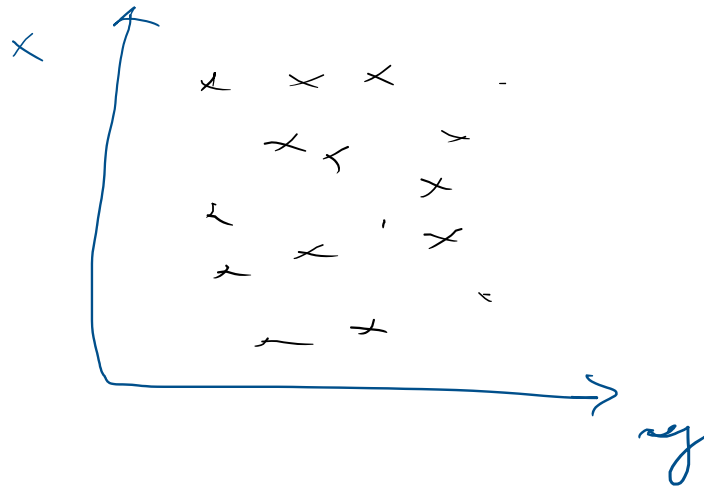① Between 0 & 1 → Positive correlation



② 0     No correlation

No relationship



③ Between 0 & −1    Negative correlation