# Probability & Statistics :-

## Random variable :-

ex:- rolling dice → 6 sides $= \{1,2,3,4,5,6\}$ ⎫
when rolled ⎬ Random Experiment
any one of them : ⎭
equal outcome

→ Sample Space.

random variable → $\boxed{X} = \boxed{\{1,2,3,4,5,6\}}$

tossing a coin → $Y = \boxed{\{H,T\}}$

$\boxed{P(x=1)} = \frac{1}{6}$    $P(x=2) = \frac{1}{6} \cdots$

$P(X \text{ is even}) = \frac{3}{6} = \frac{1}{2}$

$\begin{pmatrix} \text{probability} \\ \text{of } X \text{ being} \\ \text{even} \end{pmatrix}$

$\left( P(x=2) + P(x=4) + P(x=6) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} \right)$

$P(x \text{ is odd}) = \frac{1}{2}$

$P(x = x_i) \longrightarrow P(x_i)$   Same thing diff notation.

Finite set of values → Discrete random value

→ Height of randomly picked student.
  y could be 162, 180, 120, 140, ----- → infinite values. → Continuous Random Variable

## Outliers :-

Y : height of student.
$\{122.2, 146.4, 132.5, \cdots, \boxed{12.2}, 156.3, \boxed{92.7} \cdots \}$

→ could be an outlier

outlier → could be human error
          (or)
          actual height

→ Outliers can corrupt data.

→ A discrete value is obtained by counting
→ A continuous value is obtained by measuring.

Sample Space :- Set of all possible outcomes of an experiment.
→ A random variable value depends on the outcome of a random phenomenon.
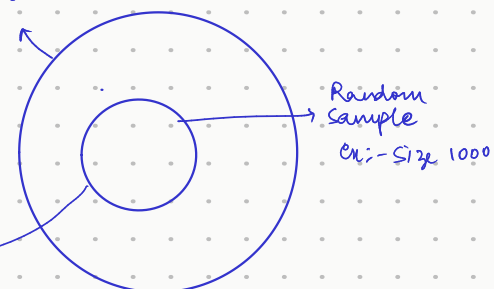
## Population & Sample :-

→ Estimating the average height of human

$\longrightarrow \mu = \frac{1}{Pop} \sum_{i=1}^{Pop} h_i$  (IMPOSSIBLE)

So we estimate.

Set of all humans in the world.

Random Sample
ex:- size 1000

often Represented by $\longleftarrow \bar{x} = \frac{1}{1000} \sum_{i=1}^{1000} h_i \longleftarrow$
$\bar{x}$

→ As sample size increases $\quad$ $\overline{x} = \mu$

$\overline{x}$ ← Sample mean

$\mu$ → original mean

Sampling is of two types:-

(i) Simple Sampling

(ii) Stratified Sampling $\quad\longrightarrow$ Unbiased Sampling & more accurate results.

ex:- 1000 people $\cdots\cdots\cdots\longrightarrow$ original sample

400 have cars $\quad$ 300 have bikes $\quad$ 200 have cycles $\quad$ 100 have nothing

Sample Size $\longrightarrow$ 250

**Simple Random Sampling**

250 could have cars
(or)

250 could have bikes
(or)

100 bikes + 150 cars.

**Stratified Random Sampling**

Cars $\longrightarrow$ 100
bikes $\longrightarrow$ 75
cycles $\longrightarrow$ 50
nothing $\longrightarrow$ 25

} These are random but equal imp to all classes

**Gaussian Distribution :- (AKA Normal Distribution)**

→ If X is a continuous random Variable & X has a PDF curve ( $\boxed{\mathord{\sim}}$ ), then we say X has a Gaussian distribution.

PDF = Probability Density Function

**why learn?**

→ Stuff in nature tends to follow G.D.
→ heights, weights of people follow it. (Natural Phenomenon)
→ They are simple models that summarize R.V.

ex:-



$\mu$ = mean
$\sigma^2$ = variance $\quad$ } → parameters

[ if this is not known, then we calc ]

[ if we are given $\mu$, $\sigma^2$ & told that X follows GD, we can plot PDF. We don't need the whole data ]

Variance is a measure of spread.

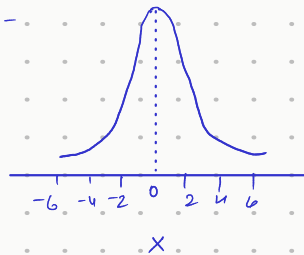PDF → $\varphi_{\mu,\sigma^2}(x)$

X → Continuous Random Variable

Red, Yellow, Blue have $\mu = 0$ but varying variances.

The peak of curve is usually at '$\mu$'.

→ The parameters of Gaussian Distribution are $\mu$ & $\sigma^2$

$$X \sim N(\mu, \sigma^2) \quad (\Longrightarrow X \text{ follows Gaussian Distribution with } \mu \text{ & } \sigma^2)$$

ex:-



$$\longrightarrow X \sim N(0, 2)$$

$-6 \quad -4 \quad -2 \quad 0 \quad 2 \quad 4 \quad 6$

$X$

$$\rightarrow P(X = x) = p(x) = \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left\{\frac{-(x-\mu)^2}{2\sigma^2}\right\}$$

↑
Probability Density at a point $(x)$ : PDF at any given point gives the probability density at that point. Probability of getting a single discrete value is 0.

ex:- If $\mu = 0$, $\sigma^2 = 1$, $\sigma = 1$

$$f(x) = \left(\frac{1}{\sqrt{2\pi}}\right) \exp\left\{\left(\frac{-1}{2}\right)x^2\right\} = y$$
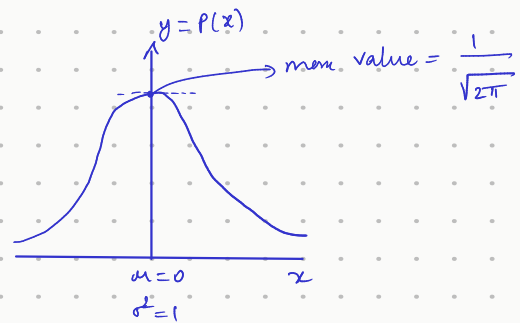
constants

further simplifying

$$y = \exp(-x^2)$$

↳ when plotted

$y = P(x)$

→ max. value $= \dfrac{1}{\sqrt{2\pi}}$

as $x$ increases
$y$ decreases.
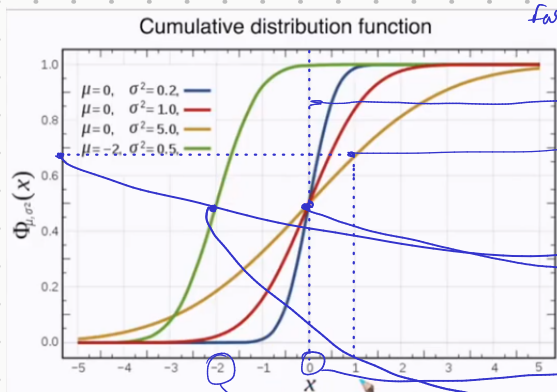as $x$ decreases
$y$ decreases.

$\mu = 0$
$\sigma^2 = 1$

conclusions :-
① $x$ moves away from $\mu$, $y$ decreases.
② Graph is symmetric.
③ In this particular graph, $y$ is reducing exponential squared. $\left(e^{-x^2}\right)$

→ If $p(x)$ is the probability density at a point '$x$', the probability can be obtained by computing the integral of $p(x)$ over a given interval.
i.e., probability of getting $X \in [a, b]$ is $\displaystyle\int_a^b p(x)\,dx$

## Cumulative Distribution Function (CDF) of Gaussian Distribution/Normal Distribution:-



Cumulative distribution function

$\mu = 0, \quad \sigma^2 = 0.2,$ ——
$\mu = 0, \quad \sigma^2 = 1.0,$ ——
$\mu = 0, \quad \sigma^2 = 5.0,$ ——
$\mu = -2, \sigma^2 = 0.5,$ ——

As $\sigma^2$ increases, CDF goes far from centre line

CDF of a random variable looks like ⌐⌐√

$\rightarrow P(X \le 1) = 0.65$

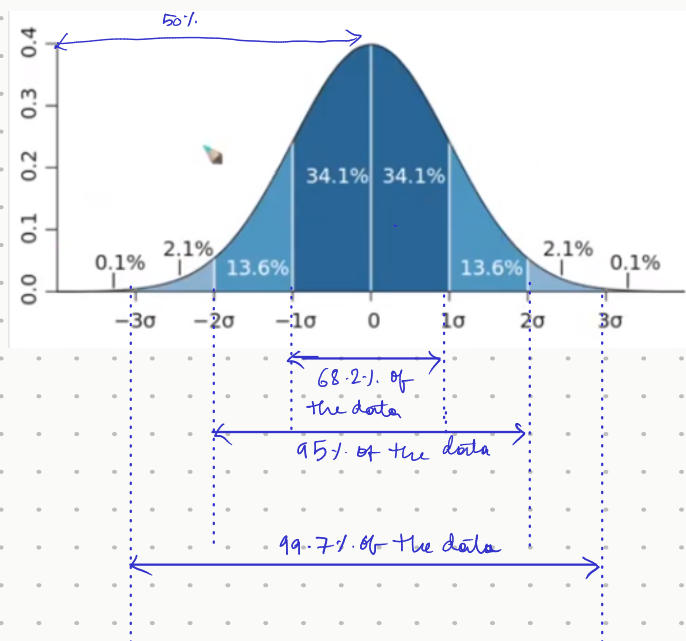$\rightarrow \mu = 0$, centre of CDF is at 0.

$\rightarrow \mu = -2$, centre of CDF is at -2

$$CDF = \frac{1}{2}\left[1 + erf\left(\frac{x - \mu}{\sigma\sqrt{2}}\right)\right] \longrightarrow \text{No need to memorize}$$

**68 - 95 - 99.7 rule :-**

if $\mu = 0$, $\sigma^2 = 4 \implies \sigma = 2$      $X \sim N(0, 4)$



- 50%
- 34.1%  34.1%
- 0.1%  2.1%  13.6%  13.6%  2.1%  0.1%
- $-3\sigma$  $-2\sigma$  $-1\sigma$  0  $1\sigma$  $2\sigma$  $3\sigma$
- 68.2% of the data
- 95% of the data
- 99.7% of the data

. How is this useful?

  eg :- if human population weight

  $$X \sim N(150, 25)$$
  $$\downarrow\mu \qquad \downarrow\sigma$$

  $\implies$ 68.2% of human populations lies b/w
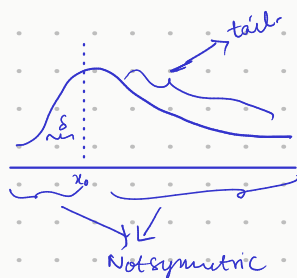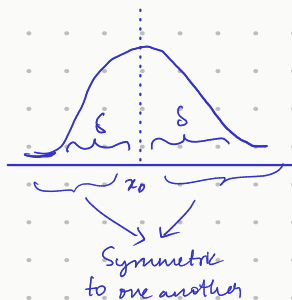  $$(150 - 25, \ 150 + 25)$$
  95% of people $(150 - 50, \ 150 + 50)$
  99.7% of people $(150 - 75, \ 150 + 75)$

$\rightarrow$ A standard guassian distribution always has a mean of 0 & Variance 1.
If it has other mean & variance, it's a non standard guassian distribution.

**Symmetric Distribution , Skewness & Kurtosis :-**

$\rightarrow$ They help understand shape of PDF.



tail

Symmetric to one another

Not symmetric

$\rightarrow$ A probability distribution is said to be symmetric if and only if there exists a value $x_0$ such that
$$f(x_0 - \delta) = f(x_0 + \delta) \text{ for all real numbers } \delta$$
$f(x)$ is the height of PDF at any point 'x'