

# Probability & Statistics :-

## Random variable :-

ex:- rolling dice  $\rightarrow$  6 sides =  $\{1, 2, 3, 4, 5, 6\}$   
when rolled

any one of these equal outcome

Random Experiment

random variable

$$X = \{1, 2, 3, 4, 5, 6\}$$

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

Sample Space

(probability of X being even)

$$P(X=1) = \frac{1}{6} \quad P(X=2) = \frac{1}{6} \dots$$

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

$$(P(X=2) + P(X=4) + P(X=6)) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6}$$

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

$$P(X=x_i) \rightarrow P(x_i) \text{ same thing diff notation}$$

Finite set of values  $\rightarrow$  Discrete random value

$\rightarrow$  Height of randomly picked student

Y could be 162, 180, 120, 140, ...

$\rightarrow$  infinite values  $\rightarrow$  Continuous Random Variable

## Outliers :-

Y: height of student

$\{122.2, 146.4, 132.5, \dots, 12.2, 156.3, 92.7, \dots\}$

outliers  $\rightarrow$  could be human error (or) actual height

could be an outlier

$\rightarrow$  Outliers can corrupt data

$\rightarrow$  A discrete value is obtained by counting

$\rightarrow$  A continuous value is obtained by measuring

Sample Space:- Set of all possible outcomes of an experiment

$\rightarrow$  A random variable value depends on the outcome of a random phenomenon

## Population & Sample :-

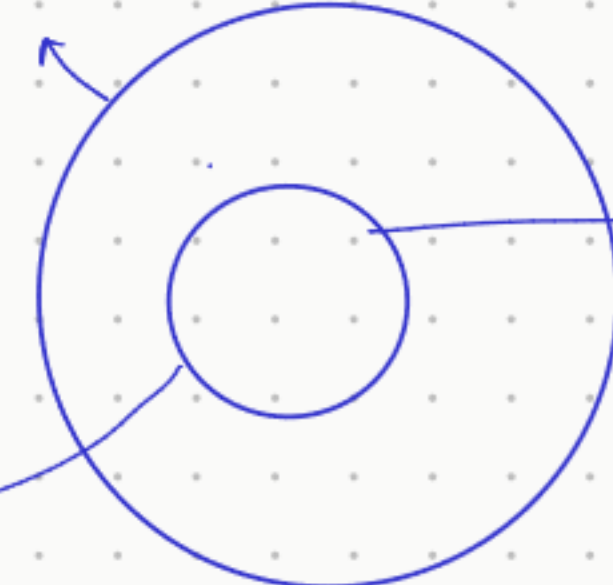
$\rightarrow$  Estimating the average height of human

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

So we estimate

often represented by  $\bar{x} = \frac{1}{1000} \sum_{i=1}^{1000} h_{is}$

Set of all humans in the world



Random Sample  
ex:- Size 1000



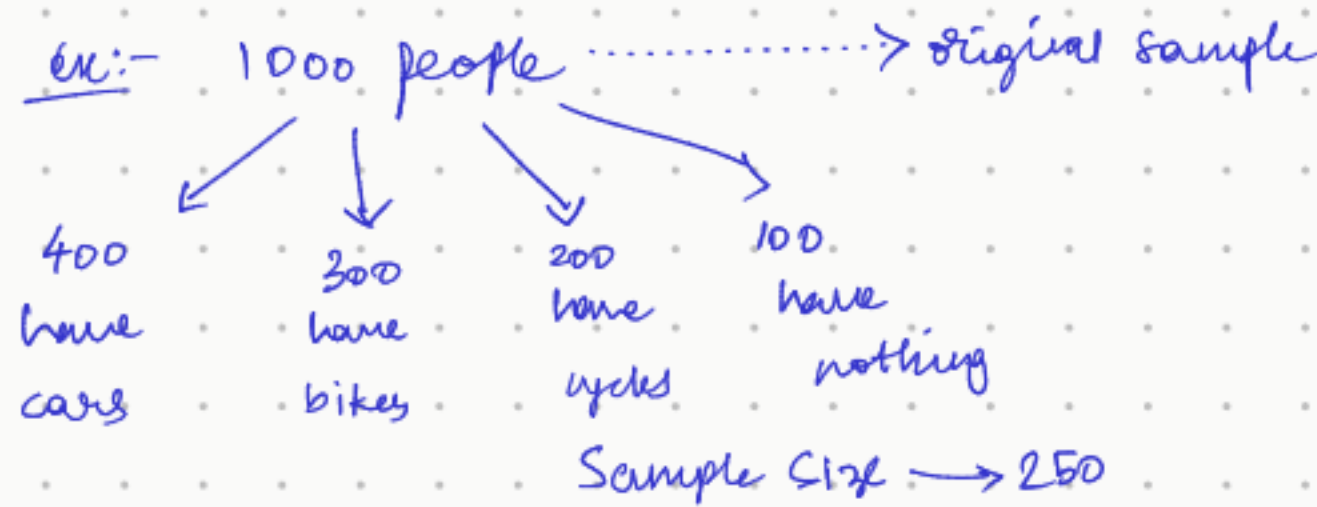
→ As sample size increases



Sampling is of two types:-

(i) Simple Sampling

(ii) Stratified Sampling → Unbiased sampling & more accurate results.



Simple Random Sampling

250 could have cars (a)

250 could have bikes (b)

100 bikes + 150 cars

Stratified Random Sampling

Cars → 100  
bikes → 75  
cycles → 50  
nothing → 25

There are random but equal imp to all classes

Gaussian Distributions:- (AKA Normal Distribution)

→ If  $X$  is a continuous random variable &  $X$  has a PDF curve (graph of a bell curve), then we say  $X$  has a Gaussian distribution.

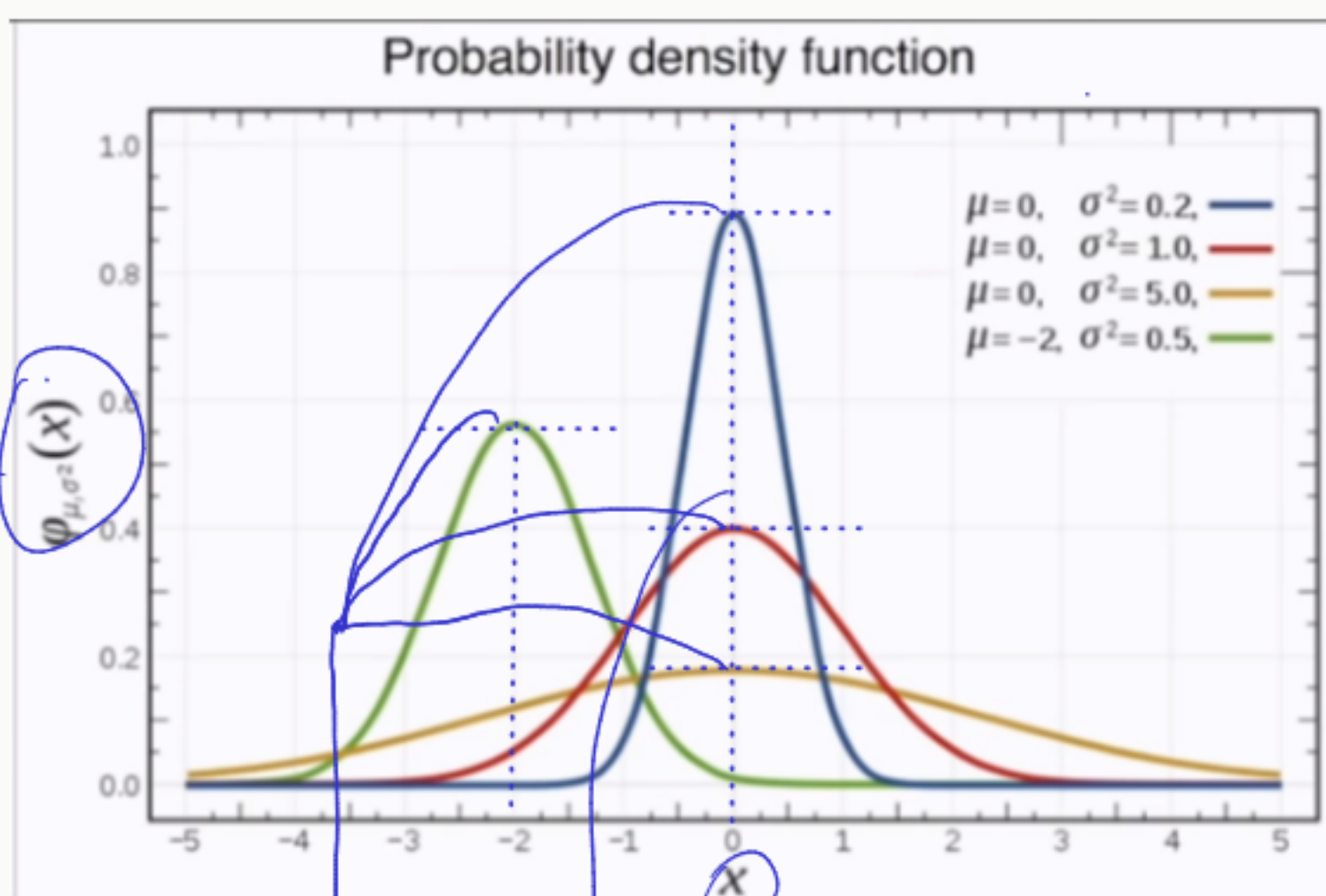
Why learn?

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

PDF = Probability Density Function

Mean, Median & Mode are all same for G.D. There are equal number of measurements above & below the mean

ex:-



$\mu$  = mean  
 $\sigma^2$  = variance

→ parameters

if this is not known, then we count

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

Variance is a measure of spread.

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

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



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

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

ex:-



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

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



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

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

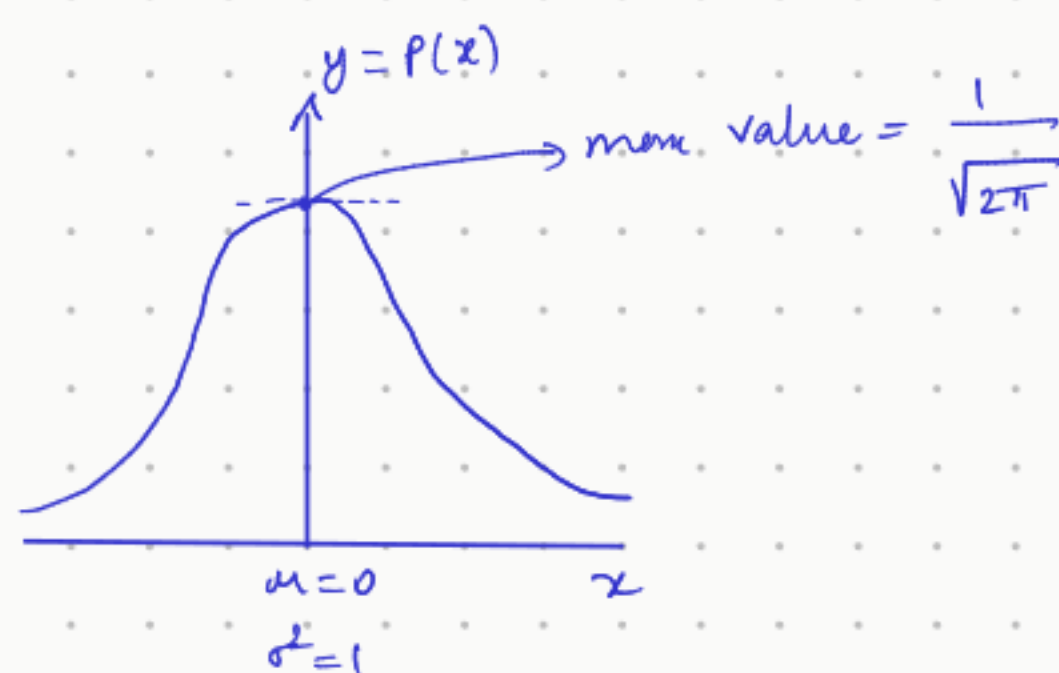
constants

further simplifying

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

when plotted

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



conclusions:-

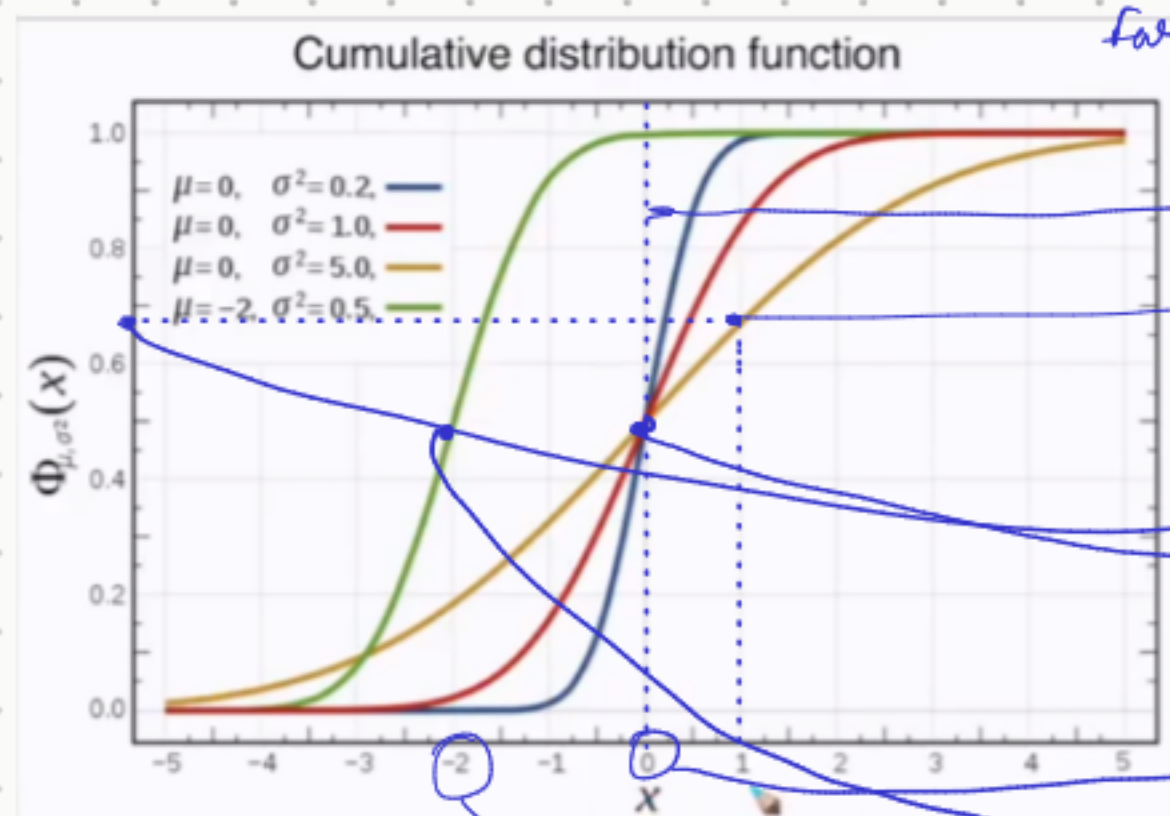
- ①  $x$  moves away from  $\mu$ ,  $y$  decreases.
- ② Graph is symmetric.
- ③ In this particular graph,  $y$  is reducing exponential squared. ( $e^{-x^2}$ )

→ 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  $\int_a^b p(x) dx$

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

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

CDF of a random variable looks like



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

$\mu=0$ , center of CDF is at 0.

$\mu=-2$ , center of CDF is at -2

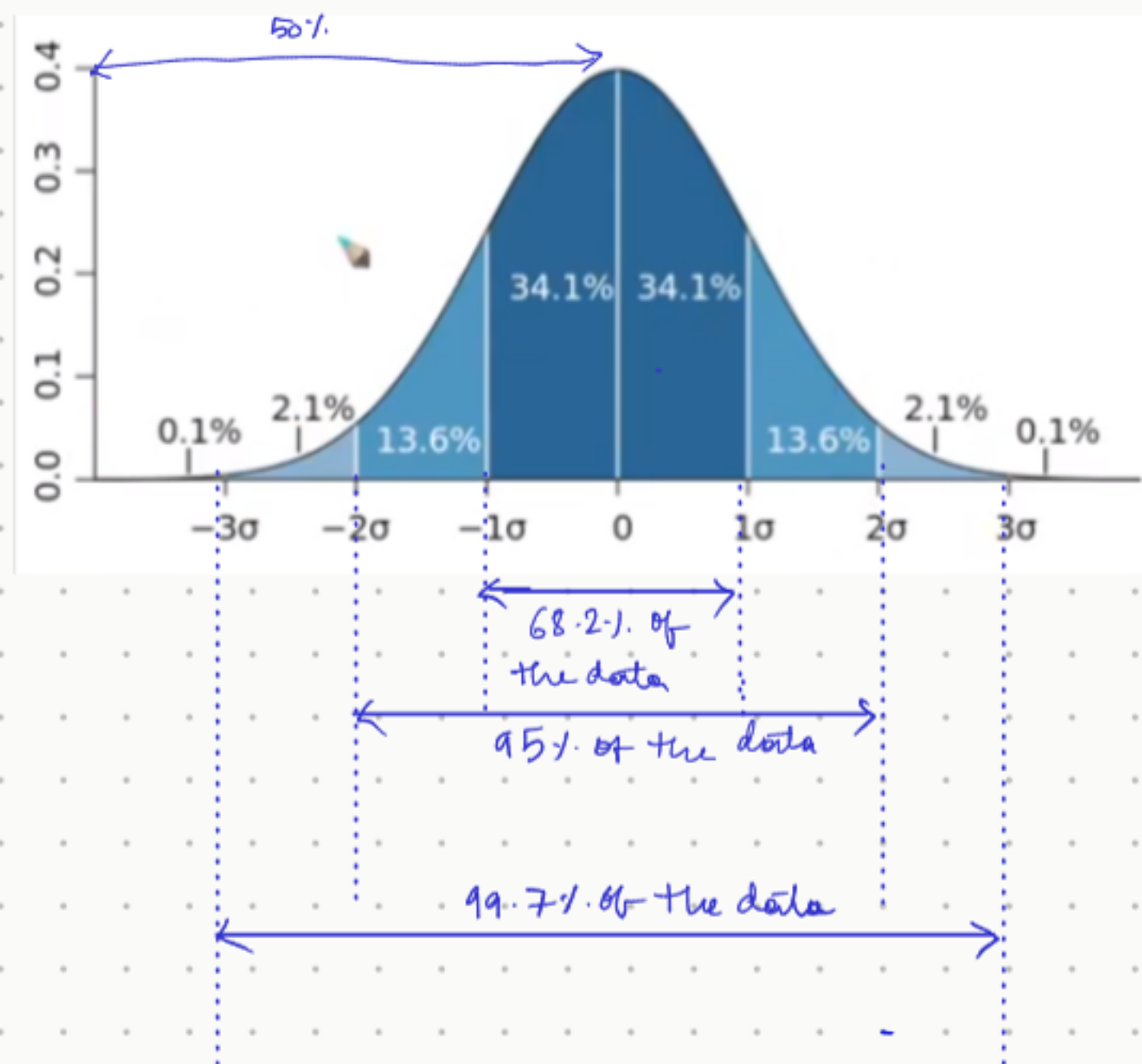


$$CDF = \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{x - \mu}{\sigma \sqrt{2}} \right) \right] \longrightarrow \text{No need to memorize}$$

68-95-99.7 rule:-

$$\text{if } \mu = 0, \sigma^2 = 4 \Rightarrow \sigma = 2$$

$$X \sim N(0, 4)$$



How is this useful?

ex:- if human populations height

$$X \sim N(150, 25)$$

$\downarrow$                        $\downarrow$   
 $\mu$                        $\sigma$

$\Rightarrow$  68.2% of human populations lies b/w  
(150-25, 150+25)

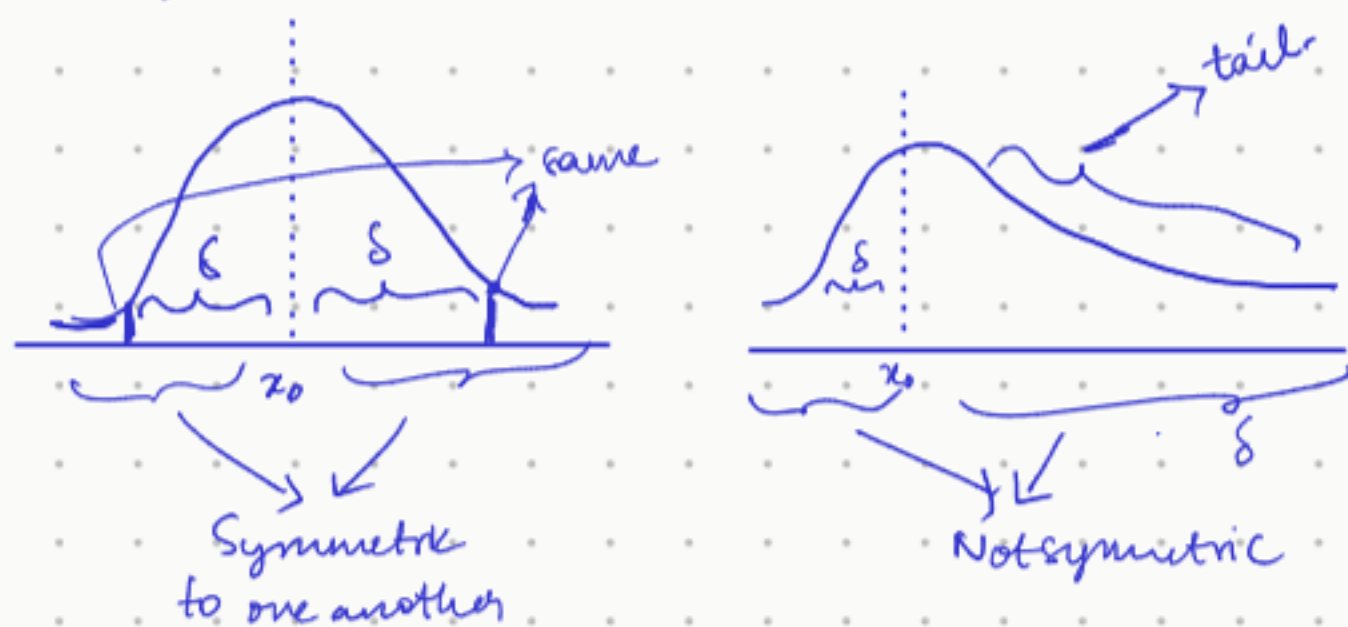
95% of people (150-50, 150+50)

99.7% of people (150-75, 150+75)

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

Symmetric Distribution, Skewness & Kurtosis:-

$\rightarrow$  They help understand shape of PDF.



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

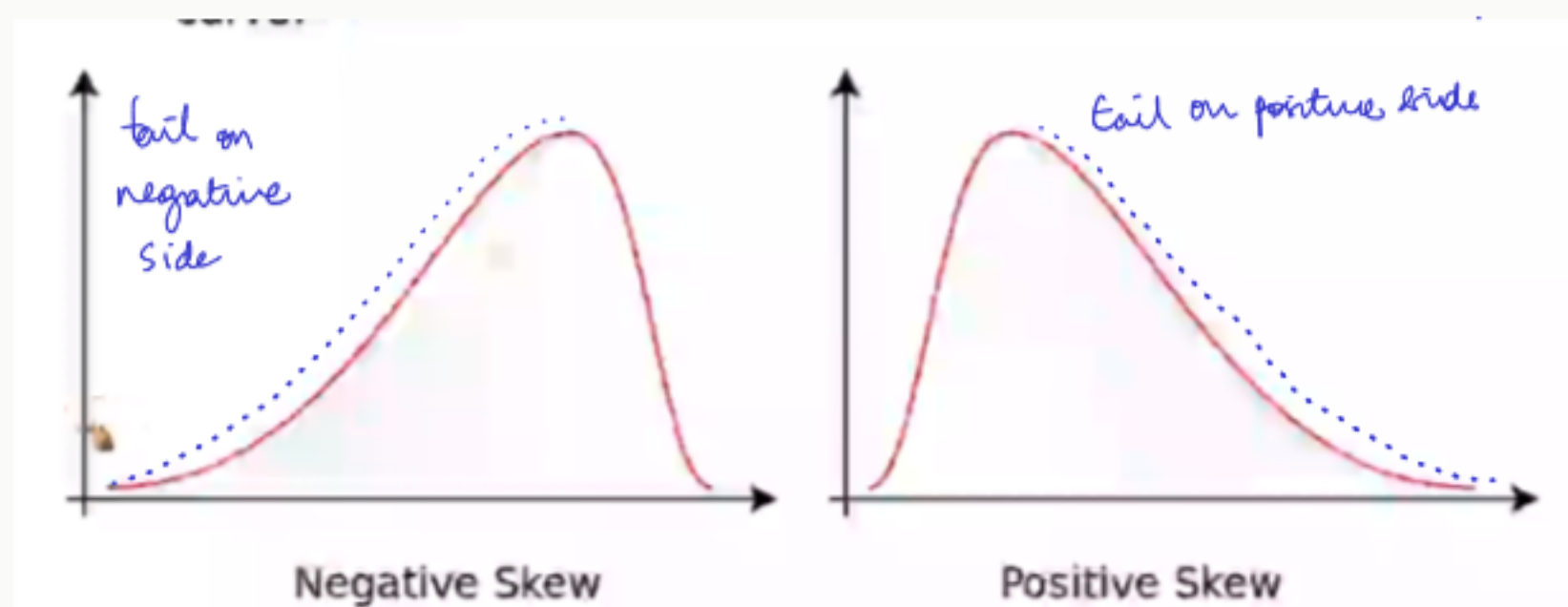
Skewness:-

Skewness is a measure of asymmetry -

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\text{Skewness} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\left[ \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{3/2}}$$

$\xrightarrow{\text{if } 2 = \text{variance}}$   
 $\xrightarrow{\text{sample std-deviation}}$





Kurtosis :-

Kurtosis of Gaussian Distribution = 3

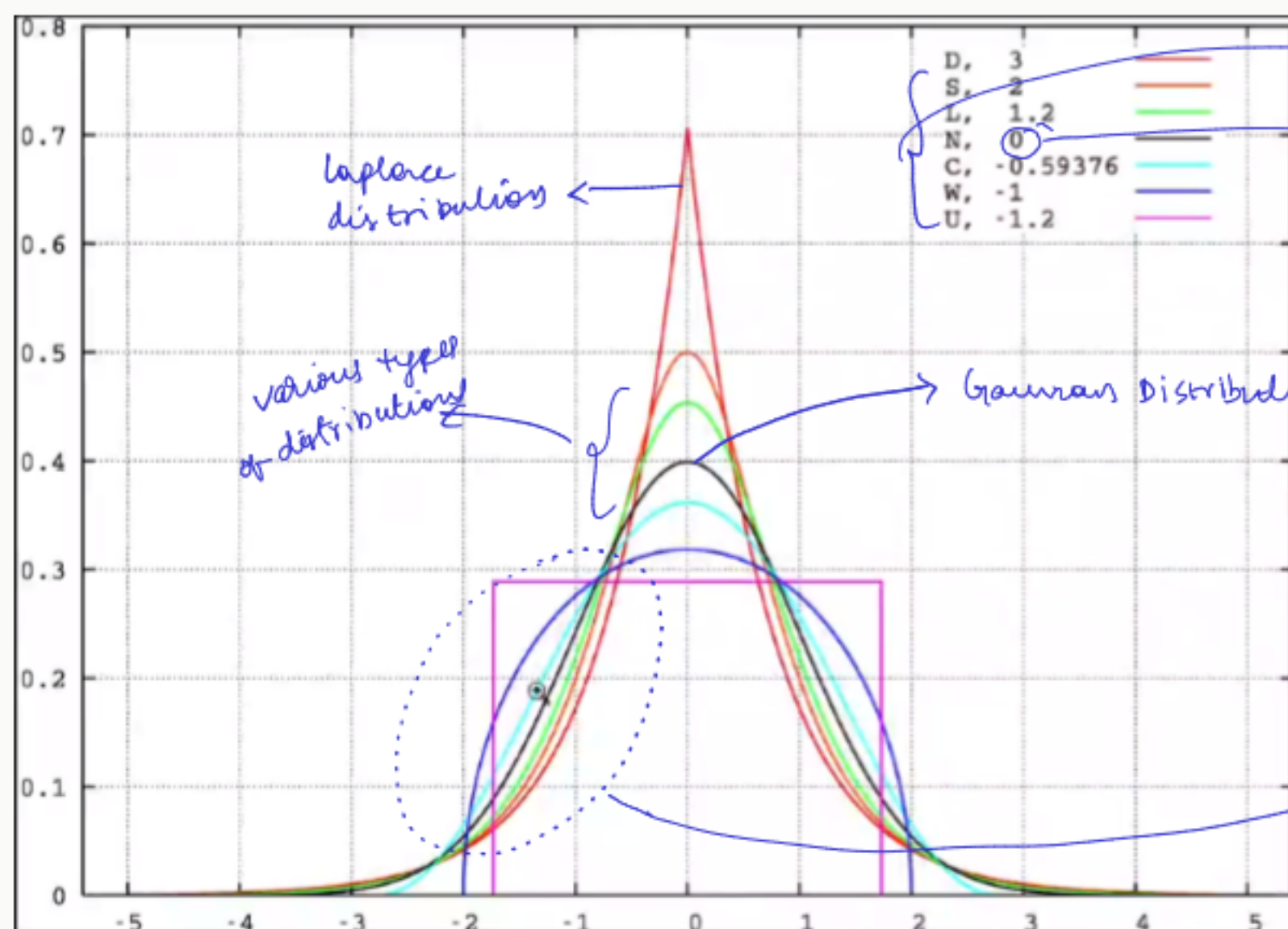
→ Measure of tailedness

excess kurtosis = kurtosis - 3

$$\text{excess kurtosis} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{\left( \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} - 3$$

Variance

→ Kurtosis is not a measure of peakedness. Might look like it -



→ excess kurtosis values.

Gaussian Distribution, excess kurtosis = 0 ⇒ Kurtosis = 3

→ The rate at which the tails touch axis is different

→ Kurtosis tells us if there are outliers are there in the data.

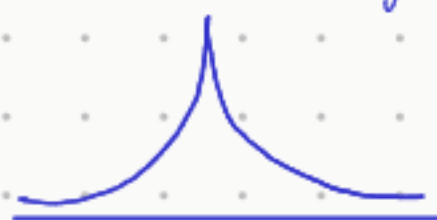
→ Kurtosis is used in analysis of trading, finance etc.

→ Kurtosis measures how much density/weight is in tails compared to middle parts

→ If there are some stocks. Some stocks fall ↓  
Some go up ↑  
Some go up enormously. ↑

To estimate risk after buying some stocks, we model data - It's generally assumed that they follow gaussian distribution for ease of calculations (Pre 2008 financial crisis). As to gaussian distributions, chances of losses are less & has good chance of high profits. But now, instead of assuming the data is a gaussian distribution, we look at the actual data & calculate the kurtosis

If kurtosis is  $> 3$  ⇒



⇒ chances of extreme profits are high & chances of losses are also high.

If kurtosis = 3 ⇒ Follows gaussian distribution

If kurtosis  $< 3$  ⇒



⇒ chances of extreme profits are low & chances of losses are also fairly low.

→ Kurtosis has a degree of 4. So if there is large deviation, then it will lead to higher kurtosis value.

→ High kurtosis is caused by infrequent extreme deviations rather than frequent modestly sized deviations.

→ In order to compare kurtosis b/w two numbers, they need to have the same variance.



## Standard Normal Variate:- (z)

$$Z \sim N(0, 1)$$

$$\mu = 0$$

$$\sigma^2 = 1$$

Let  $X \sim N(\mu, \sigma^2)$

$$X = [x_1, x_2, x_3, \dots, x_{50}]$$

Standardization:-  $x_i' = \frac{x_i - \mu}{\sigma}$

$$\Rightarrow x_i' \sim N(0, 1)$$

Standard Normal Variate.

Given any random variable  $X$ , where  $X \sim N(\mu, \sigma^2)$

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

$$\Rightarrow z \sim N(0, 1)$$

why?:-

① After standardization, PDF becomes



This is also called as basic z-score formula. This basically tells how many  $\sigma$  is  $x$  away from  $\mu$ .

68.2% of the data  
( $\because$  68-95-99.7 rule)

② If we are comparing multiple GDs with different  $\mu$  &  $\sigma^2$ , doing this helps understand better & interpret better

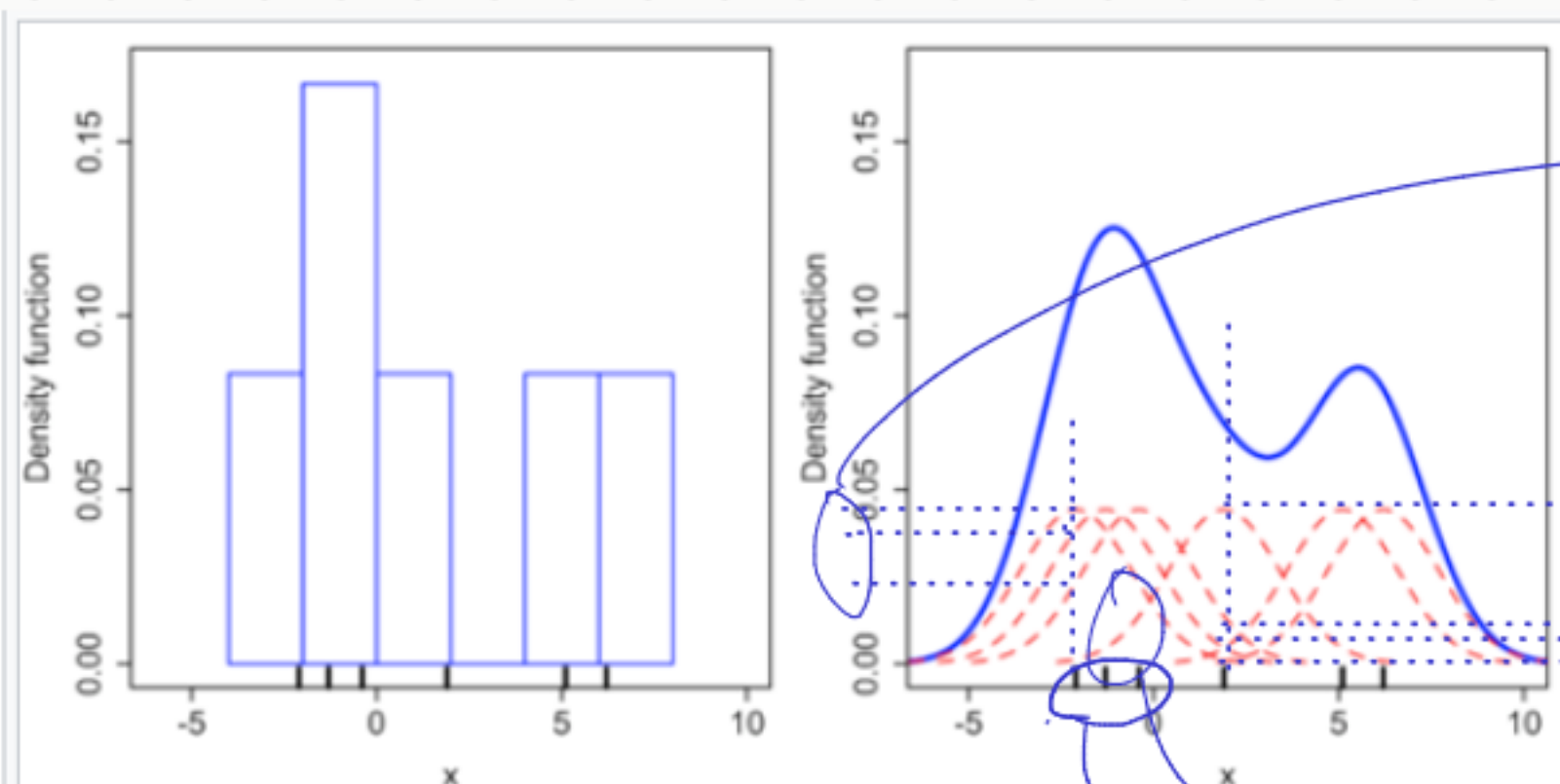
$\rightarrow$  We use StandardScaler for standardization & MinMaxScaler for Normalization.

Two techniques of feature scaling-

Normalization  $\rightarrow x_i = \frac{x_i - x_{\min}}{(x_{\max} - x_{\min})}$

## Kernel Density Estimation:-

$\rightarrow$  Used for smoothing histograms to obtain PDFs.



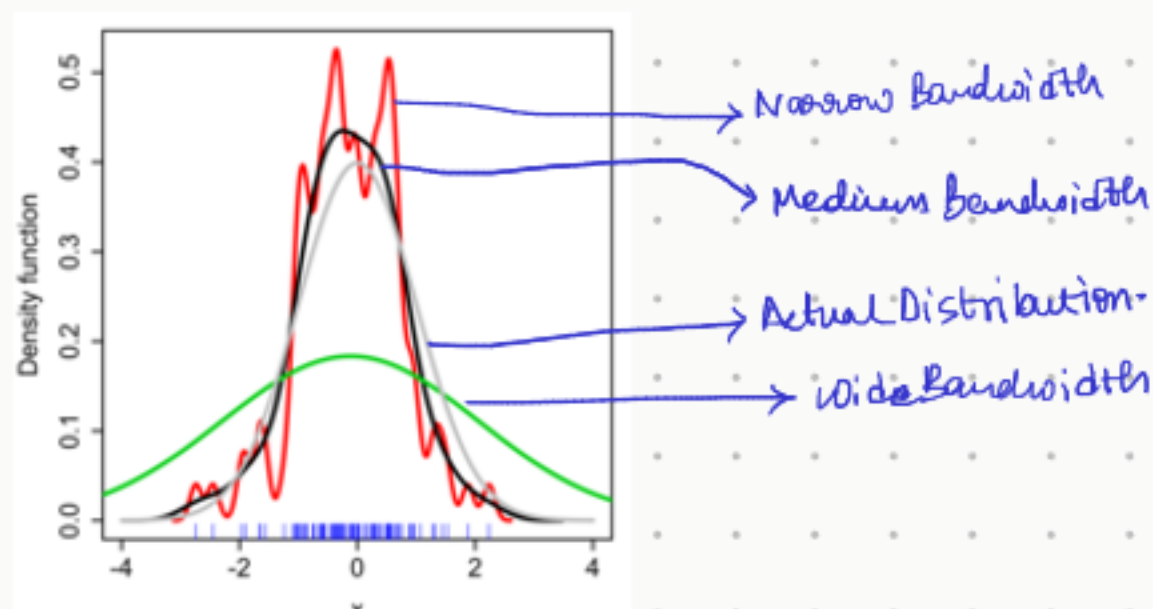
Comparison of the histogram (left) and kernel density estimate (right) constructed using the same data. The six individual kernels are the red dashed curves, the kernel density estimate the blue curves. The data points are the rug plot on the horizontal axis.

gives height at that point

$\rightarrow$  Variance of the Gaussian kernel is known as Bandwidth.

Gaussian kernel.

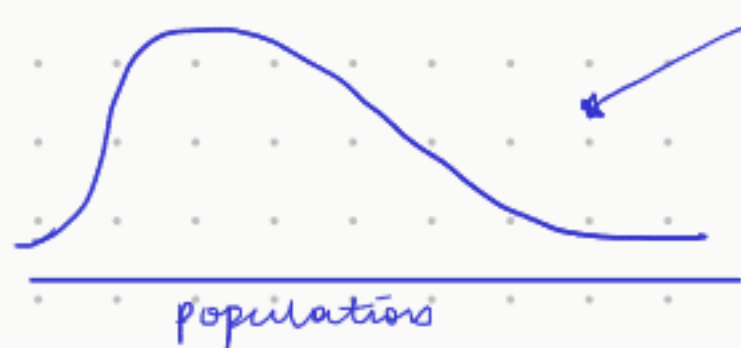
$\rightarrow$  Since there are many points here, the number of means will be high. So the PDF will be high. In normal histogram,





## Sampling Distributions & Central Limit Theorem

Let  $X$  be distributions of incomes over populations [Not necessarily Gaussian]



Let  $S_1$  be random sample of size  $n$  (let  $n=30$ )  
 $n$   $S_2$   $n$   $n$   $n$   
 $n$   $n$   $n$   $n$   $n$

mean =  $\bar{x}_1$  → Sample mean.  
 mean =  $\bar{x}_2$   
 mean =  $\bar{x}_n$

These will also have a distribution

The distribution  $\bar{x}_i$  = Sampling distributions of sample mean.

Central Limit Theorem :- If original distributions ' $X$ ' has finite mean (there can be infinite mean ex: pareto) & variance & 'm' samples are created of size ' $n$ ' whose sample means are  $\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots$  whose distribution is called sampling distribution of sampling mean, central limit theorem states that

$$\bar{x}_i \rightarrow N\left(\mu, \frac{\sigma^2}{n}\right) \text{ as } n \rightarrow \infty$$

(But IRL if  $n \geq 30$ , then it becomes gaussian. Rule of thumb)

mean is same as original or that of distribution

Gaussian Distributions

Why? → By using just  $m \times n$  datapoints, we are able to estimate  $\mu$  &  $\sigma^2$  of any distribution, if we just know that they are finite.

## Quantile Quantile Plot (Q-Q-Plot) :-

Let  $X$  be a random variable

$X: x_1, x_2, x_3, \dots, x_{500}$

Question :- Is  $X$  Gaussian Distributed. Q-Q-Plot helps in identifying other techniques such as statistical testing (KS Testing) etc, also exist.

graphical method

more powerful.

How? :- ① Sort  $X$  & calculate percentile.

$x'_1, x'_2, x'_3, \dots, x'_{500}$   
 $x'_6 \rightarrow 1^{\text{st}} \text{ percentile} \rightarrow x^{(1)}$   
 $x'_{10} \rightarrow 2^{\text{nd}} \text{ percentile} \rightarrow x^{(2)}$   
 $\vdots$   
 $x'_{500} \rightarrow 100^{\text{th}} \text{ percentile} \rightarrow x^{(100)}$

② take  $Y \sim N(0,1)$

$Y: y_1, y_2, y_3, \dots, y_{1000}$

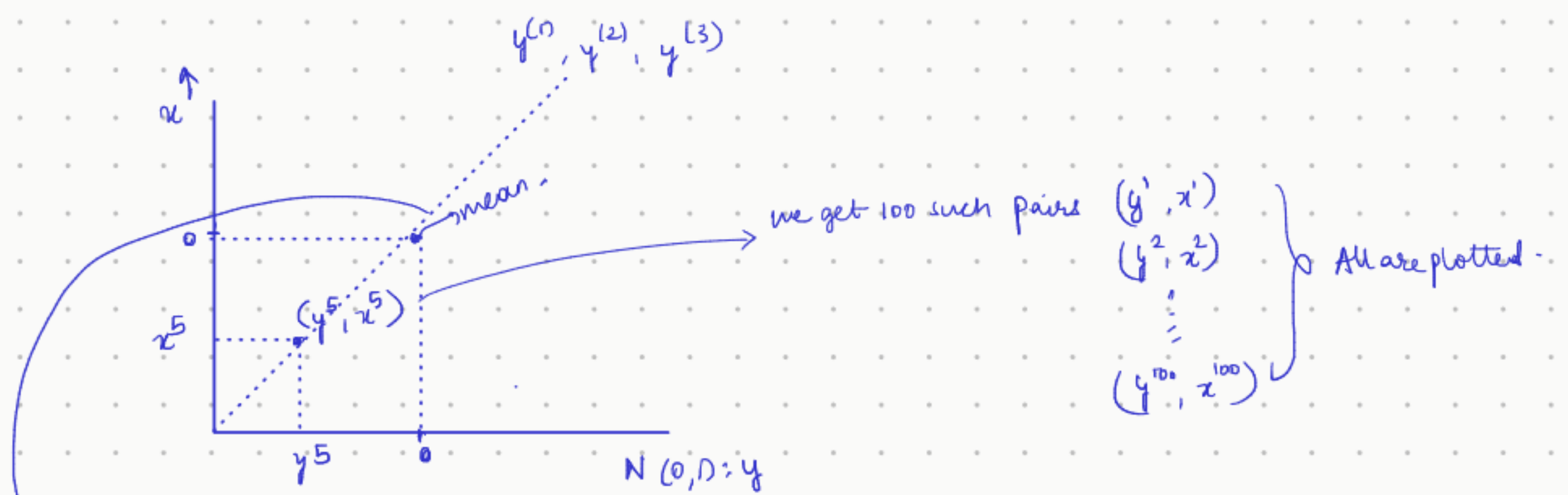
$y'_1, y'_2, y'_3, \dots, y'_{1000}$

$y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(100)}$

These are called as Theoretical Quantiles



③ Plot Q-Q-Plot using  $x^{(1)}, x^{(2)}, x^{(3)}, \dots$



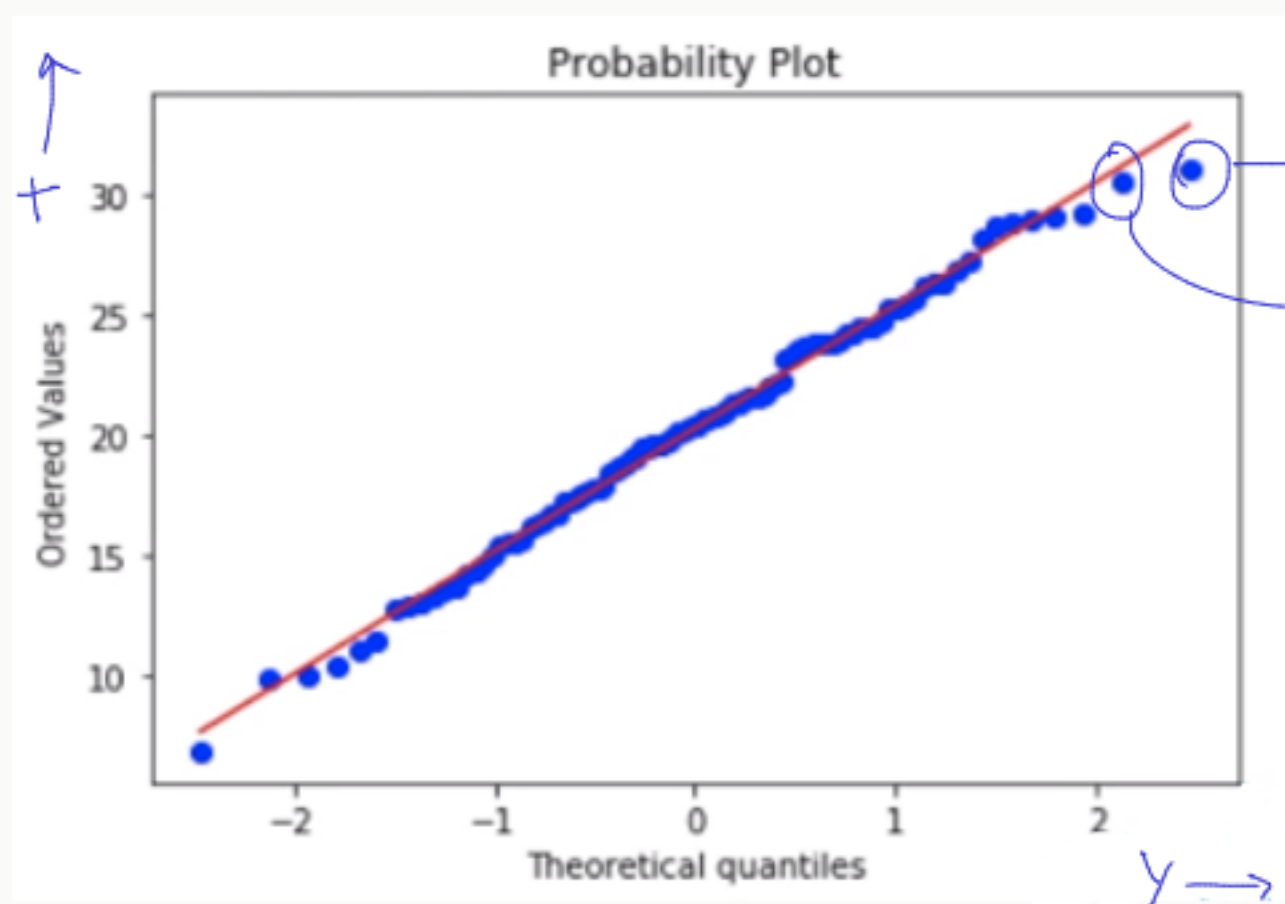
→ If  $(y^i, x^i) \forall i \Rightarrow 100$ , lie roughly on a straight line, then  $y$  &  $x$  have similar distributions.

⇒  $x$  also has gaussian distributions.

→ `stats.probplot()`

→ If size of  $x$  is small, it is hard to interpret Q-Q plot. Won't be a straight line.

→ Another use of Q-Q plot is: given  $x, y$ , does  $x$  &  $y$  have the same distribution?



Generating  $y$  :- (theoretical Quantiles)

```
#Q-Q plot
import numpy as np
import pylab
import scipy.stats as stats

# N(0,1)
std_normal = np.random.normal(loc = 0, scale = 1, size=1000)
```

mean std-dev

Generating  $x$  :-

```
# generate 100 samples from N(20,5)
measurements = np.random.normal(loc = 20, scale = 5, size=100)
# try size=1000

stats.probplot(measurements, dist="norm", plot=pylab)
pylab.show()
```

How & where to use distributions:-

→ All probability concepts are used for EDA most of the time.

→ If data is gaussian Distributed, many assumptions can be made by plotting PDF & CDF

Margin of error:- How much deviation from original numbers is allowed in samples.

Chebyshev's Inequality :-

→ If we know that data is gaussian Distributed, we can apply 68-95-99.7 rule.

→ But if we don't know the distribution, but we know  $\mu, \sigma$  can we make assumptions like this?

$\mu$  (finite)  $\sigma$  (nonzero/finite)



→ Chebyshev's Inequality states that

if  $X$  is a random variable with finite mean  $\mu$  and non zero & finite  $\sigma$

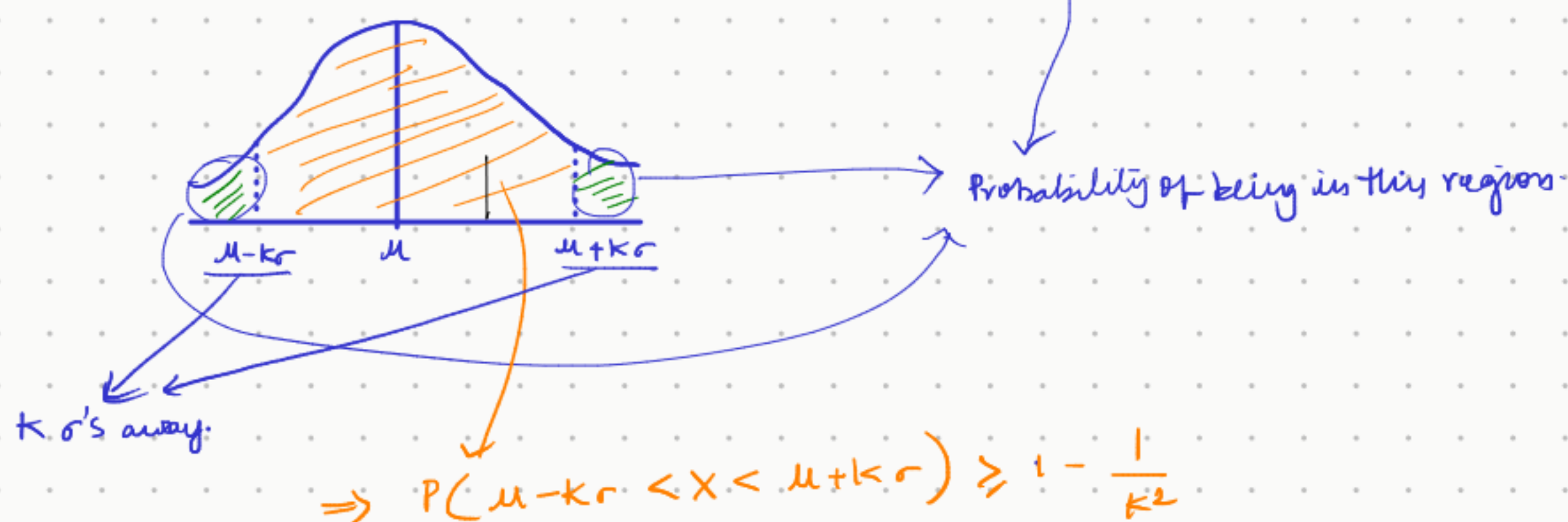
$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$$

$k$  std deviations

68 → 1 std dev away

95 → 2 std dev away

99.7 → 3 std dev away



→ From Chebyshev's inequality we can conclude that 75% of any distributions lies within  $[(\mu - 2\sigma) \text{ \& } (\mu + 2\sigma)]$  and 90% within  $[(\mu - 3\sigma) \text{ and } (\mu + 3\sigma)]$

→ In case the range is not symmetric, there are asymmetric variations of Chebyshev as well.

## Uniform Distributions

Discrete } Just like random variables  
Continuous }

PDF (Probability Density Functions) are drawn for continuous random variables.

PMF (Probability Mass Functions) are drawn for discrete random variables.

→ Probability of getting a value in dice is same for all values. This is called equiprobable.

→ CRV has  $N(\mu, \sigma)$  as parameters.

DRV has  $(a, b)$  as parameters.

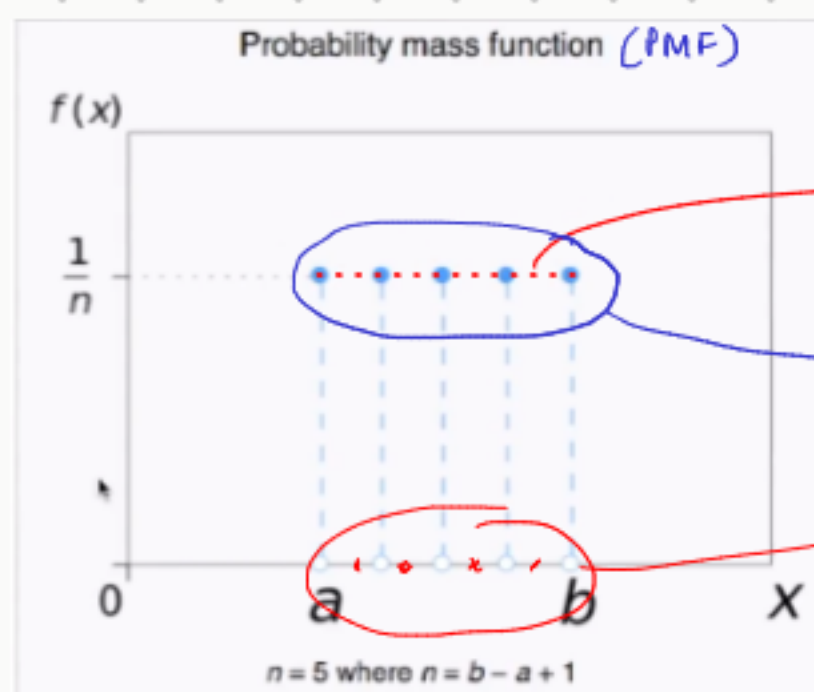
↳ both are numbers.

and

$n = b - a + 1$ . Another formula for calculating  $n = \frac{(b - a + k)}{k}$  where  $k$  is the common difference.

↳ number of outcomes.

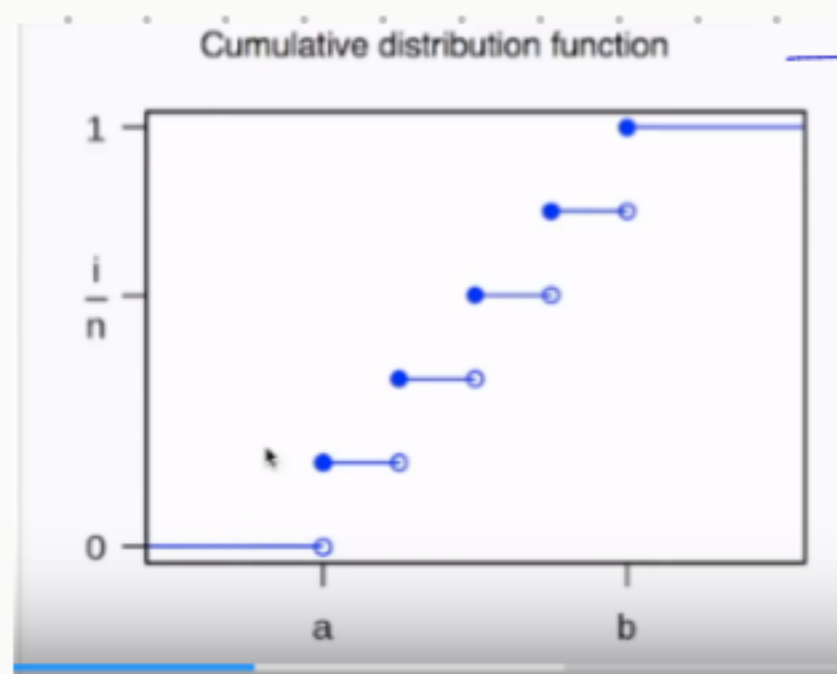
→ In uniform distributions all the variables are equiprobable. i.e.  $\frac{1}{n}$



This line can't be drawn because these points don't exist in uniform distribution

∴ only this does exist in PMF





Props of URV:-

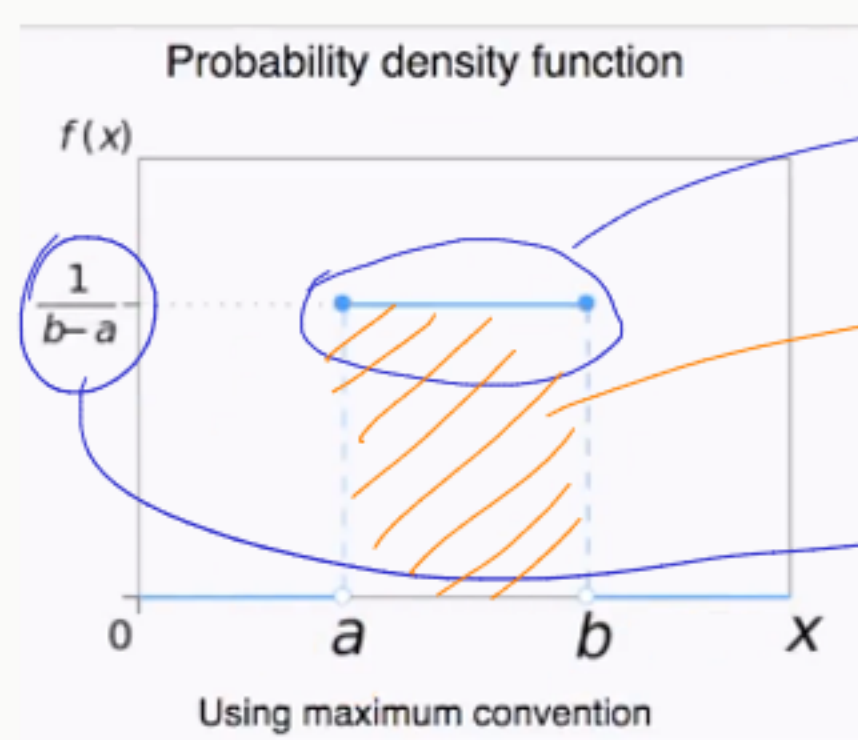
$$\text{Mean} = \frac{a+b}{2}$$

$$\text{Median} = \frac{a+b}{2}$$

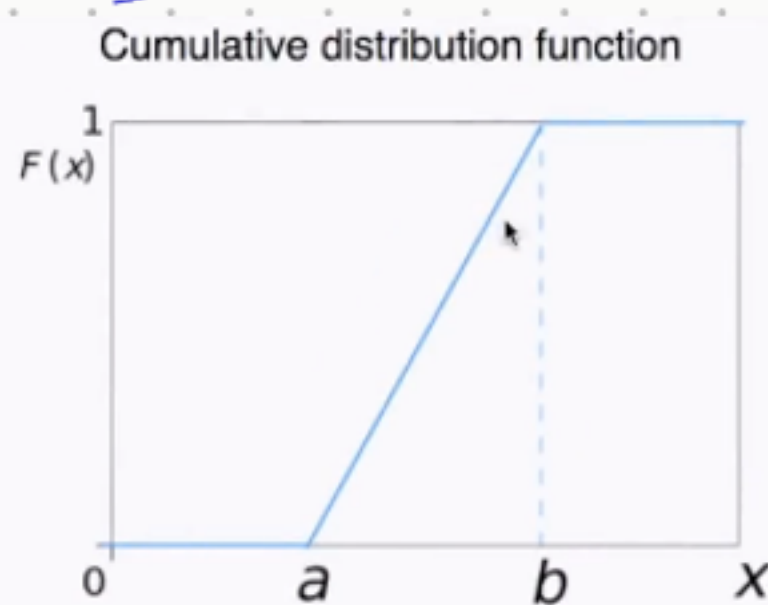
$$\text{Variance} = \frac{(b-a+1)^2 - 1}{12}$$

$$\text{Skewness} = 0$$

Continuous Random Variable:-



CDF:-



Properties of URV:-

Params:  $a, b$

$$\text{PDF} = \begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{for everything else.} \end{cases}$$

$$\text{Mean} = \frac{a+b}{2}$$

$$\text{Median} = \frac{a+b}{2}$$

$$\text{Variance} = \frac{(b-a)^2}{12}$$

PDF  $\rightarrow$  Density of data at the point

PMF  $\rightarrow$  Probability of finding that point there

Random Number generators:-

$\rightarrow$  Most random number generators generate uniform random variables unless explicitly specified.

$\rightarrow$  Python's `random.random()` picks a number b/w 0 & 1 with uniform probability.

ex:-

$\rightarrow$  Picking 30 random values from iris



$n=150$  // Length of this DS

$m=30$  // sample size

$p = m/n$  // 0.2

sample\_d = []

for i in range(0, n):

if random.random()  $\leq p$ :

sample\_d.append(d[i,:])

len(sample\_d)

$\sim 30$

Not exactly 30

random.random() picks value b/w 1 & n with equal probability.  
Chances of getting val  $\leq 0.2 = 20\%$

$\downarrow$   
20% of 150 vals = 30

→ This is an application of continuous uniform random variables.

### Bernoulli & Binomial Distributions:-

→ Discrete Distributions:-

→ Bernoulli has only two outcomes. One has a probability of  $p$  while the other has a prob of  $1-p$

$\downarrow$   
value 1

$\downarrow$   
value 0

ans:-

$X \sim \text{Bernoulli}(p=0.5)$

$$\text{PMF} = \begin{cases} q = (1-p) & \text{for } k=0 \\ p & \text{for } k=1 \end{cases}$$

Mean =  $p$

Variance =  $pq$

### Binomial Random Variable :-

Let  $X$  be tossing coin

$X \sim \text{Bernoulli}(p=0.5)$

Let  $Y = \{ \text{Number of heads when coin is tossed } n \text{ times (let } n \leq 10) \}$

$Y \in \{0, 1, 2, 3, 4, \dots, 10\}$

Now  $Y \sim \text{Binomial}(n, p)$   
 $\swarrow$  number of trials  
 $\searrow$  probability of getting 1

Params:-

$$\text{PMF} = \binom{n}{k} p^k (1-p)^{n-k} \quad \text{for } P(Y=k)$$

Not used often in Machine Learning.