

Random Variables

- * What are Algebraic Variables?

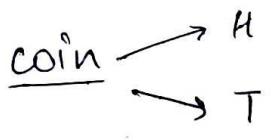
In Algebra a variable, like x , is an unknown value

$$x + 5 = 10$$

$$\boxed{x=5}$$

- * What are Random Variables in stats and Probability?

A Random Variable is a set of possible values from a random experiment.



$$\underbrace{H=1, T=0}_{\text{denote}}$$

$$X = \{1, 0\}$$

↓
sample space

* Random Variable denote with Capital letter and Algebraic variable denote small letter

Types of Random Variable

Discrete
Random Variable

$$\{H, T\}$$

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

Continuous
Random Variable

$$x = (0, 10)$$

Range of Value or
any value like 8.5, 2.4...

Probability Distribution

- * What are Probability Distributions?

A probability distribution is a list of all of the possible outcomes of a random variable along with their corresponding probability value.

coin Toss	1(H)	0(T)
Probab	1/2	1/2

Dice	1	2	3	4	5	6
Prob	1/6	1/6	1/6	1/6	1/6	1/6

2 Dice Rolled

	1	2	3	4	5	6
1	2	3	4	5	6	7
2	3	4	5	6	7	8
3	4	5	6	7	8	9
4	5	6	7	8	9	10
5	6	7	8	9	10	11
6	7	8	9	10	11	12

	<u>Prob distribution</u>	
2 →	$1/36$	9 → $4/36$
3 →	$2/36$	10 → $3/36$
4 →	$3/36$	11 → $2/36$
5 →	$4/36$	
6 →	$5/36$	12 → $1/36$
7 →	$6/36$	
8 →	$5/36$	

Problems with Distribution

In many scenarios, the number of outcomes can be much larger and hence a table would be tedious to write down. Worse still, the number of possible outcomes could be infinite, in which case, good luck writing a table for that.

Example:- Height of people, Rolling 10 dice together

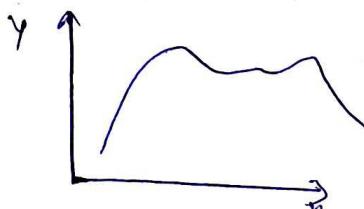
Solution - Function?

What if we use a mathematical function to model the relationship between outcome and probability?

$x \rightarrow$ outcome	1	2	3	4	5	6
$y \rightarrow$ Prob	$1/6$	$1/6$	$1/6$	$1/6$	$1/6$	$1/6$

$$y = f(x)$$

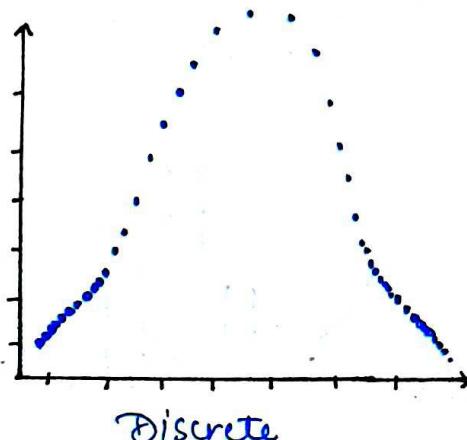
Prob dist func



Note - A lot of time Probability Distribution and probability distribution functions are used interchangeably. ↳ function

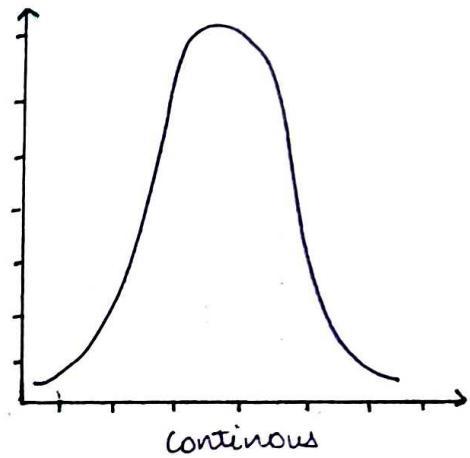
Types of Probability Distribution (PDF)

Discrete Random Variable



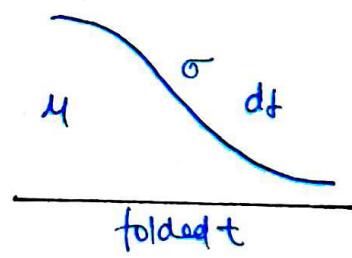
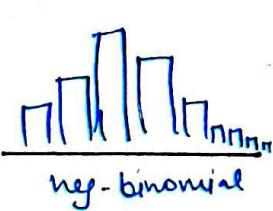
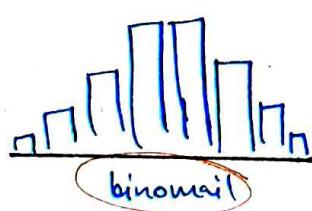
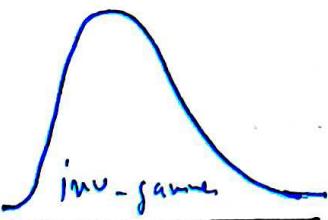
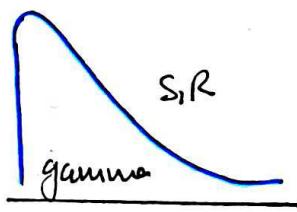
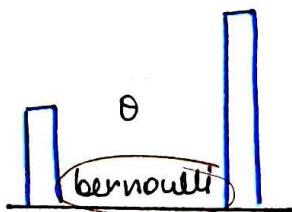
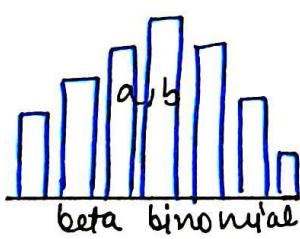
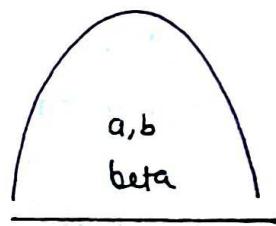
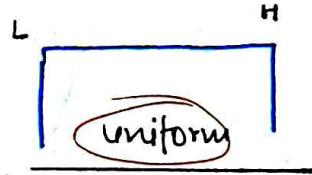
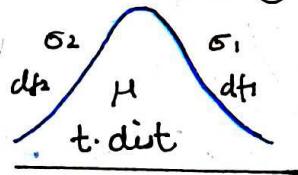
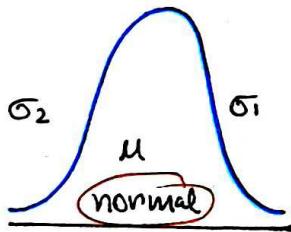
Discrete

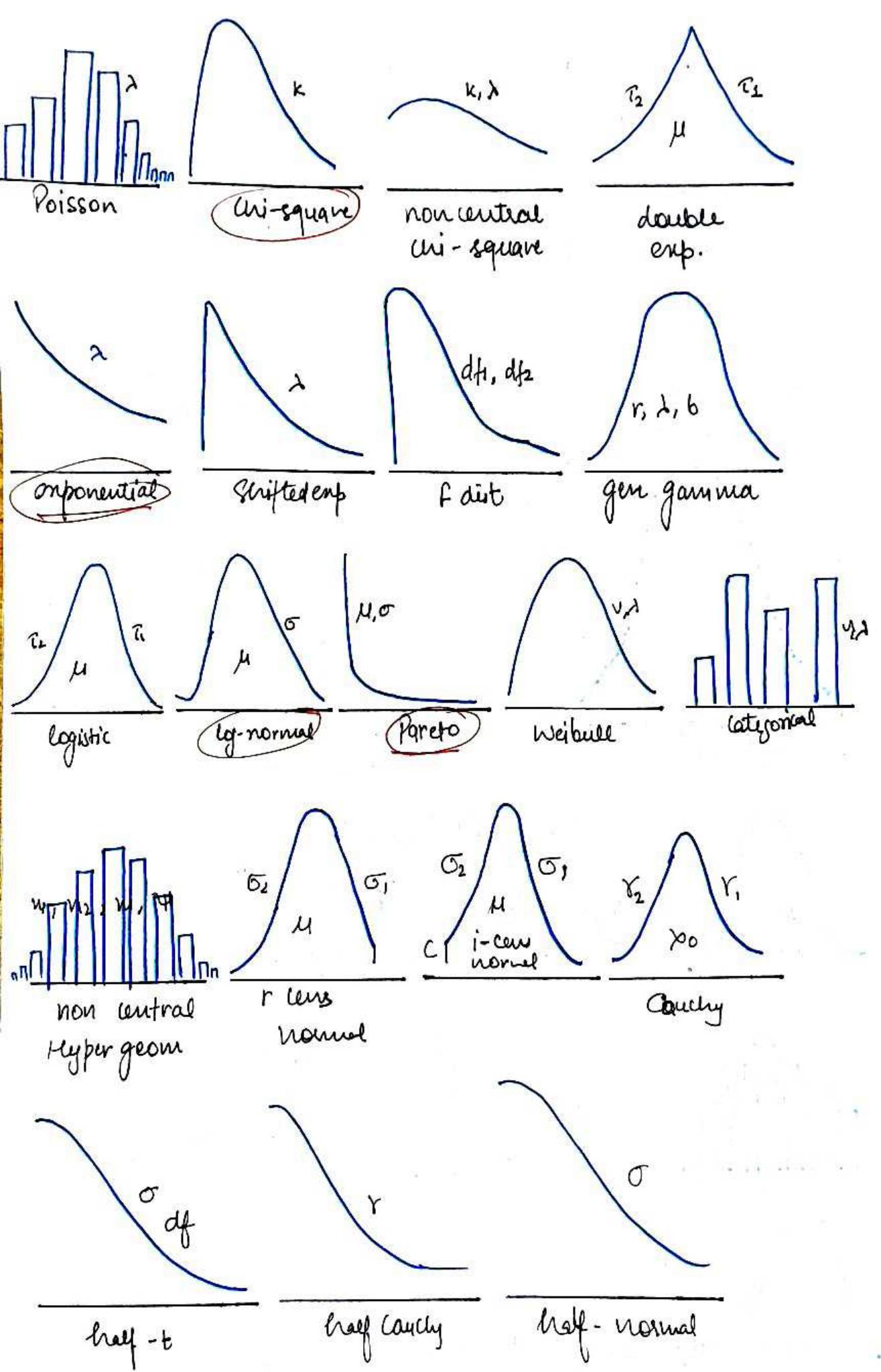
Continuous Random Variable



Continuous

Probability Distribution





* Why are Probability Distribution important?

→ Gives an idea about the shape/distribution of the data.

→ And if our data follows a famous distribution then we automatically know a lot about the data.

A note on Parameter (PDF)

Parameter in probability distribution are numerical values that determine the shape, location, and scale of the distribution.

Different probability distribution have different sets of parameters that determine their shape and character. Statistics and understanding these parameters essential in statistical analysis and inference.

Probability Distribution Functions

A probability distribution function is a mathematical function that describes the probability of obtaining different values of a random variable in a particular probability distribution.

$$y = f(x)$$

Probability distribution functions

Probability Mass function (PMF)

Creating Probability Distribution function of rolling Dice (Discrete Random Variable) is called PMF.

Probability Density function (PDF)

Creating probability dist. function of continuous Random Variable is called PDF.

Cumulative Distribution Function (CDF)

PDF \rightarrow CDF

PMF \rightarrow CDF

Probability Mass function (PMF)

It is a mathematical function that describes the probability distribution of a discrete random variable.

The PMF of a discrete random variable assigns a probability to each possible value of the random variable. The probabilities assigned by the PMF must satisfy two conditions:

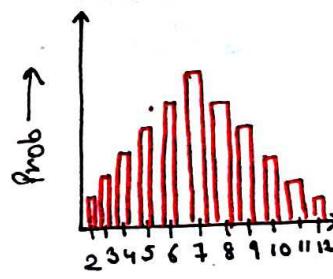
- The probability assigned to each value must be non-negative i.e., greater than or equal to zero)

One die rolled

$$Y = \begin{cases} \frac{1}{6} & \text{if } x \in \{1, 2, 3, 4, 5, 6\} \\ 0 & \text{otherwise} \end{cases}$$

Two dice rolled

$$Y = \begin{cases} \frac{1}{36} & \text{if } x \in \{2, 12\} \\ \frac{2}{36} & \text{if } x \in \{3, 11\} \\ \vdots & \\ 0 & \text{otherwise} \end{cases}$$



- * we can find the PMF and draw the graphs.

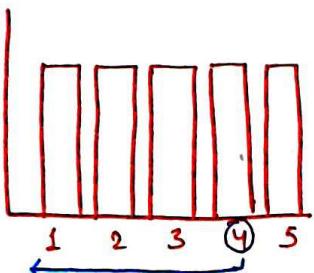
Cumulative Distribution function (CDF) of PMF

The cumulative distribution function (CDF) $F(x)$ describes the prob that a random variable X with a given prob dist. will be found at a value less than or equal to x

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

Let a dice is rolled

$$\{1, 2, 3, 4, 5\}$$



In PMF if $x=4$

then Prob of 4.

but in CDF if $x=4$

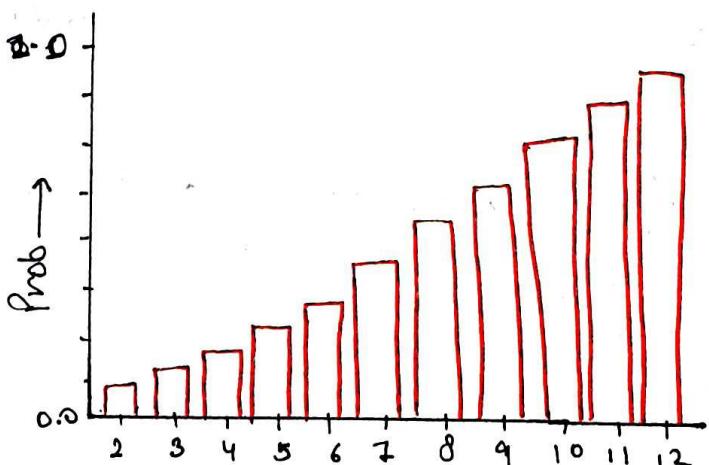
then find prob is equal to and less than 4.

$$f(x \leq 4) = f(x=4) + f(x=3) + f(x=2) \\ + f(x=1)$$

$$= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{4}{6}$$

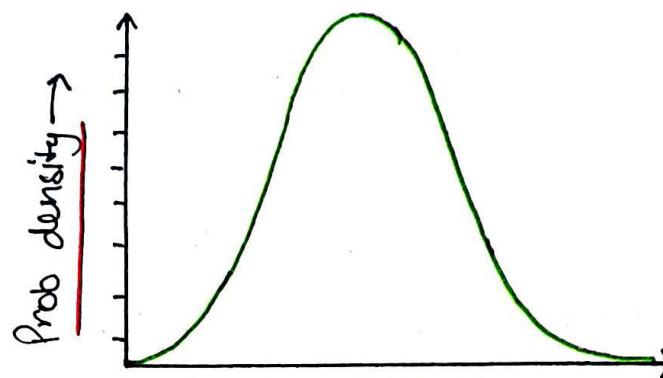
Two dice are rolled

	1	2	3	4	5	6
1	2	3	4	5	6	7
2	3	4	5	6	7	8
3	4	5	6	7	8	9
4	5	6	7	8	9	10
5	6	7	8	9	10	11
6	7	8	9	10	11	12



Probability Density Function

It is a mathematical function describe the prob dist. of a continuous random variable.



why Prob. density and why no probability?

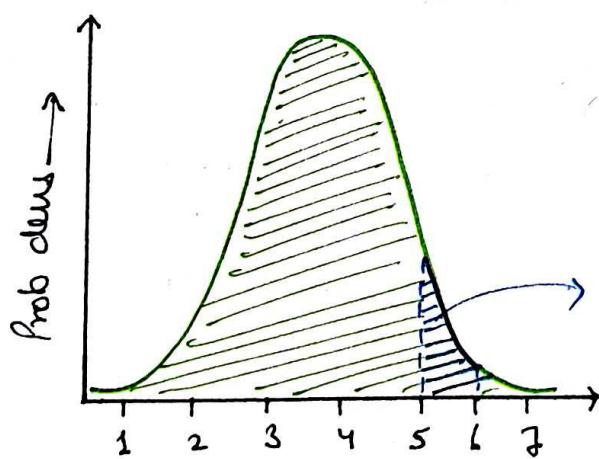
let assume Gpa = 7.912 ..

we cannot able to find prob at a particular level like 7.912... something. Prob. of this may be near to zero or zero. We can dealing with infinite no. of values. That's why we can not use prob at y-axis.

What does the area of this graph represent?

PDF helps to find the prob betn the range like (5 to 6).

$$P(0 \leq x \leq 10) = 1 \Rightarrow \text{Total area}$$



Find the area is the prob.

$$\int_5^6 f(x) dx$$

How is Graph calculated?

Density Estimation

Density estimation is a statistical technique used to estimate the prob density func (PDF) of a random variable based on a set of observation or data. In simpler terms, it involves estimating the underlying distribution of a set of data points.

Density estimation can be used for a variety of purposes, such as hypothesis testing, data analysis, and data visualization. It is particular useful in areas such as machine learning, where it is often used to estimate the probability distribution of input data or to model the likelihood of certain events or outcomes.

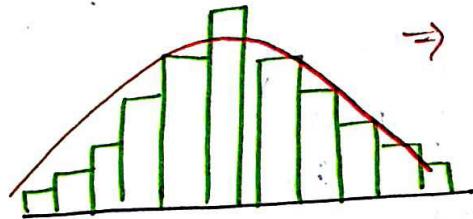
There are various methods for density estimation, including parametric and non-parametric approaches. Parametric methods assume that the data follows a specific probability distribution (such as a normal distribution), while non-parametric methods do not make any assumption about the distribution and instead estimate it directly from the data.

Commonly used techniques for density estimation include kernel density estimation (KDE), histogram estimation and Gaussian mixture models (GMMs). The choice of method depends on the specific characteristics of the data and the intended use of the density estimate.

Parametric Density function

Parametric density estimation is a method of estimating the prob density function (PDF) of random variable by assuming that the underlying distribution belongs to specific parametric family of probability distribution such as the normal, exponential or poisson distribution.

- I took 1000 data $4.3, 9.1, 6.5, \dots, 1000$
- and plot a histogram



→ It seems like Normal distribution
So, we use Pdf of normal dist.
$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{\sigma-\mu}{\sigma})^2}$$

we have to find mean and standard dev.

After plotting the graph and graph does not match with any other graph then we use Non parametric Density Estimation (KDE).

Non-Parametric Density Estimation

But sometimes the distribution is not clear or its not one of the famous distributions.
Non parametric density estimation is a statistical technique used to estimate the probability density function of a random variable without making any assumption about the underlying distribution. It is also referred to as non-parametric density estimation because it does not require the use of a predefined

Prob distribution function, as opposed to parametric method such as the Gaussian distribution.

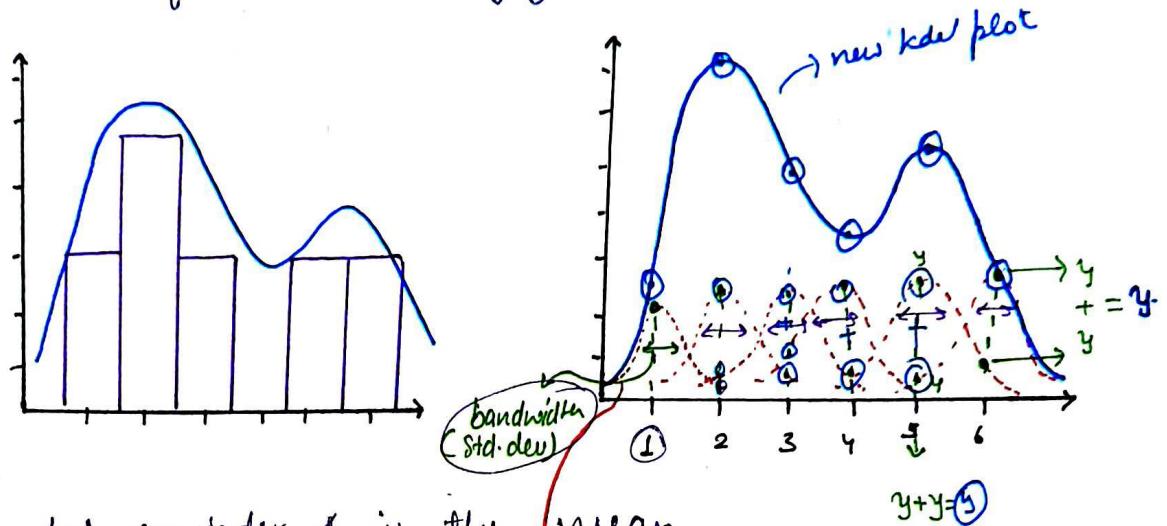
The non-parametric density estimation technique involves constructing an estimate of the prob density func. using the available data. This is typically done by creating a kernel density estimate.

Non-parametric density estimation has several advantages over parametric density estimation.

One of the main advantages is that it does not require the assumption of a specific distribution, which allows for more flexible and accurate estimation in situations where the underlying distribution is unknown or complex. However, non-parametric density estimation can be computationally intensive and may require more data to achieve accurate estimate compared to parameter methods.

Kernel Density Estimate (KDE)

The KDE technique involves using a kernel function to smooth out the data and create a continuous estimate of the underlying density functions.



Step 1:- let consider s is the mean
and create normal dist on s and same as
for $2, 3, 4, 5, 6$

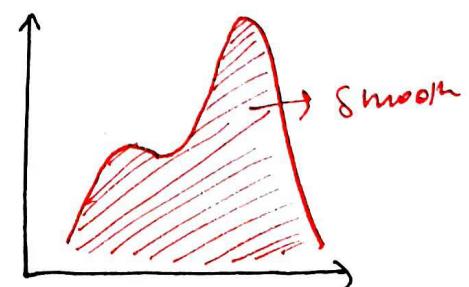
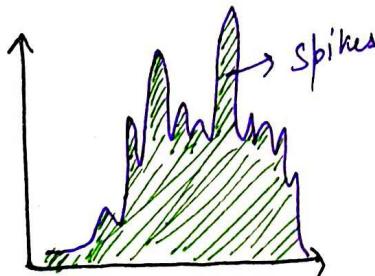
Step 2:- Add the value of y which is intercept
to normal dist of s .

Step 3:- Make new kde plot with sum of y and original x .

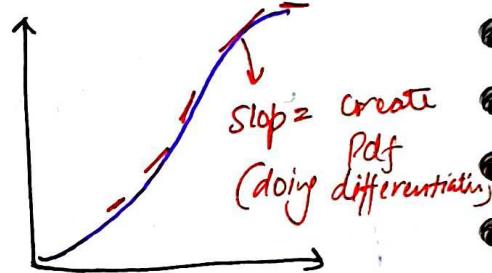
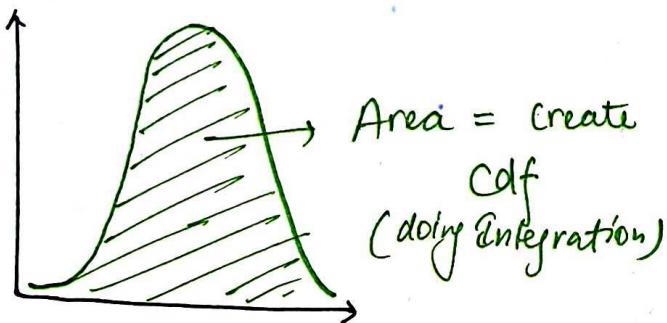
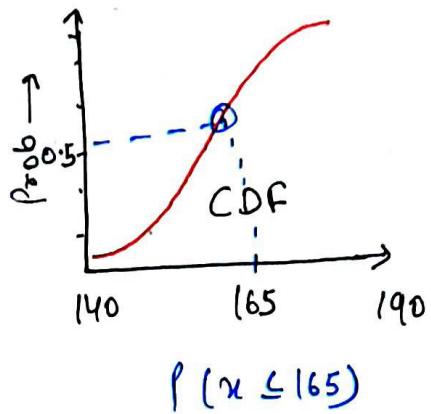
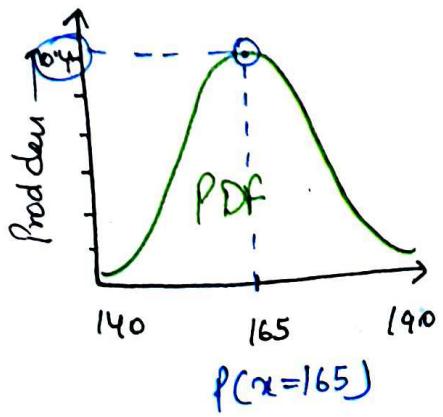
Step 4:- Band width is the hyperparameter.
Kernel = Gaussian mostly used.

bandwidth ≈ 0

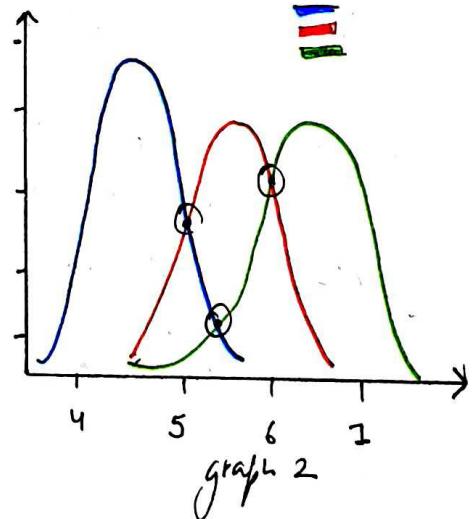
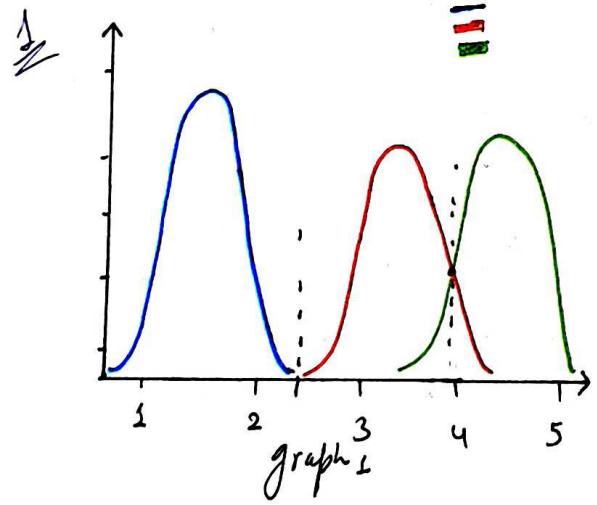
bandwidth $\uparrow \uparrow$



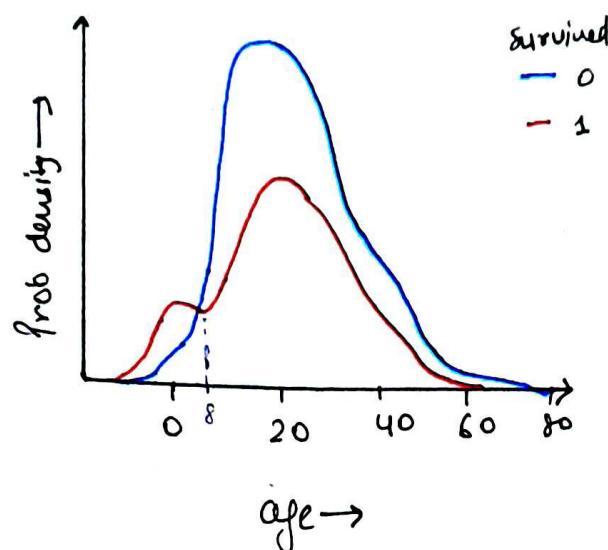
Cumulative Density function → PDF



How to use PDF in Data Science?

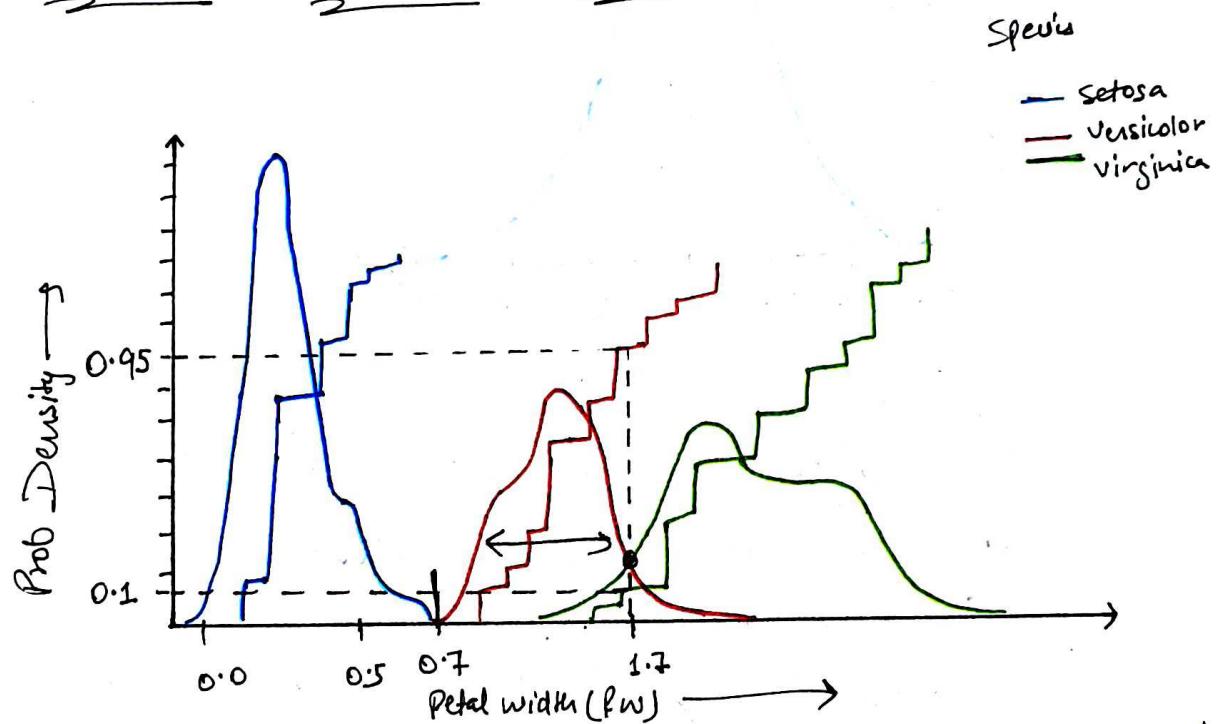


PDF helps in the features selection. In Graph 1 we can easily differentiate betⁿ the categories. but in Graph 1 we can see the kde plot near to each other and difficult to differentiate.



In this graph,
0-8 age have higher chance to survive
but after 8-80 have higher chance for prob density to not survive.

How to use CDF in data science?



If $0.7 < \text{Pw} < 1.7$ is versicolor? \Rightarrow 95% of versicolor under 1.7 and between 0.7 and 1.7.
5% can be mistake

$0.1 \rightarrow 10\%$ of ~~data~~ virginica data in versicolor data \Rightarrow 80 $\left[\begin{array}{l} \text{Pw} > 1.7 \\ \downarrow 90\% \text{ of data} \end{array} \right]$ virginica

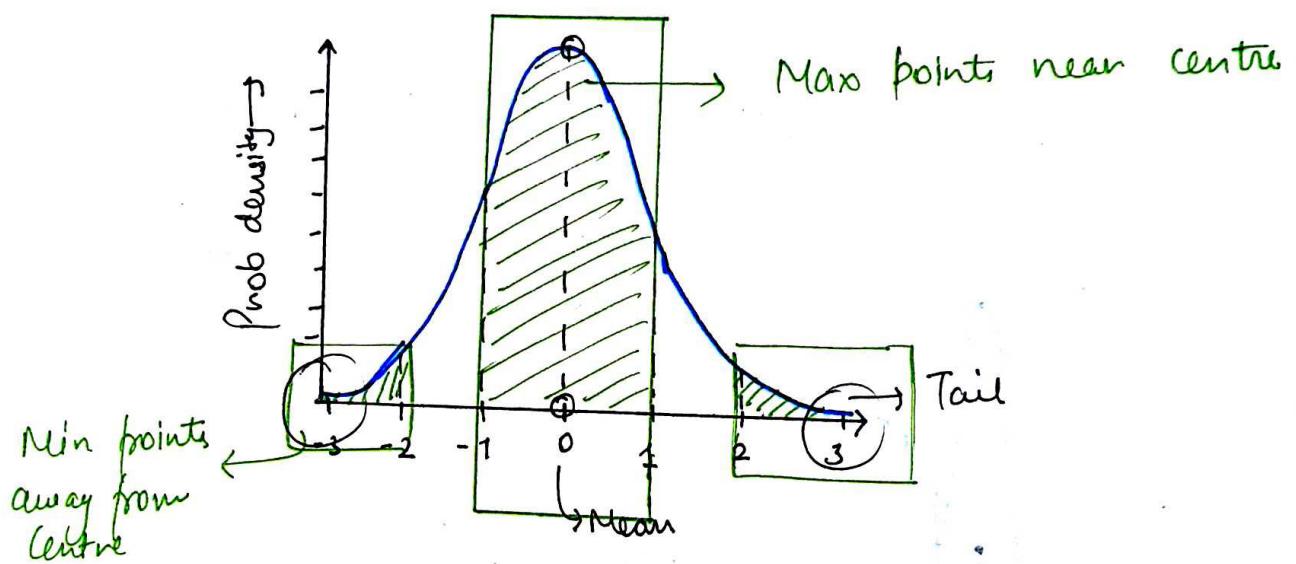
10% mistake can be

* In decision making PDF ke through kya use ek quantitative label do skta hu (95% right 5% wrong)

Normal Distribution

1. What is Normal Distribution?

Normal Distribution, also known as Gaussian distribution is a prob distribution that is commonly used in statistical analysis. It is a continuous prob distribution that is symmetrical around the mean, with a bell-shaped curve.



The normal distribution is characterized by two parameters: the mean (μ) and the standard deviation (σ). The mean represent the centre of the distribution, while the standard deviation represent the spread of the distribution.

Why is it so important?

Commonality in Nature Many natural phenomena follows a normal distribution, such as the height of people, the weights of objects, the IQ scores of a population and many more. Thus, the normal distribution provides a

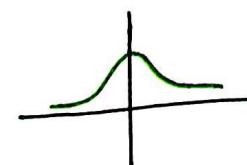
convenient way to model and analyse such data.

PDF eqn of Normal distⁿ

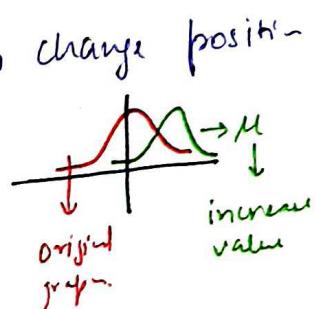
$$Y = f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

Eqn in detail:

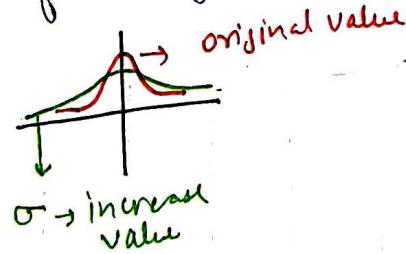
$$y = e^{-x^2} \Rightarrow$$



$$y = e^{-(x-\mu)^2} \Rightarrow \text{Normal distribution change position left to right}$$

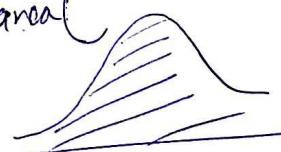


$$y = e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \text{help to spread of the graph}$$



$$Y = e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2} \rightarrow \text{Integration} = \frac{\sigma\sqrt{2\pi}}{\text{Area}} \text{ for cancelling this divide by the area}$$

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



Standard Normal Variate (z)

or

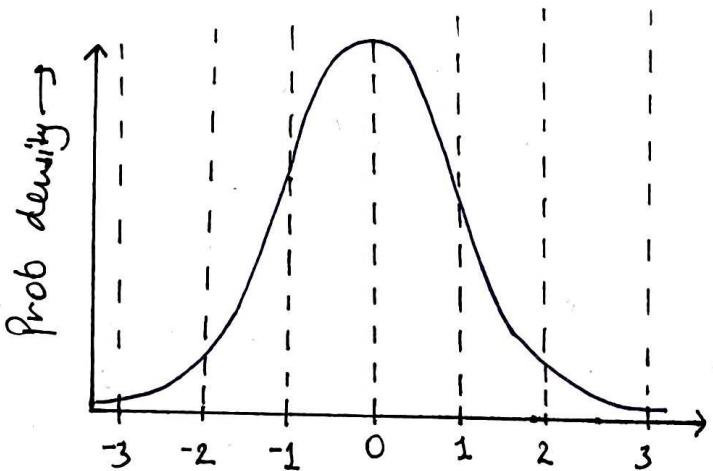
Standard Normal Distribution

* what is standard Normal variate

A Standard Normal Variate (z) is a standardized form of the normal distribution with mean = 0 and standard deviation = 1.

$$x \sim N(\mu, \sigma) \rightarrow \text{Normal dist}^n$$

$$z \sim N(0, 1) \rightarrow \text{Standard Normal Variate}$$

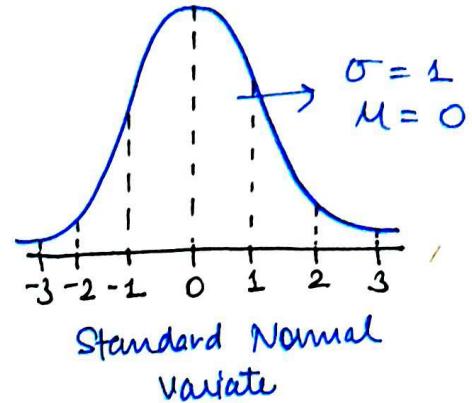
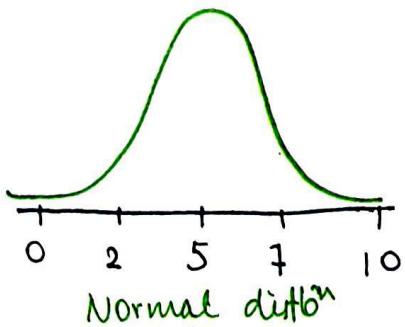


Standardizing a normal distribution allows us to compare different distributions with each other, and to calculate probabilities using standardized tables or software.

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

Standard
normal variate

* How to transform a normal distribution to Standard Normal variate.



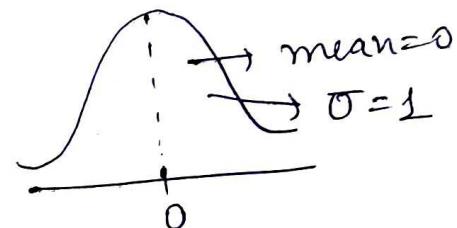
Age

27	$Z = \frac{27-\mu}{\sigma}$
25	$Z = \frac{25-\mu}{\sigma}$
61	$Z = \frac{61-\mu}{\sigma}$
73	$Z = \frac{73-\mu}{\sigma}$

$Z = \frac{61-\mu}{\sigma}$

$Z = \frac{73-\mu}{\sigma}$

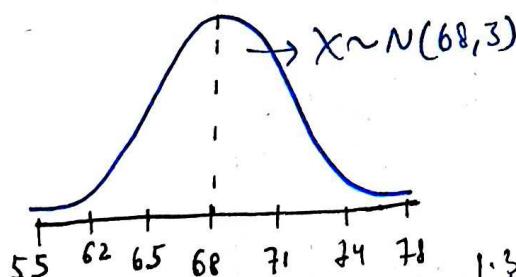
Now Plot the graph



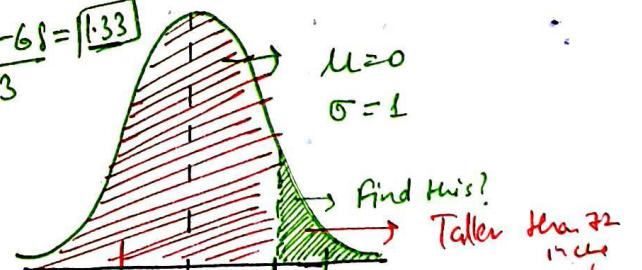
Suppose the heights of adult males in a certain population follow a normal distⁿ with a mean of 68 inches and a standard deviation of 3 inches. What is the prob that a randomly selected adult male from this population is taller than 72 inches?

* 72 inches or use kam
hone ki prob = 90.8

Given → population \Rightarrow Normal distⁿ



$$Z = \frac{72-68}{3} = 1.33$$



using Z table \rightarrow 1.3 find on the table and find 0.3 and the norm is 0.90824

$$\begin{aligned} &0.90824 \\ &(90.81\%) \downarrow \text{prob} \\ &72 \text{ inch} \rightarrow 1.33 \end{aligned}$$

$$\begin{aligned} &\text{Total area - find area} \\ &100 - 90.8 \\ &\approx 9.2\% \quad \swarrow \end{aligned}$$

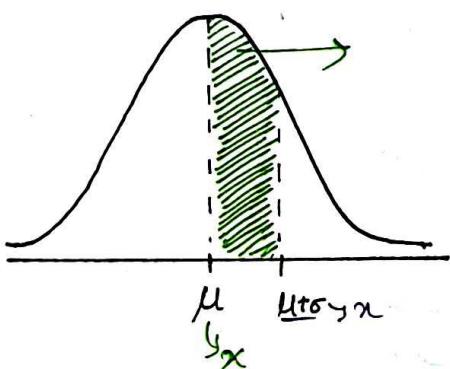
* What are Z-tables?

A Z-table tells you the area underneath a normal distribution curve, to the left of the Z-score.

Check table! - <https://www.ztable.net/>

- * For a Normal Distribution $X \sim (\mu, \sigma)$ what percent of population lie betⁿ mean and 1 std, 2 std, 3 std?

$$X \sim N(\mu, \sigma)$$



$$\text{SNV}(Z) : Z = \frac{X - \mu}{\sigma}$$

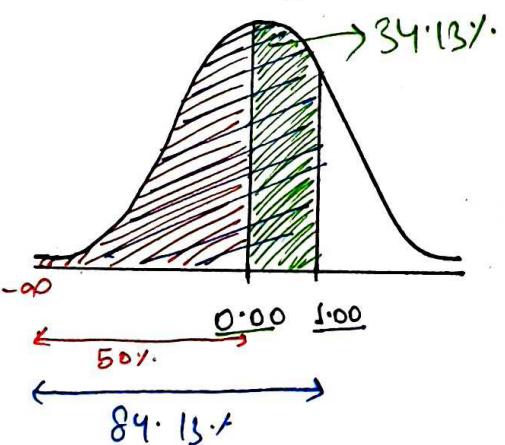
$$Z = \frac{\mu - \mu}{\sigma} = 0.00$$

$$\text{SNV}(Z) Z = \frac{\mu + \sigma - \mu}{\sigma} = 1.00$$

Z table of 0.00 = 50%.

Z table of 1.00 = 84.13%.

$$84.13 - 50 = 34.13$$



$$Z = \frac{\mu + 2\sigma - \mu}{\sigma} = 2$$

Z table of 2 = 90.972

$$97.1 - 84.13 = 13.57$$

