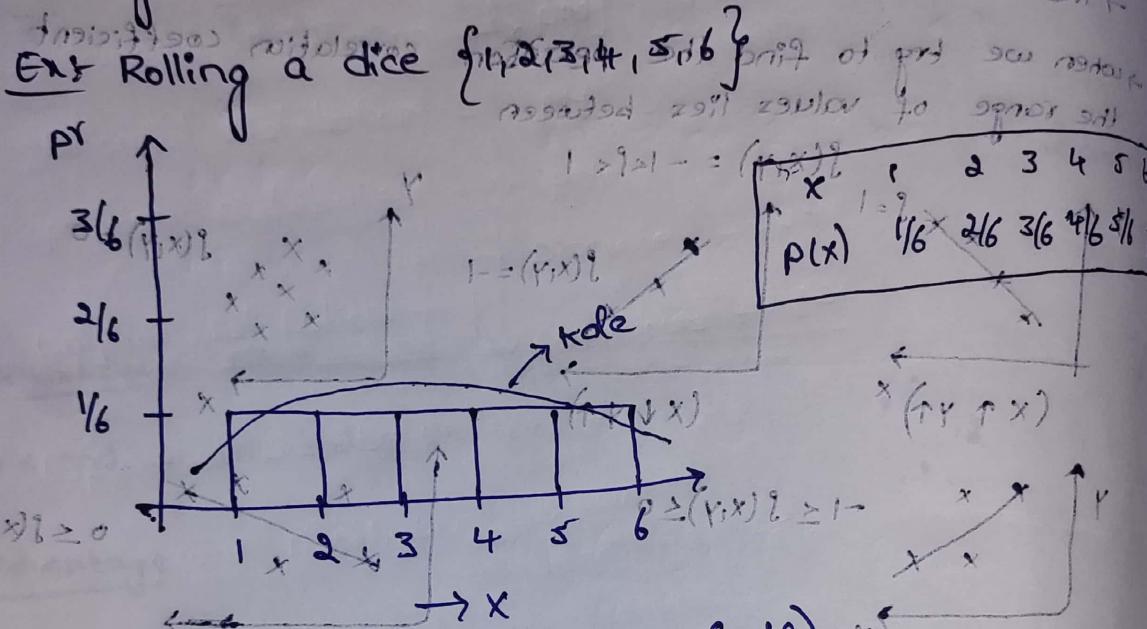
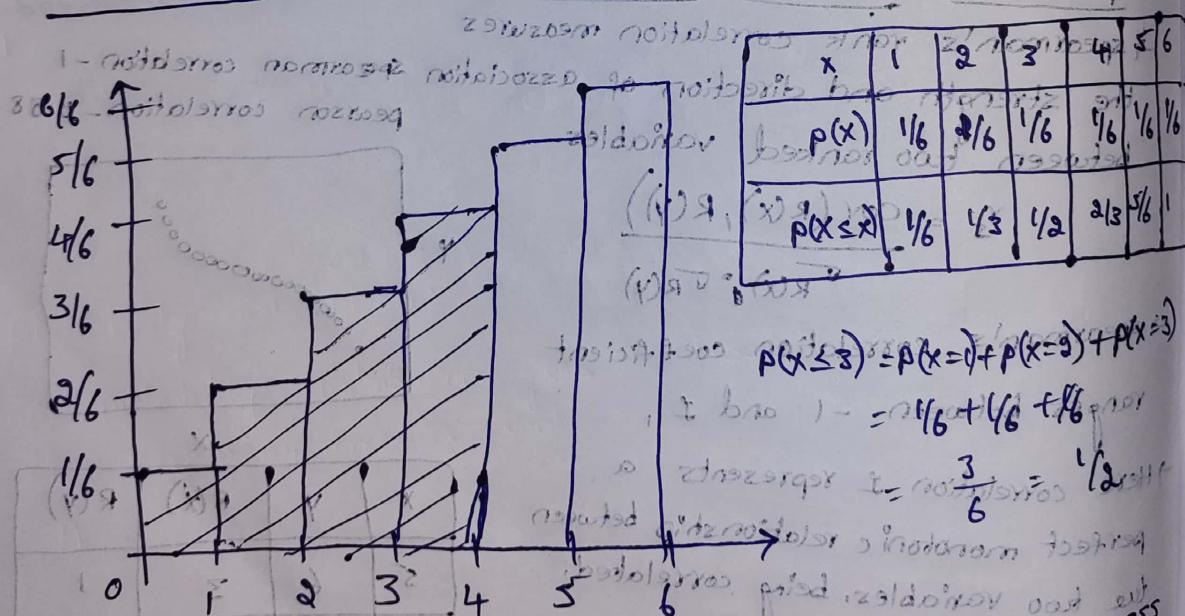


## probability distribution functions

probability mass function:  
 → A probability mass function is a function that gives the probability that a discrete random variable is exactly equal to some value.



cumulative distribution function (cdf)



→ The cumulative distribution function (cdf) for a probability mass function (pmf) is a function that describes the probability of a discrete random variable taking a value less than or equal to a specific value.

→ cdf at a given value of x is equal to sum of all possible outcomes up to and including a certain value.

## 2) probability density function:-

In probability theory, the probability density function (pdf) is a function that describes the probability distribution of a continuous random variable. It gives the probability that a random variable  $x$  takes on a value in a certain interval denoted by  $p(a \leq x \leq b)$ .

Fig: pdf

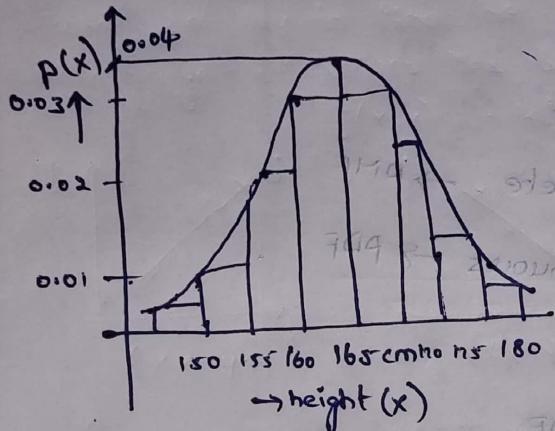
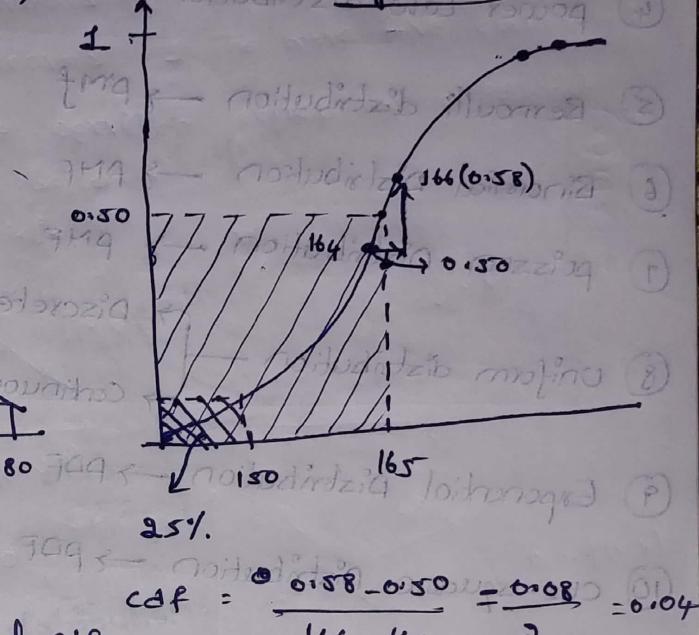


Fig: CDF



Relation between pdf and cdf:

- Here the cdf at any given point  $x$  is equal to the integral of the pdf from negative infinity upto  $x$ .
- In other words, the cdf gives the probability that a random variable  $x$  is less than or equal to a certain value  $x$ .
- while pdf gives the probability density of  $x$  at that value.
- The pdf can be obtained from the cdf by differentiating the CDF.

$$PDF(x) = \frac{d}{dx} CDF(x)$$

Ex:- The pdf is used find the probability of  $x$  falling within a certain range of values

- The probability of students being present in 160 - 170 cm tall. tall is ~~a class of students~~.

## Probability Distributions

- ① Normal / Gaussian Distribution  $\rightarrow$  PDF
- ② Standard Normal distribution  $\rightarrow$  PDF or Z-distribution
- ③ Log Normal distribution  $\rightarrow$  PDF
- ④ Power Law distribution  $\rightarrow$  PDF
- ⑤ Bernoulli distribution  $\rightarrow$  PMF
- ⑥ Binomial distribution  $\rightarrow$  PMF
- ⑦ Poisson distribution  $\rightarrow$  PMF
- ⑧ Uniform distribution  $\rightarrow$  Discrete PMF or Continuous PDF
- ⑨ Exponential distribution  $\rightarrow$  PDF
- ⑩ Chi-square distribution  $\rightarrow$  PDF

11. F Distribution

The probability distribution of ratio of two independent standard normal variables is called F distribution.

Probability distribution:

The probability distribution gives the possibility of each outcome of a random experiment or event.

Example:  $X \sim \text{Bin}(n=10, p=0.5)$

Notation:  $P(X=x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$

where  $x \in \{0, 1, 2, \dots, n\}$

Probability distribution function (PDF):

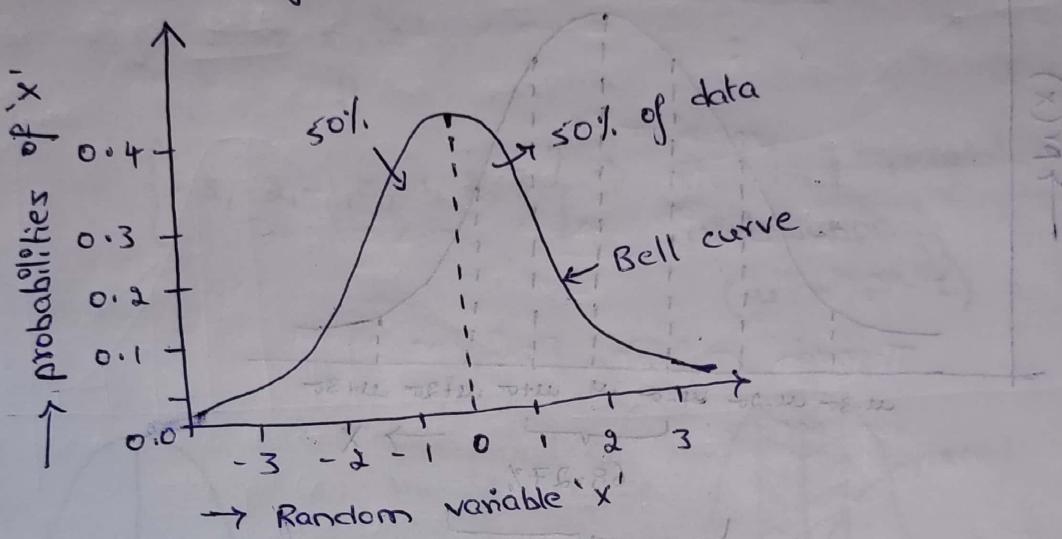
$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$


## ① Normal/Gaussian Distribution:

Let  $X$  be the random variable follows Gaussian distribution with mean ( $\mu$ ) and standard deviation ( $\sigma$ )

Here,  $x$  is a continuous random variable.

Fig: Normal distribution.

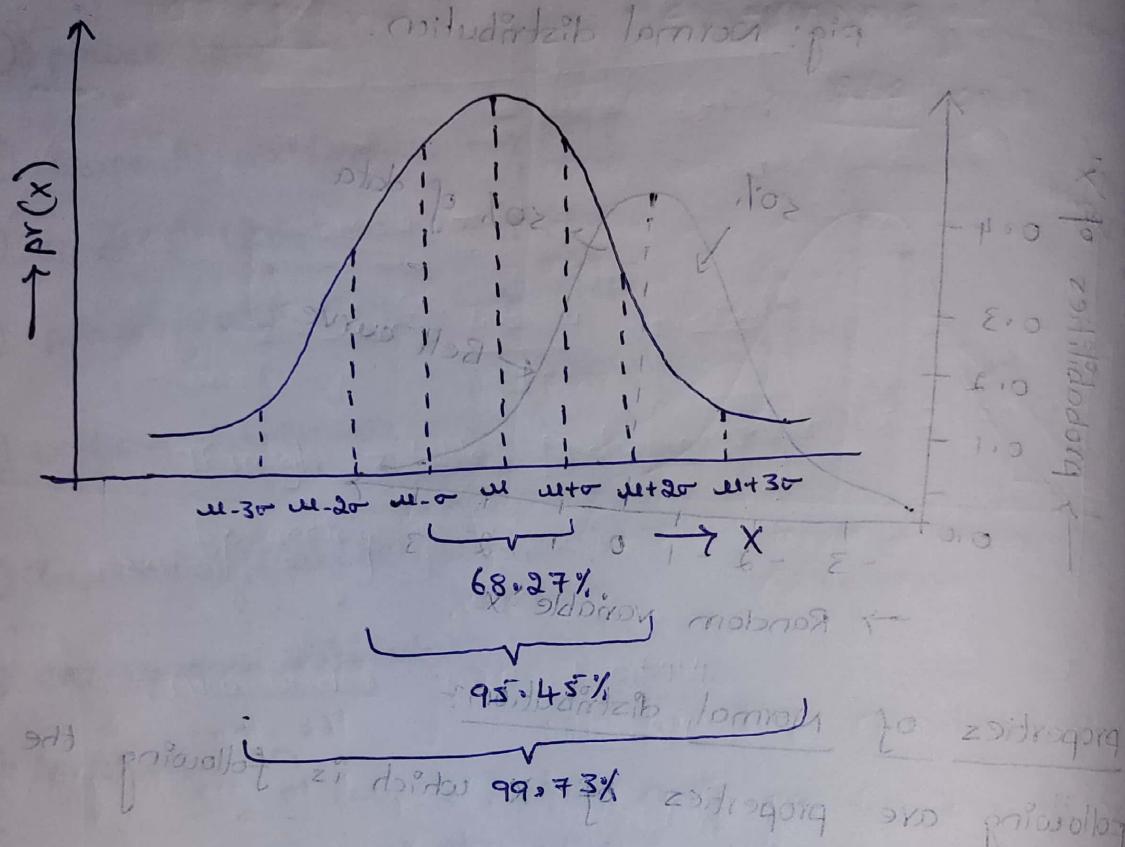


### Properties of Normal distribution:

- Following are properties of  $X$  which is following the normal distribution:
- In a normal distribution, the mean, median and mode are equal (i.e. Mean = Median = Mode).
  - The total area under the curve should be equal to 1.
  - The normally distributed curve should be symmetric at the centre.
  - There should be exactly half of the values are to the right of the centre and exactly half of the values are to the left of the centre.
  - The normal distribution should be defined by the mean and standard deviation.
  - The normal distribution curve must have only one peak.
  - The curve approaches the x-axis, but it never touches, and it extends farther away from the mean.

## Empirical Rule:

The empirical rule is the statistical rule that almost all normal distributed data fall within three standard deviations of the mean or  $\mu \pm 3\sigma$ . Here,



$Pr(\mu - \sigma \leq x \leq \mu + \sigma) \approx 68.27\%$ , within  $\sigma$  standard deviations.

$Pr(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 95.45\%$ , within  $2\sigma$  standard deviations.

$Pr(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 99.73\%$ , within  $3\sigma$  standard deviations.

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

## Applications:

- ① Height
- ⑤ Blood pressure
- ② Rolling a dice
- ⑥ Shoe size
- ③ Tossing a coin
- ⑧ Weight
- ④ IQ

## 2) Standard Normal distribution:-

A normal distribution with zero mean ( $\mu = 0$ ), unit variance ( $\sigma = 1$ ) is called as standard Normal Distribution.

$$\{0, 1, 2, 3, 4, 5, 6, 7\} \Rightarrow \text{Normal distribution.}$$

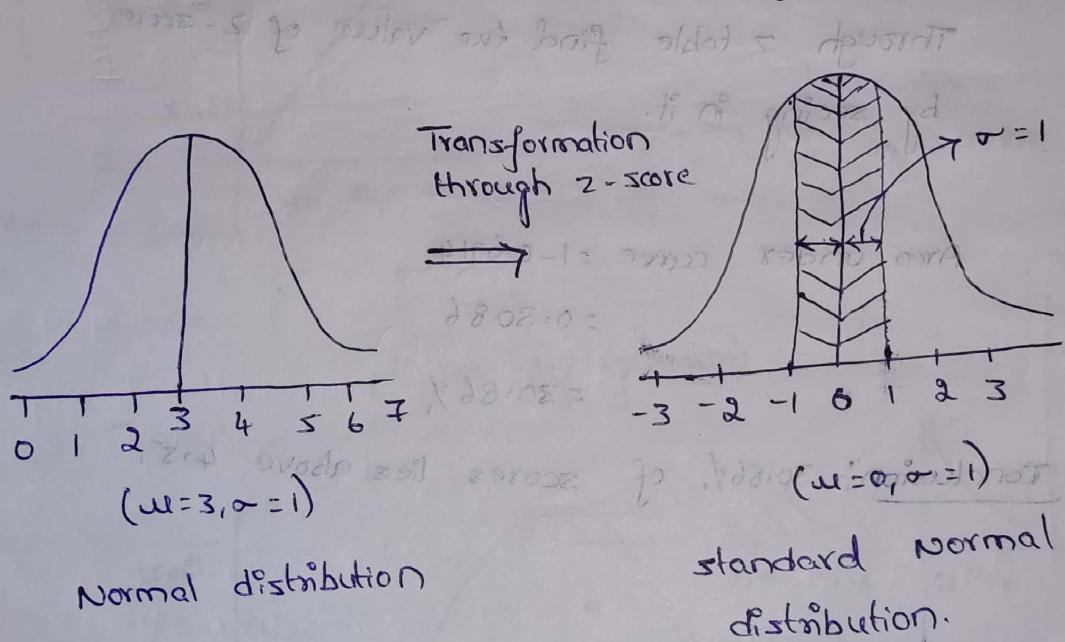
$(\mu = 3, \sigma = 1)$

↓

z-score

$$\{-3, -2, -1, 0, 1, 2, 3\} \Rightarrow \text{standard normal distribution.}$$

$(\mu = 0, \sigma = 1)$



### Z-score :-

Z-score tells you about a value how many standard deviation away from the mean.

$$Z_{\text{score}} = \frac{x_i - \mu}{\sigma}$$

$$= \frac{1-3}{1} = -2$$

$$\frac{4-3}{1} = 1$$

$$= \frac{2-3}{1} = -1$$

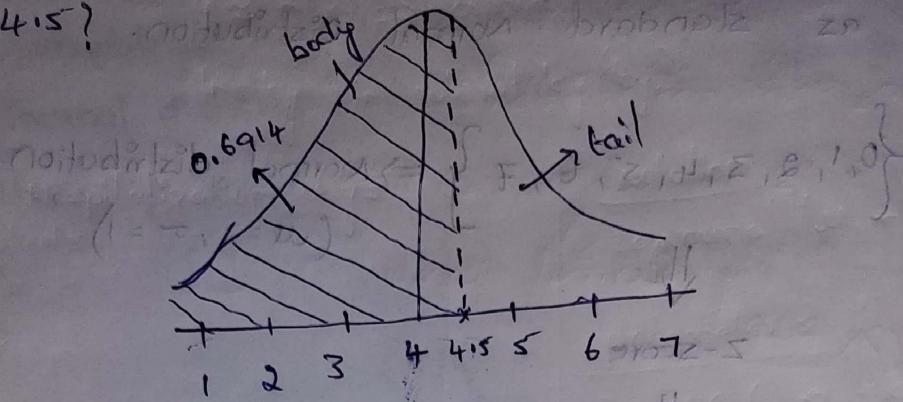
$$\frac{5-3}{1} = 2$$

$$\frac{3-3}{1} = 0$$

$$\frac{6-3}{1} = 1$$

Ex: What is the percentage of scores which lies above

4.5?



$$\mu = 4, \sigma = 1$$

Normal distribution  $\left\{ \varepsilon, \varepsilon, 1, 0, 1, 6, \varepsilon, \varepsilon \right\}$

$$z\text{-score} = \frac{4.5 - 4}{1} = 0.5$$

$$(1-\varepsilon, 6-\varepsilon)$$

Through z-table find the value of z-score

by seeing in it.

$$\text{Area under curve} = 1 - 0.6914$$

$$= 0.3086$$

$$= 30.86\%$$

Conclusion: 30.86% of scores lies above 4.5.

Normal distribution

$$\frac{\mu + 1.6}{\sigma} = 0.7025$$

$$L = \frac{\Sigma + 1}{\sigma}$$

$$B = \frac{\Sigma - 1}{\sigma}$$

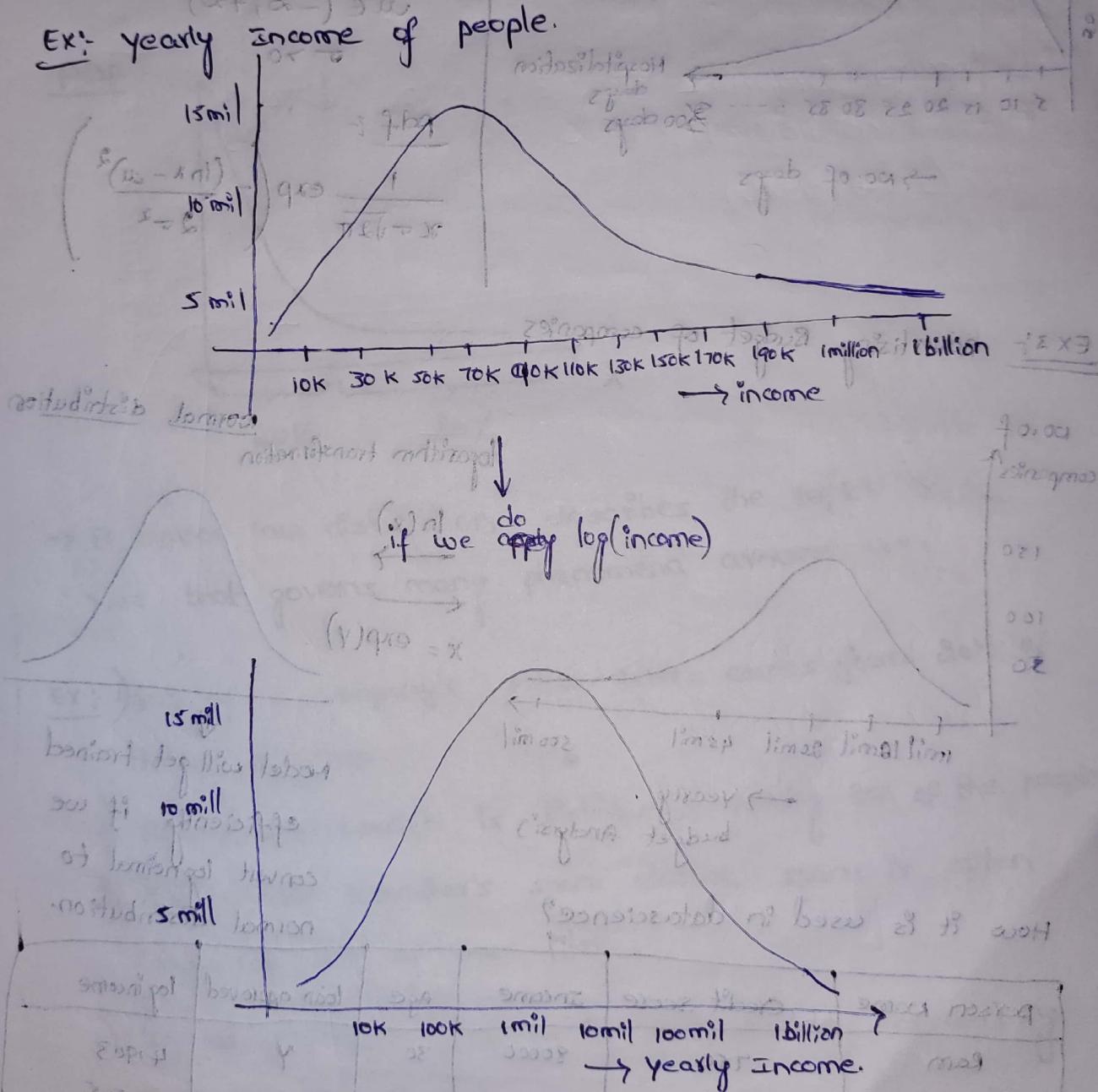
$$I = \frac{\Sigma - 3}{\sigma}$$

$$T = \frac{\Sigma - 6}{\sigma}$$

### ③ Log Normal Distribution: $(\mu, \sigma)$ for $\ln(x)$

In probability theory, a log-normal distribution is a continuous probability distribution of a random variable whose logarithm is normally distributed. Thus if the random variable  $x$  is log-normally distributed, then  $y = \ln(x)$  has a normal distribution, and the exponential function of  $y, x = \exp(y)$ , has a log-normal distribution.

Ex:- Yearly income of people.

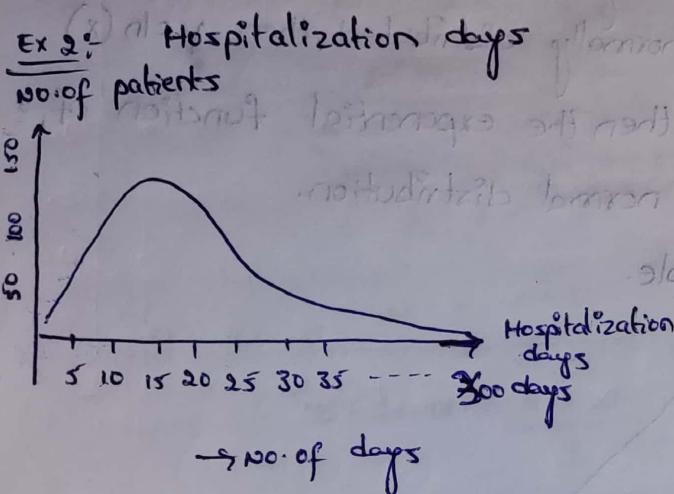


→ If you get a normal distribution by applying a log function to a dataset then dataset is log normally distributed.

$$x \sim \text{lognormal}(\mu, \sigma^2)$$

$$y \approx \ln(x) \Rightarrow \text{normal distribution}(\mu, \sigma^2)$$

$$\ln(x) = \ln = \text{natural log}(\log_e)$$



notation

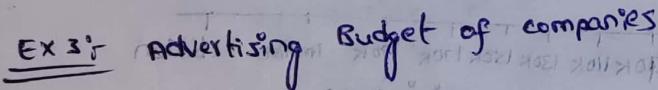
$$\text{Lognormal}(\mu, \sigma^2)$$

parameter s

$$\mu \in (-\infty, +\infty)$$

pdf :

$$\frac{1}{x \sqrt{2\pi}} \exp\left(-\frac{(\ln x - \mu)^2}{2\sigma^2}\right)$$



logarithm transformation

$$\ln(x)$$

$$x = \exp(y)$$

normal distribution

Model will get trained  
efficiently if we  
convert lognormal to  
normal distribution.

How it is used in datascience?

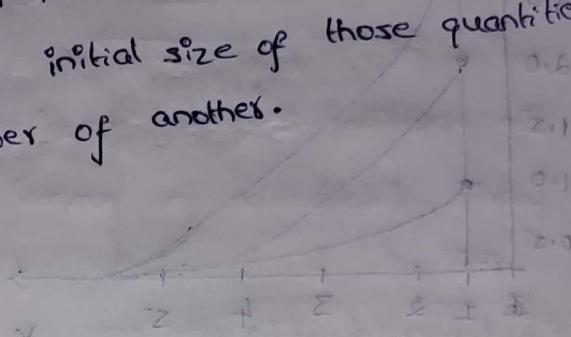
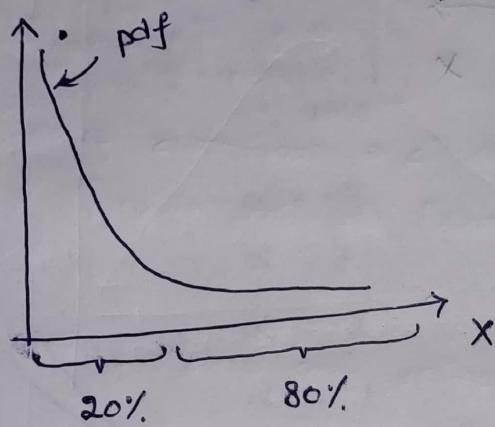
person name	Credit score	income	Age	loan approved	log income
Ram	750	80000	30	Y	4.1903
sam	310	32000	32	N	4.1505
Tom	475	77000	40	Y	4.886
Dinesh	600	65000	35	N	4.812
Hari	820	85000	37	Y	5.1740
Varun	780	73000	22	Y	4.8750

quite different

#### ④ power law distribution:

In statistics, power law is a functional relationship between two quantities, where a relative change in one quantity results in a proportional relative change in other quantity, independent of initial size of those quantities. One quantity varies as power of another.

PDF



→ A power law distribution describes the ~~80-20~~ 80-20% rule that governs many phenomena around us.

Ex: i) 80% of a company's sales often comes from 20% of their customers.

ii) 80% of the wealth is distributed among 20% of the people.

iii) 80% of the computer's space storage space is often taken by 20% of the files.

iv) In ipl, 20% of the team is responsible for winning 80% of match.

v) 80% of project complete by 20% of team.

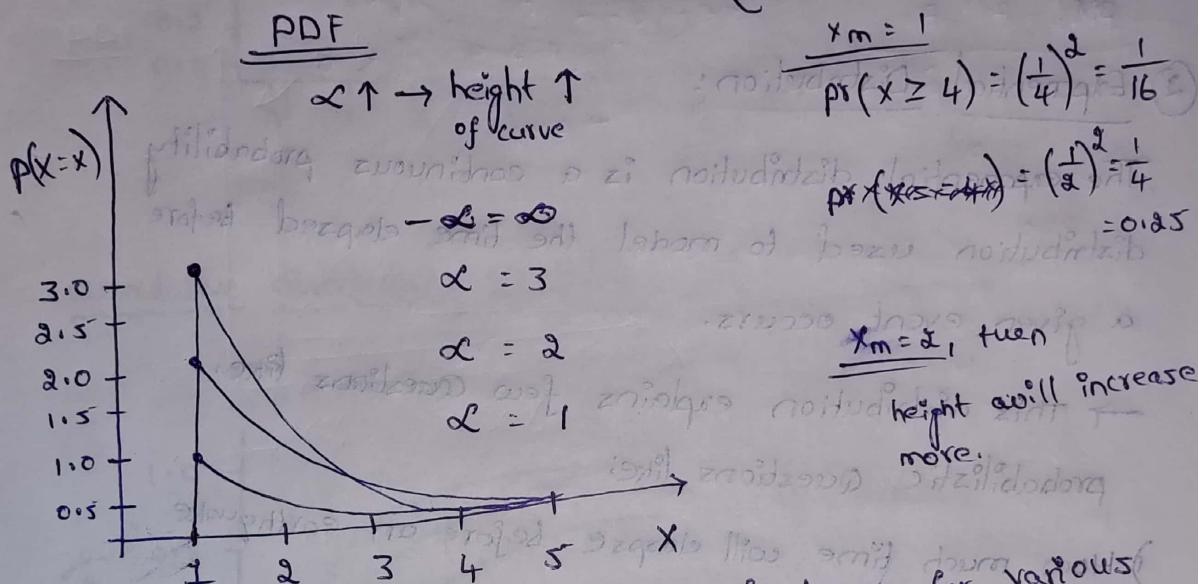
Types of power law distribution:-

- ① pareto Distribution
- ② Exponential distribution.

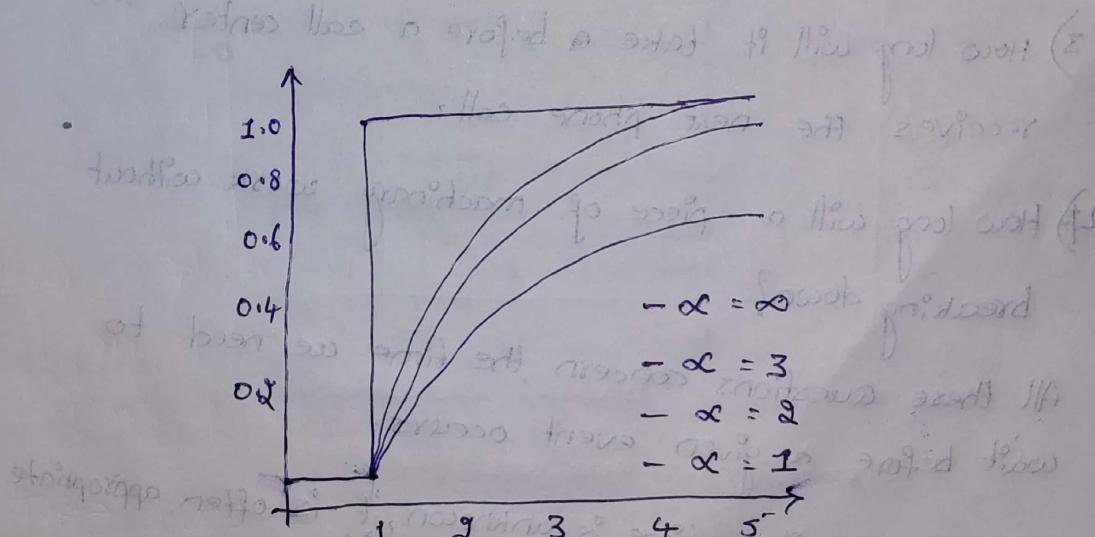
## ① Pareto Distribution: [Continuous random variable]

If  $x$  is a random variable with a pareto (Type I) distribution, then the probability that  $x$  is greater than some number  $x$ , i.e. the survival function (also called tail function), is given by.

$$\text{pdf} : F(x) = \Pr(X > x) = \begin{cases} \left(\frac{x_m}{x}\right)^{\alpha} & x \geq x_m \\ 1 & \text{otherwise} \end{cases}$$



→ Pareto type I probability density functions for various  $\alpha$  with  $x_m = 1$ . As  $\alpha \rightarrow \infty$ , the distribution approaches  $\delta(x - x_m)$  where  $\delta$  is dirac delta function.



Pareto Type I with cumulative distribution functions

→ Pareto distribution follows power law (80-20)% rule.

[old data (not having a normal distribution)  $\rightarrow$  modified data (normal distribution)]  
 Box-cox transformation  
 $y = \frac{1}{\lambda} \ln(x + \lambda)$   
 $\lambda = 0 \Rightarrow y = \ln(x)$   
 $\lambda = 1 \Rightarrow y = x$

By applying Box-cox transformation we can convert  
 pareto distribution to normal distribution.

### ② Exponential Distribution:

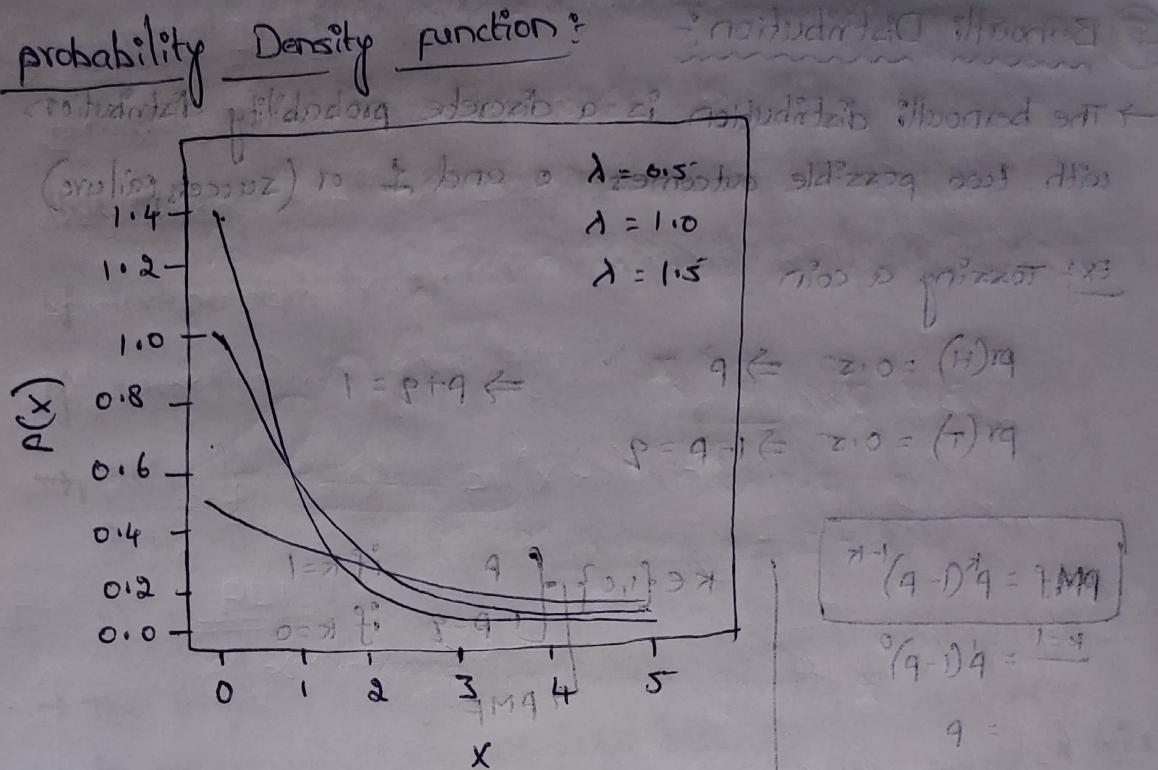
The exponential distribution is a continuous probability distribution used to model the time elapsed before a given event occurs.

→ This distribution explains few questions like:  
 probabilistic Questions like:

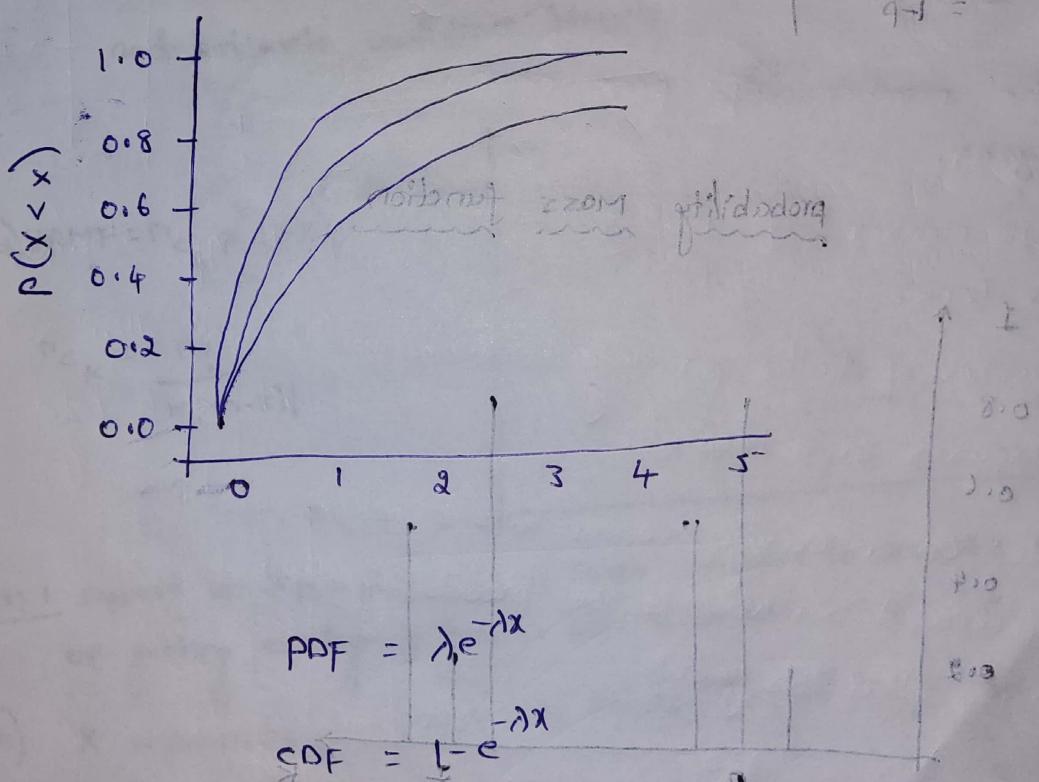
- 1) How much time will elapse before an earthquake occurs in a given region?
- 2) How long do we need to wait until a customer enters our shop?
- 3) How long will it take before a call center receives the next phone call?
- 4) How long will a piece of machinery work without breaking down?

All these questions concern the time we need to wait before a given event occurs.

→ If this waiting time is unknown, it is often appropriate to think of it as a random variable having an exponential distribution.



Cumulative Distribution Function:



Illustrated for reference series

$$8.0 = (1-x)q \quad 6.0 = (2-x)q$$

$$6.0 = (1-x)q \quad 8.0 = (3-x)q$$

$$7.0 = (1-x)q \quad 7.0 = (4-x)q$$



## 5 Bernoulli Distribution:

→ The bernoulli distribution is a discrete probability distribution with two possible outcomes 0 and 1 or (success, failure)

Ex: Tossing a coin

$$Pr(H) = 0.5 \Rightarrow p \quad \Rightarrow p+q=1$$

$$Pr(T) = 0.5 \Rightarrow 1-p=q$$

$$PMF = p^k(1-p)^{1-k}$$

$$\underset{k=1}{=} p^1(1-p)^0$$

$$= p$$

$$\underset{k=0}{=}$$

$$= p^0(1-p)^{1-0}$$

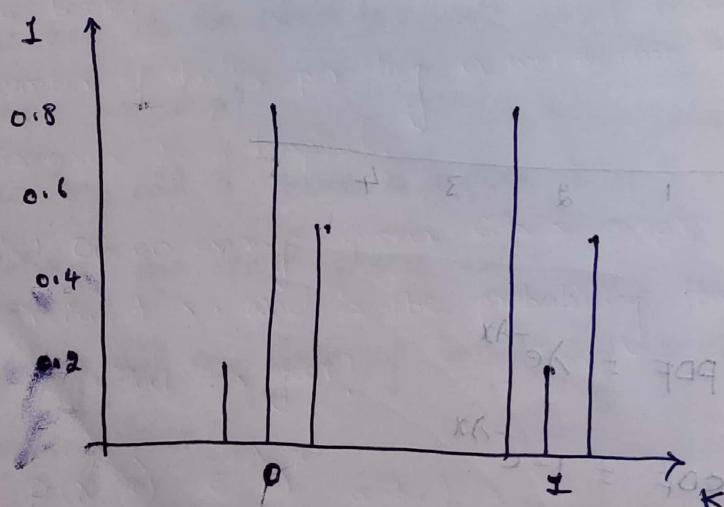
$$= 1 \times (1-p)$$

$$= 1-p$$

$$K \in \{1, 0\} \quad \begin{cases} p & \text{if } k=1 \\ 1-p = q & \text{if } k=0 \end{cases}$$

PMF

probability mass function



Three examples of bernoulli distribution.

$$P(X=0) = 0.2 \text{ and } P(X=1) = 0.8$$

$$P(X=0) = 0.8 \text{ and } P(X=1) = 0.2$$

$$P(X=0) = 0.5 \text{ and } P(X=1) = 0.5$$

⑥ Binomial distribution:

Binomial distribution is a discrete probability distribution that describes the number of successes in a fixed number of independent trials, where each trial has the same probability of success.

→ This distribution is characterized by two parameters.

1) The probability of success in each trial (denoted by  $p$ )

2) The number of trials (denoted by  $n$ ).

→ The binomial distribution is often used in data science

→ The binomial distribution is often used for outcomes such as success or failure, yes or no, or heads or tails.

→ It is used to calculate probabilities of various outcomes and can be used to perform hypothesis testing, make predictions and estimate confidence intervals.

$$P(X=k) \text{ PMF} = \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k}$$

$$\frac{n!}{k!(n-k)!} = \frac{n(n-1)(n-2)\dots(n-k+1)}{k(k-1)(k-2)\dots(2)(1)}$$

Ex: 1 Suppose we flip a fair coin 10 times and want to know the probability of getting exactly 3 heads. Here  $n=10$  and  $p=0.5$ .

A)  $X$  represents the number of heads obtained in 10 coin flips.

$$k=3$$

$$n=10$$

$p=0.5$  since the coin is fair

$$\frac{10!}{3!7!} = \frac{10 \times 9 \times 8 \times 7 \times 6 \times 5 \times 4 \times 3 \times 2 \times 1}{3 \times 2 \times 1 \times 7 \times 6 \times 5 \times 4 \times 3 \times 2 \times 1}$$

$$P(X=3) = \frac{10!}{3!7!} \times 0.5^3 \times 0.5^{10-3}$$

$$= 120 \times 0.5^10$$

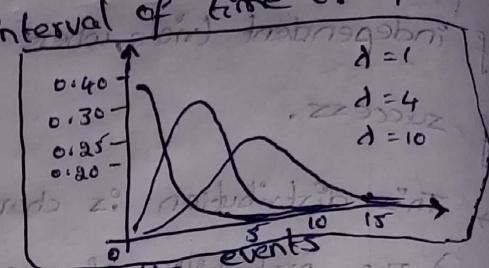
$$\approx 0.117$$

probability of getting exactly 3 heads in 10 flips is 11.7%

## 7 Poisson distribution :- (discrete random variable)

→ The poisson distribution is a probability distribution that models the number of times an event occurs within a specific interval of time or space.

$$P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}$$



$\lambda$  - Is the average rate of events per interval.

$x$  - no. of events that occur in a given interval.

→ poisson distribution is used in cases where the chances of any individual event being success very small.

Ex:- The number of plane crash in India in a year.

The number of printing mistakes in each page of a book.

→ Number of deaths per day or week due to rare diseases in a hospital.

Example:- On an average cancer kills 5 people each year in India,  $\lambda = 5$ . What is the probability that only one person is killed this year?

$$P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!} \text{ with } \lambda = 5$$

$$P(X=1) = \frac{e^{-5} 5^1}{1!} \approx 0.033$$

Thus

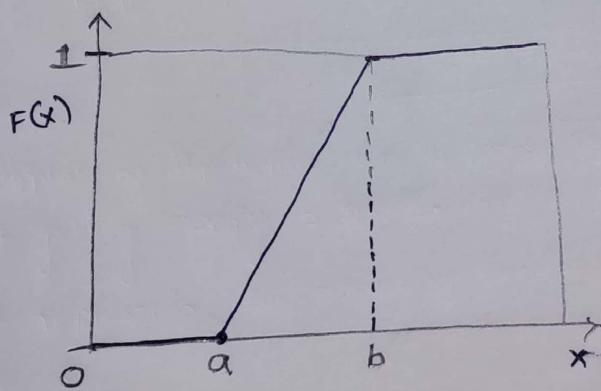
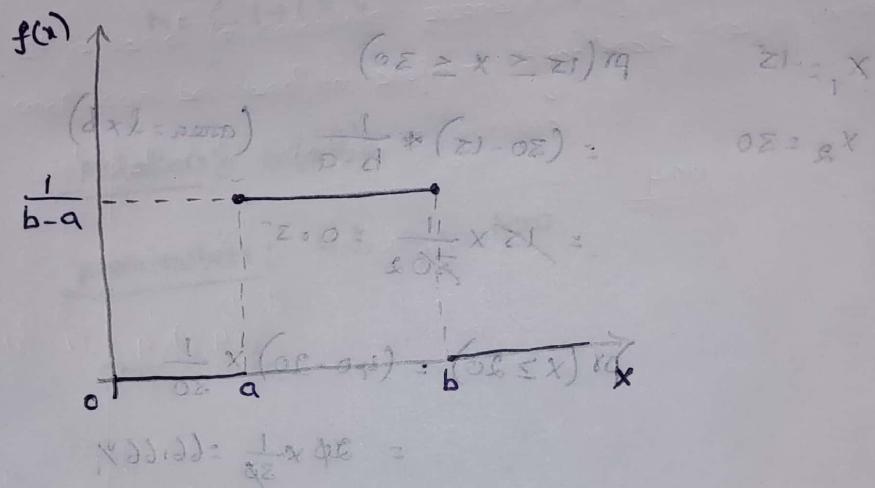
Ex:- The number of people visits bank in an hour.

- ⑧ Uniform distribution is a probability distribution function for continuous random variables.
- ① Continuous uniform distribution. (pdf)
- ② Discrete uniform distribution. (pmf)

### i) Continuous uniform distribution:-

A continuous uniform distribution is a probability that is used to describe situations where all values within a certain range are equally likely to occur.

→ In other words, the probability density function of a continuous uniform distribution is a constant within a specified interval and zero outside that interval.



→ The probability density function (pdf) of a continuous uniform distribution with parameters  $a$  and  $b$  is given by:

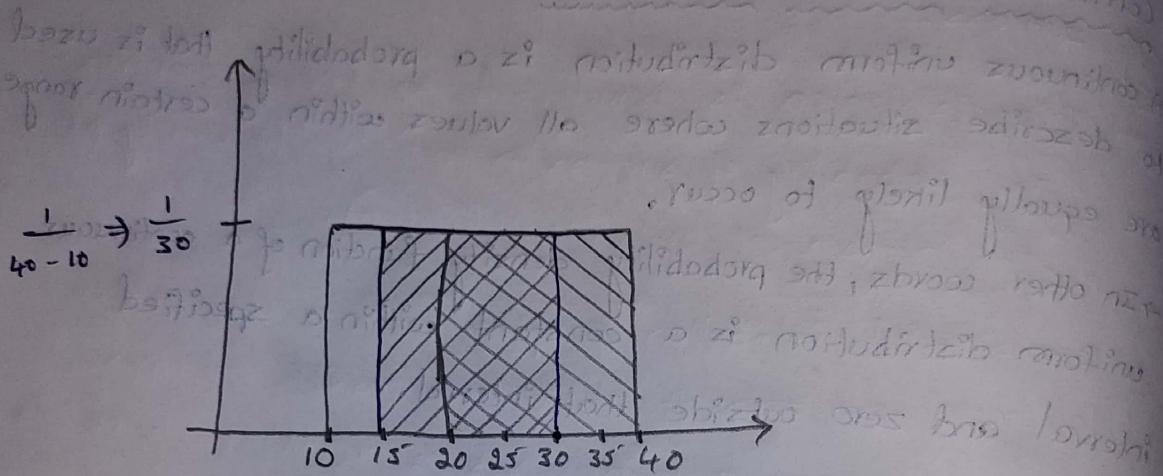
Notation:  $[max, min] \Rightarrow$  interval  $u(a, b)$

parameters:  $-\infty < a < b < \infty$   $b > a$

$$\text{pdf} = \begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

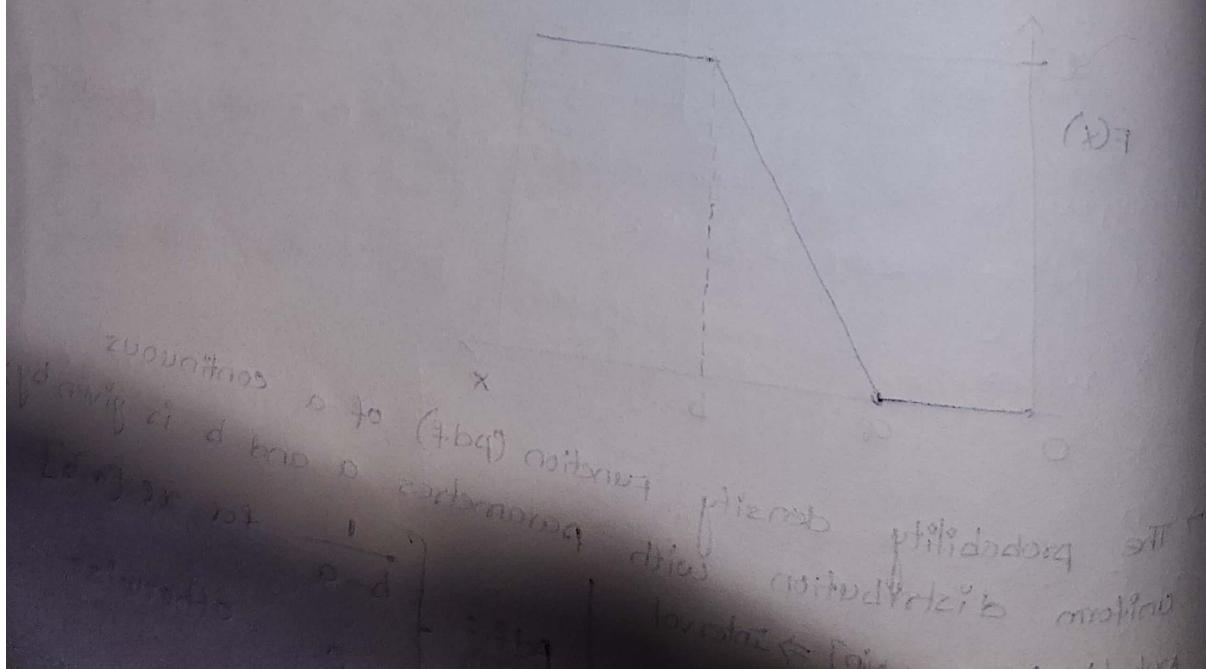
Eg:- The no. of candidates sold daily at a shop is uniformly distributed with a maximum of 40 and minimum of 10.

i) probability of daily sales to fall between 15 and 30.



$$\begin{aligned}
 x_1 &= 15 & \Pr(15 \leq x \leq 30) \\
 x_2 &= 30 & = (30 - 15) * \frac{1}{b-a} \quad (\text{area} = l \times b) \\
 & & = 15 * \frac{1}{30} = 0.5
 \end{aligned}$$

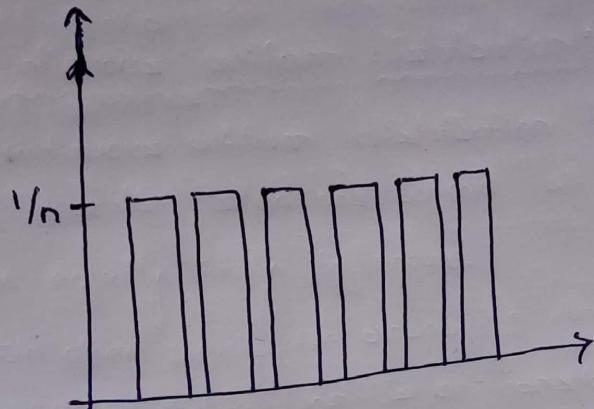
$$\begin{aligned}
 \Pr(x \geq 20) &= (40 - 20) * \frac{1}{30} \\
 &= 20 * \frac{1}{30} = 66.66\%
 \end{aligned}$$



ii) Discrete uniform distribution :-

Rolling a dice =  $\{1, 2, 3, 4, 5, 6\}$

$$pr(1) = \frac{1}{6} \quad pr(n) = \frac{1}{6}$$



$$n = b - a + 1$$

$$n = 6 - 1 + 1 = 6$$

Notation :-  $u(a,b)$

parameters :-  $a, b$  with  $b > a$

$$\text{PMF} := \frac{1}{n}$$