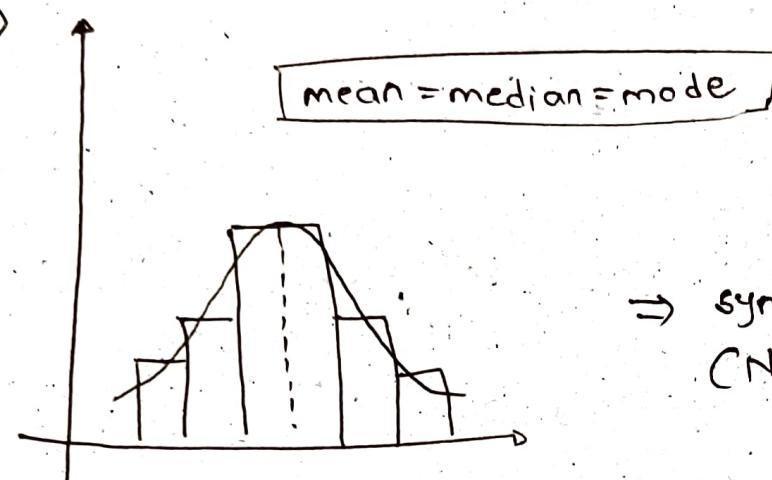


4

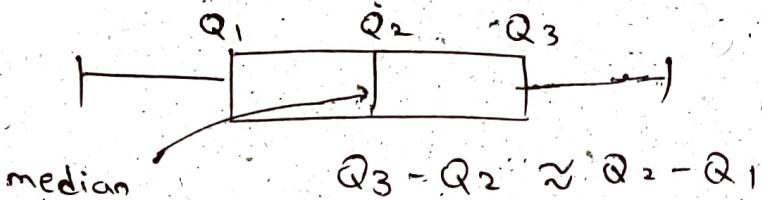
Skewness



⇒ symmetrical distribution
(NO SKEWNESS)

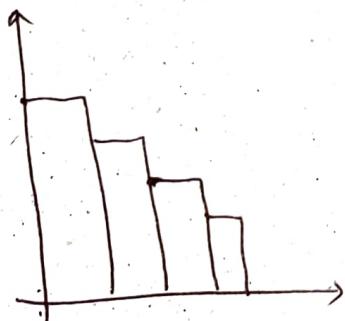
The mean, median and mode all are perfectly at the center.

Hence, boxplot for above distribution will look like:



e.g. → Normal / Gaussian Distribution

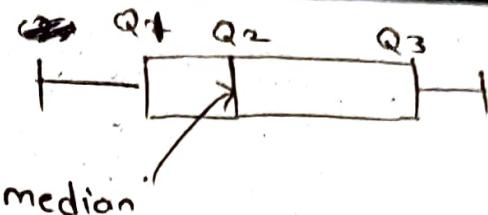
2) Right skewed



⇒ positive skewed
(right)



Boxplot for right skewed \rightarrow



$$Q_3 - Q_2 > Q_2 - Q_1$$

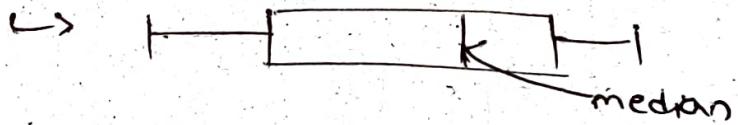
mean > median > mode

3) Left skewed (neg. skewed)



mean < median < mode

Box-plot for left skewed



$$Q_2 - Q_1 > Q_3 - Q_2$$

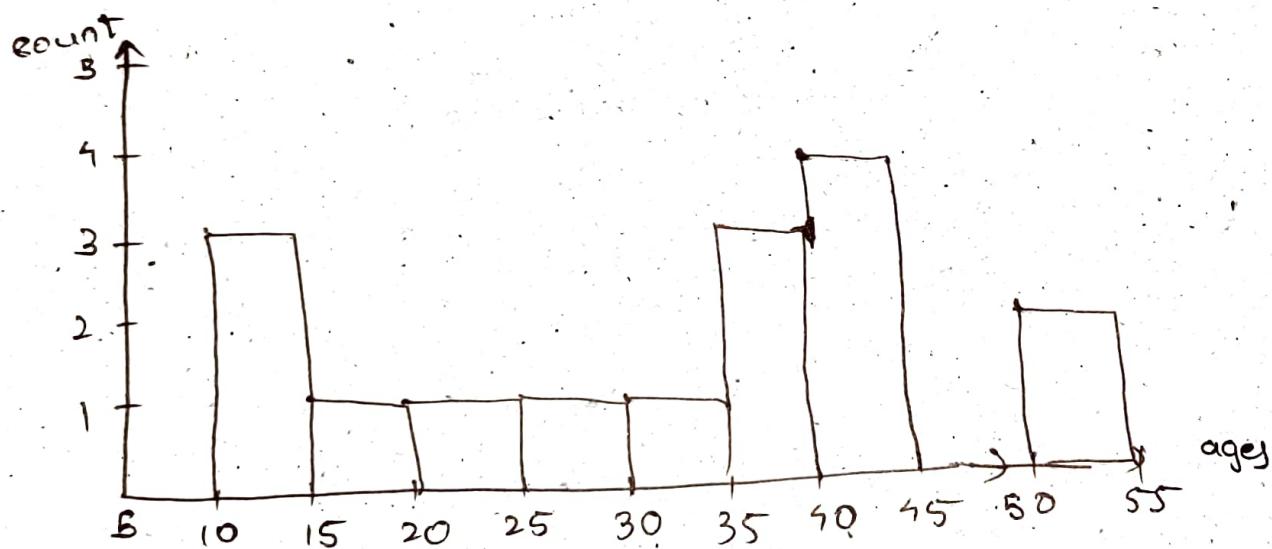
Q) Histograms - It is nothing but representation of frequency of elements.

ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51}

Suppose bins = 10

$$\therefore \text{bin size} = \frac{50 - 0}{10} = ⑥/1$$

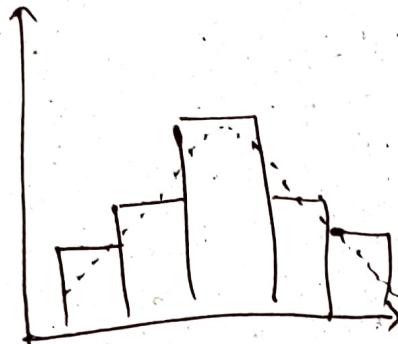
Note: while checking range (say 10 to 15), we check for ≥ 10 but < 15 .



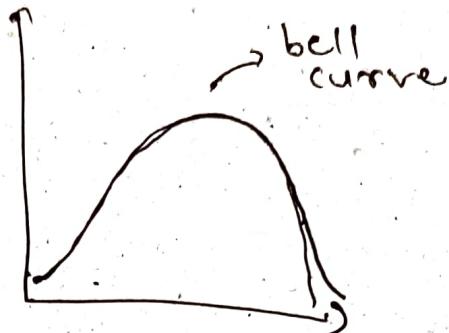
Smoothening above histograms



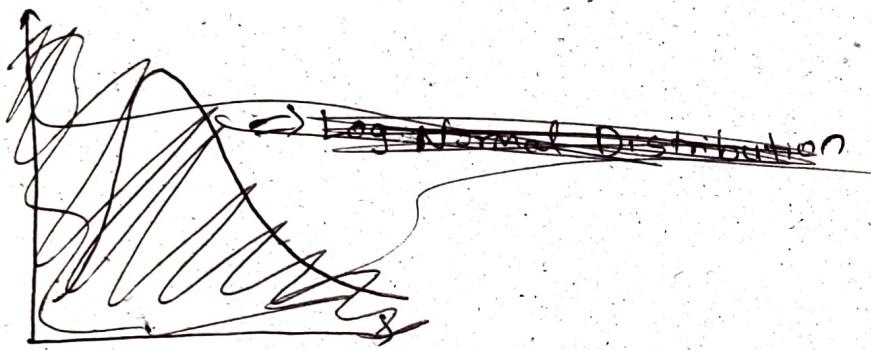
histogram



smoother



normal / gaussian
distribution



Log Normal Distribution

Probability Distribution Functions

discrete

Probability mass
function (pmf)

cumulative
distribution
function (cdf)

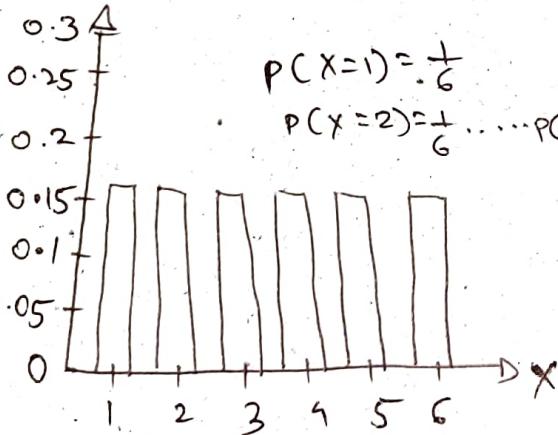
Probability density
function (pdf)

continuous

① Probability mass function (pmf) : Distribution
of discrete random variable.

e.g. → rolling a die {1, 2, 3, 4, 5, 6}

prob mass. $\frac{1}{6} = 0.16$



$$P(x) = P(X=x)$$

conditions:

- $P(x) \geq 0$
- $\sum P(x) = 1$

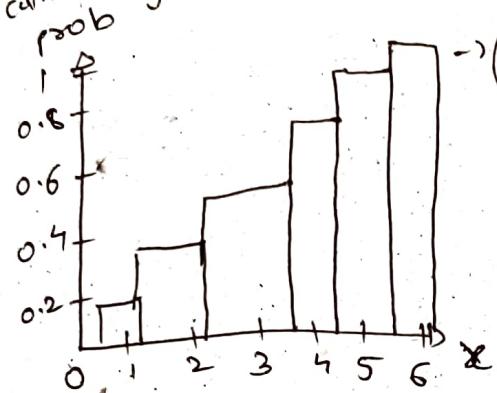
NOTE: ① Value of pmf at specific point indicated
the probability of the random variable
taking on that value

② Sum of pmf over all values of random
variable is 1

② cumulative distribution function (cdf)

For Discrete Data

eg → rolling a die $\{1, 2, 3, 4, 5, 6\}$



→ cumulative sum

→ Last bar will always be 1.

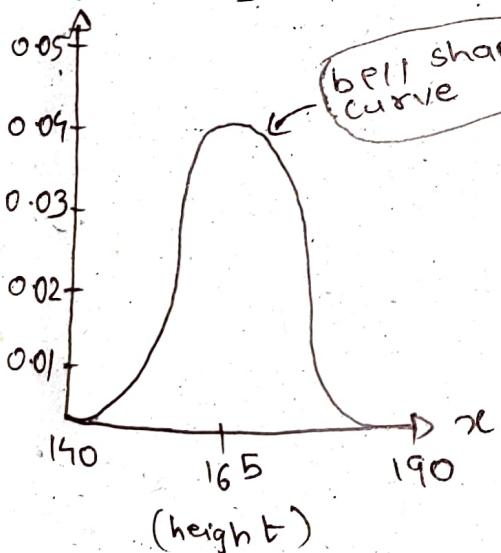
$$\begin{aligned} P(X \leq 4) &= P(X=1) + P(X=2) \\ &\quad + P(X=3) + P(X=4) \end{aligned}$$

Note: ① Value of Cdf at a specific point indicates probability that random variable will take on a value less than or equal to that point.

③ probability density function (pdf): Distribution of continuous data / random variable.

eg → height of students

prob density



- Not that students with height of 140 or 190
- max students near height of 165.

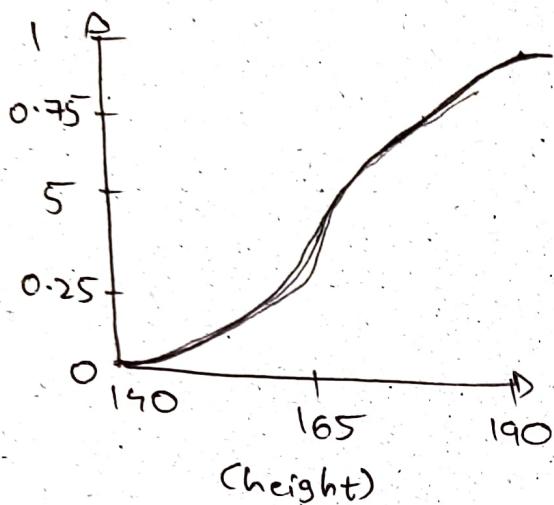
Note:

You can see that the maximum height of curve is at 0.04.
what does that mean?
will see in some time

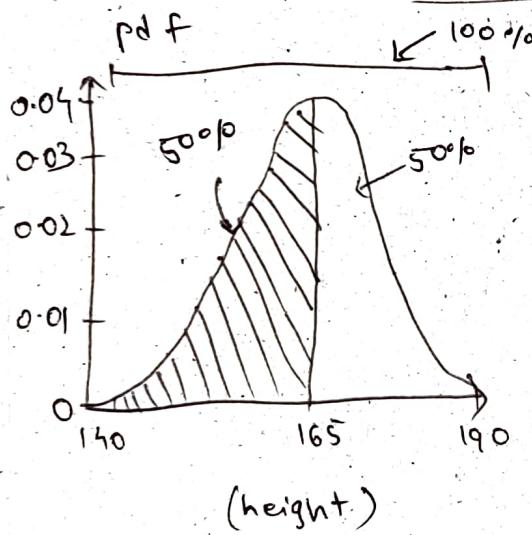
④ Cumulative distribution function: (cdf)

For continuous data

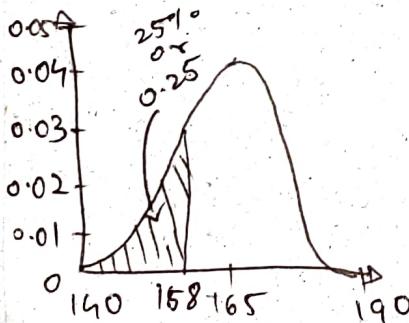
e.g. → height of students
Prob



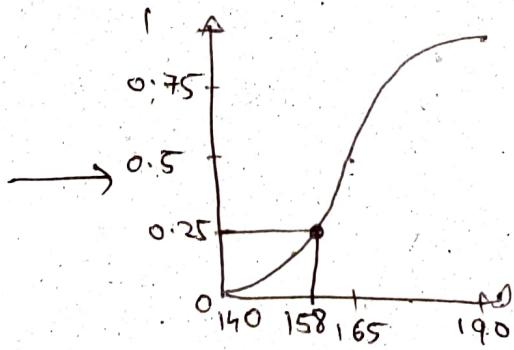
Comparing pdf & cdf



Till 165, area under curve is 50% → 0.5



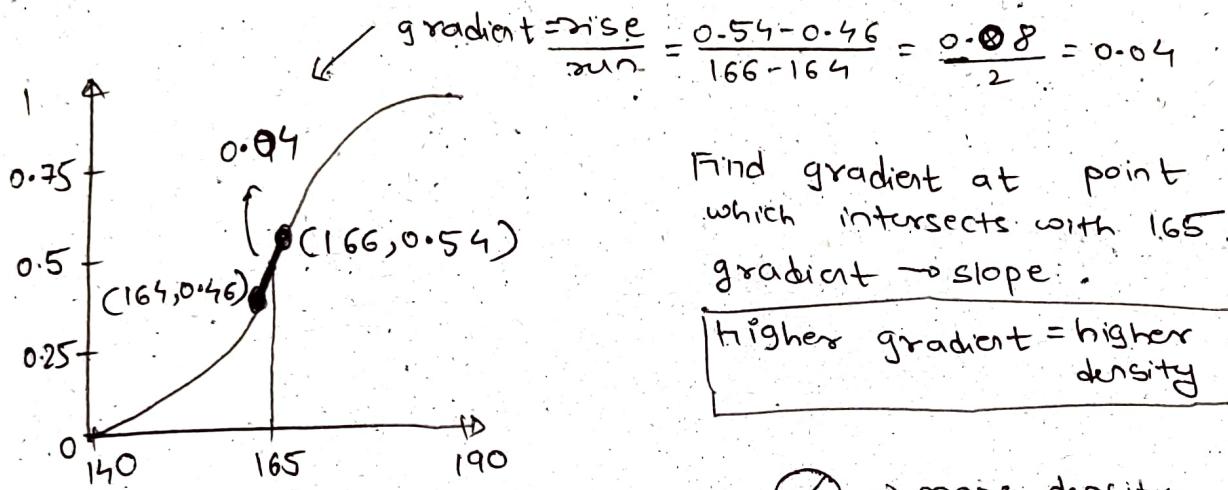
Here, at 165, we can see it matches to 0.5 at y-axis.



Hence, we can see that, the y-axis in cdf represents the area under curve for a specific value of random variable in pdf.

Hence, in this way, we can derive cdf from pdf (pdf $\xrightarrow{\text{area under curve}} \text{cdf}$)

Now, for cdf \rightarrow pdf.



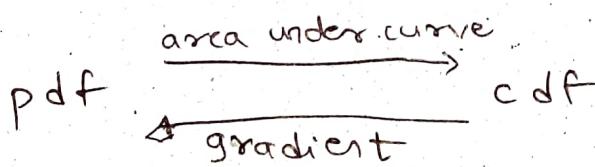
0.04 is approx gradient at 165.

You can verify this from pdf graph.

maximum height (at 165) is at 0.04 on y-axis.

This means, the y-axis in pdf represents the gradient for a specific value of random variable using the S-curve in cdf.

Hence, in this way we can derive pdf from cdf.



Probability density function $\rightarrow f(x)$

Cumulative distribution function $\rightarrow F(x)$

Note: ① differentiation of cdf = pdf

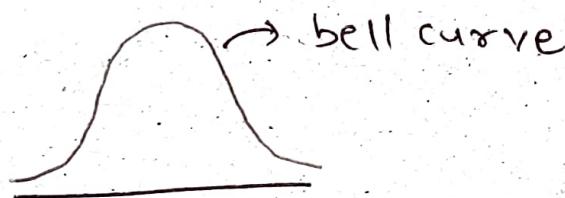
② integration of pdf = cdf.

$$\frac{d}{dx} F(x) = f(x)$$

$$\int_{-\infty}^x f(x) dx = F(x)$$

(*) Types of Probability Distribution.

① Normal / Gaussian distribution. (pdf)



② Bernoulli Distribution (pmf):

Outcomes are only

two $\rightarrow \{ \text{success}, \text{failure} \}$

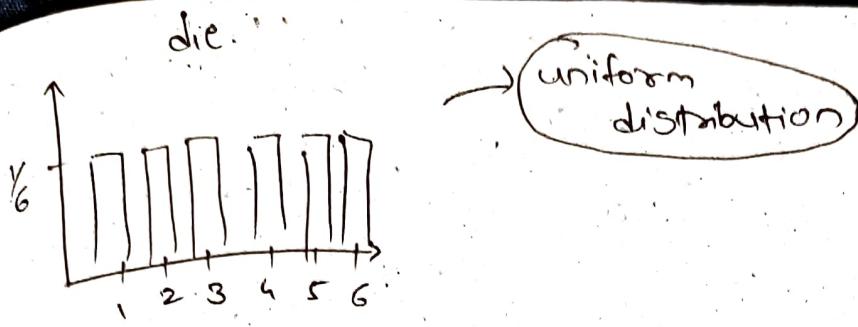
③ Uniform distribution

④ Log normal distribution (pdf)

⑤ Poisson distribution. (pmf)

⑥ Power law distribution. (pdf) $\rightarrow 80-20\% \text{ rule}$

⑦ Binomial Distribution (pmf)



uniform
distribution

① Bernoulli Distribution (pmf) (Binary Outcomes)

It is the discrete probability distribution of a random variable which takes the value $\textcircled{1}$ with probability P and the value $\textcircled{0}$ with probability $q = 1 - P$.

For a set of possible outcomes, it asks a Yes - No question.

success / yes / true / one $\rightarrow p$

failure / no / false / zero $\rightarrow q = 1 - p$

eg \rightarrow Tossing a fair coin $\{0, 1\}$

$$P(1+) = \boxed{0.5} \rightarrow p$$

$$P(0-) = 1 - 0.5 = 1 - p = \boxed{0.5} \rightarrow q$$

probability mass function

prob



$$\textcircled{1} P(x=0) = 0.2 \& P(x=1) = 0.8$$

$$\textcircled{2} P(x=0) = 0.8 \& P(x=1) = 0.2$$

$$\textcircled{3} P(x=0) = 0.5 \& P(x=1) = 0.5$$

eg \rightarrow whether the person will pass / fail.

mathematically;

pmf \rightarrow $k=0 \text{ or } 1$

$$P(X=k) = p^k(1-p)^{1-k}$$

$$P(X=1) = p^1(1-p)^{1-1} = P$$

$$P(X=0) = p^0(1-p)^{1-0} = 1-p \rightarrow q$$

$$\therefore \text{PMF} = \begin{cases} q & \text{if } k=0 \\ p & \text{if } k=1 \end{cases}$$

Mean, Variance, Std dev for Bernoulli distribution

mean:

$$\begin{aligned} E(k) &= \sum_{i=1}^k k \cdot P(k) \\ &= (1 \times 0.6) + (0 \times 0.4) \\ &= 0.6 \rightarrow p \end{aligned}$$

$k=1 \text{ or } 0$

say $P(k=1) = 0.6 \rightarrow p$

$P(k=0) = 0.4 \rightarrow q$

$$\therefore E(k) = p$$

median:

$$\begin{cases} 0 & \text{if } p < \frac{1}{2} \\ [0,1] & \text{if } p = \frac{1}{2} \\ 1 & \text{if } p > \frac{1}{2} \end{cases}$$

Variance: f std

$$\text{Variance} = p(1-p) = pq$$

$$\text{std} = \sqrt{pq}$$

↖ PMF

Binomial Distribution: Binomial distribution with parameters ' n ' and ' p ' is [discrete] probability distribution of number of successes in a sequence of n independent experiments, each asking a yes-no question and with its own Boolean-valued outcome: (success with probability P) or (failure with probability $q = 1 - P$). A single experiment is also called a Bernoulli trial or Bernoulli experiment. A sequence of outcomes is called Bernoulli process.

For a single trial/experiment i.e. $n = 1$, the binomial distribution is a bernoulli distribution.

e.g. → Tossing a coin {Bernoulli distribution}

$$P(H) = 0.5 \quad P(T) = 0.5$$

$$\begin{matrix} \downarrow & \downarrow \\ 0 & 1 \end{matrix} \quad \begin{matrix} (1-P) & (P) \end{matrix}$$

Now, e.g. → Tossing a coin [for 10 times.]

$$\left. \begin{array}{l} \textcircled{1} \quad P(T) = P \quad P(H) = 1 - P \\ \textcircled{2} \quad P(T) = P \quad P(H) = 1 - P \\ \vdots \\ \textcircled{10} \quad P(T) = P \quad P(H) = 1 - P \end{array} \right\} \Rightarrow \begin{array}{l} \text{Binomial} \\ \text{distribution.} \end{array}$$

(Combine various Bernoulli)

Parameters:

$$n \in \{0, 1, 2, \dots\} \rightarrow \text{no. of trials/exp.}$$

$P \in [0, 1] \rightarrow 0 \text{ to } 1 \rightarrow \text{prob. for each trial (success)}$

$$q = 1 - P$$

PMF for binomial distribution!

$$P(X) = {}^n C_x P^x (1-p)^{n-x}$$

or

$$P(K) = {}^n C_k P^k (1-p)^{n-k}$$

$K \in \{0, 1, 2, \dots\} \rightarrow$ no. of successes

Mean of Binomial dist:

$$\text{mean} = np$$

Variance & std:

$$\text{Var} = npq$$

$$\text{Std} = \sqrt{npq}$$

Poisson Distribution

pmf

① Discrete distribution.

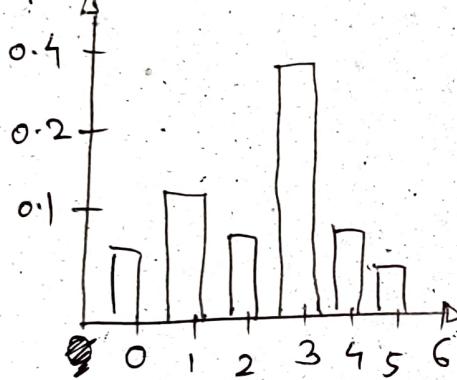
② Describes the number of [events] occurring in a fixed time interval.

e.g. → no. of people visiting hospital every hour.

- no. of people visiting banks every hour

- no. of people visiting airport every hour.

prob.



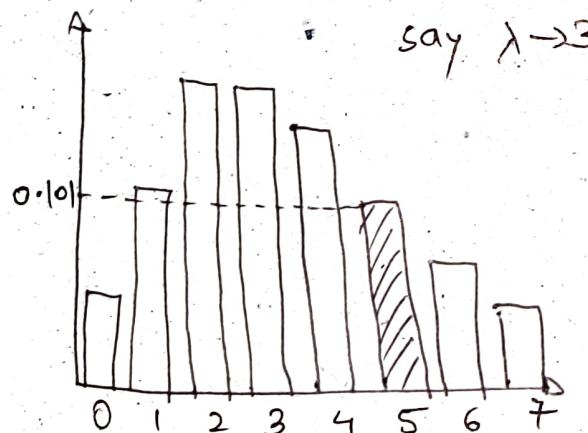
$$\lambda = 3$$

mean

~~Probability~~ no. of events occurring every hour.

$\lambda \rightarrow$ no. of events within a given event/interval of time or space

prob.



say $\lambda \rightarrow 3$

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

$$\therefore P(X=5) = \frac{e^{-3} \cdot 3^5}{5!} = 0.101 \Rightarrow 10\%$$

\therefore prob that event 5 will occur

Similarly, you can find

$$P(X=3) + P(X=5) =$$

Mean of Poisson distribution.

$$\boxed{\text{mean} \Rightarrow E(X) = \mu = \lambda * t}$$

$\lambda \rightarrow$ expected no. of events to occur at every time interval.

$t \rightarrow$ time interval.

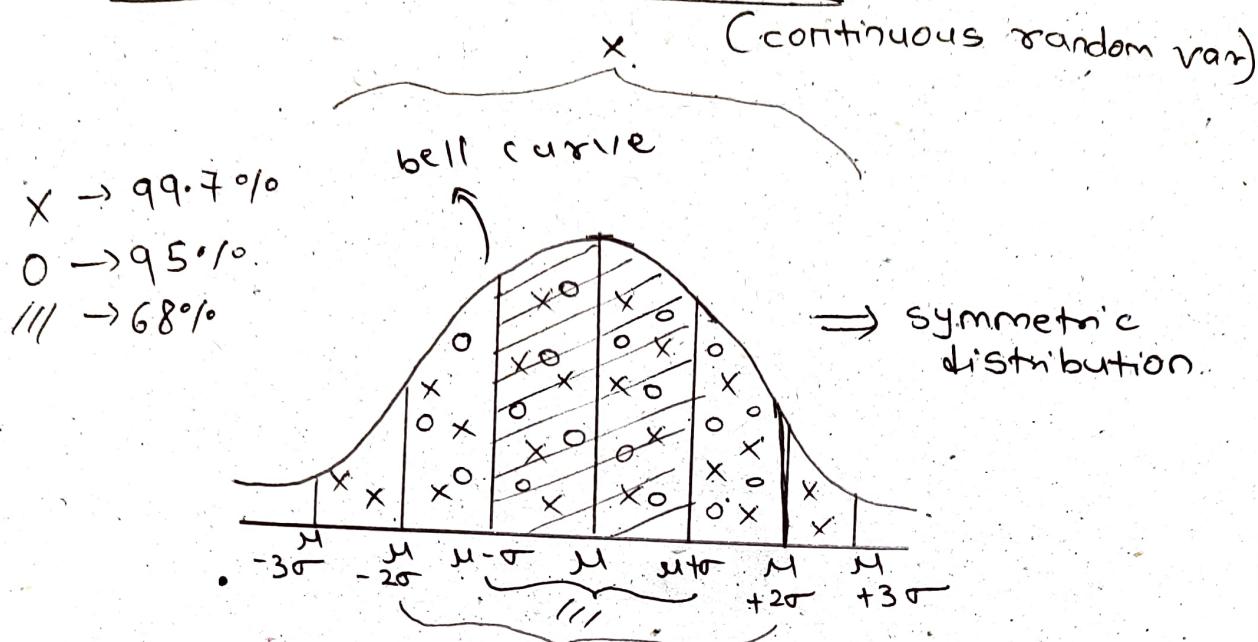
Variance of Poisson distribution.

$$\boxed{\text{Variance} \Rightarrow E(X) = \mu = \lambda * t}$$

Hence, In a Poisson distribution,

Mean & Variance are equal.

(*) Normal / Gaussian Distribution: ! pdf.



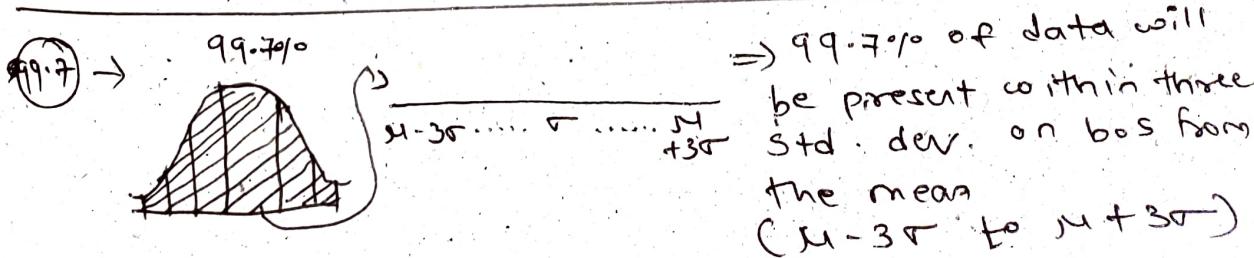
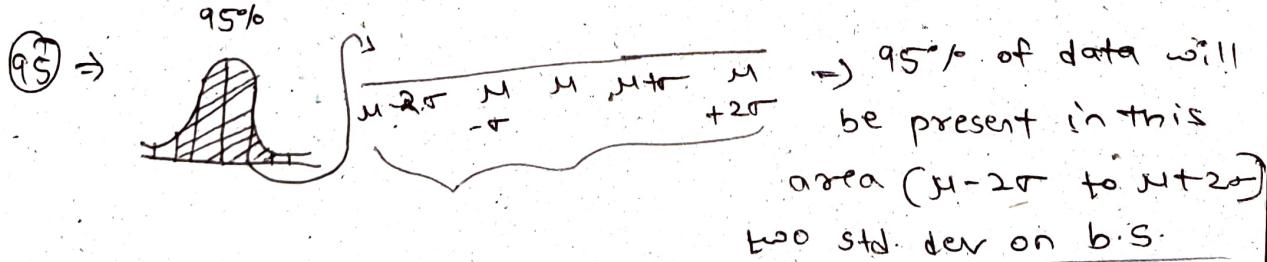
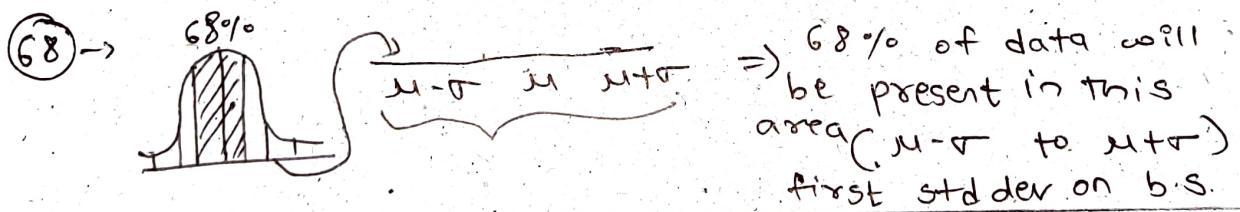
In symmetric distribution, μ is in center of x-axis
hence, area on left = area on right.

$\mu + \sigma \rightarrow \mu$ + one std dev. (right) \rightarrow

$\mu - 3\sigma \rightarrow \mu$ - three std dev. (left) \leftarrow

Empirical Rule: [3 sigma rule]

$$68 - 95 - 99.7\%$$



How to find out if a distribution is normal dist or not?

Q-Q plot \Rightarrow we can find out whether a dist is normal/gaussian distribution or not.

Probability

$X \rightarrow$ random var

$$P(\mu - \sigma \leq X \leq \mu + \sigma) \approx 68\%$$

$$P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \approx 95\%$$

$$P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) \approx 99.7\%$$

e.g. weight, height, iris dataset

Note: Anything beyond $\mu - 3\sigma$ & $\mu + 3\sigma$, we generally consider as outliers