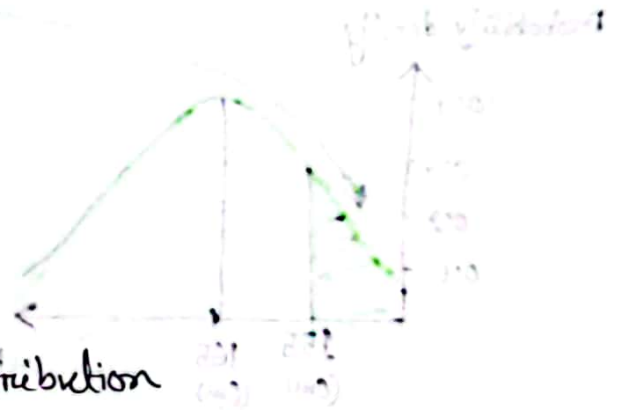


Advanced Statistics 01

Topics:

- (*) Probability Density / Distribution Function
- (*) PDF, PMF and EMF CDF
- (*) Type of probability Distribution
- (*) Bernoulli Distribution
- (*) Binomial Distribution
- (*) Poisson Distribution
- (*) Normal or Gaussian Distribution
- (*) End



(PDF)

Probability Distribution Function / Density Function:

- Actually it denotes the distribution of data

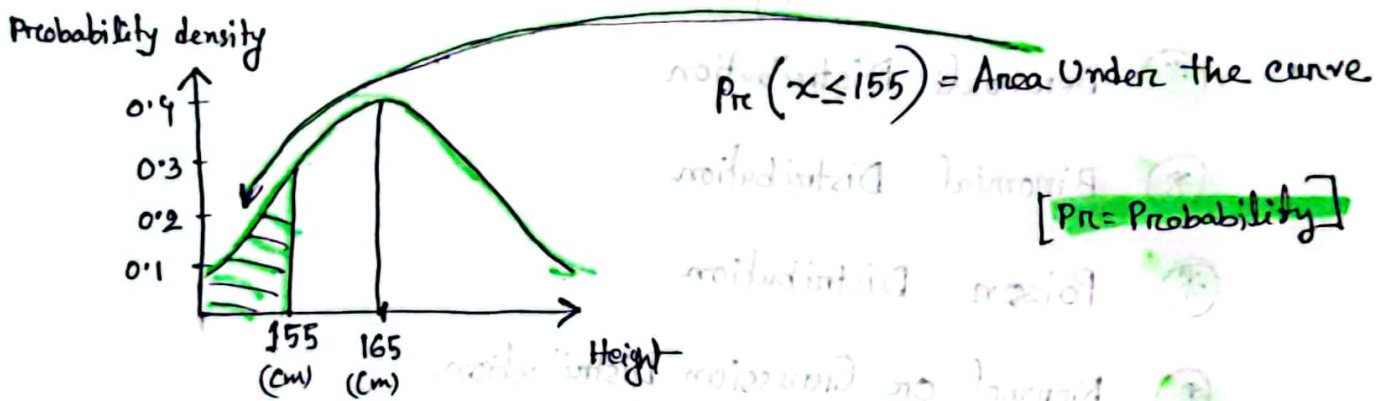
Types

→ ① Probability Density Function: (PDF)

→ Continuous Random variable

(Their distribution denotes by PDF)

Example → Height of students
in classroom



→ ② Probability Mass Function: (PMF)

→ When the variable is discrete random variable, their distribution denotes by PMF

Example → Rolling a dice

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

Histogram:



What is σ

$$Pr(x \leq 4)?$$

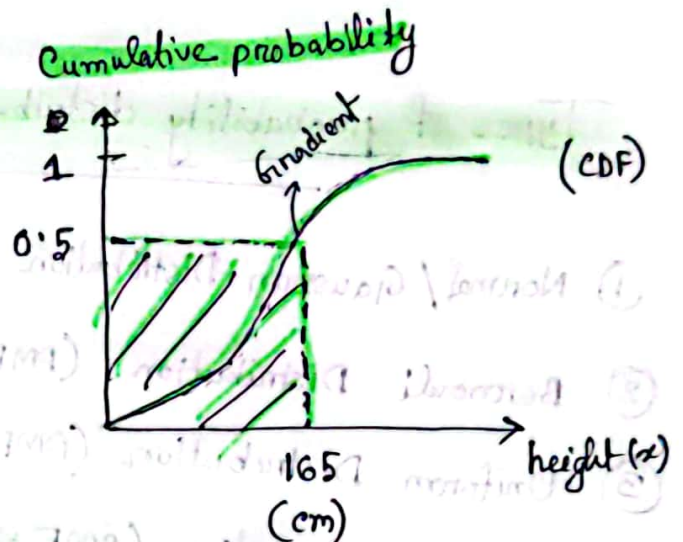
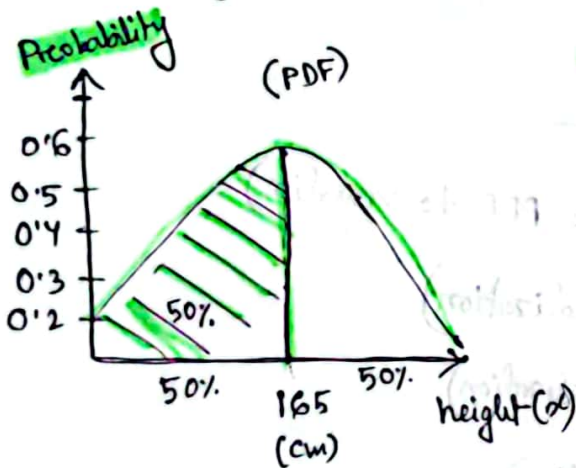
$$= Pr(x=1) + Pr(x=2) + Pr(x=3) + Pr(x=4)$$

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

$$= \frac{4}{6} = \frac{2}{3}$$

Cumulative Distribution Function:

Again the height example from PDF \rightarrow



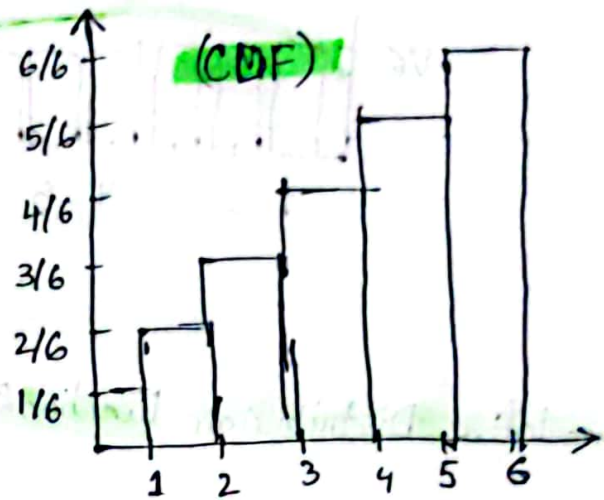
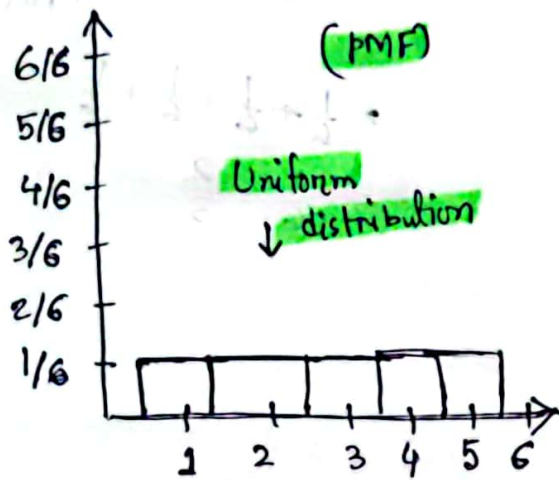
So, using the cumulative Distribution Function, we denote actually how much area height (x) has taken in the first half distribution.

Here in the left chart in the y axis, the values are denoting Gradients (steepness) of the cumulative curve (right) chart

In the right chart the y axis values are denoting area.

PMF relation with CDF:

Dice rolling probability



Types of probability distribution:

- ① Normal / Gaussian Distribution (use PDF to visualize)
- ② Bernoulli Distribution (PMF visualization)
- ③ Uniform Distribution (PMF visualization)
- ④ Poisson Distribution (PMF visualization)
- ⑤ Log normal Distribution (PDF visualization)
- ⑥ Binomial Distribution (PMF visualization)

Bernoulli Distributions

Discrete probability distribution of a random variable which takes the value 1 with probability p and the value 0 with probability $q = 1 - p$

In a easy way it can be said that,

It is a model for the set of possible outcomes of any single experiment that asks a yes-no question.

Key things to remember in Bernoulli Distribution:-

① Discrete Random Variable {PMF}

② Outcomes are binary $[(0, 1), (\text{Head}, \text{tail}), (\text{yes}, \text{no})]$

③ Example - Tossing a coin.

$$Pr(\text{Head}) = 0.5 \quad (\text{suppose } p)$$

$$Pr(\text{Tail}) = 1 - 0.5 = (1 - p)$$

④ Example \rightarrow Pass or fail in exam

$$Pr(\text{Pass}) = p = 0.7$$

$$Pr(\text{Fail}) = q = 1 - p = 1 - 0.7 = 0.3$$

Parameters:

$$0 \leq p \leq 1$$

$$q = 1 - p$$

$$K = \{0, 1\} \quad \text{or } [\text{yes or no}]$$

① PMF: (Probability Mass Function)

$$PMF = P^K * (1-P)^{1-K}$$

Hence,

$$K \in \{0, 1\}$$

$$\text{if } K=1, P_K(K=1) = P^1 * (1-P)^{1-1}$$

$$= P * 1$$

$$= P$$

$$\text{if, } K=0 \text{ then } P_K(K=0) = P^0 * (1-P)^{1-0}$$

$$= 1-P$$

PMF simplified:

$$PMF = \begin{cases} 1-P & \text{if, } K=0 \\ P & \text{if, } K=1 \end{cases}$$

Mean of Bernoulli Distribution:

$$E(K) = \sum_{k=0}^1 K \cdot P(K)$$

if, suppose

$$P_K(K=1) = 0.6$$

$$P_K(K=0) = 0.4$$

$$= [0 \times P(0) + 1 \times P(1)]$$

$$= [0 \times 0.4 + 1 \times (0.6)]$$

$$= 0.6$$

$$= P$$

So, p is the mean of Bernoulli Distribution

Median of Bernoulli Distribution:

$$\text{Median} \begin{cases} 0 & \text{if } p < 0.5 \\ [0,1] & \text{if } p = 0.5 \end{cases}$$

Based
on problem
Statement

$$1 \text{ if } p > 0.5$$

Variance:

Std:

$$\begin{aligned} \text{Var} &= p * (1-p) \\ &= p * q \quad [q = 1-p] \end{aligned} \quad \text{Std} = \sqrt{pq}$$

Binomial Distribution: $B(n, p)$

The binomial distribution with parameters n and p is the discrete probability distribution of the number of successes in a sequence of n independent experiments, each asking yes-no question, and each with its own boolean valued outcome: Success (p) or failure ($q = 1-p$).

Parameters:

n = number of trials

Suppose I am tossing a coin 7 times

So, trial $n = 7$

p = probability $[0, 1] \rightarrow$ Success probability for each trial

$$q = 1 - p$$

Key things to note:

⑥ For discrete random variable (PMF Function)

① Every outcome is binary

② This experiment is performed for n trials.

where each trail is a Bernoulli distribution

③ Every single trial from the n trial is called

Bernoulli distribution.

Support:

$K = \{1, 2, 3, \dots, n\} \rightarrow$ Number of success for n trials

PMF:

$$P_K(K, n, p) = {}^n C_K p^K (1-p)^{n-K} \quad \left\{ \begin{array}{l} K \in 0 \rightarrow n \end{array} \right.$$

Mean:

$$np,$$

Variance:

$$npq$$

Std:

$$\sqrt{npq}$$

Poisson Distribution:

\rightarrow For Discrete Random variable (PMF Function)

\rightarrow Describes the number of events in a fixed time interval

Example: Number of people (n) visiting hospital every hour
 \rightarrow fixed time



$$\lambda = 3$$

λ = Number of people expected in medical in each hour

Some statistical question that can be achieved:

1) What is the ^{probability of} number of people at 5th hour

2) What is the probability of a particular person to come at 5th hour?

PMF: $P_n(x=5) = \frac{e^{-\lambda} \lambda^x}{x!}$

$$= \frac{e^{-3} 3^5}{5!}$$

$$= 10.1\%$$

$x = 5$ th hour
 $\lambda = 3$

Mean: $\mu = \lambda t$ $\left[\begin{array}{l} \lambda = \text{Number of expected events occurred at every} \\ \text{time interval} \\ t = \text{Time interval} \end{array} \right]$

Variance: $\text{Var} = \lambda t$ (same)

Normal or Gaussian Distribution: $N(\mu, \sigma^2)$

→ Continuous probability distribution for a real-valued random variable

→ Mean = Median = Mode

Parameters:

$\mu \in \mathbb{R} = \text{mean}$

[$\mathbb{R} = \text{Real numbers}$]

$\sigma^2 \in \mathbb{R} > 0 = \text{variance}$

$x \in \mathbb{R}$

PDF:

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

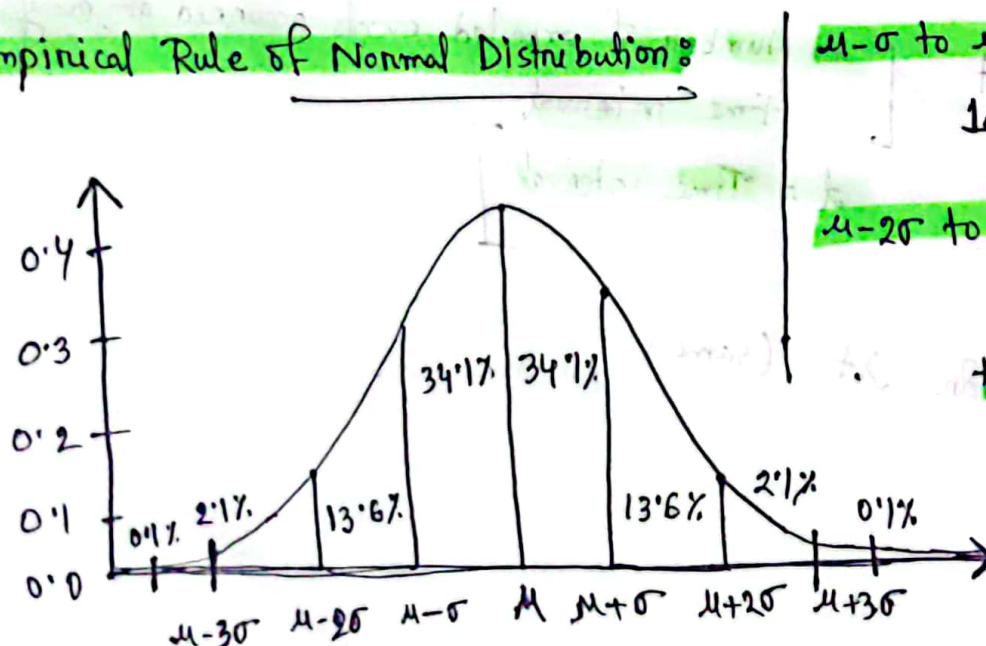
Mean:

$\mu = \text{Average value}$

Variance: $\text{Var } \sigma^2$

Std: $\sqrt{\sigma^2}$

Empirical Rule of Normal Distribution:



$\mu - \sigma$ to $\mu + \sigma \rightarrow$

1st std Region

$\mu - 2\sigma$ to $\mu + 2\sigma \rightarrow$

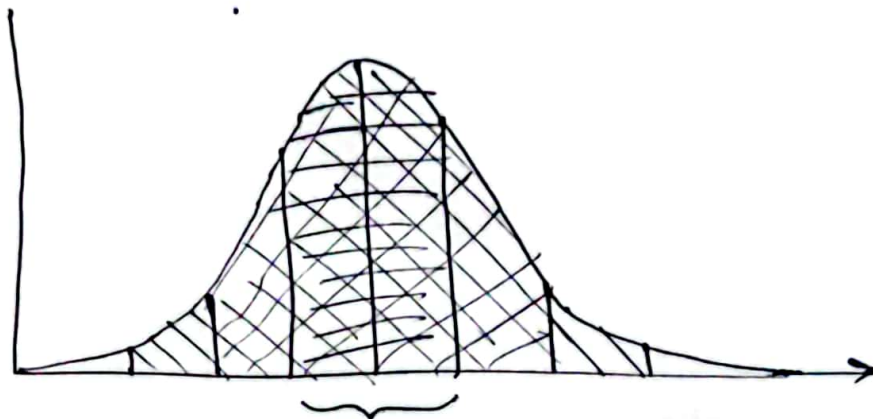
2nd std Region

then 3rd std Region

In 1st std Region \rightarrow There is 68% data distribution available

In 2nd std Region \rightarrow There is 95% of data distribution available

In 3rd std Region \rightarrow There is 99% of data distribution available



1st Region \rightarrow (—) This sign area

2nd Region \rightarrow (/) This sign area

3rd Region \rightarrow (\) This sign area

In terms of probability

for 1st std Region $\rightarrow \Pr(\mu - \sigma \leq x \leq \mu + \sigma) \approx 68\%$

2nd std Region $\rightarrow \Pr(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 95\%$

3rd std Region $\rightarrow \Pr(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 99\%$

Data that follows Normal ~~std~~ Distribution:

- 1) Height, weight
- 2) Irish