



Statistics

Through some tools
to a problem
statement

Defn: Statistics is the science of collecting,
organising and analysing the data.

Decision making process

Data: "facts or pieces of information".

e.g. heights of student in the class

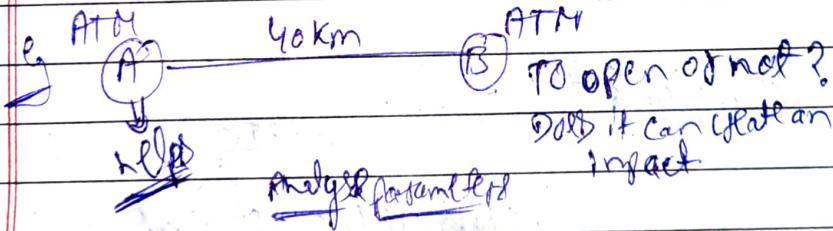
{ 175cm, 180cm, 190cm, ... }

if Data, it can be gathered
record and it can
be measured.

IQ of me student

{ 85, 90, 100, ... }

Conclusion



Types of statistics:

① Descriptive stats

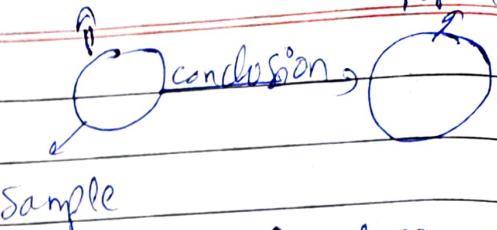
② Inferential stats

Defn: It consists of organizing
and summarizing the data

Defn: It consists of using data
you have measured to form
conclusion.

use sample to make conclusion
on population data
Population

① Measure of central tendency
[Mean, Median, Mode]



② Measure of Dispersion
[Variance, SD]

① Z-test [conclusions]
② t-test
hypothesis testing, p value,
significance.

③ Histograms, Bar chart,
Pie charts, Scatterplotlik.
distribution of data

→ All the industries used these techniques for conclusions

If let say there are 50 students in a math class in the university. we have collected the height of the students in the class

[Data]

[175cm, 180cm, 160cm, 140cm, 130cm, 140cm, 140cm, ...]

Descriptive Question

$$\frac{175 + 180 + 160 + 140 + 130 + 140 + 140}{50} = \text{mean}$$

What is the average height of the students in the class?

What is the common height of the students? = 140cm

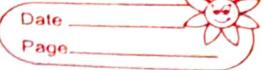
Inferential Question

Sample data

Are the average height of the students in the classroom similar to what you expect in the entire college?

Population

work on sample data to come up with a soln for the population data.

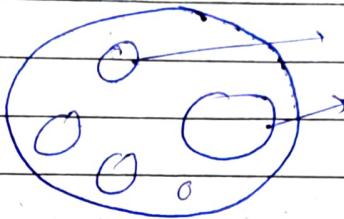


Sample Data And Population DATA

e.g. Exit poll

party A 56%

party B 46%



locations

ASK people to whom they vote.

Acc to that collection of data
make conclusions.

Types of DATA

DNC	DC	LSI	BVI	FWI	classes	Region
3.4	7.6	6.3			not fire	0
9.1	7.6	1			not fire	0
2.5	7.1	0.3			not fire	0
10.3	6.9	0			not fire	0
3	14.2	1.2			not fire	0
5.8	22.2	3.1			fire	0
9.9	30.8	6.4			fire	0
12.1	38.3	5.6			fire	0

1st step in EDA

numerical

(DATA)

categorical data

e.g., $\{x_1, x_2, \dots\}$ Quantitative

Qualitative

(Discrete)

(continuous)

(Nominal)

(Ordinal)

whole no.
specific range

Any value

e.g. gender

gr no. of bank acc

e.g. weight,
height

H, F

gr good @

No. of children in

temp,

blood grp

③ Bad

a family

speed

color

② satisfactory



Assignment 8

- ① what kind of variable Marital status is ?
- ② what kind of variable Nile River length is ?
- ③ what " " " " is Movie duration ?

Scale of Measurement of Data

- ① Nominal scale Data
- ② Ordinal scale Data
- ③ Interval scale Data
- ④ Ratio scale Data

① Nominal Scale Data

- ① Qualitative/categorical variable
- ② e.g. Gender, colors, labels
- ③ Order does not matter

↳ Survey/Logs

$$\left\{ \begin{array}{l} \text{Red} \rightarrow 5 \rightarrow 50\% \\ \text{Blue} \rightarrow 3 \rightarrow 30\% \\ \text{Yellow} \rightarrow 2 \rightarrow 20\% \end{array} \right.$$

② Interval Scale Data

③ Ordinal Scale Data

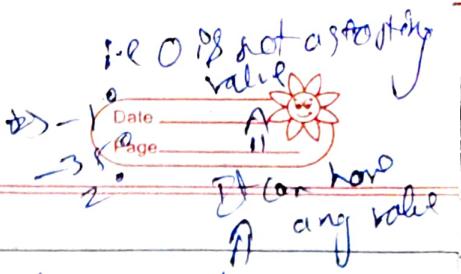
- ① Ranking and order matters
- ② Difference cannot be measured

↳ Qualification

PHD	1 st	→ 100
BSc	2 nd	→ 90
Masters	2 nd	→ 90
BCom	3 rd	→ 80
SSC	4 th	→ 70

If ranking is important

We can convert this into ranks.



③ Interval Scale Data

① The rank and order matters

② Difference can be measured (excluding ratio)

③ Does not have "0" starting value

Whatever you age feeling on 30°

That does not mean you will fell

double on 60° [Not considering Ratio]

④ Ratio Scale Data

Crudest

① order and Rank matter

0, 100, 90, 60, 30, 45

② Difference and ratio are measurable

③ It does have a "0" starting point

100, 90, 160, 150, 45

$$100:50 \Rightarrow 2:1$$

Ex:

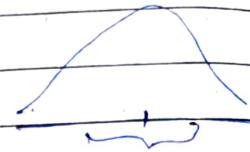
(a) length of different rivers in the world?

(b) residential status?

(c) ICL measurement?

① Measure of central tendency :-

- ① Mean
 - ② Median
 - ③ Mode
- { FDA and
feature engineering }



① Mean

Population (N)

Sample (n)

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

$(n=10)$

$$\text{population Mean } (\mu) = \frac{\sum_{i=1}^N X_i}{N}$$

sample mean (\bar{x})

$$= \frac{\sum_{i=1}^n X_i}{n}$$

$$\mu = \frac{1+1+2+2+3+3+4+5+5+6}{10}$$

$$= \frac{32}{10} = 3.2$$

②

Median

→ central element of the distribution

4, 5, 2, 3, 2, 1

soot → 1, 2, 2, 3, 4, 5

Median even counts

1, 2, 3, 4, 5

Odd count

1, 2, 2, 3, 4, 5, 7

$$\frac{2+3}{2} = 2.5$$

$$\text{Median} = 2.5$$

(Median = 3)

Why median?

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

$$S = \frac{(1+2+3+4+5)}{5} = \frac{15}{5} = 3 \text{ [Mean]}$$

f outliers

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

$$\text{Median} = 3$$

$$\{1, 2, 3, 4, 5, \underline{100}\} \rightarrow \text{outlier}$$

$$\{1, 2, \underline{3}, 4, \underline{5}, \underline{100}\}$$

$$S = \frac{1+2+3+4+5+100}{6} = \frac{115}{6} \approx 18.3 \text{ [Mean]}$$

$$\text{Median} = 3.5$$

It is a very big shift
in case of outlier

∴ use median in case of outliers.

(It is better to use median in case of outliers)

③

mode frequency maximum

$$\{2, 1, 1, 4, 5, 7, 8, 9, 10\}$$

$$\text{Mode} = 1$$

EDA and feature engineering

type of flower

Age

[to avoid loss of data]

Lily

Rose

16

3

5

mean or median
outliers

rose — sunflower

Rose



② Measure of Dispersion

① Variance

② Standard Deviation

③ Variance

Population variance

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

x_i = Data points

μ = population mean

N = population size

Sample variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Bessel's
correction

Why does the sample variance have $n-1$ in the denominator?
→ The reason we use $n-1$ rather than n is so that the sample variance will be what is called an unbiased estimator of the population variance.

x_i = Data points

\bar{x} → Sample Mean

n → Sample Size

Gr {1, 2, 3, 4, 5} → Sample

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\bar{x} = \frac{1+2+3+4+5}{5} = 3$$

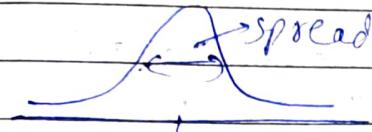
$$\begin{array}{ccc} x_i & \bar{x} & (x_i - \bar{x})^2 \\ 1 & 3 & 4 \\ 2 & 3 & 1 \\ 3 & 3 & 0 \\ 4 & 3 & 1 \\ 5 & 3 & 4 \end{array}$$

$$s^2 = \frac{10}{4} = 2.5$$

$$s^2 = \frac{10}{5} = 2$$

Variance: Spread of the data.

$$S^2 = 2.5$$



$$S^2 = 6.5 = 1.9$$



② standard Deviation

population std

$$\sigma = \sqrt{\text{variance}}$$

$$S^2 = 2.5$$

sample std

$$S = \sqrt{\text{sample variance}}$$

$$\sqrt{S^2} = \text{sample std}$$

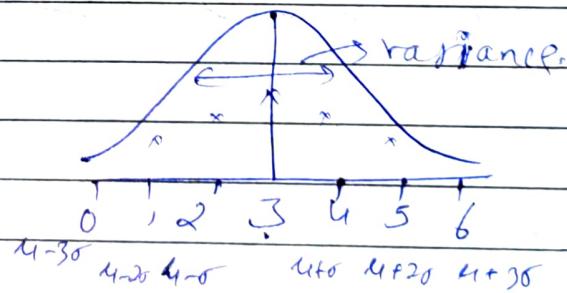
Consider

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

$$\hookrightarrow \mu = 3$$

$$\hookrightarrow \sigma = 1$$

Std



③

→ 1 std to the left & and right



③ Random variables

$$\begin{cases} x+5=7 \\ y+x=10 \end{cases}$$

value
 $x=2$
 $y=8$

But in random variable
↓

Random variable is a process of mapping the output of a random process or experiment to a number.

e.g. Tossing a coin

Rolling a dice

Predicting the temperature for the next day.

CAPITAL LETTERS

$$X = \begin{cases} 0 & \text{if H} \\ 1 & \text{if T} \end{cases}$$

Quantifying a Random

Process.

Converting a process
output or outcome into a specific number

$Y = \text{sum of the rolling of dice } T \text{ times}$

$$\{4, 5, 6, 1, 2, 2\} = 20$$

$$P(Y \geq 15), P(H)$$

(1) Sets

$$A = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

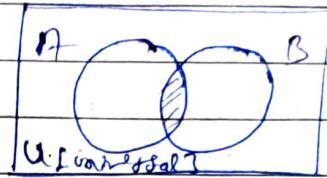
$$B = \{3, 4, 5, 6, 7\}$$

operations

(1) Intersection \cap

$$A \cap B = \{3, 4, 5, 6, 7\}$$

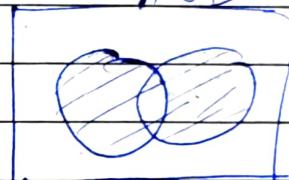
Venn diagram



(2) Union \cup

$$A \cup B = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

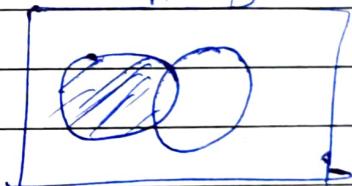
$A \cup B$



(3) Difference \setminus

$$A - B = \{1, 2, 8\}$$

$A - B$



(4) Subset \subseteq

$$A \rightarrow B \Rightarrow \text{FALSE}$$

$$B \rightarrow A \Rightarrow \text{TRUE}$$

(5) Superset \supseteq

$$A \rightarrow B \Rightarrow \text{TRUE}$$

$$B \rightarrow A \Rightarrow \text{FALSE}$$



① Covariance And Correlation

Covariance indicates the relationship of two variables whenever one variable changes.

If an increase in one variable results in an increase in the other variable, both variables are said to have a positive covariance.

Decrease in one variable also cause a decrease in the other.

X	Y
2	3
4	5
6	7
8	9

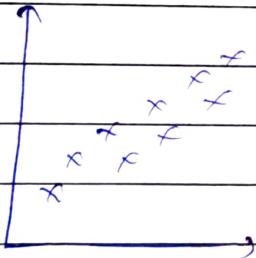
Relationship b/w X and Y

$X \uparrow$	$Y \uparrow$
$X \uparrow$	$Y \downarrow$
$X \downarrow$	$Y \uparrow$
$X \downarrow$	$Y \downarrow$

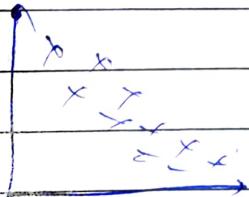
size location price

ML model

Helps in



$X \uparrow$	$Y \uparrow$
$X \downarrow$	$Y \downarrow$



$X \downarrow$	$X \uparrow$
$X \uparrow$	$X \downarrow$

Diff. b/w covariance and variance

→ It's the same thing
⇒ relationship strength

$$\text{Covariance } (X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

x^1	y^9
x^6	y^1

+ve cov.

var. of x is nothing
but its own
relationship

$$\text{var}(n) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

x^1	x^1
x^1	y^1

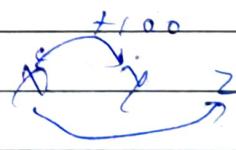
-ve cov.

$$\text{spread of the data.} \Leftarrow \text{cov}(x, x) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\text{cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\begin{array}{cc} x & y \\ 2 & 3 \\ 4 & 5 \\ 6 & 7 \end{array} \Rightarrow (2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5) = 2$$

$$\bar{x}=4, \bar{y}=5 \Rightarrow \frac{4+0+4}{3} = 4 \text{ +ve}$$



x & y are having a positive variance.

Advantages

- ① Relationship b/w x & y
+ve or -ve value

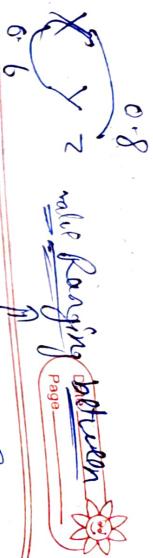
Disadvantage

- ② Covariance does not specific limit value

It can be any +ve or -ve value

To solve this disadvantage we use another kind of correlation technique

$r = 2$ is moderately correlated
man x.y.



② Pearson correlation coefficient $r \in [-1, 1]$

The Pearson coefficient is a type of correlation coefficient that represents the relationship b/w two variables that are measured on the same interval or ratio scale. The Pearson coefficient is a measure of the strength of the association b/w two continuous variables

$$r_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- ① The more the values move towards +1 the more the correlated it is (x, y)
- ② The more the values towards -1 the more we correlated it is (x, y) .

use case

Dataset - 1000 features (ML Models)

Independent Variables

Dependent Variable

size of the house	no. of rooms	no. of people	Price
large	inf	inf	nothing

re corr.

Feature selection near to 0 \Rightarrow Drop non-informative feature

feature

(3) Spearman Rank correlation

$$\gamma_s = \text{cov}(R(x), R(y))$$

$$\sigma_{R(x)} \times \sigma_{R(y)}$$

X	Y	$R(x)$	$R(y)$
1	2	5	5
3	4	4	4
5	6	3	3
7	8	2	1
0	7	6	2
8	1	1	6

Here we don't focus on values

Instead we compute the rank

of X & Y. and use them

to calculate cov.

import Seaborn as sns

```
df = sns.load_dataset('taxis')
df.head()
```

import numpy as np

np.cov

df.cov()

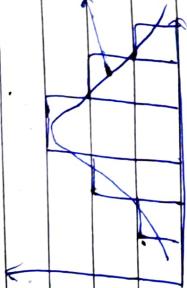
| df.cov(method='spearman')

df.cov(method='pearson')

df = sns.load_dataset('penguins')

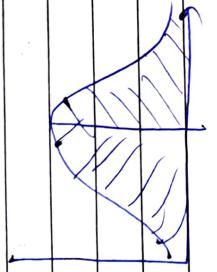
df.cov()

Skewness



→ Normal Gaussian Distribution

Symmetrical Distribution



① No skewness

For Symmetrical Distribution no skewness

Box plot

⇒ Normal Gaussian Dist.

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

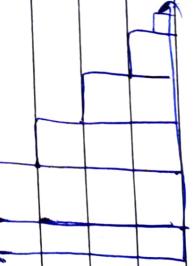
⇒ Median = Median mode

② Right Skewed or

Log normal
mode < median < mean

→ positive skewed

→ trend towards right



Relationship between median, mean, → Median > Mean
mode

mode > median
median > mode
mode > median > mean

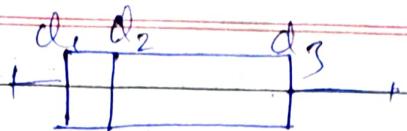
mean > median > mode

mode > peak

median > mode > peak



Box plot



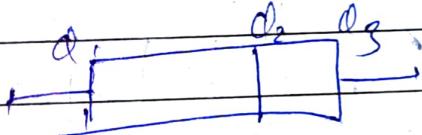
$$Q_3 - Q_2 \geq Q_2 - Q_1$$

(3) left skewed distribution:



→ negative skewed

Relationship b/w mean ≤ median ≤ mode



$$Q_2 - Q_1 \geq Q_3 - Q_2$$

Q) what is a relationship b/w mean, median & mode

histogram & [frequency] \Rightarrow count

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

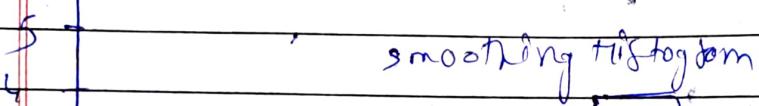
no. of \Rightarrow Buckets
Bins = 10

$$\text{binsize} = 5 = \frac{50 - 0}{10} = 5$$

\Rightarrow How many elements are present in a Bucket

$$\frac{50}{20} = 2.5 \Rightarrow \text{binsize} \\ \therefore \{20 \text{ no. of bins}\}$$

Count n

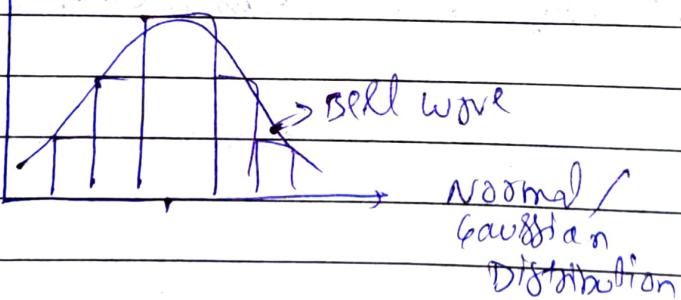


PDF (Probability distribution function)

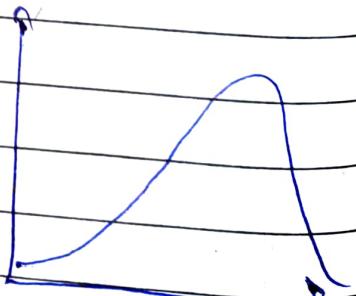
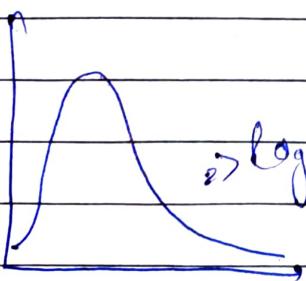
[Probability density function]

\Rightarrow Smoothening histogram

Kernel density estimator.



Further you know about data (you can make a lot of assumption about data)



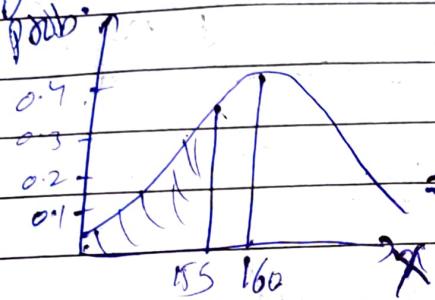
\Rightarrow When we smoothen it we can get different distribution.

Prob. Density funcn is a part of Prob. Distribution function



- ① probability Distribution Function / Density Function
- ② Probability Density Function (PDF) \Rightarrow Distribution of the continuous data

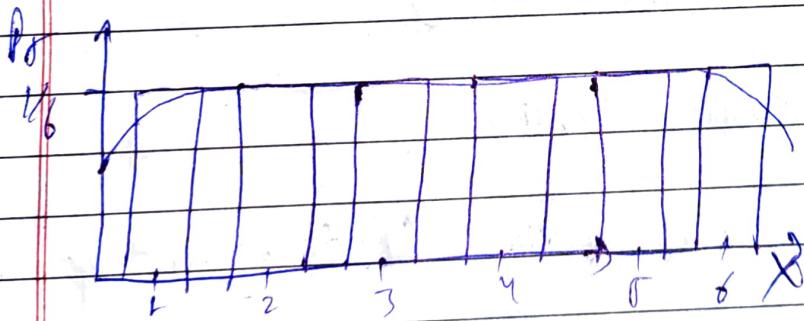
e.g. Height of Students



$P(X \leq 155)$ = Area Under the curve
0.3, 0.4, 0.5

\Rightarrow Kernel Density estimator

- ③ probability Mass Function (PMF) of Distribution of a Discrete Random variable
e.g. Rolling a Dice {1, 2, 3, 4, 5, 6}



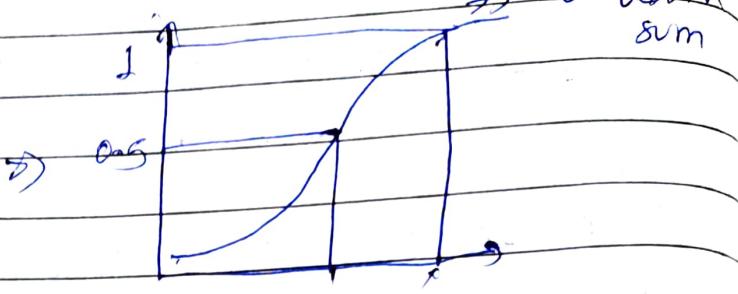
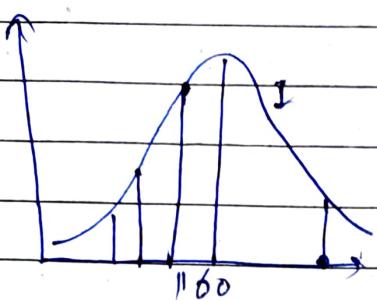
$$P(X \leq 4) = P(X=1) +$$

$$P(X=2) + P(X=3) + \\ P(X=4)$$

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

$$\Rightarrow 4 \cdot \frac{1}{6} = \frac{2}{3}$$

(3)

Cumulative Distribution Function (CDF) f

→ we are doing a cumulative sum of all the values w.r.t. the area under the wave

Types of the Probability Distributionmostly follow
in industry

(1) Normal / Gaussian Distribution (PdF)

(2) Bernoulli Distribution (PMF) → outcomes are 2
success or failure

(3) Uniform Distribution

(4) Log normal Distribution (PdF)

(5) Poisson Distribution (PMF)

(6) power law Distribution (PdF) → 80-20% rule

(7) Binomial Distribution (PMF) → discrete values

→ probabilities of various types, it may belong to a PdF or PMF.



Probability Density Function And probability mass function.

cumulative distribution function (CDF)

using the PDF & PMF we specifically create the realization graphically

PMF

[what is the formula used in the actual dist.]

○ Discrete Random variable

PMF

CDF

Ex: Rolling a dice

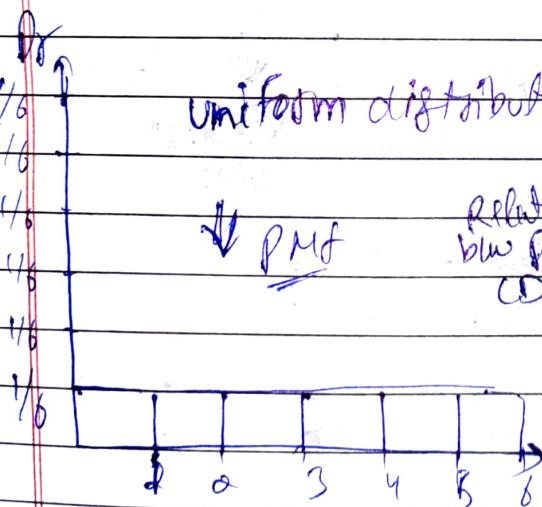
$$\{1, 2, 3, 4, 5, 6\}, P(X=1) = \frac{1}{6}$$

cumulative basic
means combining

Fair dice

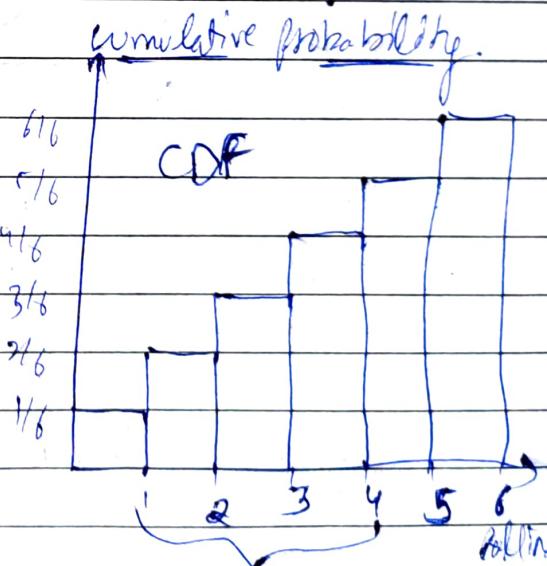
$$P(X=2) = \frac{1}{6}$$

cumulative probability.



$$P(X=1) = \frac{1}{6} \quad \left. \begin{array}{l} \text{so for } \\ \text{single} \end{array} \right\}$$

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



\Rightarrow we try to add next events
 \Rightarrow we try to cumulatively add next events

$$\text{The } P(X \leq 4) = P(X=1) + P(X=2) + P(X=3) + P(X=4)$$

we have to add here,

$$\Rightarrow \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{4}{3}$$

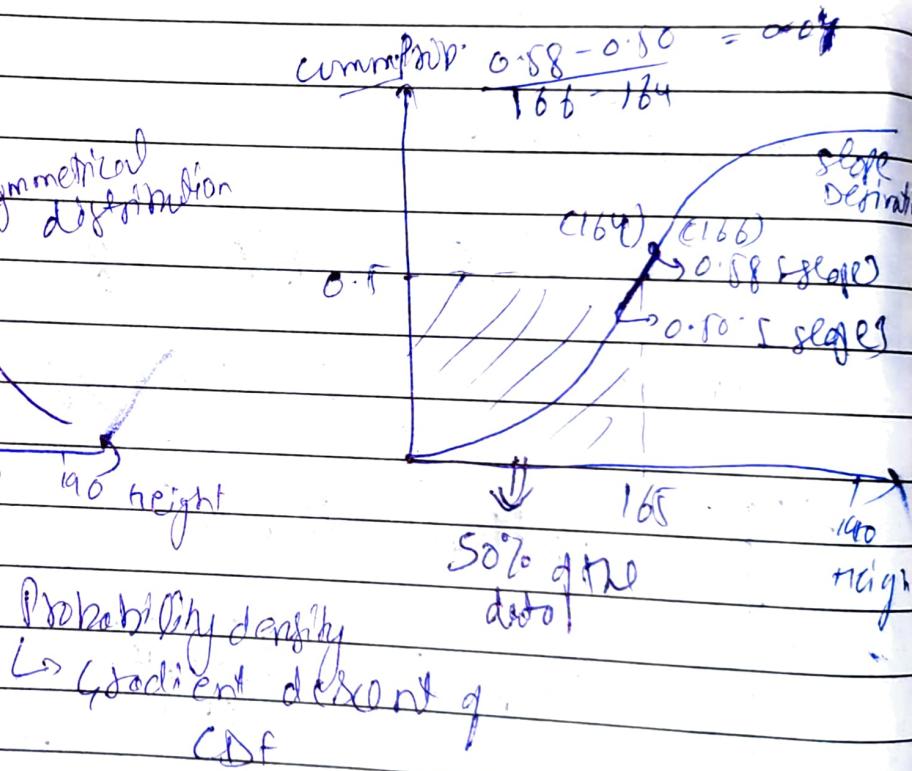
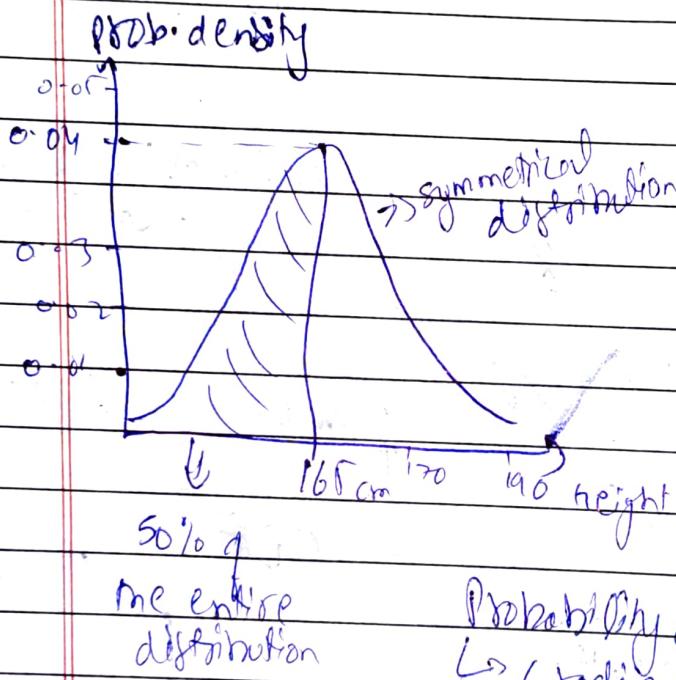
① $P(X \leq 6) = P(X=1) + P(X=2) + P(X=3)$
 $\quad \quad \quad + P(X=4) + P(X=5) + P(X=6)$

↓
cumulative prob.
 $= 1$

② Probability Density function (PDF)

application

① continuous random variable



Probability density at 185 cm = 0.06



Bernoulli Distribution

[Binary outcomes]

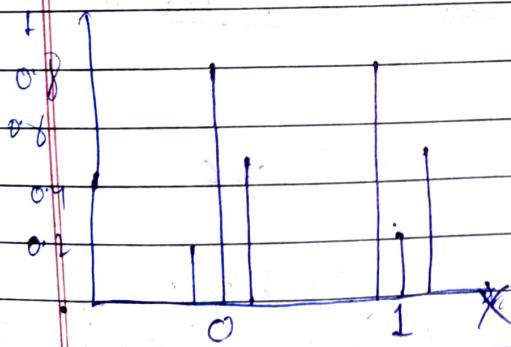
In probability theory and statistics, the Bernoulli distribution, named after the Swiss mathematician Jacob Bernoulli, is the discrete probability distribution of a random variable which takes the value 1 with probability p and the value 0 with probability $q = 1-p$.

Less formally, it can be thought of as a model for the set of possible outcomes of any single experiment that asks a yes-no question. Such questions lead to outcomes that are boolean-valued: a single bit whose value is success (yes/one) with probability p and failure/no/false/zero with probability q .

Bernoulli Distribution

Prob. Mass function.

outcomes are binary



e.g., tossing a fair coin H T

$$P(T) = 0.5 = p$$

$$P(H) = 1 - 0.5 = 1 - p = q$$

e.g., whether the person

Pass/fail.

Three examples of Bernoulli dist.

PMF

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

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

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

Q PMF vs PDF

Discrete
Random
variable

Continuous
random
values

PMF vs PDF

P

P

Discrete

continuous

random

random variable

variable

mathematically - Pdt funn eqn

PDF

K=0 or 1

$$P(X=K) = p^K (1-p)^{1-K} \rightarrow \text{PMF}$$

simplified way

$$\textcircled{1} \quad P(X=1) \Rightarrow p^1 (1-p)^0$$

$$P(X=0) \Rightarrow p^0$$

$$\text{PMF} = \begin{cases} q = 1-p & \text{IF } K=0 \\ p & \text{IF } K=1 \end{cases}$$

$$\textcircled{2} \quad P(X=0) \Rightarrow p^0 (1-p)^{1-0}$$

$$= (1-p) = q$$

\textcircled{3} Mean, variance and standard Deviation

Mean

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

K=1 or 0

Expected value q.k

$$E(K) = 1 \times 0.6 + 0 \times 0.4$$

$$\Rightarrow 0.6 = p$$

$$P(K=1) = 0.6 = p$$

$$P(K=0) = 1 - 0.6 = q = 0.4$$



Median of Bernoulli Distribution

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

variance $\Sigma \text{std } \epsilon$ $p=0.5 \quad q=0.5$

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

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