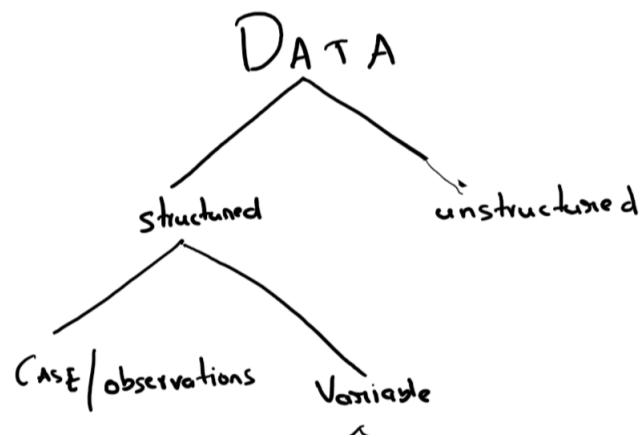
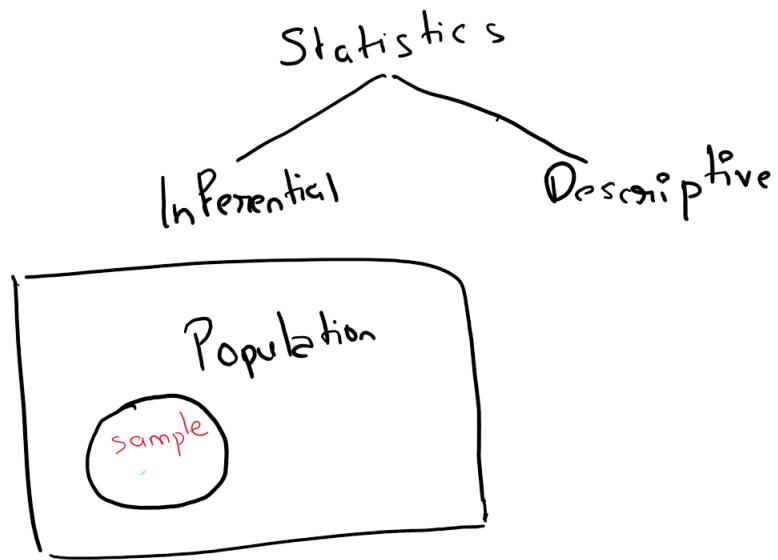
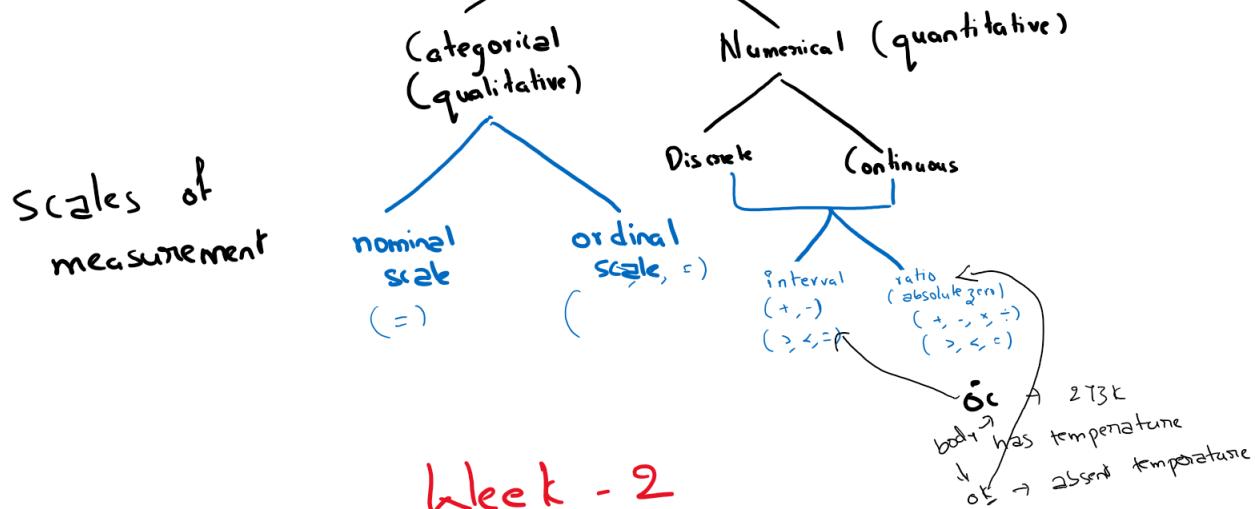


Statistics for Data Science - 1

Notes for end term



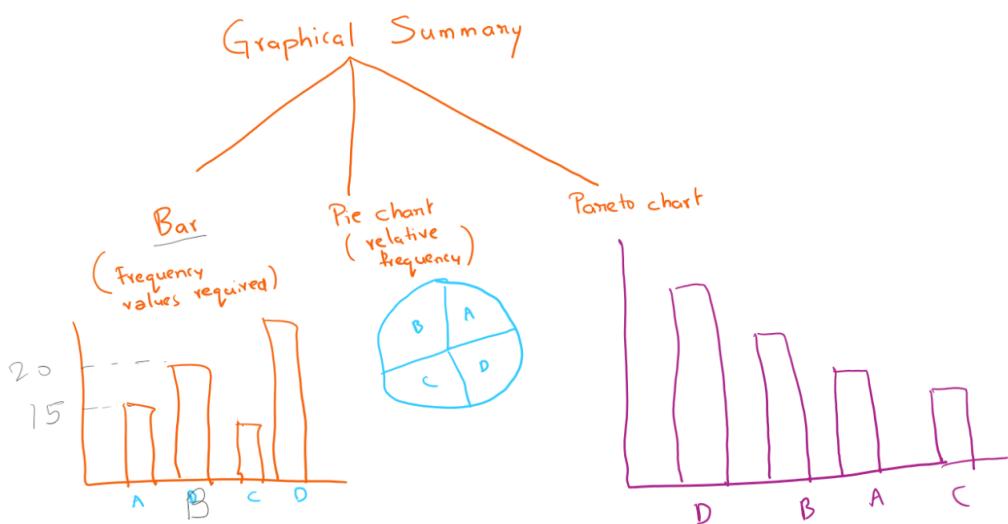


Week - 2

CATEGORICAL VARIABLE

Frequency?
Relative Frequency
Frequency Table

Categories	Frequency	relative frequency
A	15	15/60
B	20	20/60
C	15	15/60
D	10	10/60
	$\frac{60}{60}$	
Total Frequency		50



Describing categorical data

Mode
(ordinal, nominal)
highest frequency category

Median
(ordinal)

Middle value in ordered data ascending or descending order

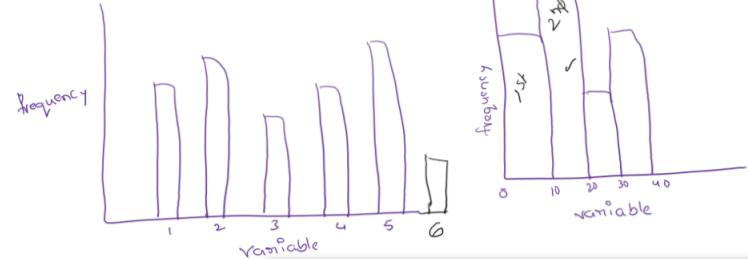
Week 3: Describing numerical variable



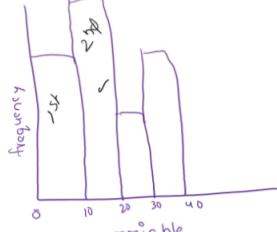
Sign in

Graphical summary

Bar chart
(for discrete less numbered variable)



Histogram
(for longer data)
Data divided into classes



Stem and Leaf plot

Stem	leaf
1-	2 9
2-	1
3-	5
4-	1

(12, 41, 21, 19, 35)

{0, 10)
{10, 20)

Descriptive measures



Measures of central tendency

Mean

$$\text{Mean } \bar{x} = \frac{\sum x_i}{n}$$

Frequency: $f_1 = 2, f_2 = 3, f_3 = 4, f_4 = 5, f_5 = 1$

$$\bar{x} = \frac{1 \times 1 + 2 \times 3 + 3 \times 4 + 4 \times 5}{1+3+4+5} = \frac{40}{15} = \frac{8}{3}$$

Median

middle value in ordered data.
median: $\frac{n+1}{2}$ th observation for n-odd
average of $\frac{n}{2}, \frac{n}{2}+1$ observations n-even

Mode

highest frequency value

Measures of dispersion

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$

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

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

$n-1 \rightarrow$ sample standard deviation (s)

$n \rightarrow$ population " " " (σ)

$$\text{Range} = \max - \min$$

$$\text{IQR: } Q_3 - Q_1$$

Inter Quartile Range

Percentiles

pth percentile is npth observation if np is exact integer
" " " next highest integer & np observation if np is not exact integer

ordered ascending

Q₁ - 25th percentile value

n → total observations

Q₃ - 75th percentile value

n × p

Q₁ = n × 0.25th observation

Q₃ = n × 0.75th observation

Interguatile range (IQR): Q₃ - Q₁

Outliers \rightarrow extreme values less than $Q_1 - 1.5 \text{ IQR}$ and values greater than $Q_3 + 1.5 \text{ IQR}$

-8 2 3.5 4.3 5.1 6.2 greater 8.5 9.1 19

\downarrow
outlier

$$n = 10,$$

$$nP = 10 \times 0.25 = 2.5$$

$Q_1 = 3^{\text{rd}}$ observation in ordered data = 3.5

$Q_3 = np = 10 \times 0.75 = 7.5$ next highest integer = 8.5

$Q_3 = 8^{\text{th}}$ observation = 8.5

$$\text{IQR} = 8.5 - 3.5 = 5$$

values $< Q_1 - 1.5 \text{ IQR} = 3.5 - 1.5(5) = -4$

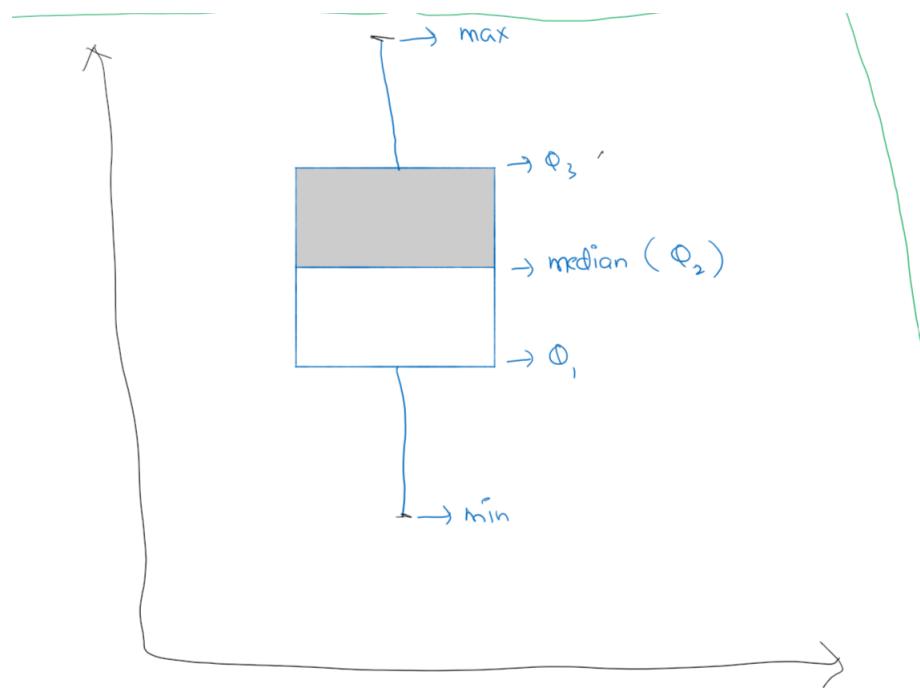
values $> Q_3 + 1.5 \text{ IQR} = 8.5 + 1.5(5) = 16$

outliers

0
0

-8, 19

are outliers



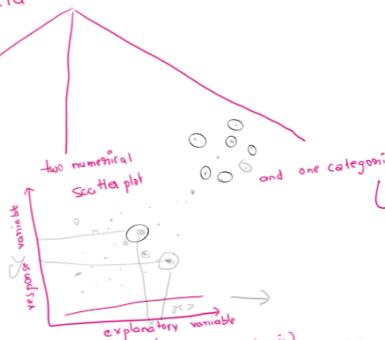


Week 4 Association between two variables

two categorical

Variable

	employed	unemployed	Total
female	100	50	150
Male	500	150	650
Total	600	200	800



and one categorical \rightarrow two categories

point, bi-variate correlation coefficient (r_{bp}):

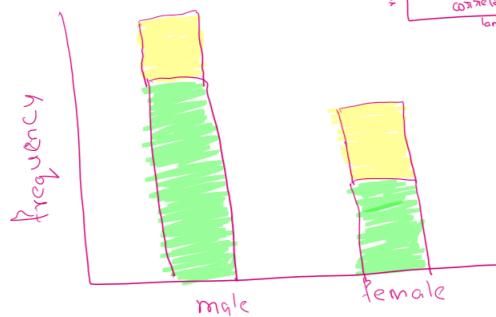
$$r_{bp} = \frac{(\bar{y}_0 - \bar{y}_1) \sqrt{P_0 P_1}}{S_d}$$

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

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

$$\text{Correlation coefficient } r = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y} = \frac{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2} \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2}}$$

r is same for sample, population



Stacked bar chart

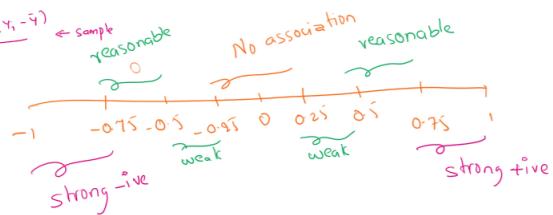
correlation coefficient (r)

Explanatory variable

$$r = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

$$r = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y} = \frac{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2} \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2}}$$

r is same for sample, population



Week 5 Permutations and combinations



Basic rules of counting

Addition rule

event A → n_1 ways

event B → n_2 ways

A or B in $n_1 + n_2$ ways

Multiplication rule

event A → n_1 ways

event B → n_2 ways

$A \in B$ in $n_1 \times n_2$ ways

... ... \vdots

A or B "

There are $n_1, n_2, n_3, n_4, \dots, n_i$ in which i events occurs

Any i events occurs in $n_1 + n_2 + n_3 + n_4 + \dots + n_i$ ways (Addition rule)

All i events occurs can occur in $n_1 \times n_2 \times n_3 \times n_4 \times \dots \times n_i$ different ways. (Multiplication rule)

Factorials

$n! = n \times (n-1) \times (n-2) \times \dots \times 1$ ways

$n! = n \times (n-1)!, \quad (n-1)! = (n-1) n!$

" " " " " in distinct objects = $n_{Pr} = \frac{n!}{(n-r)!}$

$$0! = 1$$

ROTC

All & "

Factorials

$n! = n \times (n-1) \times (n-2) \times \dots \times 1$ ways

$n! = n \times (n-1)!, \quad (n-1)! = (n-1) n!$

$$0! = 1$$

Number of possible permutations i objects from n distinct objects = $n_{Pr} = \frac{n!}{(n-r)!}$

" " " " " =

$n_{Pr} = n!$

ROOTAG → 4 letters

{ not allowed
if repetition
is not allowed

All & "

Factorials

$n! = n \times (n-1) \times (n-2) \times \dots \times 1$ ways

$n! = n \times (n-1)!, \quad (n-1)! = (n-1) n!$

Number of possible permutations i objects from n distinct objects = $n_{Pr} = \frac{n!}{(n-r)!}$

" " " " " =

$n_{Pr} = n!$

Number of possible permutations

of n objects of n objects are of 1 kind: $\frac{n!}{n_1! n_2! \dots n_r!}$

from n objects = n^r → repetition allowed.

Number of possible permutations

from n objects = n^r → repetition allowed.

$\frac{n_{Pr}}{r!}$ ← repetition

" " " " " =

from n objects = n^r → repetition allowed.

Number of objects: n
 Number of possible permutations of n distinct objects = n^r → repetition allowed.
 " " " "
 Number of possible circular permutation of n distinct objects: $\frac{n P_r}{n} = \frac{n!}{(n-r)!} = (n-1)!$
 " " "

COMBINATIONS:
 Number of ways of choosing r objects from n objects is $n C_r = \frac{n!}{(n-r)! r!}$
 $n C_r = n C_{n-r}$

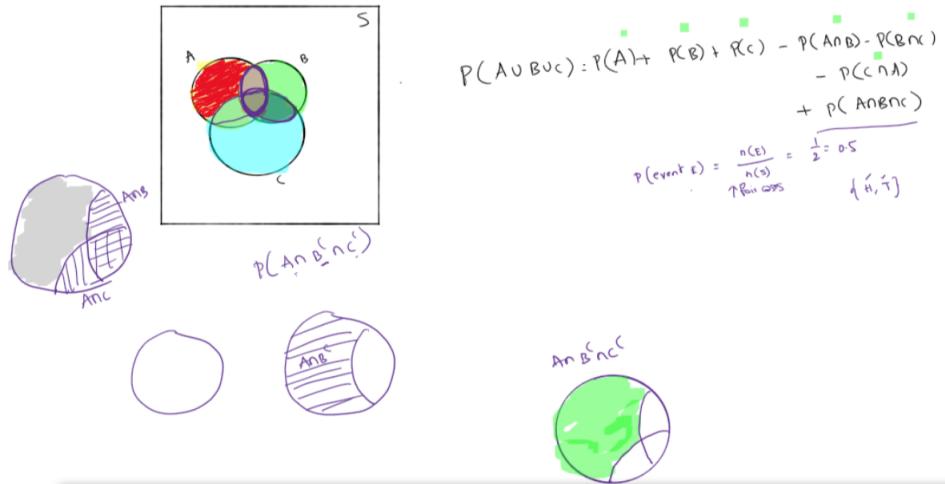
Probability

Random experiment

\downarrow
 All outcomes → sample space (S) $P(S)=1$
 $(E \text{ event } \subset \text{ sample space } S) \Rightarrow P(E) \leq 1$
 tossing a coin → All outcomes $\{H, T\} \Rightarrow S = \{H, T\}$
 event $E \rightarrow$ getting head $E = \{H\} \quad E \subset S$.

Union of events: (\cup)

$A \in B$ are two events,
 probability that either A or B occurs is $P(A \cup B)$
 $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ \cup
 $\Rightarrow P(A \cap B) = P(A) + P(B) - P(A \cup B)$
 null event → \emptyset $P(\emptyset) = 0$
 $\emptyset: \{\}$ empty set



Week - F

① Contingency table

		Employed		Unemployed
		G	H	Row total
Male	E	a	b	a+b
	F	c	d	c+d
Column total		a+c	b+d	a+b+c+d = N

$$\begin{aligned}
 N &= \text{total no. of ways} \\
 a &= E \cap G \\
 b &= E \cap H \\
 c &= F \cap G \\
 d &= F \cap H
 \end{aligned}$$

Joint Probabilities

$$\begin{aligned}
 P(E \cap G) &= a/N \\
 P(E \cap H) &= b/N \\
 P(F \cap G) &= c/N \\
 P(F \cap H) &= d/N
 \end{aligned}$$

Marginal Probabilities

$$\begin{aligned}
 P(E) &= (a+b)/N \\
 P(F) &= (c+d)/N \\
 P(G) &= (a+c)/N \\
 P(H) &= (b+d)/N
 \end{aligned}$$

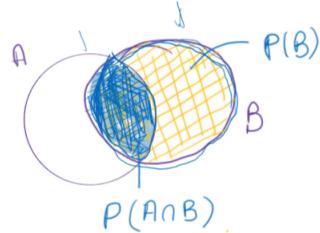
Conditional Probability

$$\begin{aligned}
 P(E|G) &= a/(a+c) \\
 P(F|G) &= c/(a+c) \\
 P(G|E) &= a/(a+b) \\
 P(H|E) &= b/(a+b) \\
 P(H|F) &= ? \quad d/(c+d) \\
 P(G|F) &= ?
 \end{aligned}$$

② Conditional Probability formula

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$P(B) > 0$



$$\left. \begin{aligned} P(A \cap B) &= P(A|B) P(B) \\ \downarrow \\ \text{Multiplication rule} \end{aligned} \right\}$$

③ Mutually Exclusive & Independent Events

\downarrow
disjoint

$$A \cap B = \emptyset$$

$$\Rightarrow P(A \cap B) = 0$$

\downarrow

$$\underline{P(A \cap B) = P(A)P(B)}$$

Let A & B be
& events for
a sample
space S.

→ Mutually exclusive but not independent:

① Random exp.: Tossing a coin

$$S = \{H, T\} \quad \checkmark$$

$$P(A) = \{H\}, B = \{T\} \quad \cancel{\checkmark}$$

$$A \cap B = \emptyset$$

$\Rightarrow P(A \cap B) = 0 \Rightarrow A \& B$ are mutually exclusive.

$$P(A) = \frac{1}{2}, P(B) = \frac{1}{2}, P(A)P(B) = \frac{1}{4} \neq 0 = P(A \cap B)$$

\Downarrow

not independent

Remark:

If A & B are 2 mutually exclusive events for a sample space S, A & B will be independent if either

$$P(A) = 0 \quad \text{or} \quad P(B) = 0.$$

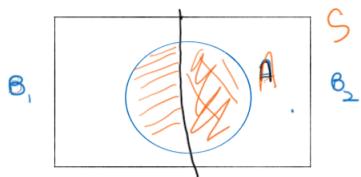
$$P(A) \times P(B) = \underline{\underline{P(A \cap B)}}$$

③ Law of Total Probability

Mutually Exclusive
Exhaustive events

① For 2 events

Let B_1 & B_2 be mutually exclusive & exhaustive events for a sample space S .



$$A = (A \cap B_1) \cup (A \cap B_2)$$

mutually exclusive events

For any event A ,

$$P(A) = P(A \cap B_1) + P(A \cap B_2)$$

(Using addition rule)

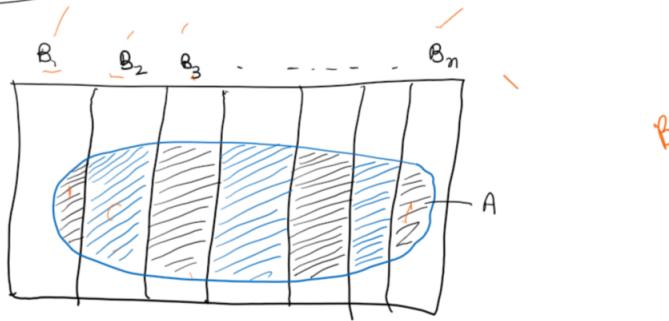
$$\Rightarrow P(A) = P(A|B_1)P(B_1) + P(A|B_2)P(B_2)$$

(Using Multiplication rule)

{ Law of Total Probability

h ↗

② For n events



$$A = (\emptyset \cap B_1) \cup (A \cap B_2) \cup \dots \cup (A \cap B_n)$$

$$\Rightarrow P(A) = P[(A \cap B_1) + P(A \cap B_2) + \dots + P(A \cap B_n)]$$

$$P(A) = P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + \dots + P(A|B_n)P(B_n)$$

④ Bayes' Theorem :

$$P(B_i | A) = \frac{P(B_i \cap A)}{P(A)}$$

$$= \frac{P(A|B_i)P(B_i)}{\sum_{i=1}^n P(A|B_i)P(B_i)}$$

Week 8

" RV "

Random Variable (Real valued function
on the sample space)

(Takes at most
countable number
of values)

Random
Discrete

Continuous

(When outcomes of random event
cannot be counted)

Eg. Experiment :- Rolling a fair die

Let X be a discrete random variable which takes value greater than 3.

$$S = \{ \underline{1}, \underline{2}, \underline{3}, \underline{4}, \underline{5}, \underline{6} \}$$

\times takes values 4, 5, 6.



Sign out



Probability mass function (p.m.f.)

Let X be a discrete random variable taking n possible values.

$x_1, x_2, x_3, \dots, x_n$

p.m.f., $p(x)$ of X is given by
$$p(x_i) = P(X=x_i)$$

x	x_1	x_2	x_n
$p(x=n)$	$p(x_1)$	$p(x_2)$	$p(x_n)$

$p(x=x_i) = 0 \rightarrow$ for other values of n

Properties of p.m.f :-

- ✓ ① $p(x_i) \geq 0$ for $i = 1, 2, \dots, n$
- ✓ ②
$$\sum_{i=1}^n p(x_i) = 1$$

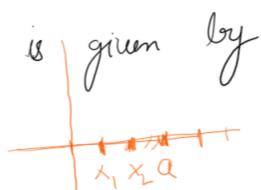
Graph of p.m.f. :

- (1) Negative skewed
- (2) Positive skewed
- (3) Symmetric
- (4) Constant

Cumulative distribution function (CDF)

CDF, F , of a random variable X is given by

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



Let $X \rightarrow$ discrete random variable taking values x_1, x_2, x_3, \dots
 $x_1 < x_2 < x_3 < \dots$

↑
 CDF, F of X is a step function.

Expectation & Variance of a Random Variable

Expectation

$$E(X) = \sum_{i=1}^{\infty} x_i P(X=x_i)$$

μ

Variance

$$\begin{aligned} \text{Var}(X) &= E[(X-\mu)^2] \\ \text{Var}(X) &= E(X^2) - [E(X)]^2 \end{aligned}$$

→ also called "long-run-average" value
 ↓
 in repeated independent observation

$$\text{① } \text{Var}(cx) = c^2 \text{Var}(X)$$

$$\text{Var}(x+c) = \text{Var}(X)$$

$$\text{Var}(ax+b) = a^2 \text{Var}(X),$$

a, b, c are constants

⇒ Expectation of $g(x)$,

$$E[g(x)] = \sum_{i=1}^{\infty} g(x_i) P(X=x_i)$$

$$\Rightarrow E[ax+b] = a E(x) + b \quad \text{where } a, b \text{ are constants}$$

$$\Rightarrow E[b] = b$$

→ X, Y are 2 random variables

$$E(X+Y) = E(X) + E(Y)$$

X = discrete random variable

$$g(x) = 3x+2 \\ = x(x-1)$$

$$\rightarrow \text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$$

X, Y are independent Random variables

$$\text{Var}(X+Y) = \sqrt{V(X+Y)} = \sqrt{V(X-Y)} = \sqrt{V(X)} + \sqrt{V(-Y)} \\ = \sqrt{V(X)} + \sqrt{V(Y)}$$

Discrete Random Variables

