

# Probability & Statistics

Estimate

Estimation — procedure

Estimator — Rule

Data

Primary

Secondary

|

|

First  
time

Existing  
sources

Des

Statistics

Descriptive

Infrential

|

|

Data into

Estimation

Graphical or

testing or

Tabular  
form.

Decision  
making

Nature of  
Variable

Qualitative  
(Attributes  
labels/names)

Quantitative  
/ /  
Discrete Continuous

Scales of Measurement:

- Nominal :

(Qualitative without  
Ranking) religion, gender,  
etc.

- Ordinal :

(Qualitative with  
Ranking) grades, Military

ranks : ~~ordinal~~

- Interval : ~~ordinal~~

(Quantitative and zero  
presence)

- Ratio : ~~ordinal~~

(Quantitative and zero  
absence).

True zero : jab kuch bhi  
na mile.

e.g.  $0^{\circ}\text{C}$  is not absolutely  
nothing.

# Probability & Statistics

Cross Sectional Data:

- Data collected for a specific time of individuals

Time Series Data:

- Data collected over a series of time.

$$R.F = \frac{\text{Freq}}{\text{Total Freq}}$$

Date

Page

Soft Drink

Freq (f)

$$\text{Angle} = R.F \times 360^\circ$$

Pepsi

5

$$\frac{5}{42} \times 360^\circ =$$

Coke

16

Sprite

5

Fanta

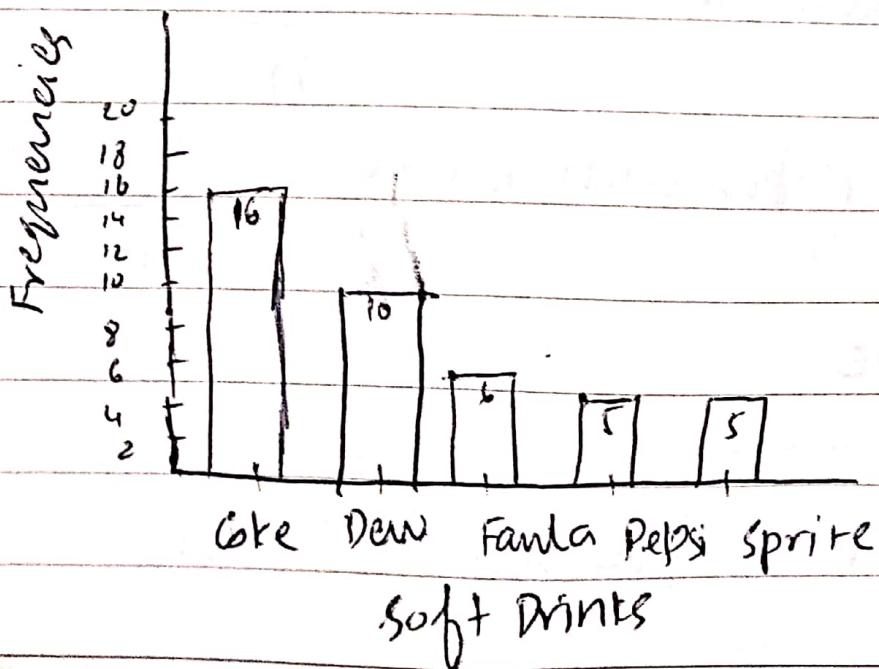
6

Dew

10

42

### Bar Chart



$\hat{\mu}$  = Population mode

$\hat{x}$  = Sample mode

## Measure of Central Tendencies:

### 1) Mean,

$n$  = sample size

$N$  = Population size

$\bar{x}$  = Sample mean =  $\frac{\sum n}{n}$

$\mu_x$  = Population mean =  $\frac{\sum n}{N}$

### 2) Mode

Most frequent value.

$\hat{\mu}$  = Pop. mode

$\hat{x}$  = Sample mode.

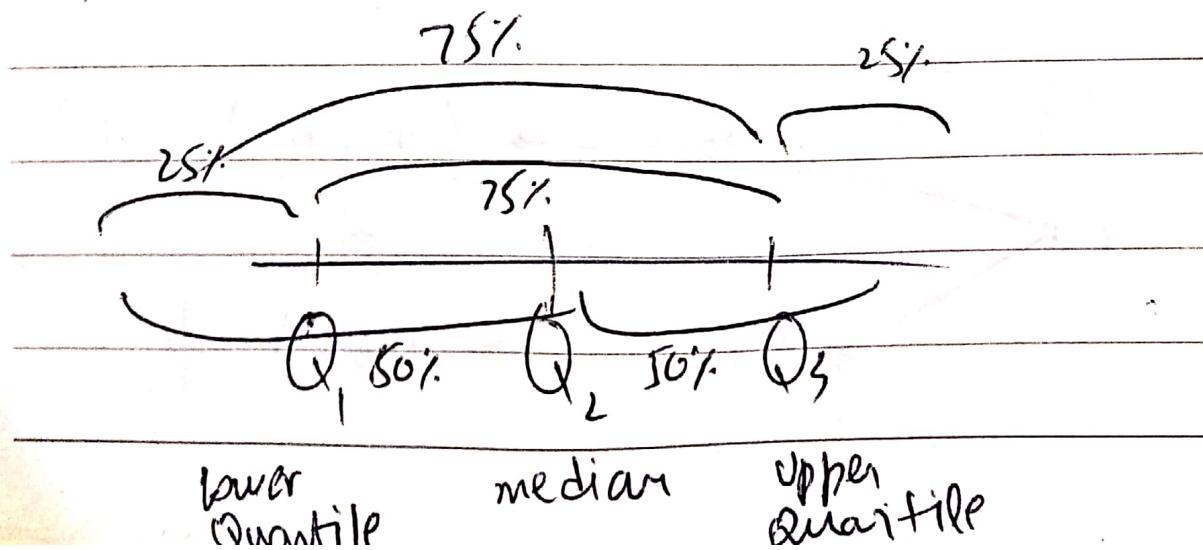
### 3) Median :

Middle value for  
sorted data set.

Even : 
$$\frac{\left(\frac{n}{2}\right)^{th} + \left(\frac{n+1}{2}\right)^{th}}{2}$$

Odd : 
$$\left(\frac{n+1}{2}\right)^{th}$$

### 4) Quartiles / Percentiles.



$P_{10}$   $P_{25}$   $P_{50}$   $P_{60}$

$D_i$  means Decile 10%.

$$i = \frac{P}{100} \times \text{total no. of obs.}$$

$$P_{25} = Q_1 = \frac{25}{100} \times n = 0.25 \times n$$

$$P_{50} = Q_2 = \text{median} = 0.5 \times n$$

$$P_{75} = Q_3 = 0.75 \times n$$

$$P_{80} = 0.80 \times n$$

$$\text{integer} \rightarrow \frac{i + (i+1)}{2}$$

$i$    
 not integer  $\rightarrow$  round-up

①      ②      ③      ④      ⑤  
 $3310, 3355, 3450, 3480, 3480 \dots$   
 $3490, 3520, 3540, 3550, 3650, 3730$   
 $3925$

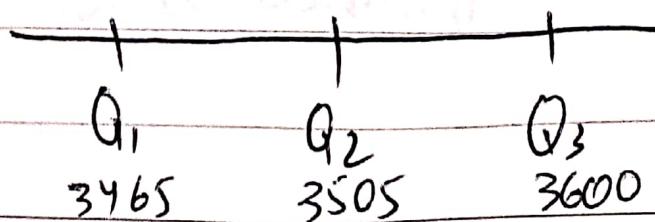
$$n = 12$$

$$P_{25} = 0.25 \times 12 = 3 \Rightarrow \frac{3^{\text{rd}} + 4^{\text{th}}}{2} \\ = 3465$$

$$P_{50} = 0.50 \times 12 = 6 \Rightarrow \frac{6^{\text{th}} + 7^{\text{th}}}{2} = 3505$$

$$P_{75} = 0.75 \times 12 = 9 \Rightarrow \frac{9^{\text{th}} + 10^{\text{th}}}{2} = 3600$$

$$P_{80} = 0.80 \times 12 = 9.6 \Rightarrow 10^{\text{th}} = 3650$$



Simple Mean = 3540

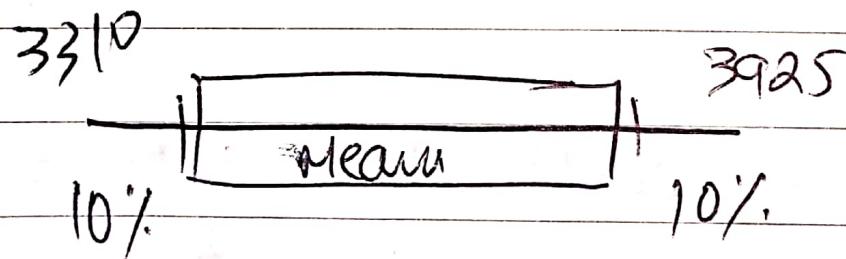
5) Trimmed Mean:

Total trim = 20%.

means 10% trim from  
one side (20% total).

trim means = 10%.

means 10% from one  
side.



trimmed by 10%

# Probability & Statistics

## Measures of Dispersion

When Average is same we  
can not analyze which data  
is better.

1) Range =  $R = X_m - X_n$

When there are outliers.

## 2) InterQuartile Range

$$\text{e.t. } IQR = Q_3 - Q_1$$

upper - lower  
quantile

## 3) Variance ~~or~~ Standard Deviation

Population:  $\sigma^2$

Sample:  $s^2$  or  $S^2$

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

$$s^2 = \frac{\sum (X_i - \bar{X})^2}{n-1} \quad (\text{Small population})$$

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n} \quad (\text{close to population})$$

Average of squared deviation of from the mean is variation.

Standard Deviation:

$$\sigma = \sqrt{\frac{\sum (x-\mu)^2}{N}}$$

$$s = \sqrt{s^2}$$

$$S = \sqrt{S^2}$$

:- For small population we divide by  $n-1$  because we have to eliminate biasness from the data set because we didn't take data from the other people

:- Small Sample ( $n < 30$ )

How S.D is related to mean?

For this:

Coefficient of Variation:

$$CV = \frac{S.d}{\text{Mean}} \times 100$$

lower is better

# Prob. & Stats

a) 68, 70, 71, 71, 71, 72, 73, 73, 73  
73, 74, 74, 74, 74, 75, 75, 75, 76  
77

$$\text{Median} = \frac{\left(\frac{n}{2}\right)^{\text{th}} + \left(\frac{n}{2} + 1\right)^{\text{th}}}{2}$$

$$= \frac{73 + 74}{2} = 73.5$$

b)  $\text{Max} = 77$

$\text{Min} = 68$

$\text{Median} = 73.5$

$$Q_1 = \frac{25}{100} \times 20 = 5$$

$$\Rightarrow \frac{5^{\text{th}} + 6^{\text{th}}}{2} = \frac{71 + 72}{2} = 71.5$$

$$Q_3 = \frac{75}{100} \times 20 = 15 \Rightarrow \frac{15^{\text{th}} + 16^{\text{th}}}{2} = \frac{74 + 75}{2} = 74.5$$

c) For outliers we calculate  
IQR

$$\text{IQR} = Q_3 - Q_1 = 74.5 - 71.5 = 3$$

Outliers are not there.

d)

~~AT&T~~

Range

$$\text{AT&T} : R = 75 - 66 = 9$$

$$\text{Sprint} : R = 69 - 63 = 6$$

$$\text{T-mobile} : R = 77 - 68 = 9$$

$$\text{Verizon} : R = 81 - 75 = 6$$

Sorting others

AT&amp;T

66, 66, 68, 68, 68, 68, 69, 69, 70,

71, 71, 72, 72, 72, 73, 73, 73, 74, 75,  
75

Sprint

63, 64, 69, 69, 65, 65, 65, 65, 66, 66,

66, 66, 66, 67, 67, 68, 68, 68, 69,  
69,

## Verizon

75, 76, 76, 77, 77, 77, 77, 77, 78, 78, 78,  
 79, 79, 79, 79, 79, 80, 80, 81, 81, 81

$Q_1$  ~~68~~

$$Q_1 = 0.25 \times 20 = 5 = 5^{\text{th}} + 6^{\text{th}} / 2$$

AT&T:  ~~$Q_1 = 0.25 \times 68 + 60 / 2 = 68$~~

Sprint:  $65 + 65 / 2 = 65$

T-mobile:  ~~$72 + 71 + 72 / 2 = 71.5$~~

Verizon:  $77 + 77 / 2 = 77$

$$Q_3 = 0.75 \times 20 = 15 \Rightarrow 15^{\text{th}} + 16^{\text{th}} / 2$$

AT&T:  $73 + 73 / 2 = 73$

Sprint:  $67 + 68 / 2 = 67.5$

T-mobile:  $74 + 75 / 2 = 74.5$

Verizon:  $79 + 80 / 2 = 79.5$

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

AT&T:  $73 - 68 = 5$

Sprint:  $67.5 - 65 = 2.5$

$$\text{T-mobile: } 74.5 - 71.5 = 3$$

$$\text{Verizon: } 79.5 - 77 = 2.5$$

$$\text{Mean } \mu = \frac{\sum X}{n}$$

$$\text{AT&T: } 1413/20 = 70.65$$

$$\text{Sprint: } 132.2/20 = 66.1$$

$$\text{T-mobile: } 1463/20 = 73.15$$

$$\text{Verizon: } 1567/20 = 78.35$$

~~Variance Standard Deviation~~

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

$$\text{AT&T: } 148.55/19 = 7.81$$

$$\text{Sprint: } 53.8/19 = 2.831$$

$$\text{T-mobile: } 85.66/19 = 4.508$$

$$\text{Verizon: } 58.55/19 = 3.081$$

Stand- Deviation

$$S = \sqrt{S^2}$$

$$\text{AT&T: } \sqrt{7.81} = 2.794$$

$$\text{T-mobile: } \sqrt{4.508} = 2.123$$

$$\text{Sprint: } \sqrt{2.831} = 1.682$$

$$\text{Verizon: } \sqrt{3.081} = 1.755$$

$$C.V = \frac{s.d}{\text{mean}} \times 100$$

$$\text{AT&T: } 2.794 / 70.65 \times 100 = 3.95\%$$

$$\text{Sprint: } 1.682 / 66.1 \times 100 = 2.54\%$$

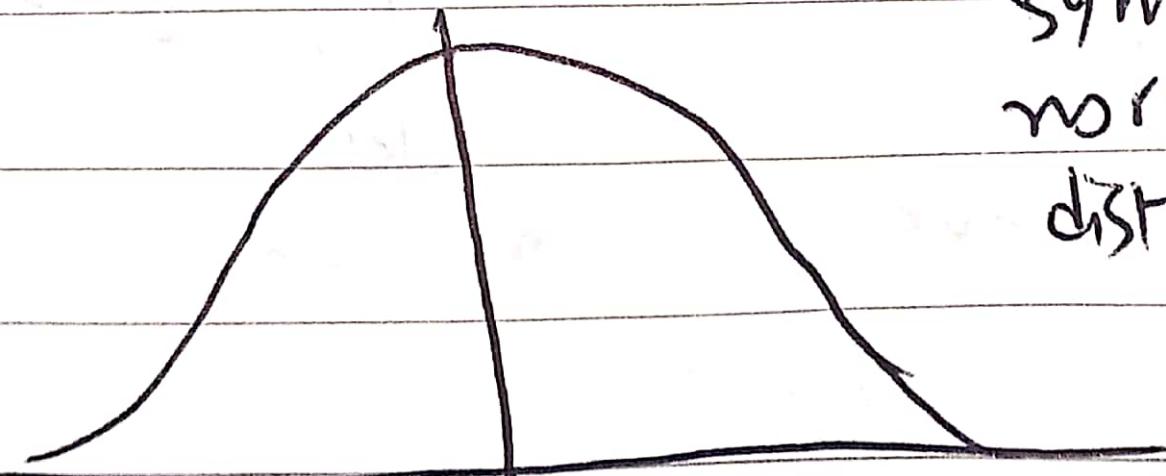
~~$$\text{T-mobile: } 2.123 / 73.15 \times 100 = 2.90\%$$~~

$$\text{Verizon: } 1.755 / 78.35 \times 100 = 2.239\%$$

# Prob & Stats

Ideal Scenario

Symmetric  
normal  
distribution



Mean = mode = median

- Everything apart from normality is skewed.
- Skewness can be positive or negative
- $\text{Mean} > \text{med} > \text{mode}$   
 $\Rightarrow$  +vely / right skewed
- $\text{Mean} < \text{med} < \text{mode}$   
 $\Rightarrow$  -vely / Left skewed
- To achieve normality we need transformation

How much the data is skewed?



Pearson's coeff. of skewness

$-3$ highly -ve	$0$ moderate	$+3$ highly +ve
--------------------	-----------------	--------------------

$$S_k = \frac{3(\text{Mean} - \text{mode})}{\text{s.d.}}$$

↓

amount of skewness

$$\text{T-mobile} = S_k = \frac{3(73.15 - 73.5)}{2.15}$$

$= -0.4884$  (slightly skewed to the left)

## Box Whisker Plot

Tmobile: Five point Summary

$X_0$	$Q_1$	med	$Q_3$	$X_m$
min	lower Quantile		upper Quantile	Max

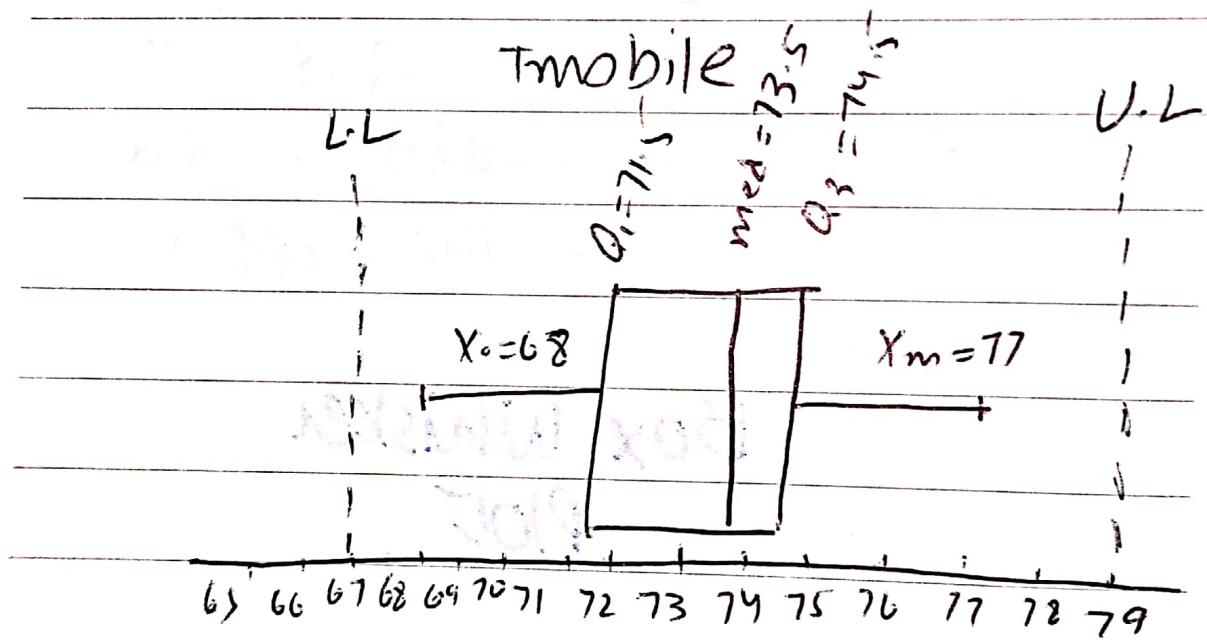
68 71.5 73.5 74.5 77

$$U.L = Q_3 + 1.5(IQR)$$

$$U.L = 74.5 + 1.5(3) = 79$$

$$L.L = Q_1 - 1.5(IQR)$$

$$L.L = 71.5 - 1.5(3) = 67$$



(CH3 41, 63, 20 David)

20 Dawson: 11 10 9 10

11 11 10 11

10 10

Clark : 8 10 13 7

10 11 10 7

15 12

Sorting data:

Dawson : 9, 10, 10, 10, 10, 10,

11, 11, 11, 11

Clark : 7, 7, 8, 10, 10, 10, 11, 12, 13

15

$$\text{Range (Dawson)} = 11 - 9 = 2$$

$$\text{Range (Clark)} = 15 - 7 = 8$$

# Five number Summary (Dawson)

$X_0$	$Q_1$	Med	$Q_3$	$X_m$
9	$0.25 \times 10$	$\frac{\left(\frac{n}{2}\right)^m + \left(\frac{n}{2}+1\right)^m}{2}$	$0.75 \times 10$	
	$= 2.5$		$= 7.5$	11
	$= 3^{\text{th}}$	$= \frac{5^{\text{th}} + 6^{\text{th}}}{2}$	$= 8^{\text{th}}$	
	$= 10$		$= 11$	
		$= \frac{10+10}{2}$		
		$= 10$		

$$IQR = Q_3 - Q_1 = 11 - 10 = 1$$

$$U.L = Q_3 + 1.5(IQR)$$

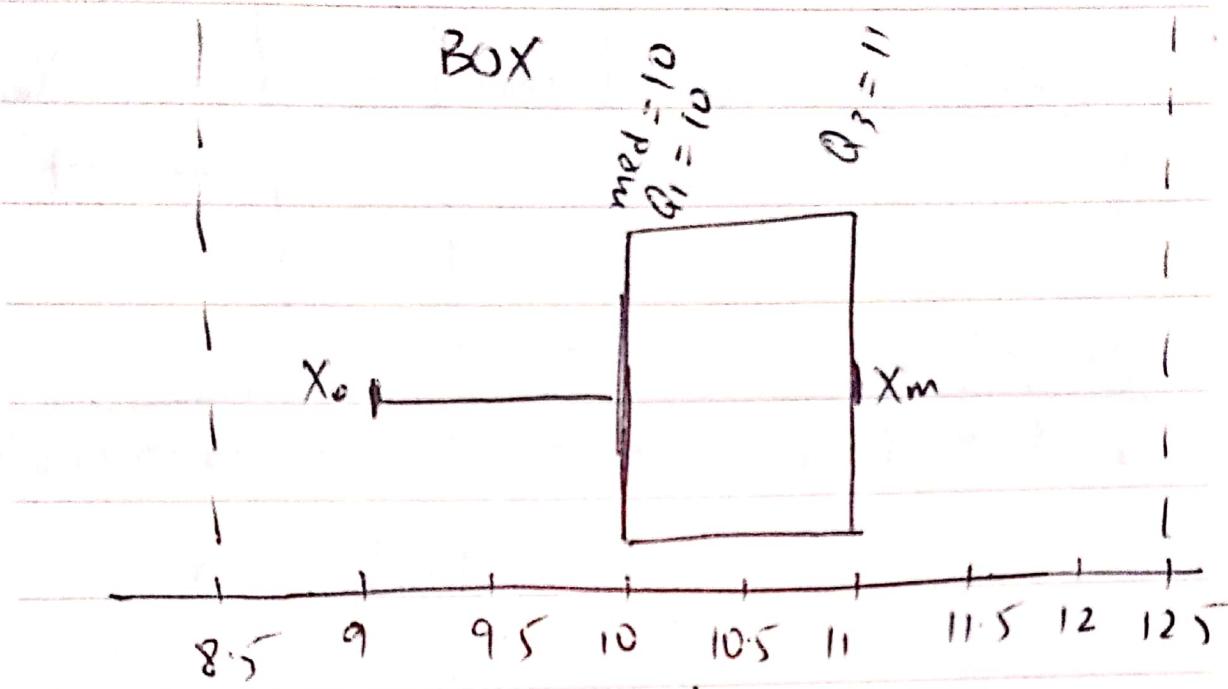
$$= 11 + 1.5(1) = 12.5$$

$$L.L = Q_1 - 1.5(IQR)$$

$$= 10 + 1.5(1) = 11.5$$

L.L

U.L



Rightly Skewed  
(dark)

$X_o$	$Q_1$	Med	$Q_3$	$X_m$
7	$0.25 \times 10$ $= 2.5$ $= 3^m$ $= 8$	$\frac{(\frac{n}{2})^m + (\frac{n}{2} + 1)^m}{2}$ $= \frac{5^m + 6^m}{2}$ $= \frac{10 + 10}{2}$ $= 10$	$0.75 \times 10$ $= 7.5$ $= 8^m$ $= 12$	15

$$IQR = Q_3 - Q_1 = 12 - 8 = 4$$

$$U.L = Q_3 + 1.5(IQR) = 12 + 1.5(4) = 18$$

$$L.L = Q_1 - 1.5(IQR) = 8 - 1.5(4) = 2$$

Date

Page

66

$Q_1, Q_3$

med

$Q_2 = 16$

U.L

$X_3$

$X_m$

2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

Normal Distribution

6/09/23 (cont.)

## Prob &amp; Stats

Year &amp; Audit time in days

12 14 19 18

15 15 18 17

20 27 22 23

n=20

22 21 33 28

14 18 16 13

① Classes	Tally	Freq(F)	RF	PF	CF	CRF	C.PF	CB
① 10-14		4	0.2	20	4	$\frac{4}{20} = 0.2$	20	9.5-14.5
② 15-19		8	0.4	40	12	$\frac{12}{40} = 0.3$	60	14.5-19.5
③ 20-24		5	0.25	25	17	0.85	85	19.5-24.5
④ 25-29		2	0.1	10	19	0.95	95	24.5-29.5
⑤ 30-35	1	1	0.05	5	20	1	100	29.5-34.5
$\Sigma$	-	$\sum F = 20$	1	100				

$$\text{No. of Classes} = 1 + 3.3 \log_{10} (\text{tot. obs})$$

Relative Freq = R.F	$X_i$	$X_i f_i$	$f_i(x_i - \bar{x})$
$= f_i / \sum f_i$	① 12	48	
Percent Freq = P.F	② 17	136	
$= R.F \times 100$	③ 22	110	
Commutative Freq = C.F	④ 27	54	
$\leftarrow P$	⑤ 32	32	
" R.F = C.R.F		$\sum$	$\sum$

$$" P.F = C.P.F$$

=

$$\text{Class Width} = \frac{\text{Max} - \text{Min}}{\text{No. of Classes}} = \frac{33 - 12}{5} = 4.2$$

No. of Classes 5

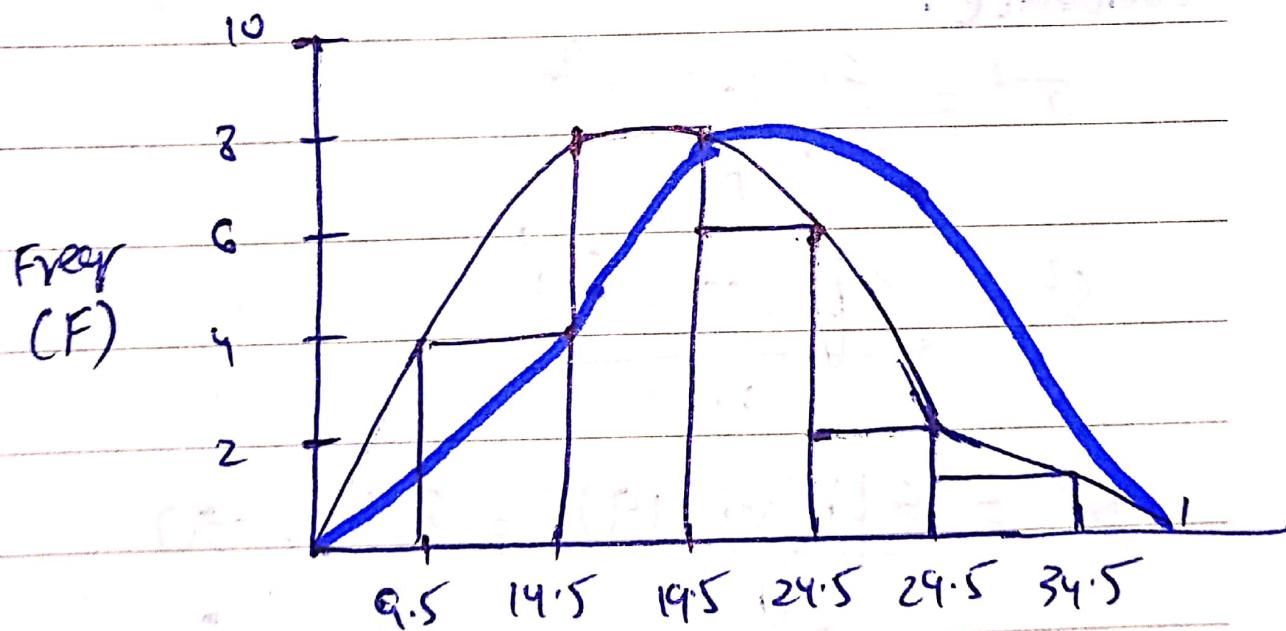
$X_i$  = Mid points

# Dot Plot

10 15 20 25 30 35 40 45

Freq  $\rightarrow$  Histogram

CF  $\rightarrow$



Class Boundaries

The data is very skewed

- Tend and audit  
- Normal curve

Mean:

$$\mu = \frac{\sum x_i f_i}{N} \text{ OR } \frac{\sum x_i f_i}{\sum f_i}$$

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

$$\sum f_i = n \text{ OR } \sum f_i = N$$

$$\bar{x} = \frac{\sum x_i f_i}{n} = \frac{380}{20} = 19 \text{ days}$$

Variance:

$$\sigma^2 = \frac{\sum f_i (x_i - \mu)^2}{N}$$

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

$$s^2 = \frac{\sum f_i (x_i - 19)^2}{20-1} = \frac{570}{19} = 30$$

$$\text{Standard Deviation} = \sigma = \sqrt{30}$$

$$= 5.48$$

# Prob & Stats

## Probability Basics

Experiment

Trial / Outcome

Random Experiment

Sample Space

Events

Simple / Compound

The <sup>Two</sup> Events

Mutually Exc / Non-M Exc

Equally likely Events

## Counting Principle

→ Rule of Multiplication

→ " " combination

→ " " permutation

# Assigning Probabilities:

- Classical Approach
- Subjective Approach
- R.F Approach

→ Repeating Experiments



Monte Carlo Runs.

Mutually Exclusive → At once

only one

category

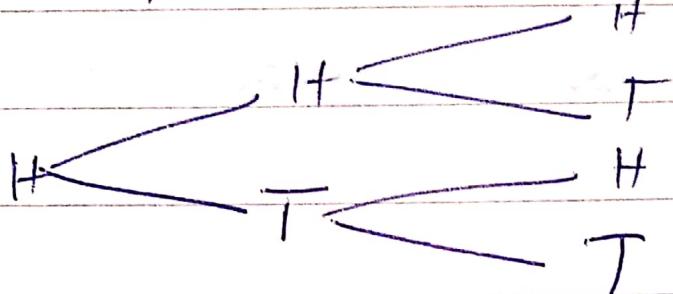
is present

and other is  
negated.

Equally likely → Probability of  
each event is  
same.

Sample space for a dice

$$\{1, 2, 3, 4, 5, 6\}$$



Rule of Multiplication:

one dice & one coin

$$\text{Sample space} = 6 \times 2 = 12$$

26 26 26 10 10 10 10

$$26^3 \times 10^4$$

order  $\rightarrow$  Permutation

Imp

order

not

Imp

$\rightarrow$  Combination

$$n = 6$$

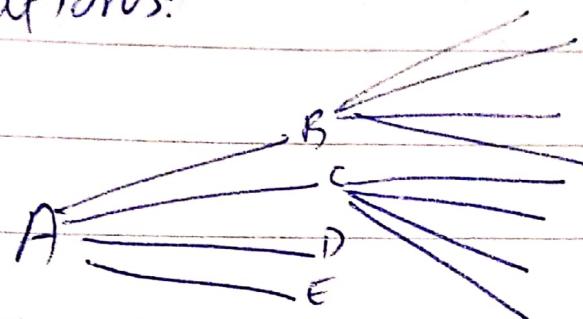
A B C D E F

$$r = 3$$

$${}^6P_3 = 120$$

$${}^6P_3 = 20$$

Permutations:



B

C

120

A B C D B C D

A B D B C E

A B E B C F

A B F B D E

A C D B D F

A C E B E F

A C F

A D E

A E F

C D E      D E F  
C D E  
C E F

$$n(S) = 6$$

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

$$E = \{2, 4, 6\}$$

$$n(E) = 3$$

$$P(E) = 3/6 = \frac{1}{2}$$

Success

$$1 - \frac{1}{2} = \frac{1}{2} \text{ (Failure)}$$

Classical is for  $\rightarrow$  Equally Likely Events

Chapter 2 Wheel pg 1

Henderson Pg 158 Q<sub>1</sub>  $\rightarrow$  Q<sub>6</sub>

Henderson Pg 169 Q<sub>22</sub>

# Prob & Stats

## Law of Complementation

$n=15$

$x=0, 1, 2, 3, \dots, 15$

Success + Failure = Total

0, 1, 2      3 ----- 15  
Failure      Success

$$P(S) + P(F) = 1$$

$$P(S) = 1 - (P(F))$$

$$= 1 - [P(x=0) + P(x=1) + P(x=2)]$$

$$P(F) = 1 - (P(S))$$

$$A^c = S - A$$

$$S = \{E_1, E_2, E_3, E_4, E_5\}$$

$$n(S) = 5$$

$$A = \{E_1, E_2\} \quad B = \{E_3, E_4\}$$

$$\Rightarrow C = \{E_2, E_3, E_5\}$$

$$n(A) = 2 \quad n(B) = 2 \quad n(C) = 3$$

$$P(A) = \frac{n(A)}{n(S)} = \frac{2}{5} \quad P(B) = \frac{2}{5}$$

$$P(C) = 3/5$$

A and B  $\Rightarrow$  Yes, A & B are

M.Ex.

- Law of Addition

$$P(A \text{ or } B)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

N.M.Ex

$$P(A \cup B) = P(A) + P(B) \quad \text{M.Ex}$$

$$P(A \cup B^c) = P(A) + P(B^c) - P(A \cap B^c)$$

$$P(A \cup B) = P(A) + P(B) = \frac{2}{5} + \frac{2}{5} = \frac{4}{5}$$

$$A \cup B = \{E_1, \bar{E}_2, \bar{E}_3, E_4\}$$

$$n(A \cup B) = 4$$

$$P(A \cup B) = \frac{4}{5}$$

$$A^c = S - A = \{E_3, E_4, \bar{E}_5\}$$

$$C^c = \{E_1, E_4\}$$

$$P(A^c) = 1 - P(A) = 1 - \frac{2}{5} = \frac{3}{5}$$

$$P(C^c) = 1 - P(C) = 1 - \frac{3}{5} = \frac{2}{5}$$

$$P(B \cup C) = P(B) + P(C) - P(B \cap C)$$

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

Q24  $\rightarrow$  ~~WT~~

F  $\sim$  fell short of exp

R  $\sim$  surpassed exp

M  $\sim$  Met exp

D  $\sim$  Didn't respond

$$P(F) = 0.26$$

$$P(M) = 0.65$$

$$P(D) = 4\% = 0.04$$

$$a) P(R) = 1 - [P(F) + P(M) + P(D)]$$

$$1 - 0.95 = 0.05$$

$$b) P(M \text{ or } R) = P(M) + P(R)$$

$$= 0.65 + 0.05$$

$$= 0.70$$

# Law of Multiplication

$$\textcircled{1} \quad P(A \cap B) = P(A) \cdot P(B|A)$$

$$\textcircled{2} \quad \text{or} \quad P(A \cap B) = P(A|B) \cdot P(B)$$

① A and B independent

② A and B dependent  
given

$$P(\text{Event 1 and 2}) \neq P(\text{Event 1})$$

$$\textcircled{1} \quad P(A \cap B) = P(A) \cdot P(B)$$

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

$$\text{Stats 1} \quad P(B)$$

		Y	N	
		n <sub>Y1</sub>	n <sub>N1</sub>	n <sub>(Y2)</sub>
Stats 2	Y	<input type="checkbox"/>	<input type="checkbox"/>	n <sub>(Y2)</sub>
	N	<input type="checkbox"/>	<input type="checkbox"/>	n <sub>(N2)</sub>
		n <sub>(Y1)</sub>	n <sub>(N1)</sub>	n <sub>(S)</sub>

Q-35

Pay Rent

		Yes ( $Y_R$ )	No ( $N_R$ )	
Buy	Yes ( $Y_B$ )	560.28	52	0.26
Car	No ( $N_B$ )	147	78	0.39
		70	130	200
		0.35	0.65	1

b)  $P(Y_B) = 0.54 > P(Y_R) = 0.35$

c)

$$P(Y_R | Y_B)$$

$$= \frac{P(Y_R \cap Y_B)}{P(Y_B)} = \frac{0.28}{0.54}$$

$$d) P(Y_R | N_B) = \frac{P(Y_R \cap N_B)}{P(N_B)}$$
$$= \frac{0.07}{0.46}$$

$0.52 \neq 0.35$  not ind.

$$e) P(Y_R | Y_B) = P(Y_R)$$

OR

$$P(Y_B | Y_R) = P(Y_B)$$

$$f) P(Y_B) + P(Y_R) - P(Y_B \cap Y_R)$$
$$= 0.54 + 0.35 - 0.28$$

Q38, 0

# Prob & Stats

Q-38 (Walpole (4.4))

$$M = 200$$

$$WNS = 40$$

$$W = 200$$

$$B = 280$$

$$MNS = 80$$

	B	S	
M	120	80	200
	(0.3)	(0.2)	(0.5)
W	160 (0.4)	40 (0.1)	200 (0.5)
	280	120	400
	(0.7)	(0.3)	(1)

$$a) P(B) = 200/400 = 0.5$$

$$b) P(S) = 120/400 = 0.3$$

$$c) P(M \cap S) = \frac{P(M \cap S)}{P(S)}$$

$$= \frac{80/400}{0.3}$$

$$= 0.66$$

$$P(W|S) = \frac{P(W \cap S)}{P(S)}$$

$$= \frac{40/400}{0.3}$$

$$= 0.33$$

$$d) P(M \cap S) = 80/400 = 0.2$$

$$P(W \cap S) = 40/400 = 0.1$$

$$e) P(M \cap S) = 80/400 = 0.2$$

$$f) P(W \cap S) = 40/400 = 0.1$$

g) For independent

$$P(M \cap S) = P(M)$$

By part (c)

$$0.66 \neq 0.5$$

∴ Dependant

$$P(W \cap S) = P(W)$$

By part (c)

$$0.33 \neq 0.5$$

∴ Dependant

# Prob & Stats

## BAYES' THEOREM

$$P(A \cap B) = P(A) (P(B|A))$$

$$P(A_i|B) = \frac{P(A_i)P(B|A_i)}{P(A_1)P(B|A_1) + \dots + P(A_n)P(B|A_n)}$$

→ where the email would go spam, inbox, there will be some condition

→ we can classify using Bayes theorem.

$$P(A_i|B) = \frac{P(A_i \cap B)}{\sum_{i=1}^n P(A_i \cap B)}$$

Date	Page	(1)	(2)	(3)	(4)	(5)
Events	Prior	Cond	Joint	Posterior		

Q40

Events	Prior	Cond	Joint	Posterior
$A_i$	$P(A_i)$	$P(A_i B)$		
$A_1$	.20	.50	0.1	0.286
$A_2$	.50	.40	0.2	0.512
$A_3$	.30	.30	<u>0.09</u>	0.23
			0.39	

Events	Prior	Cond	Joint	Posterior
$P_i$	$P(P_i)$	$P(D P_i)$		
$P_1$	0.3	0.01	0.003	0.157
$P_2$	0.2	0.03	0.006	0.315
$P_3$	0.5	0.02	<u>0.01</u>	<u>0.526</u>
			0.019	

$E_1$  = Cardholder will default

$$P(E_1) = 0.05$$

$E_2$  = missing monthly payment

$$P(E_2) = 0.20$$

$E_3$  = Cardholder not default

$$P(E_3) = 1 - 0.05 = 0.95$$

$M$  = missing a payment

~~E<sub>4</sub>~~

$$P(E_4 | M) = 0.2$$

$$P(E_1 | M) = 1$$

Events Prior Cond Joint Posterior

$$E_1 \quad 0.05 \quad 1 \quad 0.05 \quad 0.208$$

$$E_2 \quad 0.95 \quad 0.2 \quad \underline{0.19} \quad 0.791$$

$$0.24$$

$$A) \underline{0.791} + 0.208$$

# Prob & Stats

$W \sim$  msg contains word

$S \sim$  Spam

$H \sim$  Ham

let  $p$  be success as spam

q, failure of spam (for ham)

$$P(S|W) = \frac{P(W|S) P(S)}{P(W|S) P(S) + P(W|H) P(H)}$$

$$= \frac{P(\text{Word})}{P(\text{Word}) + q(\text{Word})}$$

$$= \frac{P(\text{Word1}) P(\text{Word2})}{P(\text{Word1}) P(\text{Word2}) + P(\text{Word1}) q(\text{Word2})}$$

Without prior knowledge

we assume equally likely

events = 0.5

$$\therefore P(H|D) = \frac{P(H \cap D)}{P(D)}$$

if it is reverse prob.

then  $P(D) \rightarrow$  total prob.

→ 1000 emails

700 Spams

300 Not Spams

$$H_1 = \text{Spams} = 0.7$$

$$H_2 = \text{Ham} = 0.3$$

$D = \text{email has word 'cash'}$

$$P(D|H_1) = 350/700 = 0.5$$

$$P(D|H_2) = 100/300 = 0.33$$

$$P(H_1|D) = P(D|H_1) \cdot P(H_1)$$

$$P(D|H_1)P(H_1) + P(D|H_2)P(H_2)$$

$$= \frac{(0.5) 0.7}{(0.5)(0.7) + (0.3)(0.3)}$$

$$= \frac{0.35}{0.35} < 0.78$$

$$0.449$$

$$P(D|H_1) = \frac{0.33 \times 0.3}{0.449} = \frac{0.1}{0.45} \\ = 0.23$$

## Example

word 'rolex' occurs  
in 250/2000 messages  
known to be spam

and 15/1000 messages  
to be ham which has  
the word rolex.

incoming message has rolex  
is spam assuming that it  
is equally likely that  
incoming is spam and not  
spam if our threshold for  
rejecting a message as spam  
is 0.9, will we reject?

$H_1 = \text{spam}$

$P(H_1) = 0.5$

$H_2 = \text{ham}$

$P(H_2) = 0.5$

$D = \text{rolex}$

$$P(D|H_1) = 250/2000 = 0.125$$

$$P(D|H_2) = 5/1000 = 0.005$$

$$P(H_1|D) = P(D|H_1) P(H_1)$$

$$\frac{P(D|H_1) P(H_1)}{P(D|H_1) P(H_1) + P(D|H_2) P(H_2)}$$

$$= \frac{(0.125)(0.5)}{(0.125)(0.5) + (0.005)(0.5)}$$

$$= \frac{0.0625}{0.065}$$

$$= 0.961$$

We will reject

$$H_1 = 2000$$

$$H_2 = 1000$$

$$P(H_1) = 2000/3000 = 0.66 \quad P(H_2) = 1000/3000$$

stalk under valued = 0.33

$$\downarrow \quad \downarrow$$

$$P(S|H_1) = \frac{400}{2000} \quad P(V|H_1) = \frac{200}{2000} = 0.1$$

$$= 0.2$$

$$P(S|H_2) = \frac{60}{1000} \quad P(V|H_2) = \frac{25}{1000} = 0.025$$

$$= 0.06$$

Both words is spam  
assuming we have no  
prior knowledge - 0.9

$$= \frac{P(S|H_1) P(V|H_1)}{P(S|H_1) P(V|H_2)}$$

$$= P(H_1|S) =$$

$$= \frac{P(S|H_1) P(H_1) P(V|H_1) P(H_1) + P(S|H_2) P(H_2) P(V|H_2)}{P(S|H_1) P(H_1) (P(V|H_1) P(H_1) + P(S|H_2) P(H_2) P(V|H_2))}$$

$$= (0.1)(0.66)(0.2)(0.66) \\ (0.2)(0.66)(0.1)(0.66) + (0.006)(0.33)(0.025)(0.33)$$

$$= 0.008712$$

$$\frac{0.0105}{0.0087}$$

$$= 0.829$$

$$= 0.98 > 0.9$$

We will reject