

## Lecture-1

# Statistics (Definition, why we need predictions) (R language was mainly designed for statistics)

### \* Descriptive Statistics

- Gives numerical and graphic procedures to summarize a collection of data in a clear and understandable way.
- Data properties and behaviour विवरण जैसे होते हैं।

=

Outlier: a data point that differs significantly from other observations.

### \* Inferential Statistics

- Provides procedures to draw inferences about a population from a sample.
- योग्य फैसला (Predictive decision).

Descriptive statistics वे होते हैं फैसला नहीं तो inferential statistics . Predictive decision .

Ex: Gmail में spam → inferential statistics वे होते हैं फैसला नहीं तो identify if it is spam mail or not. word pattern, या QR check करके या spam फैसला decide करके gmail. अतः यह important (spam → चल गया),

## Lecture-2

### # Descriptive Measures

\* Central Tendency measures: Mean  $\rightarrow$  Average.

They are computed to give a "center" around which the measurements in the data are distributed.

A : 100

B : 95

C : 7

D : 10

E : 10

F : 6

G : 9

$$\text{Average} = \frac{100 + 95 + 7 + 10 + 6 + 9}{7}$$

$$\approx 33$$

Real measure সত্ত্বাপন্ন অর্থে.

Average never handles extreme value.

Salman F. Rahaman - \$100000

Normal people - \$100

মাধ্যমিক আঁচ = Average.

Nature এবং প্রক্রিয়া center  $\Rightarrow$  প্রত্যেকেই,

\* Variation or Variability measures: standard deviation

They describe "data spread" or how far away the measurements are from the center.

एवढे डेटा कृत आऱ्याचे डेटाचे मरम्मी.

डेटा कृत असिल्याचे डेविएट राही.

स्टॅण्डर्ड डेविएट या गम तर आल्या, केंद्रीय तंद्रा

सो कोनक. डेटासाठी विशेष विशेष.

\* Relative Standing measures:

They describe the relative position of specific measurements in the data.

## # Measures of Central Tendency (or Location)

\* Mean (Extreme values handle अतिरिक्त विलम्ब)

- Arithmetic mean : अनुप्रयोगी दाता = weight नहीं

geometric mean

- Weighted mean :

A : 50 → 3 hrs

$$\text{arithmetic mean} = \frac{50 + 100 + 100}{3}$$

$$3 \times 50 = 150$$

$$2 \times 100 = 200$$

$$1 \times 100 = 100$$

$$\text{Weighted mean} = \frac{3 \times 50 + 2 \times 100 + 1 \times 100}{3}$$

B : 100 → 3 hrs class

$$\text{arithmetic mean} = \frac{100 + 50 + 50}{3}$$

50 → 2 hrs

$$\text{Weighted mean} = \frac{3 \times 100 + 2 \times 50 + 1 \times 50}{3}$$

- Geometric mean :

## \* Median (माध्यमिक Value)

data ते लोगो extreme value ना थाळते वा कम थाळते  
mean and median गाहाणा होते,

mean and median एवं रास्तें रास्त माने data खुल्ता balanced  
data सुल्ता ascending वा descending order द थाळते होते.

## \* Mode (सामान्य विषय वा वर्तने वाली विषय)

dataset द जावडेये तरिका महाप्रभु value कोने value आणे

रात्री असे ।

Median and mode can handle extreme values.

# Population :

The entire group about which information is desired.

# Sample :

sample is a subset of population.

A proportion or part of the population usually the proportion from which information is gathered.

$$n=3$$

$$p=7 \quad \underbrace{15, 25, 60}_{\text{sampling}} < 100$$

15

Randomly choose ~~any~~ data behaviour

25

विवरणों का चयन

60

Random sampling

179

125

1040

4049

## Lecture-3

### # Populations and Samples

A Good sampling is a sign where all the population can be represented.

#### \* Population Mean

$$\mu = \frac{\sum X}{N}$$

for ungrouped data, the population mean is the sum of all population values divided by the total number of population values.

-  $\mu$  is the population mean

- N is the total number of observations

- X is a particular value.

-  $\sum$  indicates the operation of adding .

- # A parameter is a measurable characteristic of a population.
- # A statistic is a measurable characteristic of a sample.
- # The Kiers family owns four cars. The following is the current mileage on each of the four cars:  
56000, 23000, 42000, 73000

$$\mu = \frac{\sum X}{N} = \frac{56000 + 23000 + 42000 + 73000}{4}$$
$$= 48500$$

## # Mean vs Average

## # Sample Mean

for ungrouped data, the sample mean is the sum of all the sample values divided by the number of sample values.

- n is the total number of values in the sample.

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

# A sample of five executives received the following bonus last year (\$k):

14.0, 15.0, 17.0, 18.0, 15.0

$$\bar{x} = \frac{\sum x}{n} = \frac{14.0 + 15.0 + 17.0 + 16.0 + 15.0}{5} = 16.4$$

Definition of Mean:  $\bar{x} = \frac{\sum x}{n}$

## # Arithmetic Mean Properties:

- \* Every set of interval-level and ratio-level data has a mean.
- \* All the values are included in computing the mean.
- \* A set of data has a unique mean.
- \* The mean is affected by unusually large or small data values.
- \* The arithmetic mean is the only measure of central tendency where the sum of the deviations of each value from the mean is zero.
- \*\*\* Affected by extreme values (outliers).

## # Consider the set: 3, 8, 4

The mean is 5

$$\sum (x - \bar{x}) = [(3-5) + (8-5) + (9-5)] = 0$$

## # Weighted Mean

The weighted mean of a set of numbers  $x_1, x_2, \dots, x_n$ , with corresponding weights  $w_1, w_2, \dots, w_n$  is computed, from the following formula:

$$\bar{X}_w = \frac{w_1x_1 + w_2x_2 + \dots + w_nx_n}{w_1 + w_2 + \dots + w_n}$$

Course	Grade	Points ( $x$ )	Credits ( $t$ )
Math	A	4	5
History	B	3	3
Health	A	4	2
Art	C	2	2

$$\bar{X}_w = \frac{5x_4 + 3x_3 + 2x_4 + 2x_2}{5+3+2+2}$$

$$= 3.42$$

## # Geometric Mean

The geometric mean (GM) of a set of  $n$  numbers is defined as the  $n$ th root of the product of the  $n$  numbers.

The geometric mean is used to average percents, indexes, and relatives.

$$GM = \sqrt[n]{x_1 \cdot x_2 \cdot x_3 \cdots \cdots \cdots \cdot x_n}$$

## # Properties of root

# The interest rates on three bonds were 5, 21, 4 percent.

$$\text{Geometric Mean} = \sqrt[3]{5 \times 21 \times 4} = 7.49$$

$$\text{Arithmetic Mean} = \frac{5+21+4}{3} = 10.0$$

The GM gives a more conservative profit figure because it is not heavily weighted by 21%.  
center tendency मात्रा नपेसत असै,  
Outlier से बच्याउन्छै।  
(If dataset लाई mean जस्ती value same कर्ति दिएजाउन)

## # The Median

Mean is highly sensitive to outliers or extreme values.

The median is the mid point of the values after they have been ordered from the smallest to the largest.

- There are as many values above the median as below it (in the data array).
- for an even set of values, the median will be the arithmetic average of the two middle numbers.

# 21, 25, 19, 20, 22      # 76, 73, 80, 75  
19, 20, 21, 22, 25      73, 75, 76, 80  
Median = 21      Median = 75.5

## # Properties of the Median

- \* There is a unique median for each data set.
- \* It is not affected by extremely large or small values and is therefore a valuable measure of central tendency when such values occur.
- \* It can be computed for ratio-level, interval-level, and ordinal-level data.
- \* It can be computed for an open-ended frequency distribution if the median does not lie in an open-ended class.
- \* If the dataset is uniformly distributed then  
 $\text{mean} = \text{median} = \text{mode}$ .  
Standard deviation = 0

## Lecture-4

Descriptive Statistics - data analysis वा decision  
करना

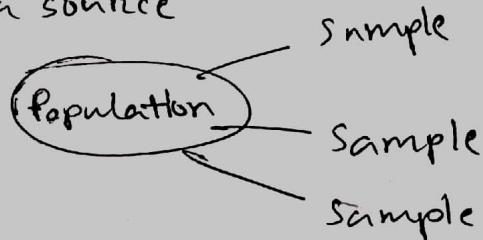
Inferential Statistics - data analysis वा mathematical  
model बनाकर future decision  
predict करना.

data analysis  
prediction       $\rightarrow$  statistics.

Statistic - value that describes samples

Parameter - value that describes a population.

Large pool of  
data source



statistical analysis एवं करना  
जी domain तरीके से जुड़ा  
इसका population

\* Sampling एवं इसे selectively and randomly.

sample mean - statistic

Population mean - Parameter.

Population mean,  $\bar{x} = \frac{\sum x_i}{n}$

Sample mean,  $\mu = \frac{\sum x_i}{n-1}$

Sample mean  $\rightarrow$  bias  $\rightarrow$   $\text{bias} = \frac{n}{n-1} - 1$  जबकि  $n-1$  नियम  
आप एवं ऐसी

# 280, 288, 300, 309

$$\bar{x} = \frac{280+288+300+309}{4}$$

$$\bar{x} = 294.25$$

$A = 280$   $d_i = \text{dispersion}/\text{केंद्रीय}$

$$d_i = x_i - A$$

$$d_i \rightarrow 0, 8, 20, 29$$

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$= \frac{\sum_{i=1}^n d_i + A}{n}$$

$$= \frac{\sum_{i=1}^n d_i + \sum_{i=1}^n A}{n}$$

$$= \frac{\sum d_i + nA}{n}$$

$$\therefore \bar{x} = A + \frac{\sum d_i}{n}$$

$$\bar{x} = 280 + \frac{0+8+20+29}{4}$$

$$= 294.25$$

$A = 67$  ( $x_i = 73$  middle)

Heights (in)	No. of Students ( $f_i$ )	Class Marks ( $x_i$ )	$\sum f_i x_i$	$d_i = x_i - A$	$d_i f_i$
60-62	5	61	305	-6	-30
63-65	18	64	1152	-3	-54
66-68	42	67	2814	0	0
69-71	27	70	1890	3	81
72-74	8	73	584	6	48
			<u>6745</u>		<u>45</u>

$$n = 100$$

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

$$= \frac{6745}{100}$$

$$= 67.45$$

$$\bar{x} = A + \frac{\sum d_i f_i}{n}$$

$$= 67 + \frac{45}{100}$$

$$= 67.45$$

## Lecture-15

Median for ungrouped data, Median is  $\frac{N+1}{2}$  th observation; N is odd.

Initial data : 20, 10, 23, 25, 100, 37, 105.

Sorted data : 10, 20, 23, 25, 37, 100, 105.

Median or  $Q_2$

Initial data : 20, 10, 23, 25, 100, 37, 105, 200

Sorted data : 10, 20, 23, 25, 37, 100, 105, 200

Median for ungrouped data, Median is average of  $\frac{N}{2}$ th and  $(\frac{N}{2}+1)$ th observation.

$$\frac{25+37}{2} = 31 \rightarrow \text{Median}$$

### Grouped data

weight	frequency
30-40	18
40-50	37
50-60	45
60-70	27
70-80	15
80-90	8

Left inclusive convention

↓  
30-39

90-49

Class Size = 10

weight	frequency	Cumulative Frequency
30-40	18	18
40-50	37	55 → $(\sum f)_e$
50-60	45 → $f_{median}$	100
60-70	27	127
70-80	15	142
80-90	8	150

$$N = 150$$

Median for grouped data:

$$\text{Median} = L + \left( \frac{\frac{N}{2} - (\sum f)_f}{f_{median}} \right) \times c$$

$L$  = Lower class boundary of median class

$N$  = Total number of items

$(\sum f)_L$  = sum of frequencies of all classes lower than the median class.

$f_{\text{median}}$  = frequency of median class.

$c$  = size of median interval.

$\frac{150}{2} \approx 75$ ,  $\frac{N}{2}$  median value  $\approx$  class 25,  $\approx$  median class.

Median class = 50-60, lower limit = 50.

$f_{\text{median}} = 45$

true lower limit  
lower class boundary  
amt 25

$(\sum f)_L = 55$

$c = 10$

$$L = \frac{50+49}{2} = 49.5$$

$L = 49.5$

$$\text{Median} = L + \left( \frac{\frac{N}{2} - (\sum f)_L}{f_{\text{median}}} \right) \times c$$

$$= 49.5 + \left( \frac{\frac{150}{2} - 55}{45} \right) \times 10$$

$$= 53.944$$

Linear interpolation method



## # Mode Calculation

5 8 10 8 5 8

एवं वैल्यु एवं फ्रेक्वेन्सी निचोड़े तरीके में mode.

5 → 2

8 → 3

10 → 1.

Mode = 8. ; Uni-modal data

5 8 10 8 5 8 5

5 → 3

8 → 3

10 → 1

Mode = 5, 8 ; Bi-modal data.

द्विघेद तरीके mode जैसा multi modal data.

## Lecture-6

1. Schaum series : statistic - Murray R. Spiegel

2. n n Probability - Seymour

Ques.

CT Topic : Inferential , population , statistic  
Descriptive sample parameter.

central tendency - arithmetic mean - grouped data.  
geometric mean

Median for grouped and ungrouped data

#

## Geometric Mean :

$$G = \sqrt[n]{x_1 \cdot x_2 \cdot x_3 \cdot \dots \cdot x_n}$$

$$= (x_1 \cdot x_2 \cdot x_3 \cdot \dots \cdot x_n)^{\frac{1}{n}}$$

$$\log G = \frac{1}{n} \log(x_1 \cdot x_2 \cdot x_3 \cdot \dots \cdot x_n)$$

$$\log G = \frac{1}{n} [\log x_1 + \log x_2 + \dots + \log x_n]$$

$$\log G = \frac{\sum f_i \log x_i}{n}$$

$$\therefore G = 10^{\frac{\sum f_i \log x_i}{n}}$$

#	Weight	Frequency
modal class →	30 - 40	18
	40 - 50	37
	50 - 60	45
	60 - 70	27
	70 - 80	15
	80 - 90	8

$$\text{Mode} = L + \frac{\Delta_1}{\Delta_1 + \Delta_2} \times c$$

L = Lower class boundary of modal class

c = class size of modal class

$\Delta_1$  = Difference between freq. of modal class and pre-modal class =  $f_m - f_1$

$\Delta_2$  = Difference between freq. of modal class and post-modal class =  $f_m - f_2$

$$L = 49.5 \quad \Delta_1, \Delta_2 \quad f_m = 45, f_1 = 37, f_2 = 27$$

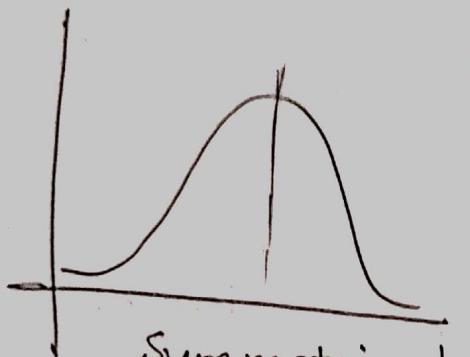
$$c = 10 \quad \Delta_1 = 45 - 37, \Delta_2 = 45 - 27 \\ = 8 \quad = 18$$

$$\therefore \text{Mode} = 49.5 + \frac{8}{8+18} \times 10 = 52.58$$

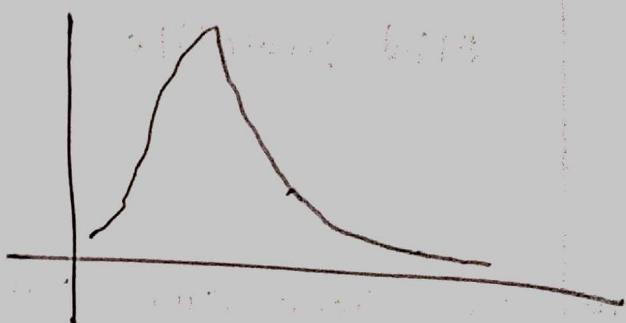
# Empirical Rule:

for a moderately skewed data

$$\text{Mode} = 3 \text{ Median} - 2 \text{ Mean}$$



Symmetric data

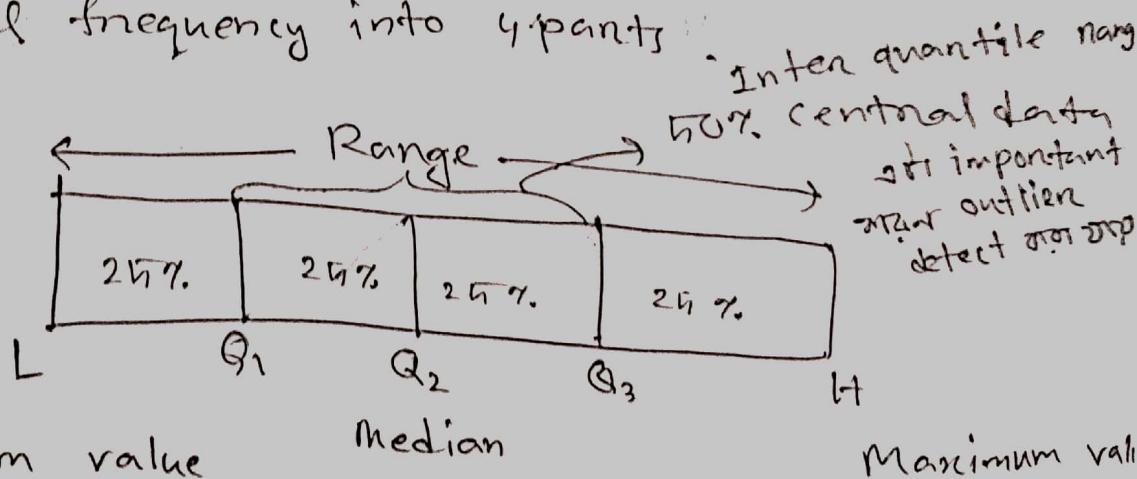


Moderately skewed data

not so symmetric  
not so asymmetric

#

Quartiles: Quartiles are values that divide the total frequency into 4 parts.



$$\text{Range} = H - L$$

$Q_1$  = At most 25% data are smaller than  $Q_1$  and at most 75% are larger.

$Q_3$  = At most 75% data are smaller than  $Q_3$   
and at most 25% are larger.

Interquartile Range (IQR) =  $Q_3 - Q_1$

Mid quartile  $\approx \frac{Q_1 + Q_3}{2} \approx Q_2$  (not always)

# If there are odd number of elements  
include median in two halves

