

Applied Statistics and Queuing Theory

↓
used in Machine Learning

Statistics is concerned with

- collecting → data collect করা-
- Organizing → unnecessary data
data Invalid data বর্তমান
- summarizing → single value রেখা
গ্রাফে distribution রেখা
- Presenting → represent করা,
Data present করা,
- analyzing → কি সিদ্ধান্ত পাইতে
পারবো, কি mean
করা।

purpose: - Make prediction

Inferential statistics → বর্তমান data র উপর ফিরি
বচতে future value predict
করা, (Machine Learning এ
use করা হবে)

Descriptive "

↓
data রে describe
করা

purpose: describing data

→ parameters

→ mean, median, std. deviation

* population এর subset sample.

* population অনেক বেশি বর্তমান
করা অবসরাম feasible হয়না।

Population

Sample → mean, median, std. deviation → statistics

↳ Randomly sampling popular method.

* Data distribution যদি biased হয় তখন analyzing
অসম্ভব নাও হতে পাবে।

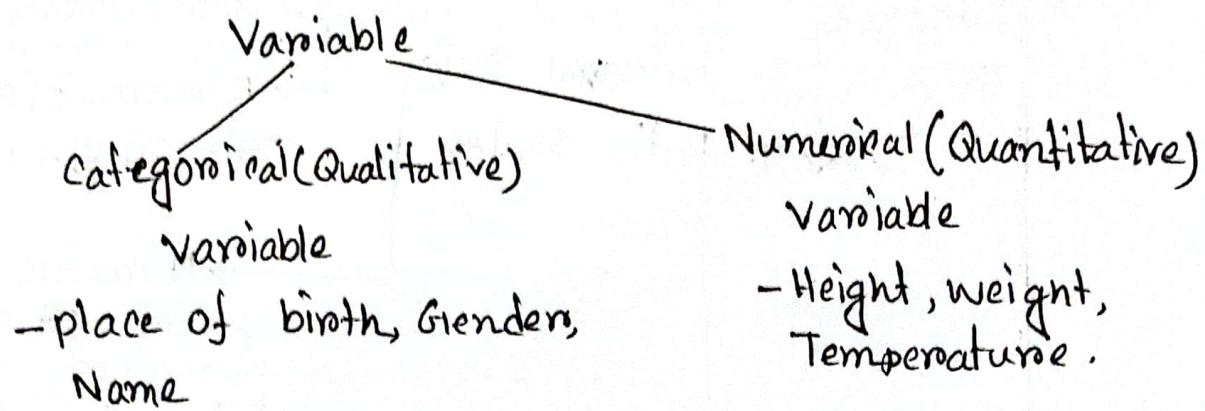
Parameden estimation :

- Sample এর statistics
state কোনো parameter state
- calculate করা, ideal
এবং আসতে উল্লেখ করা।

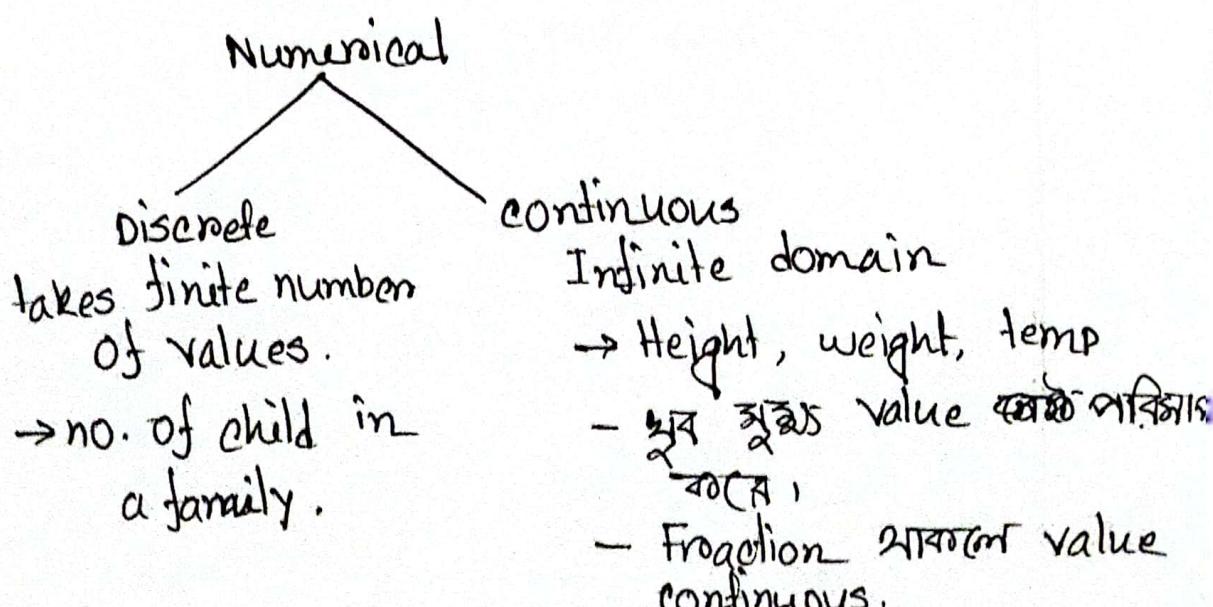
Books : statistics by Murray R. Spiegel
Probability by Seymour Lipschutz.

Variable: A variable is an attribute that describes a person, place, thing or idea.

variable changeable.



* Numerical এবং categorical এ কোথা সব বিষ্টি কোথা data loss হয়।

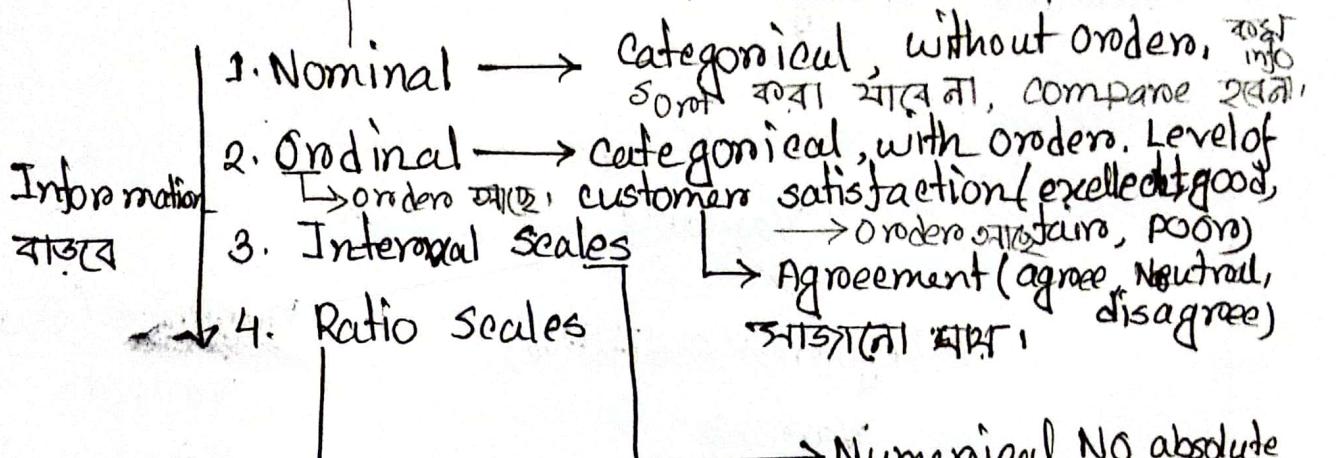


* Simplification এর ক্ষেত্রে continuous এবং discrete এ আলোচ্না,

* Level of Measurement :

* একটি variable থেকে কি পরিমাণ information

পাবো :



Numerical, absolute zero exist. Temp. in K.
একটি Reference value থাকে। যাতে compare
করা easy.

* Univariate data : looks at only variable.

একমাত্র data type নী।
কিম্বা পরিস্থিতির data নিয়ে বলতে হব।

সাধাৰণত data নিয়ে বা একটা
Feature নিয়ে বলজ কৰা হয়।

* Multivariate data : looks at multiple variable.

একাধিক data or feature নিয়ে
বলজ কৰা হয়। Classification কৰা
হাব।

Graphs

Pie chart:

* Categorical data visualization এর
জন্য use করা হয়।

Bar graphs:

* Decreasing আবশ্যক জাতান হলে
Pareto chart.

Multiple Bar Graph:

* একাধিক attribute এবং সময় মিলে chart
করতে হয়।

Line Graph:

* Progression show করার জন্য use করা হয়।

Frequency Distribution:

* যেনো particular value র frequency user কো হয়,
(ক্লাস ৫-১০ এর বিন্দু-অবিন্দু)

Left-end convention:

lower class interval
170 - 180 → upper class interval
180 - 190

a-b range & input হবে
 $a \leq x < b$

* last numbers হলো include কৰতে হবে।

* এই convention সামঞ্জ overlap এর জন্য, ambiguity র জন্য।

class size = upper class interval - lower class interval.

* 5-10 টি class interval জন্যা উচিত।

Ogive কৈথা

* Histogram and bar charts are not the same.

↓
slide দিয়ে পড়তে হবে।

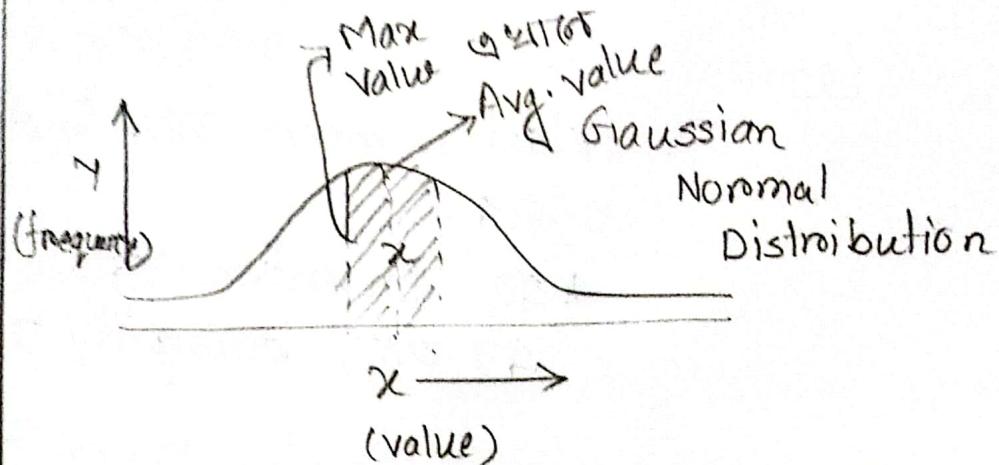
* Bin width related to class interval.

Scatterplot:

* Data distribution easily মাঝে, data র অর্ণে co-relation প্রকল্পে বুজা যায়।

Measure of Central Tendency:

* central এর মান কিরণ



* Avg এর সাথে পাঠে maximum value কানুনীভূত
থাকে।

মনসঃ 60 জন student এর marks distribution

* Arithmetic mean

* Geometricic " "

* Harmonic " "

* Quartiles " "

* Median " "

① Arithmetic Mean :

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

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

Tabular form of M.M.

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

where, $n = \sum f_i$

x	f
10	2
20	3
30	1

from ungroup frequency distribution table

* Sum of deviations from mean is equal to zero

$x_i \rightarrow x_1, x_2, \dots, x_n$

$d_i \rightarrow d_1, d_2, \dots, d_n$

$$\sum d_i = 0$$

$$d_i = x_i - \bar{x}$$

↓
বিদ্রুতি

↓
number Avg.

Formal proof :

$$\begin{aligned}
 & \sum_{i=1}^n (x_i - \bar{x}) \\
 &= \sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x} && \left| \begin{array}{l} \sum_{i=1}^n c = n c \\ \text{a constant} \end{array} \right. \\
 &= \sum_{i=1}^n x_i - n \bar{x} \longrightarrow \text{particular } \bar{x} \text{ of } n \text{ numbers} \\
 &= \sum_{i=1}^n x_i - \bar{x} \cdot \frac{\sum_{i=1}^n x_i}{\bar{x}} && \text{numbers are same avg } (\bar{x}) \\
 &= 0
 \end{aligned}$$

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

(A) → Assumed mean

$$d_i = x_i - A$$

Assumed mean actual mean
এর জন্ম না হলে কোনো d_i
0 হবে না।

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

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

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

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

$$\frac{3+4+5}{3} = 4 \rightarrow \text{True mean}$$

$A = 10 \rightarrow \text{Assume mean}$

$$\sum (3-10) + (4-10) + (5-10)$$

$$= -18$$

$$\bar{x} = 10 + \frac{-18}{3}$$

$$= 10 - 6 = 4 \quad \begin{array}{l} \text{true} \\ \text{mean} \end{array}$$

For tabular / group data:

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

Example:

Class	Frequency, f_i	x_i	$d_i = x_i - A$
0-10	7	5	-20
10-20	8	15	-10
20-30	20	25	0
30-40	10	35	10
40-50	5	45	20

$$n = \sum f = 30$$

$$\bar{x} = A + \frac{\sum f_i d_i}{n}, \quad d_i = x_i - A$$

class width, $C = 10$ * di এখন 10 রয়ে
সফল রয়ে reason
class width

→ shortest formula

$$\bar{x} = A + \frac{\sum f_i u_i}{n} \times C \rightarrow \text{class width}$$

↳ class distribution
even কর

$$u_i = \frac{di}{C}$$

$$u_i = \frac{x_i - A}{C}$$

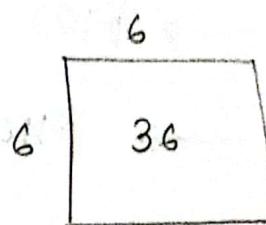
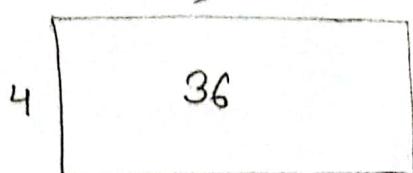
$$\bar{x} = 25 + \frac{(-14 - 8 + 0 + 10 + 10)}{50} \times 10$$

=

② Geometric Mean :

* Multiplicative in nature এবংজন data-র জন্য।

এজন একটি value নিবে যেন একার multiplication
করলে তাকে একটি সংখ্যার square root করা
যাবে।



x_1, x_2, \dots, x_n

$$G_1 = \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n}$$

$$= \left(\prod_{i=1}^n x_i \right)^{1/n}$$

$$\log G_1 = \frac{1}{n} \log (x_1 x_2 \dots x_n)$$

$$= \frac{1}{n} (\log x_1 + \log x_2 + \dots + \log x_n)$$

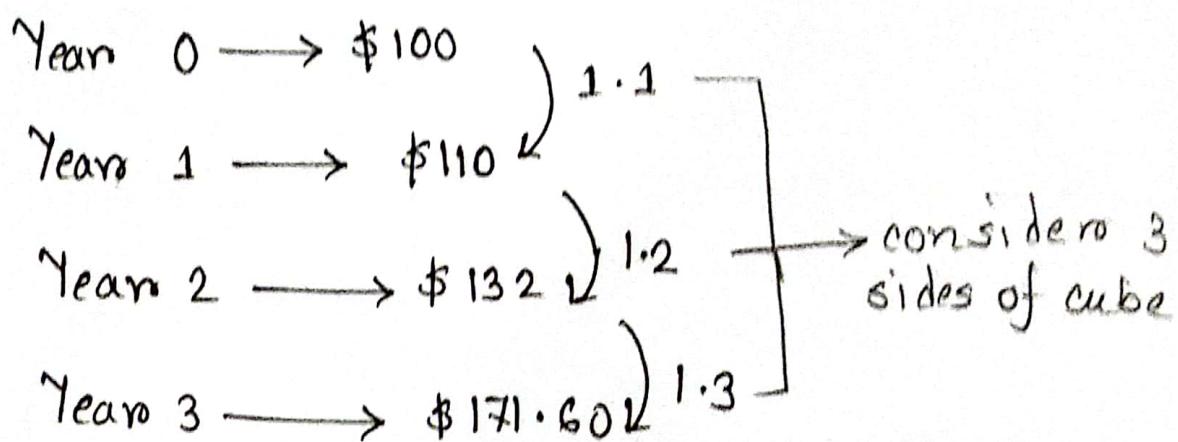
$$\log G_1 = \frac{\sum \log x_i}{n}$$

$$\log G_1 = \frac{\sum f_i \log x_i}{\sum f_i} \rightarrow \text{Group / Tabular data}$$

$$\therefore G_1 = \text{anti log} \left(\frac{\sum \log x_i}{n} \right)$$

18-03-2023

* Microsoft stock gained 10% in year 1, 20% in year 2 and 30% in year 3. What is the average yearly rate of return?



$$G = \sqrt[3]{1.1 \times 1.2 \times 1.3} = 1.1972 \rightarrow \text{Geometric mean of 3 sides of cube}$$

$$$100 \times 1.1972 = \$119.72$$

$$\$119.72 \times 1.1972 = \$143.33$$

$$\$143.33 \times 1.1972 = \$171.60$$

$$\$100 \times (1.1972)^3 \rightarrow 19.72\%$$

$$1.1972 = (1+0.1972) \rightarrow \$100 \rightarrow \$119.72$$

Arithmetic mean calculate ~~as~~ $\$100 \times (1.2)^3 = \172.8

which is quite accurate but geometric mean gives the actual rate.

Harmonic Mean :

$$\frac{1}{H} = \frac{1}{N} \sum \frac{1}{x}$$

$$H = \frac{N}{\sum \frac{1}{x}}$$

2, 3, 4 \rightarrow 3rd 31.401

$$\therefore H = \frac{3}{\frac{1}{2} + \frac{1}{3} + \frac{1}{4}} = 2.77$$

Harmonic series : $\frac{1}{1}, \frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \frac{1}{5}, \dots$

* কোনো particular frequency play করা হলে সেই
একটা play করা হলে নির্বাচিত overtone play করা
অথবা একটা value তার previous 3 next values
Harmonic mean. * Harmonic series যেকো
analysis করে Harmonic mean করা হয়।

For group / Tabular data :

$$H = \frac{\sum f}{\sum \frac{f}{x}}$$

* A car can travel 25 miles at 25 mph, 25 miles at 50 mph and 25 miles at 75 mph. What is the average velocity?

$$\text{Average velocity} = \frac{\text{Total distance}}{\text{total time}}$$

$$t = \frac{s}{v} \rightarrow = \frac{75}{\frac{25}{25} + \frac{25}{50} + \frac{25}{75}}$$

$$= \frac{3}{\frac{1}{25} + \frac{1}{50} + \frac{1}{75}} \quad \begin{array}{l} \text{Harmonic} \\ \text{Mean formula} \end{array}$$

$$= 40.9 \text{ mph}$$

Arithmetical mean : $\frac{25+50+75}{3} = 50 \text{ mph}$

অনেক মাত্রা
diff আসে।

* যখন উপরের fixed (এখানে travel distance) এবং
নিচৰা variable (sub unit এ variable), হলে
Harmonic mean use করা যাবে।

* swim 1 minute of freestyle at 3 km/h and
1 minute of butterfly at 2 km/h speed

এখানে time fixed but distance variable ($\frac{\text{distance}}{\text{time}}$)
তাহে arithmetic mean ফর্মুলা ব্যবহার করা যাবে।

$$\frac{3+2}{2} = 2.5 \text{ km/h}$$

* swim one lap of free style at 3 km/h
and one lap of butterfly at 2 km/h.
fixed distance

$$H = \frac{2}{\frac{1}{3} + \frac{1}{2}} = 2.4 \text{ km/h}$$

22. 03. 2023

Median :

Odd numbers : $\frac{N+1}{2} + h$

Even numbers : Avg. $\frac{N}{2} + h$ and $(\frac{N}{2} + 1) + h$ value

* at most 50% data median এর চেয়ে বেশি হবে।

2, 4, (5), 10, 12

Median এর চেয়ে বেশি Number আছে $= \frac{2}{5} \times 100$
 $= 40\%$

n n n ছোট n n = $\frac{2}{5} \times 100$
240%.

* Avg outline দ্বারা sensitive. but median না।

2, 4, 5, 10, 1000

Median = 5

Avg. = অনেক huge.

For even numbers

2, 4, (5, 10) 12, 25
↓
Avg যীগা হবে
Median = 7.5

Mode :

* frequency distribution এর mode বলা হয় mode value
মাত্র frequency বেশি।

2 2 2 3 3 15

x	f	
2	3	→ unimodal
3	2	mode = 2
15	1	

2 2 2 3 3 3 → uniform (সবক্ষেত্রে frequency same)

2, 2, 2, 3, 3, 3, 5, 7, 10 → Bimodal.

x	f	
2	3	
3	3	→ Bimodal.
5	1	
7	1	
10	1	

Median for group data:	weight	frequency	cumulative frequency
30 - 40	18	18	
40 - 50	37	55	P.c.f
median class → 50 - 60	45 → fm	100	
60 - 70	27	127	
70 - 80	15	142	
80 - 90	8	150	

$$N = \sum f = 150$$

$$\text{Median } \frac{N}{2} = 75$$

no. value of class

5 7.5

$$\text{Median} = L + \frac{\frac{N}{2} - \text{P.C.F}}{f_m} \times c$$

L = lower class boundary of median class

f_m = frequency of median class

p.c.f = preceding cumulative frequency
 ↙ / cumulative freq. of class before median

(Σf)₁

c = class width

Lower class boundary	lower boundary	Weight
49.5	30-39	30-40
40-49	40-49	40-50
50-59	50-59	50-60
Lower class interval		
49.5	30-39 40-49 50-59	30-40 40-50 50-60
49.5	50-59	50-60
	50 + 59	= 40.5
	2	প্রথম (50-60) এর true lower boundary
Median = 49.5 +	$\frac{150}{2} - \frac{55}{45}$	x 10
	= 53.94	

* Median class শাস্তি হতে হলে প্রথম class এর frequency
অনেক ধৈর্য হবে।

Mode for group data:

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

$\Leftrightarrow L$ = lower class boundary of modal class.

Δ_1 = Difference between modal class and pre modal class

$$= f_m - f_1$$

Δ_2 = difference between modal class and post modal class

$$= f_m - f_2$$

Weight	Frequency
30 - 40	18
40 - 50	37 $\rightarrow f_1$
-	
50 - 60	45 $\rightarrow f_m$
60 - 70	27 $\rightarrow f_2$
70 - 80	15
80 - 90	8

$$\begin{aligned}\Delta_1 &= f_m - f_1 \\ &= 45 - 37 = 8\end{aligned}$$

$$\begin{aligned}\Delta_2 &= f_m - f_2 \\ &= 45 - 27 \\ &= 18\end{aligned}$$

$$\text{Mode} = 50 + \frac{8}{8+18} \times 10$$

$$= 52.58$$

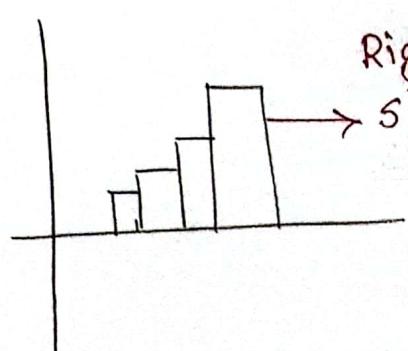
Empirical Relation of mode, Median and Mean
For moderately skewed distribution

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

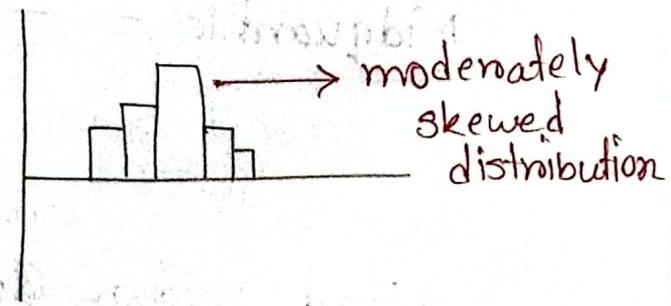
→ proof করা যাবে, only observation করে পাওয়া



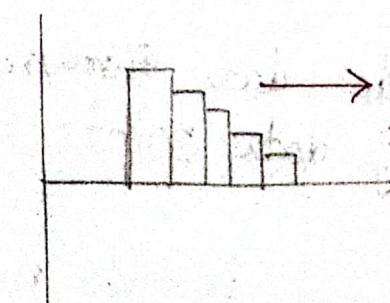
→ symmetrical distribution



Right (+ve)
skewed distribution (এর সাইডে
জড়িত)



moderately
skewed
distribution



Left (-ve)
skewed distribution

Quantiles : Quantiles are values that divide total frequency into 4 parts

25%	25%	25%	25%
Q_1	Q_2	Q_3	H

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

$Q_3 \rightarrow$ At most 75% data are smaller than Q_3 and at most 25% data are larger

$Q_2 \rightarrow 50 - 50$

$$IQR = Q_3 - Q_1$$

$$\text{Midquantile} = \frac{Q_1 + Q_3}{2}$$

↳ mid point দ্বারা এটা আবরণ Q_2 এ

Steps to find Quantiles:

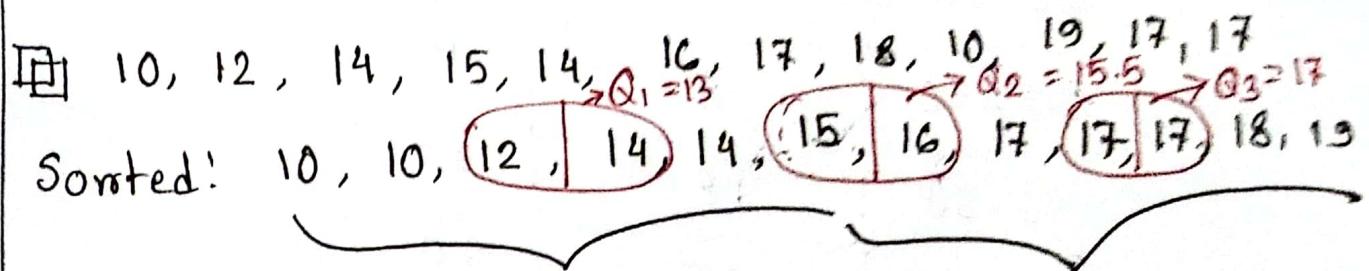
*Find median Q_2

*Use median to divide data further.

→ For odd no of data points include median in both halves.

→ For even, do not include median in either.

$2, 3, 5, 7, 8, 9, 10, 12, 15 \rightarrow$ sorted data



Grouped Data :

$$Q_i = L + \frac{\frac{ixN}{4} - P.C.f}{f} \times h$$

L = Lower class boundary of quartile class.

$$Q_2 = L + \frac{\frac{N}{2} - P.C.f}{f} \times h$$

* Quartile median හා generalization

f = frequency of quartile class

* quartile class is identified by $\frac{ixN}{4}$ th observation.

Profit	No. of companies c.f	
20 - 30	4	4
30 - 40	8	12
40 - 50	18	30
50 - 60	30	60
60 - 70	15	75
70 - 80	10	85
80 - 90	8	93
90 - 100	7	100

$$Q_2 N = 100$$

$$Q_2 = 49.5 + \frac{\frac{2 \times 100}{4} - 30}{30} \times 10$$

$$= 56.17$$

$$Q_1 = 39.5 + \frac{\frac{1 \times 100}{4} - 12}{18} \times 10$$

$$= 46.72$$

$$Q_3 = 59.5 + \frac{\frac{3 \times 100}{4} - 60}{15} \times 10$$

$$= 69.5$$

Percentiles :

$$P_i = L + \frac{\frac{ixN}{100} - P.c.f}{f} \times h$$

$\frac{ixN}{100}$ th observation

$$P_{75} = Q_3$$

$$P_{50} = Q_2$$

$$P_{25} = Q_1$$

Deciles :

$$D_i = L + \frac{\frac{ixN}{10} - P.c.f}{f} \times h$$

$$i = 1, 2, 3, \dots, 9$$

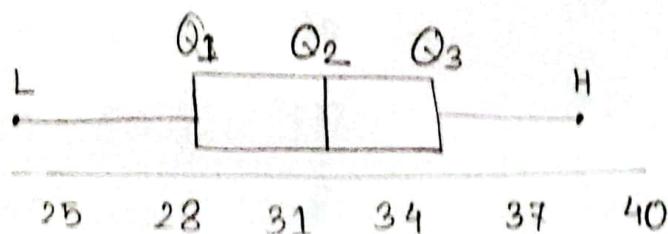
$$D_2 = P_{20}$$

Box Whiskers plot / 5 number Summary :

$L, H, Q_1, Q_2, Q_3 \rightarrow$ Quantiles

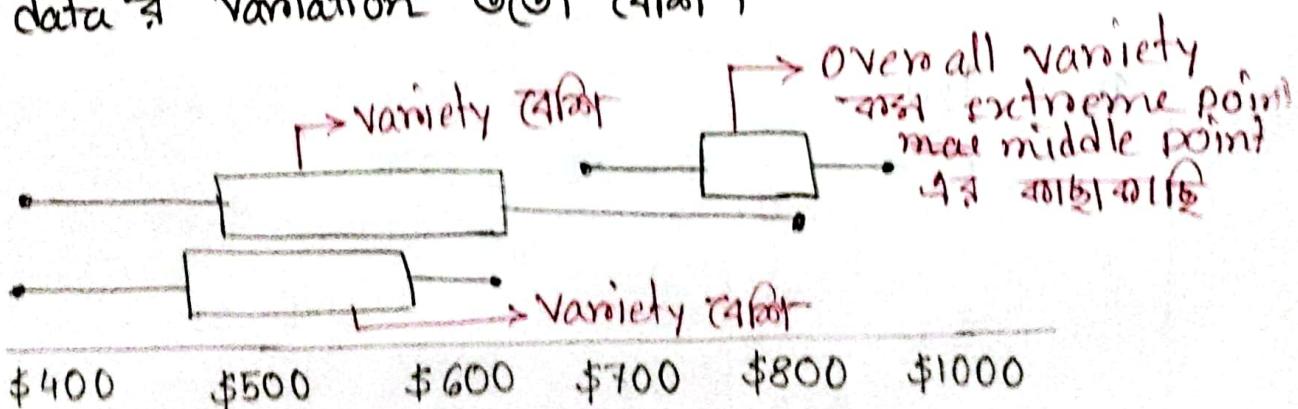
Data : 25, 28, 29, 29, 30, 34, 35, 35, 37, 38

$$L = 25, H = 38, Q_2 = 32, Q_1 = 29, Q_3 = 35$$



* Multiple box whiskers plot কে একে অপরের against
এ plot ~~করা~~ করে।

* Data টি central point দ্বারা যত দূর আছে থাকবে
data-র variation ততো হবে।



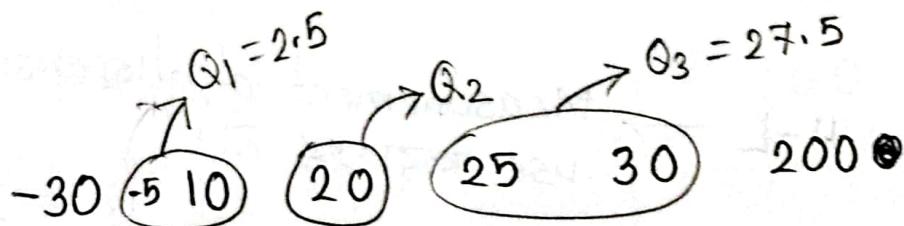
* box মতো বড় variation ততো হবে।

- * Outlier Detection using IQR
- * Outlier → Particular value distribution follows না বাবলে তা outlier's value.
Outlier highest value এবং ধোলি দের বক্রা যাবে না, comparison কীগৈ.
- * Outliers দের বাবতে Inter Quantile Range use কৰা হয়। $IQR = Q_3 - Q_1$

$$x > Q_3 + 1.5 \times IQR$$

$$x < Q_1 - 1.5 \times IQR$$

এই condition জিলে data কে outlier বলা হয়।



$$Q_1 = 2.5$$

$$Q_2 = 20$$

$$Q_3 = 27.5$$

$$IQR = 27.5 - 2.5 = 25$$

$$Q_3 + 1.5 \times IQR = 65$$

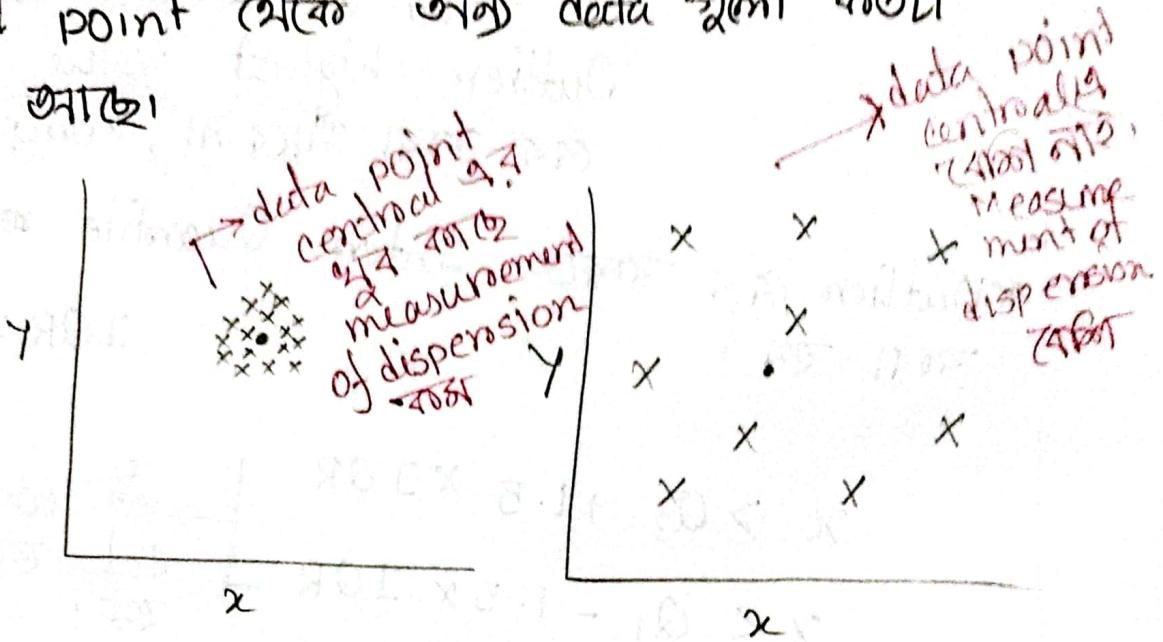
$$Q_1 - 1.5 \times IQR = -35$$

65 এর টুকু বড় ও -35 এর টুকু ছোট value outlier

$\therefore 200$ এখন outlier.

Measure of Dispersion

* Central point থেকে অন্য দোলা রেখা কতো ছড়িয়ে আছে।



* Range = H-L \rightarrow Measurement of dispersion ব্যবহার করা হয়।

* Measure of dispersion দুটি বিভিন্ন রকম। Range এর জুড়ে always একই হয়।

* Quantile Deviation, Q.D. = $\frac{Q_3 - Q_1}{2}$ IQR

A: 6, 46, 46, 46, 46

B: 6, 6, 6, 6, 46

C: 6, 10, 15, 35, 46

$$R = 46 - 6 = 40$$

A, B, C এর range
সame হলেও A & B
তে data variation কম।

* Quantile Deviation, Q.D = $\frac{Q_3 - Q_1}{2}$, $\frac{108}{2}$ km

Profit (In lakh)	No. of company	c.f
20 - 30	4	4
30 - 40	8	12
40 - 50	18	30
50 - 60	30	60
60 - 70	15	75
70 - 80	10	85
80 - 90	8	93
90 - 100	7	100

$$Q_i = L + \frac{\frac{i \times N}{4} - P.C.F}{f} \times h$$

$$Q_2 = 49.5 + \frac{\frac{2 \times 100}{4} - 30}{10} \times 10 = 56.17 \text{ lakh}$$

$$Q_1 = 46.72 \text{ lakh}$$

$$Q_3 = 69.50 \text{ lakh}$$

$$Q.D = 11.39 \text{ lakh}$$

Mean Deviation :

$$M.D = \frac{\sum_{i=1}^N |(x_i - \bar{x})|}{N}$$

For grouped data, $M.D = \frac{\sum f_i |x_i - \bar{x}|}{N}$

10 68 90 40

$$\bar{x} = \frac{208}{4} = 52$$

$$MD = \frac{|10-52| + |68-52| + |90-52| + |40-52|}{4}$$

$$= 27$$

Variance,

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{N} \rightarrow \text{MD এর বাইকালটি just square করা হয়।}$$

Standard Deviation,

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

[Grouped data র প্রক্রিয়া নিতে হবে]

- * standard deviation এ unit রাখে প্রয়োজন নির্দিষ্ট।
এখন তাই use করা হয়ে বিভিন্ন operation. Analysis
করার পরামর্শ দেওয়া হবে।
- * বিভিন্ন problem এর distribution এ calculation
এর অন্ত মেরা হবে variance.

Variance | এর shortcut method for table,

$$\begin{aligned}
 \sigma^2 &= \frac{\sum(x_i - \bar{x})^2}{n} \\
 &= \frac{\sum(x_i^2 - 2x_i\bar{x} + \bar{x}^2)}{n} \\
 &= \frac{\sum x_i^2}{n} - \frac{2\bar{x}\sum x_i}{n} + \frac{\sum \bar{x}^2}{n} \\
 &= \frac{\sum x_i^2}{n} - 2\bar{x}^2 + \frac{n\bar{x}^2}{n} \\
 &= \frac{\sum x_i^2}{n} - 2\bar{x}^2 + \bar{x}^2 \\
 &= \frac{\sum x_i^2}{n} - \bar{x}^2 \\
 \sigma^2 &= \frac{\sum x_i^2}{n} - \left(\frac{\sum x_i}{n}\right)^2
 \end{aligned}$$



→ Assumed mean
in instance

$$\sigma^2 = \frac{\sum f_d^2}{n} - \left(\frac{\sum f_d}{n} \right)^2$$

$d_i = x_i - A$
 $u_i = \frac{x_i - A}{c}$

$$= \left[\frac{\sum f_u^2}{n} - \left(\frac{\sum f_u}{n} \right)^2 \right] \times c^2$$

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

→ Bessel's correction

* ~~Best~~ Biased estimation: $(n-1)$ फूर्म अवश्यक बिसेड होता है।

Coefficient of variation:
 ↳ unitless.

Absolute measurement: একবা ধারণা (unit)

Relative " : Comparison type বর্তুল
 & Different elements এর

Group 1 (lbs) : 120, 130, 140, 150, 160
 $\mu_1 = 140 \quad \sigma_1 = 14.14$

Group 2 (kg) : 60, 65, 70, 75, 80 $\mu_2 = 70 \quad \sigma_2 = 7.07$

কোন group এ variation যেতি?

dispersion same হও, unit আলাদা বলে দেখে
 মনে হচ্ছে double.। এখানে dispersion ফর্মে measure
 করলে correct comparison আবশ্যে না as unit
 আলাদা। এজন্য আলাদা unit এর data র জন্য
 coefficient of variation use করা হব।

$$CV = \frac{\sigma}{\bar{x}} \times 100$$

$$CV_1 = \frac{14.14}{100} \times 100\% = 10.10\% \quad \rightarrow \text{group}$$

$$CV_2 = \frac{7.07}{70} \times 100\% = 10.10\% \quad \rightarrow \text{এর data র dispersion same.}$$

68-95-99.7 Rule, Empirical Rule for Normal distribution!

Z-score: standardization/standard score.

$$Z = \frac{x - \mu}{\sigma}$$

* একটি particular value mean এর টেক্স বাত স্ট্যান্ডার্ড deviation এভিয়ে বা পিছিয়ে আসা measure করবে, Z-score এটি relative measurement

Q Suppose mean of an exam is 80 and std. deviation 6.3. A student achieved 92.5. Student এর standard score কতো?

$$Z = \frac{92.5 - 80}{6.3} = \frac{12.5}{6.3} \approx 1.97$$

Avg এর টেক্স এ পরিমাণ std deviation এভিয়ে

$$\mu = 80$$

$$\sigma = 6$$

student A \rightarrow 87

$$Z_A = \frac{x - \mu}{\sigma}$$

$$= \frac{87 - 80}{6} = 1.17$$

$$= 1.50$$

$$\mu = 73$$

$$\sigma = 6$$

student B \rightarrow 82

$$Z_B = \frac{x - \mu}{\sigma}$$

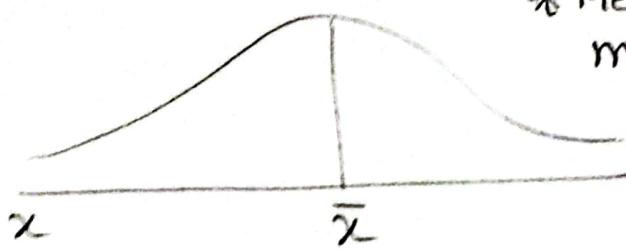
$$= \frac{82 - 73}{6} = 1.50$$

* comparatively student B result তুলি করে।

$$x_i \rightarrow 5, 7, 10, 12, 15, 16, 20, 25 \quad \bar{x} = 13.75 \\ \sigma = 6.24$$

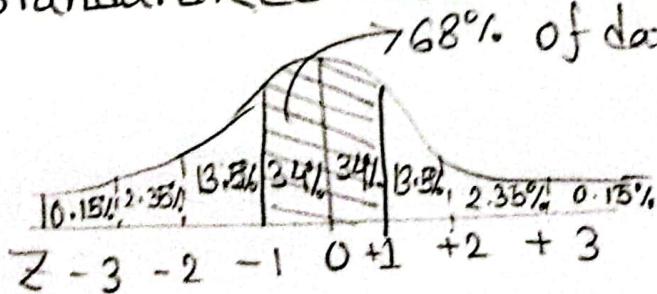
$$z_i \rightarrow -1.40, -1.08, -0.60, -0.28, 0.28, 0.36, \\ 1.00, 1.80$$

Normal distribution!



* Mean, median and mode অসমীয়া

যদি normal distribution এর z scale এ ৩ খণ্ড ২২
সময় standardized normal distribution
68% of data



* কিছি data যদি এখন ২২ এর normal distribution
এর pattern follow করে তাহলে mean মধ্যে
1 standard deviation data (-1 to +1) র ৬৮%
সহজে 68% data থাকবে।

* $-2 \leq z \leq +2$ range এ 95% data থাকবে,

* $-3 \leq z \leq +3$ range এ 99.7% data বা area পাইবে।

* $68 - 95 - 99.7$ Rule

* $z > +3$, $z < -3$ হলে সে data extreme outliers.

data standardized করার পর (Z-score এ convert করার পর)

mean $\rightarrow 0$

Variance / SD = 1

$$N(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

μ & σ \rightarrow parameters

$P(x)$ = Probability density function

standardized Normal distribution,

$$N(x; 0, 1) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

$\int_{-\infty}^0 N(x; 0, 1) dx = 0.5 \rightarrow \infty$ to 0 পর্যন্ত area under curve.

Q: Mean grade on final exam is 72 and standard deviation is 9. The top 10% student will get A+. What is the minimum mark required to get A+?

$$Z = \frac{x - \mu}{\sigma}$$

$$\therefore x = \mu + Z \sigma$$

$$= 72 + Z \times 9$$

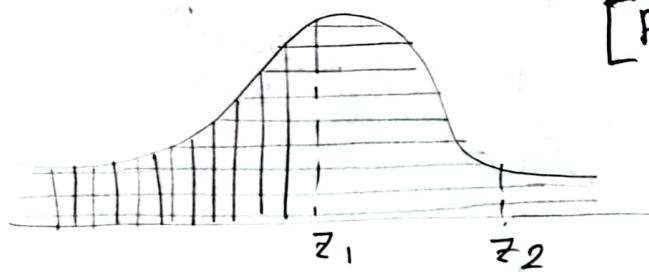
$$= 72 + (1.28 \times 9) \quad [Z = 1.28 \text{ from } Z\text{-score table}]$$

$$\text{now } \rightarrow 1.2 \quad \text{column } \rightarrow 0.08 \quad \left. \right\} \rightarrow Z = 90\% \text{ almost}$$

\downarrow
left sided Z-score table

H.W:

range of marks in central 50%



[Percentile টথকে
Z-score দ্বারা
calculation এ নাই]
বাকিটুলো আর

বাকি সব পর্যায়ে অন্তর্ভুক্ত
করা হয়েছে

এটা ZF এর মধ্যে থাকে

বাকি কোথায় থাকে? এটা কোথায় থাকে?

20.05.2023

Left sided Z-value \rightarrow percentile [Left, side ଓ
যতକ୍ଷେତ୍ର ଆଜି]

\downarrow
calculator ଏ ବରା ଯାଏ

$$1.0 \rightarrow 84\%$$

$$1.1$$

$$1.2$$

:

increment ହୁଏ ଏବଂ ବାରେ

Chebyshev's Theorem

The proportion of any distribution that lies within k standard deviation of the mean is at least

$$1 - \frac{1}{k^2}, k > 1$$

\downarrow lower bound

K	1.7	2	2.5	3
$1 - \frac{1}{k^2}$	0.65	0.75	0.84	0.89

Normal distribution ଏ 2 ଏବଂ 3 କୁ 95% ହଲେଓ,
ଏଥାଣେ ବଳେଇ at least 75% ଏବଂ ମଧ୍ୟେଇ ଥାବେ,
ଅଛାତା ମବୁ ସିନ୍ଦରିନ୍ଦର ଏବଂ ତଥା empirical
rule ନାହିଁ, ମେଣ୍ଡେତେ ଏହି theorem ବାବରେ ପ୍ରତି ଆଜି.

Moments:

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

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

$$\text{Central moment } \mu_r = \frac{\sum_{i=1}^N (x_i - \bar{x})^r}{N}$$

$$r=0, \mu_0 = 1$$

$r=1 \dots n \rightarrow$ more important

$r=1, \mu_1 = 0 \rightarrow$ sum of dispersion

$r=2, \mu_2 = \sigma^2 \rightarrow$ Variance

$r \geq 3, \rightarrow$ data এর symmetry / Assymetry
measure কৰে।

$r=4 \rightarrow$ data এর tail এর peak এর
ratio মাত্র measure কৰে।

$$\text{Raw moment, } \mu'_r = \frac{\sum_{i=1}^N (x_i - A)^r}{N}$$

$A = \text{Assumed mean}$

A reduces calculation

for grouped data

$$\mu'_r = \frac{\sum f_i (x_i - A)^r}{N}$$

$$\mu_r = \frac{\sum f_i (x_i - \bar{x})^r}{N}$$

$$\# \mu_1 = 0$$

$$\mu'_1 = \frac{\sum(x_i - A)}{N} = \frac{\sum x_i}{N} - \frac{\sum A}{N}$$

$$= \bar{x} - \frac{NA}{A}$$

$$= \bar{x} - A$$

$$\mu_2 = \frac{\sum(x_i - \bar{x})^2}{N}$$

$$= \frac{\sum\{(x_i - A) - (\bar{x}_i - A)\}^2}{N}$$

$$= \frac{\sum\{(x_i - A)^2 - 2(x_i - A)(\bar{x} - A) + (\bar{x} - A)^2\}}{N}$$

$$= \frac{\sum(x_i - A)^2}{N} - 2(\bar{x} - A) \frac{\sum x_i - A}{N} + \frac{\sum(\bar{x} - A)^2}{N}$$

$$= \mu'_2 - 2\mu'_1 \cdot \mu'_1 + \frac{\sum \mu'^2_1}{N} \rightarrow \frac{\sum \mu'^2_1}{N}$$

$$= \mu'_2 - 2\mu'^2_1 + \mu'^2_1$$

$$= \mu'_2 - \mu'^2_1$$

$$\therefore \boxed{\mu_2 = \mu'_2 - \mu'^2_1} \rightarrow \text{Relation bet^n raw and and central moment}$$

$$M_3 = M'_3 - 3M'_2 M'_1 + 2M'^3$$

$$M_4 = M'_4 - 4M'_3 M'_1 + 6M'_2 M'^2 - 3M'^4$$

$$M'_1 = \frac{\sum f_i (x_i - A)}{N}$$

$$= \frac{\sum (f_i u_i - h)}{N}$$

$$u_i = \frac{x_i - A}{h}$$

$$\Rightarrow x_i - A = u_i h$$

$$= \frac{\sum f_i u_i}{N} \times h$$

$$\therefore M'_2 = \frac{\sum f_i u_i^2}{N} \times h^2$$

$$M'_3 = \frac{\sum f_i u_i^3}{N} \times h^3$$

$$\therefore M'_n = \boxed{\frac{\sum f_i u_i^n}{N} \times h^n}$$

Skewness, moments and kurtosis:

$$\frac{\sum f_i u_i}{N} \times h, \quad \frac{\sum f_i u_i^2}{N} \times h; \quad u_i = \frac{x_i - A}{h}$$

Assumed mean = raw moment

about the mean = central moment

आलोचित असमेय विश्लेषण आलोचित नाव मोमेंट बत आलोचित केंद्रीय मोमेंट सामै भए।

1st central moment = 0 \rightarrow Always.

Standardized moment:

* Transition Invariant:

* value য় মাত্র constant

from আজ কোন

value shift হবে but

dispersion (ভৱিধি মান)

হবেনা, standard deviation

change হবেনা।

* স্ট্যান্ডার্ড বিভক্তি

* standard deviation ফর্ম অনেক সময় data কে আসা
করা যায়।

* distribution এর ক্ষেত্রে e shape এর কিছু property

Follow করে ~~standard~~ moment, skewness

* standard deviation ফর্ম দোলা data আস
ফর্ম unit less হবে, which will make the
comparison easy and Dimensionless.

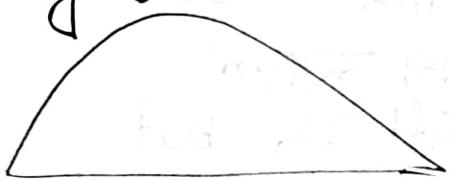
standardized moment:

$$\frac{\mu}{\sigma}, \frac{\mu^2}{\sigma^2}, \frac{\mu}{\sigma^3} \rightarrow \text{std deviation ফর্ম}-
এর করা হলো-$$

Skewness:

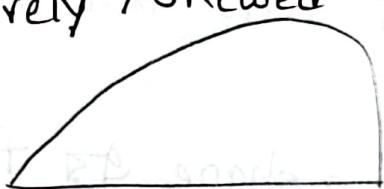
অসূত্র ফিকে দাটা, distribution asymmetric মানবে।

*Positively / Right Distribution:



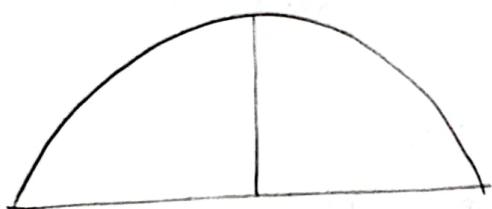
*Mean অবচেত্য বড় হয় median and mode

*Negatively / Skewed Distribution:



*Mean অবচেত্য ছোট mode বেশি।

*Symmetrical Distribution



*mean, mode, median same.

Coefficient of Skewness:

$$Skp = \frac{\text{Mean} - \text{Mode}}{\sigma}$$

→ Positive distribution (+)
→ Negative " (-)

↳ 1st skewness coefficient

* Skewness Factors same মানতে হবে, size দ্বিগুণ
বড় করলেও, shape same মানবে just scale up
বা down হবে।

$$Skp = \frac{3(\text{mean} - \text{median})}{\sigma}$$

↳ 2nd coefficient
of skewness

* Quantile কিভাবে

* Quantile and Percentile ফরি calculate করা

মান

skewness range $-\infty$ to ∞ , Ideally value (-3 to 3) মানবে
because of std. deviation

* Purely symmetric value হলে $Skp = 0$

3rd central moment $\rightarrow (+)$ হলে (+)ve skewness
 $\rightarrow (-)$ " " (-) " "

$$\frac{\mu_3}{\sigma^3}$$

Karl Pearson skewness :

$$\beta_1 = \frac{\mu_3}{\mu_2^3} \quad \mu_2 \rightarrow \text{Variance}$$

কার্ল পেরসনের পোসিটিভ না

β_1 always (+)ve হবে, রেফেন্স তুম্ব মাপের মধ্যে,
direction refers করতে পারে না।

$$\gamma_1 = \frac{\mu_3}{\sigma^3} = \frac{\mu_3}{(\sigma^2)^{3/2}} = \frac{\mu_3}{\mu^{3/2}} = \sqrt{\beta_1}$$

↓
Karl Pearson
moment
coefficient of
skewness

↙
2nd central moment
 $\sigma^2 = \text{Variance}$

$\gamma_1 \rightarrow (+)$ হলে (+)ve skewed
 $\gamma_1 \rightarrow (-)$ " (-)ve "

Kurtosis \rightarrow 4th standard central distribution.

$$\gamma_1 = \frac{\sqrt{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^{3/2}} = \frac{u_3}{\sigma^3}$$

* raw moment নির্ণয় করা easy

Kurtosis:

* Skewness distribution এর ছুটা জিনিয় নির্ণয় concern

① Peak.

② Tail

Peak and Tail এর উপর kurtosis depend করে।

kurtosis distribution

① leptokurtic:

Shapener than normal distribution, Fat tail than normal distribution

② Mesokurtic:

Normal distribution এর মতো peak & tail

③ Platykurtic :

Tail thinner than normal distribution,
Peak flatter than normal distribution

* moment এর square powers নির্ণয়

* μ_1 μ_2 μ_3 μ_4 μ_5 μ_6 Skewness বুজাই

* μ_1 μ_2 μ_3 μ_4 μ_5 μ_6 Kurtosis (Peak & tail) বুজাই

$$\beta_2 = \frac{\mu_4}{\sigma^4} = \frac{\mu_4}{\mu_2^2}$$

$$= \frac{n \sum_{i=1}^n (x_i - \bar{x})^4}{\left[\sum (x_i - \bar{x})^2 \right]^2}$$

Ø

mesokurtic $\rightarrow \beta_2 = 3$

leptokurtic $\rightarrow \beta_2 > 3$

platykurtic $\rightarrow \beta_2 < 3$

* moment generating function \rightarrow polynomial এর powers জিলো

Excess Kurtosis, $\gamma_2 = \beta_2 - 3$

\hookrightarrow platy (-) হবে

\rightarrow kurtosis দ্বারা

বর্ণনা করে এটা

lepto (+)

হবে
(using 4th and 2nd
central moment)

Correlation

* একটা বুদ্ধি করলে আবেক্ষণ্য কি change করে।

Positive Negative:

simple & multiple.

① simple \rightarrow দুই variable এর মধ্যে

② Multiple \rightarrow কয়েকটা মধ্যে

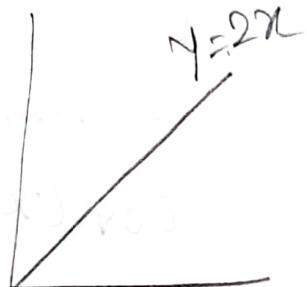
Linear: Linear অর্থাৎ আছে কিনা, একটা change
করলে আবেক্ষণ্য এই scale এই change রেখে
কিনা

Method of estimation Correlation:

* perfect positive correlation:

perfectly একটা line এ fit করবে

But negative



* perfect negative correlation:

perfectly একটা line এ fit

করবে তবে টাল (-) হবে

* High degree of Positive

* High degree of Negative

* No correlation!

যদো অন্তর নাই এবং বাড়লে আবেগ ঠিক
বজে না বাড়বে ঠিক নাই,

* Correlation হলো cause and effect হবে
এখন না।

* আবার causation হলো correlation হবে।

Covariance:

$$\text{Cov}(X, Y) = \sigma_{XY} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{N}$$

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

* দুটি variable এর জন্য

calculate করা হয়।

Variance

* Covariance দুটি ~~বিনামূল~~ variable এবং
correlation নির্ণয় করে।

