

Unit - II

DATE: 13/10/23
PAGE:

Data Exploration:

It is the first in data analysis involving the use of data visualization tools and statistical techniques to uncover data set characteristics and initial patterns.

Why is data exploration important?

- Humans are visual learners, able to process visual data much more easily than numerical data. usually
- We can represent data with the help of data exploration.

Data profiling is a process of reviewing and cleaning data to better understand how its structured and maintain data quality standards beneath an organisation.

Fundamental of statistics for data Analytics.
Statistics is a branch of mathematics, it deals with the analysis, interpretation, presentation of organisation of data.

Types of statistics: Statistics is divided into two category -

- ① Descriptive statistics - provide information on
- ② Infer^{the} distribution, central tendency and variability of data.
It includes mean, median, mode, variance, standard deviation, range.

② Inferential statistics: It is used to make prediction or draw conclusions about a population based on ^{sample} inferential data. It includes hypothesis testing, parameter estimation and relationship b/w two variables.

Benefit of statistics for data analytics

- ① Summarise and describe data.
- ② Identify pattern trend & relationship b/w data.
- ③ Make future prediction and find conclusion.
- ④ Hypothesis test.

Fundamental terms

- i) Probability -
- ii) Population & sample -
entire group data is, it (2 out of 10)
- iii) Distribution of data
has data is spread out

Normal Uniform Skewed

(iv) Measure of central tendency
It describes the central value of data set by using mean, median and mode.

v) Variability It refers how data is spread out. It find variability by using 1st deviation

vi) Central limit theorem,
conditional probability (event already occurred)
vii) Co-variance & co-relation
viii) hypothesis testing
ix) Student test, T-test and F
x) Type 1 error & Type 2 error

① find mean, median and mode

Number	7	8	9	10	11	12	13
frequency	3	8	12	15	14	12	9

$$1 \times 7 + 3 \times 8 + 12 \times 9 + 15 \times 10 + 14 \times 11 + 12 \times 12 + 9 \times 13$$

$$\text{Mean} = 10.48 \quad 12.51$$

$$\frac{n+1}{2} = \frac{73+1}{2} = 37$$

$$\text{Median} = 10$$

$$\text{Mode} = 12$$

$$CF (\text{Cumulative freq}) = 78 \quad (\text{Even})$$

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

$$= \frac{11 + 11}{2} = \frac{22}{2} = 11$$

wt	46	40	50	52	53	54	55
no of people	1	5	8	12	10	2	1

CF = 45 (odd)

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

$$= 23^{\text{th}} \text{ term}$$

$$\text{Median} = \left(\frac{12+1}{2} \right)^{\text{th}} \text{ term} = 52$$

$$\text{Mean} = \frac{46 \times 1 + 40 \times 5 + 50 \times 8 + 52 \times 12 + 53 \times 10 + 54 \times 2 + 55 \times 1}{45} = 50.64$$

$$\text{Mode} = 52$$

$$\text{Variance} = \frac{13, 5, 8, 13}{29.65}$$

$$\text{Mean} = \frac{17}{4} = 4.25$$

$$\text{Variance} = \left[\frac{\sum (\text{each value} - (\text{mean}))^2}{n} \right]$$

$$= \frac{(3-4.25)^2 + (5-4.25)^2 + (8-4.25)^2 + (13-4.25)^2}{4}$$

$$= 1.5625 + 0.5625 + 14.0625 + 10.5625$$

$$= 6.6875$$

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

$$\sigma^2 = \text{Variance} = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n} \quad (\text{Population})$$

$$\sigma (\text{Standard deviation}) = 2.586$$

$$\text{Co-Variance} =$$

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

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

$$\text{Correlation} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

$$\text{Cor}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

find mean, variance, covariance and correlation

Score	freq	Midpoint (m)
40-45	3	42.5
34-39	1	36.5
28-33	4	30.5
22-27	5	24.5
16-21	2	18.5
10-15	5	12.5
$N = 20$		

Mean = $\frac{\sum fx}{N} = \frac{3 \times 42.5 + 1 \times 36.5 + 4 \times 30.5 + 5 \times 24.5 + 2 \times 18.5 + 5 \times 12.5}{20} = 23.66$

Median = $Cf = 20$

$\frac{(N+1)}{2}$ th term

$\left(\frac{20+1}{2}\right)$ th term = $\left(\frac{21}{2}\right)$ th term

10th term + 11th term

$= \frac{24 + 24.5}{2}$

$= \frac{48.5}{2} = 24.25$

Mode = 24.5, 12.5

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

$= \frac{1}{19} [3(292.4) + (123.2) + (26.0) + (0.8) + (47.4) + (24.4) + (16.4)]$

$= 107.15$

$87.23 + 123.21 + 104.04 + 4.05 + 95.22 + 832.05$

X	Y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$
2	38	-5	10.6	-53
4	36	-3	8.6	-25.8
8	24	1	-3.4	-3.4
9	17	2	-10.4	-20.8
12	22	5	-5.4	-27

Sum =

Cov(x,y) = $\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N}$

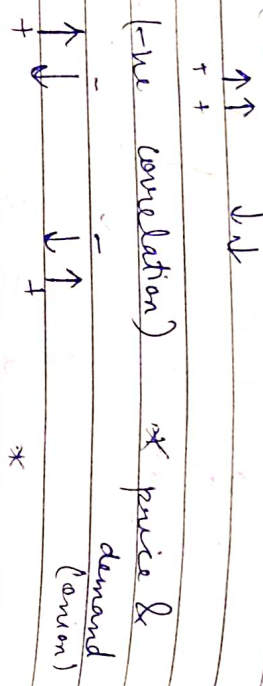
$= \frac{130}{5} = -26$

Covariation means relation b/w two variables.
 It describes the things (↑ ↓) (↑ ↑) (↓ ↓)
 1) Direction of relation (means strong or weak relation)
 2) Degree of relation

(↑ ↓) negative correlation
 correlation coefficient (r)
 correlation values lies b/w -1 to 1. except 0



Type of correlation
 There are three types of correlation
 (+ve correlation) * Age & Premium



Subtype of correlation coefficient (r)

perfect the correlation value either -1 or 1

No correlation (If the value is 0)

strong the correlation (lie b/w 0.8 to 1)

Moderate the correlation (lie b/w 0.5 to 0.8)
 low the correlation (lie b/w 0.1 to 0.5)
 -ve just above 0

Strong -ve correlation (-0.8 to -1)

Proof

$$r = \frac{n \sum xy - \sum x \sum y}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}}$$

Rank Pearson Correlation coefficient formula

X	Y
9	15
8	16
7	14
6	13
5	11
4	12
3	10
2	8
1	9

$$r = \frac{n \sum xy - \sum x \sum y}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}}$$

$$= \frac{0.85}{0.9} = 0.9$$

$$\text{Corr}(x,y) = \frac{\text{Cov}(x,y)}{\sqrt{\text{var}(x)} \sqrt{\text{var}(y)}} = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y}$$

Basic Probability

Probability is the chance that uncertain event will occur.

0 to 1
 impossible event confirm occurrence of event.
 (Sure to occur) (certain event)

There are three approaches to assessing the probability:

- ① Prior - based on prior knowledge or guess
- ② Empirical probability - based on experiment
- ③ Subjective probability - based on combination of subjective an individual's past experience (personal opinion & analysis of particular situation)

$$P = \frac{X}{T}$$

X = No. of ways in which event occurs.
 Total no. of event (Possible outcomes)

$P = \frac{\text{No. of ways in which the event occurs}}{\text{Total no. of possible outcomes}}$

$$P = \frac{X}{T}$$

What is

Q. Randomly selecting a day from the year 2015 what is the probability the day is in January.

$$P = \frac{31}{365}$$

	opt ML	Not opt ML	Total
Male	84	145	229
Female	76	134	210
Total	160	279	439

find the probability of selecting a male. opting ML from the population described in the following table.

$$P = \frac{84}{439}$$

Q. Events - there are three type of event.

- ① Simple event (Single characteristics)
- ② Joint event (more than two characteristics)

Complement of event A - All events that are not part of event A.

Sample space - collection of all possible events.

All six faces of dice is called sample space.

Q4

	Jan	Not Jan	
used	4	48	52
not used	27	286	313
Total	31	334	365

$$P_{\text{used}} = \frac{52}{365}$$

$$P_{\text{of Jan}} = \frac{31}{365}$$

$$P_{\text{of Jan \& used}} = \frac{4}{365}$$

$$P_{\text{of Jan \& not used}} = \frac{281}{365}$$

Mutually exclusive events - Event that cannot occur simultaneously. e.g. day year.
Ex: Randomly choosing a year 2015.

Event A = day in Jan
Event B = day in Feb
Event A and B are mutually exclusive event.

Collectively exhaustive event:

One of the event must occur.
The set of events covers the entire sample space.

Randomly chose a day from year 2015.

Event A = weekday
Event B = weekend.

Event C = January
Event D = Spring

Event C, B, D are collectively exhaustive but not mutually exclusive.

Event A and B are collectively exhaustive and also mutually exclusive.

Weekday cannot be Jan or Spring.
collectively exhaustive.

Joint and Marginal probability.

$P(A \text{ and } B) = \frac{\text{No. of outcome satisfying A and B}}{\text{Total no. of outcomes}}$

Q4 Addition Rule

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

Total no. of outcomes

If A and B are mutually exclusive then

$$P(A \text{ and } B) = 0$$

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

Q8

$$\begin{aligned}
 P(\text{Jan or Wed}) &= P(\text{Jan}) + P(\text{Wed}) - P(A \text{ and } B) \\
 &= \frac{31}{365} + \frac{52}{365} - \frac{4}{365} \\
 &= \frac{83-4}{365} = \frac{79}{365}
 \end{aligned}$$

Conditional Probability

$$P(A/B) = \frac{P(A \text{ and } B)}{P(B)}$$

$$P(B/A) = \frac{P(A \text{ and } B)}{P(A)}$$

Q9 In a used carlot 70% have AC and 40% have GPS. 20% have both. What is the probability that car has a GPS given that it has AC.

$$P(\text{GPS} | \text{AC}) = \frac{20}{70} = 0.2$$

even

chase

Independence

If two events are independent then $P(A)$ is not affected by $P(B)$.

$$P(A \text{ and } B) = P(A)$$

Multiplication rule of two event A and B

$$P(A \text{ and } B) = P(A/B) \cdot P(B)$$

If A and B are independent

$$P(A \text{ and } B) = P(A)P(B)$$