

# Exploratory Data Analysis and visualization

Mohaha Dharshin  
23BAID066

## UNIT - 1

### Analysis Techniques - Assignment

#### Question - 1

Scenario : Online Food Delivery performance.

A food delivery company recorded delivery times (in min) for one hour.

order ID	Delivery time (min)
01	22
02	24
03	21
04	23
05	25

tasks :

- a) Identify the level of measurement of the variable.

**Ratio Scale**

- \* Differences and ratios are meaningful.
- \* Because delivery time has a true zero (0 min means no time)

b) Compute mean and median

$$\text{Mean} = \frac{21 + 23 + 23 + 24 + 95}{5} = 37$$

\* Median (Middle value) = 23 (23 min)

c) Identify if an outlier exists.

\* 95 min is an outlier

\* It is far away from the other delivery times (19-23)

d) State which measure best represents typical delivery times and justify.

\* Median (23 min) best represents typical delivery time

\* Mean (37) is pulled upward due to the outlier (95)

\* Median is not affected by extreme values.

Question : Q : Student feedback Analysis

Scenario : Student feedback Analysis :

A college collected course feedback ratings.

Student	Rating (1-5)
S <sub>1</sub>	5
S <sub>2</sub>	4
S <sub>3</sub>	5
S <sub>4</sub>	2
S <sub>5</sub>	1

task:

- a) Identify the level of Measurement.

### Ordinal Scale

- \* Ratings (1-5) show order
- \* But the gap between 1 and 2 may not equal the gap between 4 and 5.

- b) Compute mode and median:

Mode (most frequent value): 5

Median = 4

- c) Is mean appropriate? explain.

\* Mean is not the best choice for ordinal data.

\* Ratings are categories with order, not exact numeric distances.

\* Mean assumes equal intervals between values, which may not be true.

d) Suggest a suitable visualization.

Bar chart (Frequency Ratings)

X-axis : Rating (1 to 5)

Y-axis : Number of Students.

Question : 3.

Scenerio : Temperature and Equipment Failure.

A factor records machine failures at different temperatures

Day	Temperature (°C)	Failures
D <sub>1</sub>	35	0
D <sub>2</sub>	32	1
D <sub>3</sub>	30	2
D <sub>4</sub>	28	4
D <sub>5</sub>	25	6

table:

a) Identify independent and dependent variables.

\*Independent variable : temperature (°C)

\*Dependent Variable : Number of failures

b) Compute mean temperature and mean failures.

\*Mean temperature :  $\frac{35 + 32 + 30 + 28 + 25}{5} = \frac{150}{5} = 30$

\*Mean Failures =  $\frac{0 + 1 + 2 + 4 + 6}{5} = \frac{13}{5} = 2.6$

c) Identify the trend

\* Negative relationship (inverse trend).

\* As temperature decreases, failures increase.

d) Comment on risk as temperature decreases.

\* Risk increases as temperature drops.

\* At  $35^{\circ}\text{C} \rightarrow 0$  failures

\* At  $25^{\circ}\text{C} \rightarrow 6$  failures. So machines are more likely to fail in colder temperatures, meaning higher breakdown risk.

#### Question 4

Scenario : Data cleaning in Banking

A bank records transaction amounts (₹)

Transaction ID	Amount
T <sub>1</sub>	5000
T <sub>2</sub>	7000
T <sub>3</sub>	-300
T <sub>4</sub>	7000
T <sub>5</sub>	(blank)

Tasks :

a) Identify data quality issues:

i. Negative amount (-300)  $\rightarrow$  may be invalid.

2. Missing value (blank) for T5
3. Duplicate amount (7000 repeated) → may be valid but needs checking.
4. Possible data entry error / inconsistent recording

b) Suggest cleaning actions

1) Check - 300

- \* If refund allowed → keep it and label as refund.
- \* if not allowed → correct / remove after verification.

2) Handle blank value.

- \* Fill using correct source record (blank log)
- \* if not available → mark as missing (Null) or remove row.

3) Verify duplicate amounts.

- \* Ensure T2 and T5 are real transactions, not duplicate entry.

4) Add validation rules:

\* Amount should be numeric

\* No blanks allowed

\* Negative only allowed for refund type

c) State the impact of not cleaning this data.

If not cleaned, it can cause:

- \* wrong total transaction Value.
- \* Wrong Customer balance reports.
- \* Incorrect fraud detection results.
- \* Bad business decisions ( profit / loss )
- \* Issues in auditing and compliance.

Question : 5 : Website traffic

Scenerio : Website traffic

Daily visitors recorded for a website

Day	Visitors
Mon	120
Tues	130
Wed	125
Thu	500
Fri	135

task:

a) Identify the outlier.

\* Thursday = 500 visitors is the outlier. ( because others are around 120 - 135 )

b) Compute Mean and Median

$$\text{Mean} = \frac{120 + 130 + 125 + 500 + 135}{5} = \frac{1010}{5} = 202$$

Median = 130

c) Which value is more representative and why?

\* Median (130) is more representative.

\* Mean (200) is inflated because of the outlier (500).

\* Median shows the normal daily traffic level.

d) What could cause such an outlier?

Possible reasons for 500 visitors:

\* Special offer / discount campaign, viral social media post, influencer promotion, festival / holiday traffic spike, email marketing campaign, website featured in news trending, Bot traffic / spam visits.

### Question 6 : Hospital patient Monitoring :

A hospital monitors patient heart rate.

Patient	Heart Rate (bpm)
P <sub>1</sub>	
P <sub>2</sub>	72
P <sub>3</sub>	75
P <sub>4</sub>	70
P <sub>5</sub>	180
P <sub>6</sub>	74
P <sub>7</sub>	73.

Tasks :

a) Identify the level of measurement

## Ratio Scale

- \* Heart rate has a true zero ( $0 \text{ bpm} = \text{no heartbeat}$ )
- \* Differences and ratios are meaningful.

b) Compute mean, median and range.

Mean:  $\frac{72 + 75 + 70 + 180 + 74 + 73}{6} = \frac{544}{6} = 90.67$

Median: average of 3rd and 4th values.

$$\frac{73 + 74}{2} = 73.5 \quad \text{Range} = 180 - 70 = 110$$

c) Identify outliers:

180 bpm (Patient P<sub>4</sub>) is an outlier because it is far higher than normal resting heart rates (around 60-100 bpm)

d) Interpret ~~clinical~~ clinical risk

\* P<sub>4</sub> = 180 bpm indicates high clinical risk.

\* Severe tachycardia, Stress / panic attack, Fever / infection, Heart rhythm problem (arrhythmia), Emergency condition needing immediate attention.

e) Why EDA is critical before prediction.

\* Detect outliers, identify data errors or abnormal readings  
Understand distribution, choose correct model, Improve prediction accuracy by cleaning and preprocessing

Question 7 : Sales data for Retail chain.

Monthly Sales (₹) of a Store.

Months	Jan	Feb	Mar	Apr	May	Jun
Sales	25	27	26	28	90	29

tasks:

a) Identify measurement scale.

Ratio Scale: Sales have true zero ( $0$  sales = no revenue)

Ratios make sense ( $90$  is  $3 \times$  of  $30$ )

b) Compute mean, median, standard deviation.

$$\text{Mean} = 27.5, \text{Median} = \frac{27+28}{2} = 27.5, \text{SD} = 23.52$$

$$\underline{\text{SD}} \text{ or Variance} = \frac{2317.5}{6} = 386.25 \quad \text{SD} = \sqrt{386.25} = 23.52 \text{ lacs}$$

c) Detect outliers

May = 90 lacs is an outlier

d) How outlier affects decision-making

\* Mean becomes very high (37.5) even though normal sales are 27-29.

\* Company may assume store is performing better than reality.

\* Wrong inventory planning.

\* Wrong sales targets.

\* Incorrect budgeting and forecasting.