# Phase 2 – Data Analyst Theory Tasks

## 1. End-to-End Data Analysis Lifecycle (Real Business Example)

The end-to-end data analysis lifecycle is a structured process used by data analysts to convert raw data into meaningful business insights.
Consider an e-commerce company that wants to reduce product return rates.

First, the company performs business understanding, where the main issue—high return rate—is clearly identified.
Next, in problem definition, a measurable goal is set, such as reducing returns by 15% within three months.

**After defining the goal, the analyst performs data collection from different sources like:**

- Sales transaction records

- Return history logs

- Customer feedback and reviews

- Product catalog information

**The collected data is then cleaned by:**

- Removing duplicate entries

- Handling missing values

- Correcting incorrect formats

Once the data is clean, data integration combines multiple datasets into a single structured dataset.

**The analyst then performs Exploratory Data Analysis (EDA) to find:**

- Frequently returned products

- Regions with higher returns

- Customer behavior patterns

Next comes feature selection and statistical analysis, where important variables influencing returns are identified.

The results are presented through visualization and reporting, such as dashboards and charts, so business managers can easily understand them.

Finally, insights are generated, decisions are taken (like improving product quality or descriptions), and results are continuously monitored for improvement.

This complete workflow represents the real-world data analysis lifecycle.

---

## 2. Descriptive vs Diagnostic Analytics; Correlation vs Causation

<u>**Descriptive Analytics**</u>

Descriptive analytics explains **what has happened in the past** by summarizing historical data.

**Characteristics:**

- Uses totals, averages, and summaries

- Presented through reports and dashboards

- Helps track past performance

**Example:**
Total sales achieved in the previous month.

---

## Diagnostic Analytics

Diagnostic analytics explains **why something happened**.

**Characteristics:**

- Identifies root causes of problems

- Uses comparisons, drill-downs, and correlations

- Helps in problem solving and decision making

**Example:**
Finding the reason for a sudden drop in sales.

---

## Correlation

Correlation shows a **relationship between two variables**, meaning they move together.

- Does **not prove cause**

- Can be positive or negative

- Used mainly in exploratory analysis

**Example:**
Ice cream sales increase when temperature rises.

---

## Causation

Causation means **one variable directly causes another**.

- Shows true cause-and-effect

- Requires logical proof or experimentation

- More reliable for business decisions

**Example:**
Higher temperature **causes** increased ice cream demand.

Understanding the difference between correlation and causation is important to **avoid wrong conclusions**.

---

## 3. Short Notes

### Data Bias

Data bias occurs when a dataset **does not properly represent reality**, leading to incorrect insights.

**Causes:**

- Poor sampling
- Limited demographic coverage
- Human assumptions

**Effects:**

- Misleading results
- Unfair or inaccurate decisions
- Reduced reliability of analysis

Bias can be reduced by **using diverse and balanced data sources**.

---

### Missing Data Strategies

Missing data is common in real datasets. Analysts handle it using:

- Removing incomplete rows or columns
- Filling values with **mean, median, or mode**
- Forward or backward filling
- Interpolation or prediction methods
- Treating missing values as a separate category

The correct method depends on **data importance and context**.

---

### KPIs vs Metrics

**KPIs (Key Performance Indicators):**

- Measure **critical business goals**
- Directly linked to success
- Few in number and strategic

**Metrics:**

- Measure **general performance**

- Support KPIs but are not always strategic

- Larger in number

**Example:**

- KPI → Revenue growth rate

- Metric → Daily sales count

All KPIs are metrics, but **not all metrics are KPIs**.

---

## 4. Case Study: Why Dashboards Fail Even with Correct Data

Even when data is accurate, dashboards may fail due to **design and usability issues**.

**Common reasons:**

- Not aligned with business goals

- Too many charts causing confusion

- Wrong choice of visualization

- KPIs not clearly defined

- Data not updated regularly

- Poor layout and complex design

- Users not trained to interpret data

- Insights not explained in simple language

- No clear action points or recommendations

- Dashboard not interactive or role-specific

A successful dashboard should be:

- Simple and clear

- Focused on important KPIs

- Easy for non-technical users to understand

- Able to provide actionable business insights