



ISDM (INDEPENDENT SKILL DEVELOPMENT MISSION)



OVERVIEW OF DATA ANALYTICS: IMPORTANCE & APPLICATIONS



CHAPTER 1: UNDERSTANDING DATA ANALYTICS

1.1 What is Data Analytics?

Data Analytics is the **systematic computational analysis of data** that helps organizations and individuals extract insights from raw data. It involves various processes, including **data collection, cleaning, transformation, visualization, and interpretation**, to support **decision-making, forecasting, and business intelligence**.

Data Analytics is widely used across industries to improve **efficiency, productivity, and profitability**. With advancements in **Big Data, Artificial Intelligence (AI), and Cloud Computing**, Data Analytics is becoming more powerful, enabling **real-time analysis and predictive decision-making**.

- ◆ **Key Disciplines in Data Analytics:**
- ✓ **Statistics & Mathematics** – Provides a foundation for probability, predictive modeling, and statistical inference.
- ✓ **Computer Science & Programming** – Uses **Python, R, SQL, and AI tools** to manipulate and analyze data.

- ✓ **Business Intelligence (BI)** – Converts raw data into **actionable insights and strategic decisions**.
- ✓ **Machine Learning & AI** – Automates processes and helps **predict trends, classify data, and enhance decision-making**.

📌 Example:

A **banking institution** uses Data Analytics to detect **fraudulent transactions** by analyzing past patterns and real-time activity.

💡 Conclusion:

Data Analytics is **critical** in today's data-driven world, enabling businesses and individuals to make informed, efficient, and optimized decisions.

📌 CHAPTER 2: IMPORTANCE OF DATA ANALYTICS

Data Analytics is **revolutionizing industries** by turning data into **valuable insights** that drive innovation, efficiency, and automation. With increasing data availability, businesses rely on analytics to stay competitive.

2.1 Why is Data Analytics Important?

- ✓ **Better Decision-Making** – Helps organizations make **data-driven, evidence-based** decisions rather than relying on guesswork.
- ✓ **Operational Efficiency** – Improves **workflows, automation, and cost reduction** through optimized data usage.
- ✓ **Customer Insights** – Enables businesses to **personalize user experiences** and improve customer satisfaction.
- ✓ **Fraud Detection & Risk Management** – Identifies **unusual patterns** to prevent cyber fraud in banking and finance.
- ✓ **Improved Healthcare Outcomes** – Supports **disease prediction, treatment planning, and medical research** using patient data.

✓ **Predictive Capabilities** – Forecasts future trends, market behavior, and economic conditions.

📌 **Example:**

Hospitals use **predictive analytics** to analyze patient symptoms and **forecast disease outbreaks**, improving medical response.

💡 **Conclusion:**

Data Analytics **empowers businesses and industries** by making processes more **efficient, cost-effective, and forward-looking**.

📌 **CHAPTER 3: TYPES OF DATA ANALYTICS**

There are **four major types of Data Analytics**, each serving different analytical purposes.

3.1 Categories of Data Analytics

- ◆ **Descriptive Analytics** – Answers "What happened?" by summarizing past data. It helps in **trend analysis, business reports, and dashboards**.

✓ **Example:** A retail store tracks **monthly sales performance** using historical sales data.

- ◆ **Diagnostic Analytics** – Answers "Why did it happen?" by identifying relationships and patterns. It is useful in **identifying root causes of problems**.

✓ **Example:** An HR department analyzes why **employee turnover increased in a specific quarter**.

- ◆ **Predictive Analytics** – Answers "What is likely to happen?" by using past data to forecast trends.

✓ **Example:** An e-commerce platform predicts **which products will be in high demand next season.**

◆ **Prescriptive Analytics** – Answers "What should be done?" by providing actionable recommendations based on AI and simulations.

✓ **Example:** Airlines use **Prescriptive Analytics** to optimize **ticket pricing** based on demand, competitor pricing, and historical trends.

 **Conclusion:**

Each type of analytics plays a role in helping organizations **understand past performance, predict future outcomes, and take data-driven actions.**

CHAPTER 4: APPLICATIONS OF DATA ANALYTICS IN INDUSTRIES

Data Analytics is applied in diverse industries to **improve efficiency, drive innovation, and enhance decision-making.**

4.1 How Different Sectors Use Data Analytics

 **Retail & E-Commerce**

 **Recommendation engines** suggest products based on user behavior.

 **Market Basket Analysis** helps businesses understand purchase patterns.

 **Dynamic Pricing Models** adjust pricing in real-time based on demand.

 **Healthcare & Pharmaceuticals**

 **AI-driven diagnostic tools** enhance early disease detection.

 **Electronic Health Records (EHRs)** help predict patient health

risks.

Clinical trial analytics improve **drug discovery efficiency**.

 **Finance & Banking**

Fraud detection algorithms analyze patterns to detect fraud.

Credit risk models assess customer eligibility for loans.

Algorithmic trading makes **real-time stock predictions**.

 **Manufacturing & Logistics**

Predictive maintenance prevents machine failures.

Supply chain analytics optimizes inventory and reduces costs.

Automated quality control detects product defects.

 **Example:**

Netflix uses **Data Analytics** to recommend **personalized content** based on viewing history.

 **Conclusion:**

Data Analytics **transforms industries by reducing costs, increasing efficiency, and improving customer experiences**.

 **CHAPTER 5: DATA ANALYTICS PROCESS**

A **structured approach** ensures **data accuracy, reliability, and insight generation**.

5.1 Steps in a Data Analytics Project

✓ Step 1: Problem Identification – Define business goals.

✓ Step 2: Data Collection – Gather structured and unstructured data.

✓ Step 3: Data Cleaning & Preparation – Handle missing values and format inconsistencies.

- ✓ **Step 4: Exploratory Data Analysis (EDA)** – Analyze trends, outliers, and relationships.
- ✓ **Step 5: Feature Engineering** – Convert raw data into **structured variables**.
- ✓ **Step 6: Model Building & Machine Learning** – Develop predictive and prescriptive models.
- ✓ **Step 7: Data Visualization & Reporting** – Present insights in an easy-to-understand format.
- ✓ **Step 8: Decision-Making & Optimization** – Implement findings for business improvements.

 **Example:**

A **retail chain** optimizes supply chain operations by analyzing **customer purchasing patterns**.

 **Conclusion:**

A well-defined **Data Analytics process** ensures **accurate, efficient, and impactful** decision-making.

 **CHAPTER 6: TOOLS & TECHNOLOGIES IN DATA ANALYTICS**

6.1 Popular Data Analytics Tools

- ✓ **Python & R** – Used for statistical analysis and machine learning.
- ✓ **SQL** – Extracts and manages structured data.
- ✓ **Excel & Google Sheets** – Basic data analysis tools.
- ✓ **Power BI & Tableau** – Business Intelligence dashboards.
- ✓ **Apache Spark & Hadoop** – Big Data frameworks.
- ✓ **Google Cloud AI & AWS ML** – Cloud-based data analytics platforms.

📌 Example:

Banks use **Power BI dashboards** to monitor financial performance in real-time.

💡 Conclusion:

Mastering **Data Analytics tools** is essential for working professionals and data enthusiasts.

📌 CHAPTER 7: EXERCISES & ASSIGNMENTS

7.1 Multiple Choice Questions

✓ Which tool is commonly used for data visualization?

- (a) SQL
- (b) Power BI
- (c) MySQL

✓ What type of analytics helps predict future trends?

- (a) Descriptive
- (b) Predictive
- (c) Diagnostic

7.2 Practical Assignment

📌 Task: Conduct **Exploratory Data Analysis (EDA)** on a dataset using **Python** and **Power BI** and present findings.

🌟 CONCLUSION: THE FUTURE OF DATA ANALYTICS

As data becomes more valuable, **AI, Big Data, and real-time analytics** will continue to shape the future. Data-driven insights will

power **better decision-making, automation, and smarter businesses.**



SETTING UP PYTHON ENVIRONMENT (JUPYTER NOTEBOOK, GOOGLE COLAB)

CHAPTER 1: INTRODUCTION TO PYTHON ENVIRONMENTS

1.1 Why Setting Up a Python Environment is Important?

A Python environment is a **workspace** where you can write, execute, and test your Python code efficiently. A well-configured Python environment is **essential for Data Science, Machine Learning, and Data Analytics** as it provides necessary libraries and tools.

- ◆ **Benefits of Setting Up a Python Environment:**
- ✓ **Organized Coding** – Allows structured development and testing.
- ✓ **Library Management** – Install and update libraries efficiently.
- ✓ **Reproducibility** – Ensures consistency in experiments and projects.
- ✓ **Collaborative Work** – Enables easy sharing of notebooks.

Conclusion:

A properly set up Python environment **enhances efficiency, reduces errors, and enables smooth execution of projects**.

CHAPTER 2: INTRODUCTION TO JUPYTER NOTEBOOK

2.1 What is Jupyter Notebook?

Jupyter Notebook is an **open-source, interactive computing environment** that supports multiple programming languages, including Python. It is widely used for **Data Science, Machine Learning, and Data Visualization.**

◆ **Key Features of Jupyter Notebook:**

- ✓ **Cell-Based Execution** – Allows running specific blocks of code.
- ✓ **Supports Markdown** – Enables documentation alongside code.
- ✓ **Data Visualization Support** – Displays charts, graphs, and images.
- ✓ **Integration with Machine Learning Libraries** – Works seamlessly with **Pandas, NumPy, Matplotlib, and TensorFlow.**

📌 **Example:**

A **data scientist** uses Jupyter Notebook to **load a dataset, clean data, visualize trends, and build models** in a single document.

💡 **Conclusion:**

Jupyter Notebook is **one of the most powerful tools for interactive coding, documentation, and data analysis.**

2.2 Installing Jupyter Notebook

Jupyter Notebook can be installed using **Anaconda** (recommended) or **pip** (lightweight option).

Method 1: Installing Jupyter Notebook via Anaconda (Recommended)

Download & Install Anaconda

- Visit the official **Anaconda website.**

- Download and install **Anaconda Distribution** (Python 3.x version).

☒ Launch Jupyter Notebook

- Open **Anaconda Navigator** and click "Launch Jupyter Notebook".
- Alternatively, use the terminal (or Anaconda Prompt) and run:
 - `jupyter notebook`

Method 2: Installing Jupyter Notebook via pip (Lightweight Alternative)

☒ Install Jupyter Notebook using pip:

`pip install notebook`

☒ Launch Jupyter Notebook:

`jupyter notebook`

💡 Conclusion:

Installing Jupyter Notebook is simple, and using **Anaconda** is recommended for beginners due to its pre-installed libraries.

2.3 Navigating Jupyter Notebook Interface

Once launched, Jupyter Notebook opens in a web browser.

- ◆ **Key Components:**
- ✓ **Notebook Dashboard** – Displays all files and notebooks.
- ✓ **Code Cells** – Executes Python code interactively.
- ✓ **Markdown Cells** – Writes text, comments, and documentation.

✓ **Toolbar & Shortcuts** – Run, stop, restart code, and format content.

📌 **Example:**

Running a simple Python command in Jupyter Notebook:

```
print("Hello, Jupyter Notebook!")
```

💡 **Conclusion:**

Jupyter Notebook provides an **interactive** and **user-friendly interface** for Python programming and Data Science.

📌 **CHAPTER 3: INTRODUCTION TO GOOGLE COLAB**

3.1 What is Google Colab?

Google Colab (Colaboratory) is a **cloud-based Jupyter Notebook** environment that allows users to **write and execute Python code without any local installation**.

◆ **Key Features of Google Colab:**

- ✓ **No Installation Required** – Runs in the cloud via a browser.
- ✓ **Free Access to GPUs & TPUs** – Enables faster computations.
- ✓ **Integration with Google Drive** – Easily saves and shares projects.
- ✓ **Supports Python Libraries** – Pre-installed with **TensorFlow, PyTorch, OpenCV, and Pandas**.

📌 **Example:**

A Machine Learning Engineer uses **Google Colab** to train deep learning models using free GPUs without needing expensive hardware.

Conclusion:

Google Colab **removes hardware limitations** and is ideal for cloud-based Machine Learning and Data Science tasks.

3.2 How to Access Google Colab?

- ☒ Visit Google Colab – Open [Google Colab](#) in a web browser.
- ☒ Sign in with Google – Use a **Google account** to access Colab.
- ☒ Create a New Notebook – Click "**New Notebook**" to start coding.

Example:

Running a simple Python command in Google Colab:

```
print("Hello, Google Colab!")
```

Conclusion:

Google Colab is a **powerful cloud-based alternative** to Jupyter Notebook, **eliminating the need for installation**.

3.3 Key Features of Google Colab

Google Colab enhances productivity with additional features:

- ✓ Run Code in the Cloud – No need for a local setup.
- ✓ Collaborate in Real-Time – Multiple users can edit and run code simultaneously.
- ✓ Free GPU & TPU Access – Speeds up **Deep Learning & AI** computations.
- ✓ Automatic File Saving – Saves progress to **Google Drive**.

Example:

Using **Google Drive** with Google Colab:

```
from google.colab import drive
drive.mount('/content/drive')
```

 **Conclusion:**

Google Colab improves collaboration, speeds up AI models, and ensures seamless cloud integration.

 **CHAPTER 4: JUPYTER NOTEBOOK VS GOOGLE COLAB**

| Feature | Jupyter Notebook | Google Colab |
|------------------------|---------------------------|-----------------------|
| Installation Required? | Yes | No (Cloud-based) |
| Runs Locally? | Yes | No |
| GPU/TPU Support? | No | Yes (Free) |
| Best For | Exploratory Data Analysis | Machine Learning & AI |
| File Storage | Local | Google Drive |

 **Conclusion:**

Both tools have advantages:

- **Jupyter Notebook** is great for offline, structured coding.
- **Google Colab** is ideal for collaborative, cloud-based AI & Data Science projects.

 **CHAPTER 5: EXERCISES & ASSIGNMENTS**

5.1 Multiple Choice Questions

Which tool provides cloud-based Jupyter Notebooks?

- (a) VS Code
- (b) Google Colab
- (c) PyCharm

Which command starts Jupyter Notebook from the terminal?

- (a) open notebook
- (b) jupyter notebook
- (c) start notebook

Which feature makes Google Colab different from Jupyter Notebook?

- (a) Works offline
- (b) Provides free cloud GPUs
- (c) Requires installation

5.2 Practical Assignment

- 📌 **Task 1:** Install Jupyter Notebook using Anaconda or pip and run a Python script.
- 📌 **Task 2:** Access Google Colab, create a new notebook, and execute a basic Python program.
- 📌 **Task 3:** Compare execution speed of a Deep Learning Model on Jupyter vs. Google Colab GPUs.

📌 CHAPTER 6: SUMMARY

Jupyter Notebook is an **interactive Python environment** used for **Data Science & Machine Learning**.

Google Colab is a **cloud-based alternative** that offers **free GPU**

access and real-time collaboration.

- Both tools are essential for Data Analytics, AI, and Machine Learning workflows.
-

🌟 CONCLUSION: THE FUTURE OF PYTHON ENVIRONMENTS

With increasing reliance on AI and Cloud Computing, Google Colab and Jupyter Notebooks will remain the foundation for Data Science and Machine Learning projects.

ISDM-NXT



DATA HANDLING WITH PANDAS & NUMPY



CHAPTER 1: INTRODUCTION TO DATA HANDLING

1.1 What is Data Handling?

Data handling refers to the process of **loading, cleaning, manipulating, and analyzing** data efficiently. In **Data Science and Analytics**, handling structured and unstructured data correctly is essential for extracting insights.

With the help of Python libraries like **NumPy (Numerical Python)** and **Pandas (Python Data Analysis Library)**, data handling becomes **faster, scalable, and more efficient**. These libraries help in:

- ✓ Storing and manipulating large datasets efficiently.
- ✓ Performing mathematical and statistical operations on data.
- ✓ Cleaning and transforming data into usable formats.
- ✓ Handling missing values, duplicates, and outliers.



Example:

A retail company collects customer **transaction data** in CSV format. Using Pandas and NumPy, analysts **clean the data, analyze sales trends, and optimize inventory**.



Conclusion:

Effective **data handling** ensures **accurate, consistent, and structured** data for meaningful analysis.

❖ CHAPTER 2: INTRODUCTION TO NUMPY

NumPy is the **foundation for numerical computing in Python**. It provides high-performance multi-dimensional arrays and mathematical functions for **fast computations**.

2.1 Key Features of NumPy

- ✓ **Efficient Storage** – Handles large datasets with optimized memory usage.
- ✓ **Fast Computations** – Uses C-based implementation for high-speed operations.
- ✓ **Vectorized Operations** – Performs element-wise computations without loops.
- ✓ **Linear Algebra & Statistics** – Provides advanced mathematical functions.
- ✓ **Interoperability** – Works seamlessly with Pandas, SciPy, and Matplotlib.

❖ Example:

Traditional List vs. NumPy Array Computation

```
import numpy as np  
  
# Using a Python list (slower)  
list_data = [1, 2, 3, 4, 5]  
  
list_result = [x * 2 for x in list_data]
```

```
# Using NumPy Array (faster)  
  
numpy_array = np.array([1, 2, 3, 4, 5])
```

```
numpy_result = numpy_array * 2
```

```
print(numpy_result) # Output: [2 4 6 8 10]
```

💡 Conclusion:

NumPy **enhances performance** and **reduces processing time** for large-scale numerical computations.

📌 CHAPTER 3: WORKING WITH NUMPY ARRAYS

3.1 Creating NumPy Arrays

NumPy arrays are **faster, more memory-efficient, and support advanced operations** compared to Python lists.

✓ Creating a 1D Array

```
import numpy as np  
  
arr = np.array([10, 20, 30, 40, 50])  
  
print(arr)
```

✓ Creating a 2D Array (Matrix)

```
matrix = np.array([[1, 2, 3], [4, 5, 6]])  
  
print(matrix)
```

✓ Creating an Array of Zeros

```
zero_array = np.zeros((3, 3))  
  
print(zero_array)
```

✓ Creating an Array of Random Values

```
random_array = np.random.rand(3, 3)  
print(random_array)
```

 **Conclusion:**

NumPy arrays are **highly flexible** and can be **easily created, modified, and used for advanced operations.**

 **CHAPTER 4: NUMPY OPERATIONS & FUNCTIONS**

4.1 Basic Array Operations

NumPy supports **element-wise operations** and mathematical computations.

✓ Addition, Subtraction, Multiplication, and Division

```
arr1 = np.array([1, 2, 3])
```

```
arr2 = np.array([4, 5, 6])
```

```
print(arr1 + arr2) # Output: [5 7 9]
```

```
print(arr1 - arr2) # Output: [-3 -3 -3]
```

```
print(arr1 * arr2) # Output: [4 10 18]
```

```
print(arr1 / arr2) # Output: [0.25 0.4 0.5]
```

✓ Statistical Functions

```
data = np.array([10, 20, 30, 40, 50])
```

```
print(np.mean(data)) # Output: 30.0
```

```
print(np.median(data)) # Output: 30.0
```

```
print(np.std(data)) # Standard Deviation
```

 **Conclusion:**

NumPy simplifies **complex mathematical and statistical operations** with built-in functions.

 **CHAPTER 5: INTRODUCTION TO PANDAS**

Pandas is a **powerful data manipulation library** used for **analyzing and organizing structured data**.

5.1 Key Features of Pandas

- ✓ **DataFrame & Series** – Structured data handling in tabular format.
- ✓ **Data Cleaning** – Handles missing values, duplicates, and inconsistencies.
- ✓ **Data Aggregation & Grouping** – Summarizes data efficiently.
- ✓ **Integration with NumPy & Visualization Libraries** – Works seamlessly with Matplotlib and Seaborn.

 **Example:**

Loading a CSV file using Pandas

```
import pandas as pd
```

```
df = pd.read_csv("data.csv")
```

```
print(df.head()) # Displays the first 5 rows
```

 **Conclusion:**

Pandas **simplifies data manipulation** and is essential for **handling large datasets efficiently**.



CHAPTER 6: WORKING WITH PANDAS DATAFRAMES

6.1 Creating a DataFrame

A DataFrame is a **table-like structure with rows and columns**, similar to an Excel sheet.

✓ Creating a DataFrame from a Dictionary

```
data = {  
    "Name": ["Alice", "Bob", "Charlie"],  
    "Age": [25, 30, 35],  
    "Salary": [50000, 60000, 70000]  
}
```

```
df = pd.DataFrame(data)  
print(df)
```

✓ Loading Data from CSV

```
df = pd.read_csv("employee_data.csv")  
print(df.head())
```



Conclusion:

Pandas DataFrames make **data storage, manipulation, and visualization** easier.

📌 CHAPTER 7: PANDAS DATA MANIPULATION

7.1 Handling Missing Values

✓ Checking Missing Values

```
print(df.isnull().sum())
```

✓ Filling Missing Values with Mean

```
df["Salary"].fillna(df["Salary"].mean(), inplace=True)
```

✓ Dropping Rows with Missing Values

```
df.dropna(inplace=True)
```

7.2 Filtering & Selecting Data

✓ Selecting a Specific Column

```
print(df["Name"])
```

✓ Filtering Data Based on Condition

```
filtered_data = df[df["Age"] > 30]
```

```
print(filtered_data)
```

💡 Conclusion:

Pandas provides **efficient functions** to clean and filter large datasets quickly.

📌 CHAPTER 8: DATA VISUALIZATION WITH PANDAS & NUMPY

8.1 Plotting Data Using Pandas

✓ Line Plot

```
df.plot(kind="line", x="Age", y="Salary")
```

✓ Bar Chart

```
df.plot(kind="bar", x="Name", y="Salary")
```

✓ Histogram

```
df["Age"].hist()
```

📌 Example:

Analyzing **salary trends** in an organization using a **bar chart**.

💡 Conclusion:

Visualizing data makes it **easier to interpret trends and patterns**.

📌 CHAPTER 9: EXERCISES & ASSIGNMENTS

9.1 Multiple Choice Questions

✓ What does NumPy primarily deal with?

- (a) Data visualization
- (b) Numerical computations
- (c) Text analysis

9.2 Practical Assignment

📌 **Task:** Load a dataset using Pandas, clean missing values, and create a visualization using Matplotlib.

🌟 CONCLUSION: IMPORTANCE OF DATA HANDLING

NumPy and Pandas are **essential tools for handling, processing, and analyzing data efficiently**. Mastering these libraries is **crucial for data-driven decision-making** in various industries. 🚀📊

ISDM-NxT



DATA CLEANING & TRANSFORMATION TECHNIQUES

📌 CHAPTER 1: INTRODUCTION TO DATA CLEANING

1.1 What is Data Cleaning?

Data Cleaning, also known as **data cleansing** or **data scrubbing**, is the process of **identifying, correcting, or removing errors, inconsistencies, and inaccuracies** in a dataset. It is an essential step in **data preprocessing** to ensure **data quality, accuracy, and reliability** before analysis or machine learning modeling.

Data Cleaning involves **handling missing values, removing duplicates, correcting errors, and ensuring consistency** across datasets. Poor data quality can lead to **incorrect insights and flawed decision-making**, making data cleaning a **critical step in any data-driven process**.

◆ Why is Data Cleaning Important?

- ✓ Ensures **accuracy and reliability** of data for decision-making.
- ✓ Improves **model performance** in Machine Learning and AI applications.
- ✓ Reduces **errors and inconsistencies** in business reporting.
- ✓ Enhances **data usability** for analytics, dashboards, and reporting.

📌 Example:

A retail company analyzing customer purchase behavior must remove duplicate transactions and correct missing values to generate accurate insights.

Conclusion:

Data Cleaning is **essential for data integrity**, ensuring the dataset is accurate, consistent, and free from errors.

CHAPTER 2: COMMON DATA QUALITY ISSUES

Before cleaning data, it is crucial to identify the most common issues that affect **data quality and consistency**.

2.1 Types of Data Quality Issues

- ◆ **Missing Data** – When some values are missing due to **human errors, system failures, or incomplete records**.
- ◆ **Duplicate Data** – When the same record appears multiple times due to **system glitches or data entry mistakes**.
- ◆ **Inconsistent Formatting** – Variations in **date formats, currency symbols, or text values** (e.g., "New York" vs. "NY").
- ◆ **Incorrect Data Types** – Numerical data stored as **text**, or dates stored as **strings** instead of date format.
- ◆ **Outliers & Anomalies** – Extremely **high or low values** that distort data distribution.
- ◆ **Typographical Errors** – Spelling mistakes or **incorrect labeling** in categorical variables.
- ◆ **Irrelevant Data** – Unnecessary **columns or records** that do not contribute to analysis.

Example:

A dataset with customer transactions may contain missing values in the "Age" column, duplicates in the "Customer ID" column, and inconsistent date formats.

💡 Conclusion:

Identifying **data quality issues** is the first step in ensuring that datasets are prepared for accurate analysis and decision-making.

📌 CHAPTER 3: HANDLING MISSING DATA

3.1 Why Does Data Go Missing?

- ✓ Data entry errors and **manual omissions**
- ✓ System failures or **data migration issues**
- ✓ Sensors failing to **capture data properly**
- ✓ Respondents skipping questions in **surveys**

3.2 Techniques for Handling Missing Data

- ✓ **Remove Missing Values** – If only a small portion of data is missing, dropping the missing rows or columns can be a viable solution.
- ✓ **Mean/Median/Mode Imputation** – Replace missing numerical values with the **mean (average)**, **median**, or **mode (most frequent value)** of the column.
- ✓ **Forward Fill (ffill) & Backward Fill (bfill)** – Fill missing values using the previous or next available value in a sequence (e.g., time-series data).
- ✓ **Interpolate Missing Values** – Use mathematical techniques to estimate missing values based on trends.
- ✓ **Predictive Imputation** – Use **Machine Learning models** such as K-Nearest Neighbors (KNN) or Regression to predict and fill missing values.

📌 Example:

In a medical dataset, missing **blood pressure values** can be filled using the **median** value of the dataset to ensure consistency.

💡 Conclusion:

Choosing the right method to handle missing data **depends on the dataset size, type, and business context**.

📌 CHAPTER 4: REMOVING DUPLICATE DATA

4.1 What Causes Duplicate Records?

- ✓ Data migration from **multiple sources**
- ✓ **Repeated entries** due to system errors
- ✓ Copy-pasting errors in **Excel or CSV files**

4.2 How to Remove Duplicates?

- ✓ **Identify Duplicates** – Use SQL queries (SELECT DISTINCT) or Pandas (df.duplicated()) to find duplicate records.
- ✓ **Remove Exact Duplicates** – Use .drop_duplicates() in Pandas or DELETE FROM table WHERE row_number > 1 in SQL.
- ✓ **Partial Duplicates Handling** – Merge similar records by **grouping identical entries** while keeping key attributes.
- ✓ **Verify Unique Identifiers** – Ensure columns like **Customer ID, Order ID, and Transaction ID** are unique.

📌 Example:

In an **e-commerce dataset**, a customer may have **two identical purchase records**, which should be removed to prevent **double-counting sales**.

💡 Conclusion:

Duplicate records **distort analysis and lead to incorrect conclusions**; removing them improves **data accuracy**.

📌 CHAPTER 5: HANDLING OUTLIERS & INCONSISTENCIES

5.1 What Are Outliers?

An outlier is a data point **significantly different** from the rest of the dataset. Outliers can occur due to **errors in data entry, fraud, or natural variations**.

5.2 Methods to Detect Outliers

- ✓ **Box Plots & IQR (Interquartile Range)** – Detects extreme values beyond the **1.5x IQR range**.
- ✓ **Z-score Method** – Identifies values that are **more than 3 standard deviations away from the mean**.
- ✓ **Histogram & Scatter Plots** – Visualizes distribution and unusual spikes in data.

5.3 How to Handle Outliers?

- ✓ **Remove Outliers** – If caused by errors, deleting outliers ensures data integrity.
- ✓ **Cap Values (Winsorization)** – Limits extreme values to a defined percentile (e.g., **top 1% and bottom 1%**).
- ✓ **Transform Data** – Use log transformation or scaling to reduce the effect of extreme values.

📌 Example:

A dataset analyzing **house prices** might have a few **extremely high values** due to rare luxury properties, which should be adjusted for accurate modeling.

💡 Conclusion:

Handling outliers ensures that **data distribution remains normal and results are not skewed.**

📌 CHAPTER 6: DATA TRANSFORMATION TECHNIQUES

6.1 What is Data Transformation?

Data Transformation is the process of **converting raw data into a structured, consistent format** to improve usability and model performance.

6.2 Common Data Transformation Techniques

- ✓ **Normalization & Standardization** – Scales numerical values to ensure **fair comparison** across variables.
- ✓ **One-Hot Encoding & Label Encoding** – Converts categorical variables into **numerical format** for machine learning models.
- ✓ **Feature Engineering** – Creates new meaningful variables from existing data (e.g., extracting **year** from a date column).
- ✓ **Data Aggregation** – Summarizes data by **grouping and aggregating key values** (e.g., total sales per month).

📌 Example:

Converting "Yes" and "No" responses into **binary 1s and 0s** helps models interpret categorical data effectively.

💡 Conclusion:

Data Transformation **enhances data usability, consistency, and prepares it for advanced analytics.**

📌 CHAPTER 7: EXERCISES & ASSIGNMENTS

7.1 Multiple Choice Questions

☒ Which of the following is a method to handle missing data?

- (a) Deleting the entire dataset
- (b) Filling with median values
- (c) Ignoring it completely

☒ What is the purpose of data normalization?

- (a) Convert numbers to text
- (b) Scale data for uniformity
- (c) Remove duplicates

7.2 Practical Assignment

- 📌 Task: Perform Data Cleaning & Transformation on a dataset:
- Handle missing values, duplicates, and outliers using Python or SQL.
 - Convert categorical data into numerical format for machine learning.

🌟 CONCLUSION

Data Cleaning & Transformation ensure that data is accurate, complete, and ready for analysis. Properly processed data improves insights, reduces errors, and enhances decision-making.





INTRODUCTION TO SQL FOR DATA ANALYTICS

📌 CHAPTER 1: UNDERSTANDING SQL FOR DATA ANALYTICS

1.1 What is SQL?

SQL (**S**tructured **Q**uery **L**anguage) is a standardized programming language used for managing and manipulating **relational databases**. It allows users to **store, retrieve, update, and delete** data efficiently. SQL is widely used in **Data Analytics, Business Intelligence, and Data Science** to extract meaningful insights from databases.

SQL works with **Relational Database Management Systems (RDBMS)** such as:

- ✓ **MySQL** – Open-source database widely used in web applications.
- ✓ **PostgreSQL** – Advanced open-source database known for scalability.
- ✓ **Microsoft SQL Server** – Enterprise-level database for business applications.
- ✓ **Oracle Database** – Large-scale database used in financial and enterprise systems.

📌 Example:

An e-commerce website uses SQL to store **customer orders, product details, and transactions**, enabling fast retrieval of sales data.

💡 Conclusion:

SQL is the **foundation of data storage and retrieval**, making it an

essential tool for **data analysts, engineers, and business professionals.**

📌 CHAPTER 2: IMPORTANCE OF SQL IN DATA ANALYTICS

SQL is **widely used in Data Analytics** because it enables users to **query large datasets efficiently** and extract valuable business insights.

2.1 Why is SQL Important for Data Analytics?

- ✓ **Efficient Data Retrieval** – Fetches specific data quickly using queries.
- ✓ **Data Filtering & Aggregation** – Helps in **sorting, grouping, and summarizing** large datasets.
- ✓ **Data Integration** – Combines multiple data sources into a unified view.
- ✓ **Automation & Reporting** – Enables **automated reports and dashboards** for decision-making.
- ✓ **Scalability** – Handles **millions of records** efficiently in real-time.

📌 Example:

A **retail company** uses SQL queries to analyze **monthly sales trends** and identify top-selling products.

💡 Conclusion:

SQL is a **powerful tool** for **Data Analytics**, allowing businesses to make data-driven decisions efficiently.

📌 CHAPTER 3: SQL DATABASE STRUCTURE

SQL databases are structured in a **tabular format**, making data organization and retrieval systematic and efficient.

3.1 Key Components of SQL Databases

- ✓ **Tables** – Store data in rows (records) and columns (fields).
- ✓ **Columns (Fields)** – Define the type of data stored (e.g., Name, Age, Sales).
- ✓ **Rows (Records)** – Individual entries in a table representing data instances.
- ✓ **Primary Key** – Unique identifier for each record in a table.
- ✓ **Foreign Key** – Links records from one table to another for **data relationships**.

📌 Example:

| Customer_ID | Name | Age | City | Total_Spend |
|-------------|-------|-----|-----------|-------------|
| 101 | John | 28 | New York | \$500 |
| 102 | Sarah | 32 | Chicago | \$700 |
| 103 | Alex | 25 | San Diego | \$450 |

This table stores **customer data**, and SQL queries can retrieve **customer names who spent over \$500**.

💡 Conclusion:

Understanding SQL **database structure** is crucial for efficiently querying and analyzing data.

📌 CHAPTER 4: BASIC SQL COMMANDS FOR DATA ANALYTICS

SQL provides **powerful commands** to interact with databases.

4.1 Data Querying Commands

- ✓ **SELECT** – Retrieves data from a table.
- ✓ **WHERE** – Filters records based on conditions.
- ✓ **ORDER BY** – Sorts results in ascending/descending order.
- ✓ **GROUP BY** – Groups data based on a column.
- ✓ **HAVING** – Filters grouped data.
- ✓ **JOIN** – Combines data from multiple tables.

📌 Example:

Retrieve **customers who spent more than \$500**:

```
SELECT Name, Total_Spend  
FROM Customers  
WHERE Total_Spend > 500;
```

💡 Conclusion:

Mastering **basic SQL commands** allows analysts to **efficiently filter and manipulate** data.

📌 CHAPTER 5: SQL DATA MANIPULATION & AGGREGATION

SQL provides commands for **modifying and summarizing data**.

5.1 Data Manipulation Commands

- ✓ **INSERT INTO** – Adds new records to a table.
- ✓ **UPDATE** – Modifies existing records.
- ✓ **DELETE** – Removes records from a table.

📌 **Example:**

Update a customer's **total spend** in the database:

```
UPDATE Customers
```

```
SET Total_Spend = 800
```

```
WHERE Customer_ID = 101;
```

💡 **Conclusion:**

SQL modification commands allow **dynamic updates** to datasets.

5.2 Data Aggregation & Summary Functions

- ✓ **COUNT()** – Counts the number of rows in a table.
- ✓ **SUM()** – Calculates the total of numeric columns.
- ✓ **AVG()** – Finds the average value of a column.
- ✓ **MAX() & MIN()** – Returns the highest and lowest values.

📌 **Example:**

Find the **total sales** across all customers:

```
SELECT SUM(Total_Spend) AS Total_Sales FROM Customers;
```

💡 **Conclusion:**

SQL aggregation functions help **summarize large datasets efficiently**.

📌 CHAPTER 6: SQL JOINS & RELATIONAL DATA ANALYSIS

SQL Joins allow combining multiple tables to extract meaningful insights.

6.1 Types of SQL Joins

- ◆ **INNER JOIN** – Returns matching records from both tables.
- ◆ **LEFT JOIN** – Returns all records from the left table and matching records from the right.
- ◆ **RIGHT JOIN** – Returns all records from the right table and matching records from the left.
- ◆ **FULL OUTER JOIN** – Returns all records from both tables.

📌 Example:

Combine **Customers** and **Orders** tables to retrieve customer order details:

```
SELECT Customers.Name, Orders.Order_ID, Orders.Amount  
FROM Customers  
INNER JOIN Orders  
ON Customers.Customer_ID = Orders.Customer_ID;
```

💡 Conclusion:

SQL Joins **enhance Data Analytics** by combining multiple data sources for richer insights.

📌 CHAPTER 7: SQL FOR BUSINESS INTELLIGENCE & REPORTING

SQL is widely used in **Business Intelligence (BI)** to generate **automated reports and dashboards**.

7.1 Business Applications of SQL in Reporting

- ✓ **KPI Monitoring** – Tracks business performance (e.g., total sales, average revenue).
- ✓ **Trend Analysis** – Evaluates historical data to forecast trends.
- ✓ **Customer Segmentation** – Groups users based on spending behavior.
- ✓ **Revenue & Profitability Reports** – Helps financial decision-making.

📌 Example:

Generate a report showing **monthly revenue trends**:

```
SELECT MONTH(Order_Date) AS Month, SUM(Amount) AS  
Total_Revenue  
  
FROM Orders  
  
GROUP BY MONTH(Order_Date);
```

💡 Conclusion:

SQL is indispensable for business reporting, providing structured and automated insights.

CHAPTER 8: EXERCISES & ASSIGNMENTS

8.1 Multiple Choice Questions

- Which SQL command is used to retrieve data?**
- (a) UPDATE
- (b) DELETE
- (c) SELECT
- (d) INSERT

Which SQL function returns the total count of rows?

- (a) AVG()
- (b) COUNT()
- (c) SUM()
- (d) GROUP BY

Which SQL JOIN returns only matching records from two tables?

- (a) LEFT JOIN
- (b) RIGHT JOIN
- (c) INNER JOIN
- (d) FULL JOIN

8.2 Practical Assignment

- 📌 Task 1: Write an SQL query to extract customers who made purchases above \$1,000.
- 📌 Task 2: Use SQL aggregation to calculate the **total revenue by product category**.
- 📌 Task 3: Perform an **INNER JOIN** to combine sales data with customer details.

🌟 CONCLUSION: THE FUTURE OF SQL IN DATA ANALYTICS

SQL remains the **gold standard for managing structured data**, with advancements in **cloud databases, AI integration, and real-time analytics**. Mastering SQL is essential for **anyone pursuing a career in Data Analytics, Business Intelligence, or Data Science**.  

📌 **ASSIGNMENT:**

**✓ PERFORM DATA CLEANING &
PREPROCESSING ON A DATASET USING
PANDAS & SQL**

ISDM-Nxt

🔧 SOLUTION: DATA CLEANING & PREPROCESSING USING PANDAS & SQL

🎯 Objective:

The goal of this assignment is to perform **data cleaning and preprocessing** using **Pandas (Python)** and **SQL** on a real-world dataset. This step-by-step guide will cover how to:

- ✓ Load a dataset
- ✓ Identify and handle missing values
- ✓ Detect and remove duplicates
- ✓ Convert data types
- ✓ Handle outliers
- ✓ Normalize and standardize data
- ✓ Perform SQL operations for data cleaning

◆ STEP 1: Load the Dataset

Before performing data cleaning, we need to **import the dataset** into both **Python (Pandas)** and **SQL**.

1.1 Using Pandas in Python

📌 Install required libraries (if not installed):

```
pip install pandas numpy mysql-connector-python
```

📌 Import Libraries and Load the Dataset:

```
import pandas as pd
```

```
import numpy as np
```

```
# Load dataset (Example: A CSV file)  
df = pd.read_csv("customer_data.csv")
```

```
# Display first 5 rows
```

```
print(df.head())
```

1.2 Loading Data into SQL

If you're using **MySQL**, follow these steps:

- 📌 Create a database and table in SQL:

```
CREATE DATABASE DataCleaning;
```

```
USE DataCleaning;
```

```
CREATE TABLE customers (  
    customer_id INT PRIMARY KEY,  
    name VARCHAR(100),  
    email VARCHAR(255),  
    age INT,  
    city VARCHAR(100),  
    purchase_amount FLOAT,  
    signup_date DATE  
)
```

📌 **Load Data into MySQL Table (Using MySQL Workbench or Command Line):**

```
LOAD DATA INFILE 'customer_data.csv'  
INTO TABLE customers  
FIELDS TERMINATED BY ''  
ENCLOSED BY ""  
LINES TERMINATED BY '\n'  
IGNORE 1 ROWS; -- Ignores header row
```

◆ **STEP 2: Identify and Handle Missing Values**

2.1 Checking for Missing Values

📌 **In Pandas:**

```
# Check for missing values  
print(df.isnull().sum())
```

📌 **In SQL:**

```
SELECT
```

```
    SUM(CASE WHEN name IS NULL THEN 1 ELSE 0 END) AS  
missing_names,
```

```
    SUM(CASE WHEN email IS NULL THEN 1 ELSE 0 END) AS  
missing_emails,
```

```
    SUM(CASE WHEN age IS NULL THEN 1 ELSE 0 END) AS  
missing_ages
```

FROM customers;

2.2 Handling Missing Values

📌 Filling Missing Values in Pandas:

```
# Fill missing values with default values
```

```
df['age'].fillna(df['age'].mean(), inplace=True) # Fill with mean age
```

```
df['city'].fillna('Unknown', inplace=True) # Fill missing cities with  
'Unknown'
```

📌 Handling Missing Values in SQL:

```
UPDATE customers
```

```
SET age = (SELECT AVG(age) FROM customers)
```

```
WHERE age IS NULL;
```

◆ STEP 3: Remove Duplicates

Duplicate records can cause inconsistencies in data analysis.

📌 Checking for Duplicates in Pandas:

```
# Find duplicate records
```

```
print(df.duplicated().sum())
```

```
# Drop duplicate rows
```

```
df.drop_duplicates(inplace=True)
```

📌 Checking for Duplicates in SQL:

```
SELECT customer_id, COUNT(*)  
FROM customers  
GROUP BY customer_id  
HAVING COUNT(*) > 1;
```

📌 **Remove Duplicates in SQL:**

```
DELETE FROM customers  
WHERE customer_id IN (  
    SELECT customer_id FROM (  
        SELECT customer_id, ROW_NUMBER() OVER (PARTITION BY  
customer_id ORDER BY customer_id) AS row_num  
        FROM customers  
) AS temp WHERE row_num > 1  
);
```

◆ **STEP 4: Convert Data Types**

Incorrect data types can lead to calculation errors.

📌 **Convert Data Types in Pandas:**

```
# Convert 'signup_date' to datetime  
df['signup_date'] = pd.to_datetime(df['signup_date'])  
  
# Convert 'age' to integer
```

```
df['age'] = df['age'].astype(int)
```

📌 Convert Data Types in SQL:

```
ALTER TABLE customers
```

```
MODIFY COLUMN signup_date DATE,
```

```
MODIFY COLUMN age INT;
```

◆ STEP 5: Handling Outliers

Outliers are extreme values that can distort statistical analysis.

📌 Detect Outliers in Pandas (Using IQR Method):

```
Q1 = df['purchase_amount'].quantile(0.25)
```

```
Q3 = df['purchase_amount'].quantile(0.75)
```

```
IQR = Q3 - Q1
```

```
# Define lower and upper bounds
```

```
lower_bound = Q1 - 1.5 * IQR
```

```
upper_bound = Q3 + 1.5 * IQR
```

```
# Remove outliers
```

```
df = df[(df['purchase_amount'] >= lower_bound) &  
(df['purchase_amount'] <= upper_bound)]
```

📌 Detect and Remove Outliers in SQL:

```
WITH purchase_stats AS (
    SELECT
        customer_id, purchase_amount,
        PERCENTILE_CONT(0.25) WITHIN GROUP (ORDER BY
        purchase_amount) AS Q1,
        PERCENTILE_CONT(0.75) WITHIN GROUP (ORDER BY
        purchase_amount) AS Q3
    FROM customers
)
DELETE FROM customers
WHERE purchase_amount < (Q1 - 1.5 * (Q3 - Q1))
    OR purchase_amount > (Q3 + 1.5 * (Q3 - Q1));
```

◆ STEP 6: Normalize & Standardize Data

6.1 Normalize Data (Min-Max Scaling)

📌 In Pandas:

```
from sklearn.preprocessing import MinMaxScaler
```

```
scaler = MinMaxScaler()
```

```
df[['purchase_amount']] =
scaler.fit_transform(df[['purchase_amount']])
```

📌 In SQL:

UPDATE customers

SET purchase_amount =

(purchase_amount - (SELECT MIN(purchase_amount) FROM customers)) /

((SELECT MAX(purchase_amount) FROM customers) - (SELECT MIN(purchase_amount) FROM customers));

6.2 Standardize Data (Z-Score Scaling)

📌 In Pandas:

```
from scipy.stats import zscore
```

```
df['purchase_amount'] = zscore(df['purchase_amount'])
```

📌 In SQL:

UPDATE customers

SET purchase_amount =

(purchase_amount - (SELECT AVG(purchase_amount) FROM customers)) /

(SELECT STDDEV(purchase_amount) FROM customers);

◆ STEP 7: Save the Cleaned Data

After cleaning, the dataset should be **saved** for future use.

📌 Save Cleaned Data in Pandas:

```
df.to_csv("cleaned_customer_data.csv", index=False)
```

📌 Save Cleaned Data in SQL:

```
CREATE TABLE cleaned_customers AS
```

```
SELECT * FROM customers;
```

Summary of Steps Performed

1. **Loaded dataset** using Pandas and SQL.
2. **Checked for missing values** and filled them appropriately.
3. **Removed duplicate entries** for data consistency.
4. **Converted incorrect data types** (e.g., string to date, float to integer).
5. **Handled outliers** using the IQR method.
6. **Normalized & standardized** numerical data for better model performance.
7. **Saved the cleaned dataset** for future analysis.

Next Steps:

- ◆ Use the cleaned data for **Exploratory Data Analysis (EDA)** and visualization.
- ◆ Apply **Machine Learning models** to make predictions.
- ◆ Create **interactive dashboards** using **Power BI or Tableau**.