

# Data Analytics

AD 23411

DA → Data Analytics  
DAO → Data Analyst

## Unit-1 (Notes)

- \* It's all about how to turn the raw numbers to understandable information.

### Topic to cover

- \* Overview of DA (Data Analytics)
- \* Types of DA (Data Analysis)
- \* Steps in DA's Process
- \* Data Repositories
- \* ETL (Extract, Transform, Load)
- \* Roles, Responsibilities, Skillsets of DA

### Data Analytics

Data Analysis

### Overview of DA

- \* It is the process of collecting, organizing and studying data to find or create useful information

### Importance:

- \* Helps in decision making
- \* Helps in problem solving
- \* Helps identify opportunities
- \* Improved efficiency

### DA Process

- \* Data Collection (Collecting raw data)
- \* Data Cleaning (Preparing the data)
- \* Data Analysis : Using programming language to perform tasks
- \* Data Visualization (Results are visualized and analyzed)

### Types of DA

- 1) Descriptive DA : Analytics of past data to understand them
- 2) Diagnostic DA : It finds the reason for the past data, uses correlation and plot graph (focuses on Why?)
- 3) Predictive DA : Predicts future based on past

at **Prescriptive DA**: not only predicting but finding the best solution to solve the problem.

### Methods of DA:

#### \*1 Qualitative Analytics

Source: words, Images, symbols, interviews  
Methods: Content Analytics Analysis, Grounded Theory, narrative analysis

Goal: deriving meaning and context from non-statistical inputs

#### \*1 Quantitative Analytics

Used **Numerical Data**

Methods: Hypothesis testing, Regression, Mean Average Calculations

Goal: Measuring variables and testing theories statistically

### Example of Types of Data Analytics

#### 1) Descriptive Data Analytics

code: (Pandas, Matplotlib)

```
data = {'Product': ['A', 'B'],
        'January': [200, 180],
        'February': [250, 210],
        'March': [300, 240]}
```

df = pd.DataFrame(data)

df.set\_index('Products').plot(kind='bar', color=['red', 'blue'])

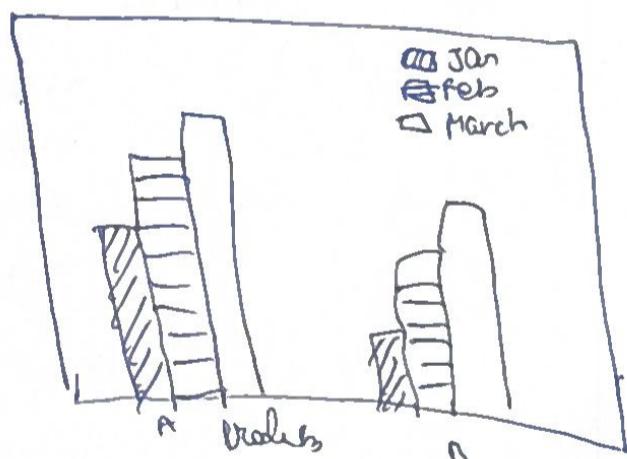
plt.title("Monthly Sales Data")

plt.xlabel("Products")

plt.ylabel("Sales")

plt.show()

Output:



## 2) Diagnostic Analytics:

Code: pandas as pd, seaborn as sns, matplotlib.pyplot as plt

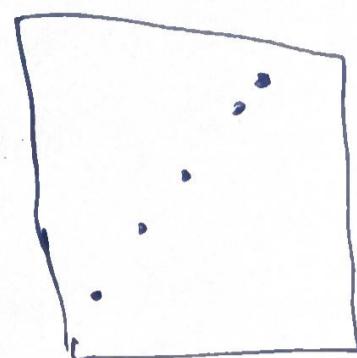
data = {'Ad\_Spend': [1000, 1500, 2000, 2500, 3000],  
 'Sales': [12000, 18000, 18000, 21000, 28000]}

df = pd.DataFrame(data)  
cor = df.corr()

print(cor)

sns.scatterplot(x='Ad\_Spend', y='Sales', data=df, color='green')  
plt.title

Output:



Ad Spend  
Show +ve correlation

## 3) Predictive Analytics

pandas, matplotlib, seaborn, from sklearn import linear\_model import

Linear Regression

df = pd.DataFrame({'Month': [1, 2, 3, 4, 5], 'Sales': [2000, 2500, 3000, 3500, 4000]})  
X, Y = df[['Month']], df['Sales']

model = LinearRegression().fit(X, Y)

nextMonth = [6]

prediction = model.predict([nextMonth])[0]

sns.set\_style("darkgrid")

plt.scatter(X, Y, color='blue', label='Actual')

plt.plot(range(1, 7), model.predict([[i]] for i in range(1, 7)), 'g--',  
 label='Trend')

plt.scatter(nextMonth, prediction, color='red', label='Forecast')

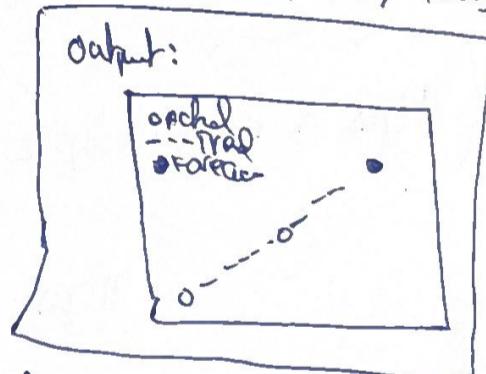
plt.title('Month 6 Forecast: \$ {Prediction:,.2f}'), \$=100, (label = Forecast))

plt.show()

Next ('Prediction:,.2f'))

Output:

Actual  
Trend  
Forecast



## 1) Predictive DA

Code: pandas, matplotlib.pyplot as plt, seaborn &, sklearn.linear\_model

```
df = pd.DataFrame({'sqft' : [1000, 1500, 2000, 2500, 3000], 'Price' : [20000, 25000, 30000, 350000, 400000]})
```

```
, df['sqft']] = [1800]
```

```
model = LinearRegression().fit(x,y)
```

```
target_sqft = [1800]
```

```
prediction = model.predict(target_sqft)[0]
```

```
sns.set_theme(style = "whitegrid")
```

```
plt.scatter(x,y, ...)
```

```
plt.plot(x, model.predict(x), color = 'green', label = 'Pred')
```

```
plt.scatter(..., ...)
```

```
:
```

## Output:

Same as Predictive code, we would just draw some conclusion, that's the difference.

## Steps in Data Analysis Process:

### Steps :

- 1) Define the Problem (understand before solving)
- 2) Data Collection (Extracting/Collecting the necessary data for the process)

Code:

```
titanic = sns.load_dataset('titanic')
titanic.head()
```

### Shortcuts

Sns → Seaborn  
Pd → Pandas  
Plt → Matplotlib

Output: Dataset is displayed

④

③

3) **Data Cleaning**: It is the longest process, here we prepare the data for the further process i.e.: Removing null, duplicates, fixing inaccuracy  
Code: for filling

- \*1 titanic['embarked'].fillna(titanic['Embarked'].mode()[0], inplace=True)
- \*1 df['filled['alone']] = df['filled['alone']].fillna(df['alone'].mean())

4) **Analyzing Data**: Here you will visualize and find the patterns, trends behind the given data, also if there is any inaccurate value it is analyzed in this step

5) **Visualizing Data**:

- \*1 converting to graph or any chart to easily understand the data

6) **Interpret And Make Decisions**: This is the final step

\*1 Explain the result and more

\*1 This is 100% dependent on the user

Ex code: for accuracy

$$\text{accuracy} = \text{accuracy\_score}(y_{\text{val}}, y_{\text{pred}})$$

similar

### Data Repositories in DA

1) **Data warehouse**: It is a centralized data repository, just like real life warehouse, they are structured historical data. schema-on-write  
key features:

o) **Subject-oriented**: organized around major business subjects, rather than applications or processes.

o) **Integrated**: Data from various sources are integrated

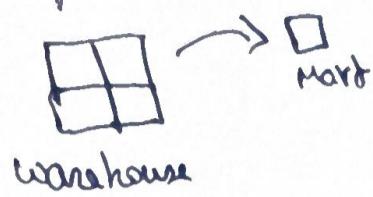
o) **Time-variant**: Trend analysis and long term reporting

o) **Non-volatile**: Read only (not modified)

o) **Optimized for querying and analysis**

## Architecture :

- 1) Source Systems : External Data Sources (CRM, ERP)
- 2) ETL Process (Extract, Transform, Load)
- 3) Data warehouse
- 4) Data marts
- 5) BI Tools

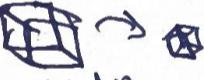


## Advantages:

- 1) Centralized Data access
- 2) Improved Data Quality
- 3) Historical Analysis
- 4) Efficient query Performance
- 5) Support Decision-making

2. Data cube: Multi dimensional Data structure used in data warehousing and online analytical processing (OLAP)

## Operations:

- \* Slice : 3D  $\rightarrow$  2D view of particular field
- \* Cube : Smaller Sub Cube 
- \* Roll up : Aggregates data by climbing up a dimension hierarchy
- \* Drill down : Breaks down aggregated data into finer levels
- \* Pivot : Rotate the cube to view different perspective

## 3. Data Marts:

- \* Small subdivision of Data warehouse

## Key Features:

- o Subject oriented
- o Smaller Scope
- o Optimized for specific use
- o Derived from a Data warehouse
- o Improved performance

## Types:

- o Dependent DM
- o Independent DM
- o Hybrid DM

## Advantages:

- o Faster
- o cost effective
- o de-centralized control
- o Improved performance

4. **Data Lakes**: It's like a real life lake, unstructured, elements of different types are there, raw, unstructured, flexible, schema-on-read model  
**Key features**:

- \* 1 Raw Data storage
- \* 1 Scalable
- \* 1 Flexible Schema
- \* 1 Cost effective
- \* 1 Supports diverse datatypes
- \* 1 Accessible for Big Data Analytics

### **Advantages**

- 1) Flexible
- 2) Scalability
- 3) Cost Saving
- 4) Supports Advanced Analytics

### **Challenges :**

- 1) Extra Governance
- 2) Performance
- 3) Security

### **Advantages of Data Repos**:

- \* 1 Centralized
- \* 1 Data Accessibility
- \* 1 Cost efficiency
- \* 1 Enhanced Security
- \* 1 Supports Advanced Analytics

### **ETL (Extract, Transform, Load)**

#### **1. Extract: Getting the Data**

#### **Data sources**

- 1) Relational databases, APIs, 5) cloud storage
- 2) flat file
- 3) NoSQL databases
- 4) Real-time data streams

## Steps in Extractions:

- 1) Identify the data sources and their formats
- 2) Establish a connection to the sources
- 3) Extract the data efficiently, ensuring minimal disruption to the source system

## challenges:

- \* Handling Large Data volumes
- \* Managing inconsistent formats
- \* Minimizing Impact on source systems

## 2. Transform

- 1) Data cleaning
- 2) Data validation
- 3) Data Aggregation (grouping)
- 4) Data Maping (choose where the data can be stored )
- 5) Data enrichment (improving its accuracy)
- 6) Data formating (making it in proper format)
- 7) Data Hierarchy creation (creates hierarchies for drill down analysis)

## 3. Load:

### Types:

- \* full load (Transforms all data)
- \* Incremental load (Transforms new or updated data)
- \* Real-Time load (live)

## Advantages of ETL

- \* Centralized data integration
- \* Improved data quality
- \* Historical Analysis
- \* Enhanced decision Making
- \* Scalability

## ETL Tools:

- \* Informatic power centre
- \* Talend
- \* Acos glue
- \* Apache Nifi , (data flow automation)
- \* Google data cloud

## Brülage:

### DA Problems:

- \* 1) Define the problem
- \* 1) Data collection
- \* 1) Data cleaning and preprocessing
- \* 1) Exploiting Data Analytics (Visualizing)
- \* 1) Data transformation
- \* 1) Data Modeling
- \* 1) Evaluation of results
- \* 1) Visualization and reporting
- \* 1) Decision Making and action

### Roles and responsibilities of DA:

→ roles " -> responsibilities

- \* 1) Collection and "Integration of data"
- \* 1) "Data cleaning" and preprocessing "and transformation"
- \* 1) "EDA"
- \* 1) "Data Analysis and modeling"
- \* 1) "Reporting and communication", Visualization
- \* 1) "Visualization"
- \* 1) "Collaboration"
- \* 1) "Maintain data quality"

### Skills required:

#### Tools and Tech:

- \* 1) Data Base
- \* 1) Data cleaning
- Libraries,
- \* 1) Visualization tool
- \* 1) Big data tool

#### Technical skill:

- \* 1) Data wrangling
- \* 1) EXCEL
- \* 1) Data visualization
- \* 1) Statistical analytical tools SAS, SPSS  

↓

Statistical  
package of  
Social science

#### Machine learning

- \* 1) Problem solving
- \* 1) Communication
- \* 1) Attention to detail
- \* 1) Collaboration

#### Business skill:

- \* 1) Domain
- \* 1) Critical thinking