# Exploratory Data Analysis (EDA) & Data Visualization

# Introduction to EDA

- Definition of EDA:-Understanding the data before making any decisions

- Importance in Data Science:-Detecting Missing value,Finds Outliers,Improves Data Quality for Better Models,Saves Time and Prevents Costly Mistakes

# cont…..

- Goals of EDA: Understanding Data, Detecting Anomalies

# Volume of Data

- Small Data
  Ex:- excel sheet,CSV,JSON,TEXT(Kaggle)
- Medium Data
  Ex:- MySQL, PostgreSQL, or NoSQL(HDFC,SBI,BOI,UNION BANK)
- Big Data
- Ex. Hadoop, Apache Spark, and AWS

# variables required for the analysis(Types of Variable)

- **Dependent Variable (Target Variable)**

   Def:- Where the output is available is called dependent Variable

   Ex.House Price Prediction

# cont…………

- **Independent Variables (Predictor Variables)**
- Def:- Where the Input is available is called dependent Variable

    Ex.**Marketing Spend , Ad Clicks**

# open and paid data sources

- Free Data Sources
    Ex. Kaggle.com, NASA Open Data
- Paid Data Sources(Premium)
    Ex. IBM MarketScan for Health Industry data

# Describe the metadata

- **Def**:-Data About data is called as metadata
- **Data Identification**

    :- **what the dataset contains**.

- **Data Quality**

    :-Missing values, data format**. Data Processing & Automation**

- **Data Processing & Automation**

    **:- ETL pipelines** for understand data

# Data Validation: Tools and Processes

- Data validation:- the process of ensuring that data is **accurate, complete, and consistent** before it is used for analysis
- Tools:- EDA Technique

# Role of EDA in Data Science

- Data Understanding
  Ex. EDA

- Feature Engineering

  Ex. PCA

- Model Preparation

  Ex. ML Algorithm

# Types of Data

- Structured vs. Unstructured Data
  Ex.
   Structured:-Tabular data
  Unstructured:- Image,Video,Audio

# cont………..

- Numerical vs. Categorical Data

- Numerical:-

  :-**Continuous:** Height (170.5 cm), Temperature (36.8°C)

  :-**Discrete:** Number of employees (50), Number of orders (120)

# Categorical Data

- **Nominal:** Eye color (Blue, Brown, Green), Car brand (Toyota, Honda, BMW)

- **Ordinal:** Shirt size (S, M, L, XL), Exam Grades (A, B, C, D, F)

# Data Cleaning in EDA

- Handling Missing Values

- Removing Duplicates

- Dealing with Outliers

# Handling Missing Data

- Mean, Median, Mode Imputation
- Dropping Missing Values

# Outlier Detection Techniques

- Boxplot Method

- Z score Method

- IQR Method

# Summary Statistics

- Mean, Median, Mode

- Variance

# Measuring Data Distribution

- Skewness & Kurtosis

- Normal Distribution

# Introduction to Data Visualization

- Importance of Data Visualization
- Storytelling with Data

# Types of Data Visualization

- Univariate, Bivariate, Multivariate Analysis

# Univariate Analysis

- Histograms

- Box plots
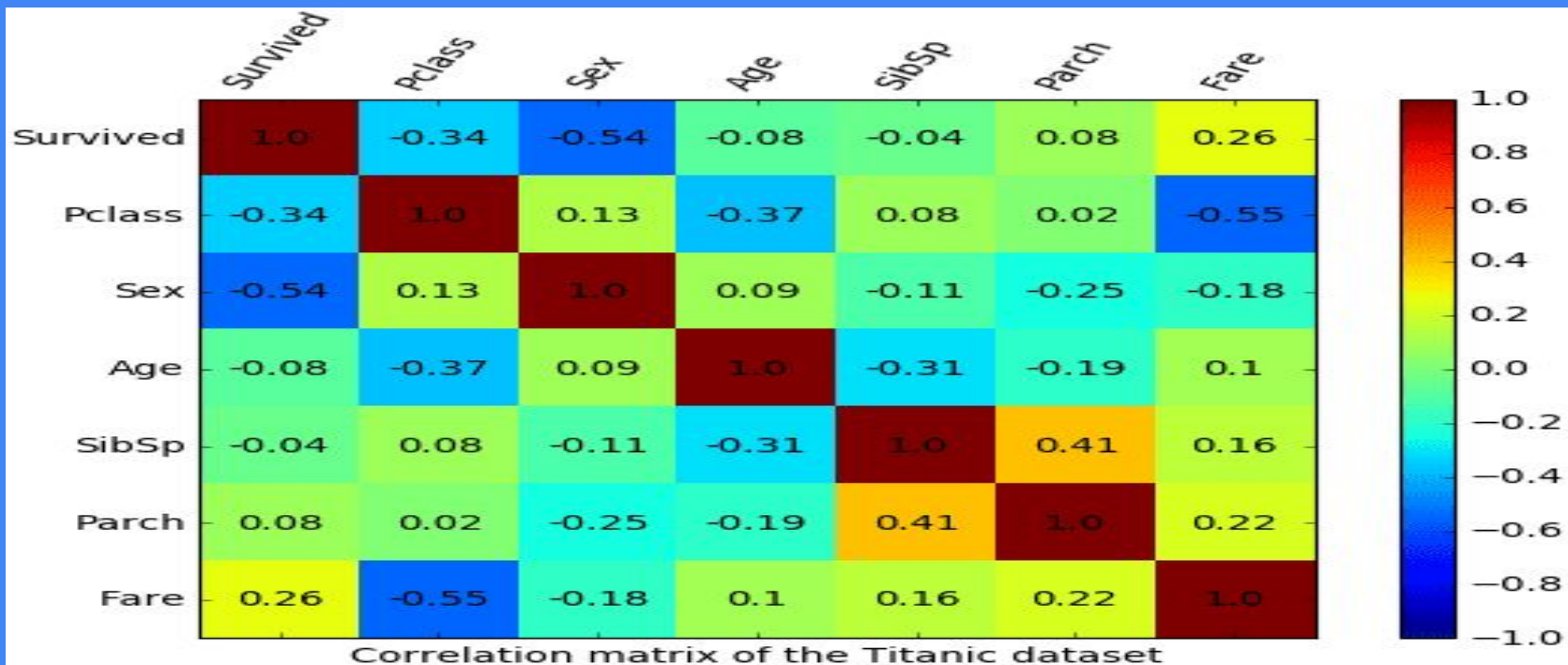
- KDE Plots

# Bivariate Analysis

- Scatter Plots

# cont……

- Line Graphs

## 6-Month sales report and forecast

IN THOUSANDS (USD)

# cont….

- Correlation Matrix



Correlation matrix of the Titanic dataset

# Multivariate Analysis

- Pair Plots

# cont………

- Heatmaps



Correlation matrix of the Titanic dataset
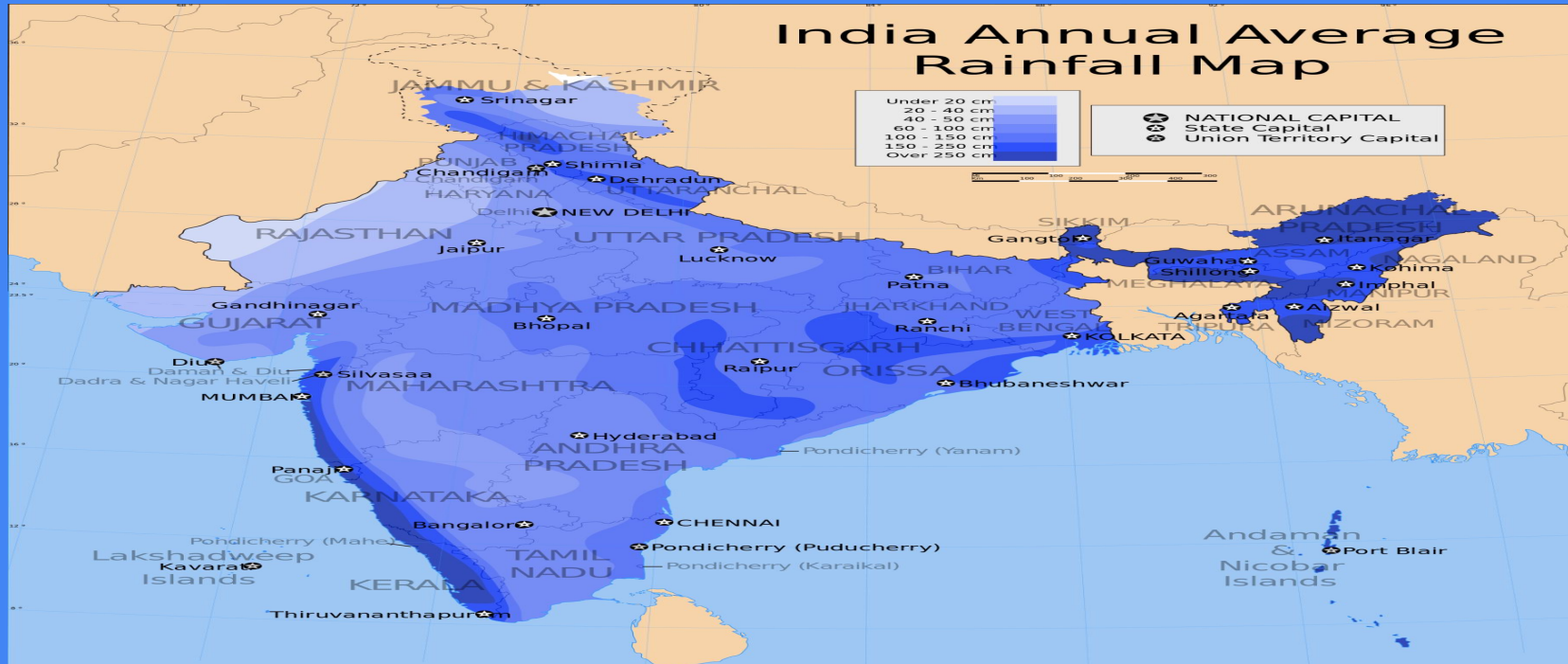
# cont..........

- 3D Scatter Plots

# Advanced Visualization Techniques

- Interactive Visualizations
    - Ex. Power Bi S/w

# cont…..
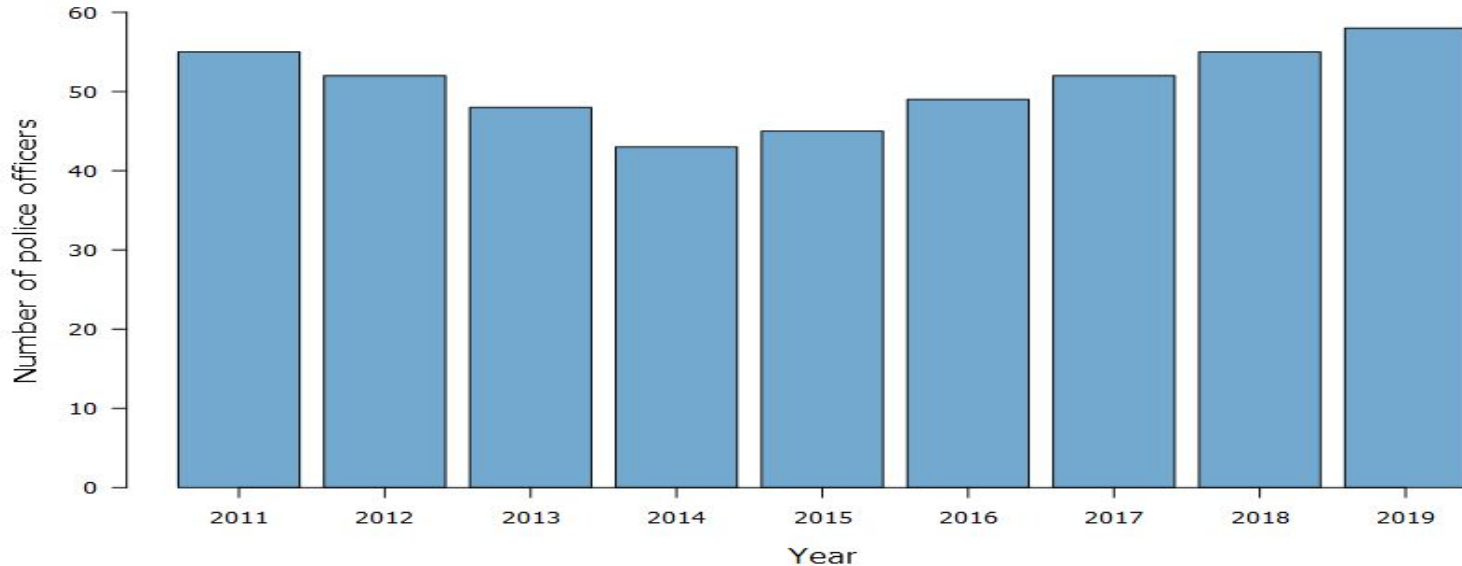
- Geospatial Data Visualization

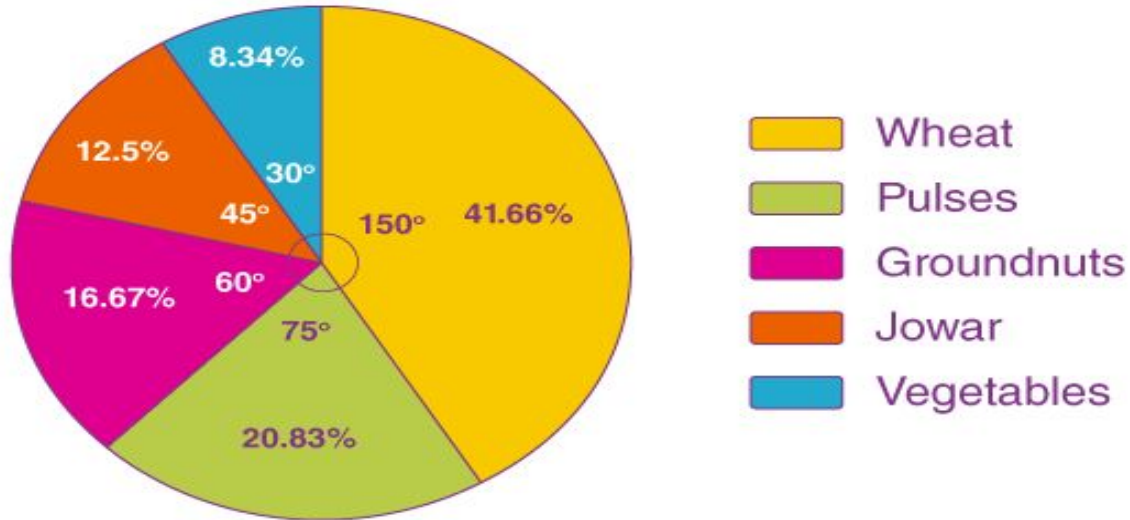# Visualizing Categorical Data

- Bar Charts



**Chart 5.2.1**
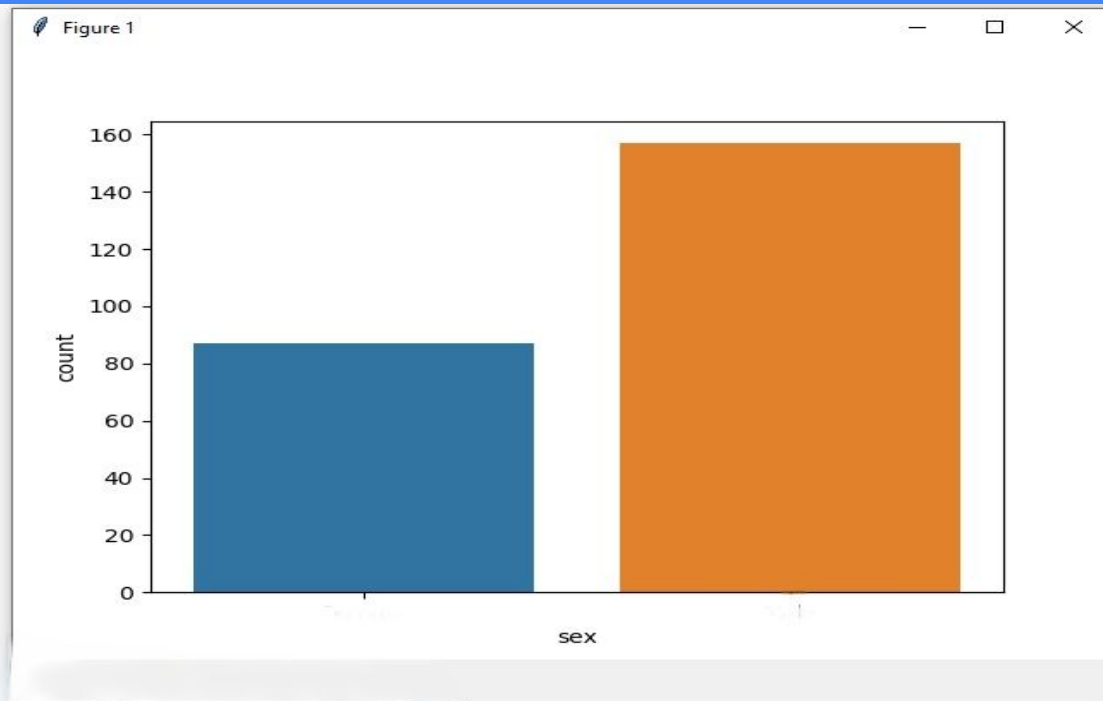**Number of police officers in Crimeville, 2011 to 2019**
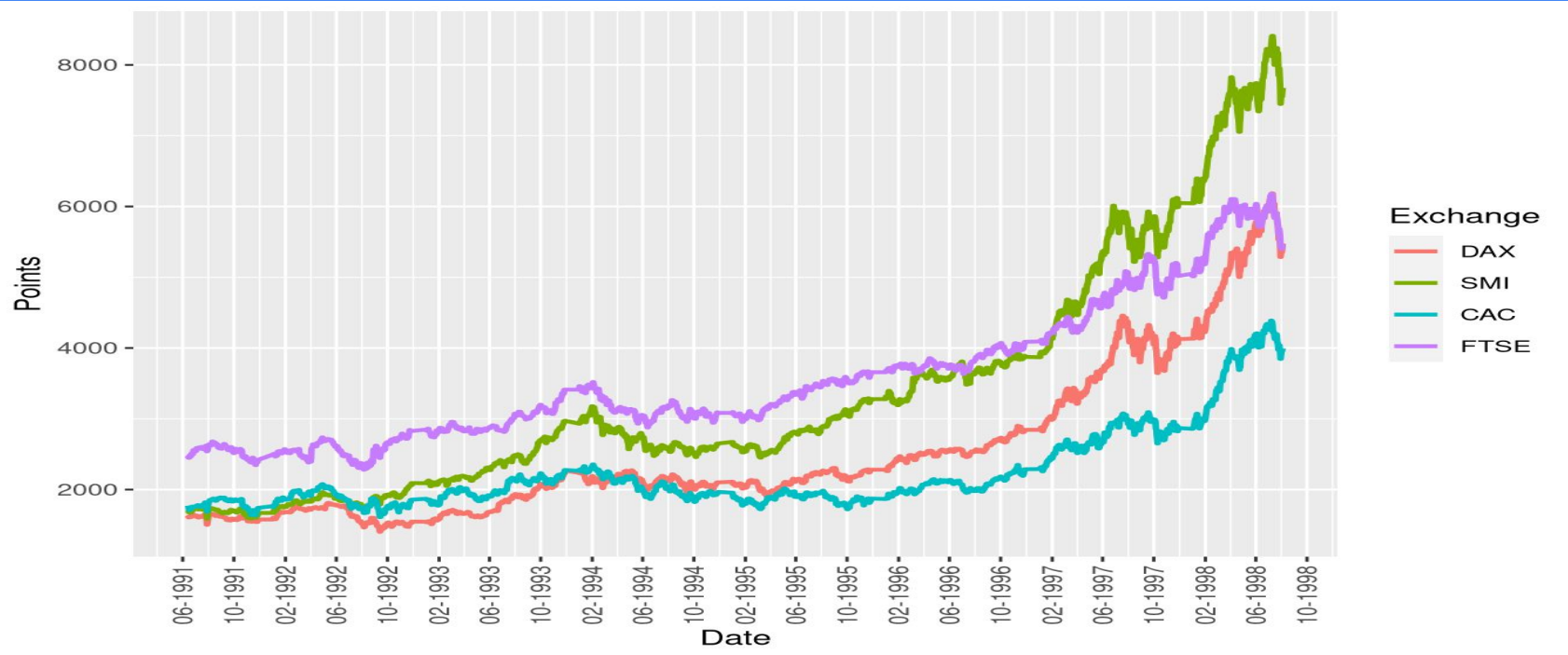
# cont..........
- Pie Charts

# cont..

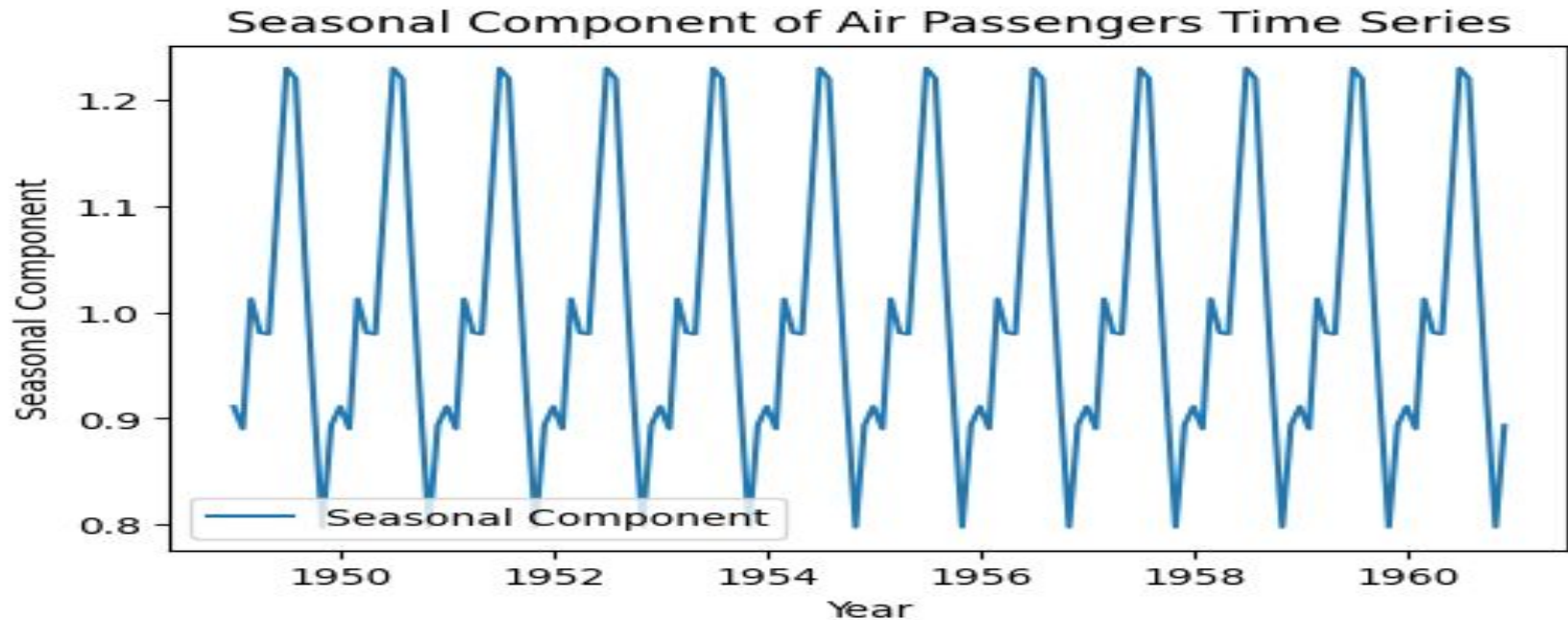- Count Plots

# Time Series Data Visualization

- Line Plots

# Rolling Averages

# cont….

- Seasonality Trends



Seasonal Component of Air Passengers Time Series

# Data Visualization Libraries in Python

- Pandas Visualization
- Matplotlib
- Seaborn

# Advanced Libraries

- Plotly
- Tableau

# Creating Interactive Dashboards

- Using Streamlit & Plotly

# Case Study: EDA on a Real Dataset

- Data Cleaning
- Feature Engineering
- Visualization Insights

# Common Challenges in EDA

- Handling Large Datasets

- Dealing with Noisy Data

- Choosing the Right Visuals

# Best Practices for EDA & Visualization

- Use Meaningful Visuals

- Avoid Misleading Graphs

- Keep It Simple & Clear

# Q&A Session

- Open Floor for Questions