



<https://pianalytix.com/>

# DATA SCIENCE MASTERY CHEAT SHEET

: <https://instagram.com/pianalytix>

: <https://www.linkedin.com/company/pianalytix/>

(140+ Code Snippets to 10x your Productivity)

## 1 Introduction

### Components

1. Data Collection and Preparation:
  - Gathering relevant data from various sources.
  - Cleaning, preprocessing, and transforming the data for analysis.
2. Statistical Analysis and Modeling:
  - Applying statistical techniques to identify patterns and relationships in the data.
  - Creating mathematical models and algorithms for predictive and prescriptive analysis.
3. Machine Learning and Data Mining:
  - Utilizing machine learning algorithms to automatically extract insights from data.
  - Uncovering hidden patterns and making predictions or classifications.
4. Data Visualization and Communication:
  - Presenting data findings through visualizations, charts, and graphs.
  - Communicating complex concepts and insights to stakeholders effectively.
5. Business Impact and Decision-making:
  - Using data-driven insights to inform strategic decisions and solve problems.
  - Maximizing business value and driving innovation through data-driven approaches.

### Libraries

1. NumPy:
  - Numerical computation and array manipulation for scientific computing.
2. pandas:
  - Data manipulation and analysis for scientific data processing.
3. scikit-learn:
  - Machine learning algorithms and tools for scientific modeling.
4. Matplotlib:
  - Data visualization library for scientific plotting and graphing.
5. TensorFlow and PyTorch:
  - Deep learning frameworks for scientific research and applications.
6. SciPy:
  - Scientific computing library providing a wide range of mathematical functions.
7. Seaborn:
  - Statistical data visualization library for exploratory data analysis.

## 2 Data Collection

### 2.1 Definitions - Web Scraping

#### Web Scraping

- Web Scraping refers to the extraction of data from a website. This information is collected and then exported into a format that is more useful for the user. Be it a spreadsheet or an API.
- Web Scraping Techniques:**
1. Fetching Webpage Content:
    - Use the requests library to send a GET request and retrieve the HTML content of the webpage.
  2. HTML Parsing:
    - Utilize a parsing library like BeautifulSoup or lxml to parse the HTML structure and extract data.
  3. Locating Relevant Data:
    - Inspect the HTML structure and identify appropriate HTML tags and attributes to locate the desired data.
  4. Data Extraction:
    - Use methods such as find, find\_all, or CSS selectors to extract the relevant data from the HTML.
  5. Data Processing and Cleaning:
    - Process and clean the extracted data by removing unwanted characters, formatting, or converting data types.
  6. Handling Pagination:
    - Implement techniques to handle webpages with multiple pages, such as iterating through page numbers or following next page links.
  7. Handling Dynamic Content:
    - Address dynamic content loading by using headless browsers or JavaScript rendering libraries like Selenium or Scrapy-Splash.
  8. Error Handling and Retrying:
    - Implement error handling mechanisms to handle exceptions, network errors, and connection timeouts. Include retry logic if necessary.
  9. Respecting Website Policies:
    - Follow the website's terms of service, including rate limits, crawling restrictions, and respecting robots.txt guidelines.
  10. Data Storage:
    - Store the scraped data in a structured format like CSV, JSON, or a database for further analysis integration with other systems.

### 2.1.1 Code Snippets - Sequence Padding

#### Sequence Padding

```
import requests
from bs4 import BeautifulSoup

# Send a GET request to the webpage
response = requests.get(url)

# Parse the HTML content using BeautifulSoup
soup = BeautifulSoup(response.content, 'html.parser')

# Extract relevant data from the webpage
data = soup.findAll('img')

# Optimization
# Define the URL to scrape
url = "https://example.com"

# Set headers to mimic a web browser
headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/91.0.4472.124 Safari/537.36'}

# Send a GET request to the webpage
response = requests.get(url, headers=headers)

# Check if the request was successful
if response.status_code == 200:
    # Parse the HTML content using BeautifulSoup
    soup = BeautifulSoup(response.content, 'html.parser')

    # Extract relevant data using specific CSS selectors
    data = soup.select('.class_name')

    # Process and return the extracted data
    for item in data:
        print(item.text)
else:
    print('Error!', response.status_code)

# Add a delay between requests to avoid overwhelming the server
time.sleep(1)
```

### Strategies

#### 1 Pagination:

```
import requests

base_url = "http://example.com/page"
```

```
for page in range(1, 6):
    url = base_url + str(page)
    response = requests.get(url)
    # Process the page data
```

#### 2 Handling Dynamic Content:

```
from selenium import webdriver
```

```
driver = webdriver.Chrome('path_to_chromedriver')
driver.get('https://example.com')
```

```
# Wait for dynamic content to load
time.sleep(3)
```

```
# Extract data from the loaded page
# Close the driver
driver.quit()
```

#### 3 Login and Session Handling:

```
import requests
```

```
login_url = "https://example.com/login"
data = {'username': 'your_username', 'password': 'your_password'}
```

```
# Perform login
session = requests.Session()
session.post(login_url, data=data)
```

```
# Use the session for subsequent requests
response = session.get('https://example.com/dashboard')
```

```
# Process the response
```

#### 4 Captcha Handling:

```
import requests
```

```
# Fetch the captcha image
response = requests.get('https://example.com/captcha')
```

```
# Save the captcha image
with open('captcha.png', 'wb') as f:
    f.write(response.content)
```

```
# Manually solve the captcha and get the solution
```

```
# Send a POST request with the captcha solution
data = {'captcha': 'your_solution'}
response = requests.post('https://example.com/solve', data=data)
```

```
# Process the response
```

## 2.2 Definitions - API Data Retrieval

### API Data Retrieval

An API (Application Programming Interface) is a standardized and secure interface that allows applications to communicate and work with each other.

The process of API data retrieval involves several steps:

1. Identification and Authentication:
  - Identify API and obtain access credentials (API key, tokens, etc.).
2. Making the Request:
  - Formulate HTTP request to API endpoint with required parameters.
3. API Server Processing:
  - API server processes the request and validates permissions.
4. Receiving the Response:
  - API sends structured response (JSON, XML) back to the client.
5. Data Extraction:
  - Extract and parse relevant information from the API response.
6. Error Handling:
  - Gracefully handle API errors with informative messages and codes.
7. Data Processing:
  - Manipulate and transform retrieved data for application requirements.
8. Utilization of Data:
  - Use data in the application's logic, display, or other functionalities.

### 2.2.1 Code Snippets - API Data Retrieval

#### API Data Retrieval

```
import requests

# Make a GET request to the API endpoint
response = requests.get(api_url, params=params)

# Get the JSON data from the response
data = response.json()

# Extract relevant information from the JSON data
extracted_data = data['key']
```

### Optimization

```
import requests
```

```
def retrieve_data_from_api(api_url, params=None, headers=None):
    try:
        # Use connection pooling to improve performance
        session = requests.Session()
    except:
        # Set session headers if needed
        if headers:
            session.headers.update(headers)
```

```
# Make the API request with appropriate parameters
response = session.get(api_url, params=params)
```

```
# Check for successful response (status code 200)
if response.status_code == 200:
    # Process the response data
    data = response.json()
    # Perform additional data processing as needed ...
    return data
else:
    # Handle API errors with appropriate error handling
    response.raise_for_status()
```

```
except requests.exceptions.RequestException as e:
    # Handle connection errors and exceptions gracefully
    print(f"Error occurred: {e}")
    return None
```

```
# Example Usage:
api_url = "https://api.example.com/data"
parameters = {"param1": "value1", "param2": "value2"}
headers = {"Authorization": "Bearer YOUR_TOKEN_HERE"}

results = retrieve_data_from_api(api_url, params=parameters, headers=headers)

# Process and utilize the retrieved data
print(results)
```

### Strategies

#### 1 Caching using requests\_cache:

```
import requests
import requests_cache
```

```
requests_cache.install_cache("api_cache", expire_after=3000)
```

```
data = requests.get("https://api.example.com/data").json()
```

#### 2 Asynchronous Requests using asyncio and aiohttp:

```
import asyncio
import aiohttp
```

```
async def get_data():
    sync_with_aiohttp.ClientSession() as session:
        sync_with_session.get("https://api.example.com/data") as response:
            await response.read()}
```

```
loop = asyncio.get_event_loop()
task = loop.run_until_complete(get_data())
print(task.result())
```

#### 3 Rate Limiting and Retry using ratelimit library:

```
from ratelimit import ratelimit
```

```
@ratelimit(max_per_second=1000, wait_exponential_multiplier=1000, wait_exponential_max=10000, stop_max_attempt_number=3)
def get_data():
    return requests.get("https://api.example.com/data").json()
```

```
data = get_data()
print(data)
```

#### 4 Handling Errors and Retry with retrying:

```
from retrying import retry
```

```
import requests
```

```
@retry(wait_exponential_multiplier=1000, wait_exponential_max=10000, stop_max_attempt_number=3)
def get_data():
    return requests.get("https://api.example.com/data").json()
```

```
data = get_data()
print(data)
```

### Strategies

#### 4 Pagination Handling for Large Data:

```
import requests
```

```
data = []
url = "https://api.example.com/data"
```

```
while url:
    response = requests.get(url)
    data.extend(response.json())
    url = response.links['next'].url
```

```
print(data)
```

#### 5 Utilizing Webhooks for Real-Time Updates:

```
from flask import Flask, request, jsonify
```

```
app = Flask(__name__)
```

```
app.route('/webhook', methods=['POST'])
def webhook_listener():
    data = request.json
    # Process and store the received data in the application
    # ...
    return jsonify({'message': 'Webhook data received successfully'}), 200
```

```
if __name__ == '__main__':
    app.run(port=5000)
```

## 2.3 Definitions - Database Query

### Database Query

A database query for data collection is a command or a set of commands that are used to retrieve specific information or data from a database. In other words, it is a request made to a database management system (DBMS) to fetch records or data that match certain criteria.

A database query allows you to interact with the database and obtain the exact data you need, such as:

**Data Retrieval:** Fetching data from one or multiple database tables based on specified conditions.

**Data Filtering:** Selecting specific rows or columns that meet certain criteria.

**Data Aggregation:** Calculating summary information, like sums or averages, over groups of data.

**Data Modification:** Inserting, updating, or deleting records within the database.

**Data Joins:** Combining data from multiple tables based on shared columns.

### 2.3.1 Code Snippets - Database Query

#### Database Query

```
import sqllite3
```

```
# Connect to the database
conn = sqllite3.connect('database.db')
cursor = conn.cursor()
```

```
# Execute SQL query to fetch data
cursor.execute('SELECT * FROM table')
```

```
# Fetch all rows of data
data = cursor.fetchall()
```

```
# Close the database connection
conn.close()
```

### Optimization

#### 1 Use Indexing:

```
# Create an index on the "age" column to speed up queries
CREATE INDEX idx_age ON Customers (age);
```

```
# Example query using the index
SELECT name FROM Customers WHERE age > 20;
```

#### 2 Limit the Number of Rows:

```
# Fetch only the top 10 rows from the "Customers" table
SELECT * FROM Customers LIMIT 10;
```

#### 3 Avoid SELECTing unnecessary Columns:

```
# Select only the necessary columns from the "Customers" table
SELECT name, email FROM Customers WHERE age > 30;
```

#### 4 Use Proper Joins:

```
# Avoid using SELECT *; for the required columns explicitly
SELECT name, email FROM Customers;
```

#### 5 Use Subqueries and Use JOINS:

```
# Avoid subquery for better performance
SELECT name, email
FROM Customers
WHERE customer_id IN (SELECT customer_id FROM Orders WHERE order_amount > 100);
```

```
# Use JOIN for the same result
SELECT DISTINCT name, email
FROM Customers
JOIN orders ON Customers.customer_id = Orders.customer_id
WHERE order_amount > 100;
```

#### 6 Use Aggregate Functions Wisely:

```
# Avoid unnecessary aggregation and filtering
SELECT MAX(salary) FROM Employees WHERE department = 'IT';
```

```
# Prefer filtering before aggregation
SELECT MAX(salary) FROM Employees WHERE department = 'IT';
```

#### 8 Analyze Query Execution Plan:

```
# Explain query to analyze execution plan
EXPLAIN SELECT * FROM Customers WHERE order_date BETWEEN '2023-01-01' AND '2023-06-30';
```

#### 9 Optimize Date Comparisons:

```
# Use BETWEEN instead of > and <
SELECT * FROM Orders WHERE order_date BETWEEN '2023-01-01' AND '2023-06-30';
```

#### 10 Partitioning:

```
# Create partitions for large databases
CREATE TABLE Orders (
    order_id INT PRIMARY KEY,
    order_date DATE,
    order_amount DECIMAL
) PARTITION BY RANGE (YEAR(order_date));
```

### Strategies

```
# Create an index on the "age" column in the "Customers" table
CREATE INDEX idx_age ON Customers (age);
```

```
# Select only the necessary columns for the query
SELECT name, email FROM Customers WHERE age > 30;
```

```
# Inner join for matching records between "Customers" and "Orders" tables
SELECT Customers.name, Orders.order_id
FROM Customers
INNER JOIN Orders ON Customers.customer_id = Orders.customer_id
WHERE Orders.order_id > 100;
```

```
# Use proper aggregation and filtering
SELECT MAX(salary) FROM Employees WHERE department = 'IT';
```

```
# Analyze the execution plan to optimize the query
EXPLAIN SELECT * FROM Customers WHERE age > 30;
```

```
# Use proper date comparisons
SELECT * FROM Orders WHERE order_date BETWEEN '2023-01-01' AND '2023-06-30';
```

## 2.4 Definitions - File Reading

### File Reading

File reading access data from files. It's vital in data collection to ingest, preprocess, and integrate information from diverse sources.

**Use of File Reading in Data Collection:**

**Data Ingestion:** File reading enables the transfer of data from files to applications or databases, making it accessible for further processing or analysis.

**Data Interoperability:** Different file formats are used for data storage, and file reading provides a standardized way to access data regardless of the file type.

**Data Preprocessing:** File reading allows you to parse and preprocess data before using it in applications, ensuring that it is in a usable format.

**Data Integration:** In data collection, data often comes from various sources in different formats. File reading helps combine and integrate data into a uniform structure for analysis and reporting.

**Automation:** File reading can be incorporated into automated data collection pipelines to streamline the process of extracting and processing data.

### File Reading

```
import pandas as pd
```

```
# Read data from a CSV file
data = pd.read_csv('file.csv')
```

```
# Read data from an Excel file
data = pd.read_excel('file.xlsx')
```

```
# Read data from a JSON file
data = pd.read_json('file.json')
```

**Optimization**

```
1 Reading Large Files in Chunks:
def read_large_file(filename):
    with open(filename, 'r') as file:
        while True:
            chunk = file.read(1024) # Read in 1KB chunks
            if not chunk:
                break
            # Process the chunk

2 Using Buffered Reading:
def read_file_with_buffered(filename):
    with open(filename, 'r', buffering=8192) as file: # Use a larger buffer size
        data = file.read()
    # Process the data

3 Context Management with Iteration:
def read_file_in_chunks(filename):
    with open(filename, 'r') as file:
        for line in file:
            # Process each line

4 Memory Mapping (for large files):
import mmap

def read_large_file(filename):
    with open(filename, 'r') as file:
        mmap_obj = mmap.mmap(file.fileno(), length=0, access=mmap.ACCESS_READ)
        data = mmap_obj.read()
    # Process the data

5 Parallel File Reading Using multiprocessing:
import multiprocessing

def read_parallel_file(filename):
    with open(filename, 'r') as file:
        n_processes = multiprocessing.cpu_count()
        chunk_size = 1024 * 1024 # 1MB chunk size
        with multiprocessing.Pool(n_processes) as pool:
            results = pool.map(mmap_obj.read, [chunk_size] * n_processes)
    # Process the results
    ...


```

**Strategies****1 Reading Entire File:**

```
def read_entire_file(filename):
    with open(filename, 'r') as file:
        data = file.read()
    # Process the data


```

**2 Reading Line by Line:**

```
def read_line_by_line(filename):
    with open(filename, 'r') as file:
        for line in file:
            # Process each line
            # ...


```

**3 Reading Binary Files:**

```
def read_binary_file(filename):
    with open(filename, 'rb') as file:
        data = file.read()
    # Process the binary data
    ...


```

**2.5 Definitions – Sensor/Device Data Collection****Sensor/Device Data**

Sensor/Device Data Collection is the process of gathering information from sensors or devices for analysis, monitoring, or control purposes.

**Sensor/Device Data Collection Process:**

- Deployment:
  - Install sensors/devices in the target environment or system.
- Data Acquisition:
  - Sensors/devices measure physical parameters and generate electrical or digital signals.
- Data Transmission:
  - Transmit data from sensors/devices to a central collection point.
- Data Aggregation:
  - Gather data from multiple sensors/devices into a centralized repository.
- Data Preprocessing:
  - Cleanse and filter raw sensor data to remove noise or outliers.
- Data Storage:
  - Store collected and preprocessed data in databases or data warehouses.
- Data Integration:
  - Combine sensor/device data with other relevant sources (e.g., weather data, GIS).
- Real-time Processing (Optional):
  - Implement real-time data processing for immediate insights and control responses.
- Data Analysis and Visualization:
  - Analyze data using statistical techniques, machine learning, or domain-specific algorithms.
- Decision Making and Control:
  - Interpret data insights to make informed decisions or take appropriate actions.
- Data Archival and Retention:
  - Archive historical sensor/device data for reference or compliance.
- Data Security and Privacy:
  - Ensure data security during transmission, storage, and access.

**2.5.1 Code Snippets - Sensor/Device Data Collection****Sensor/Device Data**

```
import serial

# Connect to the serial port of the device
ser = serial.Serial('COM1', baudrate=9600)

# Read data from the device
data = ser.readline()

# Close the serial connection
ser.close()


```

**Optimization****1 Reading Large Data from Sensors in Chunks:**

```
# Read sensor data in chunks (do it at once)
def read_sensor_data_in_chunks(sensor):
    chunk_size = 1024
    while True:
        data_chunk = sensor.read(chunk_size)
        if not data_chunk:
            break
        # Process data_chunk
        # ...


```

**2 Asynchronous Data Acquisition:**

```
import asyncio

# Asynchronous data acquisition from multiple sensors
async def read_sensor_data(sensor):
    tasks = [read_sensor_data(sensor) for sensor in sensors]
    await asyncio.gather(*tasks)
    return tasks

async def main():
    sensors = [Sensor10, Sensor20, Sensor30]
    tasks = [read_sensor_data(sensor) for sensor in sensors]
    await asyncio.gather(*tasks)
    sensor_data = tasks[0].result()


```

**3 Data Preprocessing with Pandas (for CSV Data):**

```
import pandas as pd

# Read data from a CSV file
df = pd.read_csv('data.csv')

# Drop rows with missing values
df.dropna(inplace=True)

# Check for outliers
df_element = df[(df['value'] > 0) & (df['value'] < 100)]


```

**4 Efficient Data Storage with SQLite:**

```
import sqlite3

# Create a database and table
conn = sqlite3.connect('sensor_data.db')
cursor = conn.cursor()
cursor.execute("CREATE TABLE IF NOT EXISTS sensor_data (time TEXT, value REAL, ...)")
cursor.execute("SELECT * FROM sensor_data WHERE time='2023-07-01 10:00:00'")

# Insert data
sensor_data = [(time, value) for time, value in cursor.fetchall()]
cursor.executemany("INSERT INTO sensor_data VALUES (?, ?)", sensor_data)
conn.commit()


```

**5 Using a Time-Series Database:**

```
from influxdb import InfluxDBClient
client = InfluxDBClient(host='localhost', port=8086)
client.switch_database('sensors')
data = [
    {"measurement": "temperature",
     "tags": {"sensor_id": "sensor1", "location": "room1"},
     "time": "2023-07-01T10:00:00",
     "fields": {"value": 25.5}},
    ...
]
# More data points...
client.write_points(data)
result = client.query('SELECT * FROM temperature WHERE location="room1"')


```

**3 Data Preparation****3.1 Definitions – Data Preparation**

Data preparation is the process of preparing raw data so that it is suitable for further processing and analysis. Key steps include collecting, cleaning, and labeling raw data into a form suitable for machine learning (ML) algorithms and then exploring and visualizing the data.

**3.2 Definitions – Reading Data from CSV****Reading Data from CSV**

Reading data from CSV (Comma-Separated Values) is the process of loading structured tabular data stored in a CSV file format into a program or data analysis environment for further processing, analysis, or modeling.

Process of Reading Data from CSV:

- Importing Required Library:
  - Import the necessary library (e.g., Pandas) to read CSV data in Python.
- Specifying File Path:
  - Provide the path or filename of the CSV file you want to read.
- Reading CSV Data:
  - Use the library function (e.g., Pandas' `read\_csv()`) to load the CSV data.
- Data Exploration:
  - Examine the loaded data to understand its structure and basic information.
- Data Cleaning and Preprocessing:
  - Handle missing values, outliers, and perform initial data cleaning if needed.
- Data Analysis & Modeling (Optional):
  - Proceed with data analysis or modeling tasks if applicable to the project.

**3.2 Code Snippets - Reading Data from CSV****Reading Data (CSV)**

```
import pandas as pd

# Step 1: Specify the file path
file_path = "path_to_your_csv_file.csv"

# Step 2: Read the CSV file
data = pd.read_csv(file_path)

# Step 3: Explore the data
print(data.head()) # Displays the first few rows of the data
print(data.shape) # Get the number of rows and columns in the data
print(data.info()) # Get information about data types and missing values


```

**Optimization****1 Using Chunking for Large Files:**

```
import pandas as pd

# Read data in chunks (e.g., 1000 rows at a time)
chunk_size = 1000
chunks = pd.read_csv('data.csv', chunksize=chunk_size)

for chunk in chunks:
    # Process each chunk
    # ...


```

**2 Specifying Data Types:**

```
# Explicitly set data types while reading CSV
data = pd.read_csv('data.csv', dtype=[{'column': int, 'column2': float, 'column3': str}])


```

**3 Using CSV Reader:**

```
import csv

# Read data using CSV Reader
with open('data.csv', 'r') as csv_file:
    csv_reader = csv.reader(csv_file)
    headers = next(csv_reader) # Skip headers if present
    for row in csv_reader:
        # Process each row
        # ...


```

**4 Skipping Rows and Columns:**

```
# Skip the first row (header) and read specific columns only
data = pd.read_csv('data.csv', skiprows=1, usecols=['column1', 'column2'])


```

**5 Using CSV Librarian:**

```
import askaiter

# Read data using fask for parallel processing
data = pd.read_csv('data.csv')
# Process operations on 'data' direct DataFrame


```

**3.3 Definitions – Handling Missing Values****Handling Missing Values**

When dealing with missing data, data scientists can use two primary methods to solve the error: imputation or the removal of data.

Process of Handling Missing Values:

- Identifying Missing Values:
  - Check for missing values in the dataset.
- Understanding the Pattern:
  - Analyze the missing data pattern (MCAR, MAR, MNAR).
- Deletion or Dropping:
  - Remove rows or columns with missing values.
- Imputation Techniques:
  - Fill missing values with estimated data.
- Creating Missing Value Indicators:
  - Generate binary indicators for missing values.
- Advanced Imputation Techniques:
  - Use multiple imputation or probabilistic methods for complex cases.
- Consider Domain Knowledge:
  - Apply domain knowledge to handle missing values.
- Evaluate and Monitor:
  - Assess the impact of handling and monitor for potential issues.

**Missing Values**

```
import pandas as pd

# Step 1: Read data from CSV (assuming 'data' is a DataFrame)
data = pd.read_csv('data.csv')

# Step 2: Identifying missing values
missing_values = data.isnull().sum()

# Step 3: Deletion/Dropping Missing Values (if applicable)
data_cleaned = data.dropna()

# Step 4: Imputation (fill missing values) using mean (for example)
data_imputed = data.fillna(data.mean())

# Step 5: (Optional): Creating Missing Value Indicators
data['column_mean_indicator'] = data['column_mean'].isnull().astype(int)


```

**Optimization****1 Dropping Rows with Missing Values:**

```
# Drop rows with missing values in a specific column
data_cleaned = data.dropna(subset=['column_name'])

# Fill missing values in a numerical column with its median
data['numerical_column'].fillna(data['numerical_column'].median(), inplace=True)

# Fill missing values in a categorical column with its mode
data['categorical_column'].fillna(data['categorical_column'].mode()[0], inplace=True)


```

**2 Forward-Fill and Backward-Fill:**

```
# Forward-Fill (forward fill) - use the last valid observation to fill
data.fillna(method='ffill', inplace=True)

# Backward-Fill (backward fill) - use the next valid observation to fill
data.fillna(method='bfill', inplace=True)


```

**3 Forward-Fill and Backward-Fill:**

```
# Forward-Fill (forward fill) - use the last valid observation to fill
data.fillna(method='ffill', inplace=True)

# Backward-Fill (backward fill) - use the next valid observation to fill
data.fillna(method='bfill', inplace=True)


```

**4 Advanced Imputation with KNN:**

```
# Use KNNImputer to fill missing values with nearest neighbors
knn_imputer = KNNImputer(n_neighbors=5)
# Impute missing values using KNN
data_imputed = knn_imputer.fit_transform(data)


```

**5 Using Pandas' fillna() with Dictionary:**

```
# Fill missing values using a dictionary (e.g., different values for different columns)
imputation_dict = {'column1': value1, 'column2': value2, ...}
data.fillna(imputation_dict, inplace=True)


```

**6 Impute Categorical Variables with 'Unknown' Category:**

```
# Fill missing values in a categorical column with 'Unknown'
data['categorical_column'].fillna('Unknown', inplace=True)


```

**Optimization****7 Imputation with Linear Regression:**

```
from sklearn.linear_model import LinearRegression
# Separate features and target columns
data_features = data[['feature1', 'feature2', 'feature3']]
data_target = data['target']
# Train a linear regression model to predict missing values
model = LinearRegression()
model.fit(data_features, data_target)
# Predict missing values
predicted_values = model.predict(data_features[['feature1', 'feature2']])

# Fill missing values with predicted values
data.loc[data['categorical_column'].isnull(), 'categorical_column'] = predicted_values


```

**3.4 Definitions – Data Transformation****Data Transformation**

Data transformation is the process of converting, cleaning, and structuring data into a usable format that can be analyzed to support decision making processes, and to propel the growth of an organization.

Process of Data Transformation:

- Data Cleaning and Preprocessing:
  - Handle missing values and prepare data for further processing.
- Feature Scaling:
  - Normalize or standardize numerical features to a common scale.
- Feature Engineering:
  - Create new features or derive relevant information from existing ones.
- Encoding Categorical Variables:
  - Convert categorical variables into numerical representations.
- Data Discretization (Binning):
  - Group continuous values into bins or intervals for simplification.
- Log Transformation:
  - Apply log transformation to reduce the impact of extreme values.
- Data Normalization:
  - Scale data to ensure equal weightage for all features.
- Data Smoothing:
  - Apply techniques to remove noise from time series data.

**9 Principal Component Analysis (PCA):**

Reduce data dimensionality for better visualization or modeling.

**10 Data Integration and Aggregation:**

Merge data from different sources and aggregate it for analysis.

**11 Data Reduction (Sampling):**

Reduce the size of large datasets using sampling techniques.

**12 Handling Skewed Data (if applicable):**

Apply transformations to handle data with a skewed distribution.

**13 Handling Seasonality (for time series data):**

Use techniques like seasonal decomposition to extract seasonal components.

**3.4.1 Code Snippets – Data Transformation****Data Transformation**

```
# Convert categorical variables to numerical using one-hot encoding
data_encoded = pd.get_dummies(data, columns=['categorical_column'])

# Log transformation of a column (for positive values)
data['log_transformed'] = np.log(data['original_column'])


```

**Handling Missing Values with Imputation**

```
# Use SimpleImputer for numerical columns
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy='mean')
data[['numerical_feature1', 'numerical_feature2']] = imputer.fit_transform(data[['numerical_feature1', 'numerical_feature2']])

# Use most frequent value for categorical columns
data['categorical_feature'].fillna(data['categorical_feature'].mode()[0], inplace=True)


```

**Feature Scaling with Min-Max Scaling**

```
from sklearn.preprocessing import MinMaxScaler
# Use MinMaxScaler for scaling numerical features to the range [0, 1]
scaler = MinMaxScaler()
data[['numerical_feature1', 'numerical_feature2']] = scaler.fit_transform(data[['numerical_feature1', 'numerical_feature2']])


```

**One-Hot Encoding with get\_dummies()**

```
data_encoded = pd.get_dummies(data, columns=['categorical_column'], drop_first=True)


```

**Log Transformation (avoiding zero values)**

```
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
data[['numerical_feature1', 'numerical_feature2']] = scaler.fit_transform(data[['numerical_feature1', 'numerical_feature2']])


```

**Data Normalization (Z-score Scaling)**

```
from sklearn.preprocessing import StandardScaler
# Use StandardScaler for z-score normalization
scaler = StandardScaler()
data[['numerical_feature1', 'numerical_feature2']] = scaler.fit_transform(data[['numerical_feature1', 'numerical_feature2']])


```

**Principal Components Analysis (PCA)**

```
from sklearn.decomposition import PCA
# Use PCA to reduce data dimensionality
pca = PCA(n_components=2)
pca_data = pca.fit_transform(data[['numerical_feature1', 'numerical_feature2']])


```

**Data Aggregation with GroupBy**

```
grouped_data = data.groupby('group_column')[['value_column']].mean()


```

**Applying Transformations to Multiple Columns**

```
# List of selected column names
numerical_columns = ['numerical_feature1', 'numerical_feature2']

# Apply Min-Max Scaling to multiple columns
data[numerical_columns] = scaler.fit_transform(data[numerical_columns])


```

**Data Reduction (Sampling)**

```
# Use sample() to randomly sample a fraction of data (e.g., 30%)
sampled_data = data.sample(frac=0.3)


```

```
# Use sample() to randomly sample a fixed number of data points (e.g., 1000)
sampled_data = data.sample(n=1000)


```

**Handling Skewed Data with Log Transformation**

```
import numpy as np

# Use log transformation to handle skewed data (avoiding zero values)
data['log_transformed'] = data['original_column'].apply(lambda x: np.log(x) if x > 0 else 0)


```

**Handling Outliers**

```
# Use percentile-based outlier detection and removal
q1_low = data['numerical_feature'].quantile(0.05)
q1_high = data['numerical_feature'].quantile(0.95)
data = data[(data['numerical_feature'] > q1_low) & (data['numerical_feature'] < q1_high)]


```

**Handling Seasonality  
(for Time Series Data)**

```
# Assuming you have a time series dataset named 'ts_data'
# ts_data.index = pd.date_range('1990-01-01', periods=100)
# ts_data['date'] = pd.date_range('1990-01-01', periods=100)

# Use seasonal decomposition using statsmodel's seasonal_decompose
from statsmodels.tsa.seasonal import seasonal_decompose
result = seasonal_decompose(ts_data['Sales'], model='additive', period=12)

# Extract seasonal, trend, and residual components
seasonal = result.seasonal
trend = result.trend
residual = result.resid
```

**Data Smoothing**

```
# Use rolling window (e.g., 5 days) for moving average smoothing
data['smoothed_values'] = data['original_values'].rolling(window=5).mean()
```

**Feature Engineering**

```
# Create new features based on existing ones
data['new_feature'] = data['feature1'] + data['feature2']

# Extract date components from datetime feature
data['year'] = data['date'].dt.year
data['month'] = data['date'].dt.month
data['day'] = data['date'].dt.day
```

**3.5 Definitions – Data Scaling****Data Scaling**

Data scaling is the process of transforming numerical features to a common scale, ensuring that they have comparable magnitudes. It is a crucial step in data preprocessing, especially when working with algorithms that are sensitive to the scale of features, such as distance-based methods and gradient-based optimization algorithms.

**Process of Data Scaling:**

1. Select Numerical Features:  
Identify the numerical features that require scaling.
2. Choose Scaling Method:  
Decide on the scaling method based on the nature of the data and the requirements of the machine learning model. Common scaling methods include Min-Max Scaling (Normalization) and Z-score Scaling (Standardization).
3. Apply Scaling:  
Scale the numerical features based on the chosen method to bring them to a common scale.

**3.5.1 Code Snippets – Data Scaling****Data Scaling**

```
from sklearn.preprocessing import StandardScaler

# Scale numerical features to have mean=0 and variance=1
scaler = StandardScaler()
scaled_data = scaler.fit_transform(data[['numeric_features', 'non-numeric_features']])
```

**Min-Max Scaling (Normalization)  
with MinMaxScaler**

```
from sklearn.preprocessing import MinMaxScaler

# Create a DataFrame with numerical features
data = pd.DataFrame({'Feature1': [10, 20, 30, 40], 'Feature2': [100, 200, 300, 400]})

# Use MinMaxScaler for scaling numerical features to the range (0, 1)
scaler = MinMaxScaler()
scaled_data = scaler.fit_transform(data[['Feature1', 'Feature2']])
```

**Z-score Scaling (Standardization)  
with StandardScaler**

```
from sklearn.preprocessing import StandardScaler

# Create a DataFrame with numerical features
data = pd.DataFrame({'Feature1': [10, 20, 30, 40], 'Feature2': [100, 200, 300, 400]})

# Use StandardScaler for z-score normalization
scaler = StandardScaler()
scaled_data = scaler.fit_transform(data[['Feature1', 'Feature2']])
```

**Scaling Specific Features  
with MinMaxScaler**

```
from sklearn.preprocessing import MinMaxScaler

# Create a DataFrame with numerical features
data = pd.DataFrame({'Age': [25, 35, 45], 'Income': [50000, 75000, 60000]})

# Use MinMaxScaler to scale specific features to the range (0, 1)
scaler = MinMaxScaler()
data['Scaled_Age'] = scaler.fit_transform(data[['Age']])
```

**Applying Scaling to Test Data Using Scaled**

```
from sklearn.preprocessing import MinMaxScaler

# Create a DataFrame with numerical features
train_data = pd.DataFrame({'Feature1': [10, 20, 30, 40], 'Feature2': [100, 200, 300, 400]})

# Use MinMaxScaler for scaling numerical features to the range (0, 1)
scaler = MinMaxScaler()
scaler.fit(train_data)
scaled_train_data = scaler.transform(train_data[['Feature1', 'Feature2']])

# Save the scaler for future use (e.g., scaling test data)
import joblib
joblib.dump(scaler, 'scaler.pkl')

# Later, when scaling test data:
loadscaler = joblib.load('scaler.pkl')
scaled_test_data = loadscaler.transform(test_data[['Feature1', 'Feature2']])
```

**3.6 Definitions – Data Splitting for Machine Learning****Data Splitting for ML**

Data splitting is when data is divided into two or more subsets. Typically, with a two-part split, one part is used to evaluate or test the data and the other to train the model. Data splitting is an important aspect of data science, particularly for creating models based on data.

**Process of Data Splitting:**

1. Data Preparation:  
Preprocess, clean, handle missing values, and engineer features for analysis or modeling.
2. Shuffling (Optional):  
Randomly shuffle data to eliminate order bias and ensure better generalization.
3. Train-Test Split:  
Divide data into training and testing sets for model training and evaluation.
4. Optional Validation Set:  
Further split training data for hyperparameter tuning and model optimization.
5. Model Training:  
Train machine learning models using the training set for learning patterns.
6. Model Evaluation:  
Assess model performance on the testing (or validation) set for generalization testing.
7. Adjustment and Improvement:  
Fine-tune models based on evaluation results to improve accuracy and efficiency.

**3.6.1 Code Snippets – Data Splitting for Machine Learning****Data Splitting for ML**

```
from sklearn.model_selection import train_test_split

# Split data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(features, target, test_size=0.2)
```

**Optimized Data Splitting with Stratified Sampling  
(for Classification Task)**

```
from sklearn.model_selection import train_test_split

# Assuming you have a DataFrame "data" with features and labels
X = data.drop(['target'], axis=1) # Features
y = data['target'] # Labels

# Use stratified sampling to preserve class distribution in training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

**Strategies****1 Simple Train-Test Split:**

```
from sklearn.model_selection import train_test_split

# Assuming you have a DataFrame "data" with features and labels
X = data.drop(['target'], axis=1) # Features
y = data['target'] # Labels

# Split the data into a training set and a testing set (80-20 split)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

**2 Stratified Split for Classification Tasks:**

```
from sklearn.model_selection import train_test_split

# Assuming you have a DataFrame "data" with features and class labels "target"
X = data.drop(['target'], axis=1) # Features
y = data['target'] # Class labels

# Use stratified sampling to preserve class distribution in training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y, random_state=42)
```

**3 Time Series Data Splitting:**

```
# For time series data, split sequentially based on time
train_size = int(0.8 * len(data))
X_train, X_test, y_train, y_test = data[:train_size], data[train_size:]
```

**4 Cross-Validation Split:**

```
from sklearn.model_selection import KFold
```

```
# Assuming you have a DataFrame "data" with features and labels
X = data.drop(['target'], axis=1) # Features
y = data['target'] # Labels

# Use K-Fold cross-validation (e.g., K=5)
kf = KFold(n_splits=5, shuffle=True, random_state=42)
for train_index, test_index in kf.split(X):
    X_train, X_test = X.loc[train_index], X.loc[test_index]
    y_train, y_test = y.loc[train_index], y.loc[test_index]
```

**5 Group-Based Split (for Grouped Data):**

```
from sklearn.model_selection import GroupShuffleSplit

# Assuming you have a DataFrame "data" with features, labels, and groups
X = data.drop(['target'], axis=1) # Features
y = data['target'] # Labels
groups = data['group'] # Group labels

# Use GroupShuffleSplit for grouped data splitting
gss = GroupShuffleSplit(n_splits=1, test_size=0.2, random_state=42)
train_index, test_index = next(gss.split(X, y, groups))
X_train, X_test, y_train, y_test = X.loc[train_index], X.loc[test_index], y.loc[train_index], y.loc[test_index]
```

**3.7 Definitions – Feature Selection****Feature Selection**

Feature selection is the process of selecting a subset of relevant and informative features (variables or attributes) from the original set of features in a dataset. The goal of feature selection is to improve the performance of machine learning models by reducing the dimensionality of the data, removing irrelevant or redundant features, and focusing on the most important ones that have a significant impact on the target variable.

**Process of Feature Selection:**

1. Data Preparation:  
Clean, preprocess, and encode categorical variables if necessary for feature selection.
2. Feature Importance Ranking:  
Rank features based on relevance scores to prioritize their selection.
3. Filter Methods:  
Apply statistical tests to select the most informative features effectively.
4. Wrapper Methods:  
Utilize ML models to evaluate subsets of features iteratively.
5. Embedded Methods:  
Leverage model's built-in feature selection capabilities during training.
6. Validation/Evaluation:  
Selects features impact on model performance through evaluation metrics.
7. Model Adjustment:  
Iterate with different techniques or settings to optimize feature selection outcome.
8. Final Model Training:  
Train ML model with the chosen subset of features for prediction or analysis.

**3.7.1 Code Snippets – Feature Selection****Feature Selection**

```
from sklearn.feature_selection import SelectKBest, f_regression
```

```
# Select top k features based on ANOVA F-value
selector = SelectKBest(score_func=f_regression, k=5)
selected_features = selector.fit_transform(X_train, y_train)
```

**Optimization****1 Feature Importing Ranking with Random Forest:**

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.feature_selection import SelectKBest, f_classif
```

```
# Assuming you have X_train and y_train for training
# Create a Random Forest classifier
rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)
```

```
# Fit the classifier to the training data
rf_classifier.fit(X_train, y_train)
```

```
# Use Feature_importance_ to get Feature importance scores
feature_importance = rf_classifier.feature_importance_
```

```
# Use SelectKBest to select features based on a threshold
sfm = SelectKBest(f_classif, threshold=0.1)
X_train_selected = sfm.fit_transform(X_train, y_train)
```

**2 Recursive Feature Elimination with Cross-Validation (RFE):**

```
from sklearn.feature_selection import RFE
```

```
from sklearn.linear_model import LogisticRegression
```

```
# Assuming you have X_train and y_train for training
# Create a logistic regression model (for any other model)
logreg_model = LogisticRegression()
```

```
# Use RFE to recursively eliminate features with cross-validation
rfe = RFE(logreg_model, n_features_to_select=5)
X_train_selected = rfe.fit_transform(X_train, y_train)
```

**3 Feature Selection with SelectKBest using Chi-Square:**

```
from sklearn.feature_selection import SelectKBest, chi2
```

```
# Assuming you have X_train and y_train for training
```

```
# Use SelectKBest with chi2 to select the top k features
```

```
k = 10 # Set the number of top features you want to select
selector = SelectKBest(score_func=chi2, k=k)
X_train_selected = selector.fit_transform(X_train, y_train)
```

**4 Feature Selection with L1 Regularization (Lasso):**

```
from sklearn.linear_model import Lasso
```

```
# Assuming you have X_train and y_train for training
```

```
# Use Lasso with regularization for feature selection
```

```
lasso_model = Lasso(alpha=0.1)
X_train_selected = lasso_model.fit_transform(X_train, y_train)
```

**5 Feature Selection with Mutual Information:**

```
from sklearn.feature_selection import SelectKBest, mutual_info_classif
```

```
# Assuming you have X_train and y_train for training
```

```
# Use SelectKBest with mutual information for feature selection
```

```
k = 10 # Set the number of features you want to select
selector = SelectKBest(score_func=mutual_info_classif, k=k)
X_train_selected = selector.fit_transform(X_train, y_train)
```

**6 Feature Selection with Recursive Feature Elimination and Grid Search:**

```
from sklearn.feature_selection import RFE
```

```
from sklearn.linear_model import LogisticRegression
```

```
# Assuming you have X_train and y_train for training
```

```
# Create a logistic regression model (for any other model)
```

```
logreg_model = LogisticRegression()
```

```
# Use RFE to recursively eliminate features with cross-validation
```

```
rfe = RFE(logreg_model, n_features_to_select=5, step=1)
```

```
param_grid = {'n_estimators': [5, 10, 15]} # Set the number of trees you want to select
```

```
grid_search = GridSearchCV(rfe, param_grid, cv=5)
```

```
# Get the best features subset from grid search
```

```
X_train_selected = grid_search.best_estimator_.transform(X_train)
```

**Optimization****7 Feature Selection with Variance Threshold:**

```
from sklearn.feature_selection import VarianceThreshold
```

```
# Assuming you have a DataFrame "data" with features and labels
```

```
# Create a VarianceThreshold object for variance features
```

```
threshold = 0.01 # Set the threshold for variance
```

```
vt = VarianceThreshold(threshold=threshold)
```

```
X_train_selected = vt.fit_transform(X_train)
```

**8 Feature Selection with SelectFromModel and Grid Search:**

```
from sklearn.feature_selection import SelectFromModel
```

```
from sklearn.ensemble import RandomForestClassifier
```

```
from sklearn.model_selection import GridSearchCV
```

```
# Assuming you have X_train and y_train for training
```

```
# Create a Random Forest classifier
```

```
rf_classifier = RandomForestClassifier()
```

```
# Use SelectFromModel to select features based on a threshold
```

```
sfm = SelectFromModel(rf_classifier)
```

```
param_grid = {'threshold': [0.1, 0.2, 0.3]} # Set the threshold values for feature selection
```

```
grid_search = GridSearchCV(sfm, param_grid, cv=5)
```

```
# Get the best feature subset from grid search
```

```
X_train_selected = grid_search.best_estimator_.transform(X_train)
```

**3.8 Definitions – Handling Imbalanced Data (Oversampling)****Handling Imbalanced Data**

Handling imbalanced data is a critical preprocessing step in machine learning when the distribution of classes in the target variable is highly skewed. Imbalanced data can lead to biased model performance, where the model becomes more accurate in predicting the majority class but performs poorly on the minority class. The goal of handling imbalanced data is to ensure that the model can effectively learn patterns from both classes, providing a more balanced and fair performance.

**Process of Handling Imbalanced Data:**

1. Data Analysis:  
Examine class distribution to identify imbalance and its impact on model performance.
2. Resampling Techniques:  
Use oversampling, undersampling, or SMOTE to balance the dataset.
3. Class Weights:  
Assign higher weight to the minority class during model training.
4. Ensemble Methods:  
Utilize ensemble techniques to handle imbalanced data effectively.
5. Performance Metrics:  
Use precision, recall, F1-score, or ROC-AUC for evaluation.
6. Fine-Tuning:  
Continuously monitor performance and adjust techniques for better balance and accuracy.

**3.8.1 Code Snippets – Handling Imbalanced Data****Handling Imbalanced Data**

```
from imblearn.over_sampling import SMOTE
```

```
# Apply SMOTE to oversample the minority class
```

```
smote = SMOTE()
```

```
X_resampled, y_resampled = smote.fit_resample(X_train, y_train)
```

**Resampling Techniques****1 OverSampling:**

```
from imblearn.over_sampling import SMOTE
```

```
# Assuming you have X_train and y_train for training
```

```
# Create X_train and y_train to create a single dataset
```

```
train_data = pd.concat([X_train, y_train], axis=1)
```

```
# Separate majority and minority classes
```

```
majority_class = train_data[train_data['target'] == 0]
```

```
minority_class = train_data[train_data['target'] == 1]
```

```
# oversample the minority class to match the majority class size
```

```
minority_oversampled = resample(minority_class, replace=True, n_samples=len(majority_class), random_state=42)
```

```
# Create majority and minority oversampled classes
```

```
balanced_data = pd.concat([majority_oversampled, minority_oversampled])
```

```
# Split the balanced data into X_train balanced and y_train balanced
```

```
X_train_balanced = balanced_data.drop(['target'], axis=1)
```

```
y_train_balanced = balanced_data['target']
```

**Class Weights**

```
from sklearn.linear_model import LogisticRegression
```

```
# Assuming you have X_train and y_train for training
```

```
# Create a logistic regression model with balanced class weights
```

```
logreg_model = LogisticRegression(class_weight='balanced', random_state=42)
```

```
logreg_model.fit(X_train, y_train)
```

**Ensemble Methods**

```
from imblearn.ensemble import BalancedRandomForestClassifier
```

```
# Assuming you have X_train and y_train for training
```

```
# Create a Random Forest classifier with balanced class weights
```

```
rf_classifier = BalancedRandomForestClassifier(class_weight='balanced', random_state=42)
```

```
rf_classifier.fit(X_train, y_train)
```

**SMOTE (Synthetic Minority Over-sampling Technique)**

```
from imblearn.over_sampling import SMOTE
```

```
# Assuming you have X_train and y_train for training
```

```
# Apply SMOTE to oversample the minority class
```

```
smote = SMOTE(sampling_strategy='auto', random_state=42)
```

```
X_train_balanced, y_train_balanced = smote.fit_resample(X_train, y_train)
```

**3.9 Definitions – Data Normalization(Min-Max Scaling)****Data Normalization**

Data normalization is the process of reorganizing data within a database so that users can utilize it for further queries and analysis.

**Process of Data Normalization:**

1. Selecting Features:  
Identify the numerical features (attributes) that need normalization in the dataset.

2. Choosing Normalization Method:  
Decide on the appropriate normalization method based on the data distribution and the requirements of the machine learning algorithm.

**Min-Max Normalization:**

3. Standardization (Standardization):  
Standardize the features to a predefined range, often [0, 1], using the formula:

4. Scaling to [-1, 1] Range:  
Standardize the features to have a mean of 0 and a standard deviation of 1, using the formula:

5. Scaling to [0, 1] Range:  
Use a variation of Min-Max normalization to scale features to the range [0, 1], particularly if the data contains negative values:

**3.9.1 Code Snippets – Data Normalization****Data Normalization**

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

```
# Scale features to a specific range (e.g., [0, 1])
```

```
scaler = MinMaxScaler()
```

```
normalized_data = scaler.fit_transform(data[['numeric_features', 'non-numeric_features']])
```

**Min-Max Normalization**  
from sklearn.preprocessing import MinMaxScaler  
  
# Assuming you have X\_train and X\_test for training and testing  
# Create a MinMaxScaler object  
scaler = MinMaxScaler()  
  
# Fit the scaler on the training data and transform both training and testing data  
X\_train\_normalized = scaler.fit\_transform(X\_train)  
X\_test\_normalized = scaler.transform(X\_test)

**Z-Score Normalization**  
from sklearn.preprocessing import StandardScaler  
  
# Assuming you have X\_train and X\_test for training and testing  
# Create a StandardScaler object  
scaler = StandardScaler()  
  
# Fit the scaler on the training data and transform both training and testing data  
X\_train\_standardized = scaler.fit\_transform(X\_train)  
X\_test\_standardized = scaler.transform(X\_test)

**Scaling to [-1, 1] Range**  
from sklearn.preprocessing import RobustScaler  
  
# Assuming you have X\_train and X\_test for training and testing  
# Calculate the min and max values of the dataset  
min\_value = X\_train.min().max()  
max\_value = X\_train.max().min()  
  
# Scale the features in the range [-1, 1]  
X\_train\_normalized = 2 \* (X\_train - min\_value) / (max\_value - min\_value) - 1  
X\_test\_normalized = 2 \* (X\_test - min\_value) / (max\_value - min\_value) - 1

**Robust Scaling**  
from sklearn.preprocessing import RobustScaler  
  
# Assuming you have X\_train and X\_test for training and testing  
# Create a RobustScaler object  
scaler = RobustScaler()  
  
# Fit the scaler on the training data and transform both training and testing data  
X\_train\_scaled = scaler.fit\_transform(X\_train)  
X\_test\_scaled = scaler.transform(X\_test)

**Log Transformation**  
import numpy as np  
  
# Assuming you have X\_train and X\_test for training and testing  
# Apply log transformation to the features  
X\_train\_log = np.log(X\_train)  
X\_test\_log = np.log(X\_test)

**Power Transformation**  
(Box-Cox or Yeo-Johnson)  
from sklearn.preprocessing import PowerTransformer  
  
# Assuming you have X\_train and X\_test for training and testing  
# Create a PowerTransformer object with method 'box-cox' or 'yeo-johnson'  
scaler = PowerTransformer(method='boxcox')  
  
# Fit the scaler on the training data and transform both training and testing data  
X\_train\_transformed = scaler.fit\_transform(X\_train)  
X\_test\_transformed = scaler.transform(X\_test)

## 3.10 Definitions – Data Binning

**Data Scaling**  
Binning is a way to group a number of more or less continuous values into a smaller number of "bins".  
  
Process of Data Binning:  

1. Data Selection:  
Choose the numerical feature (attribute) that requires binning.
2. Determining the Number of Bins:  
Decide the number of bins or intervals to create. This can be based on domain knowledge, the desired level of granularity, or using statistical methods like Sturge's rule or Scott's rule.
3. Binning Method:  
- Select a binning method:  
- Equal-Width Binning:  
- Divide the data range into equal-sized intervals.  
- Equal-Cardinality Binning:  
- Divide the data into intervals with an equal number of data points in each bin.  
- Custom Binning:  
- Define custom bin boundaries based on specific criteria or requirements.
4. Binning Process:  
Assign each data point to its corresponding bin based on its value.
5. Data Transformation:  
Replace the original numerical values with the bin labels or bin identifiers.

### 3.10.1 Code Snippets – Data Binning

**Data Scaling**  
from sklearn.preprocessing import StandardScaler  
  
# Scale numerical features to have mean=0 and variance=1  
scaler = StandardScaler()  
scaler.fit(data[['numerical\_feature']]).transform(data[['numerical\_feature']])

**Optimization**  
**1 Equal-Width Binning:**  
import pandas as pd  
  
# Assuming you have a DataFrame named 'data' with a numerical column 'age'  
# Define the number of bins you want  
num\_bins = 5  
  
# Calculate the width of each bin  
bin\_width = (data['age'].max() - data['age'].min()) / num\_bins  
# Define the bin edges based on the bin width  
bin\_edges = [data['age'].min()] + [i \* bin\_width for i in range(num\_bins - 1)]  
# Create a new 'age\_bin' column  
bin\_labels = ['Bin{0}'.format(i) for i in range(1, num\_bins + 1)]  
# Use pd.cut() to create a new 'age\_bin' column with the bin labels  
data['age\_bin'] = pd.cut(data['age'], bin\_edges, labels=bin\_labels, right=False)  
  
**2 Equal-Frequency Binning:**  
import pandas as pd  
  
# Assuming you have a DataFrame named 'data' with a numerical column 'score'  
# Define the number of bins you want  
num\_bins = 5  
# Use pd.qcut() to create a new 'score\_bin' column with the bin labels  
data['score\_bin'] = pd.qcut(data['score'], num\_bins, labels=False)

## 3.11 Definitions – Handling DateTime Data

**Handling DateTime Data**  
Handling datetime data is a crucial part of data preparation when working with datasets that contain date and time information. The process involves converting datetime data into a format that is suitable for analysis and machine learning tasks.  
Handling datetime data typically includes the following steps:

1. Data Parsing:  
Convert datetime strings to machine-readable format.
2. Feature Extraction:  
Extract relevant components (e.g., year, month, day).
3. Datetime Indexing:  
Set datetime values as the dataset's index for time-based analysis.
4. Handling Time Zones:  
Handle time zone conversions to maintain consistency.
5. Visualizing Temporal Data:  
Create time series plots to identify trends and patterns.

### 3.11.1 Code Snippets – Handling DateTime Data

**Handling DateTime Data**  
# Converting string to datetime objects  
data['date\_column'] = pd.to\_datetime(data['date\_column'])  
  
# Extract features from datetime (e.g., year, month, day)  
data['year'] = data['date\_column'].dt.year  
data['month'] = data['date\_column'].dt.month  
data['day'] = data['date\_column'].dt.day

**Optimization**  
**1 Data Parsing and Conversion:**  
import pandas as pd  
  
# Assuming you have a DataFrame named 'data' with a column 'date\_time' containing datetime strings  
# Convert 'date\_time' column to datetime objects  
data['date\_time'] = pd.datetime(data['date\_time'])

**2 Feature Extraction:**  
# Assuming you have a DataFrame named 'data' with a column 'date\_time' of datetime objects  
# Extract year, month, day, hour, minute as new columns  
data['year'] = data['date\_time'].dt.year  
data['month'] = data['date\_time'].dt.month  
data['day'] = data['date\_time'].dt.day  
data['hour'] = data['date\_time'].dt.hour  
data['minute'] = data['date\_time'].dt.minute

**3 Datetime Indexing:**  
# Assuming you have a DataFrame named 'data' with a 'date\_time' column of datetime objects  
# Set 'date\_time' as the datetime index for time-based analysis  
data.set\_index('date\_time', inplace=True)

**4 Handling Time Zones:**  
# Assuming you have a DataFrame named 'data' with a 'date\_time' column of datetime objects  
# Convert timezone to a different time zone (e.g., UTC to Eastern)  
data['date\_time'] = data['date\_time'].dt.tz\_convert('US/Eastern')

**5 Visualization Temporal Data:**  
import matplotlib.pyplot as plt  
# Assuming you have a DataFrame named 'data' with a 'date\_time' column of datetime objects  
# Create a time series plot for visualization  
plt.figure(figsize=(10, 6))  
plt.plot(data['date\_time'], data['value'])  
plt.xlabel('Date Time')  
plt.ylabel('Value')  
plt.title('Time Series Plot')  
plt.show()

## 4 Statistical Analysis

### 4.1 Definitions – Statistical Analysis

Statistical analysis, or statistics, is the process of collecting and analyzing data to identify patterns and trends, remove bias and inform decision-making. It's an aspect of business intelligence that involves the collection and scrutiny of business data and the reporting of trends.

### 4.2 Definitions – Descriptive Statistics

#### Descriptive Statistics

Descriptive statistics is a branch of statistical analysis that involves summarizing and presenting data in a meaningful and interpretable way. It provides an overview of the main features of a dataset, allowing analysts to understand its central tendencies, variability, and distribution.  
Process of Descriptive Statistics:

1. Central Tendency Measures:  
- Calculate mean, median, and mode to determine central values representing the dataset.
2. Variability Measures:  
- Find range, variance, and standard deviation to understand the spread of dispersion of data.
3. Distribution Shape Analysis:  
- Use histograms and kernel density plots to visualize data distribution's shape.  
- Assess skewness to understand the symmetry or skewness of the data.
4. Percentiles and Quartiles:  
- Calculate percentiles to identify values below specific percentages of the data.  
- Determine quartiles to divide data into four equal parts.
5. Data Visualization and Interpretation:  
- Utilize box plots, scatter plots, and bar charts to visualize the data.  
- Interpret descriptive statistics to gain insights into data characteristics and trends.

#### 4.2.1 Code Snippets - Descriptive Statistics

##### Descriptive Statistics

import numpy as np  
  
# Assuming you have a Pandas Series 'data'  
# Calculate mean, median, and standard deviation  
mean\_value = np.mean(data)  
median\_value = np.median(data)  
std\_deviation = np.std(data)  
  
# Count occurrences of each unique value in the data  
value\_counts = np.unique(data, return\_counts=True)

##### Optimization

**1 Calculating Mean, Median and Mode with Pandas:**  
import pandas as pd  
  
# Assuming you have a DataFrame named 'data' with a column 'values'  
# Calculate mean, median, and mode using Pandas  
mean\_value = data['values'].mean()  
median\_value = data['values'].median()  
mode\_value = data['values'].mode().iloc[0]

##### 2 Variability Measures with NumPy:

import numpy as np  
  
# Assuming you have a Numpy array named 'data'  
# Calculate range, variance, and standard deviation with NumPy

data\_range = np.ptp(data) # Range: Difference between max and min values  
data.variance = np.var(data) # Variance  
data\_std.dev = np.std(data) # Standard deviation

##### 3 Quantiles and Percentiles with Pandas:

import pandas as pd  
  
# Assuming you have a DataFrame named 'data'  
# Calculate quartiles and percentiles using Pandas

quartiles = data['values'].quantile([0.25, 0.5, 0.75])  
percentiles\_75 = data['values'].quantile(0.75) # 75th percentile

##### 4 Data Visualization with Matplotlib :

import matplotlib.pyplot as plt  
  
# Assuming you have a Pandas Series named 'data'  
# Create a histogram to visualize data distribution

plt.hist(data, bins=10, edgecolor='black')  
plt.xlabel('Age')  
plt.ylabel('Frequency')  
plt.title('Histogram')  
plt.show()

##### 5 Correlation Analysis:

# Assuming you have two sets of data: 'data1' and 'data2'

# Perform a two-sample t-test to compare means  
t\_statistic, p\_value = ttest\_ind(data1, data2)

if p\_value < 0.05:  
 print("Reject null hypothesis: Means are significantly different.")

else:  
 print("Fail to reject null hypothesis: Means are not significantly different.")

#### Hypothesis Testing (T-Test)

from scipy.stats import ttest\_ind

# Assuming you have two sets of data: 'data1' and 'data2'

# Perform an independent two-sample t-test  
t\_statistic, p\_value = ttest\_ind(data1, data2)

if p\_value < 0.05:  
 print("Reject null hypothesis: Means are significantly different.")

else:  
 print("Fail to reject null hypothesis: Means are not significantly different.")

#### Confidence Intervals

from scipy.stats import t

# Assuming you have a sample data 'data' and sample mean 'sample\_mean'

# Calculate the confidence interval for a given confidence level (e.g., 95%)

confidence\_level = 0.95

n = len(data)

standard\_error = data.std() / np.sqrt(n)

margin\_of\_error = standard\_error \* t.ppf((1 + confidence\_level) / 2, df=n - 1)

lower\_bound = sample\_mean - margin\_of\_error

upper\_bound = sample\_mean + margin\_of\_error

## Regression Analysis (Linear Regression)

import statsmodels.api as sm

# Assuming you have a DataFrame named 'data' with columns 'x' and 'y'

# Fit a linear regression

X = data['x']

y = data['y']

model = sm.OLS(y, X).fit()

intercept, slope = model.params

print(f'Intercept: {intercept}, Slope: {slope}')

print(f'Predicted value for x=1: {intercept + slope \* 1}')

print(f'Predicted value for x=2: {intercept + slope \* 2}')

print(f'Predicted value for x=3: {intercept + slope \* 3}')

print(f'Predicted value for x=4: {intercept + slope \* 4}')

print(f'Predicted value for x=5: {intercept + slope \* 5}')

print(f'Predicted value for x=6: {intercept + slope \* 6}')

print(f'Predicted value for x=7: {intercept + slope \* 7}')

print(f'Predicted value for x=8: {intercept + slope \* 8}')

print(f'Predicted value for x=9: {intercept + slope \* 9}')

print(f'Predicted value for x=10: {intercept + slope \* 10}')

print(f'Predicted value for x=11: {intercept + slope \* 11}')

print(f'Predicted value for x=12: {intercept + slope \* 12}')

print(f'Predicted value for x=13: {intercept + slope \* 13}')

print(f'Predicted value for x=14: {intercept + slope \* 14}')

print(f'Predicted value for x=15: {intercept + slope \* 15}')

print(f'Predicted value for x=16: {intercept + slope \* 16}')

print(f'Predicted value for x=17: {intercept + slope \* 17}')

print(f'Predicted value for x=18: {intercept + slope \* 18}')

print(f'Predicted value for x=19: {intercept + slope \* 19}')

print(f'Predicted value for x=20: {intercept + slope \* 20}')

print(f'Predicted value for x=21: {intercept + slope \* 21}')

print(f'Predicted value for x=22: {intercept + slope \* 22}')

print(f'Predicted value for x=23: {intercept + slope \* 23}')

print(f'Predicted value for x=24: {intercept + slope \* 24}')

print(f'Predicted value for x=25: {intercept + slope \* 25}')

print(f'Predicted value for x=26: {intercept + slope \* 26}')

print(f'Predicted value for x=27: {intercept + slope \* 27}')

print(f'Predicted value for x=28: {intercept + slope \* 28}')

print(f'Predicted value for x=29: {intercept + slope \* 29}')

print(f'Predicted value for x=30: {intercept + slope \* 30}')

print(f'Predicted value for x=31: {intercept + slope \* 31}')

print(f'Predicted value for x=32: {intercept + slope \* 32}')

print(f'Predicted value for x=33: {intercept + slope \* 33}')

print(f'Predicted value for x=34: {intercept + slope \* 34}')

print(f'Predicted value for x=35: {intercept + slope \* 35}')

print(f'Predicted value for x=36: {intercept + slope \* 36}')

print(f'Predicted value for x=37: {intercept + slope \* 37}')

print(f'Predicted value for x=38: {intercept + slope \* 38}')

print(f'Predicted value for x=39: {intercept + slope \* 39}')

print(f'Predicted value for x=40: {intercept + slope \* 40}')

print(f'Predicted value for x=41: {intercept + slope \* 41}')

print(f'Predicted value for x=42: {intercept + slope \* 42}')

print(f'Predicted value for x=43: {intercept + slope \* 43}')

print(f'Predicted value for x=44: {intercept + slope \* 44}')

print(f'Predicted value for x=45: {intercept + slope \* 45}')

print(f'Predicted value for x=46: {intercept + slope \* 46}')

print(f'Predicted value for x=47: {intercept + slope \* 47}')

print(f'Predicted value for x=48: {intercept + slope \* 48}')

print(f'Predicted value for x=49: {intercept + slope \* 49}')

print(f'Predicted value for x=50: {intercept + slope \* 50}')

print(f'Predicted value for x=51: {intercept + slope \* 51}')

print(f'Predicted value for x=52: {intercept + slope \* 52}')

print(f'Predicted value for x=53: {intercept + slope \* 53}')

print(f'Predicted value for x=54: {intercept + slope \* 54}')

print(f'Predicted value for x=55: {intercept + slope \* 55}')

print(f'Predicted value for x=56: {intercept + slope \* 56}')

print(f'Predicted value for x=57: {intercept + slope \* 57}')

print(f'Predicted value for x=58: {intercept + slope \* 58}')

print(f'Predicted value for x=59: {intercept + slope \* 59}')

print(f'Predicted value for x=60: {intercept + slope \* 60}')

print(f'Predicted value for x=61: {intercept + slope \* 61}')

print(f'Predicted value for x=62: {intercept + slope \* 62}')

print(f'Predicted value for x=63: {intercept + slope \* 63}')

print(f'Predicted value for x=64: {intercept + slope \* 64}')

print(f'Predicted value for x=65: {intercept + slope \* 65}')

print(f'Predicted value for x=66: {intercept + slope \* 66}')

print(f'Predicted value for x=67: {intercept + slope \* 67}')

print(f'Predicted value for x=68: {intercept + slope \* 68}')

print(f'Predicted value for x=69: {intercept + slope \* 69}')

print(f'Predicted value for x=70: {intercept + slope \* 70}')

print(f'Predicted value for x=71: {intercept + slope \* 71}')

print(f'Predicted value for x=72: {intercept + slope \* 72}')

print(f'Predicted value for x=73: {intercept + slope \* 73}')

print(f'Predicted value for x=74: {intercept + slope \* 74}')

print(f'Predicted value for x=75: {intercept + slope \* 75}')

print(f'Predicted value for x=76: {intercept + slope \* 76}')

print(f'Predicted value for x=77: {intercept + slope \* 77}')

print(f'Predicted value for x=78: {intercept + slope \* 78}')

print(f'Predicted value for x=79: {intercept + slope \* 79}')

print(f'Predicted value for x=80: {intercept + slope \* 80}')

print(f'Predicted value for x=81: {intercept + slope \* 81}')

print(f'Predicted value for x=82: {intercept + slope \* 82}')

print(f'Predicted value for x=83: {intercept + slope \* 83}')

print(f'Predicted value for x=84: {intercept + slope \* 84}')

print(f'Predicted value for x=85: {intercept + slope \* 85}')

print(f'Predicted value for x=86: {intercept + slope \* 86}')

print(f'Predicted value for x=87: {intercept + slope \* 87}')

print(f'Predicted value for x=88: {intercept + slope \* 88}')

print(f'Predicted value for x=89: {intercept + slope \* 89}')

print(f'Predicted value for x=90: {intercept + slope \* 90}')

print(f'Predicted value for x=91: {intercept + slope \* 91}')

print(f'Predicted value for x=92: {intercept + slope \* 92}')

print(f'Predicted value for x=93: {intercept + slope \* 93}')

print(f'Predicted value for x=94: {intercept + slope \* 94}')

print(f'Predicted value for x=95: {intercept + slope \* 95}')

print(f'Predicted value for x=96: {intercept + slope \* 96}')

print(f'Predicted value for x=97: {intercept + slope \* 97}')

print(f'Predicted value for x=98: {intercept + slope \* 98}')

print(f'Predicted value for x=99: {intercept + slope \* 99}')

print(f'Predicted value for x=100: {intercept + slope \* 100}')

print(f'Predicted value for x=101: {intercept + slope \* 101}')

print(f'Predicted value for x=102: {intercept + slope \* 102}')

print(f'Predicted value for x=103: {intercept + slope \* 103}')

print(f'Predicted value for x=104: {intercept + slope \* 104}')

print(f'Predicted value for x=105: {intercept + slope \* 105}')

print(f'Predicted value for x=106: {intercept + slope \* 106}')

print(f'Predicted value for x=107: {intercept + slope \* 107}')

## 4.5 Definitions – Analysis of Variance (ANOVA)

### Analysis of Variance

**Analysis of Variance (ANOVA)** is a statistical technique used to compare the means of two or more groups to determine if there are statistically significant differences between them. It assesses whether the variations between group means are greater than the variations within the groups.

ANOVA is commonly used when dealing with categorical independent variables and a continuous dependent variable. It helps to answer questions such as:

Are there any significant differences in the means of a dependent variable across different groups?

Which group means are significantly different from each other?

ANOVA compares the variance between the group means (explained variance) with the variance within each group (unexplained variance). If the explained variance is significantly greater than the unexplained variance, it suggests that there are significant differences between the groups.

There are different types of ANOVA, depending on the number of independent variables and the design of the study:

**One-Way ANOVA:** Compares the means of a single dependent variable across multiple groups (e.g., comparing the test scores of students from different schools).

**Two-Way ANOVA:** Considers the effects of two independent variables on a single dependent variable (e.g., studying the impact of both gender and age on exam performance).

**Multivariate ANOVA (MANOVA):** Handles multiple dependent variables simultaneously, allowing the examination of group differences across several outcomes.

## 4.5.1 Code Snippets – Analysis of Variance (ANOVA)

### ANOVA

```
from scipy.stats import f_oneway
# Assuming you have multiple sets of data in a list "data_groups"
# Perform one-way ANOVA to compare means across groups
f_statistic, p_value = f_oneway(*data_groups)
if p_value < 0.05:
    print("Reject null hypothesis: Means are significantly different!")
else:
    print("Fail to reject null hypothesis: Means are not significantly different.")
```

### One-Way ANOVA

```
From statsmodels.formula.api import ols
# Assuming you have multiple sets of data in a list "data_groups"
# Perform one-way ANOVA to compare means across groups
f_statistic, p_value = f_oneway(*data_groups)

if p_value < 0.05:
    print("Reject null hypothesis: Means are significantly different!")
else:
    print("Fail to reject null hypothesis: Means are not significantly different.")
```

### Two-Way Anova

```
import pandas as pd
from statsmodels.formula.api import ols
from statsmodels.formula.concrete import ols
# Assuming you have a Dataframe named "data" with columns "dependent_var", "independent_var1", and "independent_var2"
# Create the ANOVA model using formula notation
formula = "dependent_var ~ independent_var1 + independent_var2 + independent_var1:independent_var2"
model = ols(formula, data).fit()

# Perform two-way ANOVA and extract results
anova_results = model.anova()

# Print the ANOVA table with F-statistic and p-value
print(anova_results)
```

### Multivariate ANOVA

```
import pandas as pd
import statsmodels.api as sm
from statsmodels.multivariate import MANOVA
# Assuming you have a Dataframe named "data" with columns "dependent_var1", "dependent_var2", and "group_var"
# Create a multivariate design matrix
y = data[['dependent_var1', 'dependent_var2']]
values = X * sX.add_constant(data['group_var']).values

# Add a constant column for the intercept
X = sX.add_constant(data['group_var'])

# Perform MANOVA
manova = MANOVA(y, X)

# Print the MANOVA results
print(manova.mv_test())
```

### Optimization

**1. One-Way ANOVA with DataFrames:**

```
import pandas as pd
from scipy.stats import f_oneway
# Assuming you have a Dataframe named "data" with columns "dependent_var" and "group_var"
# Group the data by "group_var" and perform one-way ANOVA
grouped_data = data.groupby('group_var')[['dependent_var']].values
f_statistic, p_value = f_oneway(*grouped_data)
```

**2. One-Way ANOVA with Groupby:**

```
import pandas as pd
from scipy.stats import f_oneway
# Assuming you have a Dataframe named "data" with columns "dependent_var" and "group_var"
# Perform one-way ANOVA using the Groupby function
grouped_data = data.groupby('group_var')['dependent_var'].values
f_statistic, p_value = f_oneway(*grouped_data)
```

**3. Two-Way ANOVA with Statsmodels:**

```
import pandas as pd
import statsmodels.api as sm
from statsmodels.formula.api import ols
# Assuming you have a Dataframe named "data" with columns "dependent_var", "independent_var1", and "independent_var2"
# Create the ANOVA model using formula notation
formula = "dependent_var ~ independent_var1 + independent_var2 + independent_var1:independent_var2"
model = ols(formula, data).fit()

# Print the ANOVA table with F-statistic and p-value
anova_table = sm.stats.anova_lm(model)
print(anova_table)
```

**4. One-Way ANOVA with Post-Hoc Test:**

```
import pandas as pd
from statsmodels.stats.multicomp import pairwise_tukeyhsd
# Assuming you have a Dataframe named "data" with columns "dependent_var" and "group_var"
# Perform one-way ANOVA
pairwise_tukeyhsd(data, 'group_var', 'dependent_var')

# Perform Tukey's HSD post-hoc test for pairwise comparisons
if p_value < 0.05:
    tukey_results = pairwise_tukeyhsd(data['dependent_var'], data['group_var'])
    print(tukey_results)
```

**5. Mixed-Design ANOVA with Pingouin:**

```
import pandas as pd
from pingouin import mixed_anova
# Assuming you have a Dataframe named "data" with columns "dependent_var", "between_group_var", and "within_group_var"
# Perform mixed-Design ANOVA (Between-Within ANOVA)
mixed_results = mixed_anova(data, dv='dependent_var', between='between_group_var', within='within_group_var')
print(mixed_results)
```

## 4.6 Definitions – Probability Distribution

### Probability Distribution

In statistical analysis for data science, a probability distribution is a mathematical function that describes the likelihood of various outcomes or events occurring in a random experiment or process. It provides a systematic way to model and understand the uncertainty associated with random variables. Commonly used probability distributions in data science include:

1. Normal Distribution (Gaussian Distribution): Often used for continuous data and is characterized by its bell-shaped curve.

2. Binomial Distribution: Applicable for discrete data with two possible outcomes, like success/failure.

3. Poisson Distribution: Suitable for discrete data that represents the number of events in a fixed interval.

4. Exponential Distribution: Useful for modeling time-to-event data, such as the time between two events occurring.

5. Uniform Distribution: Represents data where all outcomes have equal probabilities.

6. Chi-Squared Distribution: Used in hypothesis testing and confidence interval calculations.

7. Student's t-Distribution: Similar to the normal distribution but used when the sample size is small or the population variance is unknown.

## 4.6.1 Code Snippets – Probability Distribution

### Probability Distribution

```
from scipy.stats import norm
# Assuming you want to calculate probabilities for a normal distribution
# Calculate probability density function (PDF) for a specific value
pdf_value = norm.pdf(x, loc=mean, scale=sdev)
# Calculate cumulative distribution function (CDF) for a specific value
cdf_value = norm.cdf(x, loc=mean, scale=sdev)
```

### Normal Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import norm
# Generate random data from a normal distribution with mean=0 and standard deviation=1
data = np.random.normal(loc=0, scale=1, size=1000)
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='orange')
```

```
# Plot the probability density function (PDF) for the standard normal distribution
x = np.linspace(-1, 1, 1000)
plt.plot(x, norm.pdf(x, 0, 1), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Normal Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### Binomial Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import binom
# Generate random data from a binomial distribution with n=10 and p=0.5
data = np.random.binom(10, 0.5, 1000)
# Plot the histogram of the generated data
plt.hist(data, bins=10, range=(-0.5, 10.5), density=True, alpha=0.6, color='blue')
```

```
# Plot the probability mass function (PMF) for the binomial distribution with n=10 and p=0.5
x = np.arange(0, 11, 1)
plt.plot(x, binom.pmf(x, 10, 0.5), 'r', lw=2, marker='o')
```

```
plt.xlabel('Number of Successes')
plt.ylabel('Probability Mass')
```

```
plt.title('Binomial Distribution')
```

```
plt.legend(['PMF', 'Data'])
```

```
plt.show()
```

### Poisson Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import poisson
# Generate random data from a Poisson distribution with lambda=2
data = np.random.poisson(2, size=1000)
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='blue')
```

```
# Plot the probability mass function (PMF) for the Poisson distribution with lambda=2
x = np.arange(0, 15, 1)
plt.plot(x, poisson.pmf(x, 2), 'r', lw=2, marker='o')
```

```
plt.xlabel('Number of Events')
plt.ylabel('Probability Mass')
```

```
plt.title('Poisson Distribution')
```

```
plt.legend(['PMF', 'Data'])
```

```
plt.show()
```

### Exponential Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import expon
# Generate random data from an exponential distribution with scale parameter=2
data = np.random.exponential(scale=2, size=1000)
```

```
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='blue')
```

```
# Plot the probability density function (PDF) for the exponential distribution with scale parameter=2
x = np.linspace(0, 15, 1000)
plt.plot(x, expon.pdf(x, scale=2), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Exponential Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### Uniform Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import uniform
```

```
# Generate random data from a uniform distribution between 0 and 1
data = np.random.uniform(0, 1, size=1000)
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='blue')
```

```
# Plot the probability density function (PDF) for the uniform distribution between 0 and 1
x = np.linspace(0, 1, 1000)
plt.plot(x, uniform.pdf(x, 0, 1), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Uniform Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### Chi-Squared Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import chi2
```

```
# Generate random data from a chi-squared distribution with degrees of freedom=3
data = np.random.chisquare(df=3, size=1000)
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='purple')
```

```
# Plot the probability density function (PMF) for the chi-squared distribution with degrees of freedom=3
x = np.linspace(0, 20, 1000)
plt.plot(x, chi2.pdf(x, df=3), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Chi-Squared Distribution')
```

```
plt.legend(['PMF', 'Data'])
```

```
plt.show()
```

### Student's t-Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import t
```

```
# Generate random data from a t-distribution with degrees of freedom=10
data = np.random.standard_t(df=10, size=1000)
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='orange')
```

```
# Plot the probability density function (PDF) for the t-distribution with degrees of freedom=10
x = np.linspace(-5, 5, 1000)
plt.plot(x, t.pdf(x, df=10), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Student's t-Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### Log-Normal Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import lognorm
```

```
# Generate random data from a log-normal distribution with degrees of freedom=10
data = np.random.lognorm(mean=1, df=10, size=1000)
```

```
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='orange')
```

```
# Plot the probability density function (PDF) for the log-normal distribution with degrees of freedom=10
x = np.linspace(0, 10, 1000)
plt.plot(x, lognorm.pdf(x, 1, df=10), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Student's t-Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### Beta Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import beta
```

```
# Generate random data from a beta distribution with shape parameters alpha=2 and beta=5
data = np.random.beta(alpha=2, beta=5, size=1000)
```

```
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='grey')
```

```
# Plot the probability density function (PDF) for the beta distribution with shape parameters alpha=2 and beta=5
x = np.linspace(0, 1, 1000)
plt.plot(x, beta.pdf(x, alpha=2, beta=5), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Beta Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### Gamma Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import gamma
```

```
# Generate random data from a gamma distribution with shape=2 and scale=1
data = np.random.gamma(shape=2, scale=1, size=1000)
```

```
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='blue')
```

```
# Plot the probability density function (PDF) for the gamma distribution with shape=2 and scale=1
x = np.linspace(0, 15, 1000)
plt.plot(x, gamma.pdf(x, 2, 1), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Gamma Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### Poisson Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import poisson
```

```
# Generate random data from a Poisson distribution with shape=2 and scale=2
data = np.random.poisson(2, size=1000)
```

```
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='blue')
```

```
# Plot the probability density function (PDF) for the Poisson distribution with shape=2 and scale=2
x = np.linspace(0, 15, 1000)
plt.plot(x, poisson.pdf(x, 2, 2), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Poisson Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### Weibull Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import weibull
```

```
# Generate random data from a weibull distribution with shape=2 and scale=2
data = np.random.weibull(2, size=1000)
```

```
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='green')
```

```
# Plot the probability density function (PDF) for the Weibull distribution with shape=2 and scale=2
x = np.linspace(0, 5, 1000)
plt.plot(x, weibull.pdf(x, 2, 2), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Weibull Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### F - Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import f
```

```
# Generate random data from an F-distribution with degrees of freedom 2 and 10
data = np.random.f(2, 10, size=1000)
```

```
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='purple')
```

```
# Plot the probability density function (PDF) for the F-distribution with degrees of freedom 2 and 10
x = np.linspace(0, S, 1000)
plt.plot(x, f.pdf(x, dfn=2, dfd=10), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('F-Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### Pareto Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import pareto
```

```
# Generate random data from a Pareto distribution with shape=3
data = np.random.pareto(3, size=1000)
```

```
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='cyan')
```

```
# Plot the probability density function (PDF) for the Pareto distribution with shape=3
x = np.linspace(0, 10, 1000)
plt.plot(x, pareto.pdf(x, b=3), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Pareto Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### Logistic Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import logistic
```

```
# Generate random data from a logistic distribution with location=0 and scale=1
data = np.random.logistic(loc=0, scale=1, size=1000)
```

```
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='pink')
```

```
# Plot the probability density function (PDF) for the logistic distribution with location=0 and scale=1
x = np.linspace(0, 10, 1000)
plt.plot(x, logistic.pdf(x, loc=0, scale=1), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Logistic Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```

### Laplace Distribution

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import laplace
```

```
# Generate random data from a Laplace distribution with location=0 and scale=1
data = np.random.laplace(loc=0, scale=1, size=1000)
```

```
# Plot the histogram of the generated data
plt.hist(data, bins=30, density=True, alpha=0.6, color='orange')
```

```
# Plot the probability density function (PDF) for the Laplace distribution with location=0 and scale=1
x = np.linspace(-5, 5, 1000)
plt.plot(x, laplace.pdf(x, loc=0, scale=1), 'r', lw=2)
```

```
plt.xlabel('Value')
plt.ylabel('Probability Density')
```

```
plt.title('Laplace Distribution')
```

```
plt.legend(['PDF', 'Data'])
```

```
plt.show()
```



## 5.4 Definitions – Time Series Modeling (ARIMA)

### Handling Date-Time Data

ARIMA (AutoRegressive Integrated Moving Average) is a time series modeling technique that uses the following notation:

- $y(t) = C + \phi_1y(t-1) + \phi_2y(t-2) + \dots + \phi_py(t-p) + \theta_1e(t-1) + \theta_2e(t-2) + \dots + \theta_qe(t-q) + \epsilon(t)$
- $\epsilon(t)$  represents the value of the time series at time  $t$ .
- $C$  is the constant term or the intercept.
- $\phi_1, \phi_2, \dots, \phi_p$  are the AR coefficients.
- $\theta_1, \theta_2, \dots, \theta_q$  are the lagged values of the error terms.
- $e(t), e(t-1), \dots, e(t-q)$  are the moving average (MA) coefficients.
- $\epsilon(t), \epsilon(t-1), \dots, \epsilon(t-q)$  are the residual errors or the noise terms.
- $\epsilon(t)$  is the current residual error or noise term.

The equation combines the autoregressive (AR) terms, moving average (MA) terms, and the constant term to model the dependencies and patterns in the time series. The parameters  $\phi$  and  $\theta$  are estimated from the historical data using techniques like maximum likelihood estimation.

## 5.4.1 Code Snippets – Time Series Modeling

### Time Series Modeling

```
import numpy as np
import matplotlib.pyplot as plt
from statsmodels.tsa.arima.model import ARIMA

# Generate synthetic time series data
np.random.seed(42)
data = np.random.randn(1000)
# Fit an AR(1) model
model = ARIMA(data, order=(1, 1, 1))
model.fit()
# Forecast future values
forecast_steps = 10
forecast, _, conf_int = model.forecast(steps=forecast_steps)

# Plot the original data and the forecast
plt.plot(data, label='Original Data', color='blue')
plt.plot(arima_forecast, label='Forecast', steps=forecast_steps, color='red')
plt.title('ARIMA Forecasting')
plt.xlabel('Value')
plt.ylabel('Value')
plt.title('ARIMA Forecast')
plt.legend()
plt.show()
```

### Grid Search for Hyperparameter Tuning

```
import itertools
import pandas as pd
import numpy as np
from statsmodels.tsa.arima.model import ARIMA

# Define the range of values for p, d, q
q_values = range(0, 5) # Example: 0, 1, 2
d_values = range(0, 3)
p_values = range(0, 3)

# Create a list of all possible combinations of p, d, and q values
param_grid = list(itertools.product(p_values, q_values, d_values))

# Load the time series data
data = pd.read_csv('time_series_data.csv')
# Fit the ARIMA model
best_params = {}
for param in param_grid:
    model = ARIMA(data, order=param)
    results = model.fit()
    # Calculate AIC (Akaike Information Criterion)
    aic = results.aic
    if aic < best_params.get('aic'):
        best_params['aic'] = aic
        best_params['params'] = param
    else:
        continue

# Print the best parameters
print("Best Parameters: ", best_params)
```

### Rolling Window Cross-Validation

```
import pandas as pd
from statsmodels.tsa.arima.model import ARIMA
from sklearn.metrics import mean_squared_error

# Load the time series data
data = pd.read_csv('time_series_data.csv')

# Split the data into train and test sets
train_size = int(len(data) * 0.8)
train_data, test_data = data[:train_size], data[train_size:]

# Fit the ARIMA model on the training data
model = ARIMA(train_data, order=(1, 1, 1))
results = model.fit()

# Perform rolling window cross-validation
window_size = 7
mse_scores = []

for i in range(len(test_data) - window_size):
    window = test_data[i:window_size]
    # Make predictions on the window
    predictions = results.predict(start=i+window, end=i+window+6)
    # Calculate the mean squared error
    mse = mean_squared_error(window, predictions)
    mse_scores.append(mse)

# Calculate the average mean squared error
avg_mse = sum(mse_scores) / len(mse_scores)
print("Average MSE: ", avg_mse)
```

### Model Selection using Information Criteria

```
import pandas as pd
import numpy as np
from statsmodels.tsa.arima.model import ARIMA

# Load the time series data
data = pd.read_csv('time_series_data.csv')

# Split the data into train and validation sets
train_size = int(len(data) * 0.8)
train_data, val_data = data[:train_size], data[train_size:]

# Define the range of values for p, d, q
q_values = range(0, 5) # Example: 0, 1, 2
d_values = range(0, 3)
p_values = range(0, 3)

best_params = float('inf')
best_params = None

# Iterate over all parameter combinations
for p in p_values:
    for q in q_values:
        for d in d_values:
            # Fit the ARIMA model
            model = ARIMA(train_data, order=(p, d, q))
            results = model.fit()

            # Calculate AIC (Akaike Information Criterion)
            aic = results.aic

            # Update the best parameters if AIC is lower
            if aic < best_params['aic']:
                best_params['aic'] = aic
                best_params['params'] = (p, d, q)
            else:
                continue

# Fit the ARIMA model on the combined train and validation data using the best parameters
model = ARIMA(data, order=best_params['params'])
results = model.fit()

# Forecast future values
forecast_steps = 10
forecast, _, conf_int = results.forecast(steps=forecast_steps)

# Print the best parameters
print("Best Parameters: ", best_params)
```

## 5.5 Definitions – Decision Tree

### Decision Tree

The decision tree algorithm uses a flowchart-like structure to make decisions or predict values based on input features. The decision-making process involves splitting the data at internal nodes based on different features, and assigning outcomes or predicted values at the leaf nodes. The algorithm uses specific criteria to determine the best feature and the best split at each node.

For classification tasks, the decision tree algorithm uses criteria such as Gini impurity or entropy to measure the purity of each split and select the feature that provides the most information gain.

For regression tasks, the decision tree algorithm uses criteria such as variance reduction or sum of squared errors to measure the quality of each split and choose the feature that minimizes the variance or the error.

The final decision tree structure is created through a recursive process of selecting the best features and splits until a stopping criterion is met, such as reaching a maximum depth or having a minimum number of instances in each leaf.

## 5.5.1 Code Snippets – Decision Tree

### Decision Tree

Import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load\_iris
from sklearn.tree import DecisionTreeClassifier, plot\_tree

# Load the iris dataset
data = load\_iris()
X, y = data.data, data.target

# Create a decision tree classifier
model = DecisionTreeClassifier()
model.fit(X, y)

# Visualize the decision tree
plt.figure(figsize=(10, 8))
plot\_tree(model, feature\_names=iris.feature\_names, class\_names=iris.target\_names, filled=True)
plt.show()

## 5.6 Definitions – K-Means Clustering

### K-Means Clustering

K-Means Clustering is an Unsupervised Learning algorithm, which groups the unlabeled dataset into distinct clusters. K-Means is one of the most popular clustering algorithms and is used to be created in the process of K=2, there will be two clusters, and for K=3, there will be three clusters, and so on. The goal of K-Means clustering is to divide the data points into K clusters in such a way that the within-cluster variation is minimized and the between-cluster variation is maximized.

It allows us to cluster the data into different groups and a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training. It is a centroid-based algorithm, where each cluster is associated with a centroid.

The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters.

## 5.6.1 Code Snippets – K-Means Clustering

### K-Means Clustering

Import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import blobs
from sklearn.cluster import KMeans

# Generate synthetic data for clustering
X, \_ = make\_blobs(n\_samples=300, centers=4, random\_state=42)
# Train the K-Means model
model = KMeans(n\_clusters=4)
model.fit(X)

# Visualize the clustered data
plt.scatter(X[:, 0], X[:, 1], c=model.labels\_, cmap='viridis', s=50, alpha=0.7)
plt.scatter(model.cluster\_centers[:, 0], model.cluster\_centers[:, 1], c='red', marker='x', s=100)
plt.xlabel('Feature 1')
plt.ylabel('Feature 2')
plt.title('K-Means Clustering')
plt.show()

## 5.7 Definitions – Principal Component Analysis(PCA)

### PCA

PCA, or Principal Component Analysis, is a technique used to reduce the dimensionality of a dataset while retaining most of the information. It achieves this by transforming the original features into a new set of uncorrelated variables called principal components. These components are ranked based on the amount of variance they explain in the data. By selecting the top components, PCA allows for data visualization, noise reduction, and improved computational efficiency in machine learning algorithms.

PCA (Principal Component Analysis):

- Dimensionality Reduction Technique
- Capturing Data Variance
- Utilizing Eigenvectors for Transformation

## 5.7.1 Code Snippets – Linear Regression

### Linear Regression

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.decomposition import PCA

# Load the iris dataset
data = load_iris()
X, y = data.data, data.target

# Apply PCA for dimensionality reduction to 2 components
X_pca = PCA(n_components=2)
X_pca.fit(X)

# Plot the first two principal components
plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y, cmap='plasma')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.title('PCA - Dimensionality Reduction')
plt.show()
```

## 5.8 Definitions – Support Vector Machines(SVM)

### Support Vector Machine

Support Vector Machine (SVM) is a supervised machine learning algorithm used for classification and regression tasks. It is particularly effective in solving complex classification problems with a clear margin of separation between classes.

For SVM classification, the decision boundary equation for a linearly separable case can be expressed as:

$$w^T x + b = 0$$

where:

- $w$  is the weight vector perpendicular to the decision boundary
- $x$  is the input feature vector
- $b$  is the bias term

The predicted class label for a new sample  $x$  can be determined by evaluating the sign of the decision function:

$$f(x) = \text{sign}(w^T x + b)$$

For non-linear SVM classification, the decision boundary equation is transformed using a kernel function, such as the radial basis function (RBF) kernel:

$$f(x) = \text{sign}(\alpha_i \phi(x_i) \cdot \phi(x) + b)$$

where:

- $\alpha_i$  are the Lagrange multipliers (support vectors)
- $\phi$  is the non-linear kernel function that computes the similarity between the support vectors and the new sample
- $x$  is the new sample
- $b$  is the bias term

## 5.8.1 Code Snippets – Support Vector Machines(SVM)

### SVM

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.svm import SVC
from sklearn.svm import LinearSVC
from sklearn.svm import SVR
# Generate synthetic data for binary classification
X, y = make_classification(n_samples=100, n_features=2, n_informative=2, n_redundant=0, random_state=42)
# Train a Support Vector Machine (SVM) classifier
model = SVC()
model.fit(X, y)
# Plot the data and the decision boundary
plt.scatter(X[:, 0], X[:, 1], c=y, cmap='plasma')
plt.xlabel('Feature 1')
plt.ylabel('Feature 2')
# Plot the decision boundary
xx, yy = np.meshgrid(np.linspace(-5, 5, 30), np.linspace(-5, 5, 30))
Z = model.decision_function(np.c_[xx.ravel(), yy.ravel()])
Z = Z.reshape(xx.shape)
plt.contourf(xx, yy, Z, colors='c', levels=[-1, 0, 1], linestyles='--', alpha=0.5)
plt.title('Support Vector Machine (SVM)')
plt.show()
```

## 5.9 Definitions – Naive Bayes Classifier

### Naive Bayes Classifier

Naive Bayes is a probabilistic classifier that is based on Bayes' theorem with the assumption of independence between features. It is commonly used for classification tasks and is particularly suited for text classification and spam filtering.

Bayes theorem is a formula that offers a conditional probability of an event  $A$  taking happened given another event  $B$  has previously happened. Its mathematical formula is as follows:

Where:

$$P(A|B) = P(A) \cdot P(B|A) / P(B) = P(A) / P(B)$$

A and B are two events

$P(A|B)$  is the probability of event A provided event B has already happened.

$P(A)$  is the independent probability of A

$P(B)$  is the independent probability of B

## 5.9.1 Code Snippets – Naive Bayes Classifier

### Naive Bayes Classifier

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier, plot_tree
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
```

# Load the iris dataset
data = load\_iris()
X, y = data.data, data.target

# Split the data into training and testing sets
X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=42)

# Create a Gaussian Naive Bayes classifier
model = GaussianNB()
model.fit(X\_train, y\_train)

# Make predictions on the test data
y\_pred = model.predict(X\_test)

# Calculate accuracy
accuracy = accuracy\_score(y\_test, y\_pred)
print("accuracy", accuracy)

## 5.10 Definitions – Random Forest Classifier

### Random Forest Classifier

Random Forest Classifier is an ensemble learning method that combines multiple decision trees to create a powerful classification model. It is a type of supervised learning algorithm that is widely used for both classification and regression tasks.

In Random Forest, a collection of decision trees is created, where each tree is trained on a random subset of the training data and a random subset of input features. During the training process, each tree independently makes predictions, and the final prediction is determined by combining the predictions from all the trees through voting (for classification) or averaging (for regression).

Random Forest works in two-phase first is to create the random forest by combining N decision tree, and second is to make predictions for each tree created in the first phase.

The Working process can be explained in the below steps and diagram:

Step-1: Select random N data points from the training set.

Step-2: Build the decision trees associated with the selected data points (Subsets).

Step-3: Assign the N for decision trees that you want to build.

Step-4: Repeat Step 1 & 2.

Step-5: For new data points, find the predictions of each decision tree, and assign the new data points to the category that wins the majority votes

## 5.10.1 Code Snippets – Random Forest Classifier

### Random Forest Classifier

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_classification
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
```

# Generate synthetic data for binary classification
X, y = make\_classification(n\_samples=100, n\_features=2, n\_informative=2, n\_redundant=0, random\_state=42)

# Split the data into training and testing sets
X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=42)

# Create a Random Forest Classifier
model = RandomForestClassifier(n\_estimators=100, random\_state=42)
model.fit(X\_train, y\_train)

# Make predictions on the test data
y\_pred = model.predict(X\_test)

# Calculate accuracy
accuracy = accuracy\_score(y\_test, y\_pred)
print("accuracy", accuracy)

## 5.11 Definitions – Poisson Regression

### Poisson Regression

Poisson Regression is a statistical model used for analyzing count data. It models the relationship between predictor variables and the expected counts of an event. It assumes that the response variable follows a Poisson distribution.

The equation for Poisson Regression can be represented as:

$$\log(E(Y)) = B_0 + B_1 X_1 + B_2 X_2 + \dots + B_n X_n$$

Where:

- $E(Y)$  is the natural logarithm of the expected count of the event (response variable Y).
- $B_0, B_1, B_2, \dots, B_n$  are the regression coefficients associated with the predictor variables  $X_1, X_2, \dots, X_n$ .
- $X_1, X_2, \dots, X_n$  are the predictor variables.

The model assumes that the log of the expected count is a linear combination of the predictor variables, which allows for modeling the relationship between the predictors and the event count. The coefficients  $B_0, B_1, B_2, \dots, B_n$  represent the effects of the predictor variables on the log of the expected count.

## 5.11.1 Code Snippets – Poisson Regression

### Poisson Regression

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import poisson
import statsmodels.api as sm
```

# Generate synthetic count data following a Poisson distribution
np.random.seed(42)
data = np.random.poisson(3, size=100)

# Fit a Poisson regression model
model = PoissonRegressor()
model.fit(data)

# Print summary of the Poisson regression
print(result.summary())

## 5.12 Definitions – Gradient Boosting

Gradient Boosting is a machine learning technique that combines multiple weak prediction models to create a strong predictive model. It sequentially builds models, focusing on the errors of the previous models to improve predictions.

## 5.12.1 Code Snippets – Gradient Boosting

### Gradient Boosting

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_boston
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
```

# Load the Boston Housing dataset
data = load\_boston()

X, y = data.data, data.target

# Split the data into training and testing sets
X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=42)

# Create a Gradient Boosting regression model
model = GradientBoostingRegressor(n\_estimators=100, learning\_rate=0.1, random\_state=42)
model.fit(X\_train, y\_train)

# Make predictions on the test data
y\_pred = model.predict(X\_test)

# Calculate error
use = mean\_squared\_error(y\_test, y\_pred)

print("Mean Squared Error:", use)



