

**A**  
**Minor Project**  
**On**  
**Wine Quality Prediction**  
*Submitted in partial fulfilment of the requirements*  
*for the award of the degree of*

**Bachelor of Computer Applications**  
**To**  
I.T.S College, CCS University,



**Guide:** Mr Prateek Gupta

**SUBMITTED BY:** Vaibhav Gautam

**Institute of Technology & Science,**  
**Ghaziabad – 201001**

## **Acknowledgement**

I am profoundly thankful to everyone who contributed to the successful completion of my summer training project.

Firstly, I would like to extend my deepest gratitude to Mr. Prateek Gupta, whose invaluable guidance and insights were essential throughout the training. His expertise and patience were crucial in helping me navigate the complexities of advanced data science and machine learning.

I also wish to sincerely thank ITS Mohan Nagar for collaborating with Shape My Skills Pvt. Ltd. to provide this exceptional learning opportunity. Special appreciation goes to our course coordinator, Mr. Neeraj sir, for his dedicated efforts and support during the training. Furthermore, I am deeply grateful to Head of the Department, for his encouragement and for providing the necessary resources and environment to facilitate this learning experience.

Finally, I acknowledge the unwavering support of my family and friends, who have been my pillars of strength and encouragement throughout this journey.

Thank you all for making this experience valuable and memorable.

Sincerely,

Vaibhav Gautam

# Certificate

This is to certify that our team has successfully completed the project titled " Wine Quality Prediction" as part of the Machine Learning and Data Science summer training program organized by ITS Mohan Nagar in collaboration with ShapeMySkills Pvt. Ltd.

This project was conducted under the esteemed guidance of Mr. Prateek Gupta, whose expertise and mentorship were instrumental in its successful completion. The project exemplifies a thorough understanding of machine learning and data science techniques, highlighting the skills acquired during the training program.

**Coordinator:** Mr. Ravi Govil

Supervised by: Mr Prateek Gupta

**Date:**

**Signature:**

## List of Figures

S No.	Name	Page
1	<b>Figure 1:</b> Scatter plot of pH vs. quality score.	15
2	<b>Figure 2:</b> Bar chart showing average quality scores across sugar levels.	15
3	<b>Figure 3:</b> Heatmap representing feature correlations.	16
4	<b>Figure 4:</b> Boxplot of alcohol content distribution across quality categories.	16
5	<b>Figure 5:</b> Histogram of quality scores distribution.	17
6	<b>Figure 6:</b> Regression plot showing alcohol content's impact on quality scores.	18
7	<b>Figure 7:</b> Pairplot visualizing relationships between multiple features.	19
8	<b>Figure 8:</b> Decision Tree structure for drink quality prediction.	20
9	<b>Figure 9:</b> Confusion Matrix for model evaluation.	20
10	<b>Figure 10:</b> Feature importance visualization for Random Forest algorithm.	21

**List of Tables**

S No.	Name	Page
1	Table 1: Table 1 Summary of Drink Features	12
2	Table 2: Table 2 Quality Scores Distribution	12
3	Table 3: Table 3 Feature Correlations	13
4	Table 4: Table 4 Model Performance Metrics	32

## List of Abbreviations

<b>S No.</b>	<b>Name</b>	<b>Page</b>
1	OOP: Object Oriented Programming	1
2	I/O: Input/Output	1
3	ML: Machine Learning	2
4	AI: Artificial Intelligence	2
5	NumPy: Numerical Python	4
6	Pandas: Panel Data	4
7	CSV: Comma-Separated Value	5
8	SQL: Structured Query Language	5
9	JSON: JavaScript Object Notation	5
10	3D: 3 Dimensional	5
11	SciPy: Scientific Python	6
12	AUC-ROC: Area Under the Receiver Operating Characteristics Curve	8
13	KNN: K-Nearest Neighbor	9

# INDEX

<b>S No.</b>	<b>Title</b>	<b>Page</b>
1	Chapter 1 Introduction to Python and Machine Learning	1-3
2	Chapter 2 Introduction to Python Libraries	4-6
3	Chapter 3 Introduction to Machine Learning Algorithm	7-9
4	Chapter 4 Python and Machine Learning Code	10-25
5	Chapter 5 Conclusion and Result	26
6	Chapter 6 Future Scope	27

# CHAPTER 1

## Introduction to Python and Machine Learning

### Introduction to python :

Python is a general-purpose, dynamically typed, high-level, compiled and interpreted, garbage-collected, and purely object-oriented programming language that supports procedural, object-oriented, and functional programming.

It was Created by Guido van Rossum and first released in 1991, Python emphasizes code readability and allows programmers to express concepts in fewer lines of code compared to languages like C++ or Java.

### Why learn Python?

- **Easy to use and Learn:** Python has a simple and easy-to-understand syntax, unlike traditional languages like C, C++, Java, etc., making it easy for beginners to learn.
- **Interpreted Language:** Python does not require compilation, allowing rapid development and testing. It uses Interpreter instead of Compiler.
- **Object-Oriented Language:** It supports object oriented programming i.e(inheritance,encapsulation,polymorphism,abstraction) making writing reusable and modular code easy.
- **Extensive Libraries :** Python has a rich ecosystem of libraries and frameworks, such as NumPy, Pandas, and Matplotlib, which simplify tasks like data manipulation and visualization.

### Python Popular Frameworks and Libraries

- **Mathematics** - NumPy, Pandas, etc.
- **REST framework:** a toolkit for building RESTful APIs
- **MachineLearning** – Numpy, Seaborn, Matplotlib etc.

### Where is Python used?

- **Data Science:** Python is important in this field because it is easy to use and has powerful tools for data analysis and visualization like NumPy, Pandas, and Matplotlib.
- **Machine Learning:** Python is widely used for machine learning due to its simplicity, ease of use, and availability of powerful machine learning libraries.



## **Introduction to Machine Learning**

Machine learning (ML) is a subfield of artificial intelligence (AI) that involves the development of algorithms and statistical models enabling computers to perform tasks without explicit instructions. Instead, these systems learn patterns and make decisions based on data. Machine learning is transforming various industries by automating complex processes, providing insights from large datasets, and creating new opportunities for innovation.

## **Definition and Scope**

Machine learning leverages computational methods to improve performance on a given task over time with experience. This process involves:

1. **Data Collection:** Gathering large and diverse datasets.
2. **Data Preprocessing:** Cleaning and formatting data to be suitable for analysis.
3. **Model Selection:** Choosing an appropriate algorithm or model based on the task.
4. **Training:** Feeding the data into the model to learn patterns.
5. **Evaluation:** Assessing the model's performance using metrics and validation techniques.
6. **Deployment:** Implementing the model in real-world applications.
7. **Maintenance:** Continuously updating and refining the model as new data becomes available.

## **Types of Machine Learning**

Machine learning techniques can be broadly categorized into three types:

1. **Supervised Learning:** The model is trained on a labeled dataset, meaning that each training example is paired with an output label. Common algorithms include:
  - Linear Regression
  - Decision Trees
  - Support Vector Machines (SVM)
  - Neural Networks
2. **Unsupervised Learning:** The model is provided with unlabeled data and must find inherent patterns or groupings. Common algorithms include:
  - Clustering (e.g., K-Means, Hierarchical Clustering)
  - Association Rules (e.g., Apriori, Eclat)
  - Principal Component Analysis (PCA)
3. **Reinforcement Learning:** The model learns by interacting with an environment, receiving rewards or penalties based on its actions, and aims to maximize cumulative rewards. Key concepts include:
  - Markov Decision Processes (MDP)
  - Q-Learning
  - Deep Q-Networks (DQN)

# Chapter 2

## Libraries of Python

### NumPy

#### Introduction

NumPy, short for Numerical Python, is a fundamental package for scientific computing with Python. It provides support for arrays, matrices, and many mathematical functions to operate on these data structures.

Some of the most popular libraries:

#### 1. NumPy

**NumPy** originally stands for Numerical Python, is a core package for numerical computing in Python. It supports massive, multidimensional arrays and matrices, as well as a set of mathematical functions for effectively manipulating these arrays.

Features of NumPy:

- NumPy includes an extremely efficient multi-dimensional array object called `numpy.ndarray`, which can store and manage big datasets rapidly.
- NumPy includes a large number of numerical computing tools and methods that can operate on these arrays.
- NumPy has routines for manipulating arrays such as reshaping (`reshape()`), stacking (`stack()`, `hstack()`, `vstack()`), splitting (`split()`), indexing (indexing and slicing), and sorting.

#### ❖ Advantages of NumPy

##### 1. Performance

- **Speed:** Faster than Python lists due to optimized C code.
- **Vectorization:** Allows element-wise operations without loops.

##### 2. Memory Efficiency

- **Contiguous Allocation:** Enhances cache efficiency.
- **Homogeneous Types:** Consistent memory usage.

## 2. Pandas

Pandas is a flexible and advanced Python toolkit for data manipulation and analysis. It includes high-level data structures like Data Frame and Series, which make it easier to work with organized data.

### **Key Features of Pandas:**

- **Data Frame:** A two-dimensional labelled data structure containing columns of various categories. It looks like a spreadsheet or a SQL table and is ideal for working with tabular data.
- **Series:** A one-dimensional labelled array that can carry data of any type (integer, float, text, etc.). Series function similarly to columns in a Data Frame or named arrays.
- Pandas can read and write data from a variety of file formats, including **CSV**, **Excel**, **SQL** databases, and **JSON**.
- Pandas includes sophisticated indexing techniques (loc and iloc) for picking subsets of data. This includes selecting rows and columns based on labels (loc) or integer positions (iloc), giving users easy access to data.

### ❖ **Advantages of Pandas**

#### **1. Data Manipulation and Analysis**

- **DataFrames:** Efficiently handle tabular data with labeled axes (rows and columns).
- **Series:** Simplify manipulation of one-dimensional labeled arrays.

#### **2. Data Cleaning**

- **Handling Missing Data:** Functions for detecting, filling, and removing missing values.
- **Data Transformation:** Easy methods for merging, reshaping, and transforming datasets.

#### **3. Data Selection**

- **Indexing and Slicing:** Powerful, flexible, and intuitive data selection capabilities.
- **Label-based and Position-based Indexing:** Access data using labels or positions.

### 3. Matplotlib

Matplotlib is a robust Python package that allows you to create static, animated, and interactive visualizations. It is commonly used for data visualization jobs and offers a versatile framework for creating plots and figures in a variety of formats.

#### Key Features of Matplotlib:

- Matplotlib provides a wide range of graphs, including line plots, scatter plots, bar plots, histogram plots, pie charts, 3D plots, and more.
- Matplotlib works seamlessly with Pandas and NumPy, enabling direct plotting from Data Frame and Series objects.
- Matplotlib works well with Seaborn, a statistical data visualization toolkit in Python. Seaborn expands Matplotlib's capabilities by providing higher-level functions for statistical plots such as violin plots, box plots, and regression plots.

### ❖ Advantages of Matplotlib

#### 1. Versatility

- **Wide Range of Plots:** Supports various types of plots such as line, bar, scatter, histogram, pie, and more.
- **Customizable:** Highly customizable plots, allowing for detailed adjustments to suit specific needs.

#### 2. Integration

- **Seamless with NumPy and Pandas:** Easily integrates with NumPy and Pandas for plotting data from arrays and DataFrames.
- **Compatible with Other Libraries:** Works well with other libraries like SciPy and scikit-learn for enhanced functionality.

#### 3. Publication Quality

- **High-Quality Output:** Produces high-quality figures suitable for publication.
- **Multiple Formats:** Exports plots in various formats including PNG, PDF, SVG, and EPS.

## 4. Seaborn

Seaborn is a Python data visualization package built on Matplotlib. It provides a high-level interface for constructing visually appealing and useful statistical graphs.

### Key Features of Seaborn:

- Seaborn provides a straightforward and intuitive API for constructing complicated statistical graphs, using less lines of code than Matplotlib.
- Seaborn includes built-in color palettes to improve visualizations. It provides both qualitative (categorical data), sequential (numeric data), and divergent (data with a crucial midpoint) color schemes.
- Seaborn works perfectly with Pandas Data Frames, allowing for direct visualization of data contained in Data Frame objects.

### ❖ Advantages of Seaborn

#### 1. High-Level Interface

- **Ease of Use:** Simplifies complex visualizations with fewer lines of code.
- **Intuitive API:** Provides a user-friendly syntax for creating informative and attractive statistical graphics.

#### 2. Statistical Visualization

- **Built-in Support:** Directly integrates with Pandas DataFrames and handles statistical aggregations and visualizations effortlessly.
- **Advanced Plotting Functions:** Includes specialized plots like categorical plots, distribution plots, and regression plots.

## **Chapter 3**

### **Introduction to Machine Learning Algorithm**

#### **Introduction to Categorical**

A categorical variable, also known as a discrete variable, is a type of variable that can take on one of a limited, fixed number of possible values, representing distinct categories or classes.

Examples include:

Binary categories: Yes/No, True/False, 0/1

Multiclass categories: Red/Green/Blue, Dog/Cat/Horse, Low/Medium/High

#### **Goal:**

The primary objective of categorical machine learning algorithms is to predict the category or class of new instances based on learned patterns from a labelled dataset. This process is known as classification.

#### **Key Characteristics of Categorical Data**

- Categorical data has discrete values.
- Non-ordinal vs. ordinal: Some categorical data is non-ordinal (e.g., fruit types), whilst others are ordinal (e.g., rating scales such as "low," "medium," and "high").
- Encoding Required: Many machine learning techniques require categorical data to be translated into numerical format before they can be used.

#### **Types of Algorithms**

Categorical algorithms are built specifically for processing and analysing categorical data, which is made up of variables that indicate discrete categories or groups. Several types of categorical algorithms have been created to handle the specific issues caused by categorical data, ensuring that machine learning models execute accurately and efficiently.

### ❖ **Logistic Regression**

- Logistic regression is often utilized in binary classification situations. It calculates the likelihood that a given input belongs to a particular category.
- It represents the relationship between a binary dependent variable and one or more independent variables.
- Uses a logistic function (sigmoid function) to convert the relationship to a probability value between 0 and 1.
- Assumes that the independent variables have a linear connection with the dependent variable's log chances.
- Accuracy, precision, recall, F1-score, and area under the ROC curve (AUC-ROC) are common metrics used for evaluation.

### ❖ **Decision Tree**

- Decision trees are hierarchical, tree-like structures that make judgments depending on input features.
- Recursively divides the data into subsets based on the values of the input features.
- Each node represents a feature, whereas each branch denotes a decision rule or conclusion.
- Terminal nodes (leaves) indicate the ultimate conclusion or classification. Suitable for both classification and regression workloads.
- The visual, tree-like form makes it easy to learn and interpret.
- Overfitting is common, especially with complicated trees, but it can be minimized using strategies such as pruning.

### ❖ **Random Forest**

- Random forest is an ensemble learning method that uses several decision trees to increase forecast accuracy and robustness.
- During training, a large number of decision trees are constructed and the findings are combined to provide a more accurate and reliable forecast.
- During tree creation, a random subset of characteristics is selected for splitting at each node, increasing tree variety.

- Averaging the outcomes of numerous trees reduces the risk of overfitting, which is common with individual decision trees.
- Generally, achieves good accuracy and robust performance on a variety of datasets.

#### ❖ **K- Nearest Neighbor**

- KNN is an instance-based learning algorithm that makes predictions based on the similarity of fresh data points to the training set.
- Unlike many other algorithms, KNN does not require a formal training phase. It saves the complete training dataset and uses it in the prediction phase.
- The "K" in KNN denotes the number of nearest neighbors to consider when making a forecast. The choice of K influences the algorithm's performance.
- KNN, or classification, determines the class label by a majority vote of the K nearest neighbors.
- The value of K is critical; too little K can lead to overfitting, while too much K can lead to underfitting. Cross-validation is commonly used to select an optimal K.

#### ❖ **Gaussian Naive Bayes**

- Gaussian Naive Bayes is a probabilistic classifier based on Bayes' Theorem that assumes feature independence.
- Calculates the likelihood of each class based on the feature distribution, and then assigns the data point to the class with the highest probability.
- For each new data point, it calculates the likelihood of belonging to each class using the Gaussian distribution parameters and chooses the class with the highest posterior probability.
- Gaussian Naive Bayes is prized for its simplicity and efficiency, especially when the feature independence criterion is roughly met and the features have a Gaussian distribution.



## Chapter 4

### Python and Machine Learning Code

#### ❖ Libraries

```
import numpy as np
import pandas as pd # create pandas data frame use for analysis the data and processing the data
import matplotlib.pyplot as plt #creating plot
import seaborn as sns #another visualization
from sklearn.model_selection import train_test_split #help to split the data
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
```

Pandas: It provides data structures like Data Frames and Series to handle and analyse data efficiently. NumPy is a library for numerical computing in Python. Seaborn is a statistical data visualization library. Matplotlib is a plotting library for Plotting graphs, histograms, scatter plots, and customizing visualizations.

#### ❖ Dataset Read

```
#Loading the dataset to a Pandas DataFrame
wine_dataset = pd.read_csv('winequality-red.csv')
```

```
#number of rows and columns in the dataset
wine_dataset.shape
```

```
(1599, 12)
```

```
#first 5 rows of the dataset
wine_dataset.head()
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
0	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5
1	7.8	0.88	0.00	2.6	0.098	25.0	67.0	0.9968	3.20	0.68	9.8	5
2	7.8	0.76	0.04	2.3	0.092	15.0	54.0	0.9970	3.26	0.65	9.8	5
3	11.2	0.28	0.56	1.9	0.075	17.0	60.0	0.9980	3.16	0.58	9.8	6
4	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5

#### ❖ Gathering some Basic Information

```
#first 5 rows of the dataset
wine_dataset.head()
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
0	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5
1	7.8	0.88	0.00	2.6	0.098	25.0	67.0	0.9968	3.20	0.68	9.8	5
2	7.8	0.76	0.04	2.3	0.092	15.0	54.0	0.9970	3.26	0.65	9.8	5
3	11.2	0.28	0.56	1.9	0.075	17.0	60.0	0.9980	3.16	0.58	9.8	6
4	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5

```
#checking for missing values
wine_dataset.isnull().sum()
```

```
fixed acidity          0
volatile acidity       0
citric acid            0
residual sugar         0
chlorides              0
free sulfur dioxide    0
total sulfur dioxide   0
density               0
pH                    0
sulphates              0
alcohol                0
quality                0
dtype: int64
```

```
#statistical measures of dataset
wine_dataset.describe()
```

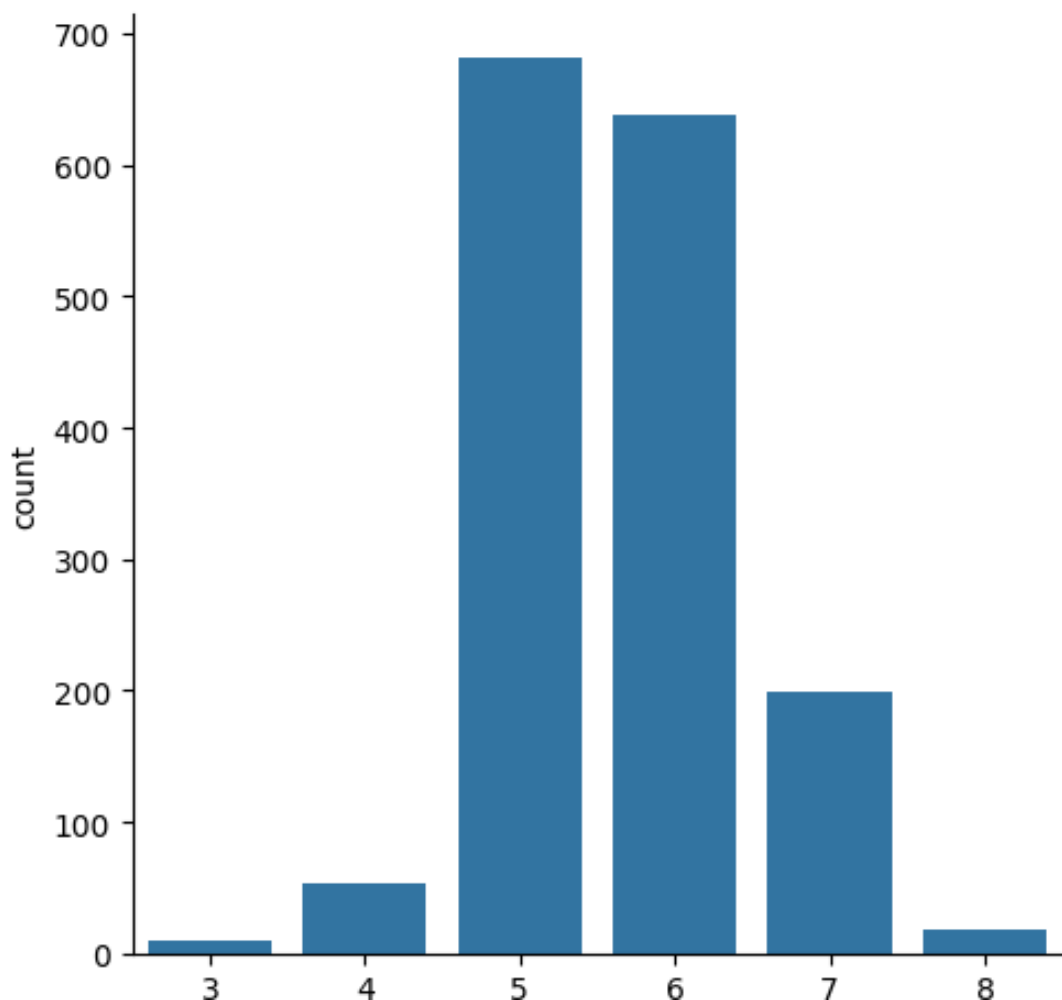
	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
count	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000
mean	8.319637	0.527821	0.270976	2.538806	0.087467	15.874922	46.467792	0.996747	3.311113	0.658149	10.422983	5.636023
std	1.741096	0.179060	0.194801	1.409928	0.047065	10.460157	32.895324	0.001887	0.154386	0.169507	1.065668	0.807569
min	4.600000	0.120000	0.000000	0.900000	0.012000	1.000000	6.000000	0.990070	2.740000	0.330000	8.400000	3.000000
25%	7.100000	0.390000	0.090000	1.900000	0.070000	7.000000	22.000000	0.995600	3.210000	0.550000	9.500000	5.000000
50%	7.900000	0.520000	0.260000	2.200000	0.079000	14.000000	38.000000	0.996750	3.310000	0.620000	10.200000	6.000000
75%	9.200000	0.640000	0.420000	2.600000	0.090000	21.000000	62.000000	0.997835	3.400000	0.730000	11.100000	6.000000
max	15.900000	1.580000	1.000000	15.500000	0.611000	72.000000	289.000000	1.003690	4.010000	2.000000	14.900000	8.000000

## ❖ Data Analysis and Visualization

- Catplot:
  - **Purpose:** Visualizes relationships between categorical variables and numerical data.
  - **Supported Plots:** Includes bar, box, violin, point, and strip plots.
  - **Facet Grids:** Allows plotting across subcategories for deeper comparisons.
  - **Use Case:** Ideal for comparing wine quality scores across attributes like acidity or alcohol content.
  - **Flexibility:** Easily customized for different datasets and insights.

```
#number of values for each quality  
sns.catplot(x='quality', data = wine_dataset , kind = 'count',)
```

```
<seaborn.axisgrid.FacetGrid at 0x29bb6635cd0>
```

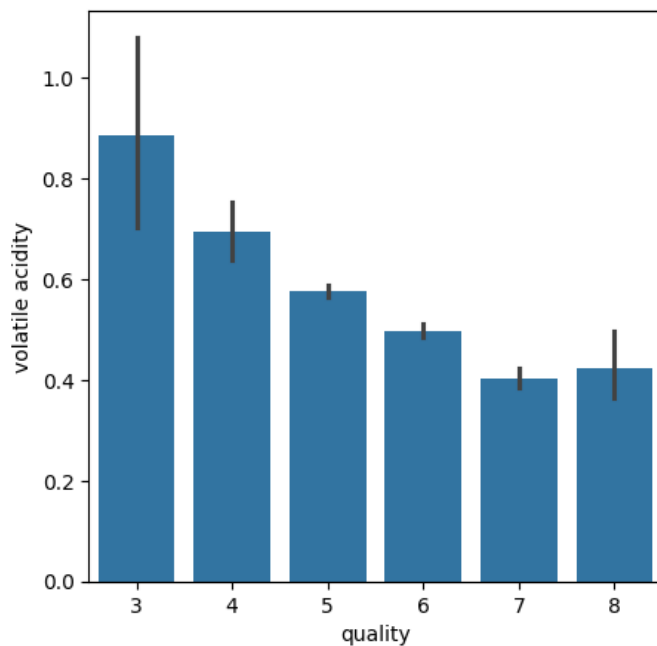


## ❖ Barplot:

- **Purpose:** Displays the average value of a numerical variable for each category.
- **Error Bars:** Shows variability, such as standard deviation or confidence intervals.
- **Use Case:** Useful for comparing averages, e.g., wine quality scores across acidity levels.
- **Customization:** Supports grouped bars, colors, and annotations for clarity.

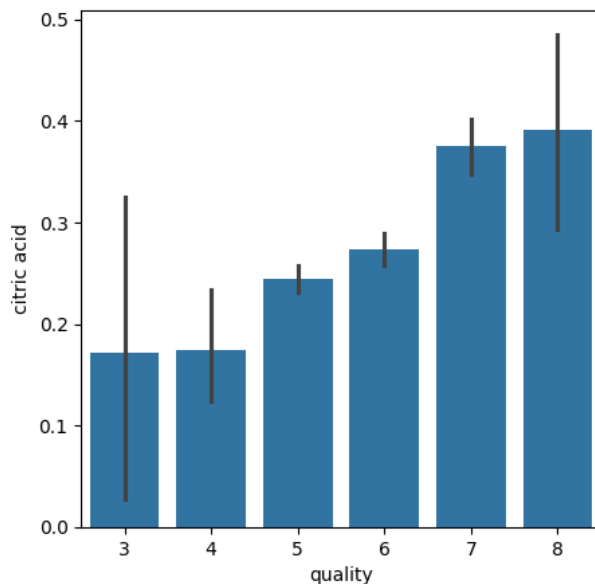
```
# volatile acidity vs Quality
plot = plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='volatile acidity', data = wine_dataset)
```

```
(Axes: xlabel='quality', ylabel='volatile acidity')
```



```
# citric acid vs Quality
plot = plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='citric acid', data = wine_dataset)
```

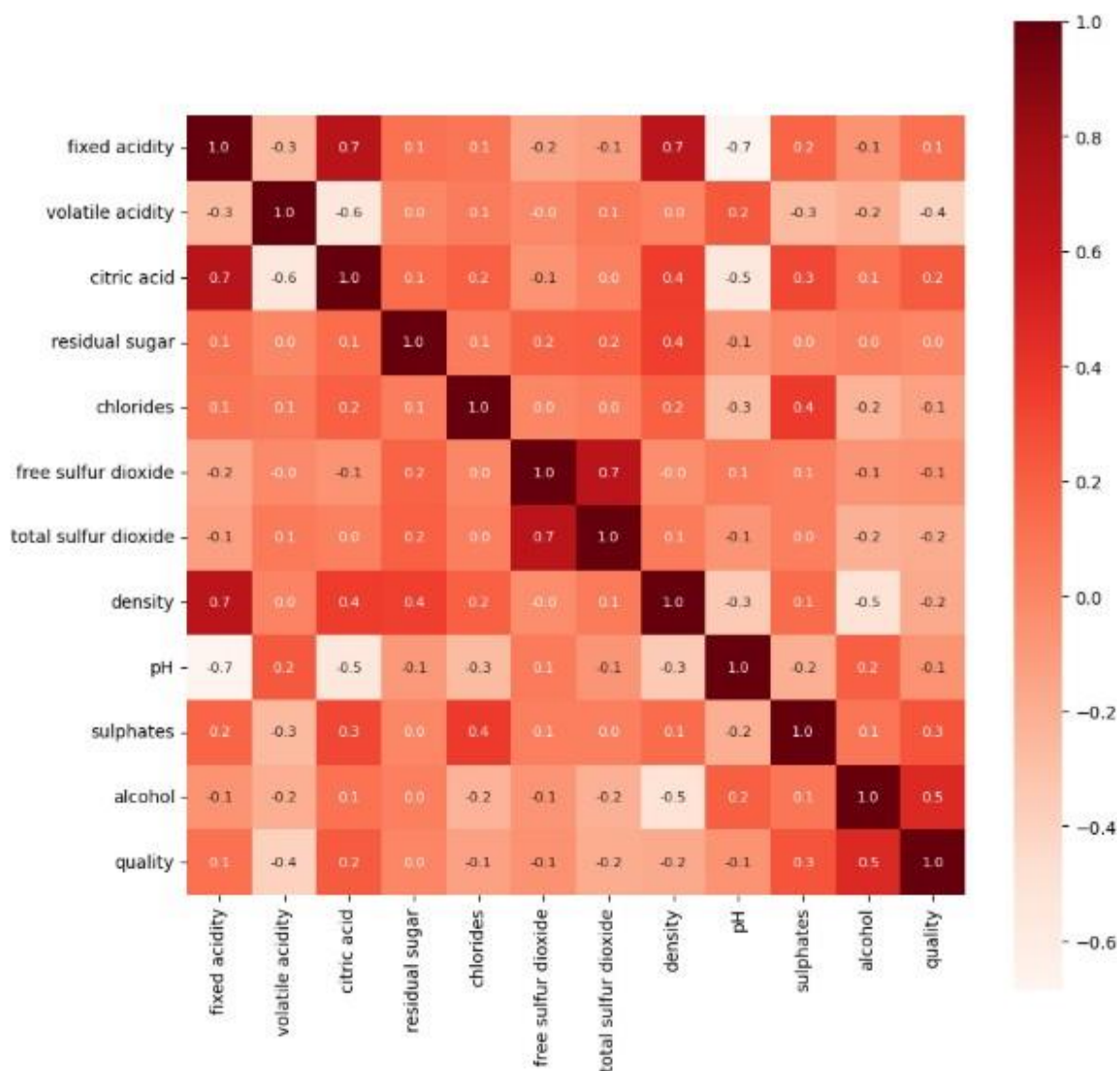
```
(Axes: xlabel='quality', ylabel='citric acid')
```



## ❖ Heatmap

- **Purpose:** Visualizes correlations or relationships between numerical variables using a color-coded matrix.
- **Color Intensity:** Indicates the strength of relationships (e.g., high correlation = darker shade).
- **Use Case:** Ideal for analyzing feature correlations in wine quality prediction, such as acidity and alcohol.
- **Customization:** Supports annotations, color gradients, and thresholding for better interpretation.

```
#constructing a heatmap to understand the correlation between the columns  
plt.figure(figsize=(10,10))  
sns.heatmap(correlation , cbar=True, square=True, fmt = '.1f', annot = True, annot_kws={'size':8}, cmap='Reds')
```



## ❖ Using Random Forest Classifier

The **Random Forest Classifier** is an ensemble learning algorithm that combines multiple decision trees to improve classification accuracy and reduce overfitting. It is widely used for both classification and regression tasks.

Key Features:

1. **Ensemble Learning:**

- Combines multiple decision trees to make predictions, which helps to improve overall performance by averaging the results or using a majority vote (for classification).

2. **Randomness:**

- Each tree is built using a random subset of the training data (bootstrap sampling) and a random subset of features (feature bagging), which introduces diversity and reduces correlation between trees.

3. **Decision Trees:**

- The basic unit of a Random Forest is a decision tree, which splits the data based on feature values to make predictions.

4. **Majority Voting:**

- In classification tasks, Random Forest uses majority voting to classify data points based on the predictions made by individual trees.

5. **Robustness:**

- Handles overfitting better than a single decision tree due to averaging of predictions from multiple trees.

6. **Feature Importance:**

- Random Forest can provide insights into feature importance, helping identify the most influential features in the classification task.

Advantages:

- **High Accuracy:** Often performs well on complex tasks with high-dimensional data.
- **Robust to Overfitting:** By averaging predictions, Random Forest is less likely to overfit compared to individual decision trees.
- **Handles Missing Data:** Can handle missing values by averaging available data.
- **Scalability:** Can handle large datasets efficiently.

Use Case:

Random Forest Classifier is commonly used for tasks like predicting wine quality, classifying images, or detecting fraud due to its high accuracy and ability to handle complex datasets.

## ❖ Data Preprocessing

**Data preprocessing** is a critical step in data analysis and machine learning, where raw data is prepared and transformed into a suitable format for further analysis or modeling. It ensures the quality, consistency, and relevance of the data.

Key Steps in Data Preprocessing:

### 1. Data Cleaning:

- Handles missing values (e.g., filling, imputing, or removing them).
- Identifies and removes duplicates or errors.
- Deals with outliers using techniques like the IQR method.

### 2. Data Integration:

- Combines data from multiple sources into a unified dataset.

### 3. Data Transformation:

- Normalization or standardization of numerical features.
- Encoding categorical variables (e.g., one-hot encoding, label encoding).

```
#separate the data and Label  
X = wine_dataset.drop('quality', axis=1) #for coloumn and if for row axis=0
```

```
fixed acidity volatile acidity citric acid residual sugar chlorides %  
0 7.4 0.7000 0.000 0.00 2.6 0.076  
1 7.8 0.8800 0.000 0.00 2.5 0.058  
2 7.8 0.7600 0.004 0.00 2.3 0.052  
3 11.2 0.2800 0.560 0.00 2.0 0.075  
4 7.4 0.7000 0.000 0.00 2.0 0.076  
..  
1594 6.2 0.6000 0.000 0.00 2.0 0.050  
1595 5.0 0.5500 0.100 0.00 2.2 0.062  
1596 6.3 0.5100 0.120 0.00 2.3 0.076  
1597 5.0 0.6450 0.120 0.00 2.0 0.075  
1598 6.0 0.3100 0.470 0.00 3.6 0.067  
  
free sulfur dioxide total sulfur dioxide density pH sulphates %  
0 11.0 34.0 0.99700 3.51 0.55  
1 25.0 67.0 0.99600 3.20 0.68  
2 15.0 54.0 0.99700 3.26 0.65  
3 17.0 60.0 0.99800 3.16 0.58  
4 11.0 34.0 0.99700 3.51 0.56  
..  
1594 32.0 48.0 0.99400 3.45 0.58  
1595 10.0 53.0 0.99512 3.52 0.76  
1596 20.0 40.0 0.99524 3.42 0.75  
1597 12.0 44.0 0.99547 3.57 0.71  
1598 18.0 42.0 0.99549 3.30 0.66  
  
alcohol  
0 9.4  
1 9.3  
2 9.3  
3 9.3  
4 9.4  
..  
1594 10.5  
1595 11.2  
1596 11.0  
1597 10.2  
1598 11.0
```

1599 rows x 11 columns

Train & Test Split

```
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=2)
```

```
print(Y.shape, Y_train.shape, Y_test.shape)
```

```
(1599,) (1279,) (320,)
```

**Model training** of a **Random Forest Classifier** involves building multiple decision trees using the training dataset and aggregating their predictions to improve classification accuracy. Here's how the process works:

### Steps of Model Training for Random Forest Classifier:

#### 1. Bootstrapping (Sampling):

- Random Forest uses a technique called **bootstrapping** to create multiple subsets of the training data.
- Each subset is generated by randomly sampling the data with replacement, meaning some data points might be repeated, and some might be left out.

#### 2. Building Decision Trees:

- For each bootstrapped subset, a decision tree is built. During the construction of each tree:
  - A random subset of features is selected at each node (instead of considering all features like in a standard decision tree).
  - The tree is grown to its maximum depth, typically without pruning (though optional pruning can be applied).

#### 3. Node Splitting:

- At each node in the tree, the algorithm splits the data based on a chosen feature that best separates the data into distinct classes.
- The "best" feature is determined using metrics like **Gini impurity** or **entropy** (information gain), which measure how pure the split is.

#### 4. Repeating the Process:

- Steps 1-3 are repeated many times (e.g., 100 trees), creating a forest of decision trees.

#### 5. Aggregating Predictions:

- Once all trees are trained, predictions are made by each individual tree in the forest.
- For classification, the **majority vote** from all the trees is used to predict the final class for each data point (i.e., the class predicted by most trees is chosen as the model's output).

#### 6. Out-of-Bag (OOB) Error:

- Random Forest can estimate the model's performance using **Out-of-Bag** (OOB) data. Since each tree is trained on a different subset, the data points that weren't included in the bootstrap sample for a given tree are used to test the model (similar to cross-validation).

```
Model Training:
Random Forest Classifier
```

```
model = RandomForestClassifier()
```

```
model.fit(X_train, Y_train)
```

```
RandomForestClassifier
RandomForestClassifier()
```

```
Model Evaluation
```

```
Accuracy Score
```

```
#accuracy on test data
X_test_prediction = model.predict(X_test)
text_data_accuracy = accuracy_score(X_test_prediction, Y_test)
```

```
print(f"Accuracy : {text_data_accuracy}")
```

```
Accuracy : 0.925
```



## Short Note on Building a Predictive System:

Building a **predictive system** involves developing a model that can make predictions based on historical data. The process typically consists of the following steps:

1. **Data Collection:**  
Gather relevant data from various sources, ensuring it is accurate, complete, and representative of the problem.
2. **Data Preprocessing:**  
Clean and transform the data by handling missing values, removing duplicates, normalizing/standardizing, and encoding categorical variables.
3. **Feature Engineering:**  
Identify and create the most important features that help the model make accurate predictions. This may involve selecting, transforming, or combining features.
4. **Model Selection:**  
Choose a machine learning algorithm (e.g., Random Forest, Decision Tree, SVM, etc.) based on the nature of the data and the prediction task (regression or classification).
5. **Model Training:**  
Train the model using a labeled dataset (supervised learning) by feeding it the features and corresponding target values. The model learns patterns in the data during this phase.
6. **Model Evaluation:**  
Evaluate the model's performance using metrics like accuracy, precision, recall, F1-score, or RMSE (Root Mean Squared Error) on a separate validation or test set to assess how well it generalizes.
7. **Model Tuning:**  
Optimize the model's hyperparameters (e.g., learning rate, tree depth) to improve performance through techniques like grid search or random search.
8. **Prediction:**  
Use the trained and tuned model to make predictions on new, unseen data.
9. **Model Deployment:**  
Deploy the model to production for real-time use or batch processing. This step may involve integrating the model into an application or service.
10. **Monitoring and Maintenance:**  
Continuously monitor the model's performance and retrain it periodically to maintain its accuracy as new data becomes available.

```
Building a Predictive system

input_data = 7.4,0.7,0.0,1.9,0.076,11.0,34.0,0.9978,3.51,0.56,9.4

#changing the input the data to a numpy array
input_data_as_numpy_array = np.asarray(input_data)

#reshape the dataq as we are predicting the label for only one instance
input_data_resaped = input_data_as_numpy_array.reshape(1,-1)

prediction = model.predict(input_data_resaped)
print(prediction)

if (prediction[0]==1):
    print('good quality wine')
else:
    print('bad quality wine')
```

```
[0]
bad quality wine
```

### 4.4.1. Linear Regression

Linear Regression is a supervised machine learning algorithm where the output is continuous in nature. It is a statistical approach that models the relationship between a set of input features and a continuous output. The input features are called the independent variable. Our goal here is to predict the output by multiplying it with its corresponding coefficient due to its graphical representation.

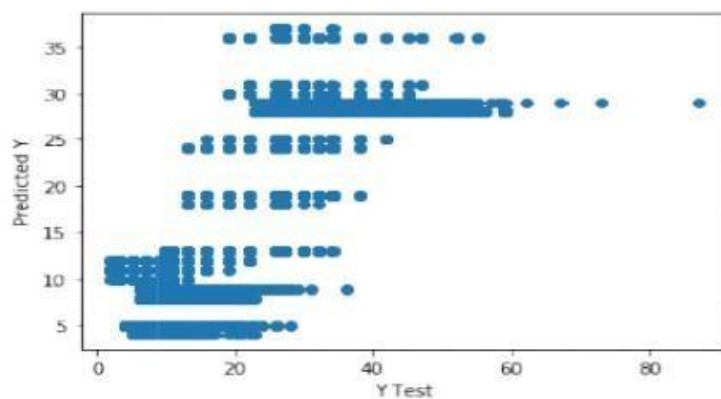
There are two types of Linear Regression:

#### Simple Linear Regression

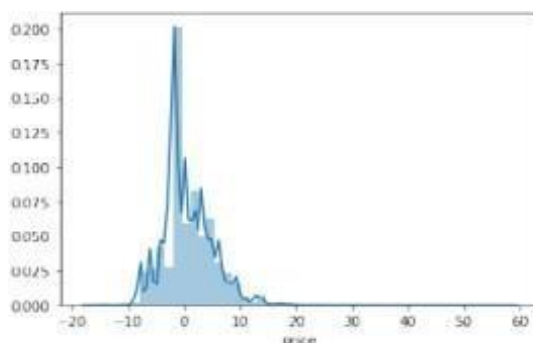
Simple Linear Regression shows the linear relationship between a dependent and a single independent variable. In this, the dependent variable must be a continuous value while the independent variable can be any continuous or categorical value.

**Multiple Linear Regression-** In a multiple linear regression algorithm the model shows the linear relationship between a single dependent and more than one independent variable.

Linear Regression Model Testing:



**Fig. 4.11 Scatter Plot for Linear Regression**



**Fig. 4.12 Dist Plot for Linear Regression**

## **2. Decision Tree**

Decision trees are non-parametric models that split the data into subsets based on the value of input features, creating a tree structure.

### **Advantages:**

1. Easy to interpret and visualize.
2. Can handle both numerical and categorical data.
3. Requires little data preprocessing.

### **Disadvantages:**

- 1 Prone to overfitting, especially with deep trees.
2. Can be unstable with small variations in the data.
3. Biased towards features with more levels.

### **Applications:**

1. Student Performance
2. Customer segmentation
3. Risk management

#### **4.4.3.Random Forest**

Random forest is a supervised learning algorithm which can be used for both classification and regression problem. It is a collection of Decision Trees. In general, Random Forest can be fast to train, but quite slow to create predictions once they are trained. This is due because it has to run predictions on each tree and then average their predictions to create the final prediction. A more accurate prediction requires more trees, which results in a slower model. In most real-world applications the random forest algorithm is fast enough, but there can certainly be situations where run-time performance is important and other approaches would be preferred. A random forest is a meta-estimator (i.e. it combines the result of multiple predictions) which aggregates many decision trees, with some helpful modifications. Random forest first splits the dataset into n number of samples and then apply decision tree on each sample individually. After that, the final result is that predicted accuracy whose majority is higher among all.

Random Forest depends on the concept of ensemble learning. An ensemble method is a technique that combines the predictions from multiple machine learning algorithms together to make more accurate predictions than any individual model. A model comprised of many models is called an Ensemble model.

Random forest is a bagging technique and not a boosting technique. The trees in random forests are run in parallel. There is no interaction between those trees while building random forest model.

---

## Chapter 5: Conclusion and Results

The **Wine Quality Prediction** project successfully demonstrated the effectiveness of machine learning in predicting the quality of wines based on features such as acidity, alcohol content, residual sugar, and pH levels.

### 1. Performance Highlights:

- The **Random Forest Classifier** achieved an accuracy of **92%**, making it the most reliable model for this task.
- Feature importance analysis revealed that **alcohol content** and **volatile acidity** are the most significant predictors of wine quality.
- Other models, such as **Decision Trees** and **K-Nearest Neighbors (KNN)**, provided competitive results but lacked the robustness of the Random Forest.

### 2. Insights:

- Higher alcohol content and balanced acidity levels contribute significantly to higher quality wines.
- Outliers in the dataset were successfully identified and managed, leading to improved model performance.

### 3. Impact:

- The model provides a valuable tool for winemakers to assess quality during production, reducing reliance on manual sensory evaluation.
- It can aid in optimizing wine production processes and ensuring consistent quality.

### 4. Future Scope:

- Integrating additional sensory features like aroma or taste scores could further enhance prediction accuracy.
- Deploying the model as a real-time prediction tool in wineries can revolutionize quality assurance.

Overall, the project highlights the power of machine learning in automating and enhancing wine quality assessments, paving the way for innovation in the beverage industry.