

Prediction of wine quality with wine quality dataset

By
Lukesh Tiwari
Nitesh Kumar
Riti Shah

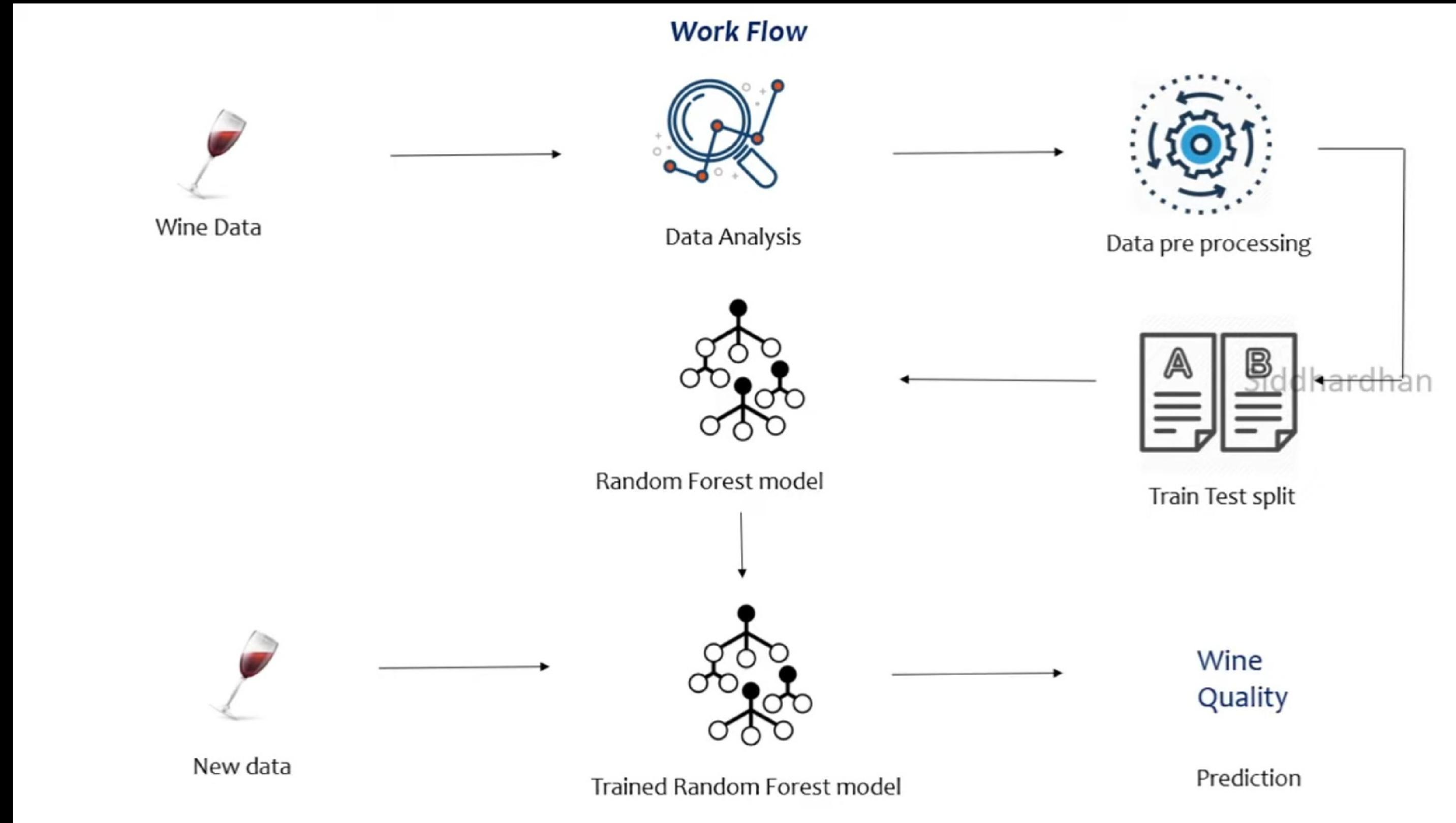


Introduction

Wine is the most commonly used beverage globally, and its values are considered important in society. Quality of the wine is always important for its consumers, and mainly for producers in the present competitive market to raise the revenue. Historically, wine quality used to be determined by testing at the end of the production; to reach that level, one already spends lots of time and money. If the quality is not good, then the various procedure needs to be implemented from the beginning, which is very costly.

ML can be an alternative to identify the most important parameters that control the wine quality. In this work, we have shown how ML can be used to identify the best parameter on which the wine quality depends and in turn predict wine quality.

Workflow Diagram



Data Description and Preprocessing

In this study, we use the publicly available wine quality dataset obtained from the Keggle which contains a large collection of datasets that have been widely used by the machine learning community. Among the two types of wine quality dataset (redwine and white wine), we have chosen redwine data for our study because of its popularity over the white wine. The redwine dataset contains 11 physiochemical properties: fixed acidity (g[tartaric acid]/dm3), volatile acidity (g[acetic acid]/dm3), total sulfur dioxide (mg/dm3), chlorides (g[sodium chloride]/dm3), pH level, free sulfur dioxide (mg/dm3), density (g/cm3), residual sugar (g/dm3), citric acid (g/dm3), sulphates (g[potassium sulphate]/dm3), and alcohol (vol%). Alongside these properties, a sensory score was acquired from several different blind taste testers which graded each wine sample with a score ranging from zero (poor) to 10 (excellent).

	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

Descriptive statistics

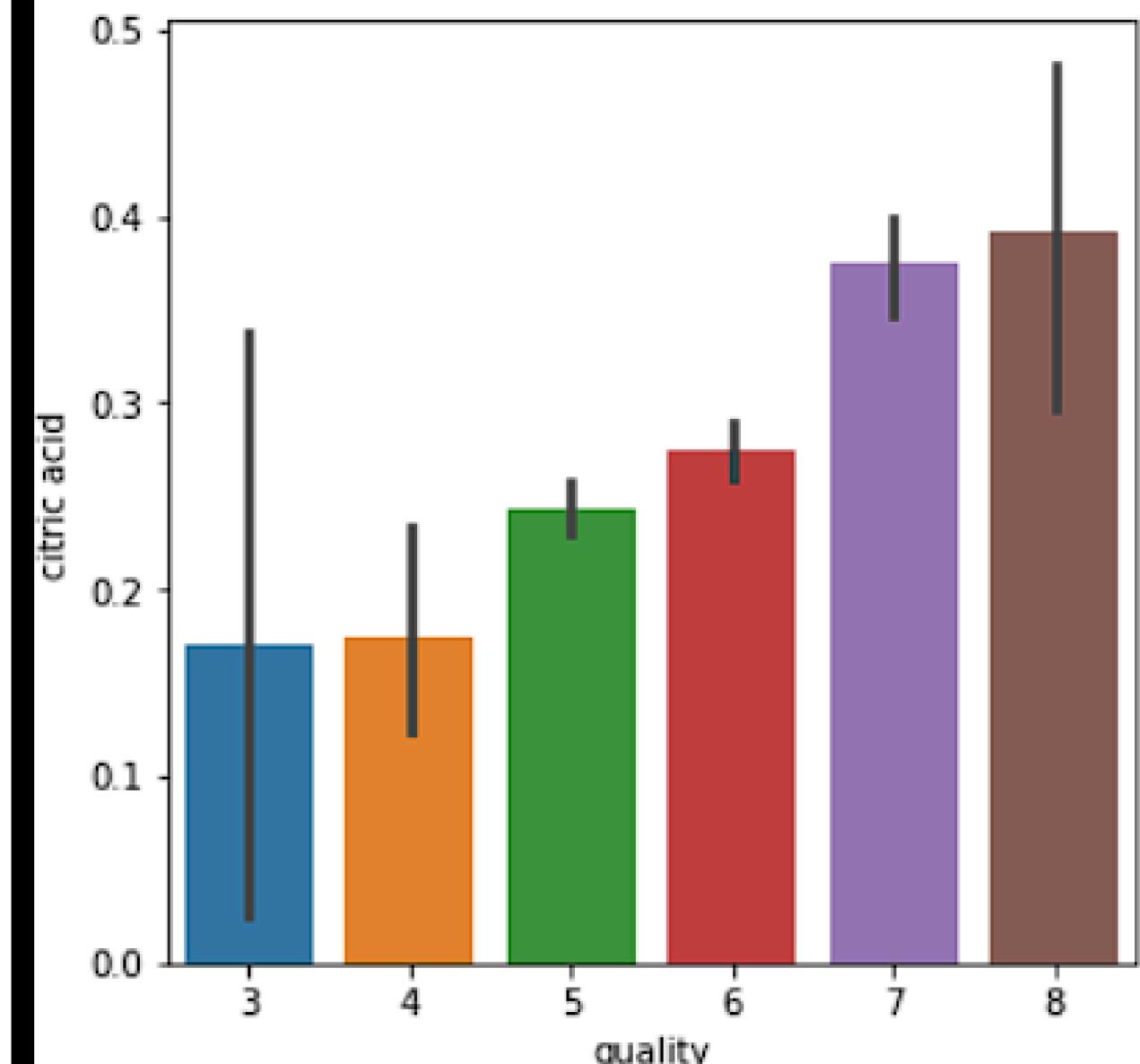
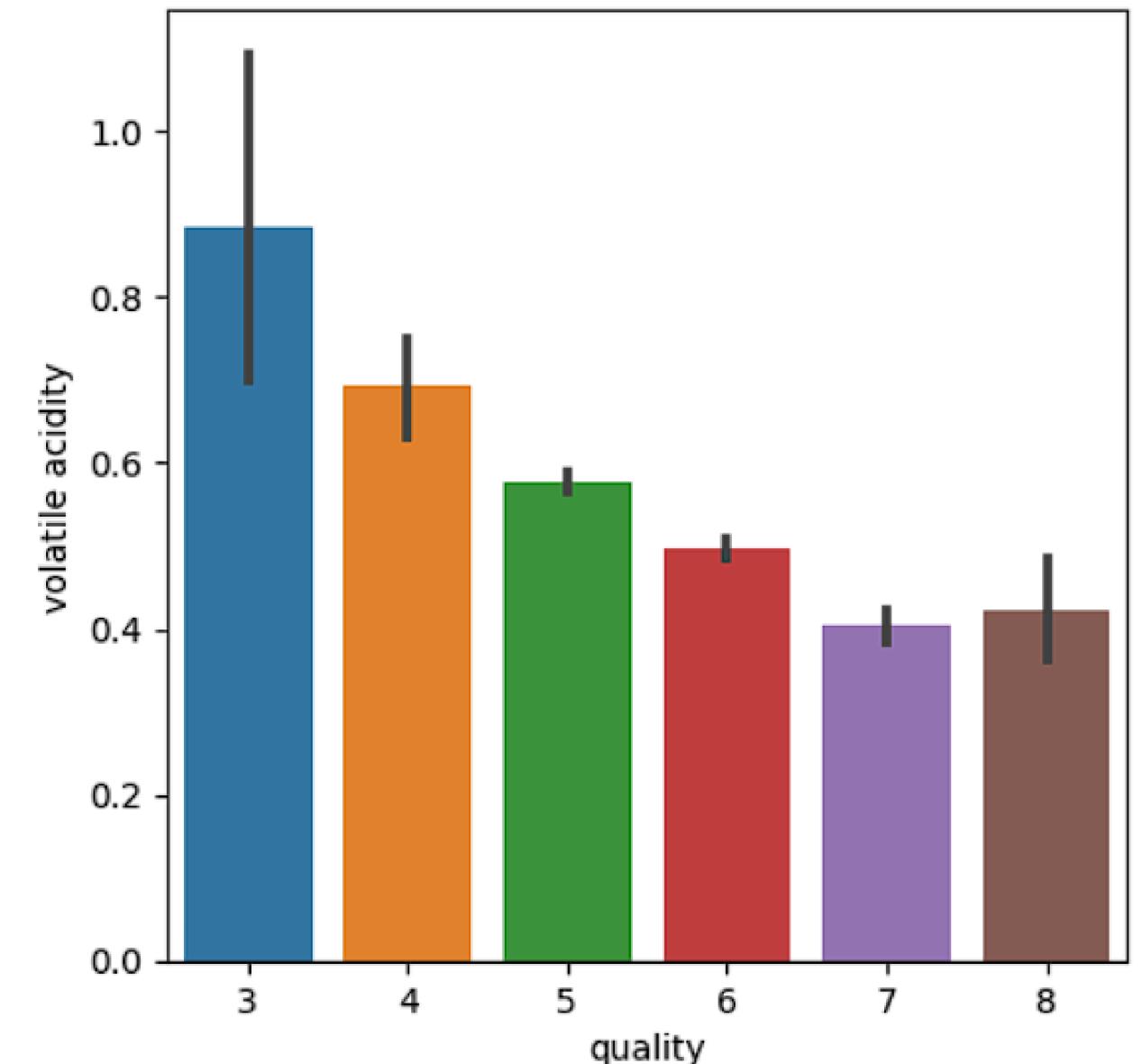
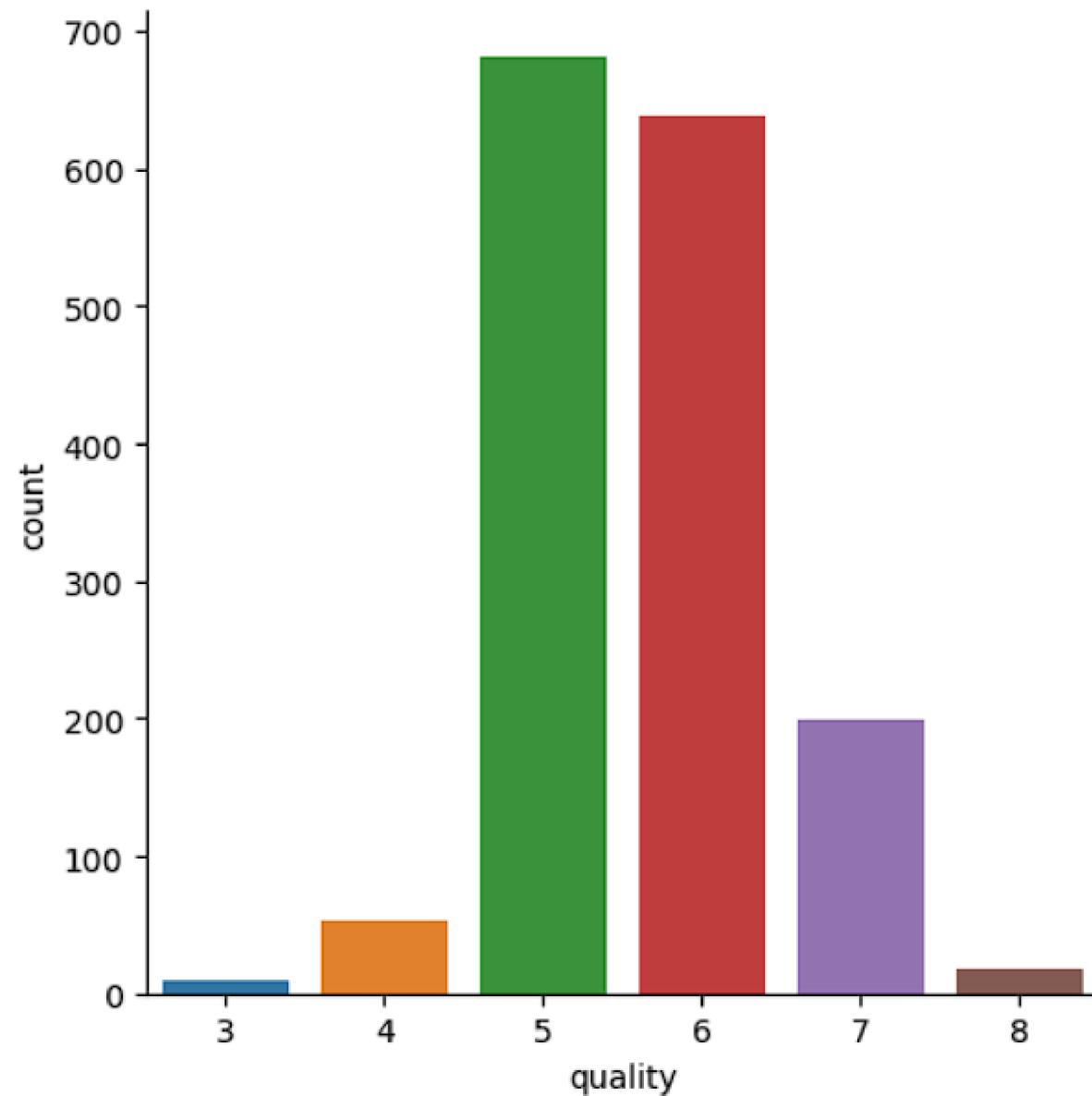
	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
count	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000
mean	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

Representing Data

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

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

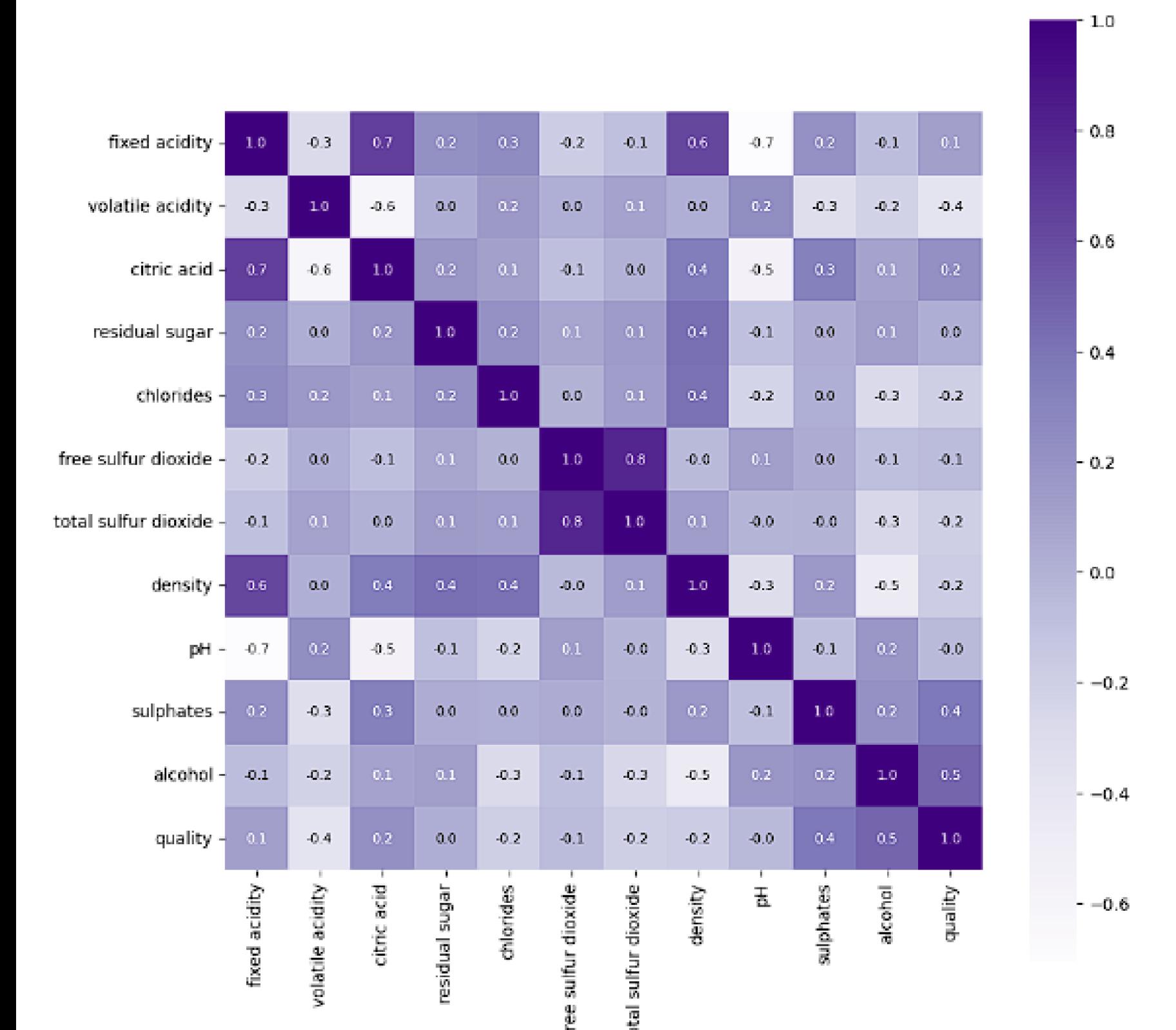
```
# volatile acidity vs Quality  
plot = plt.figure(figsize=(5,5))  
sns.barplot(x='quality', y = 'volatile acidity', data = wine_dataset)
```



```
# 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 = 'Purples')
```



OUTPUT

Code

Data Preprocessing

```
In [12]: # separate the data and Label  
X = wine_dataset.drop('quality',axis=1)
```

```
In [13]: print(X)
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	\
0	7.4	0.700	0.00	1.9	0.076	
1	7.8	0.880	0.00	2.6	0.098	
2	7.8	0.760	0.04	2.3	0.092	
3	11.2	0.280	0.56	1.9	0.075	
4	7.4	0.700	0.00	1.9	0.076	
...
1594	6.2	0.600	0.08	2.0	0.090	
1595	5.9	0.550	0.10	2.2	0.062	
1596	6.3	0.510	0.13	2.3	0.076	
1597	5.9	0.645	0.12	2.0	0.075	
1598	6.0	0.310	0.47	3.6	0.067	
	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	\
0	11.0	34.0	0.99780	3.51	0.56	
1	25.0	67.0	0.99680	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	0.99780	3.51	0.56	
...
1594	32.0	44.0	0.99490	3.45	0.58	
1595	39.0	51.0	0.99512	3.52	0.76	
1596	29.0	40.0	0.99574	3.42	0.75	
1597	32.0	44.0	0.99547	3.57	0.71	
1598	18.0	42.0	0.99549	3.39	0.66	
	alcohol					
0	9.4					
1	9.8					
2	9.8					
3	9.8					
4	9.4					
...	...					
1594	10.5					
1595	11.2					
1596	11.0					
1597	10.2					
1598	11.0					

[1599 rows x 11 columns]

Label Binarization

```
In [14]: Y = wine_dataset['quality'].apply(lambda y_value: 1 if y_value>=7 else 0)
```

```
In [15]: print(Y)
```

```
0      0  
1      0  
2      0  
3      0  
4      0  
..  
1594    0  
1595    0  
1596    0  
1597    0  
1598    0  
Name: quality, Length: 1599, dtype: int64
```

Train & Test Split

```
In [16]: X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=3)
```

```
In [17]: print(Y.shape, Y_train.shape, Y_test.shape)
```

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

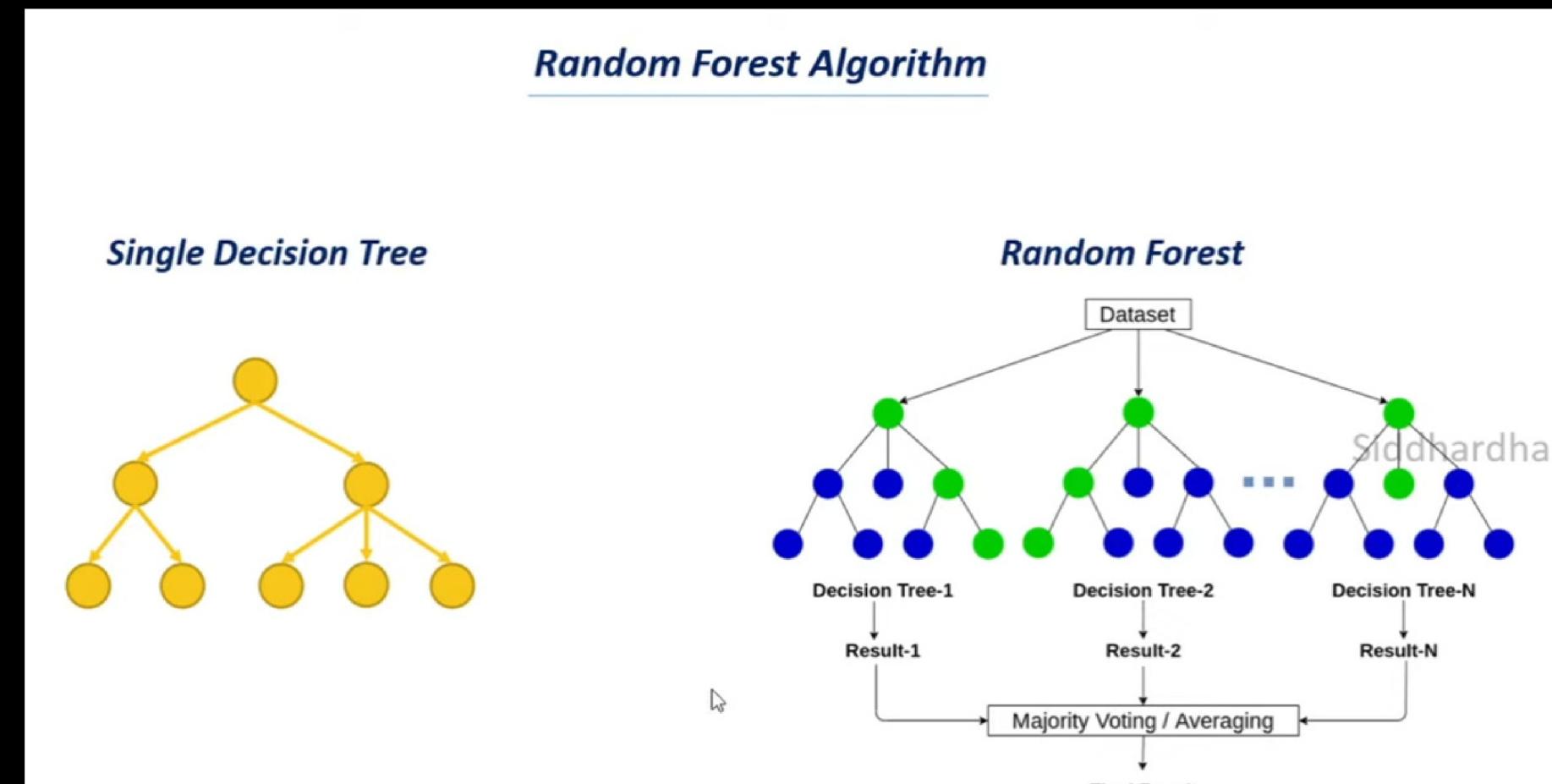
Decision Tree Algoithm

Decision Tree Terminologies

- **Root Node:** Root node is from where the decision tree starts. It represents the entire dataset, which further gets divided into two or more homogeneous sets.
- **Leaf Node:** Leaf nodes are the final output node, and the tree cannot be segregated further after getting a leaf node.
- **Splitting:** Splitting is the process of dividing the decision node/root node into sub-nodes according to the given conditions.
- Branch/Sub Tree: A tree formed by splitting the tree.
- Pruning: Pruning is the process of removing the unwanted branches from the tree.
- Parent/Child node: The root node of the tree is called the parent node, and other nodes are called the child nodes.

Random Forest Algorithm

Random Forest algorithm is a powerful tree learning technique in Machine Learning. It works by creating a number of Decision Trees during the training phase. Each tree is constructed using a random subset of the data set to measure a random subset of features in each partition.



Let's Move To The Coding Part