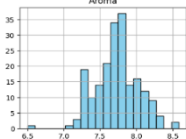
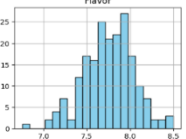
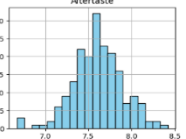
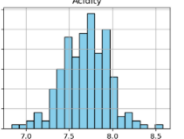
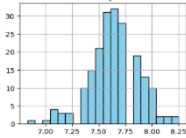
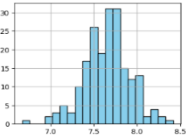
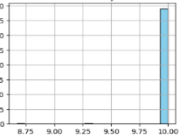
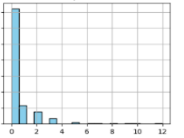


Data Collection and Preprocessing Phase

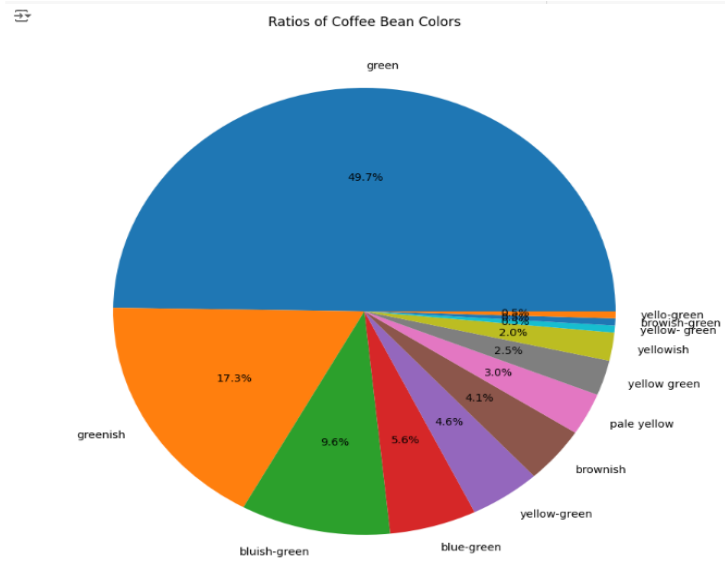
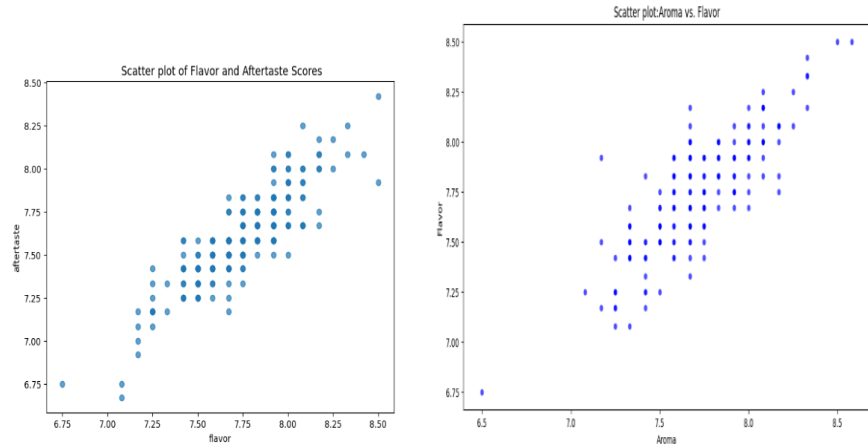
Date	9 July 2024
Team ID	team-739994
Project Title	Precise Coffee Quality Prediction
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

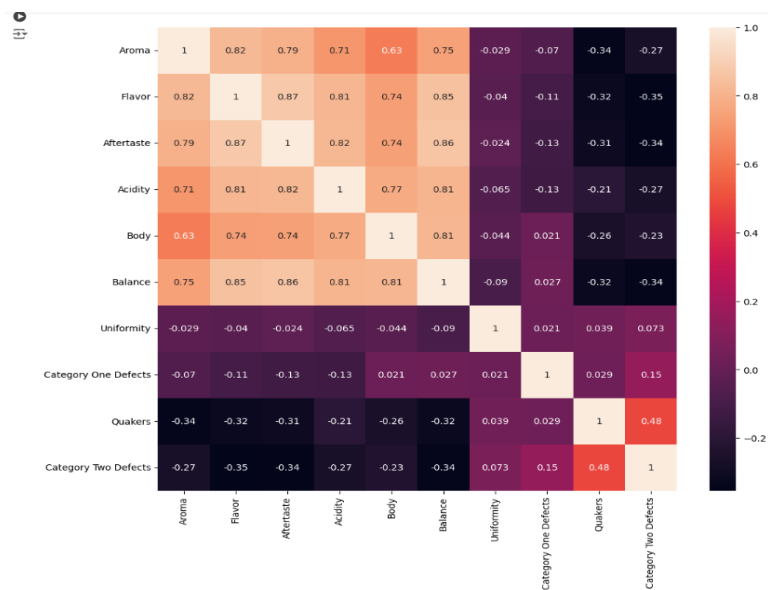
Dataset variables will be statistically analyzed to identify patterns and outliers, with python employed for preprocessing tasks like normalization and feature engineering .Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description
Data Overview	Dimensions: 207 rows x 19 columns <u>Descriptive Statistics:</u>
Univariate Analysis	<div><div></div><div>Histograms of Coffee Quality Scores</div><div><div><div>Aroma</div></div><div><div>Flavor</div></div><div><div>Aftertaste</div></div><div><div>Acidity</div></div></div><div><div><div>Body</div></div><div><div>Balance</div></div><div><div>Uniformity</div></div><div><div>Quakers</div></div></div></div>

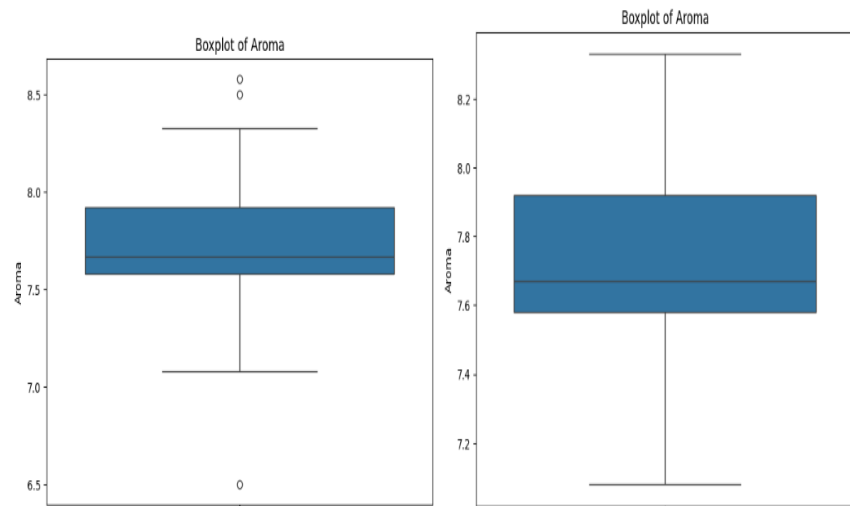
Bivariate Analysis



Multivariate Analysis



Outliers and Anomalies



Data Preprocessing Code Screenshots

Loading Data

```
[ ] import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
df = pd.read_csv("coffee_data.csv")
df
```

ID	Number of Bags	Bag Weight	Variety	Processing Method	Aroma	Flavor	Aftertaste	Acidity	Body	Balance	Uniformity	Overall	Total Cup Points	Moisture Percentage	Category One Defects	Quakers	Color	Category Two Defects	
0	0	1	35 kg	Castillo Double Arabica	Washed	8.58	8.58	8.42	8.58	8.25	8.42	10.0	8.58	88.33	11.5	0	0	green	3
1	1	1	80 kg	Gesha	Washed / Wet	8.50	8.50	7.92	8.00	7.92	8.25	10.0	8.50	87.58	10.5	0	0	blue-green	0
2	2	19	25 kg	Java	Semi Washed	8.33	8.42	8.08	8.17	7.92	8.17	10.0	8.33	87.42	10.4	0	0	yellowish	2
3	3	1	22 kg	Gesha	Washed / Wet	8.08	8.17	8.17	8.25	8.17	8.08	10.0	8.25	87.17	11.0	0	0	green	0
4	4	2	24 kg	Red Bourbon	Honey Moist	8.33	8.33	8.08	8.25	7.92	7.92	10.0	8.25	87.08	11.6	0	2	yellow-green	2
...
202	202	2240	60 kg	Mundo Novo	Natural / Dry	7.17	7.17	6.92	7.17	7.42	7.17	10.0	7.08	88.08	11.4	0	0	green	4
203	203	300	30 kg	SHG	Natural / Dry	7.33	7.08	6.75	7.17	7.42	7.17	10.0	7.08	88.00	10.4	0	2	green	12
204	204	343	60 kg	Caturra	Washed / Wet	7.25	7.17	7.08	7.08	7.08	7.08	10.0	7.08	79.87	11.6	0	9	green	11
205	205	1	2 kg	Mangrove	Natural / Dry	6.50	6.75	6.75	7.17	7.08	7.00	10.0	6.83	78.08	11.0	0	12	black-green	13
206	206	600	60 kg	Mundo Novo	SEM-LAMGO	7.25	7.08	6.87	6.83	6.83	6.87	10.0	6.87	78.00	11.3	0	0	green	1

207 rows x 19 columns

Handling Missing Data

```
[ ] df.isnull().sum()
```

```
ID
Number of Bags    0
Bag Weight         0
Variety            6
Processing Method  5
Aroma              0
Flavor             0
Aftertaste         0
Acidity            0
Body              0
Balance            0
Uniformity         0
Overall            0
Total Cup Points   0
Moisture Percentage 0
Category One Defects 0
Quakers            0
Color              0
Category Two Defects 0
dtype: int64
```

```
[ ] df.dropna(inplace=True)
```

```
[ ] df.isna().sum()
```

```
ID
Number of Bags    0
Bag Weight         0
Variety            0
Processing Method  0
Aroma              0
Flavor             0
Aftertaste         0
Acidity            0
Body              0
Balance            0
Uniformity         0
Overall            0
Total Cup Points   0
Moisture Percentage 0
Category One Defects 0
Quakers            0
Color              0
Category Two Defects 0
dtype: int64
```

Data Transformation	<pre>[] from sklearn.preprocessing import LabelEncoder label_encoder = LabelEncoder() df1["Color_Encoded"] = label_encoder.fit_transform(df1["Color"]) df1 = df1.drop(["Color"],axis=1)</pre> <p>df1</p> <table><thead><tr><th></th><th>Aroma</th><th>Flavor</th><th>Aftertaste</th><th>Acidity</th><th>Body</th><th>Balance</th><th>Uniformity</th><th>Category One Defects</th><th>Quakers</th><th>Category Two Defects</th><th>Color_Encoded</th></tr></thead><tbody><tr><td>0</td><td>8.58</td><td>8.50</td><td>8.42</td><td>8.58</td><td>8.25</td><td>8.42</td><td>10.0</td><td>0</td><td>0</td><td>3</td><td>4</td></tr><tr><td>1</td><td>8.50</td><td>8.50</td><td>7.92</td><td>8.00</td><td>7.92</td><td>8.25</td><td>10.0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr><tr><td>2</td><td>8.33</td><td>8.42</td><td>8.08</td><td>8.17</td><td>7.92</td><td>8.17</td><td>10.0</td><td>0</td><td>0</td><td>2</td><td>11</td></tr><tr><td>3</td><td>8.08</td><td>8.17</td><td>8.17</td><td>8.25</td><td>8.17</td><td>8.08</td><td>10.0</td><td>0</td><td>0</td><td>0</td><td>4</td></tr><tr><td>4</td><td>8.33</td><td>8.33</td><td>8.08</td><td>8.25</td><td>7.92</td><td>7.92</td><td>10.0</td><td>0</td><td>2</td><td>2</td><td>10</td></tr><tr><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td></tr><tr><td>202</td><td>7.17</td><td>7.17</td><td>6.92</td><td>7.17</td><td>7.42</td><td>7.17</td><td>10.0</td><td>0</td><td>0</td><td>4</td><td>4</td></tr><tr><td>203</td><td>7.33</td><td>7.08</td><td>6.75</td><td>7.17</td><td>7.42</td><td>7.17</td><td>10.0</td><td>0</td><td>2</td><td>12</td><td>4</td></tr><tr><td>204</td><td>7.25</td><td>7.17</td><td>7.08</td><td>7.00</td><td>7.08</td><td>7.08</td><td>10.0</td><td>0</td><td>9</td><td>11</td><td>4</td></tr></tbody></table>		Aroma	Flavor	Aftertaste	Acidity	Body	Balance	Uniformity	Category One Defects	Quakers	Category Two Defects	Color_Encoded	0	8.58	8.50	8.42	8.58	8.25	8.42	10.0	0	0	3	4	1	8.50	8.50	7.92	8.00	7.92	8.25	10.0	0	0	0	0	2	8.33	8.42	8.08	8.17	7.92	8.17	10.0	0	0	2	11	3	8.08	8.17	8.17	8.25	8.17	8.08	10.0	0	0	0	4	4	8.33	8.33	8.08	8.25	7.92	7.92	10.0	0	2	2	10	202	7.17	7.17	6.92	7.17	7.42	7.17	10.0	0	0	4	4	203	7.33	7.08	6.75	7.17	7.42	7.17	10.0	0	2	12	4	204	7.25	7.17	7.08	7.00	7.08	7.08	10.0	0	9	11	4
	Aroma	Flavor	Aftertaste	Acidity	Body	Balance	Uniformity	Category One Defects	Quakers	Category Two Defects	Color_Encoded																																																																																																														
0	8.58	8.50	8.42	8.58	8.25	8.42	10.0	0	0	3	4																																																																																																														
1	8.50	8.50	7.92	8.00	7.92	8.25	10.0	0	0	0	0																																																																																																														
2	8.33	8.42	8.08	8.17	7.92	8.17	10.0	0	0	2	11																																																																																																														
3	8.08	8.17	8.17	8.25	8.17	8.08	10.0	0	0	0	4																																																																																																														
4	8.33	8.33	8.08	8.25	7.92	7.92	10.0	0	2	2	10																																																																																																														
...																																																																																																														
202	7.17	7.17	6.92	7.17	7.42	7.17	10.0	0	0	4	4																																																																																																														
203	7.33	7.08	6.75	7.17	7.42	7.17	10.0	0	2	12	4																																																																																																														
204	7.25	7.17	7.08	7.00	7.08	7.08	10.0	0	9	11	4																																																																																																														
Feature Engineering	<pre>[] df1.loc[condition_healthy,'Bean_Status']='Healthy' condition_healthy=(df1["Category One Defects"]==0) & (df1["Category Two Defects"]==0) df1.loc[condition_healthy,'Bean_Status']='Healthy' condition_unhealthy=(df1["Category One Defects"]!=0) & (df1["Category Two Defects"]!=0) df1.loc[condition_unhealthy,'Bean_Status']='Unhealthy'</pre> <p>df1</p> <table><thead><tr><th></th><th>Aroma</th><th>Flavor</th><th>Aftertaste</th><th>Acidity</th><th>Body</th><th>Balance</th><th>Uniformity</th><th>Category One Defects</th><th>Quakers</th><th>Category Two Defects</th><th>Color_Encoded</th><th>Bean_Status</th></tr></thead><tbody><tr><td>0</td><td>8.58</td><td>8.50</td><td>8.42</td><td>8.58</td><td>8.25</td><td>8.42</td><td>10.0</td><td>0</td><td>0</td><td>3</td><td>4</td><td>Healthy</td></tr><tr><td>1</td><td>8.50</td><td>8.50</td><td>7.92</td><td>8.00</td><td>7.92</td><td>8.25</td><td>10.0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>Healthy</td></tr><tr><td>2</td><td>8.33</td><td>8.42</td><td>8.08</td><td>8.17</td><td>7.92</td><td>8.17</td><td>10.0</td><td>0</td><td>0</td><td>2</td><td>11</td><td>Healthy</td></tr><tr><td>3</td><td>8.08</td><td>8.17</td><td>8.17</td><td>8.25</td><td>8.17</td><td>8.08</td><td>10.0</td><td>0</td><td>0</td><td>0</td><td>4</td><td>Healthy</td></tr></tbody></table>		Aroma	Flavor	Aftertaste	Acidity	Body	Balance	Uniformity	Category One Defects	Quakers	Category Two Defects	Color_Encoded	Bean_Status	0	8.58	8.50	8.42	8.58	8.25	8.42	10.0	0	0	3	4	Healthy	1	8.50	8.50	7.92	8.00	7.92	8.25	10.0	0	0	0	0	Healthy	2	8.33	8.42	8.08	8.17	7.92	8.17	10.0	0	0	2	11	Healthy	3	8.08	8.17	8.17	8.25	8.17	8.08	10.0	0	0	0	4	Healthy																																																							
	Aroma	Flavor	Aftertaste	Acidity	Body	Balance	Uniformity	Category One Defects	Quakers	Category Two Defects	Color_Encoded	Bean_Status																																																																																																													
0	8.58	8.50	8.42	8.58	8.25	8.42	10.0	0	0	3	4	Healthy																																																																																																													
1	8.50	8.50	7.92	8.00	7.92	8.25	10.0	0	0	0	0	Healthy																																																																																																													
2	8.33	8.42	8.08	8.17	7.92	8.17	10.0	0	0	2	11	Healthy																																																																																																													
3	8.08	8.17	8.17	8.25	8.17	8.08	10.0	0	0	0	4	Healthy																																																																																																													
Save Processed Data	<pre>[] import pickle import warnings [] with open("./coffee_quality_prediction(rfc).pkl","wb") as f: pickle.dump(RFC,f)</pre>																																																																																																																								