# Data Collection and Preprocessing Phase

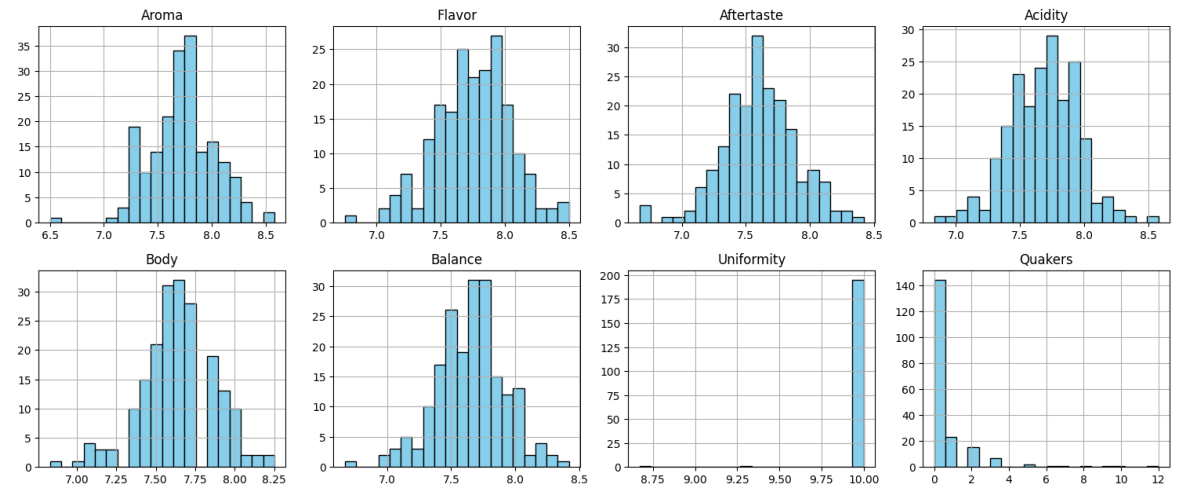| Date | 10 July 2024 |
|---|---|
| Team ID | Team-740058 |
| Project Title | Masterful Machines: Precise Coffee Quality Predictions Throogh ML |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Report**

The coffee quality dataset was explored to understand distributions, correlations, and missing data patterns. Key preprocessing steps included handling missing values, outlier treatment, feature scaling, and encoding categorical variables. The data was then split into training and testing sets, ensuring readiness for precise coffee quality predictions using machine learning models.

| Section | Description |
|---|---|
| Data Overview | Dimension:<br>207 rows × 19 columns<br>Descriptive statistics:<br> |
| Univariate Analysis | |

4o
.

## Histograms of Coffee Quality Scores



Histograms of Coffee Quality Scores (Aroma, Flavor, Aftertaste, Acidity, Body, Balance, Uniformity, Quakers)

**Bivariate Analysis**



Scatter plot of Flavor and Aftertaste Scores



Scatter plot: Aroma vs. Flavor

| | |
|---|---|
| Multivariate Analysis | 
Ratios of Coffee Bean Colors |
| Outliers and Anomalies |  |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Loading Data |  |
| Handling Missing Data |  |

**Loading Data**

```python
df = pd.read_csv("/content/beans_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 Anaerobic Washed | 8.58 | 8.50 | 8.42 | 8.58 | 8.25 | 8.42 | 10.0 | 8.58 | 89.33 | 11.8 | 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 | 11.8 | 0 | 0 | green | 0 |
| 4 | 4 | 2 | 24 kg | Red Bourbon | Honey,Mossto | 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 | 80.08 | 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 | 80.00 | 10.4 | 0 | 2 | green | 12 |
| 204 | 204 | 343 | 60 kg | Catimor | Washed / Wet | 7.25 | 7.17 | 7.08 | 7.00 | 7.08 | 7.08 | 10.0 | 7.00 | 79.67 | 11.6 | 0 | 9 | green | 11 |
| 205 | 205 | 1 | 2 kg | Maragogype | Natural / Dry | 6.50 | 6.75 | 6.75 | 7.17 | 7.08 | 7.00 | 10.0 | 6.83 | 78.08 | 11.0 | 0 | 12 | bluish-green | 13 |
| 206 | 206 | 600 | 60 kg | Mundo Novo | SEMI-LAVADO | 7.25 | 7.08 | 6.67 | 6.83 | 6.83 | 6.67 | 10.0 | 6.67 | 78.00 | 11.3 | 0 | 0 | green | 1 |

207 rows × 19 columns

**Handling Missing Data**

```python
df1.isna().sum()
```

```
Aroma                   0
Flavor                  0
Aftertaste              0
Acidity                 0
Body                    0
Balance                 0
Uniformity              0
Category One Defects    0
Quakers                 0
Color                   0
Category Two Defects    0
dtype: int64
```

```python
df.duplicated().sum()
```

```
0
```

```python
df['Color'].value_counts()
```

```
Color
green            98
greenish         34
bluish-green     19
blue-green       11
yellow-green      9
brownish          8
pale yellow       6
yellow green      5
yellowish         4
yellow- green     1
browish-green     1
yello-green       1
Name: count, dtype: int64
```

| | |
|---|---|
| Data Transformation | ```python\ndf1['Bean_Status']='Healthy'\ncondition_healthy=(df1['Category One Defects']==0) & (df1['Category Two Defects']==0)\ndf1.loc[condition_healthy,'Bean_Status']='Healthy'\ncondition_unhealthy=(df1['Category One Defects']!=0) & (df1['Category Two Defects']!=0)\ndf1.loc[condition_unhealthy,'Bean_Status']='Unhealthy'\n```<br><br>`[ ] df1`<br><br><table><tr><td></td><td>Aroma</td><td>Flavor</td><td>Aftertaste</td><td>Acidity</td><td>Body</td><td>Balance</td><td>Uniformity</td><td>Category One Defects</td><td>Quakers</td><td>Category Two Defects</td><td>Color_Encoded</td><td>Bean_Status</td></tr><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><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><td>Healthy</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><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><td>Healthy</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><td>Healthy</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><td>Healthy</td></tr><tr><td>205</td><td>6.50</td><td>6.75</td><td>6.75</td><td>7.17</td><td>7.08</td><td>7.00</td><td>10.0</td><td>0</td><td>12</td><td>13</td><td>1</td><td>Healthy</td></tr><tr><td>206</td><td>7.25</td><td>7.08</td><td>6.67</td><td>6.83</td><td>6.83</td><td>6.67</td><td>10.0</td><td>0</td><td>0</td><td>1</td><td>4</td><td>Healthy</td></tr></table>197 rows × 12 columns<br><br>`[ ] df1['Bean_Status'].value_counts()`<br><br>```\nBean_Status\nHealthy     186\nUnhealthy    11\nName: count, dtype: int64\n``` |
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | - |