
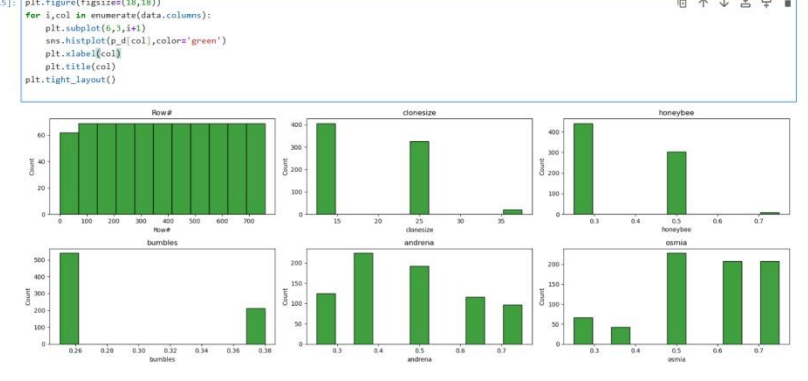


Data Collection and Preprocessing Phase

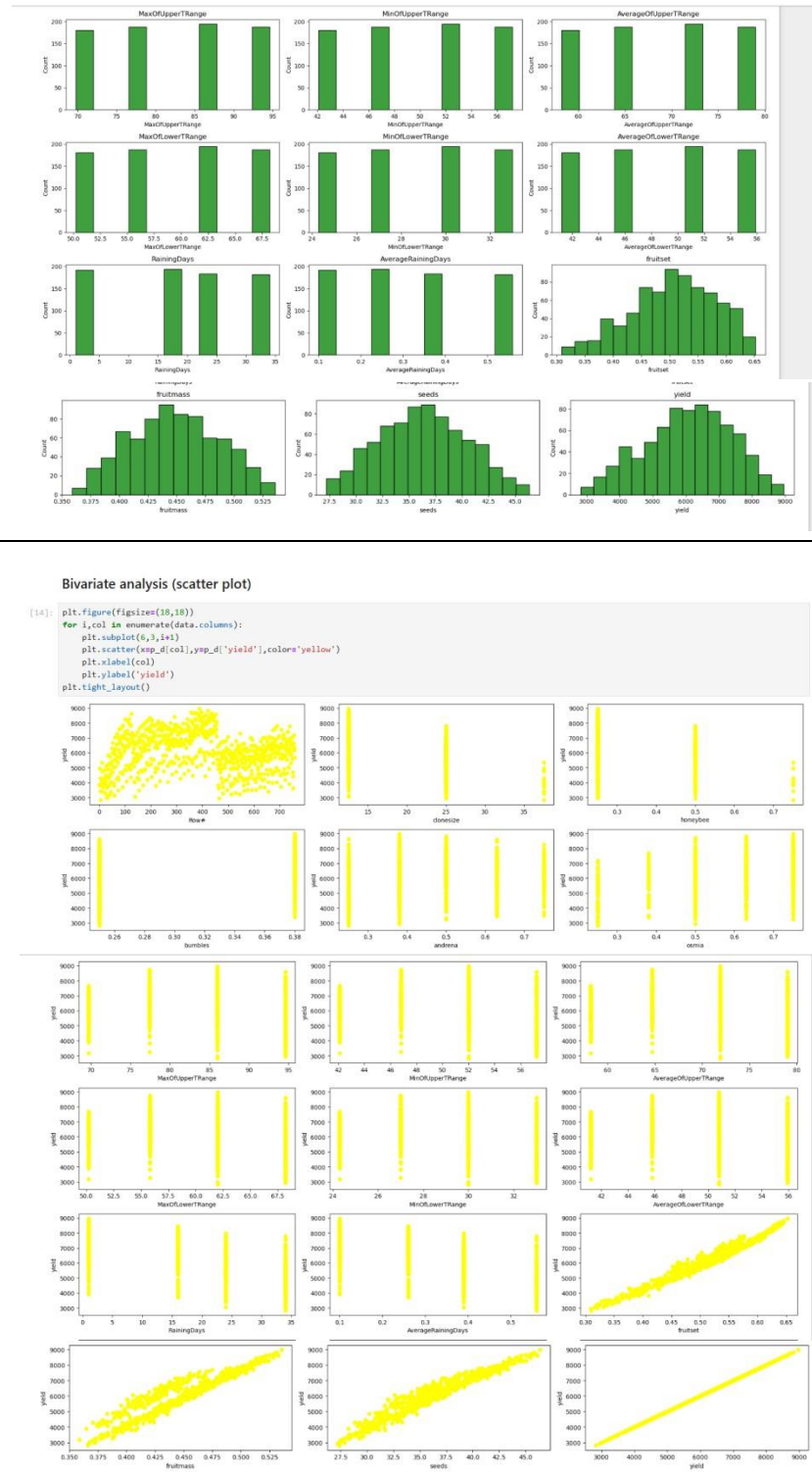
Date	06 July 2024
Team ID	739665
Project Title	BlueBerry Yield Prediction
Maximum Marks	6 Marks

Data Exploration and Preprocessing Report

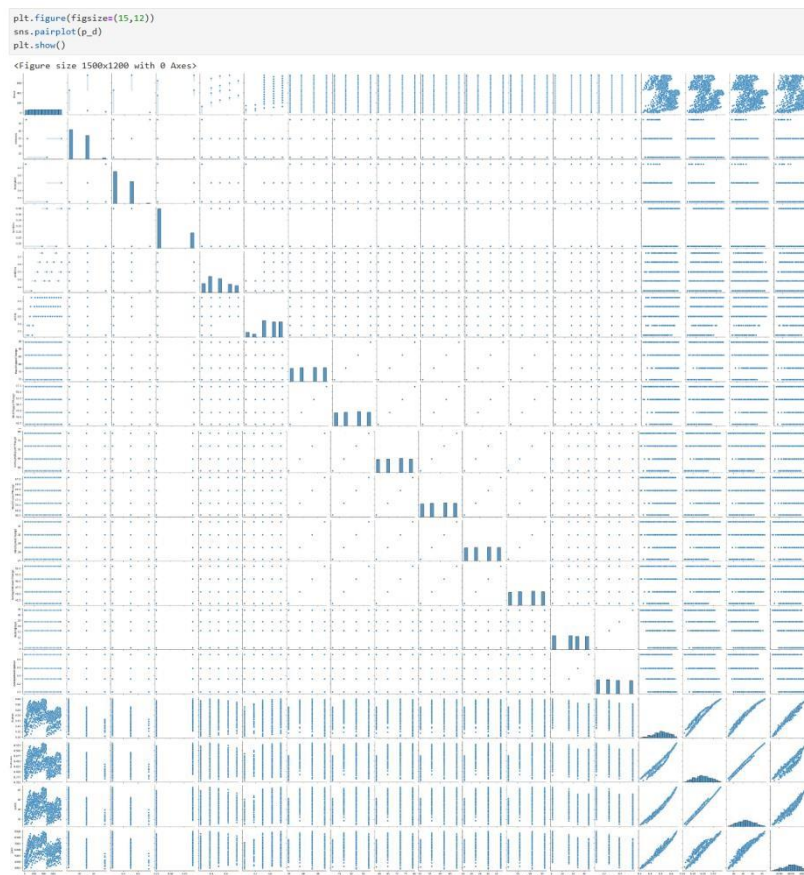
Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

Section	Description
Data Overview	<pre>[14]: p_d.describe()</pre> 
Univariate Analysis	<pre>[15]: plt.figure(figsize=(18,18)) for i,col in enumerate(data.columns): plt.subplot(6,3,i+1) sns.histplot(p_d[col],color='green') plt.xlabel(col) plt.title(col) plt.tight_layout()</pre> 

Bivariate Analysis



Multivariate Analysis



Data Preprocessing Code Screenshots

Loading Data

```
data=pd.read_csv("WildBlueberryPollinationSimulationData.csv")
data
```

	Row#	clonesize	honeybee	bumbles	andrena	osmia	MaxOfUpperTrange	MinOfUpperTrange	AverageOfUpperTrange	MaxOfLowerTrange	MinOfLowerTrange
	0	0	37.5	0.750	0.250	0.250	86.0	52.0	71.9	62.0	30.0
	1	1	37.5	0.750	0.250	0.250	86.0	52.0	71.9	62.0	30.0
	2	2	37.5	0.750	0.250	0.250	94.6	57.2	79.0	68.2	33.0
	3	3	37.5	0.750	0.250	0.250	84.6	57.2	70.8	68.2	33.0

Handling Null Values

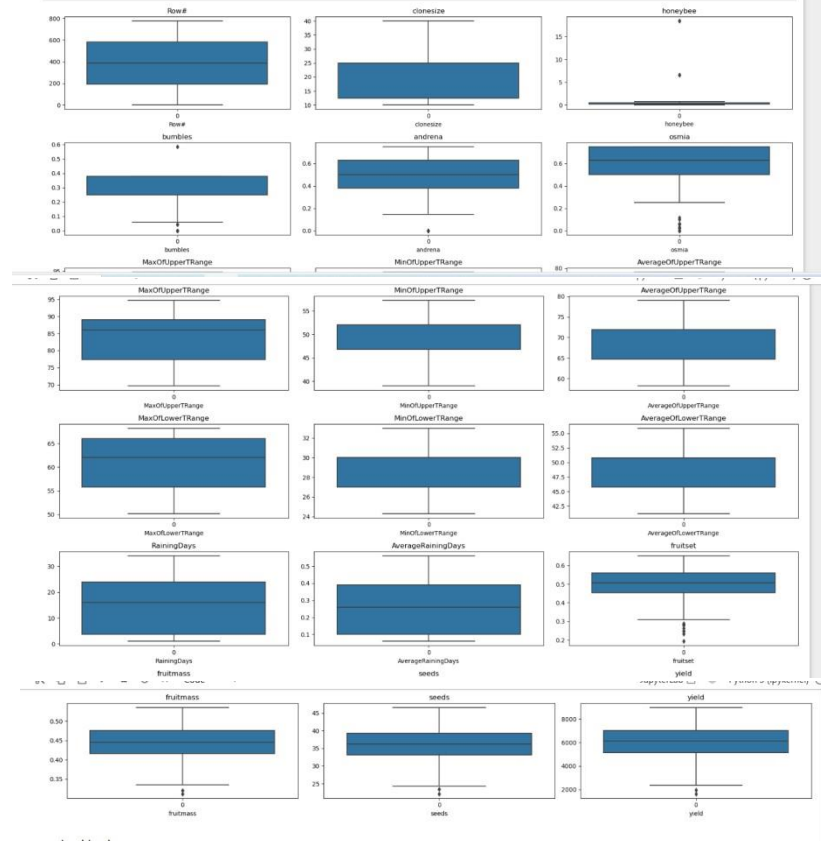
Handling null values

```
[8]: data.isnull().sum()
```

```
[8]: Row#              0
      clonesize        0
      honeybee         0
      bumbles          0
      andrena          0
      osmia            0
      MaxOfUpperTRange  0
      MinOfUpperTRange  0
      AverageOfUpperTRange  0
      MaxOfLowerTRange  0
      MinOfLowerTRange  0
      AverageOfLowerTRange  0
      RainingDays       0
      AverageRainingDays  0
      fruitset          0
      fruitmass         0
      seeds            0
      yield            0
      dtype: int64
```

veiwng imbalanced data

```
[9]: plt.figure(figsize=(18,18))
    for i,col in enumerate(data.columns):
        plt.subplot(6,3,i+1)
        sns.boxplot(data[col])
        plt.xlabel(col)
        plt.title(col)
    plt.tight_layout()
```



<p>Handling outliers</p>	<div> <h3>handling imbalance data</h3> <p>by removing outliers</p> <pre>[223]: x=data q1=x.quantile(0.25) q3=x.quantile(0.75) iqr=q3-q1 iqr</pre> <pre>[223]: Row# 388.000000 clonesize 12.500000 honeybee 0.250000 bumbles 0.130000 andrena 0.250000 osmia 0.250000 MaxOfUpperTRange 11.600000 MinOfUpperTRange 5.200000 AverageOfUpperTRange 7.200000 MaxOfLowerTRange 10.200000 MinOfLowerTRange 3.000000 AverageOfLowerTRange 5.000000 RainingDays 20.230000 AverageRainingDays 0.290000 fruitset 0.106571 fruitmass 0.059869 seeds 6.123577 yield 1897.334830 dtype: float64</pre> </div>
<p>Saved Processed Data</p>	<pre>p_d=data[~((data<(q1-1.5*iqr)) (data>(q3+1.5*iqr))).any(axis=1)] p_d.shape</pre> <p>(752, 18)</p>