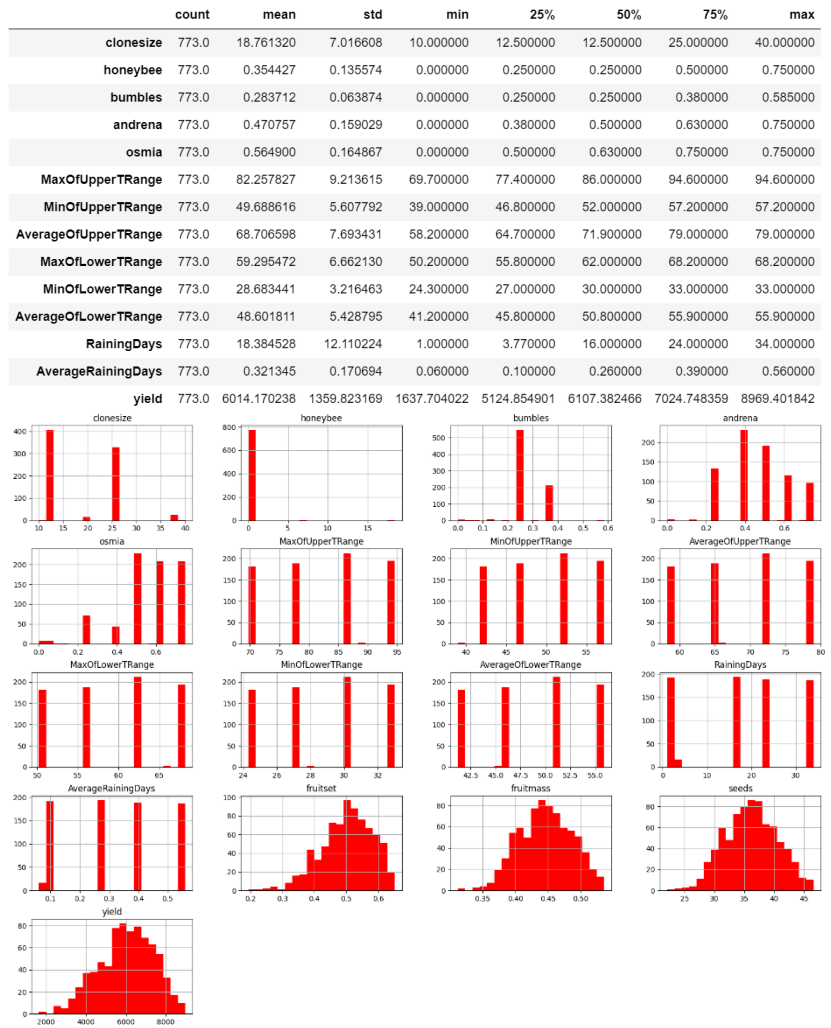


Data Collection and Preprocessing Phase

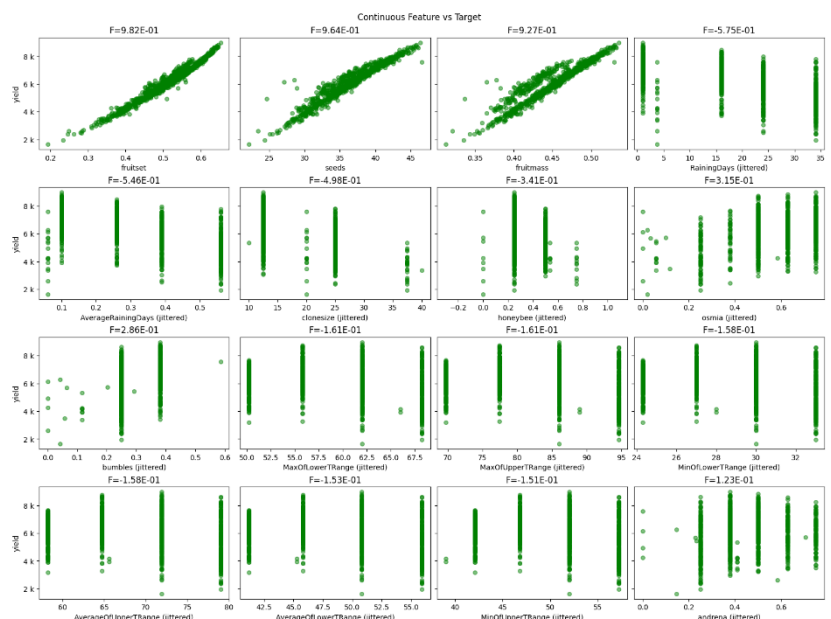
Date	12 July 2024
Team ID	SWTID1720077079
Project Title	Wild Blueberry Yield Prediction
Maximum Marks	6 Marks

Section	Description																																																																								
Data Overview	Data columns (total 17 columns):																																																																								
	<table><tr><th>#</th><th>Column</th><th>Non-Null Count</th><th>Dtype</th></tr><tr><td>0</td><td>clonesize</td><td>777 non-null</td><td>float64</td></tr><tr><td>1</td><td>honeybee</td><td>777 non-null</td><td>float64</td></tr><tr><td>2</td><td>bumbles</td><td>777 non-null</td><td>float64</td></tr><tr><td>3</td><td>andrena</td><td>777 non-null</td><td>float64</td></tr><tr><td>4</td><td>osmia</td><td>777 non-null</td><td>float64</td></tr><tr><td>5</td><td>MaxOfUpperTRange</td><td>777 non-null</td><td>float64</td></tr><tr><td>6</td><td>MinOfUpperTRange</td><td>777 non-null</td><td>float64</td></tr><tr><td>7</td><td>AverageOfUpperTRange</td><td>777 non-null</td><td>float64</td></tr><tr><td>8</td><td>MaxOfLowerTRange</td><td>777 non-null</td><td>float64</td></tr><tr><td>9</td><td>MinOfLowerTRange</td><td>777 non-null</td><td>float64</td></tr><tr><td>10</td><td>AverageOfLowerTRange</td><td>777 non-null</td><td>float64</td></tr><tr><td>11</td><td>RainingDays</td><td>777 non-null</td><td>float64</td></tr><tr><td>12</td><td>AverageRainingDays</td><td>777 non-null</td><td>float64</td></tr><tr><td>13</td><td>fruitset</td><td>777 non-null</td><td>float64</td></tr><tr><td>14</td><td>fruitmass</td><td>777 non-null</td><td>float64</td></tr><tr><td>15</td><td>seeds</td><td>777 non-null</td><td>float64</td></tr><tr><td>16</td><td>yield</td><td>777 non-null</td><td>float64</td></tr></table>	#	Column	Non-Null Count	Dtype	0	clonesize	777 non-null	float64	1	honeybee	777 non-null	float64	2	bumbles	777 non-null	float64	3	andrena	777 non-null	float64	4	osmia	777 non-null	float64	5	MaxOfUpperTRange	777 non-null	float64	6	MinOfUpperTRange	777 non-null	float64	7	AverageOfUpperTRange	777 non-null	float64	8	MaxOfLowerTRange	777 non-null	float64	9	MinOfLowerTRange	777 non-null	float64	10	AverageOfLowerTRange	777 non-null	float64	11	RainingDays	777 non-null	float64	12	AverageRainingDays	777 non-null	float64	13	fruitset	777 non-null	float64	14	fruitmass	777 non-null	float64	15	seeds	777 non-null	float64	16	yield	777 non-null	float64
	#	Column	Non-Null Count	Dtype																																																																					
	0	clonesize	777 non-null	float64																																																																					
	1	honeybee	777 non-null	float64																																																																					
	2	bumbles	777 non-null	float64																																																																					
	3	andrena	777 non-null	float64																																																																					
	4	osmia	777 non-null	float64																																																																					
	5	MaxOfUpperTRange	777 non-null	float64																																																																					
	6	MinOfUpperTRange	777 non-null	float64																																																																					
	7	AverageOfUpperTRange	777 non-null	float64																																																																					
	8	MaxOfLowerTRange	777 non-null	float64																																																																					
	9	MinOfLowerTRange	777 non-null	float64																																																																					
	10	AverageOfLowerTRange	777 non-null	float64																																																																					
	11	RainingDays	777 non-null	float64																																																																					
	12	AverageRainingDays	777 non-null	float64																																																																					
	13	fruitset	777 non-null	float64																																																																					
	14	fruitmass	777 non-null	float64																																																																					
	15	seeds	777 non-null	float64																																																																					
	16	yield	777 non-null	float64																																																																					
dtypes: float64(17)																																																																									

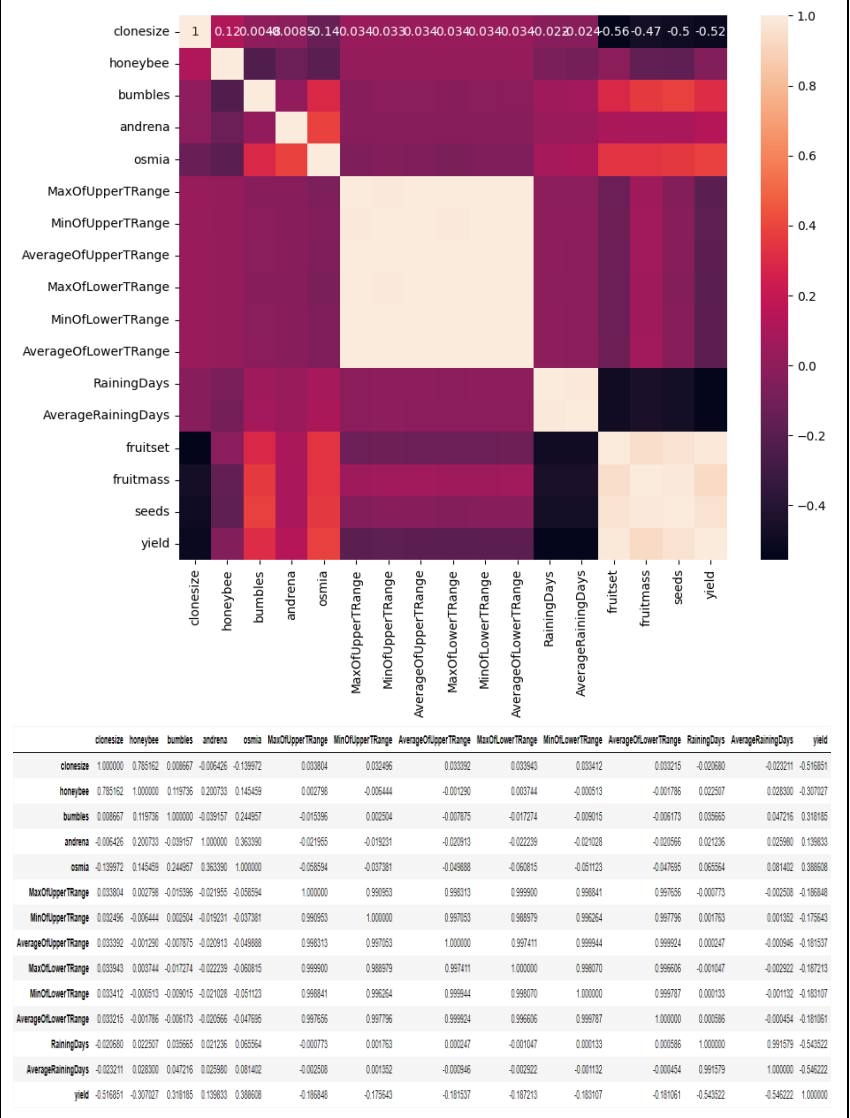
Univariate Analysis



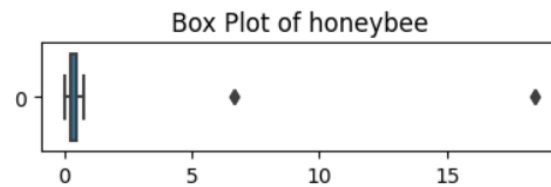
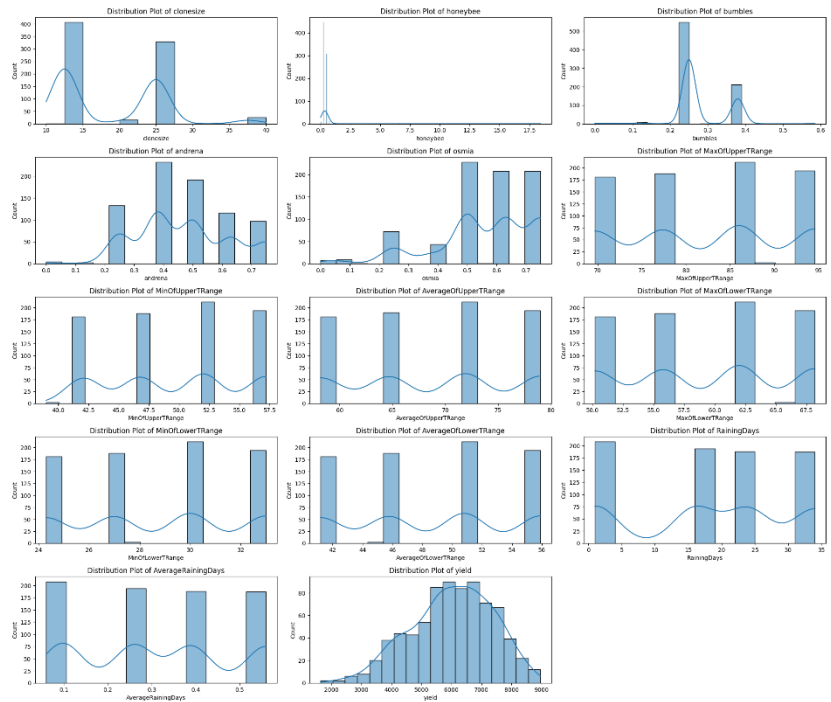
Bivariate Analysis



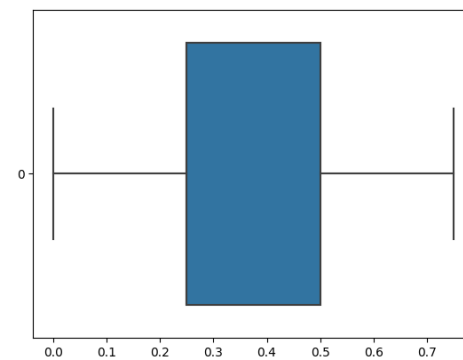
Multivariate Analysis



Outliers and Anomalies



Outlier in feature 'honeybee' found using boxplot



Handled outlier

Data Preprocessing Code Screenshots	
Loading Data	<pre>import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns import dabl</pre> <pre>data=pd.read_csv('WildBlueberryPollinationSimulationData.csv')</pre>
Handling Missing Data	<pre>data.isna().sum()</pre> <pre> 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 </pre>
Data Transformation	<pre>from sklearn.preprocessing import StandardScaler scale = StandardScaler()</pre> <pre>X_scaled=scale.fit_transform(X)</pre> <pre>X_scaled</pre> <pre> array([[2.67234719, 2.91964747, -0.52812593, ..., 0.40517505, -0.19702952, -0.35962034], [2.67234719, 2.91964747, -0.52812593, ..., 0.40517505, -1.43645427, -1.29757569], [2.67234719, 2.91964747, -0.52812593, ..., 1.34521837, -0.19702952, -0.35962034], ..., [0.17664981, 1.34753621, -2.61170931, ..., 0.40517505, 0.46399701, 0.40246839], [0.17664981, 1.34753621, -2.61170931, ..., -0.60859715, -1.20757383, -1.53206453], [0.17664981, 1.34753621, -2.61170931, ..., -0.60859715, 0.46399701, 0.40246839]]) </pre>
Feature Engineering	For handling outlier

	<pre>Q1 = data['honeybee'].quantile(0.25) Q3 = data['honeybee'].quantile(0.75) IQR=Q3-Q1 lower_limit = Q1 - 1.5 * IQR upper_limit = Q3 + 1.5 * IQR print('lower_limit: ',lower_limit) print('upper_limit: ',upper_limit) data = data[(data.honeybee>lower_limit)&(data.honeybee<upper_limit)]</pre>
Save Processed Data	<pre>X=pd.DataFrame(X_scaled, columns=names)</pre> <p>X</p> <p>Saving the scaler</p> <pre>with open('standard_scaler.pkl', 'wb') as file: pickle.dump(scale, file)</pre>