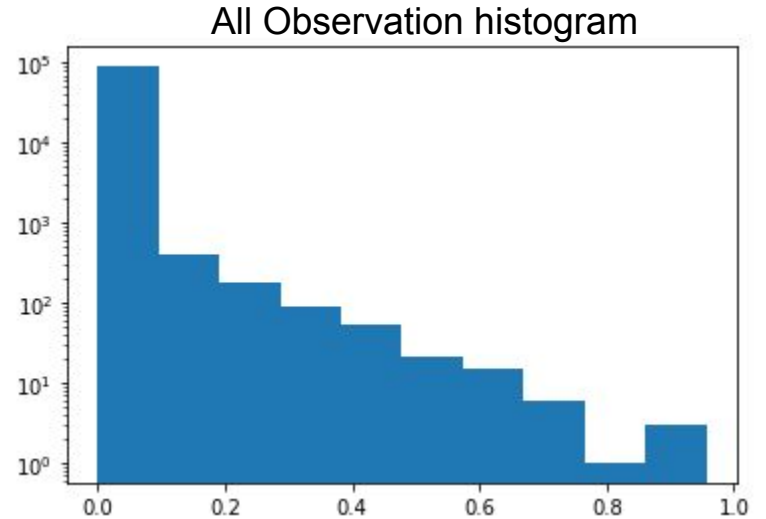- Matrix Size: 250*361
- Total Observation: 90250
- Top 0.5% CutOff: 0.16

Very Sparse Matrix, after processing, there are 451 observations of 1, the rest are set to 0.
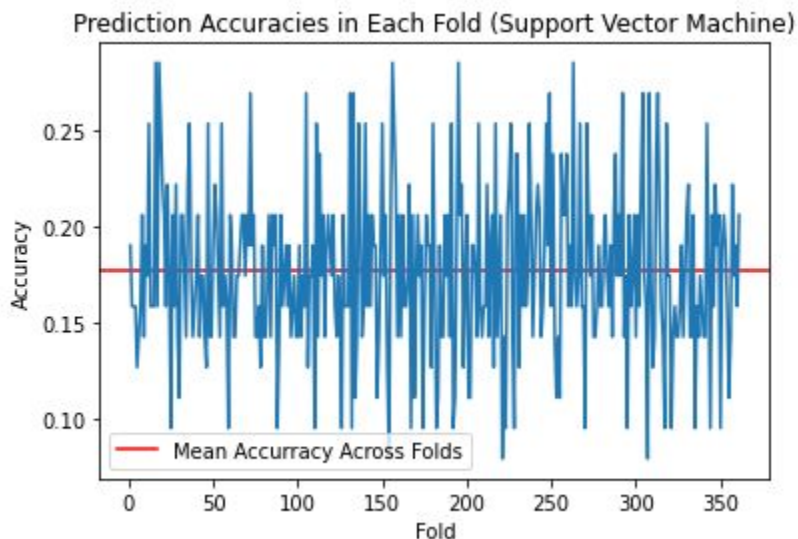


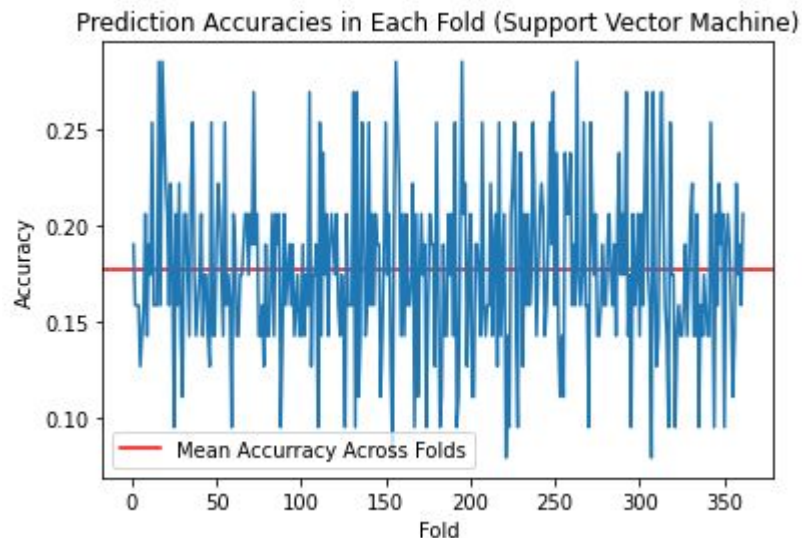All Observation histogram

# Plant Species

13 Unique Species, 250 samples total, Average 19 samples per label.

{'POPR', 'ELSC', 'POLE', 'TRSP', 'ELEL', 'ELTR', 'POST', 'FETH', 'FEBR', 'FESA', 'POAL', 'ACNE', 'ACLE'}

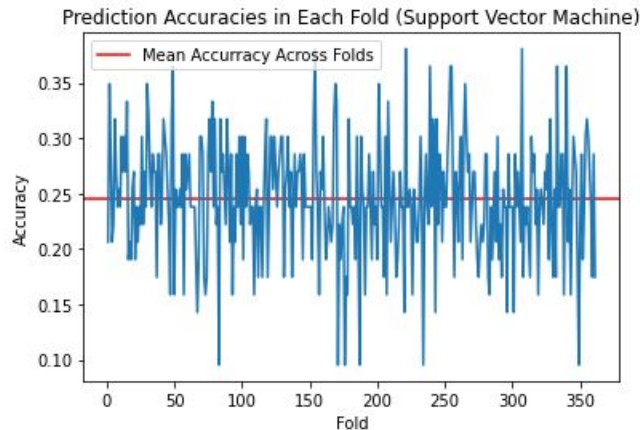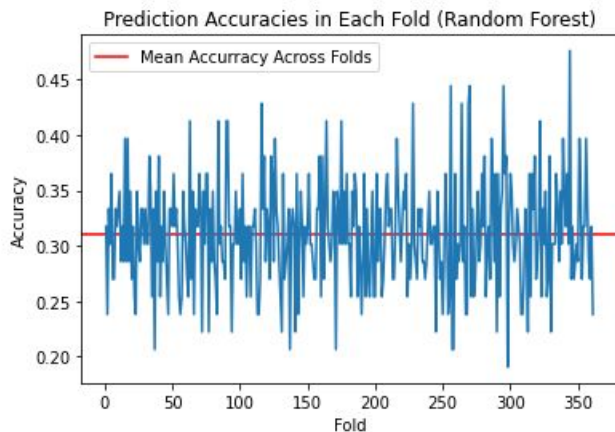Continuous Measurement

Binary Matrix

# Plant Species (grouped by genus)
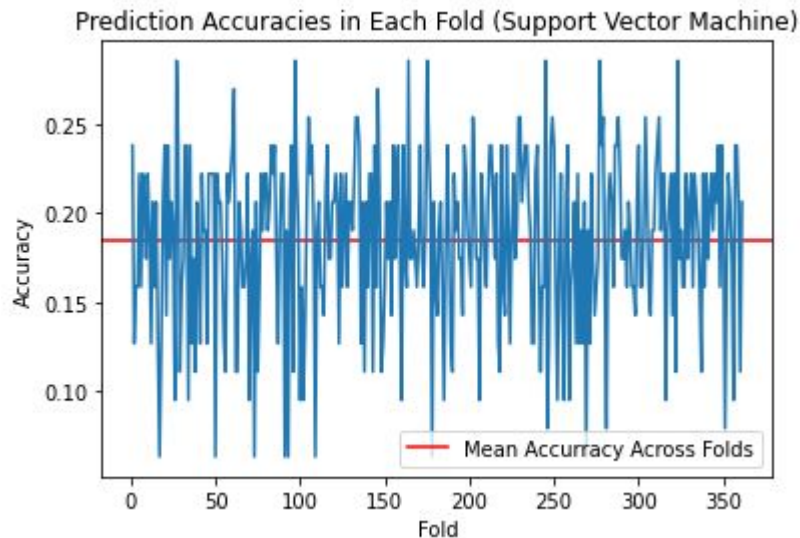
{'PO', 'TR', 'AC', 'FE', 'EL'}

Random Guess Accuracy: 20%

Mean prediction Accuracy using SVM: 24.57%

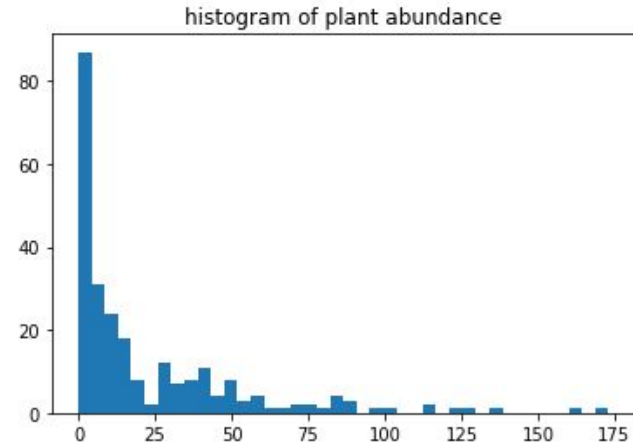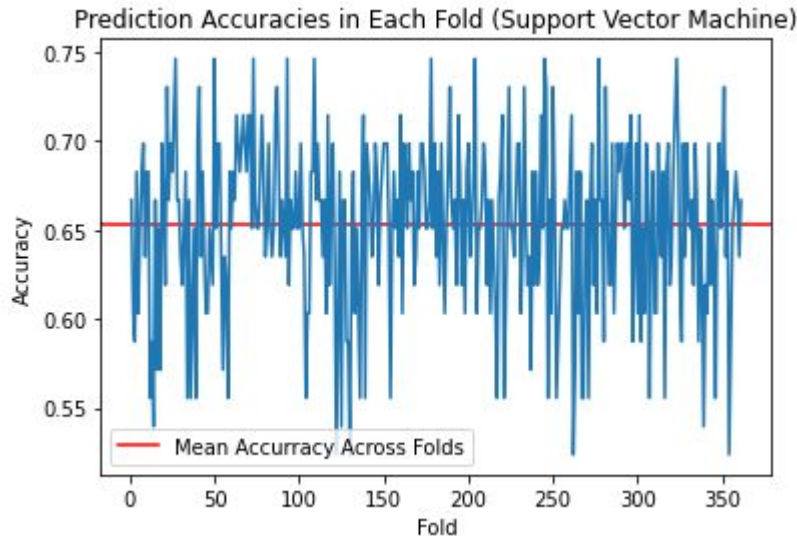Mean prediction Accuracy using Random Forest: 31%



Prediction Accuracies in Each Fold (Random Forest)



Prediction Accuracies in Each Fold (Support Vector Machine)

# Gradient

{'Avery', 'Treasury', 'Ruby', 'HunterHill', 'RP', 'Cinnamon', 'KP', 'JC2', 'Teocali', 'JC1'}



Prediction Accuracies in Each Fold (Support Vector Machine)

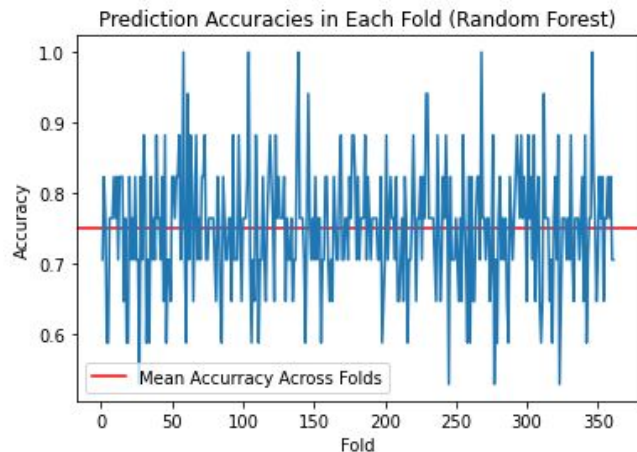# Plant Abundance (Ignoring Species and Genus)

Continuous measurements, generalized into 2 class: {Abundant, Not Abundant}
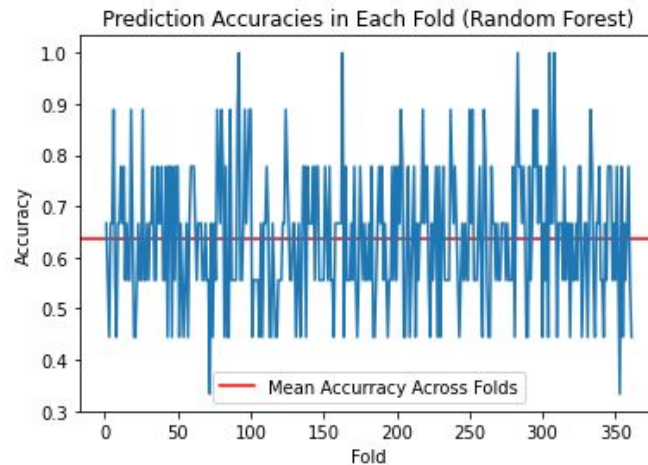
Currently using 5 as the cutoff.
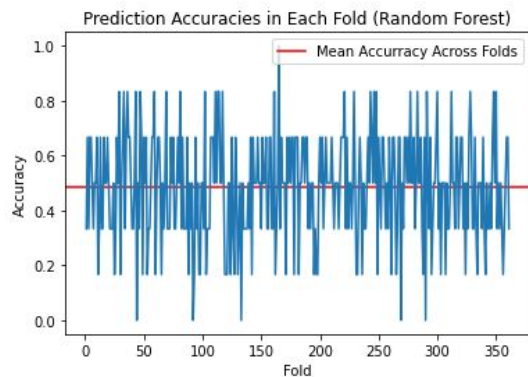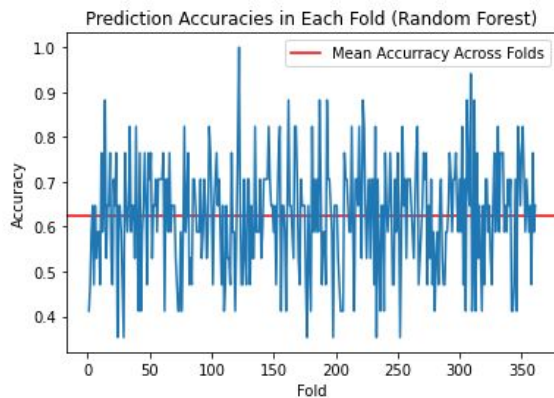
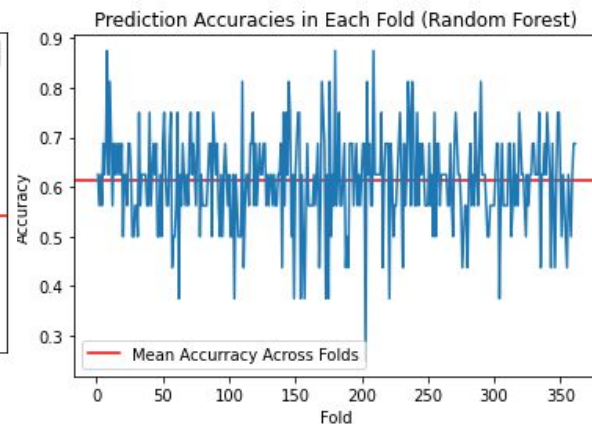# Plant Abundance by Genus

PO



TR

# Plant Abundance by Genus



AC

Prediction Accuracies in Each Fold (Random Forest)

FE

Prediction Accuracies in Each Fold (Random Forest)

EL

Prediction Accuracies in Each Fold (Random Forest)

# Plant Abundance by Genus

Mean Accuracy Prediction by Genus:

PO: 0.751

TR: 0.638

FE: 0.623

EL: 0.6125346260387812

AC: 0.483

# POPR

Sample Size is small (20), and extremely unbalanced:

SVM prediction with all ASVs:

**Real Labels:**

['Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Not Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Not Abundant']
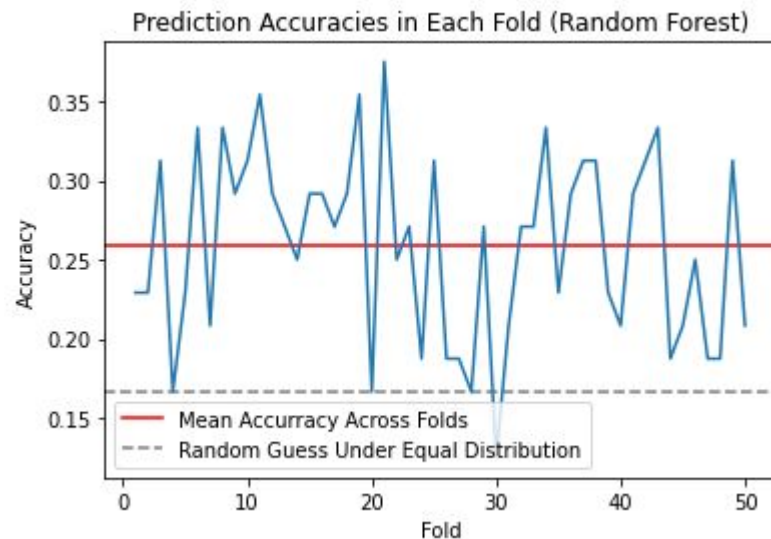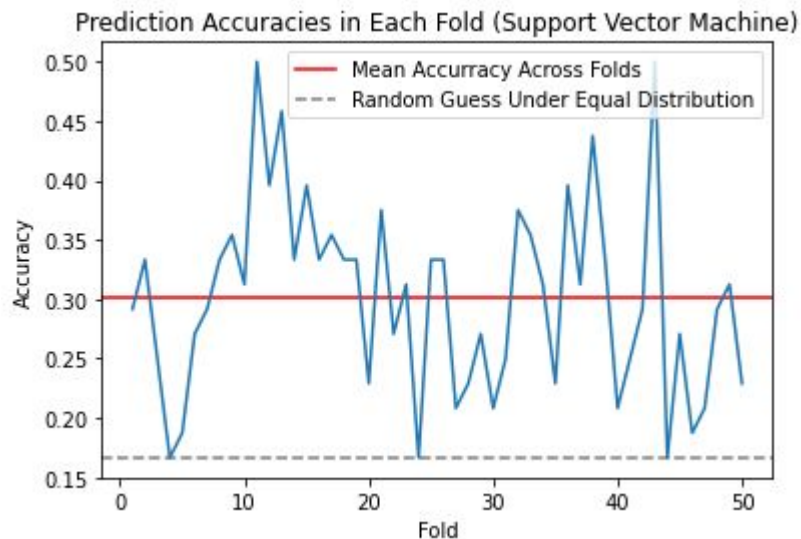
**Predictions:**

['Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant', 'Abundant']

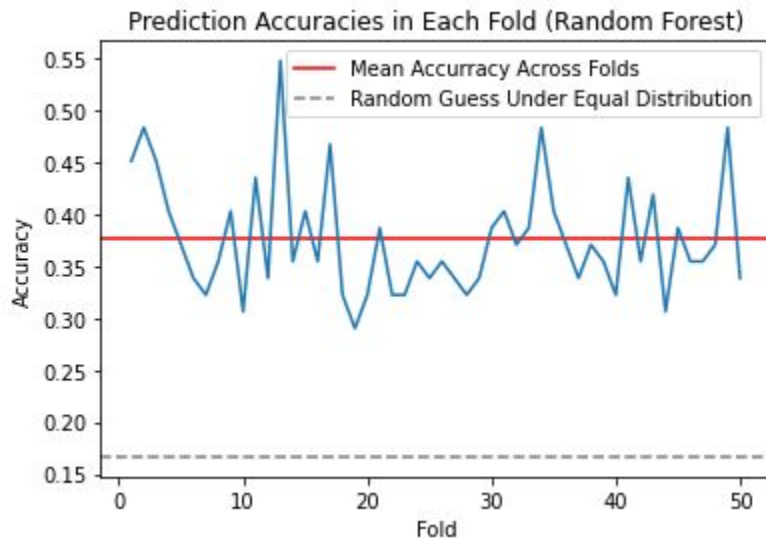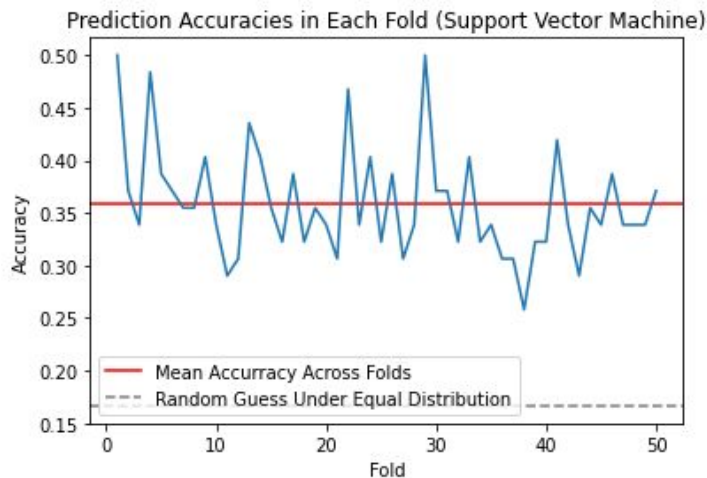# Species Prediction (One VS Others)

SVM, NB and RF all yield meaningless predictions (no true positives among all species and samples). With one false positive in the prediction of ELTR nad FETH.

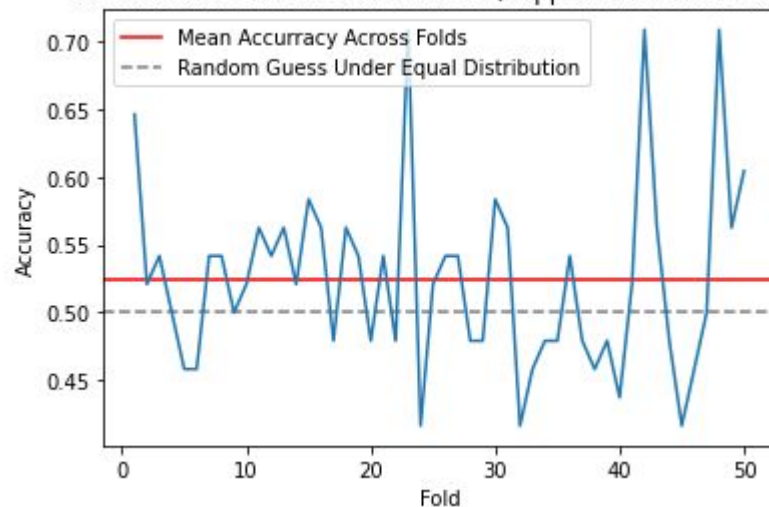Most predictions are completely all negatives.

# Leaves-Gradient



Prediction Accuracies in Each Fold (Support Vector Machine)

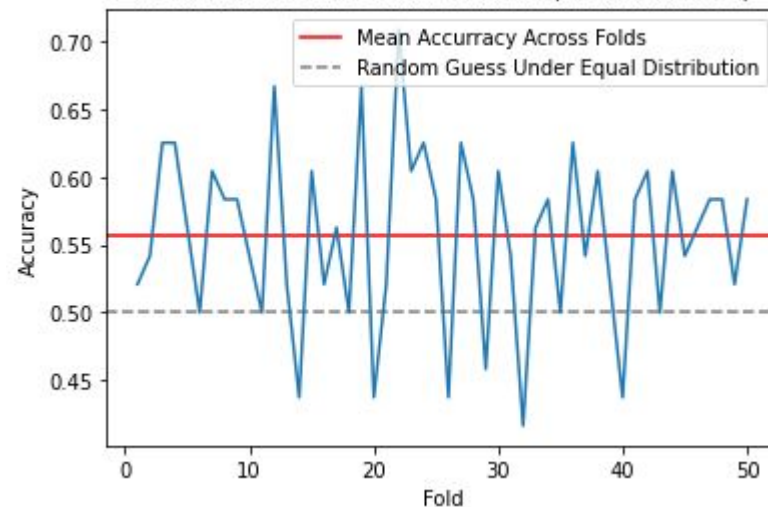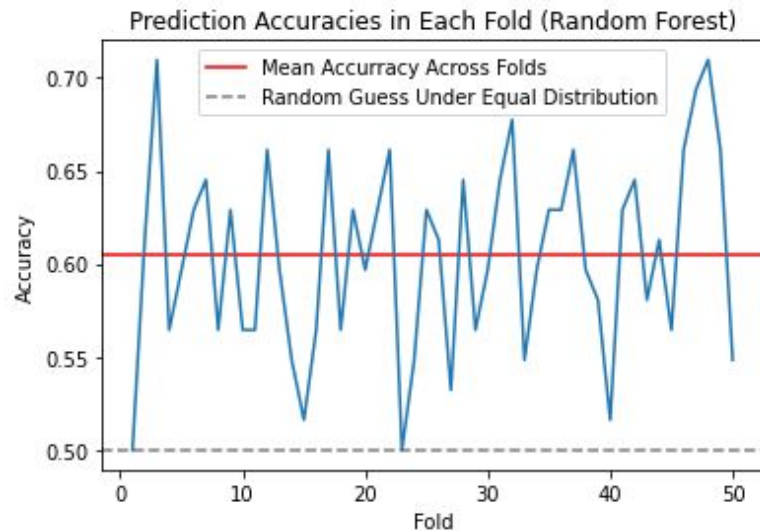Prediction Accuracies in Each Fold (Random Forest)

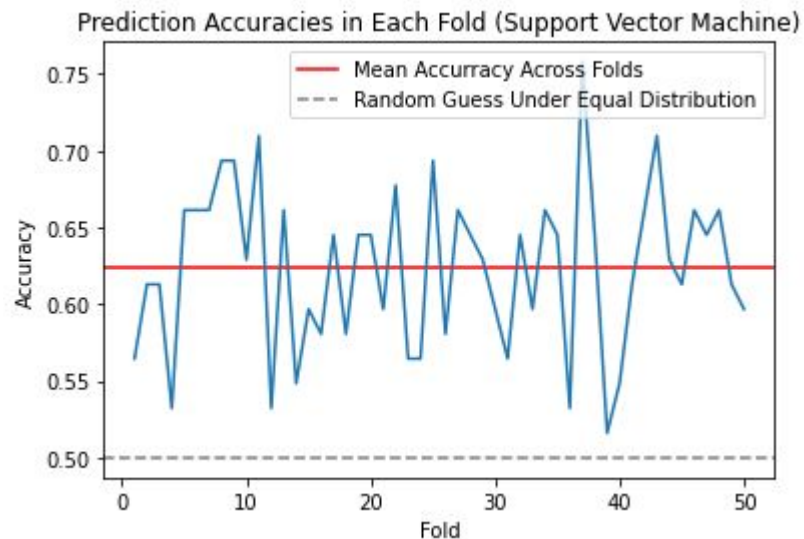# Root-Gradient

# Leave-Richness



Prediction Accuracies in Each Fold (Support Vector Machine)
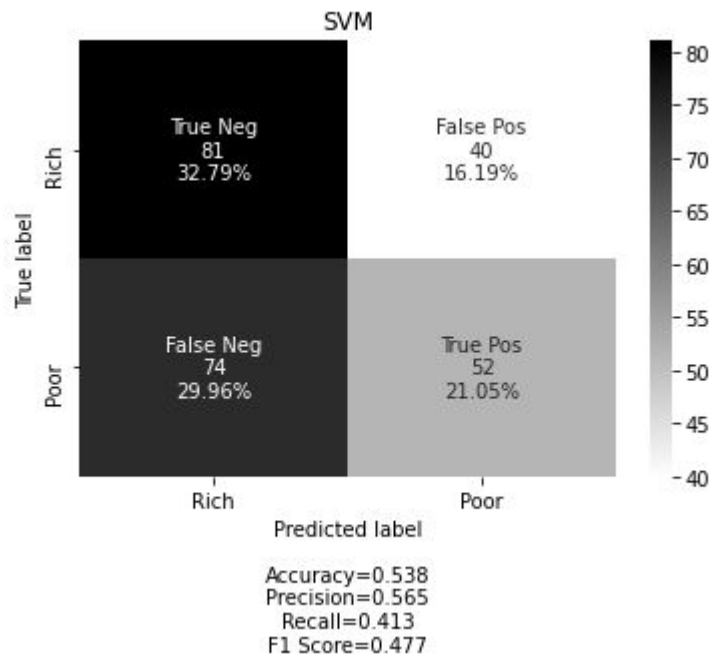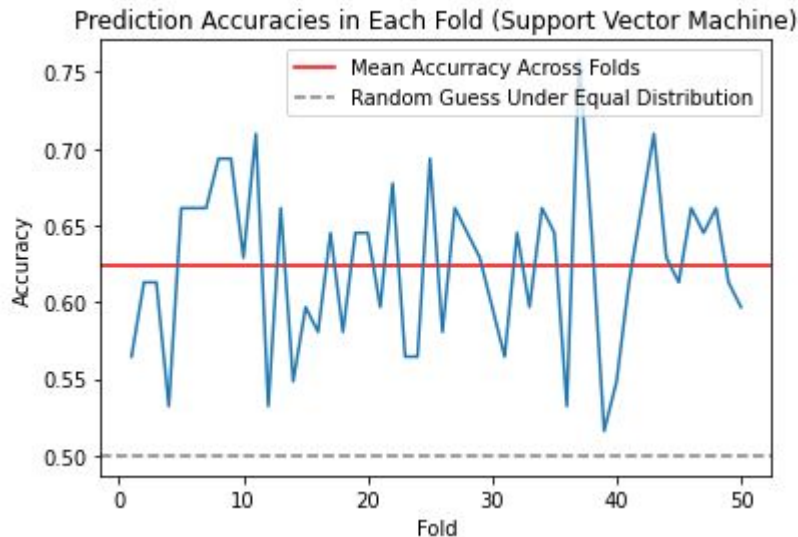
Prediction Accuracies in Each Fold (Random Forest)
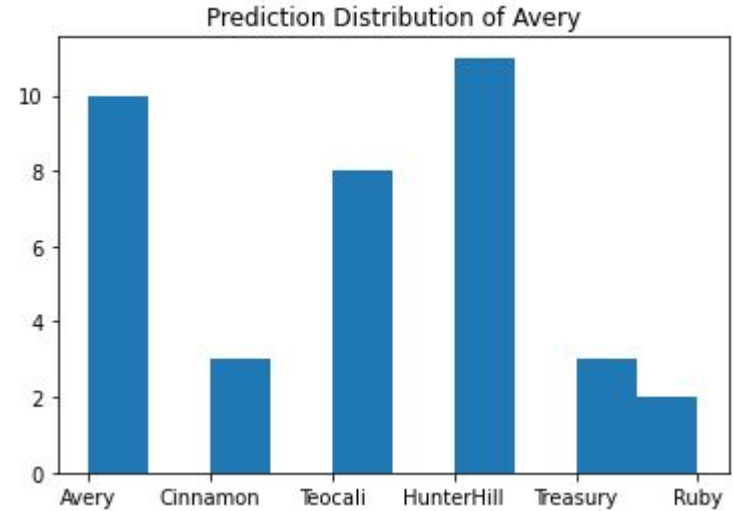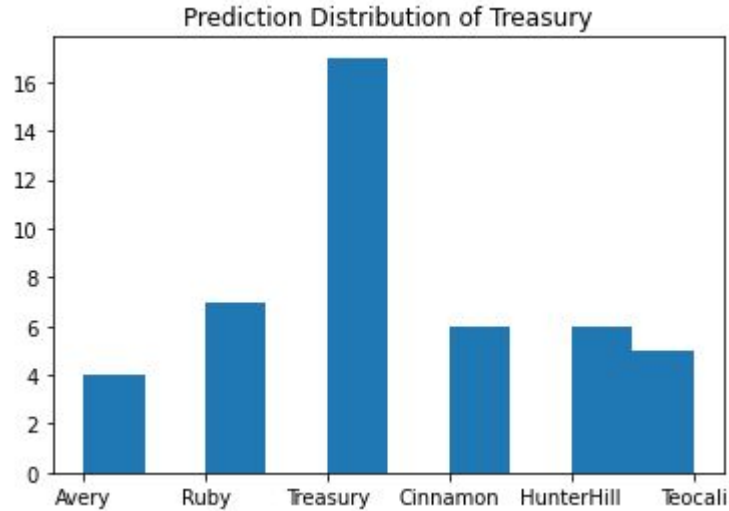
# Root-Richness

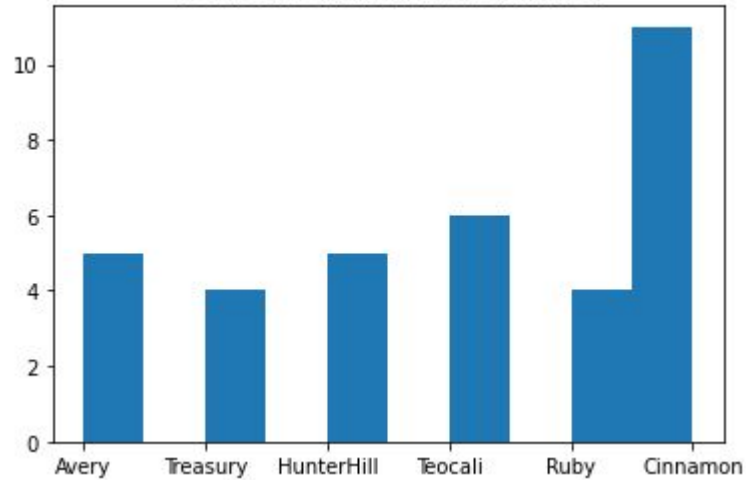# Richness using total-N (seeing minor signal so far)

Using Top 30 Features

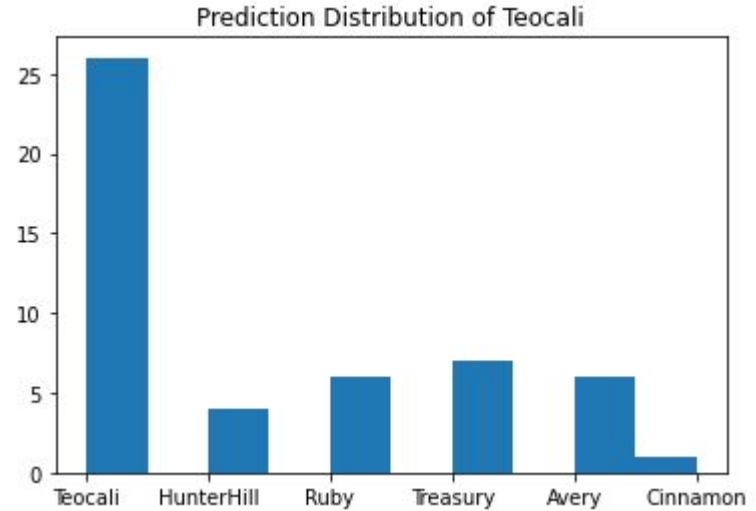# Gradient Prediction By Groups
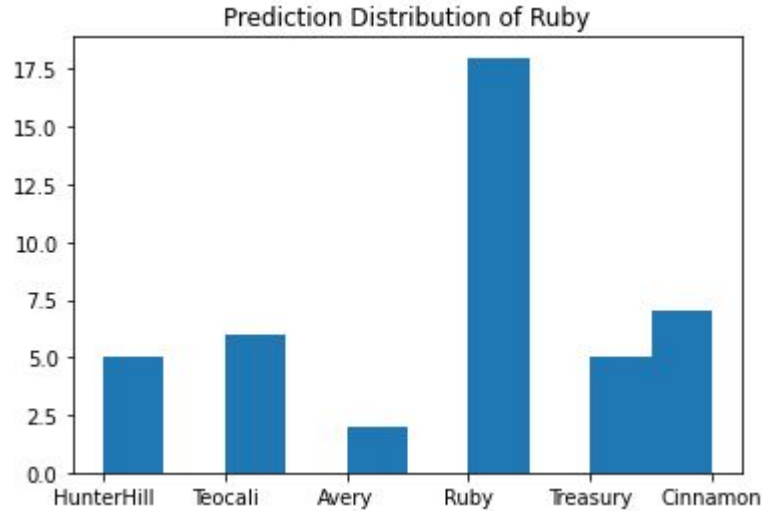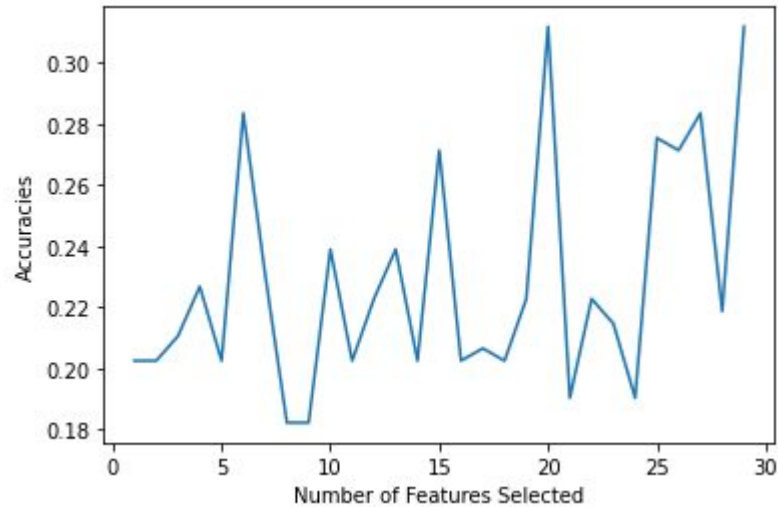
# Gradient Prediction By Species

# Gradient Prediction By Species

# Feature Selection on Gradient

# Treasury

Mean Accuracy Across Folds 0.7887096774193548



Prediction Accuracies in Each Fold (Support Vector Machine)

One vs Rest without FS, for this plot, its binary labels of Treasury vs Not Treasury

# Cinnamon :

Mean Accuracy Across Folds 0.8338709677419354



Prediction Accuracies in Each Fold (Support Vector Machine)

# Avery

Mean Accuracy Across Folds 0.7741935483870968



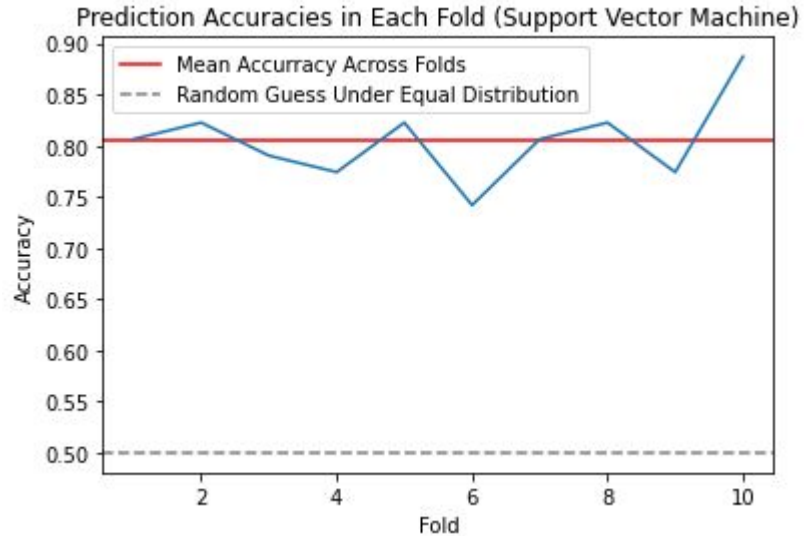Prediction Accuracies in Each Fold (Support Vector Machine)
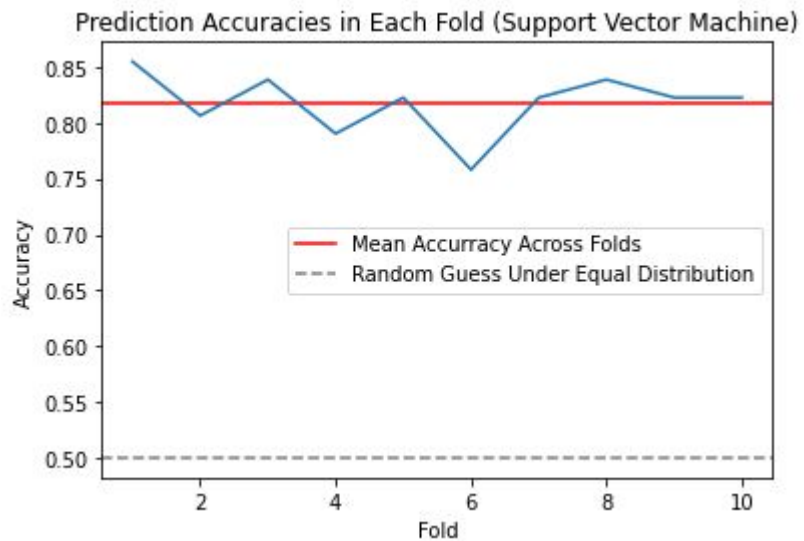
# Hunterhill

Mean Accuracy Across Folds 0.8483870967741935
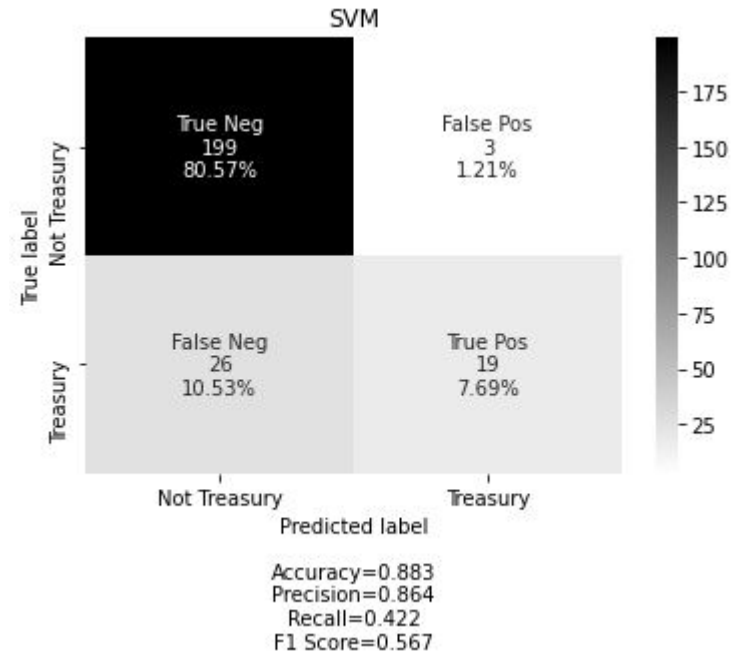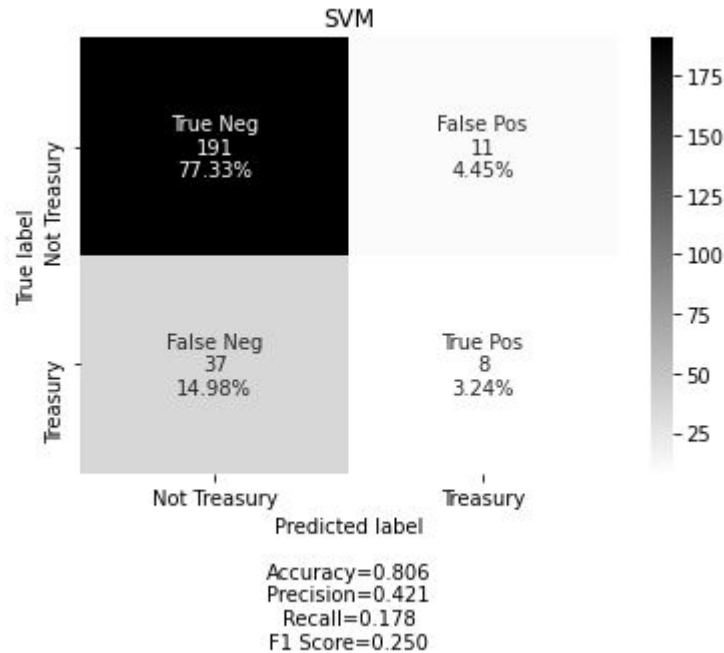
# Teocali

Mean Accuracy Across Folds 0.8048387096774194
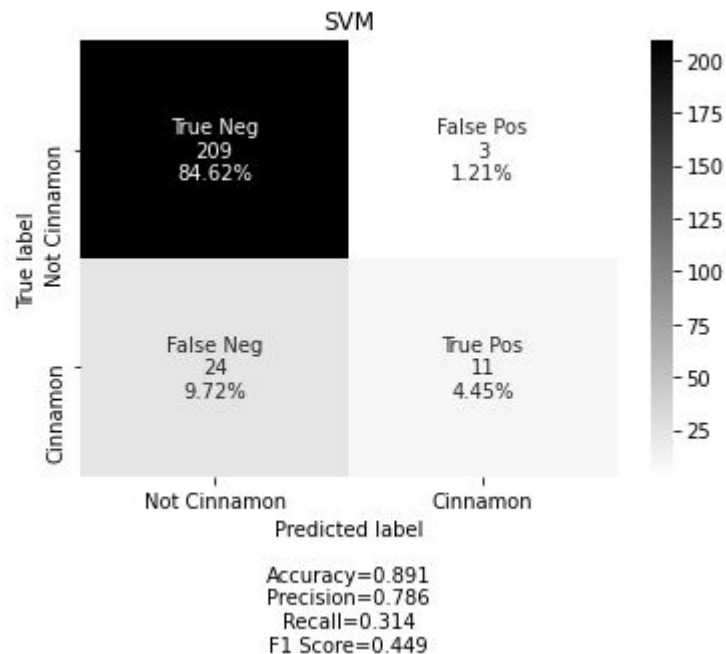
# Ruby

Mean Accuracy Across Folds 0.817741935483871



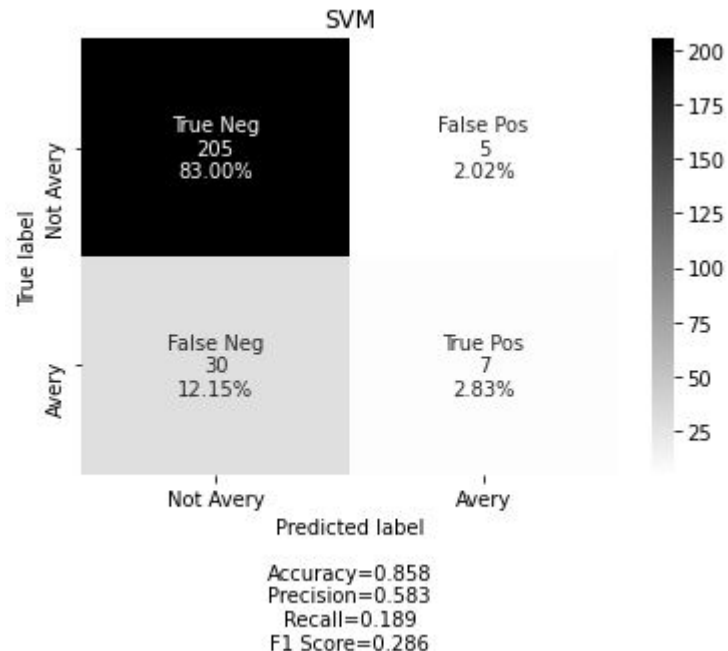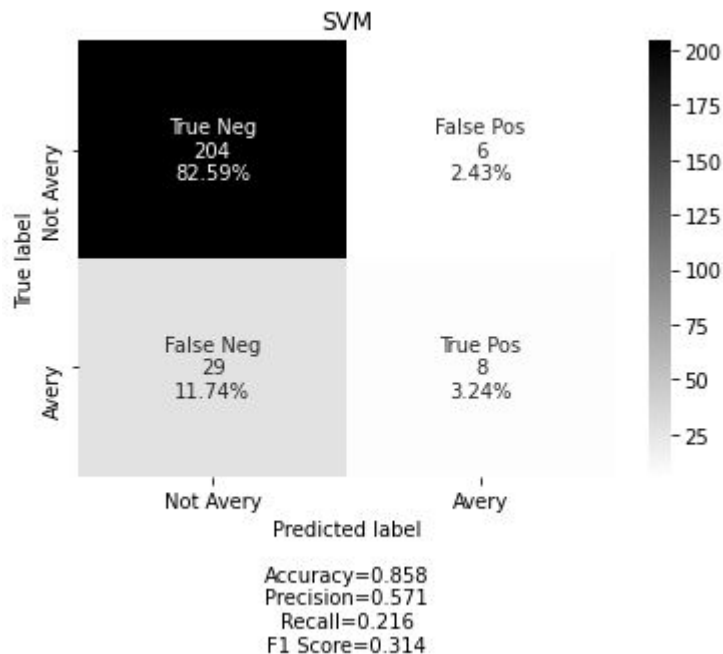Prediction Accuracies in Each Fold (Support Vector Machine)

Different Number of features impact different gradient prediction precision
Top 30 features (left) vs Top 100 features(right) for Treasury, select K best
with Chi2 scoring function.

# Cinnamon



SVM

|  | Not Cinnamon | Cinnamon |
|---|---|---|
| **Not Cinnamon** | True Neg 206 83.40% | False Pos 6 2.43% |
| **Cinnamon** | False Neg 27 10.93% | True Pos 8 3.24% |

Accuracy=0.866
Precision=0.571
Recall=0.229
F1 Score=0.327

SVM

|  | Not Cinnamon | Cinnamon |
|---|---|---|
| **Not Cinnamon** | True Neg 209 84.62% | False Pos 3 1.21% |
| **Cinnamon** | False Neg 24 9.72% | True Pos 11 4.45% |

Accuracy=0.891
Precision=0.786
Recall=0.314
F1 Score=0.449

# Avery



Left confusion matrix (SVM):

|  | Not Avery | Avery |
|---|---|---|
| **Not Avery** | True Neg 204 82.59% | False Pos 6 2.43% |
| **Avery** | False Neg 29 11.74% | True Pos 8 3.24% |

Accuracy=0.858
Precision=0.571
Recall=0.216
F1 Score=0.314

Right confusion matrix (SVM):

|  | Not Avery | Avery |
|---|---|---|
| **Not Avery** | True Neg 205 83.00% | False Pos 5 2.02% |
| **Avery** | False Neg 30 12.15% | True Pos 7 2.83% |

Accuracy=0.858
Precision=0.583
Recall=0.189
F1 Score=0.286
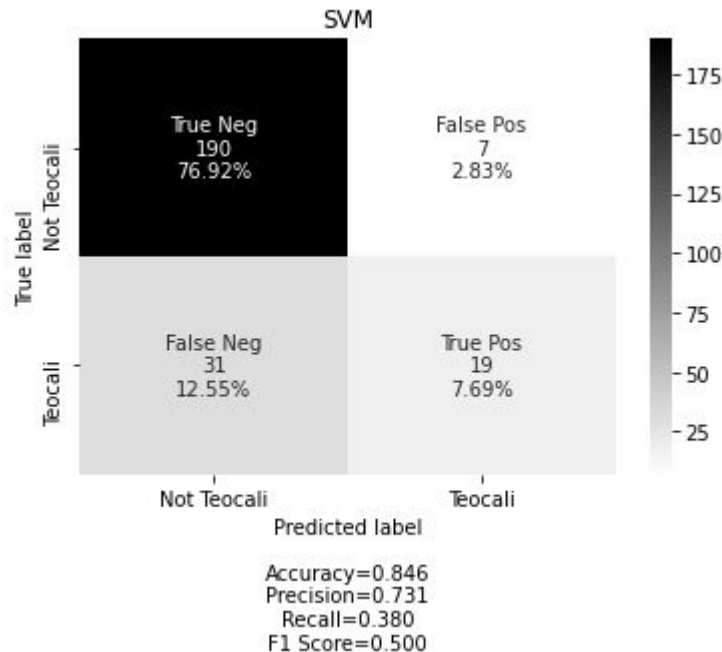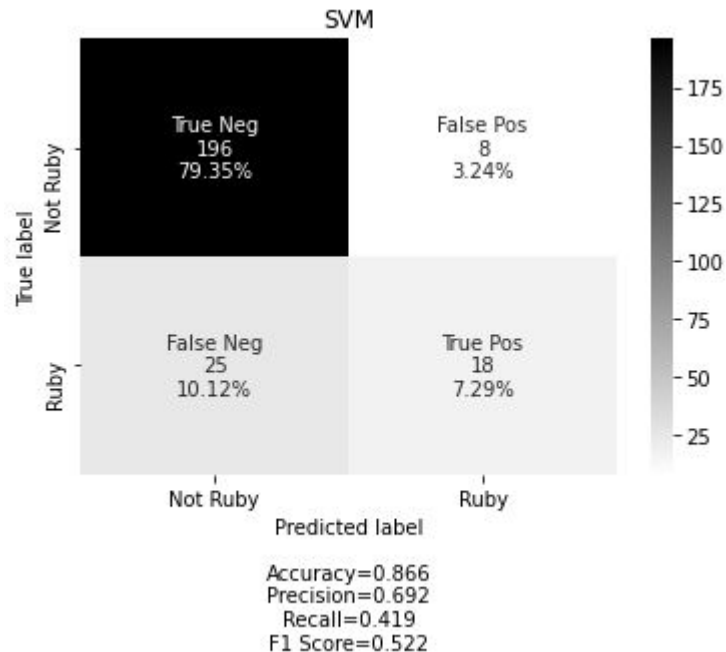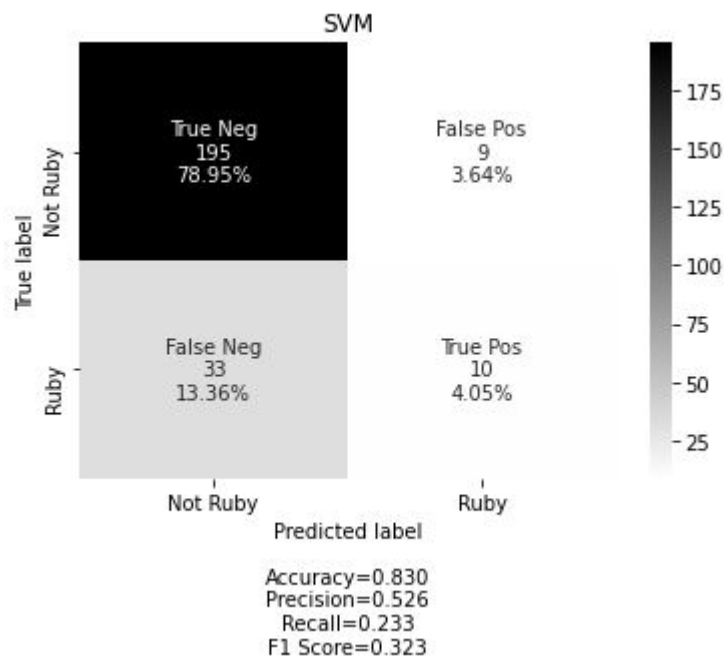
# Hunterhill

# Teocali

# Ruby



SVM

| | Predicted label | |
|---|---|---|
| | Not Ruby | Ruby |
| Not Ruby | True Neg 195 78.95% | False Pos 9 3.64% |
| Ruby | False Neg 33 13.36% | True Pos 10 4.05% |

Accuracy=0.830
Precision=0.526
Recall=0.233
F1 Score=0.323

SVM

| | Predicted label | |
|---|---|---|
| | Not Ruby | Ruby |
| Not Ruby | True Neg 196 79.35% | False Pos 8 3.24% |
| Ruby | False Neg 25 10.12% | True Pos 18 7.29% |

Accuracy=0.866
Precision=0.692
Recall=0.419
F1 Score=0.522
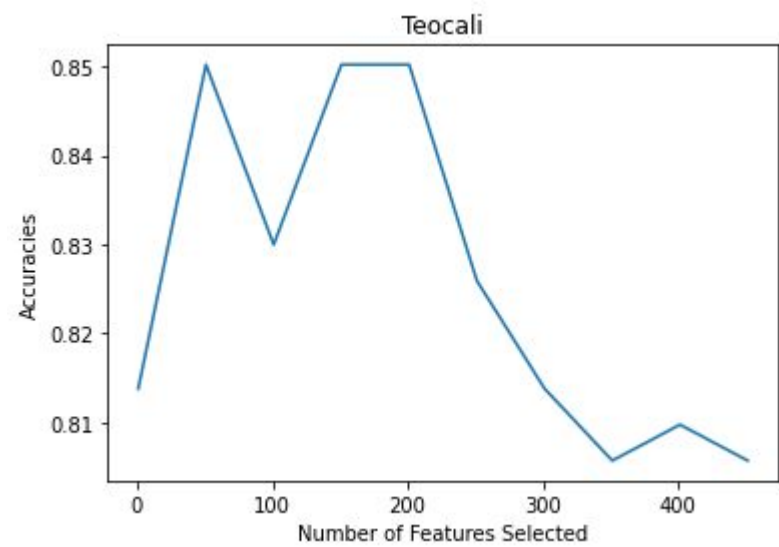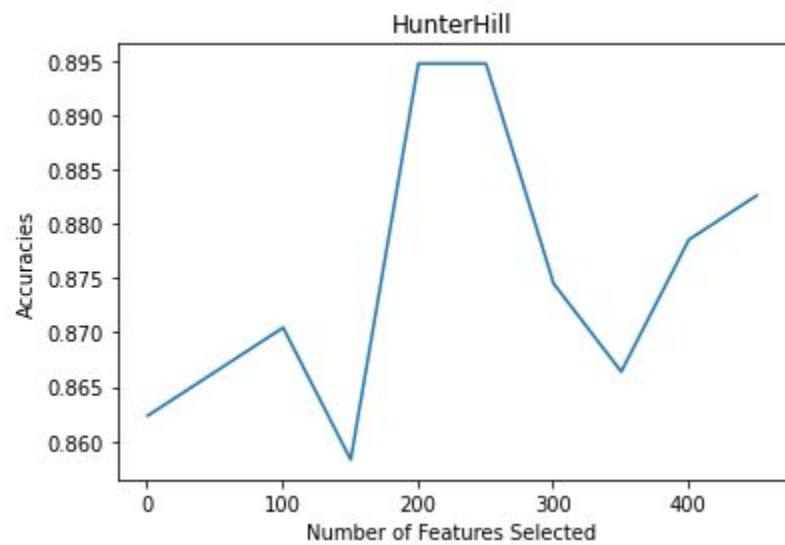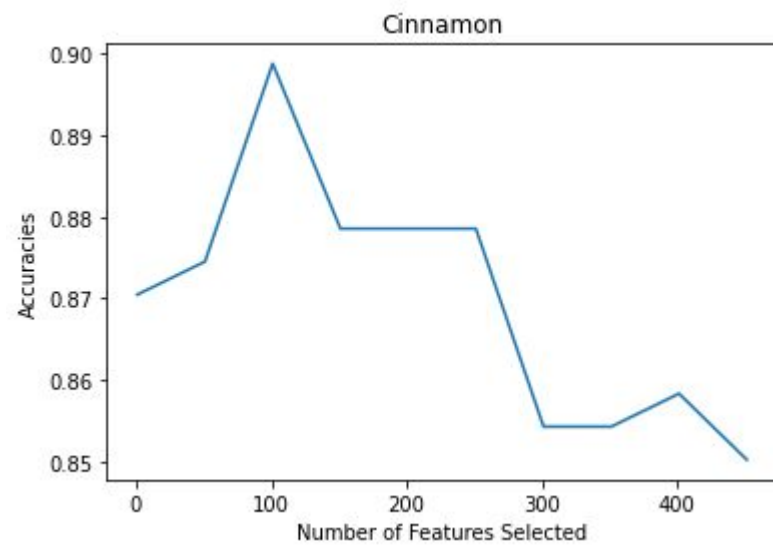
# Number of Top Features Selected (increment of 50)

Teocali

HunterHill

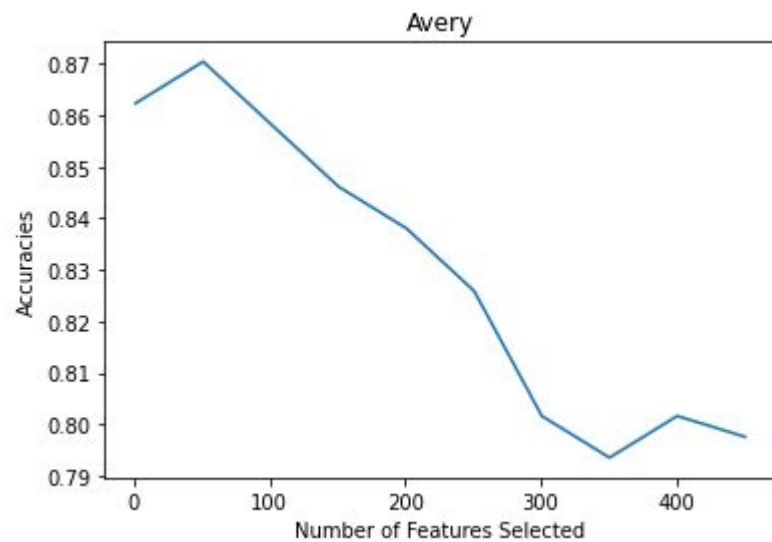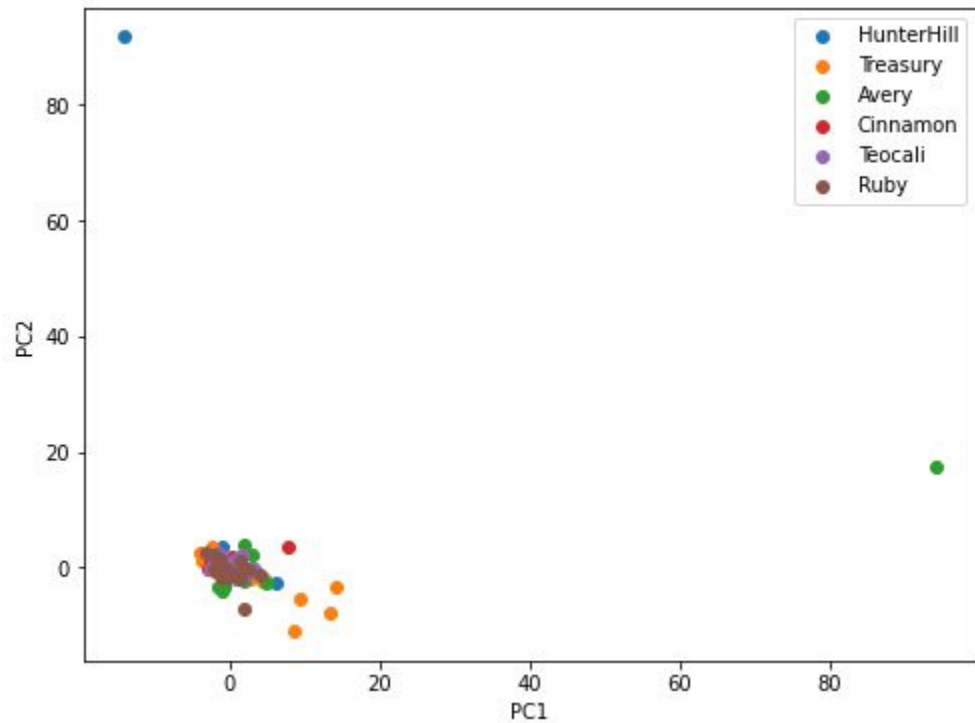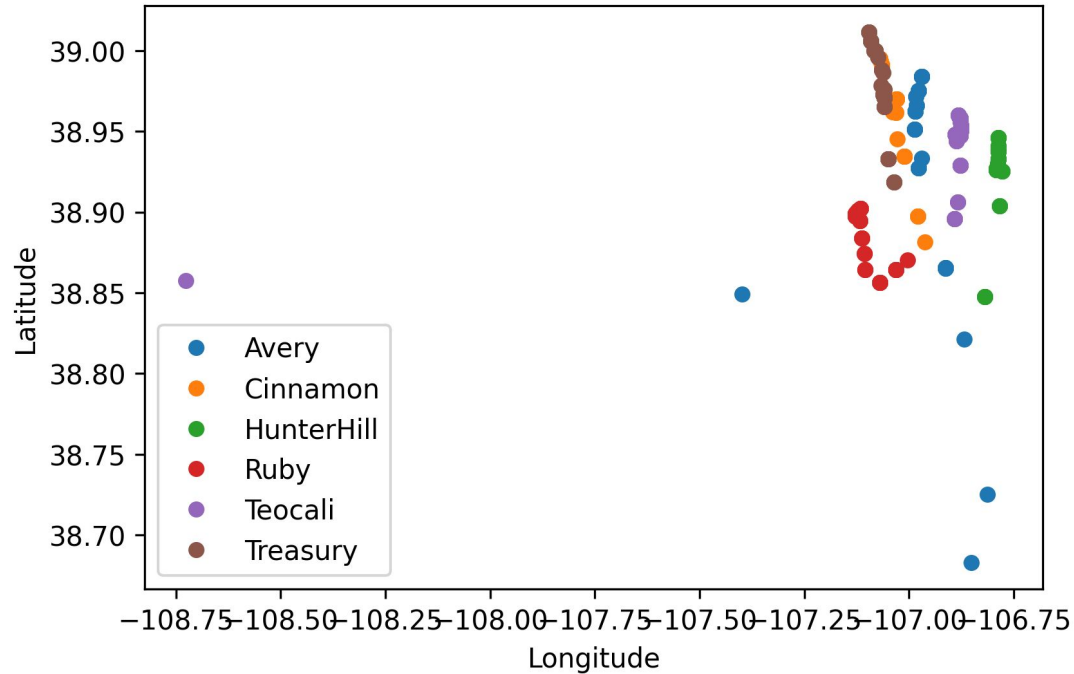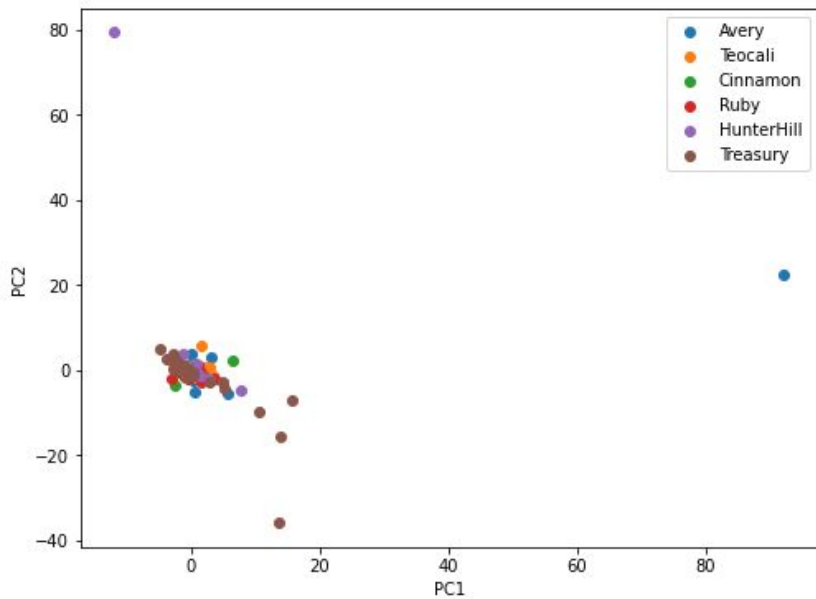Cinnamon

Avery

PCA Plot (Not much Variance Explained, as expected,about 1.5 percent variance explained on PC1 and PC2)

# Sample Map (one odd Teocali Sample) (old version with ELEL)
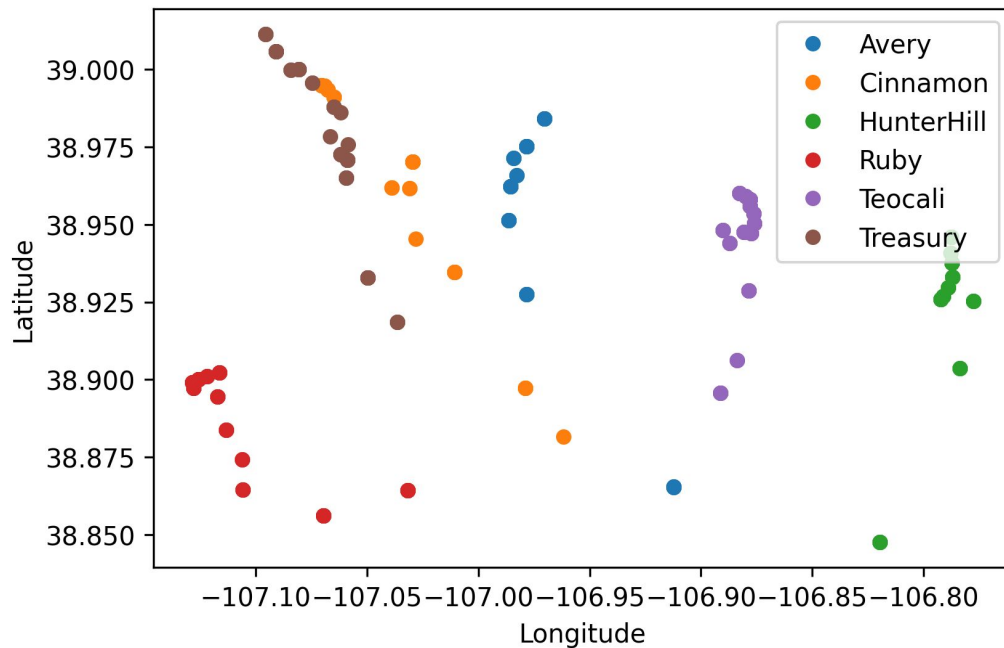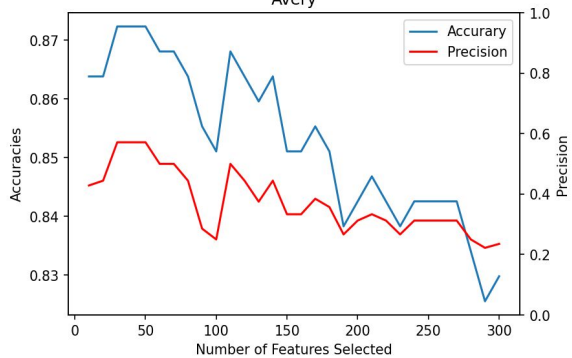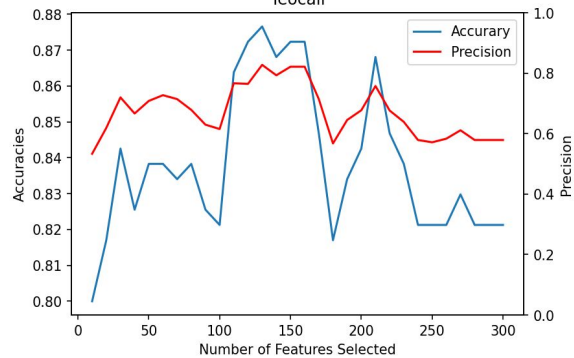
# Updated PCA (slightly better, [0.01696325 0.01507139])

# Updated Map (removed ELEL from data)

# RUBY

# CINNAMON



## RUBY — SVM
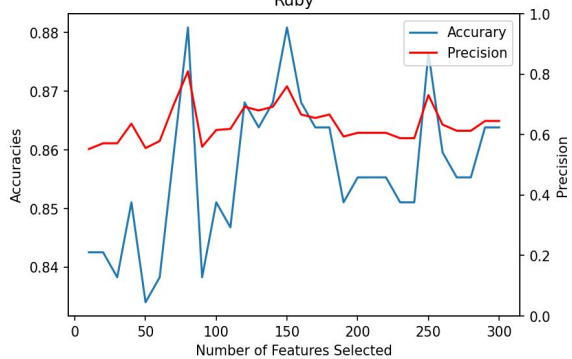
|  | Not Ruby | Ruby |
|---|---|---|
| **Not Ruby** | True Neg 32 39.02% | False Pos 9 10.98% |
| **Ruby** | False Neg 5 6.10% | True Pos 36 43.90% |

True label / Predicted label

Accuracy=0.829
Precision=0.800
Recall=0.878
F1 Score=0.837

## CINNAMON — SVM

|  | Not Cinnamon | Cinnamon |
|---|---|---|
| **Not Cinnamon** | True Neg 25 36.76% | False Pos 9 13.24% |
| **Cinnamon** | False Neg 6 8.82% | True Pos 28 41.18% |

True label / Predicted label

Accuracy=0.779
Precision=0.757
Recall=0.824
F1 Score=0.789

Avery

SVM

True label: Not Avery / Avery
Predicted label: Not Avery / Avery

True Neg 21 33.87%
False Pos 10 16.13%
False Neg 6 9.68%
True Pos 25 40.32%

Accuracy=0.742
Precision=0.714
Recall=0.806
F1 Score=0.758

Treasury

SVM

True label: Not Treasury / Treasury
Predicted label: Not Treasury / Treasury

True Neg 30 34.09%
False Pos 14 15.91%
False Neg 16 18.18%
True Pos 28 31.82%

Accuracy=0.659
Precision=0.667
Recall=0.636
F1 Score=0.651

Hunterhill

Teocali

SVM

|  | Not HunterHill | HunterHill |
|---|---|---|
| Not HunterHill | True Neg 26 35.14% | False Pos 11 14.86% |
| HunterHill | False Neg 8 10.81% | True Pos 29 39.19% |

Accuracy=0.743
Precision=0.725
Recall=0.784
F1 Score=0.753

SVM

|  | Not Teocali | Teocali |
|---|---|---|
| Not Teocali | True Neg 31 32.29% | False Pos 17 17.71% |
| Teocali | False Neg 14 14.58% | True Pos 34 35.42% |

Accuracy=0.677
Precision=0.667
Recall=0.708
F1 Score=0.687

# After including 20 percent noise



Prediction Accuracies in Each Fold (Support Vector Machine)