

# *Forest cover type predication*

A watercolor illustration on the left side of the slide. It features several blue solar panels with white grid lines, arranged in a row and slightly overlapping. Above the panels are large, billowing clouds in shades of blue and white, with some darker blue and grey tones. The background is a light, textured white.

*Hello!*

**Group 16**

**Dina Ibrahim Mohammady**

**300389383**

**Nada Mohammed Zakaria**

**300389901**

**Sema Abdelnasser Helali**

**300389914**

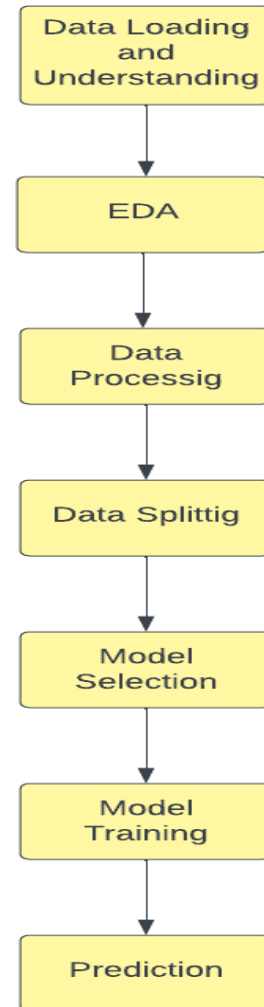


# *Problem Overview*





The main objective is to predict the dominant forest cover type in 30 x 30 meter cells of the Roosevelt National Forest using cartographic variables, without employing remotely sensed data. The forest cover types in the four wilderness areas are primarily shaped by ecological processes due to minimal human-caused disturbances, emphasizing the significance of understanding the ecological relationships between cartographic features and forest cover types in these undisturbed wilderness areas.



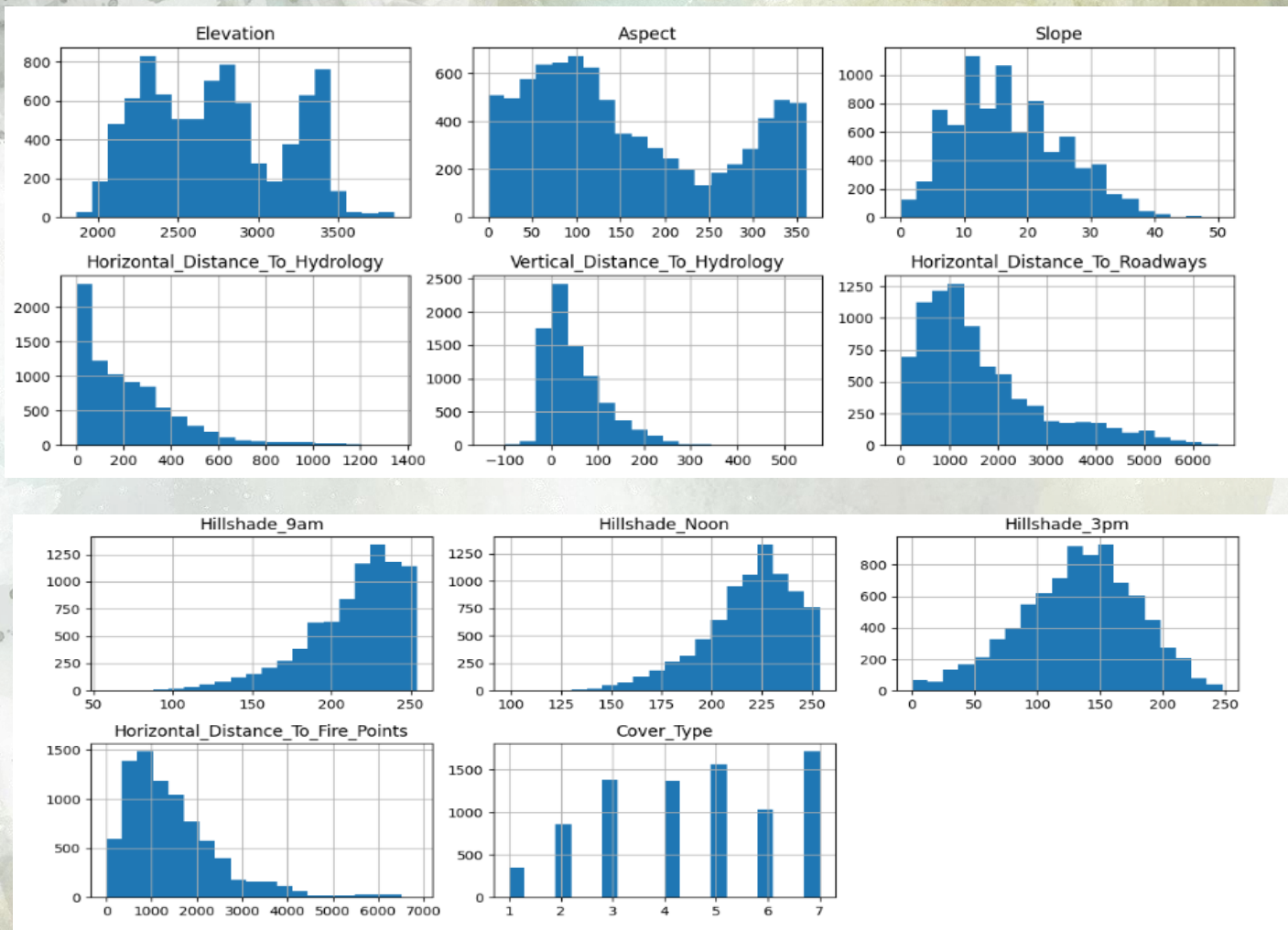


# Exploratory Data Analysis

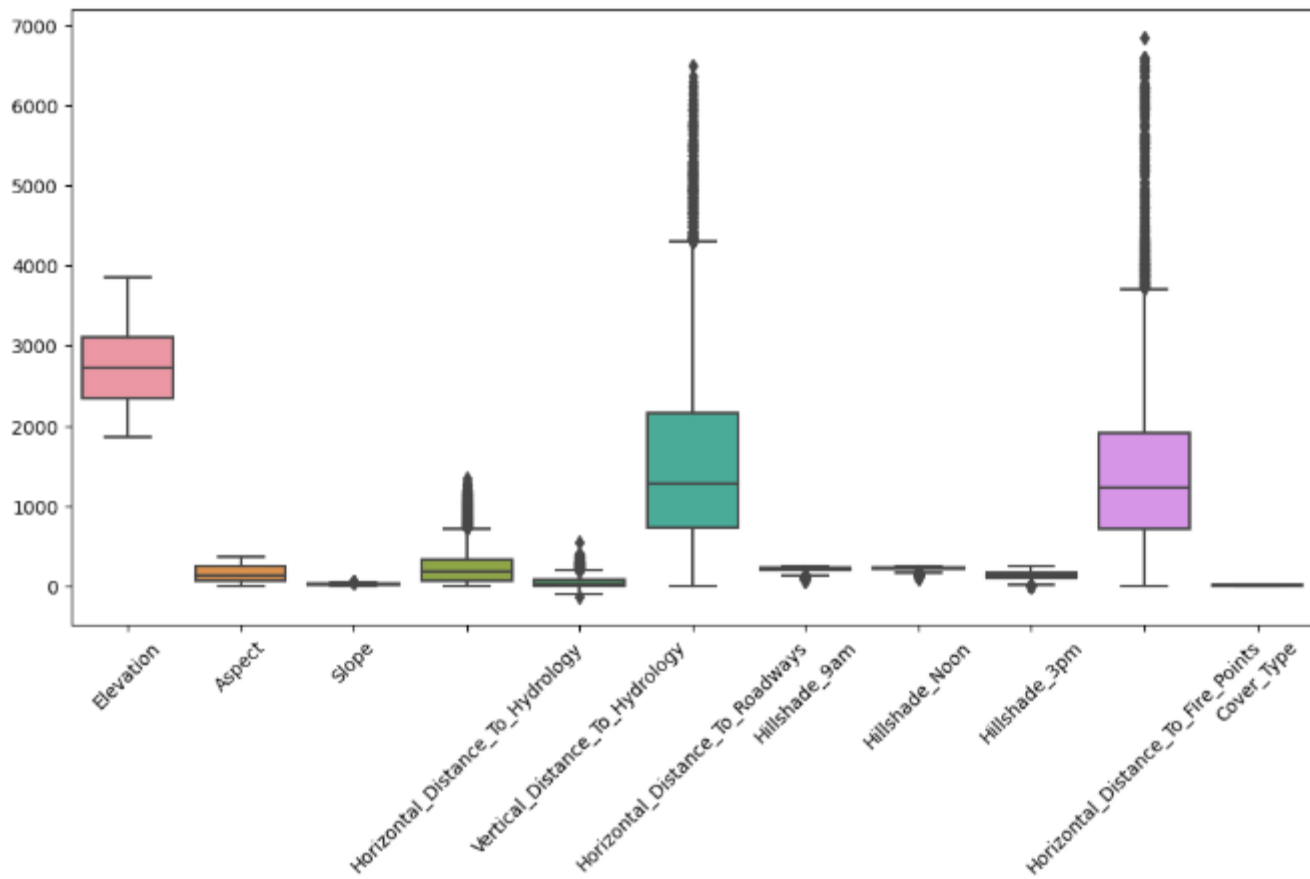
## Dataset statistics

Number of variables	56
Number of observations	8286
Missing cells	0
Missing cells (%)	0.0%
Duplicate rows	0
Duplicate rows (%)	0.0%
Total size in memory	3.5 MiB
Average record size in memory	448.0 B

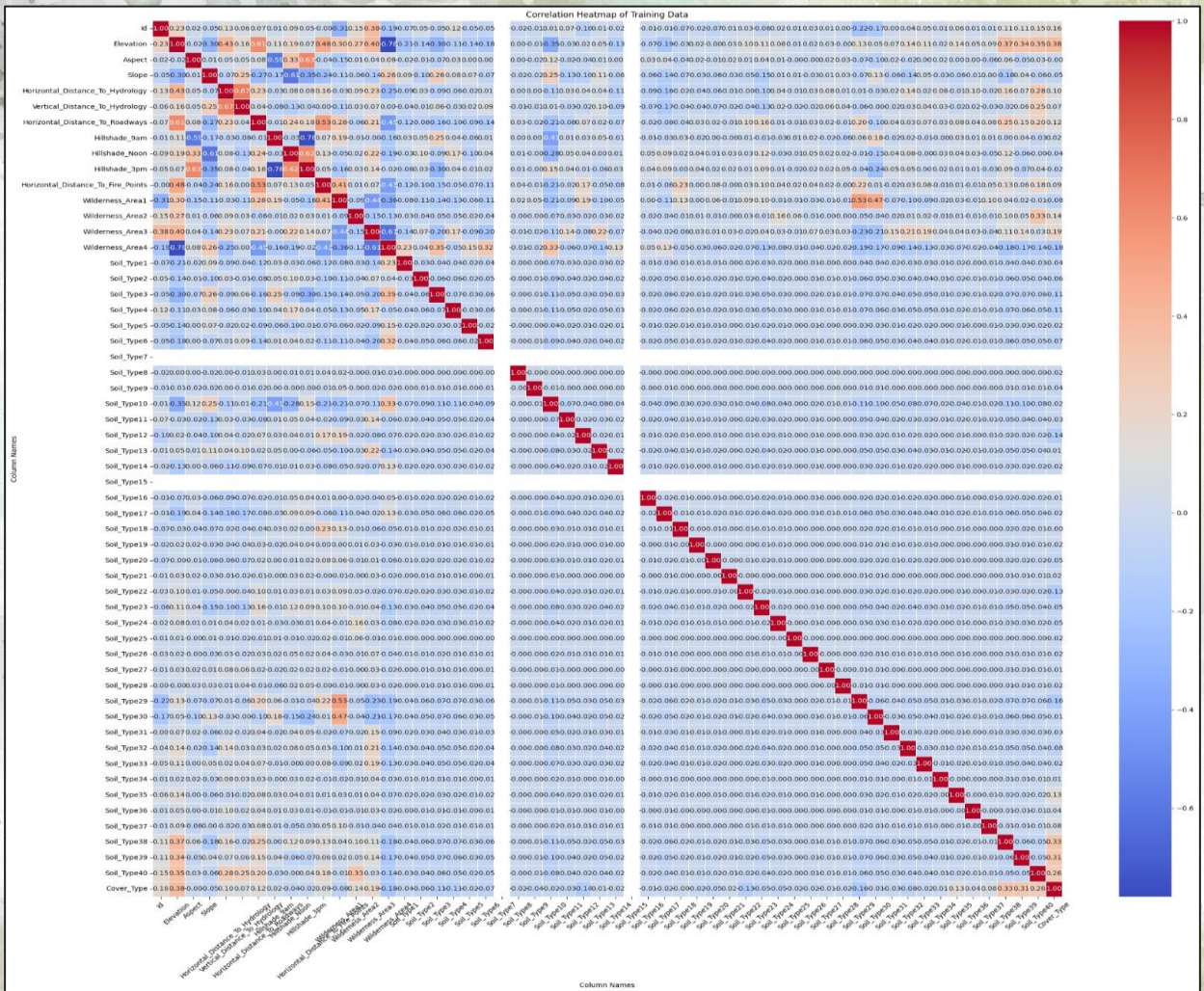


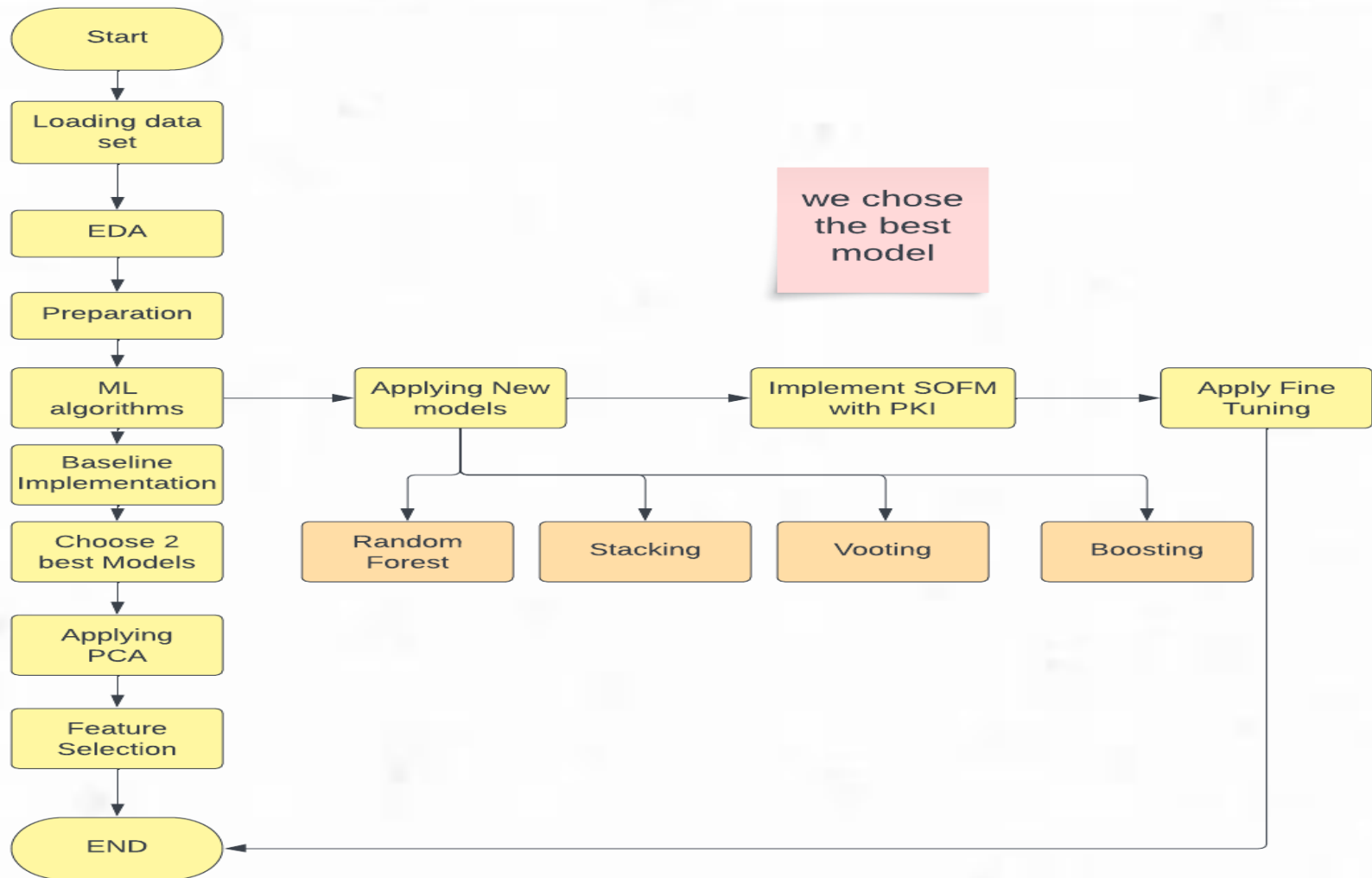


Box Plots of Numerical Features



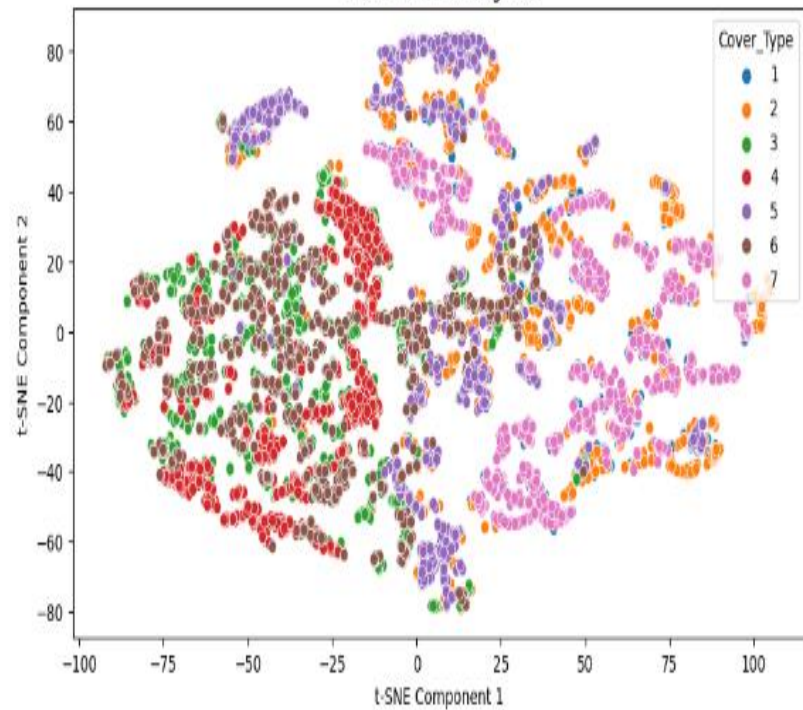




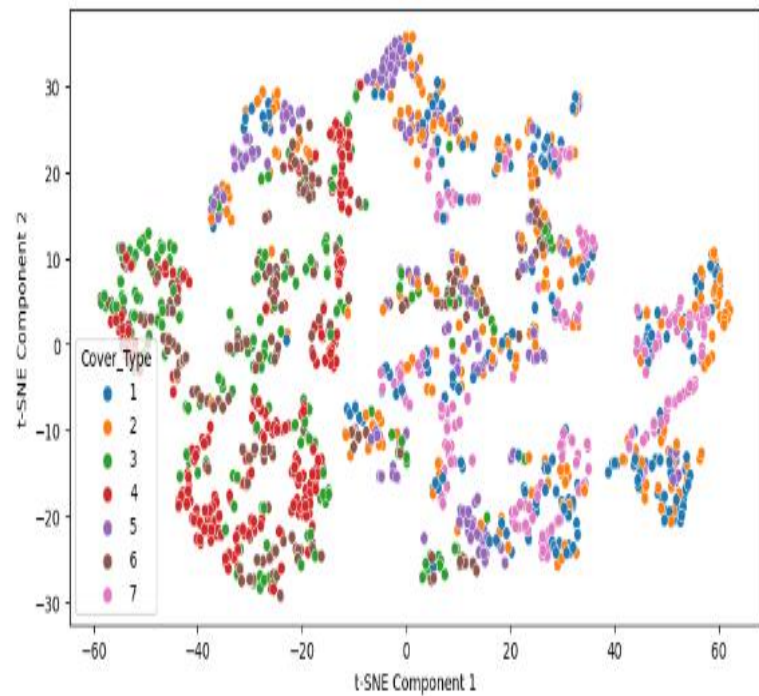




t-SNE Plot - Training Set



t-SNE Plot - Test Set

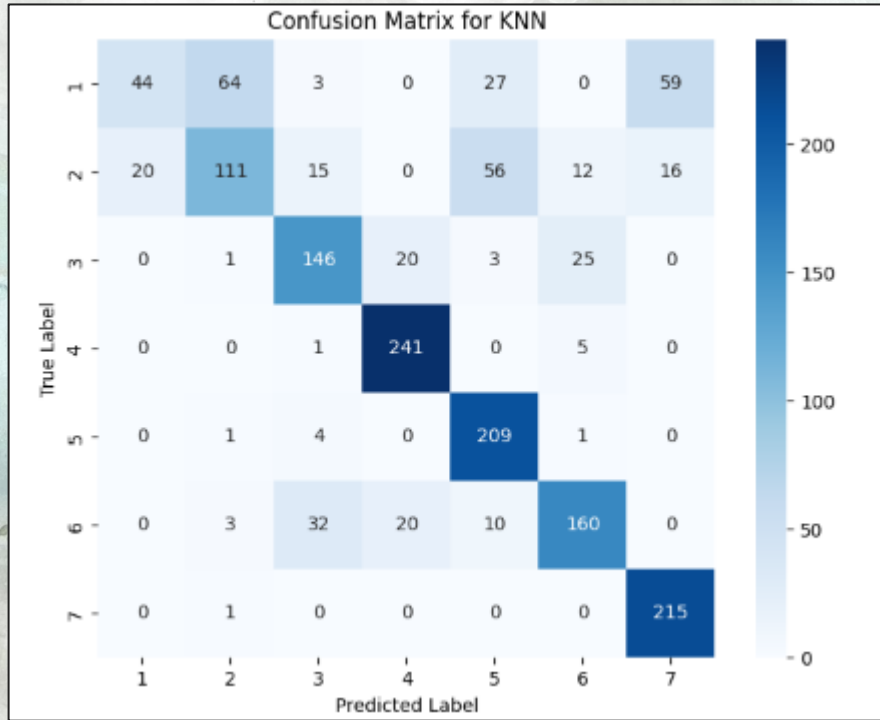




# Baseline Performance



## For KNN



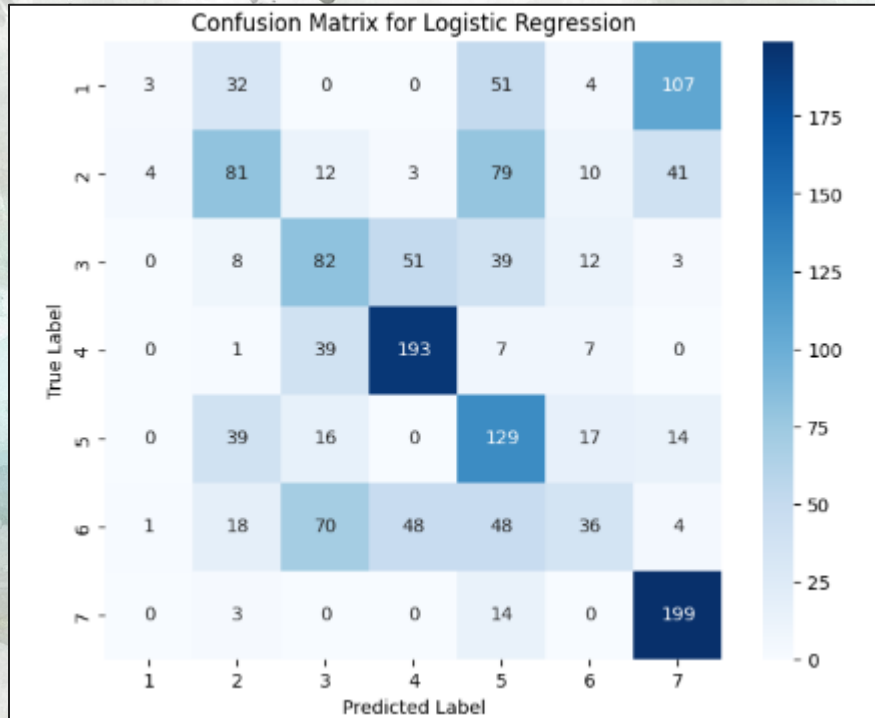
Training Accuracy for KNN: 0.8853

Testing Accuracy for KNN: 0.7384

Validation Accuracy for KNN: 0.7523



## For Logistic Regression



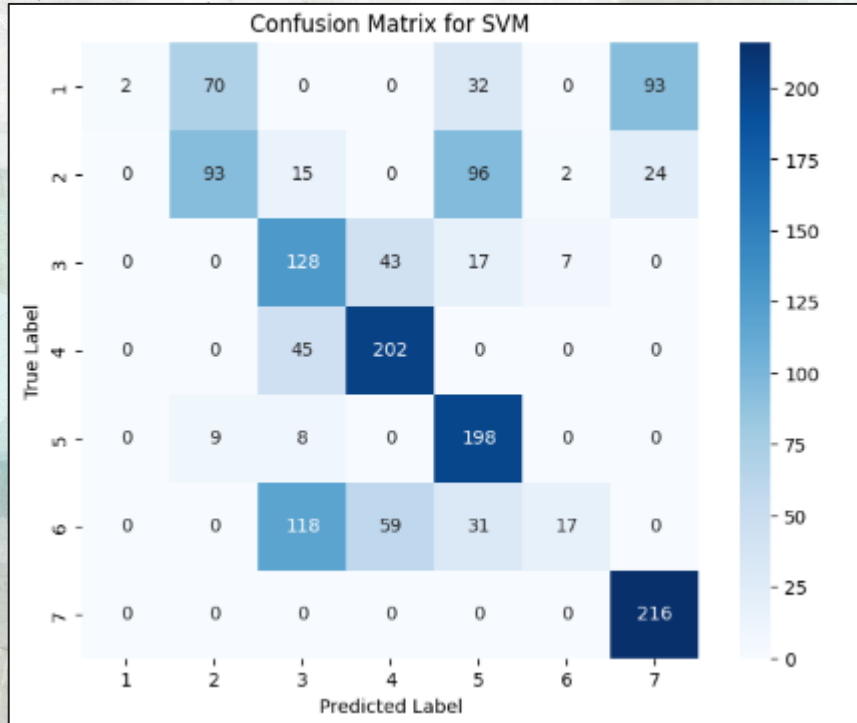
Training Accuracy for Logistic Regression: 0.5782

Testing Accuracy for Logistic Regression: 0.4741

Validation Accuracy for Logistic Regression: 0.4773



## For SVM

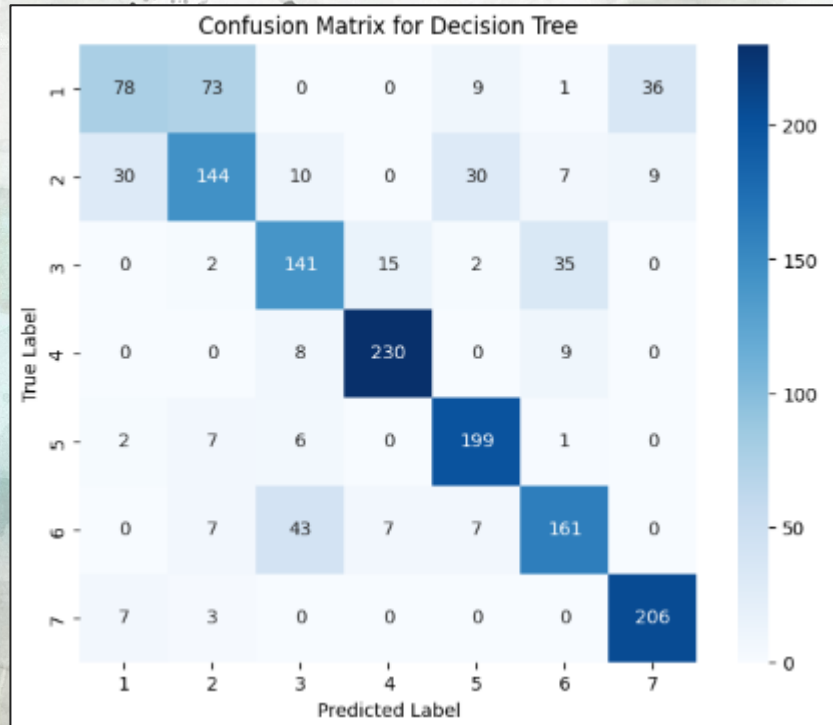


Training Accuracy for SVM: 0.6859

Testing Accuracy for SVM: 0.5613

Validation Accuracy for SVM: 0.5501

## For Decision Tree



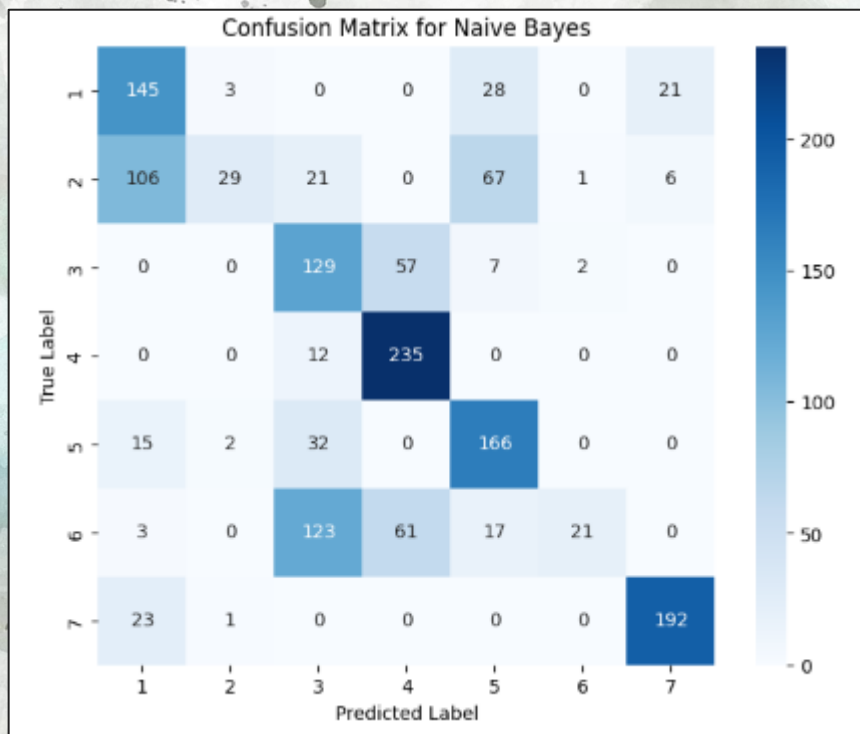
Training Accuracy for Decision Tree: 1.0000

Testing Accuracy for Decision Tree: 0.7600

Validation Accuracy for Decision Tree: 0.7530



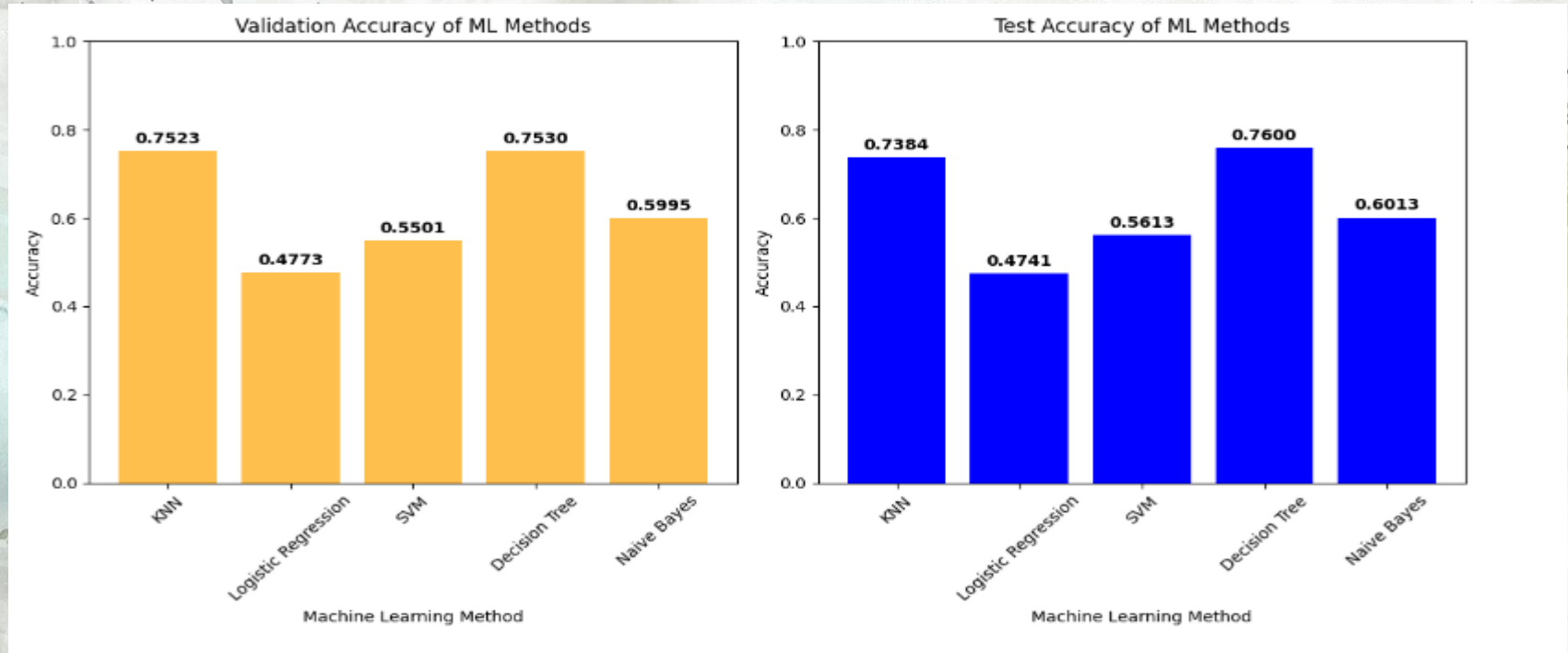
## For Naïve Bayes



Training Accuracy for Naive Bayes: 0.6417  
Testing Accuracy for Naive Bayes: 0.6013  
Validation Accuracy for Naive Bayes: 0.5995



# Accuracy Comparison

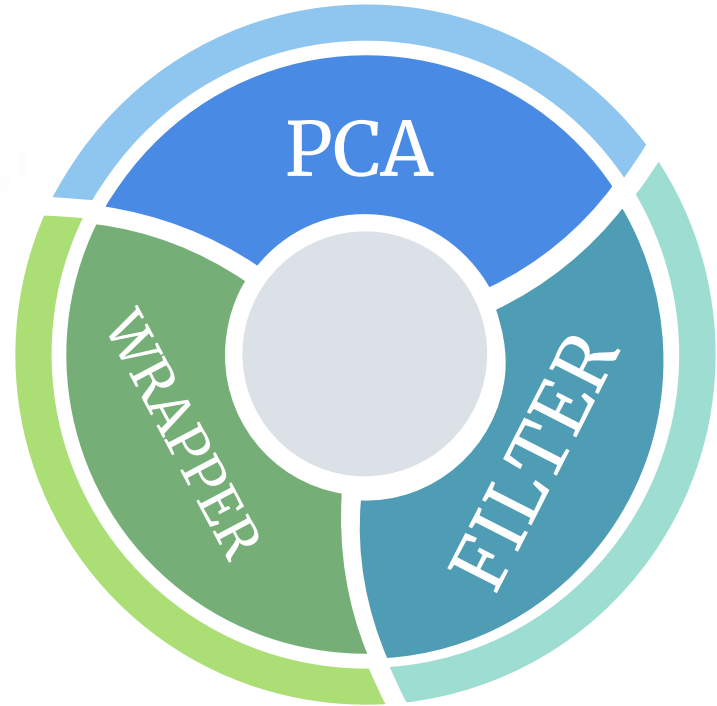


The Best 2 Baseline are Decision Tree and KNN

# *Improvement Strategy*

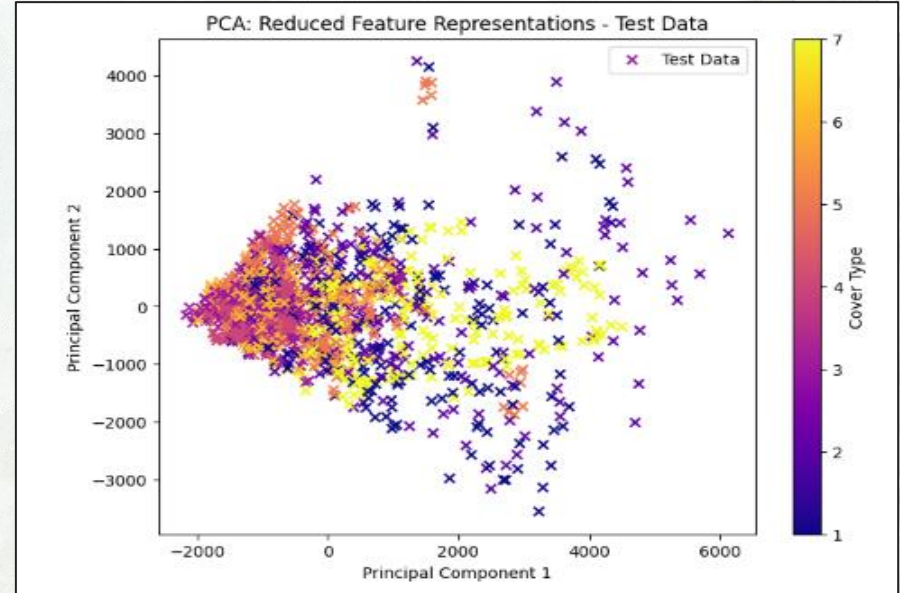
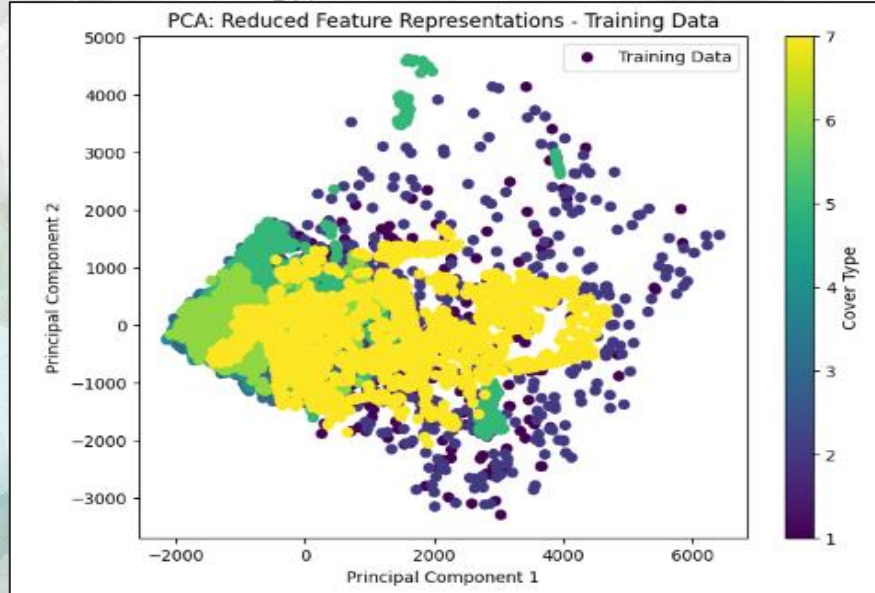


# Improvement Strategy

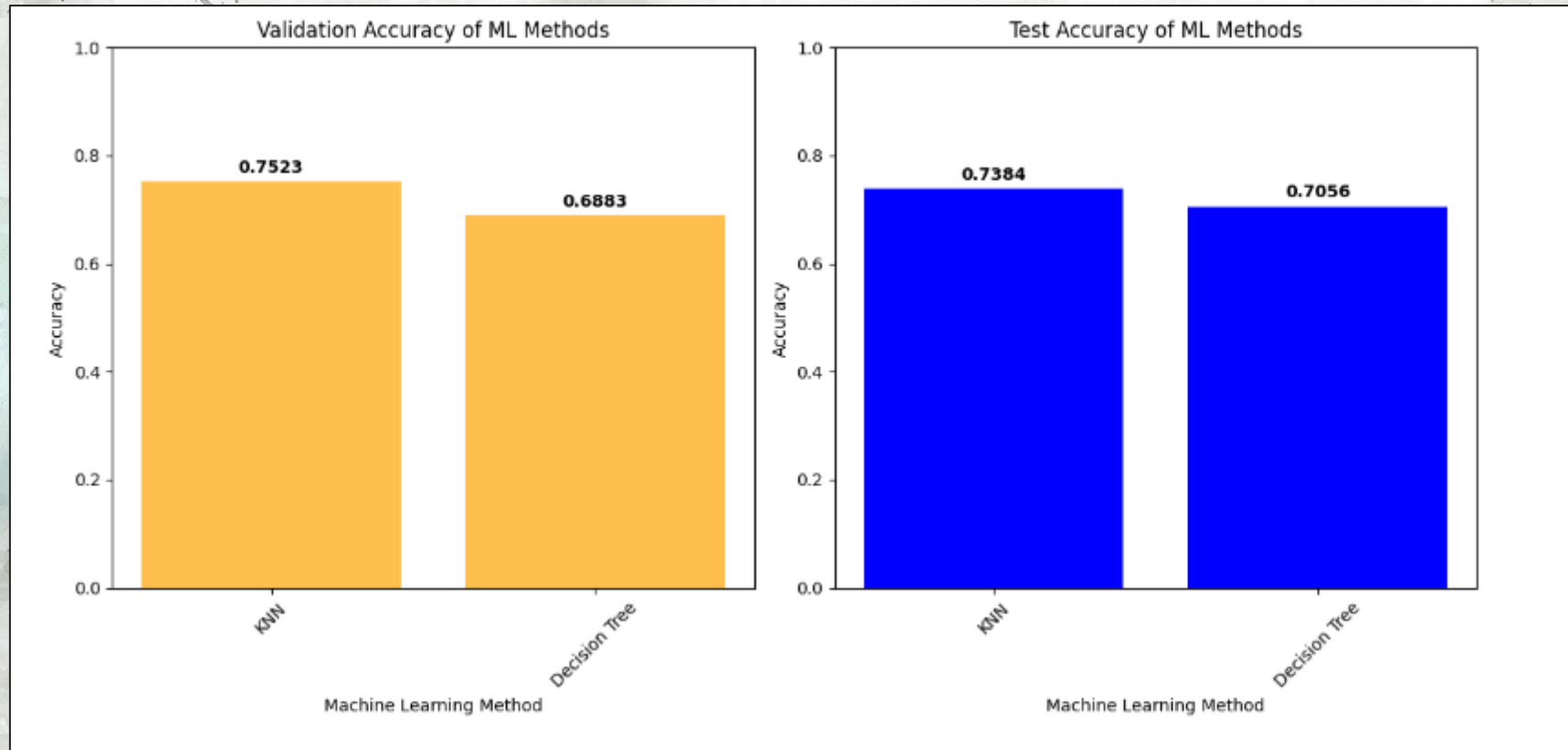




# Applying PCA

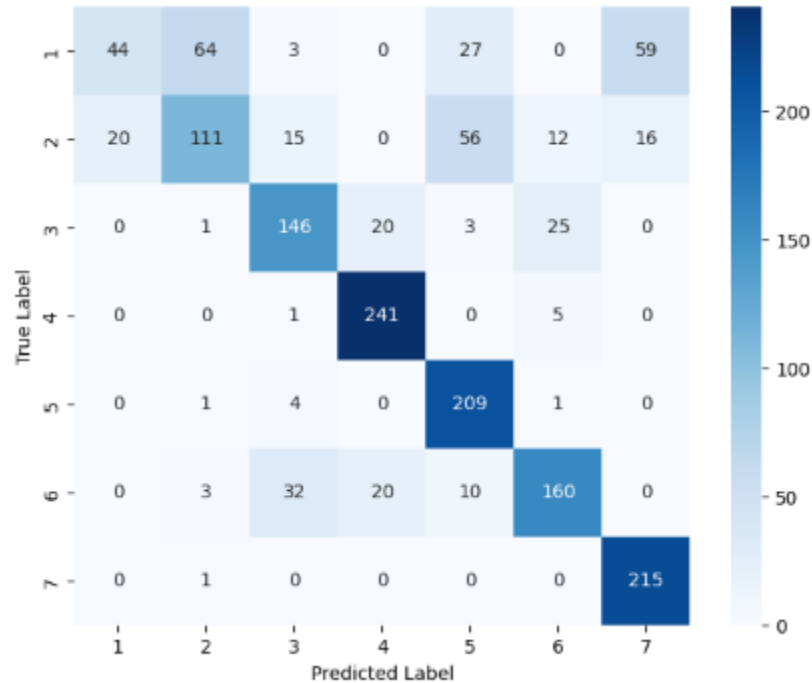


# Comparing 2 models after PCA



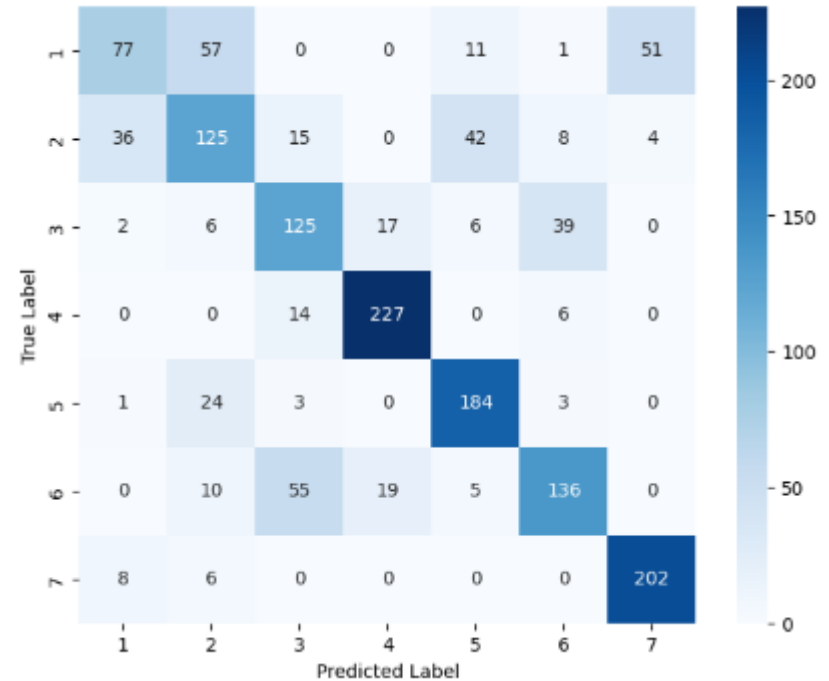
# Comparing 2 models after PCA

Confusion Matrix for KNN



Training Accuracy for KNN: 0.8853  
Testing Accuracy for KNN: 0.7384  
Validation Accuracy for KNN: 0.7523

Confusion Matrix for Decision Tree

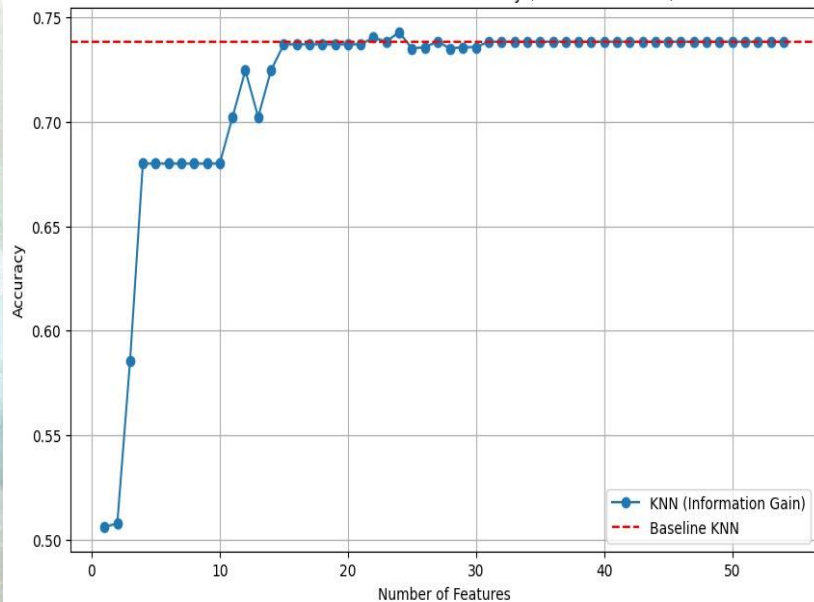


Training Accuracy for Decision Tree: 1.0000  
Testing Accuracy for Decision Tree: 0.7056  
Validation Accuracy for Decision Tree: 0.6883

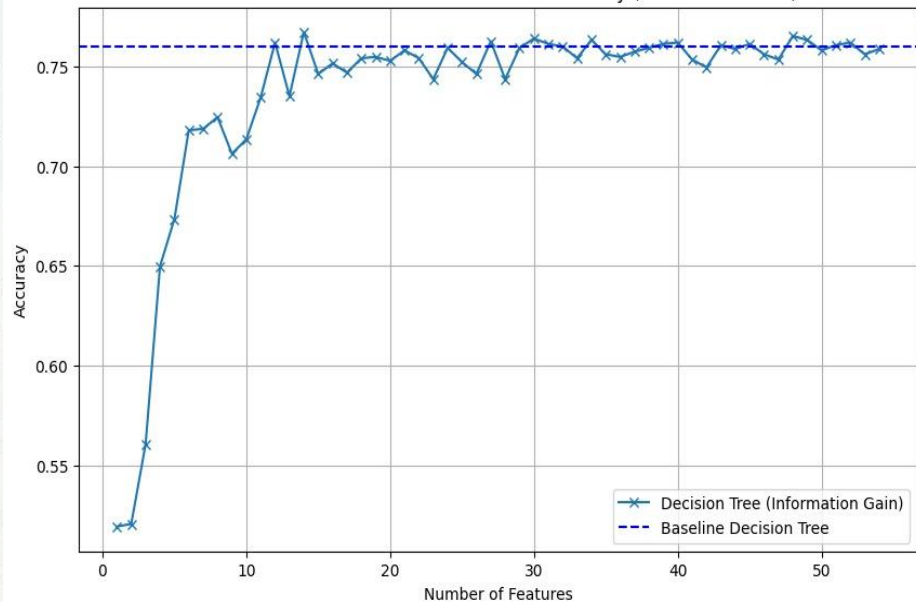


# Applying Filter

KNN: Number of Features vs. Accuracy (Information Gain)

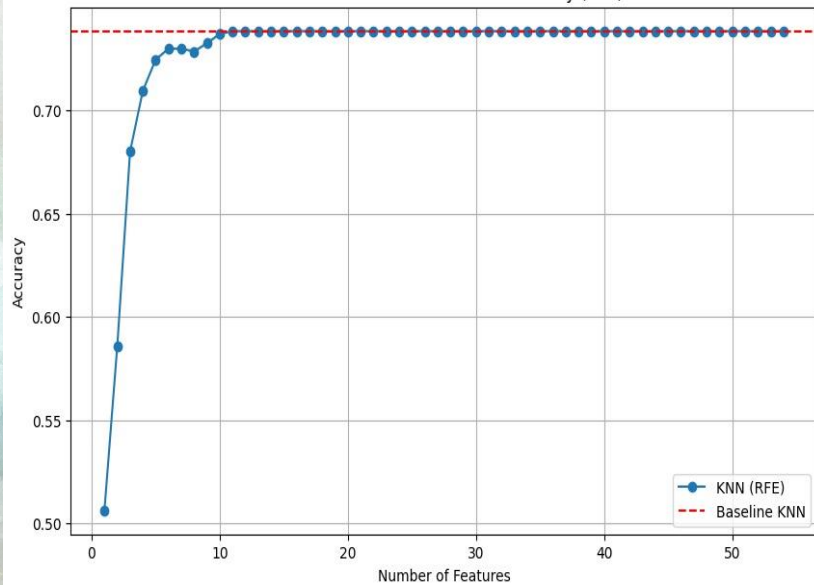


Decision Tree: Number of Features vs. Accuracy (Information Gain)

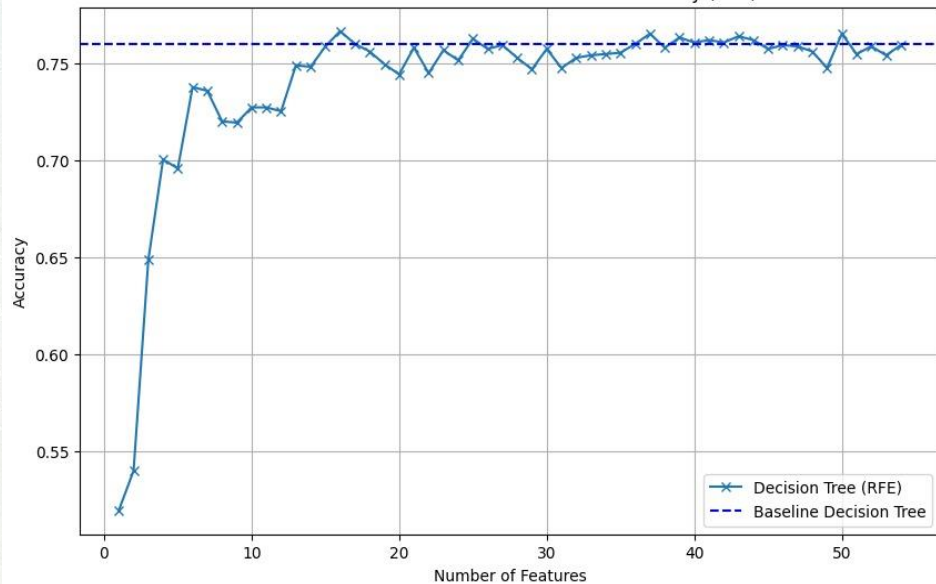


# Applying Wrapper

KNN: Number of Features vs. Accuracy (RFE)



Decision Tree: Number of Features vs. Accuracy (RFE)







# Applying New models



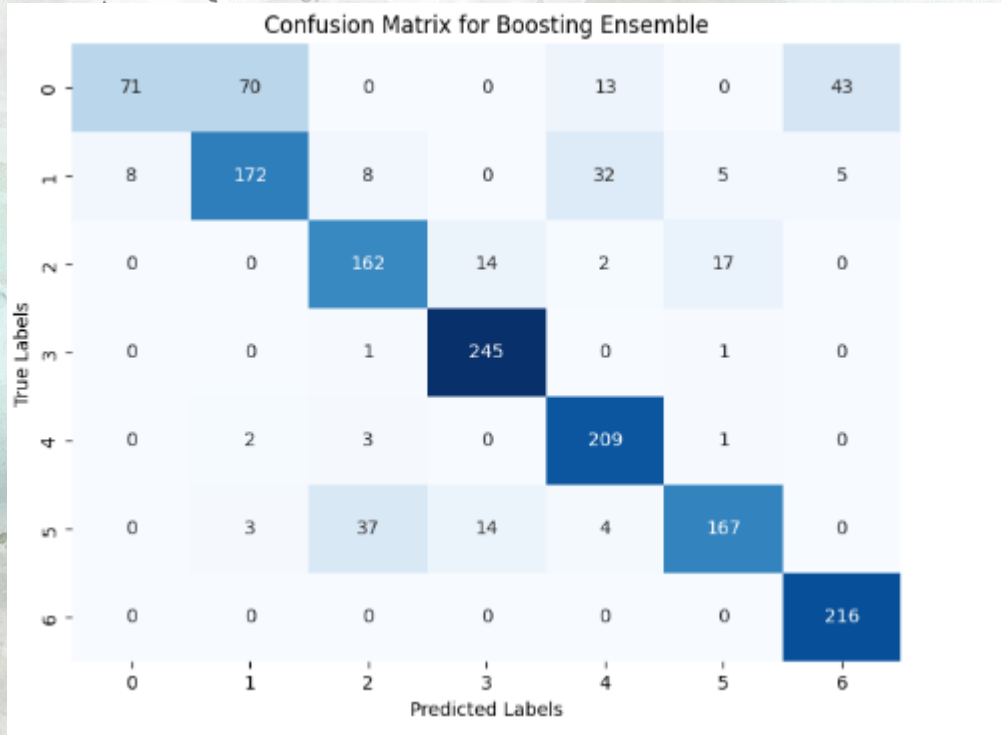
# VOTING Model

Confusion Matrix for Voting Ensemble

True Labels \ Predicted Labels	0	1	2	3	4	5	6
0	76	65	0	0	13	0	43
1	13	166	9	0	31	5	6
2	0	0	166	16	3	10	0
3	0	0	1	245	0	1	0
4	0	3	3	0	208	1	0
5	0	3	41	15	3	163	0
6	0	0	0	0	0	0	216

Accuracy of Voting Ensemble: 0.8131147540983606

# Boosting Model



Accuracy of Boosting: 0.8144262295081968

# Stacking Model

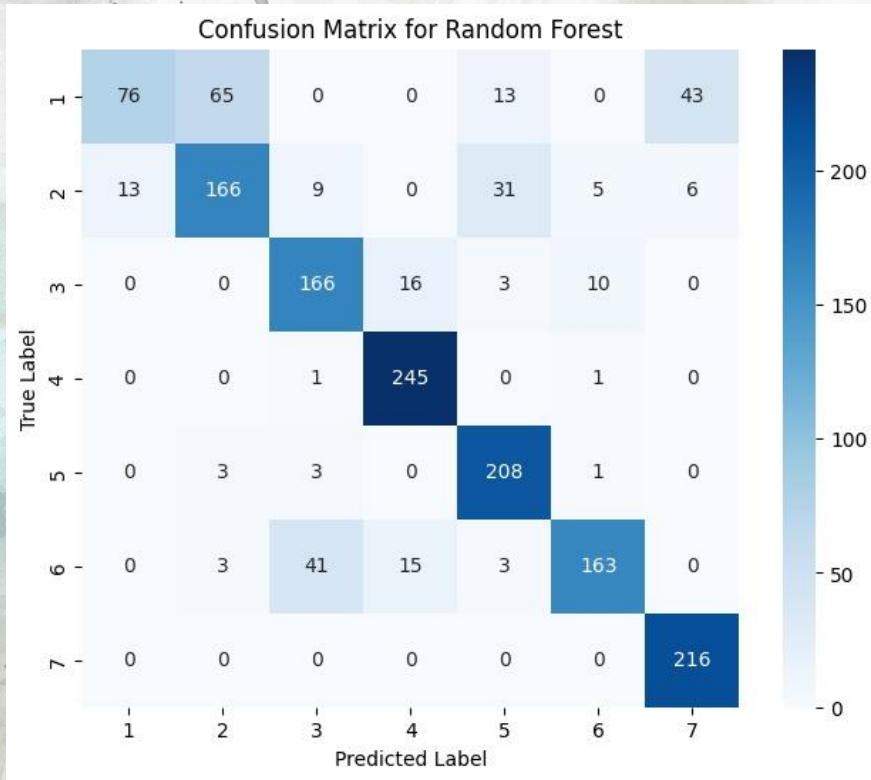
Confusion Matrix for Stacking Ensemble

True Labels	0	1	2	3	4	5	6
	95	65	0	0	9	0	28
	21	167	10	0	21	6	5
	0	0	162	12	2	19	0
	0	0	4	239	0	4	0
	0	4	3	0	208	0	0
	0	3	38	13	3	168	0
	0	1	2	3	4	5	6
		Predicted Labels					

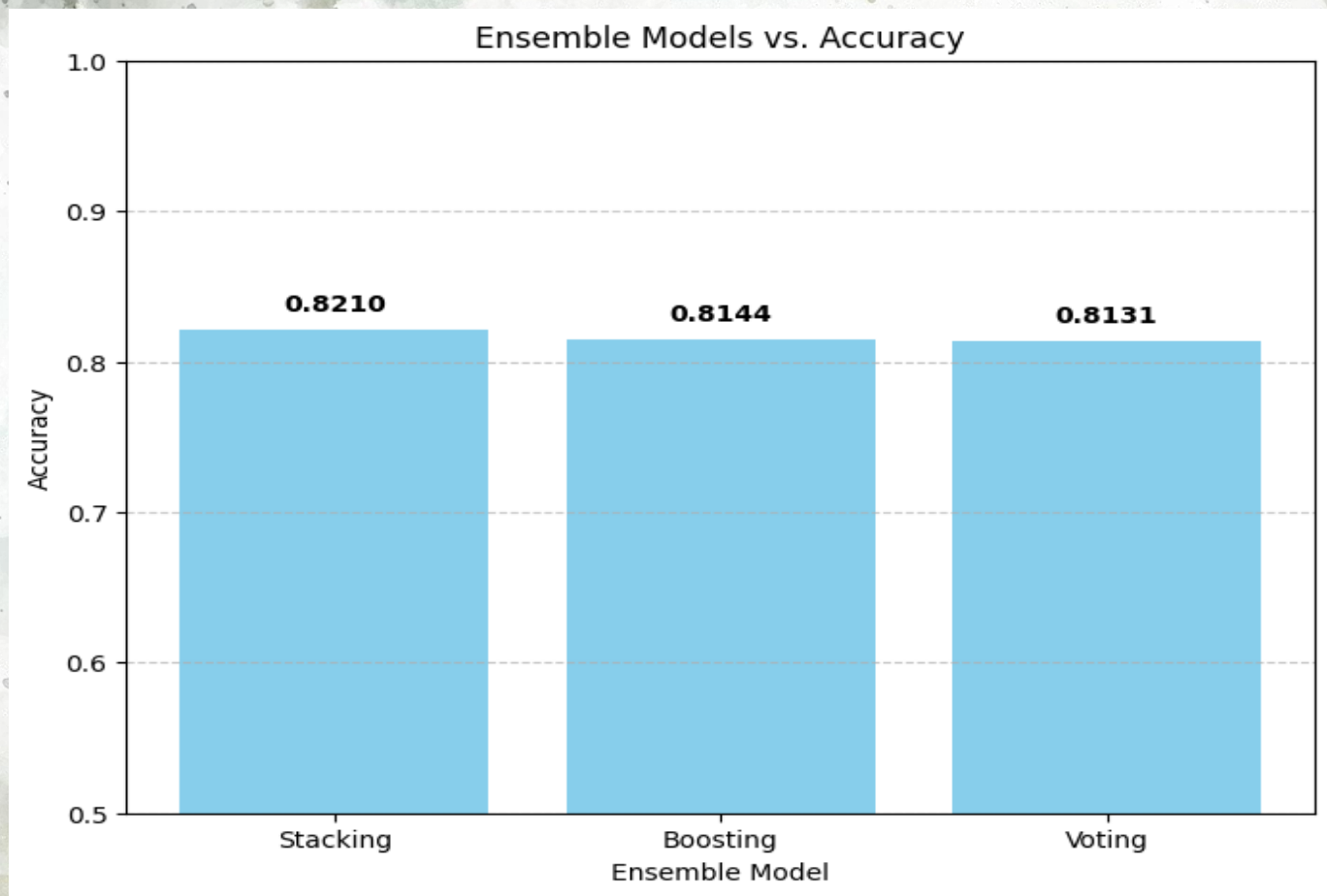
Accuracy of Stacking: 0.820983606557377



# Random Forest Model



Testing Accuracy for Random Forest: 0.8131



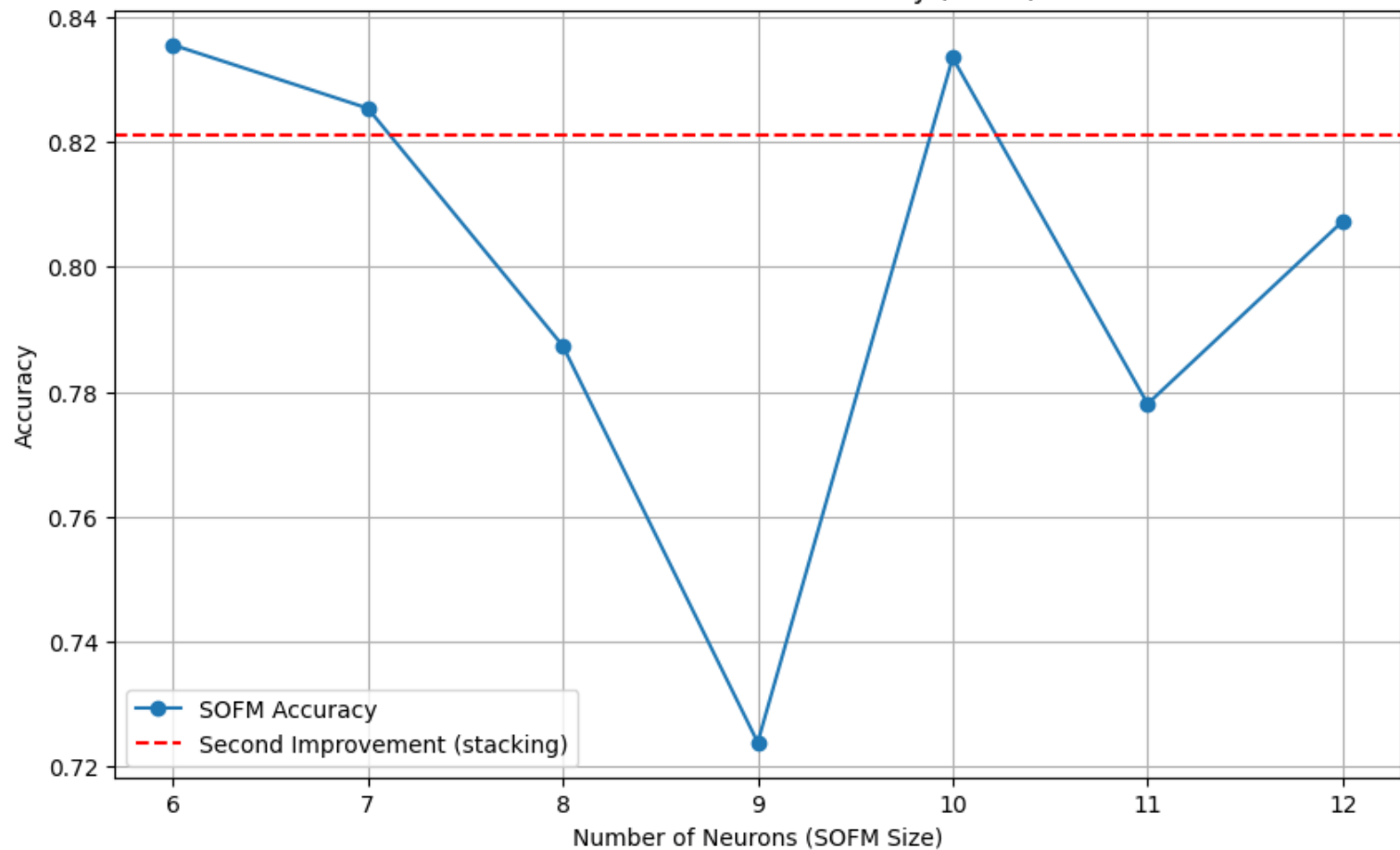


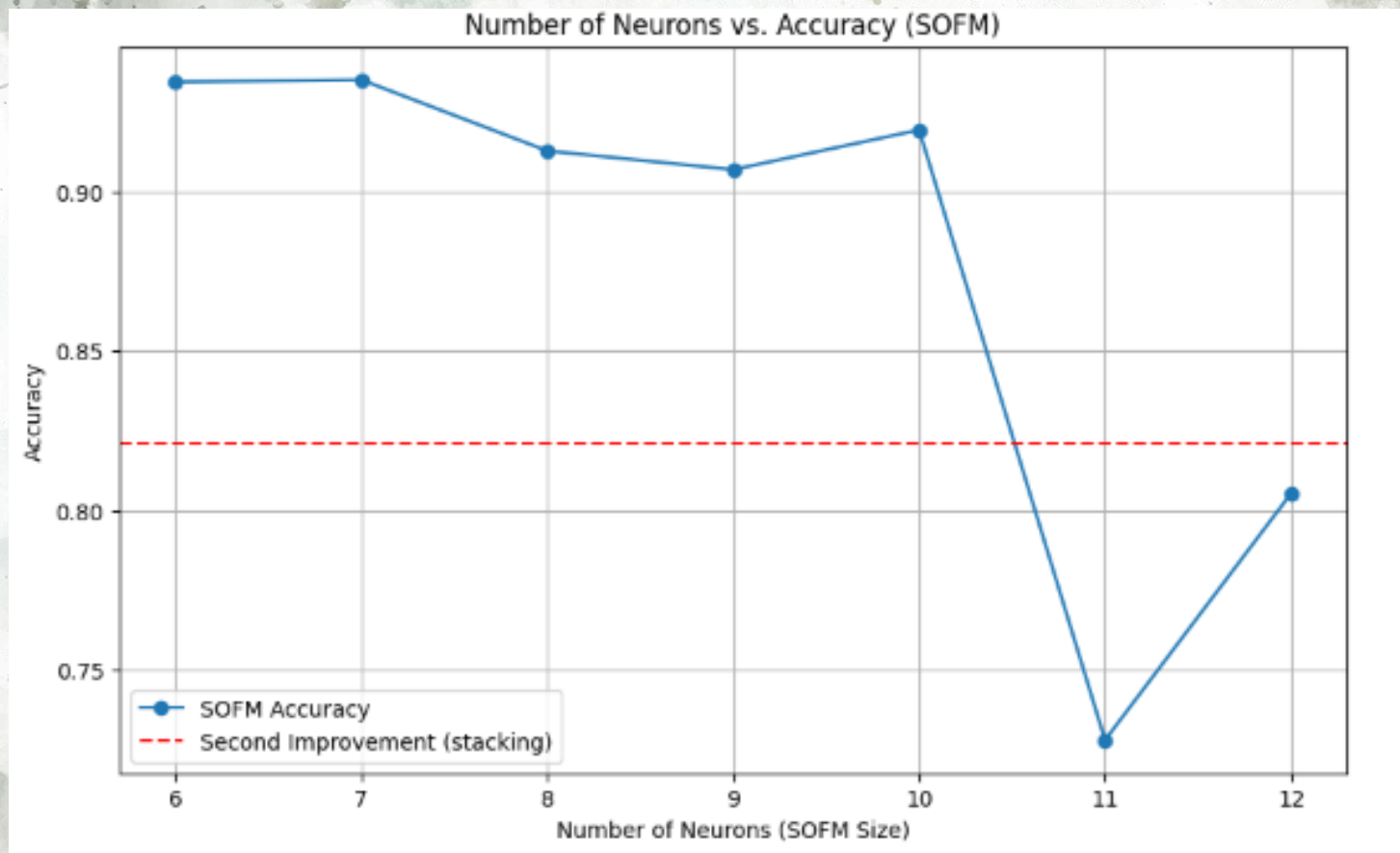
*Supervised &  
Unsupervised  
Combination via PKI*





Number of Neurons vs. Accuracy (SOFM)

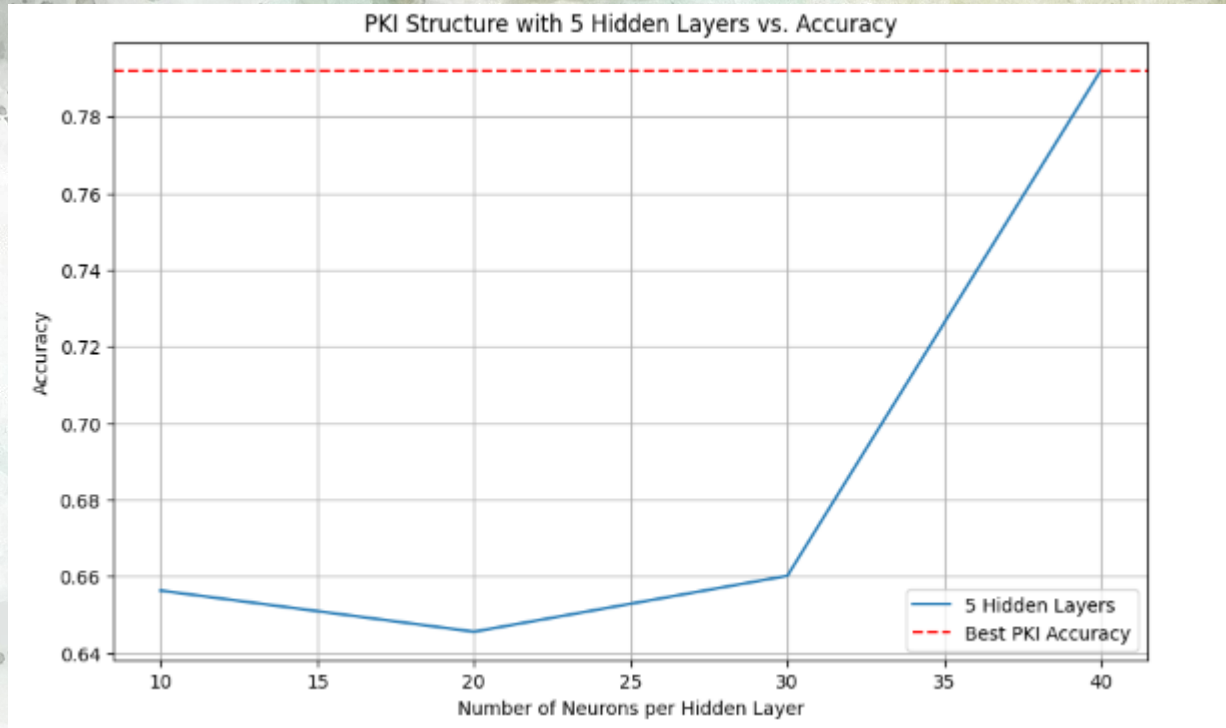




# fine tuning







Best SQFM Structure: 10x10, Accuracy: 0.8425 Best PKI Structure: 5 Hidden Layers, 40 Neurons per Layer, Accuracy: 0.7921

## conclusion

1. **First Improvement** : Decision Tree with Filter-based Feature Selection (SelectKBest): By selecting the most informative features using SelectKBest and applying Decision Tree, we improved the model's accuracy compared to the baseline.
2. **Second Improvement** :Random Forest: Utilizing Random Forest as an ensemble model further boosted the accuracy beyond the first improvement.
3. **Third Improvement** : PKI Model (SOFM + DNN): - Incorporating Self-Organizing Feature Map (SOFM) clustering information into a Deep Neural Network (DNN) model provided significant accuracy improvement compared to the previous strategies.

By sequentially applying these three improvement strategies, we achieved substantial enhancements in the supervised model's accuracy and performance, making it more capable of making accurate predictions on the target dataset.



*Thanks!*

