

Project Report

Executive Summary

This project aims at the implementation of 3 machine learning algorithms – Support Vector Machines, Decision Trees and Gradient Boosting for 2 different datasets. A variety of control parameters are looked at for experimentation purposes for the mentioned algorithms such as *penalty, gamma, training size, degree* for SVM; *training size, maximum depth and pruning* for Decision Trees, *Number of boosting stages, pruning* for Gradient Boosting etc. The relationship between the variables and the accuracy/error rate is captured using appropriate plots/learning curves. A final summary/comparison is provided at the end of the report aggregating all the results. The entire project is implemented using 'Python 3' (Jupyter Notebook).

Tasks

Task 1a – Dataset 1 – Energy Data

The Energy dataset, which was used for previous assignment is considered as the First dataset. The dataset has 29 columns and 19735 rows with 1 target variable – 'Appliances' and 28 dependent variables. Correlation matrix is plotted and visualized using a heatmap for all the independent variables. The features – 'date', 'rv1', and 'rv2' are dropped. The target variable is continuous in nature. It is converted to a binary classification problem using a median of 60 for implementing the above listed algorithms.

Task 1b – Dataset 2 – Audit data

The second dataset I have chosen is an Audit dataset from Kaggle. The Audit risk dataset consists of data of various firms and their risk factors. They belong to a multitude of sectors ranging from Irrigation, Public Health, Animal Husbandry to Fisheries, Tourism, Science and Technology. This dataset is developed by a 3rd party Audit company who wishes to calculate and assess risk by analysing the present and historical risk factors thereby facilitating the audit-planning process. The dataset has over 18 columns, with one 'target' variable – 'Risk' (binary).

The primary reason I found this dataset interesting is because I have a strong interest in fraud detection. Moreover, I consider that implementing these algorithms on this dataset will expose me to understand what are the factors of importance that significantly contribute in assessing the 'Risk' of a firm.

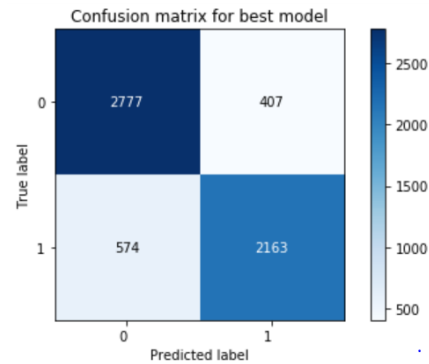
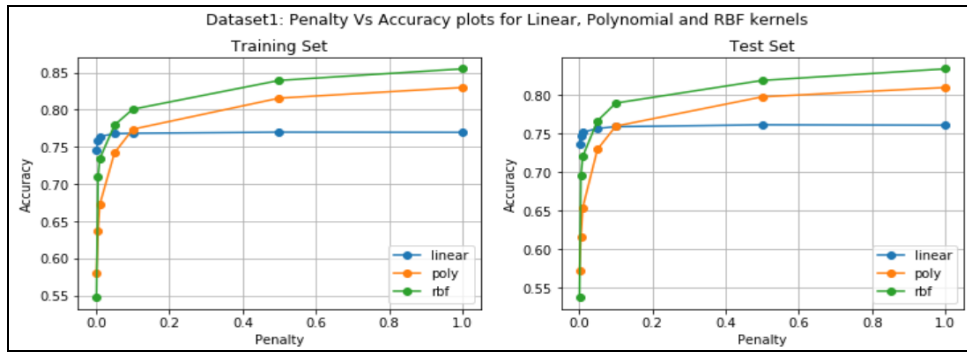
Task 2 – Support Vector Machines

This algorithm is implemented using 'svm' package from 'scikit-learn' library. I have performed 4 experiments on SVM. I have considered linear, polynomial and rbf kernels because most non-linear datasets can be made linearly separable with one of these kernels (with proper control parameters).

Experiment 2.1 – Analysis of 'penalty' parameter for various kernels

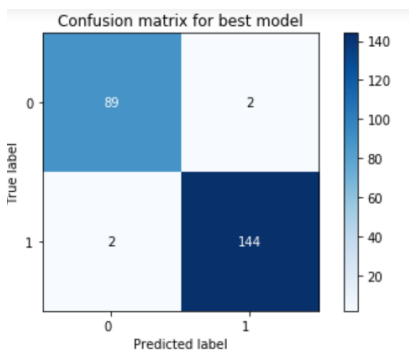
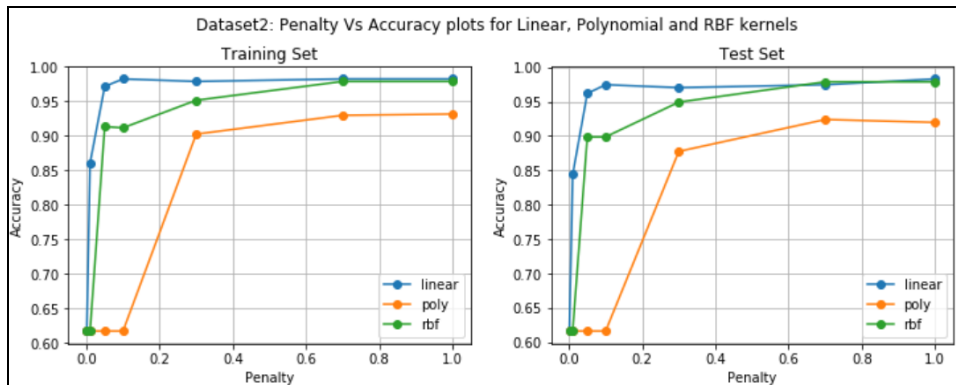
Penalty parameter decides the amount of allowable misclassification and consequently, the width of the margin. If penalty is high, then penalty for misclassification is high, hence margin is small and misclassification rate is small.

Dataset 1 – Rbf kernel is the best kernel as it records best performance. The best value of penalty is 1 as all models records higher train/test accuracies for that value of penalty parameter. (Plots below)



Metrics	Values
Accuracy	0.8343
Recall	0.8722
Precision	0.8287

Dataset 2 – The linear kernel is the best kernel as it records the best train/test accuracies for varying values of Penalty parameter. The best value for the control parameter is 1 as the test accuracies are the best for this value for all the models. The best model's classification metrics are reported below.

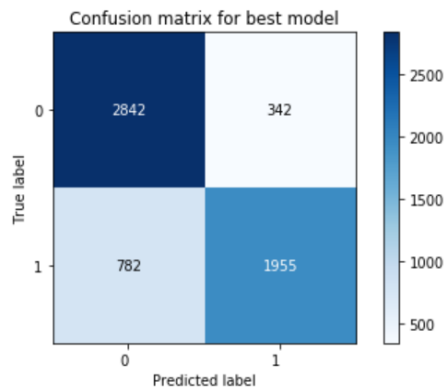
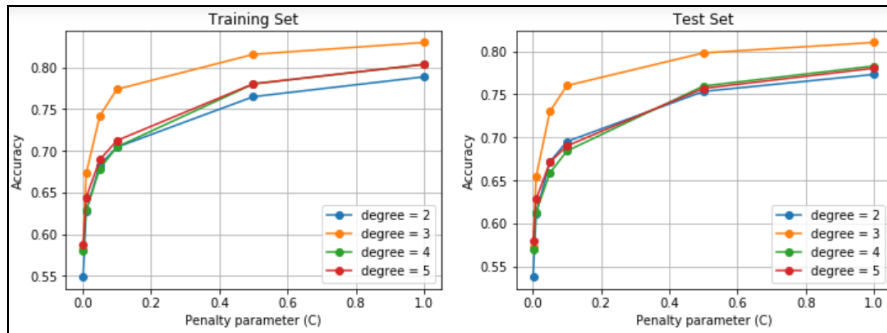


Metrics	Values
Accuracy	0.9831
Recall	0.978
Precision	0.978

Experiment 2.2 – Analysis of ‘degree’ parameter for Polynomial kernel

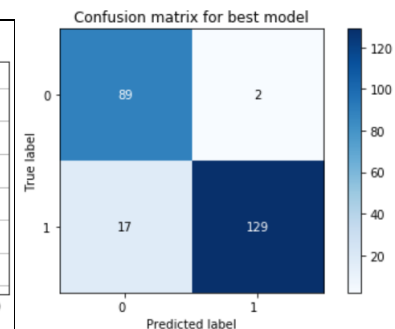
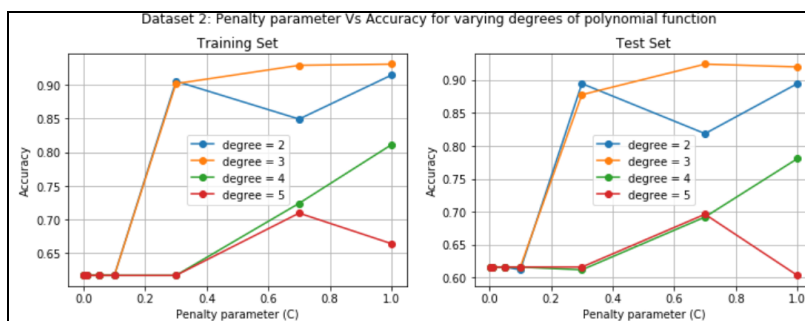
Though polynomial kernel did not record best performance in earlier experiment, I want to experiment if higher degrees of polynomial kernel can fit the data better with varying values of penalty parameter.

Dataset 1 – The polynomial with degree 2 has the best performance with respect to Train and Test accuracies for a penalty parameter value of 1. However, this does not perform better than ‘rbf’ model in experiment 2.1. The classification metrics are reported below.



Metrics	Values
Accuracy	0.8102
Recall	0.8926
Precision	0.7842

Dataset 2 – Here again, polynomial with degree 2 outperforms other models as it records best accuracy values for penalty parameter of 1. Still, this model is not better than the ‘linear’ model in Experiment 2.1. The test accuracy of the previous model is 98.31% > 91.98% which is accuracy of current model. The classification metrics are reported below for reference.

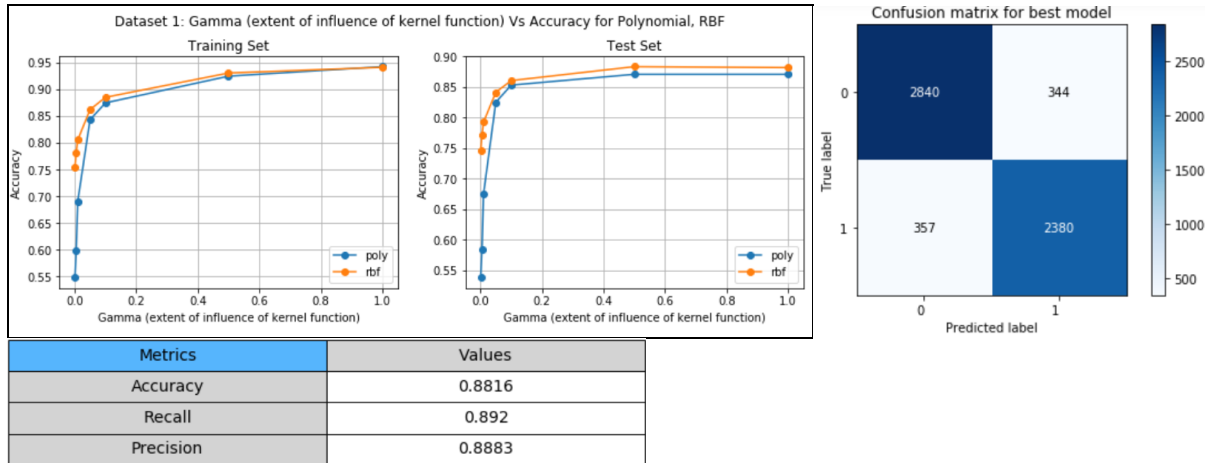


Metrics	Values
Accuracy	0.9198
Recall	0.978
Precision	0.8396

Experiment 2.3 – Analysis of ‘gamma’ parameter for Polynomial, RBF kernel

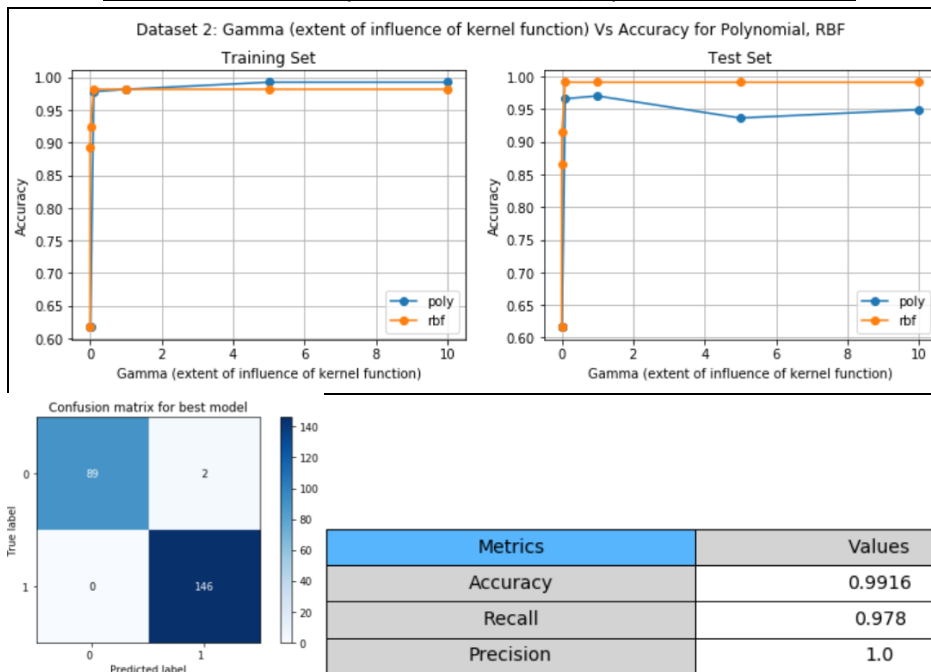
‘Gamma’ is another control parameter for SVM which signifies the ‘extent of influence of the kernel function’. My intention behind this experiment is to find if there is a model that can fit the model better than the previous models by varying the ‘gamma’ control parameter for polynomial and rbf kernels. As ‘gamma’ increases the decision boundaries become more wiggled as compared to lower gamma values.

Dataset 1 – The RBF model outperforms the polynomial model slightly looking at the test accuracy metric at gamma = 1. Here, we have a model that is better than the one obtained in Experiment 2.1.
Current test accuracy – 88.16% > 83.43% (test accuracy of model in Experiment 2.1)



Dataset 2 – The RBF kernel performs better than the polynomial kernel in terms of test accuracy. And, like in dataset 1, we have a better model as compared to the model obtained in Experiment 2.1.

- Current test accuracy 99.16% > 98.31% (Experiment 2.1 model)

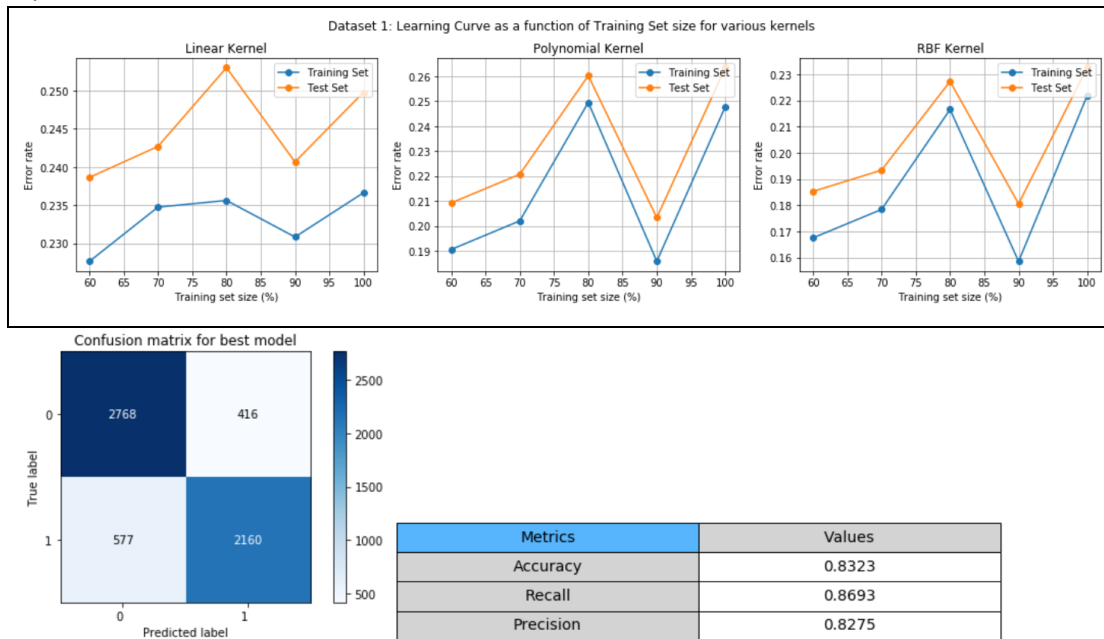


The best ‘Gamma’ value is 1 for which the best accuracy metrics are recorded in above plots for both datasets.

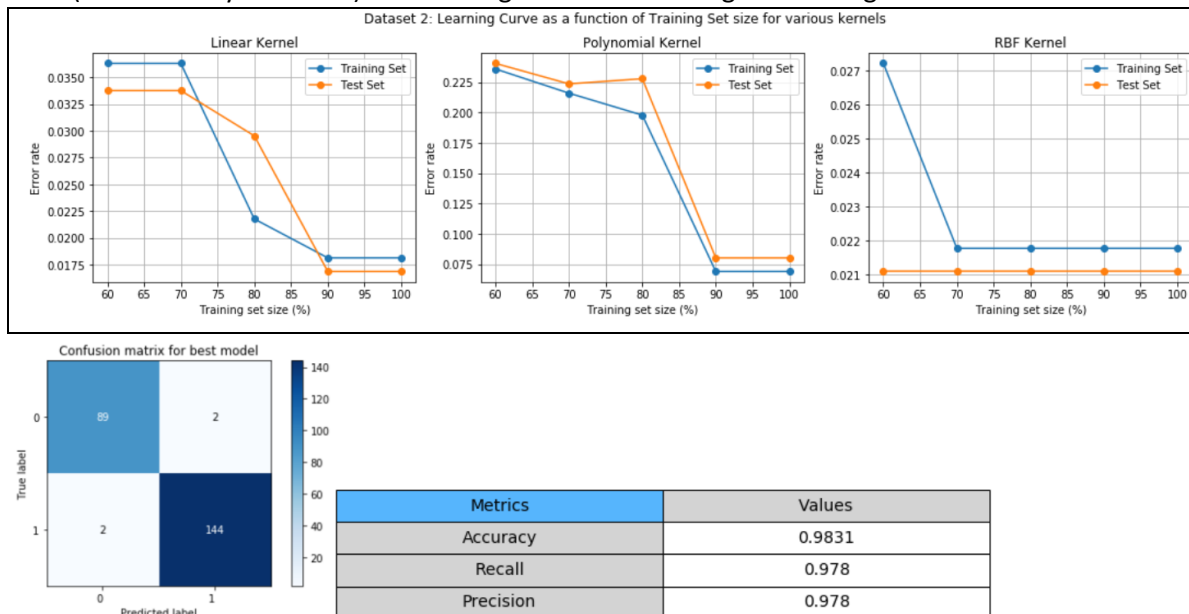
Experiment 2.4 – Learning curves as a function of Training Size for all 3 kernels

In this experiment, learning curves are plotted as a function of varying sizes of the training set. The corresponding test and training error rates are captured and reported to find if there are any models that fit the data better than the previous models. The parameters for the kernels are the best parameters that were found as part of earlier experimentation.

Dataset 1 – The RBF kernel with a penalty=1 reports the best performance at a training size of 90% of original training set in aspects of Train/Test error rates. The classification metrics of the best model is reported below.



Dataset 2 – The linear kernel with penalty=1 reports the best performance with the lowest test error of 1.6% (test accuracy – 98.31%) at a training size of 90% of original training set.



In both datasets, the best model is still the one that is obtained in Experiment 2.3.

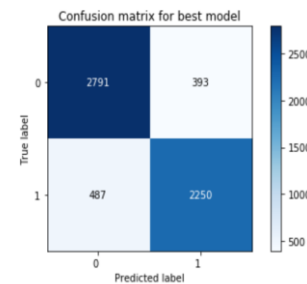
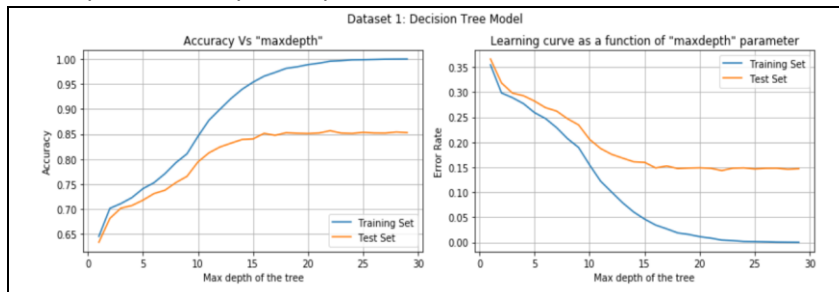
Task 3 – Decision Tree Implementation

The Decision Tree model is implemented using 'DecisionTreeClassifier' package from 'scikit-learn' library. I have chosen 'Entropy' (Information Gain) as the splitting criterion. The primary reason for going with Entropy is it gives us a better sense of impurity at each level for binary classification problems. Gini could be more useful when the target variable is continuous nature. For this task, I have performed 3 experiments involving control parameters maxdepth, pruning size and training size.

Experiment 3.1 – Learning Curves as a function of 'maxdepth'

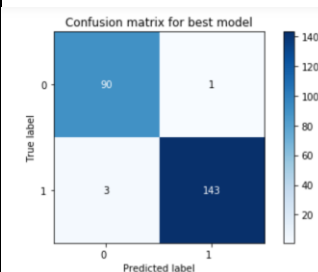
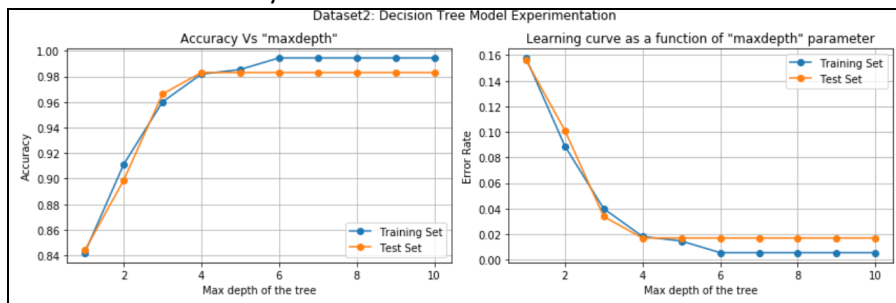
The control parameter 'maxdepth' denotes the depth of the decision tree while it is being grown. As 'maxdepth' increases, the training accuracy of the model increases, but past a point the model starts to overfit resulting in poor test accuracy. Obtaining the ideal 'maxdepth' value through learning curves is the intention behind this experiment.

Dataset 1 - Past the value of 16, even though the training accuracy increases, the validation accuracy saturates and slightly fluctuates. Hence the ideal 'maxdepth' for Dataset1 is 16. The Test accuracy is 85.14% with a True positive rate of 87.66% meaning the model has correctly classified 87.66% of the actual positive samples as positive.



Metrics	Values
Accuracy	0.8514
Recall	0.8766
Precision	0.8514

Dataset 2 - For Dataset 2, interpreting similarly, it can be seen that past the value of 4, the test accuracy saturates even though training accuracy keeps increasing indicating the condition of overfitting. Hence, the ideal 'maxdepth' value is 4 for this model. The classification matrix for this model is plotted above. It records a test accuracy of 98.31% with a Recall of 98.9%.

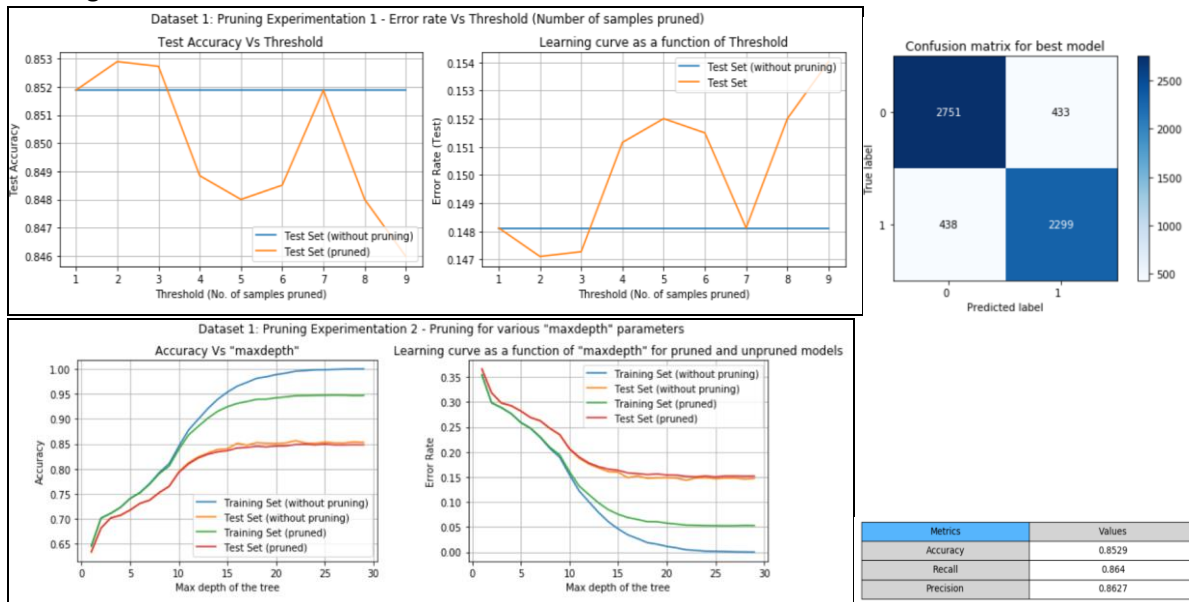


Metrics	Values
Accuracy	0.9831
Recall	0.989
Precision	0.9677

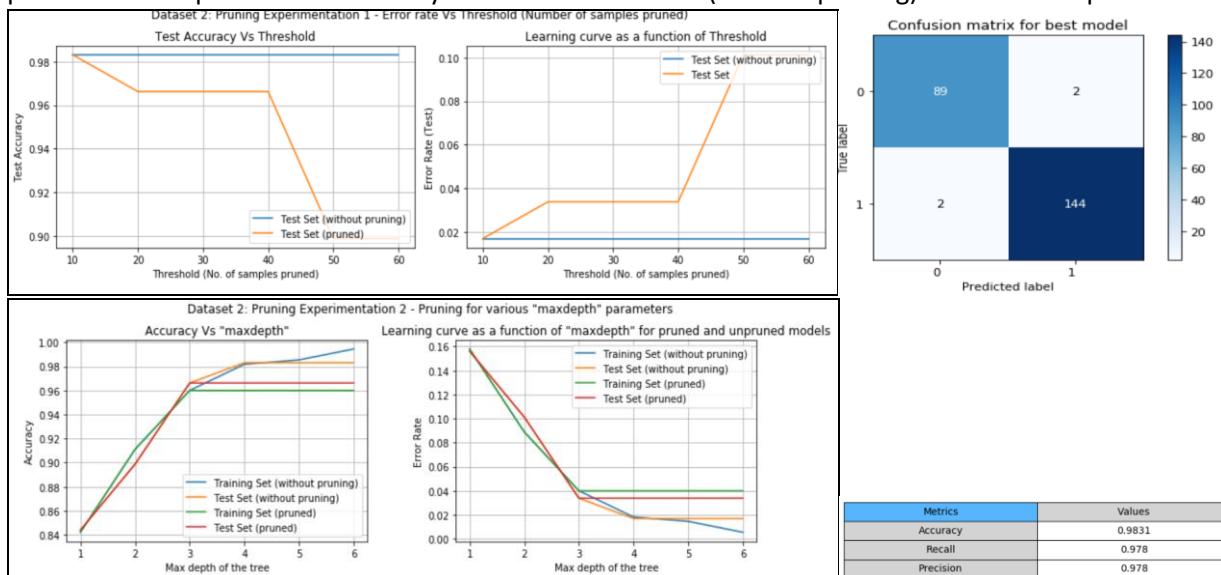
Experiment 3.2 – Pruning

In this experiment, I have chosen ‘minimum number of samples in a leaf node’ as pruning condition as it is simplistic and easy to analyze. If a leaf node has the number of samples below the “*threshold*”, it is cut and taken off the decision tree. Intention behind this experiment to understand effect of pruning on overfitted models.

Dataset 1 - At *threshold*=2, the test error is the least and it performs better than the original unpruned model. For this experiment, the *maxdepth* was not set to specify a fully grown decision tree. The experiment is performed for varying ‘*maxdepth*’ values and fixed pruning threshold of 2. The below plot tells us that as ‘overfitting’ starts occurring (past 16), the pruned model performs reasonably similar to the original model.



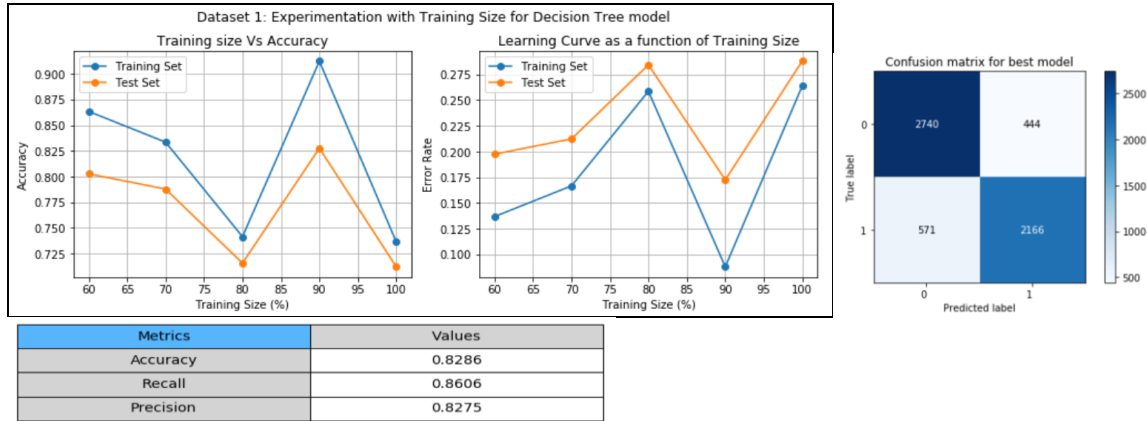
Dataset 2 - The pruned model performs reasonably better at ‘*threshold*’=10. The test accuracy is significantly high at 98.31%. Also, it can be interpreted that, for a pruning with threshold of 10, the pruned test set performs reasonably similar to the test set (without pruning) until a *maxdepth* of 3.



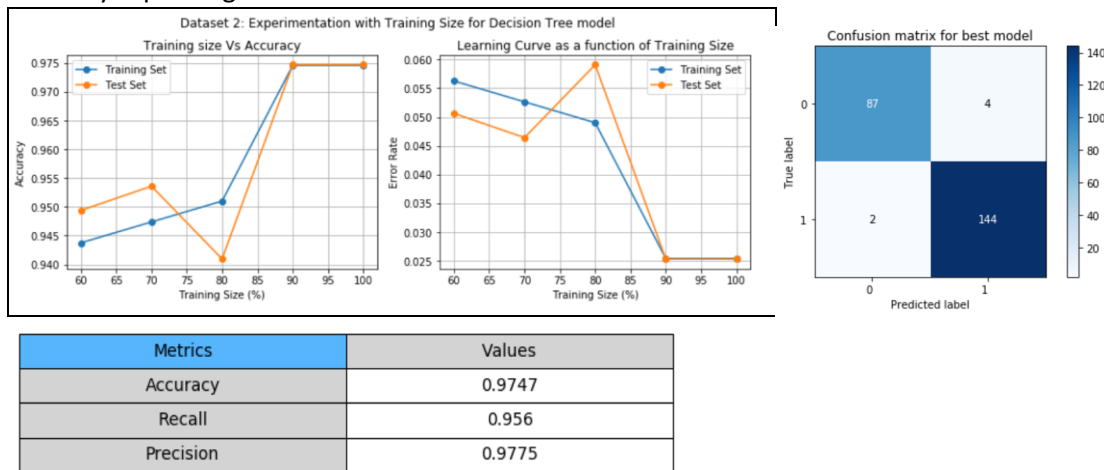
Experiment 3.3 – Experimenting with Training Size

In this experiment, the control parameter is the training size. This experiment is performed to understand the relationship between Training set size and train/test error rates.

Dataset 1 - The best model is the one with a training size of 90% of the entire train set. Looking at the confusion matrix, the test accuracy is 82.86% with a Recall rate of 86.06%.



Dataset 2 – Here also, the best model is the one with training size of 90% of original training set, the Test accuracy is peaking at 97.47%.



The main takeaway in this experiment is that it is not expected that as training size increases, test accuracy increases. The distribution of data is the key point to consider. Based on the distribution of the data, the learnings can be different, and the model can perform differently on the test set.

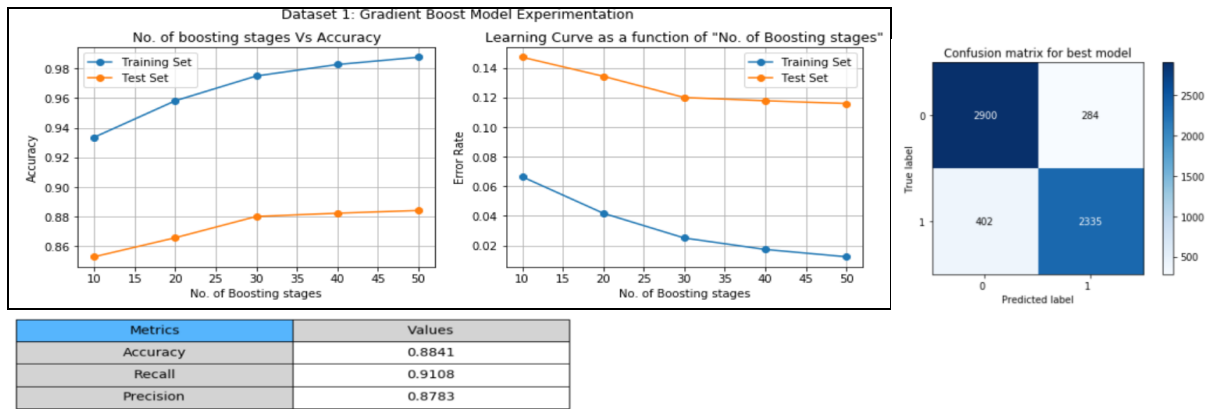
Task 4 – Boosted version of Decision Trees

In this task, I have implemented the Gradient boost model for Decision Trees using the 'GradientBoostingClassifier' package from 'scikit-learn' library. I have performed 2 experiments using control parameters – Number of boosting stages and pruning.

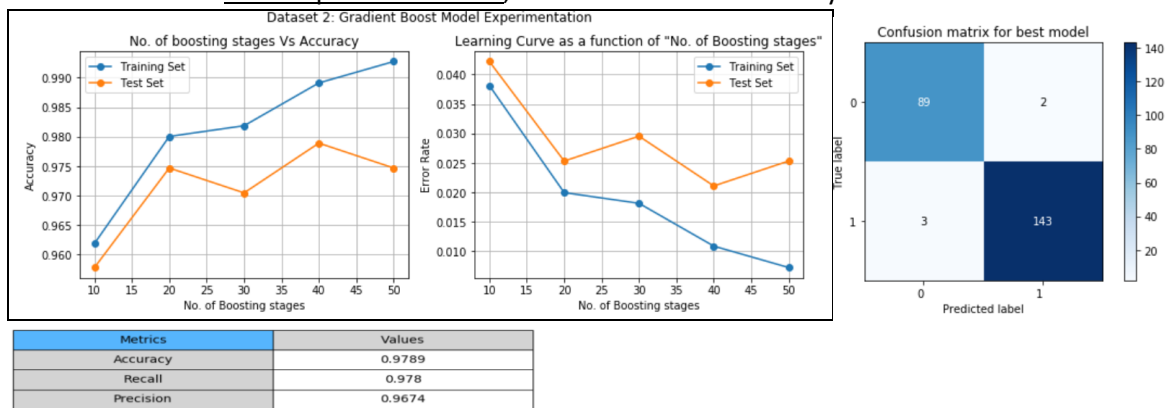
Experiment 4.1 – Experimenting with 'Number of boosting stages'

'Boosting stages' refers to the number of times the model is sequentially worked upon on its misclassification. The intention behind this experiment is to understand the relationship between the control parameter and train/test error rates.

Dataset 1 - In the considered range, 50 is the best value for the control parameter as evident by the highest test/train accuracies. The classification metrics for this model is reported below.



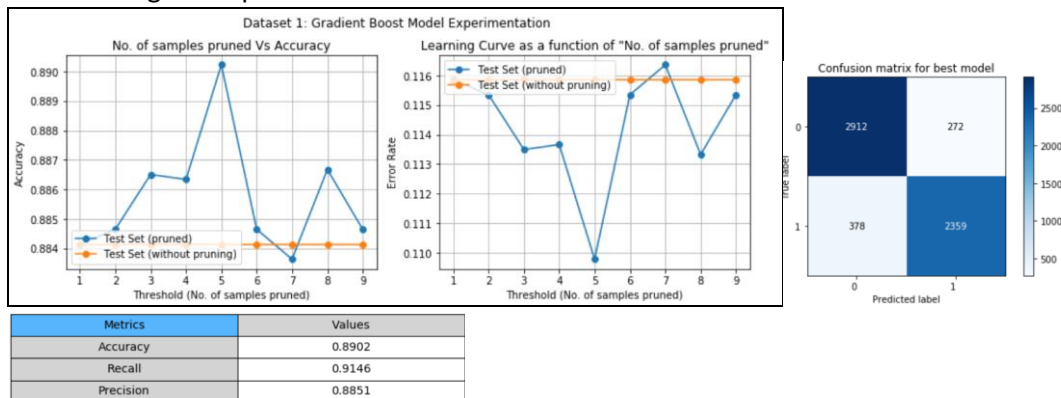
Dataset 2 - As 'No. of boosting stages' increases, the test accuracy fluctuates exhibiting no pattern. The best value for the control parameter is 40, for which the test accuracy is 97.89%.



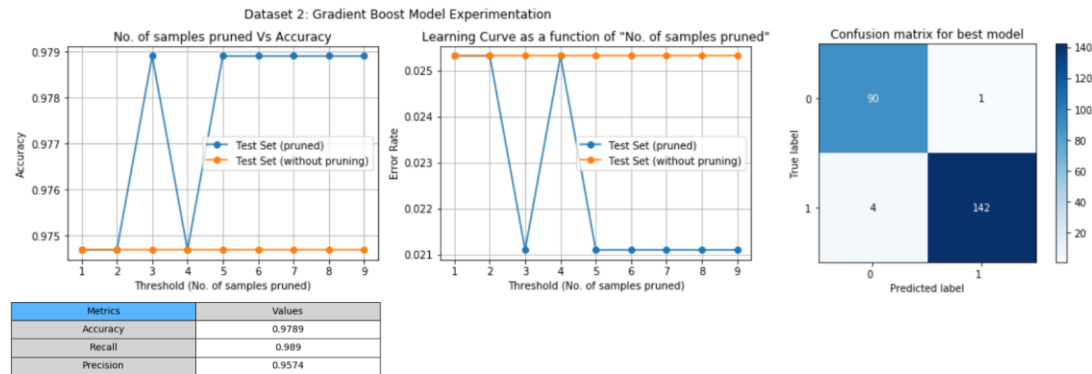
Experiment 4.2 – Experimenting with Pruning

Pruning is performed by cutting nodes with less than 'threshold' number of samples off the tree.

Dataset 1 – Best model is when threshold = 5, recording test accuracy of 89.02% and performing better than the original unpruned model.

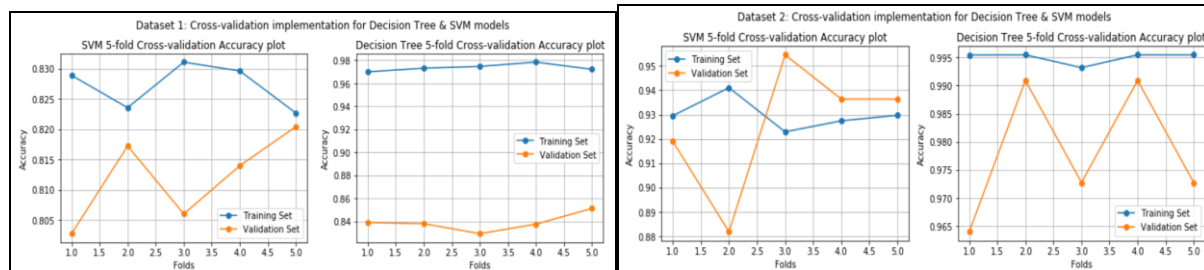


Dataset 2 - The best model is when the control parameter is 5. Key takeaway from this experiment is that pruning does help in arriving at a better model. The best model is arrived at through learning curves and experimentation.



Task 5 – Cross Validation Implementation

5-fold cross-validation is implemented for both datasets. This experiment is performed to understand the importance of cross-validation sets. For dataset 1, the 5th fold records very low train/test error and produces the best model for both SVM and Decision Tree models. For dataset 2, for SVM, the 5th fold produces the best model because the train and test accuracies are comparable and don't exhibit extremities. For Decision Tree, 2nd and 4th fold produce the best results for dataset 2.



Conclusion / Discussion

- Dataset 1 – Best model – Boosted model (boosting stages = 50) – 88.41% test accuracy**
 The boosted model is the best model. In my opinion, this is because the dataset is dense, not linearly separable and no single model can precisely capture all the variance exhibited by the dataset. Hence, an ensemble model (a series of models) that sequentially works on misclassification of previous models is the best model.
 SVM and Decision Tree models performed reasonably similar on Dataset 1. The SVM RBF kernel was the only model that neared the best model with an accuracy of 88.16% for a gamma value of 1. This is just an indication of how difficult to separate the data linearly.
- Dataset 2 – Best model – RBF SVM (gamma = 1) – 99.16% test accuracy**
 The dataset has over 1k rows. Being sparse, the data is linearly separable on a higher dimension and hence, an RBF kernel is able to classify the data almost perfectly.
 SVM performed way better than Decision Tree in the 2nd dataset. This is primarily because the dataset is linearly separable. Though an ensemble model sequentially works on misclassification errors, it could not stand out as compared to SVM, the main reason being the dataset's distribution and size.
- There is no "best" algorithm on a global basis. It solely depends on the dataset, its size, distribution, variance etc. From my perspective, the "best" algorithm would have optimal training accuracy and the best testing accuracy. The "best" model never occurs by chance, it must be worked upon incessantly. Extensive pre-processing, visualization and proper tuning using hyperparameters/control parameters of the built model is mandatory for to get the "best model".