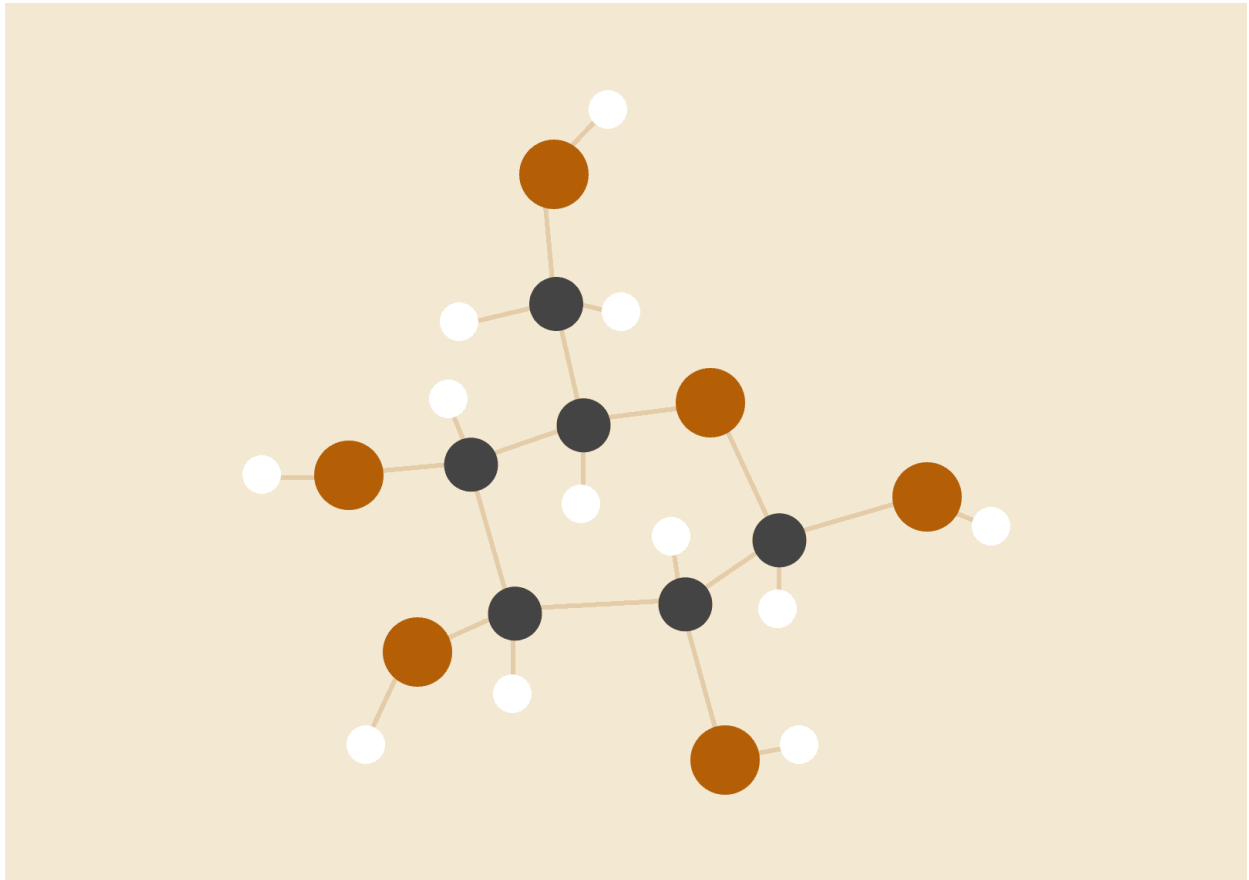


MACHINE LEARNING PROJECT REPORT

CREDIT RISK ANALYSIS USING MACHINE LEARNING



GROUP 5

RAJEEV RANJAN 20BME044
YAGNESH SHETTY 20BSM060
PULKIT SATYARTHI 20BSM043

System under study

A significant activity of the banking industry is to extend credit to customers. Credit risk management evaluates available data and decides the credibility of a customer, with the intent of protecting the financial institution against fraud.

We apply three machine learning classification algorithms. A summary of the methods are as follows:

1. Logistic Regression (LR) —Logistic Regression fits a logistic (sigmoid) function to the data. This is useful when the dependent variable is binary. To compute the weights, we will utilize the Iteratively Re-weighted Least Squares algorithm (IRLS).
2. Support Vector Machine (SVM) —SVM identifies hyperplanes to separate the data based on labels. To handle nonlinear data, projections using various kernels can be used.
3. Classification Trees –We will implement two types of classification trees: Decision Tree and Random Forest. Decision Tree is similar to a flowchart and is very easy to visualize. Every terminal node represents the output. Random Forest is a collection of decision trees, where the majority vote is taken for prediction.

Data Collection

The UC Irvine Machine Learning repository (UCI - ML) contains a collection of datasets useful for evaluating machine learning algorithms. Australian credit dataset from the UCI - ML 1 repository were used for this project. The Australian dataset has 14 characteristics and 690 instances. The response variable is a binary decision, whether a customer is credible or not. The datasets have some attributes such as the credit score, the purpose of the loan and customer information (occupation, salary, age and account duration).

Data Exploration, Feature selection and Data Cleaning

For data selection ,preprocessing,feature selection and creating dummy variable we have used various module of sklearn such as label encoders,onehotencoder and standardscaler.

```
[ ] def data_preprocessing(dataset) :  
    listt = dataset.select_dtypes(include =['category', object]).columns  
    X = dataset.iloc [: , :-1].values  
    y = dataset.iloc [: , -1].values  
    for i in range(0 , len(y)) :  
        if y[i] == 2:  
            y[i] = 0  
  
    listn = np.empty((len(listt ),1))  
    for i in range (0 , len(listt )):  
        listn [i ,0] = dataset.columns.get_loc(listt[i])  
  
    # Categorical data  
    from sklearn.preprocessing import LabelEncoder , OneHotEncoder  
    labelencoder_X = LabelEncoder ()  
    for i in range (0 , len( listn ) ):  
        X[: , int(listn [i ,0]) ] = labelencoder_X .fit_transform (X[: , int (listn [i ,0]) ] )  
        onehotencoder = OneHotEncoder(categorical_features = [int(listn[i ,0])])  
  
    # Splitting the dataset into the Training set and Test set  
    from sklearn.model_selection import train_test_split  
    X_train , X_test , y_train , y_test = train_test_split(X, y, test_size = 0.2)  
  
    # Feature Scaling  
    from sklearn.preprocessing import StandardScaler  
    sc = StandardScaler ()  
    X_train = sc.fit_transform(X_train)  
    X_test = sc.transform(X_test)  
    return X_train ,X_test , y_train , y_test ,X,y
```

Machine learning Model

LOGISTIC REGRESSION-

Suppose we have N training samples (x_i, y_i) where the response variable, y_i , is a binary variable, i.e., $y_i \in \{0, 1\}$. The log-odds then becomes,

$$\ln \left(\frac{p(y_i | x_i)}{1 - p(y_i | x_i)} \right) = w^T x_i, \quad (3.1)$$

where $p(y_i | x_i)$ is the posterior probability and w is the weight vector that we seek. Logistic Regression (LR) models are fitted via maximum likelihood. Solving eq. (3.1) for the class posterior probability,

$$\begin{aligned} Pr(y_i = 1) &= p(y_i | x_i) = \frac{e^{w^T x_i}}{1 + e^{w^T x_i}}, \\ Pr(y_i = 0) &= 1 - p(y_i | x_i) = \frac{1}{1 + e^{w^T x_i}}. \end{aligned} \quad (3.2)$$

We recall Bernoulli RVs (random variables) are discrete RVs with only two values, typically 0 and 1. The probability mass function is completely specified with one parameter: the “probability of success”, $p = P r(X = 1)$,

$$pr(k) = p^k (1 - p)^{(1-k)}, \quad k \in \{0, 1\}.$$

We note that for binary classification, each y_i can be modelled as a Bernoulli RV. Since y_i is Bernoulli RV, the (joint) likelihood is

$$\mathcal{L}(w) = Pr(\{y_1, y_2, \dots\} | w, \{x_1, x_2, \dots\}) = \prod_{i=1}^N p^{y_i} (1 - p)^{(1-y_i)}.$$

Hence,

$$\begin{aligned}
\ln \mathcal{L}(w) &= \sum_{i=1}^N [y_i \ln(p) + (1 - y_i) \ln(1 - p)] \\
&= \sum_{i=1}^N \left[y_i \ln \left(\frac{e^{w^T x_i}}{1 + e^{w^T x_i}} \right) + (1 - y_i) \ln \left(1 - \frac{e^{w^T x_i}}{1 + e^{w^T x_i}} \right) \right] \quad (3.3) \\
&= \sum_{i=1}^N \left[y_i w^T x_i - \ln \left(1 + e^{w^T x_i} \right) \right].
\end{aligned}$$

The task is to maximize the (log) likelihood, eq. (3.3). This is equivalent to minimizing the negative log likelihood (NLL)

$$\text{NLL}(w) = - \sum_{i=1}^N \left[y_i w^T x_i - \ln \left(1 + e^{w^T x_i} \right) \right].$$

Gradient descent is a first-order minimization method that uses information about the gradient. The update formula looks like

$$w^{n+1} = w^n - \gamma \left. g(w) \right|_{w^n}.$$

where γ is a stepsize to be computed, and the gradient satisfies

$$g(w) = \nabla_w \text{NLL}(w) = X^T (p - y). \quad (3.4)$$

Here, X is $N \times (p + 1)$ matrix of x_i values, p the vector of fitted probabilities has elements $p_i = p(x_i, w)$, y is a vector of y_i values. Newton's method is a second-order optimization method which uses the gradient and the Hessian (curvature). The update formula looks like

$$w^{n+1} = w^n - \tilde{\gamma} H^{-1}(w) \big|_{w^n} g \big|_{w^n}, \quad (3.5)$$

where $\tilde{\gamma}$ is a stepsize to be computed, the gradient is defined above in eq. (3.4), and the Hessian H satisfies,

$$H = \nabla^2 \text{NLL}(w) = X^T S X. \quad (3.6)$$

Here, S is a $N \times N$ diagonal matrix with entries $s_i = p_i(1 - p_i)$. Both methods converge, however, Newton's method is faster since it takes curvature into account. Substituting the definitions of H , eq. (3.6) and g eq. (3.4) into eq. (3.5), we get

$$\begin{aligned} w^{n+1} &= w^n - H^{-1} g \\ &= w^n - (X^T S_n X)^{-1} X^T (y - p_n) \\ &= (X^T S_n X)^{-1} ((X^T S_n X) w^n - X^T (y - p_n)) \\ &= (X^T S_n X)^{-1} X^T (S_n X w^n + y - p_n) \\ &= (X^T S_n X)^{-1} X^T S_n z, \end{aligned} \quad (3.7)$$

where we define the working response or adjusted response,

$$z_n = X w^n + S_n^{-1} (y - p_n).$$

Iteratively Re-weighted Least Squares (IRLS)

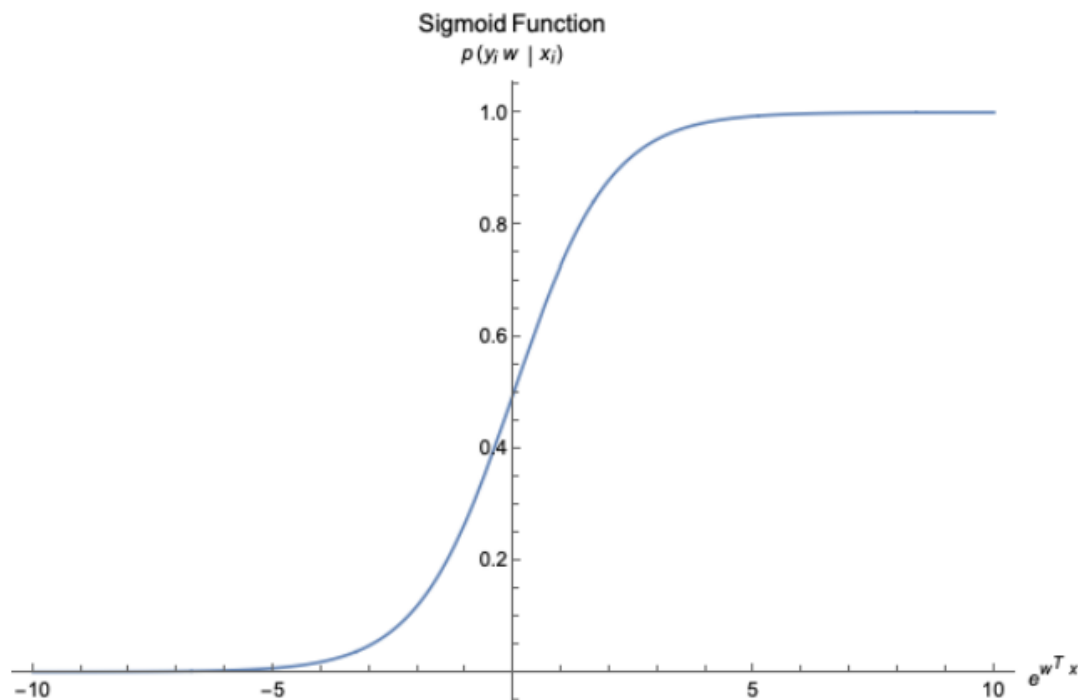
We have seen previously, eq. (3.7) that after each iteration, p_n changes which in turn changes S_n and z_n . Therefore we update the weights, s_n , after every iteration, and solve a weighted least squares problem. The algorithm for finding the weights is summarized

in algorithm 1.

Algorithm 1 Iteratively Reweighted Least Squares (IRLS)

```
1:  $w = 0$ 
2: while not converged do
3:   for  $i = 1 : N$  do
4:      $p_i = \frac{e^{w^T x_i}}{1 + e^{w^T x_i}};$ 
5:      $s_i = p_i(1 - p_i);$ 
6:      $z_i = w^T x_i + \frac{y_i - p_i}{s_i}$ 
7:   end for
8:    $S = \text{diag}(s_1 \ s_2 \ \dots \ s_N)$ 
9:    $w = (X^T S X)^{-1} X^T S z$ 
10: end while
```

Predictions with the Sigmoid Function: So for a new data point x , compute $p(y_i | w, x)$. If $p(y_i | w, x) \geq 0.5$ then classify x as Class A. If $p(y_i | w, x) < 0.5$ then classify x as Class B.



Support Vector Machine (SVM)

Non-Linear SVM (Kernel Trick)

Suppose the data is non-linearly separable. SVM can still be applied to separate data with the use of a kernel trick – a kernel is used to transform the data to a higher dimensions which is then separable with hyperplanes. Radial Basis Functions (RBF) and polynomial kernels are two popular kernel functions used in the SVM community.

$$\dagger \text{ RBF Kernel: } K(x_i, x_k) = e^{-\gamma \|x_i - x_k\|^2}$$

$$\dagger \text{ } d\text{-th degree Polynomial Kernel: } K(x_i, x_k) = (1 + x_i^T x_k)^d$$

We can rewrite the Wolfe dual minimization problem using the kernel functions,

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^N \alpha_i \alpha_k y_i y_k K(x_i, x_k) - \sum_{i=1}^N \alpha_i \\ \text{subject to} \quad & \sum_{i=1}^N \alpha_i y_i = 0, \quad \alpha_i \geq 0, \quad i = 1, \dots, N \end{aligned}$$

Decision Tree

Mathematical Framework

Suppose we view a data set in terms of feature vectors, i.e., our $N \times p$ matrix has p feature vectors,

$$X = \begin{bmatrix} \vdots & \vdots & \dots & \vdots \\ \tilde{x}_1 & \tilde{x}_2 & \dots & \tilde{x}_p \\ \vdots & \vdots & \dots & \vdots \end{bmatrix}$$

The ID3 algorithm for finding the decision tree proceeds as follows.

1. Set X to be the entire data set.
2. Compute the entropy of X, where the entropy is defined as

$$H(X) = \sum_{c \in C} -p(c) \log_2 p(c).$$

For our credit risk analysis, the response variable has binary values, $C = \{\text{Yes}, \text{No}\}$, and $p(c)$ is the respective probability.

3. Calculate the weighted sum entropy of a particular feature.

$$H(X | \tilde{x}) = \sum_x p(\tilde{x}) \sum_{c \in C} -p(c | \tilde{x}) \log_2 p(c | \tilde{x})$$

4. Calculate information gain for each feature and partition X using the maximum gain.

$$IG(X, x) = \sum_{c \in C} -p(x) \log_2 p(x) - \sum_x p(x) \sum_{c \in C} -p(c | x) \log_2 p(c | x)$$

5. Repeat steps 2-5 with X as each child node until we get the desired decision tree.

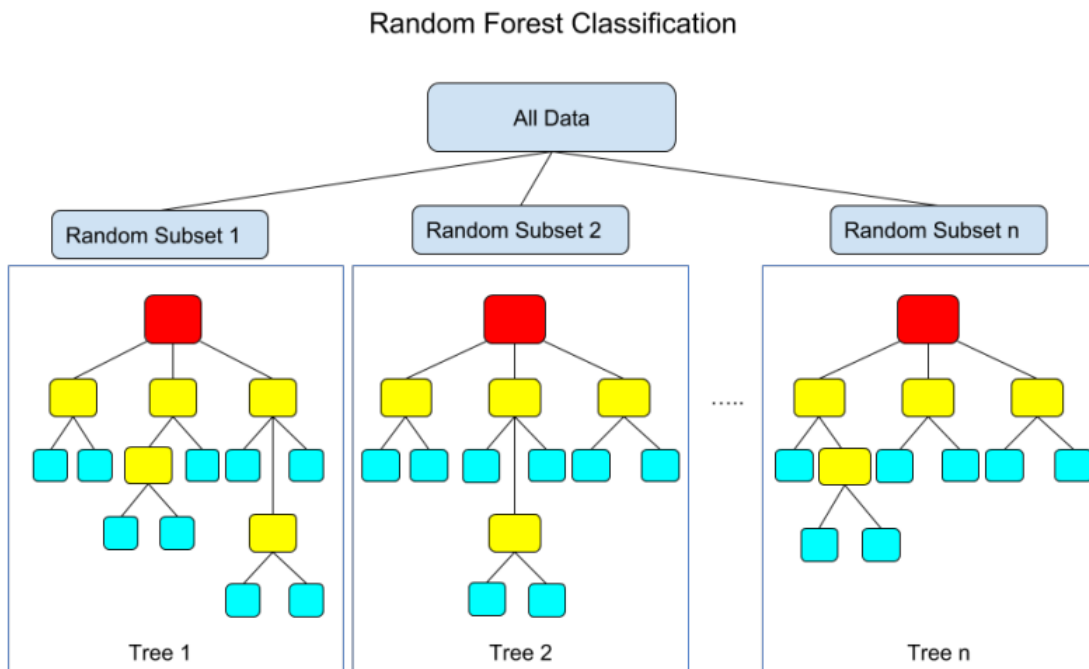
The CART algorithm is similar except that the entropy computation, eq. (6.1) is replaced by computing the gini index,

$$G(X) = \sum_{c \in C} p(c) (1 - p(c))$$

Random Forest Classification

Methodology

The Random Forest Classification is an algorithm which collects decision trees using the random subsets of features and choosing the majority vote among them for the classification. This method minimizes overfitting and increases the overall accuracy. Figure below shows how the random forest classification looks. So, for a given data each tree will make a vote for the prediction, and the majority vote will be considered as the prediction.



Feature Selection

In Logistic Regression

In SVM we will use parameters such as kernel , poly , rbf (radial basis function) along with various degrees such as 1 ,2 , 3 ,4 ,5 , 6 ,7 .

In Decision Tree the criteria used are gini and entropy along with this minimum samples split is taken into consideration with values such as 2 , 4 , 6 , 8 , 15

In Random Forest classification criteria such as gini and entropy and minimum samples split remain same and another factor of number of samples is taken into consideration.

HyperParameters Optimisation

PARAMETRIC TUNING- it is done on various value of C and tolerance.

```
[28] # Best parameters for Logistic Regression
parameters_lr = [{ 'C': [1 , 5 , 10] , 'tol': [1e-4 ,1e-5 ,1e-6 ,1e-10]}]
best_para_lr = choosing_parameters ( parameters_lr , classifier_lr , X_train , y_train )
print(best_para_lr)
C_lr = best_para_lr [0]
tole_lr = best_para_lr [1]

[1, 0.0001]
```

Parametric tuning:-

The other tuning values for degree are from 1 to 7 and for kernel, RBF is also considered as different kernel. However, for Australian data set, it gives even good accuracy by setting the polynomial kernel with degree 1.

```
# Best parameters for SVM
parameters_svc = [{ 'kernel': ['rbf','poly'] , 'degree': [1 , 2 , 3 , 4 , 5 , 6 , 7]}]
best_para_svc = choosing_parameters(parameters_svc , classifier_svc , X_train , y_train )
ker_svc = best_para_svc[1]
deg_svc = best_para_svc[0]
print(best_para_svc)

[1, 'rbf']
```

Parametric tuning:-

Parameters for criterion used are gini and entropy ; for mini sample split 2,4,6,8,10,15.

Out of all parameter best result comes from entropy(criterion) and 10(split).

```
[31] # Best parameters for Decision Tree
parameters_dt = [{ 'criterion': ['gini','entropy'] , 'min_samples_split': [2 ,4 ,6 ,8 ,10 ,15]}]
best_para_dt = choosing_parameters( parameters_dt , classifier_dt , X_train , y_train )
criteria_dt = best_para_dt [0]
min_split_dt = best_para_dt [1]
print(best_para_dt)

['entropy', 10]
```

Parametric tuning:-

Parameters for criterion used are gini and entropy ; for mini sample split 2,5,10,15,20; for estimators (10,20,30,40,50,100,150)

Out of all parameter best result comes from gini(criterion) ,15(split) and 30(estimators).

```
# Best parameters for Random Forest
parameters_rf = [{ 'criterion': ['gini','entropy'] , 'min_samples_split': [2 ,5 , 10 , 15 , 20] ,
                    'n_estimators': [10 , 20 , 30 , 40 , 50 ,100 ,150]}]
best_para_rf = choosing_parameters ( parameters_rf ,classifier_rf , X_train , y_train )
print(best_para_rf)
min_samp_spl_rf = best_para_rf [1]
n_est_rf = best_para_rf [2]
criterion_rf = best_para_rf [0]

['gini', 15, 30]
```

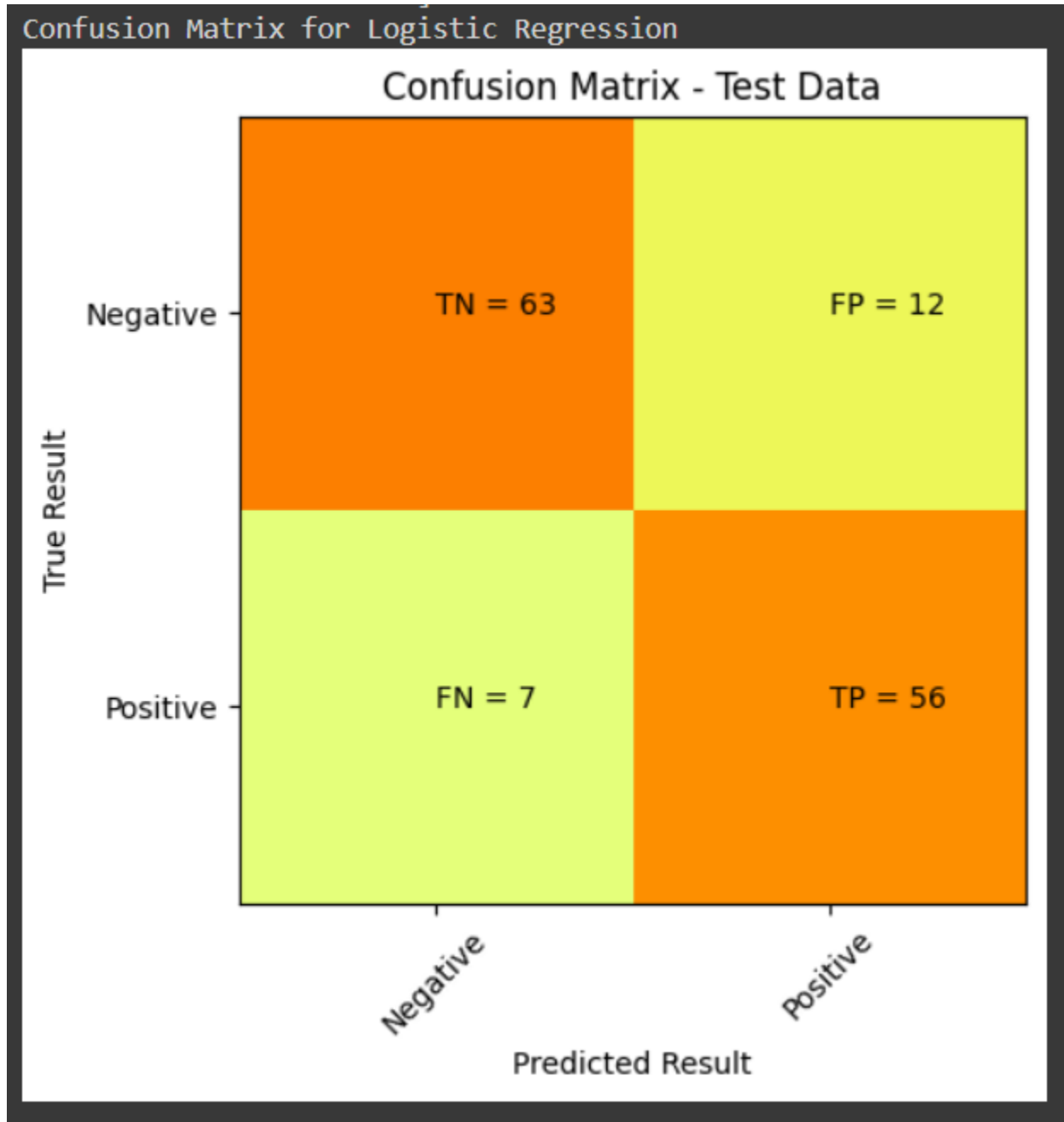
Coding the Actual System

The code of the system can be found in the attachment with the report.

Result and Analysis

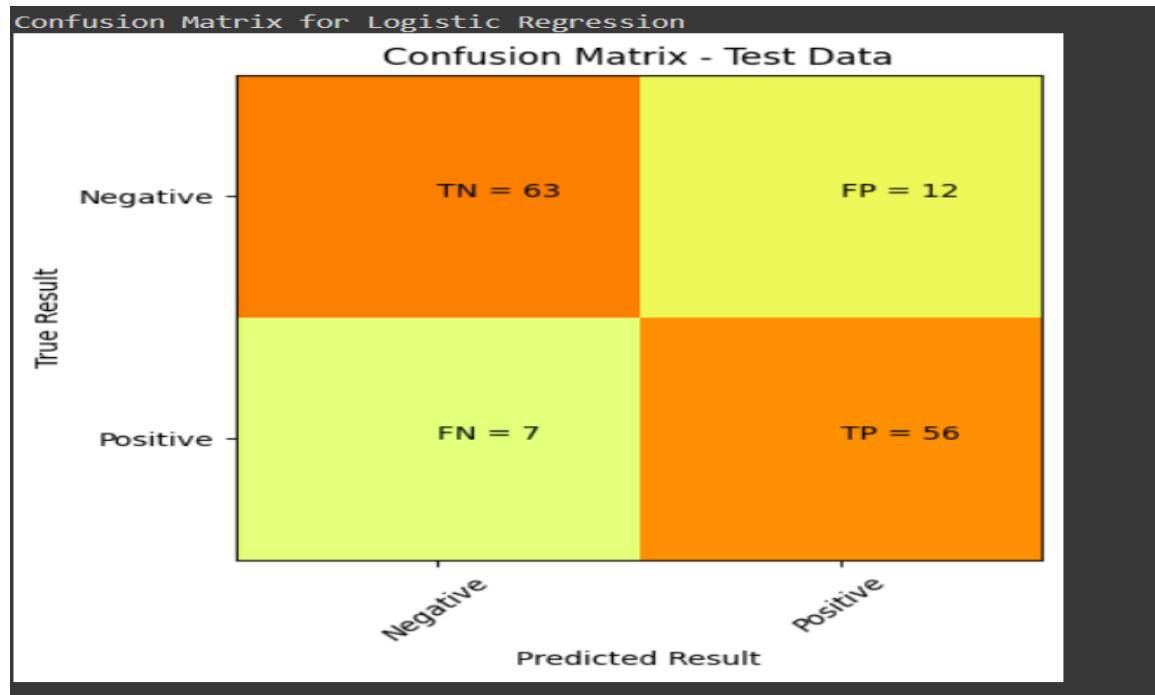
Mean accuracy and confusion matrix before tuning in logistic Regression are as follows:-

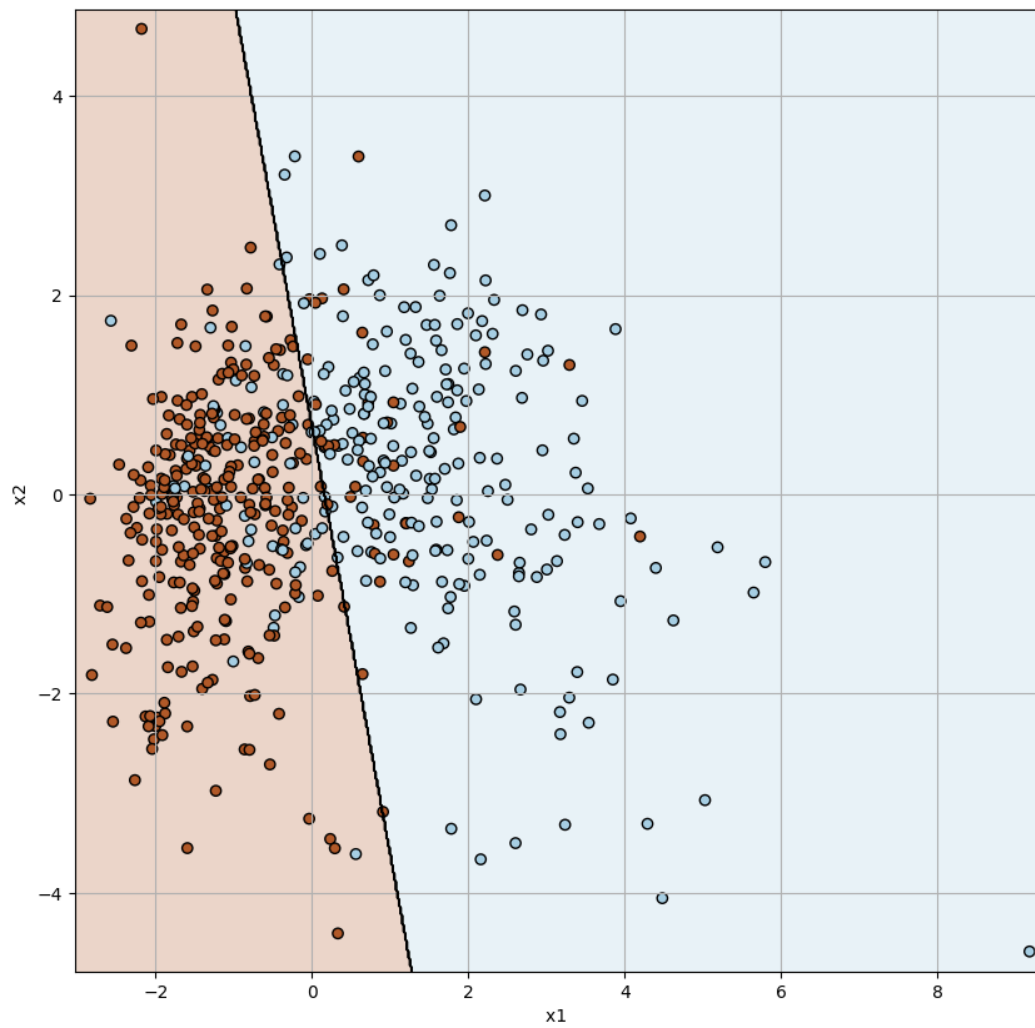
LR: Mean Accuracy before tuning 86.21428571428572



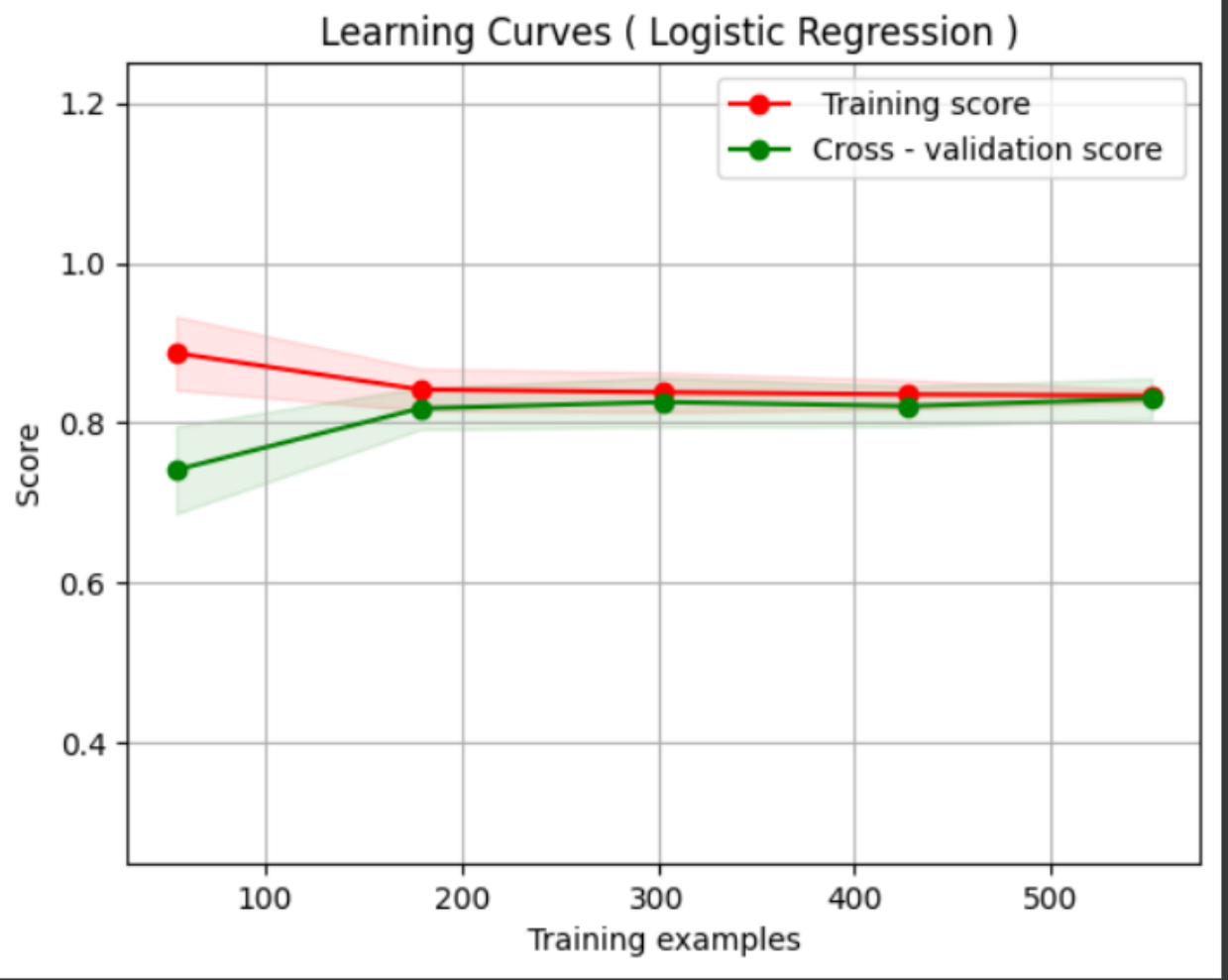
Mean accuracy and confusion matrix after tuning are as follows:-

LR : Mean Accuract after tuning 86.21428571428572





Decision boundary plot using Logistic Regression of the Australian data set.



Learning curve of Logistic Regression using the Australian data set. No high bias or high variance is observed, indicating that the LR model is not over-fitting or under-fitting this Australian data set.

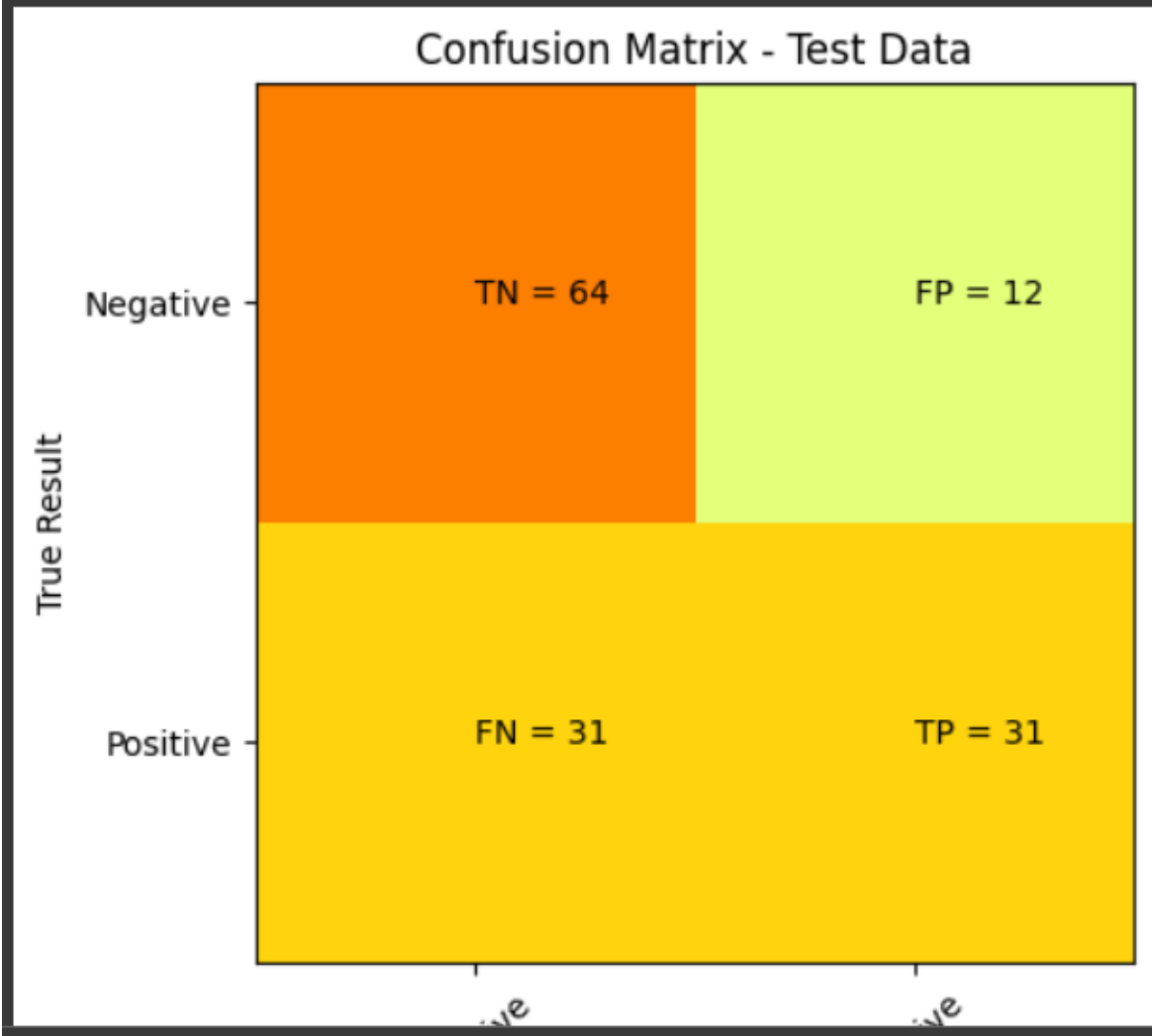
Results (both default and after tuning)

By applying Kernel SVM to Australian credit data sets, the accuracies are obtained shown in table. The default parameters for SVM is set to be polynomial kernel with degree 2.

Mean accuracy and confusion matrix before tuning are as follows:-

SVC : Mean Accuracy before tuning 71.14285714285715

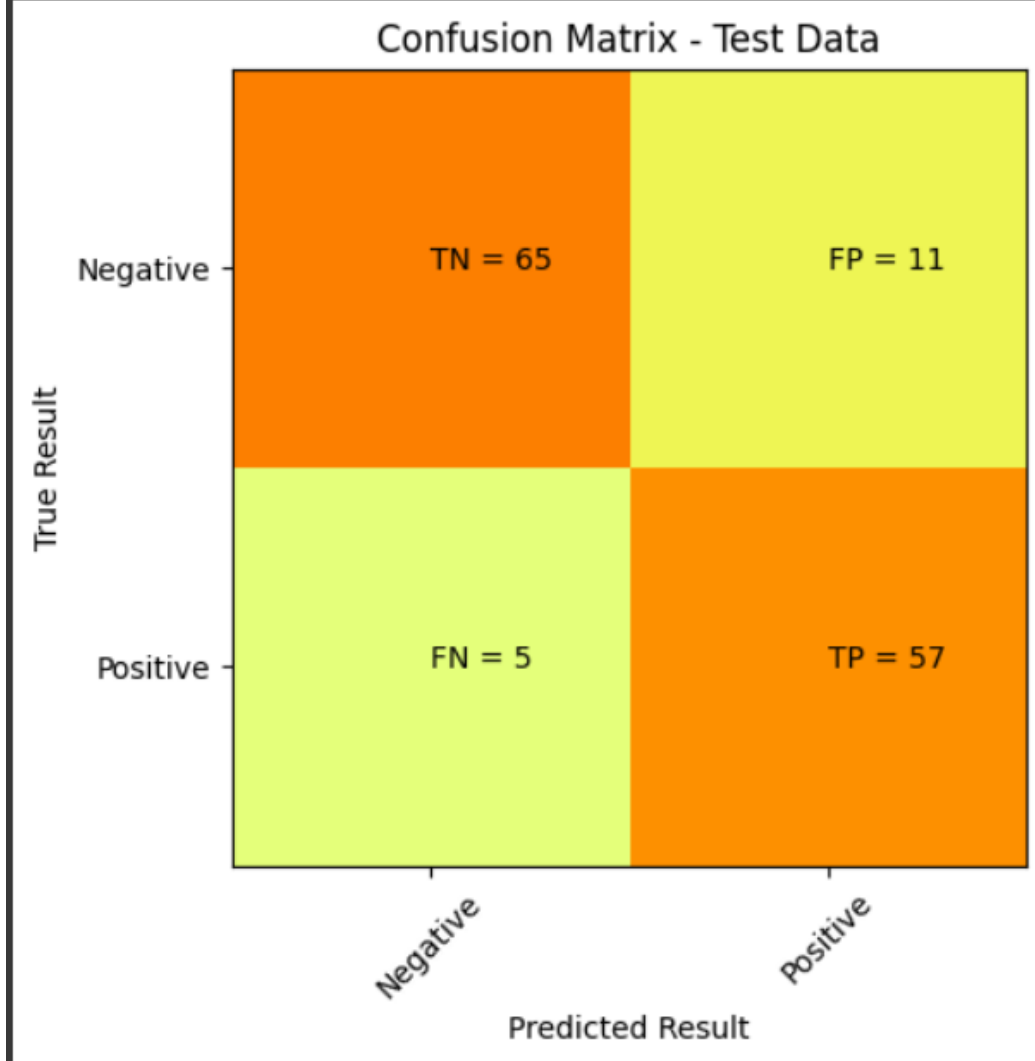
Confusion Matrix for SVM

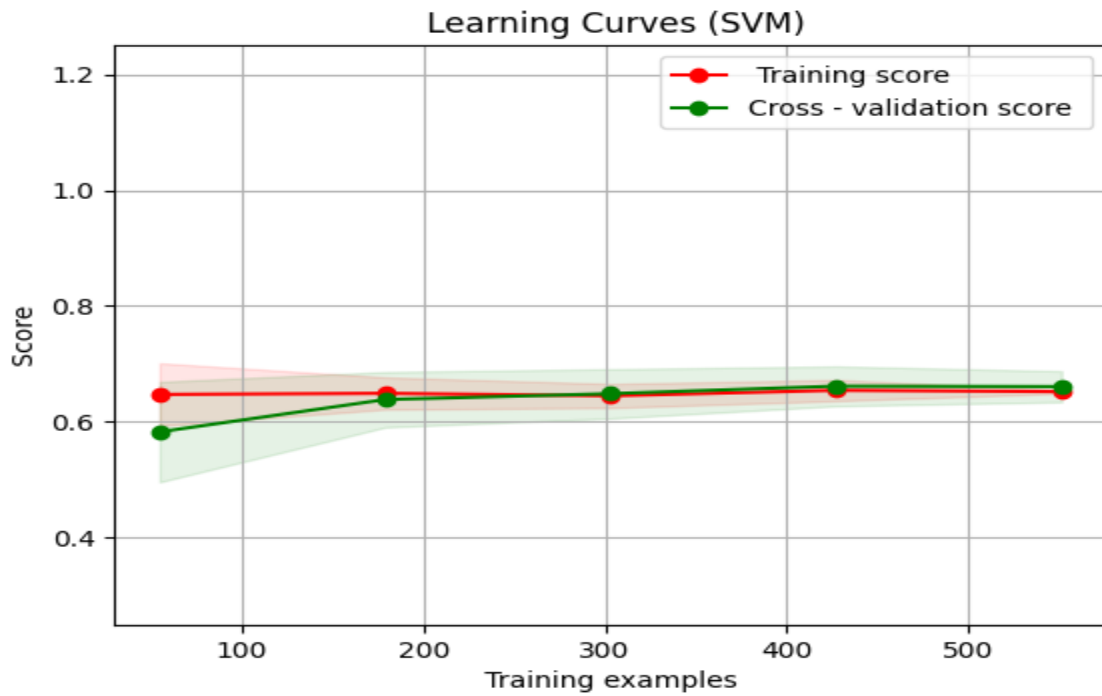


Mean accuracy and confusion matrix after tuning are as follows:-

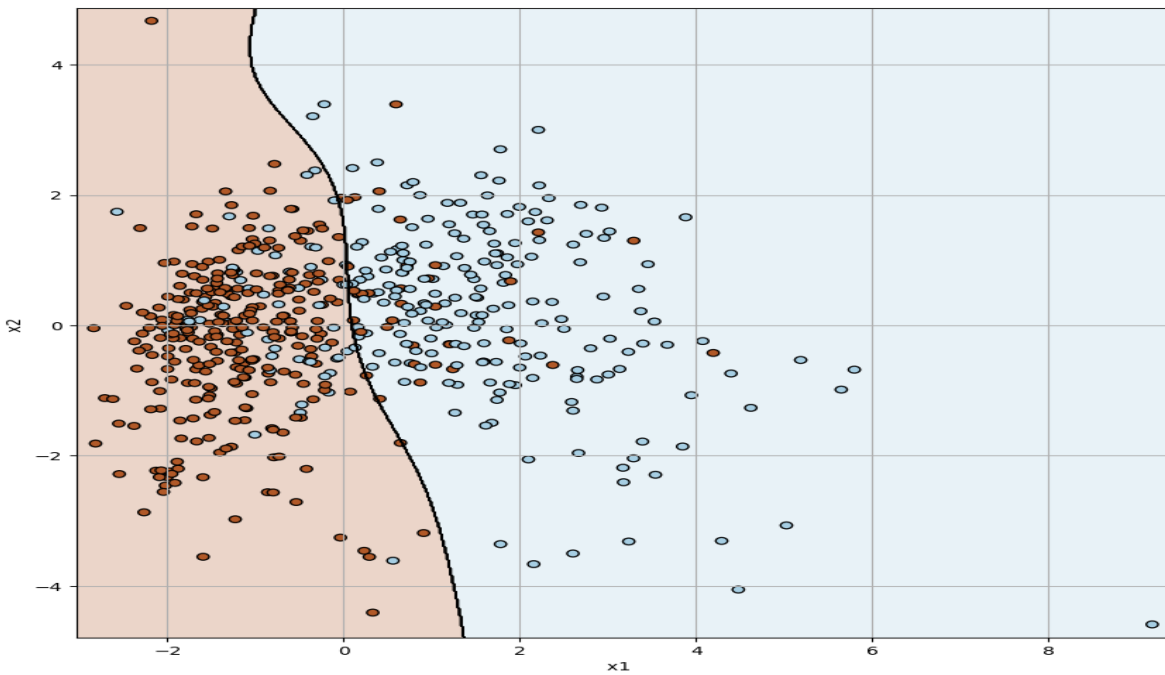
SVM : Mean Accuracy after tuning 85.30194805194805

Confusion Matrix for SVM





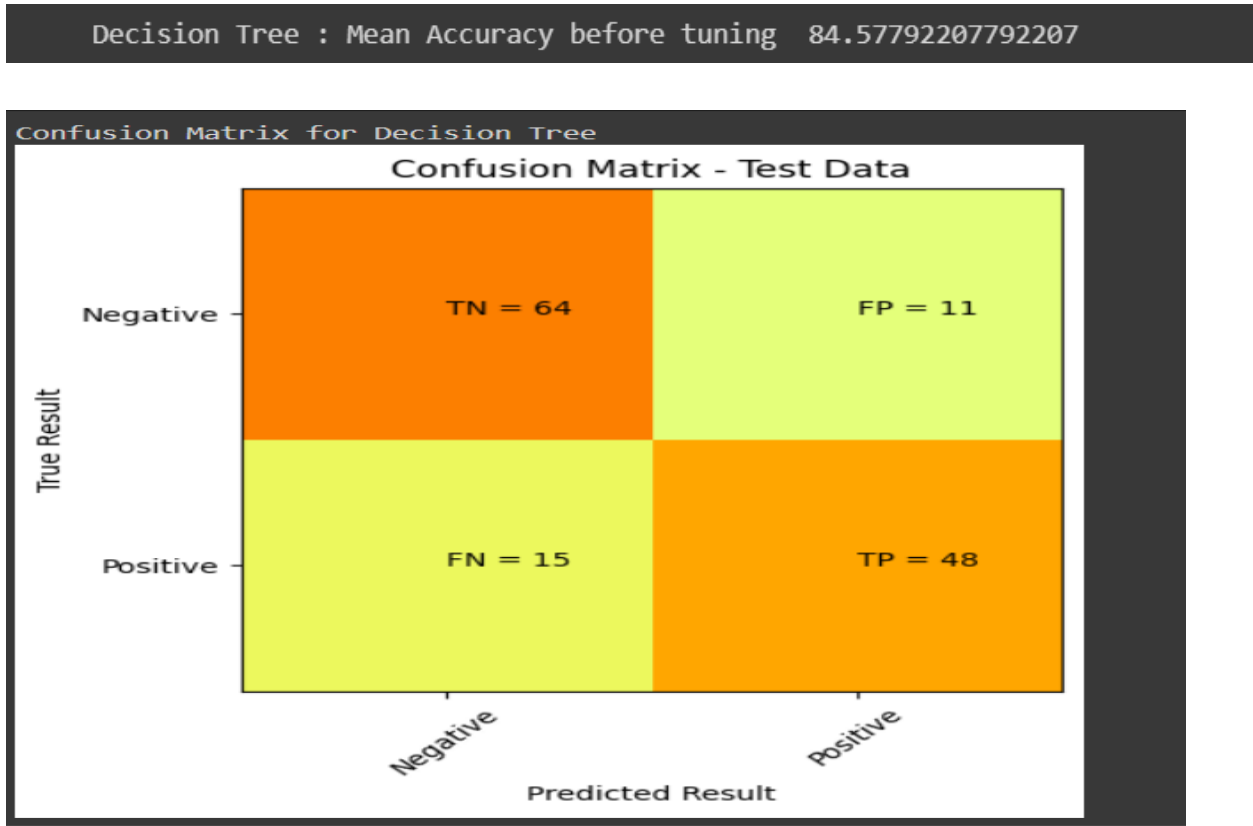
Learning curve of Support Vector Machine using the Australian data set. No gap is observed as increases the samples, indicating that the Support Vector Machine model is under-fitting this Australian data set.



Decision boundary plot using Support Vector Machine of Australian data set.

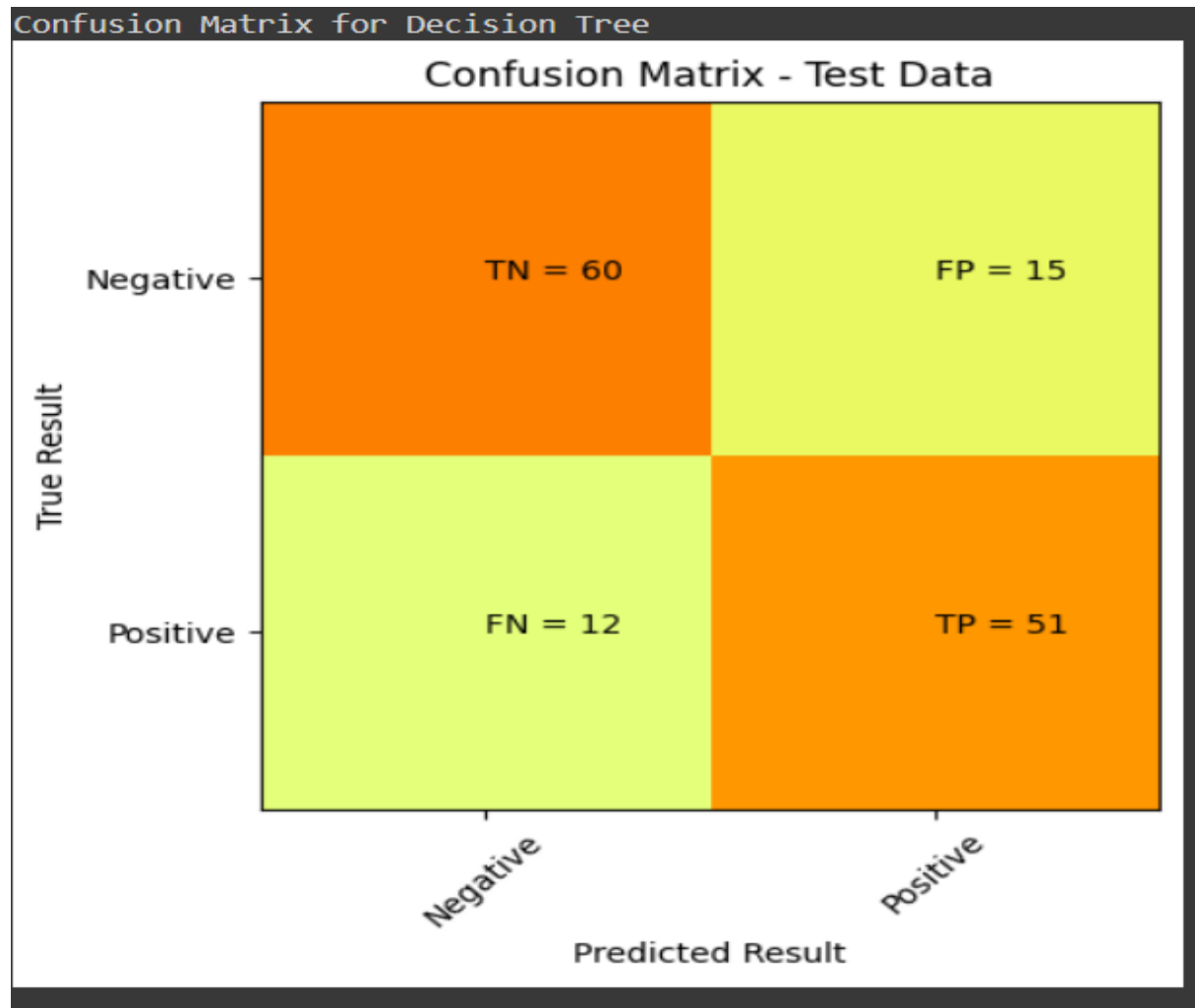
Results(both default and after tuning)-

Mean accuracy and confusion matrix without tuning:-

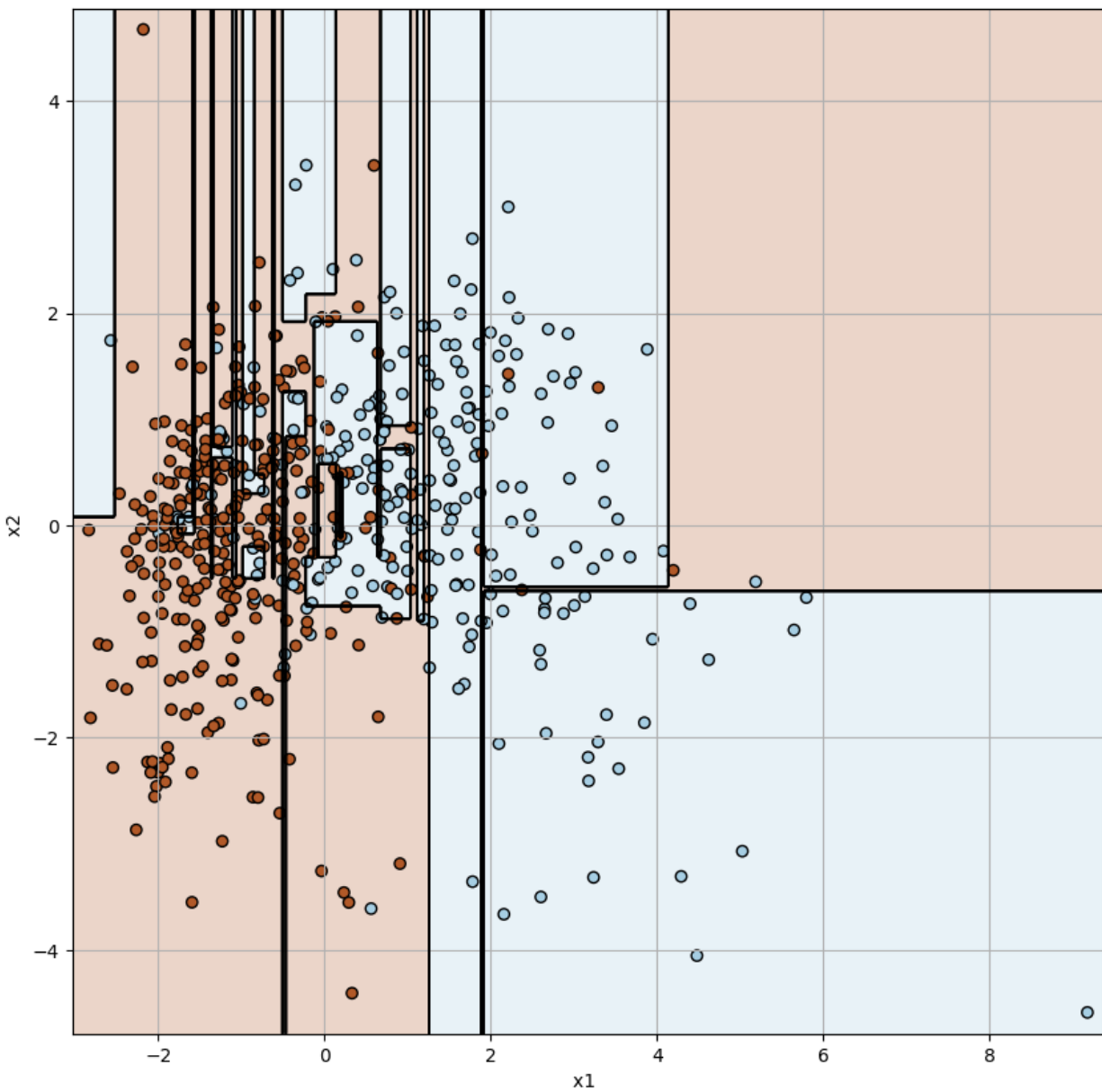


Mean accuracy and confusion matrix with tuning:-

Decision Tree : Mean Accuracy after tuning 85.47727272727272







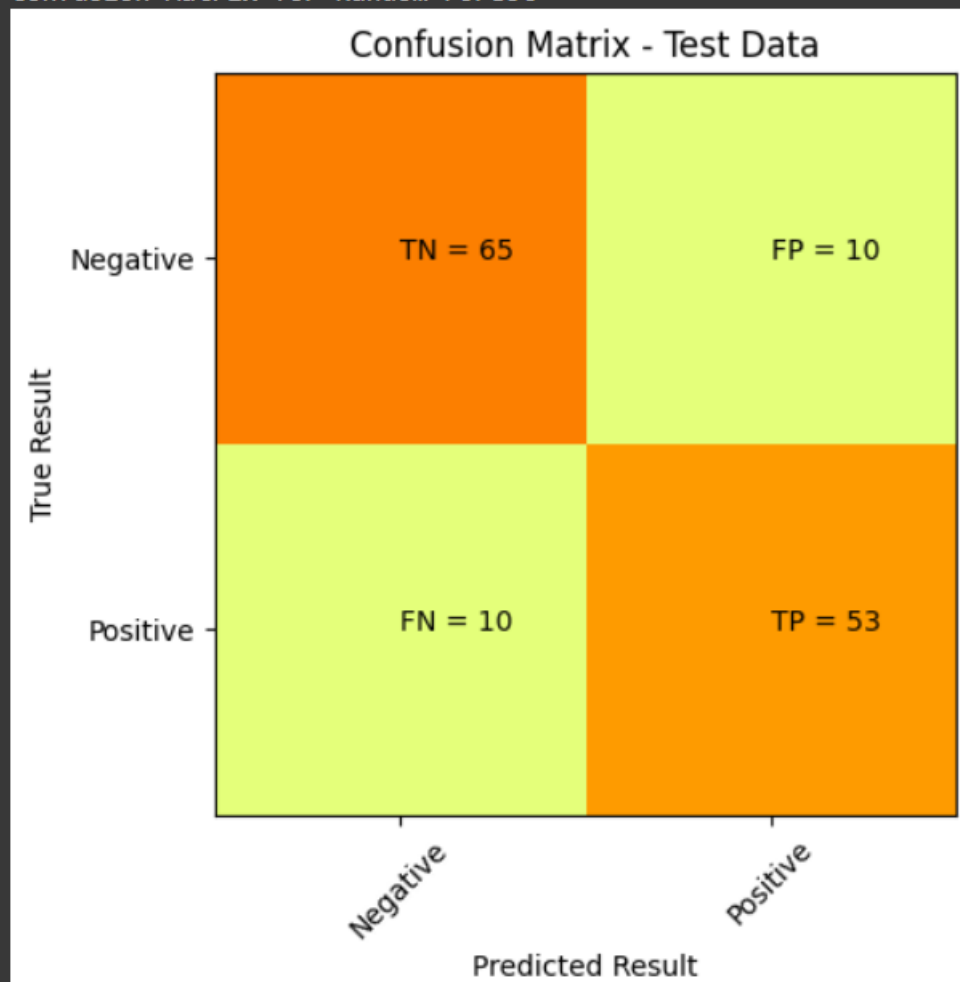
Decision boundary plot using Decision Tree of Australian data set.

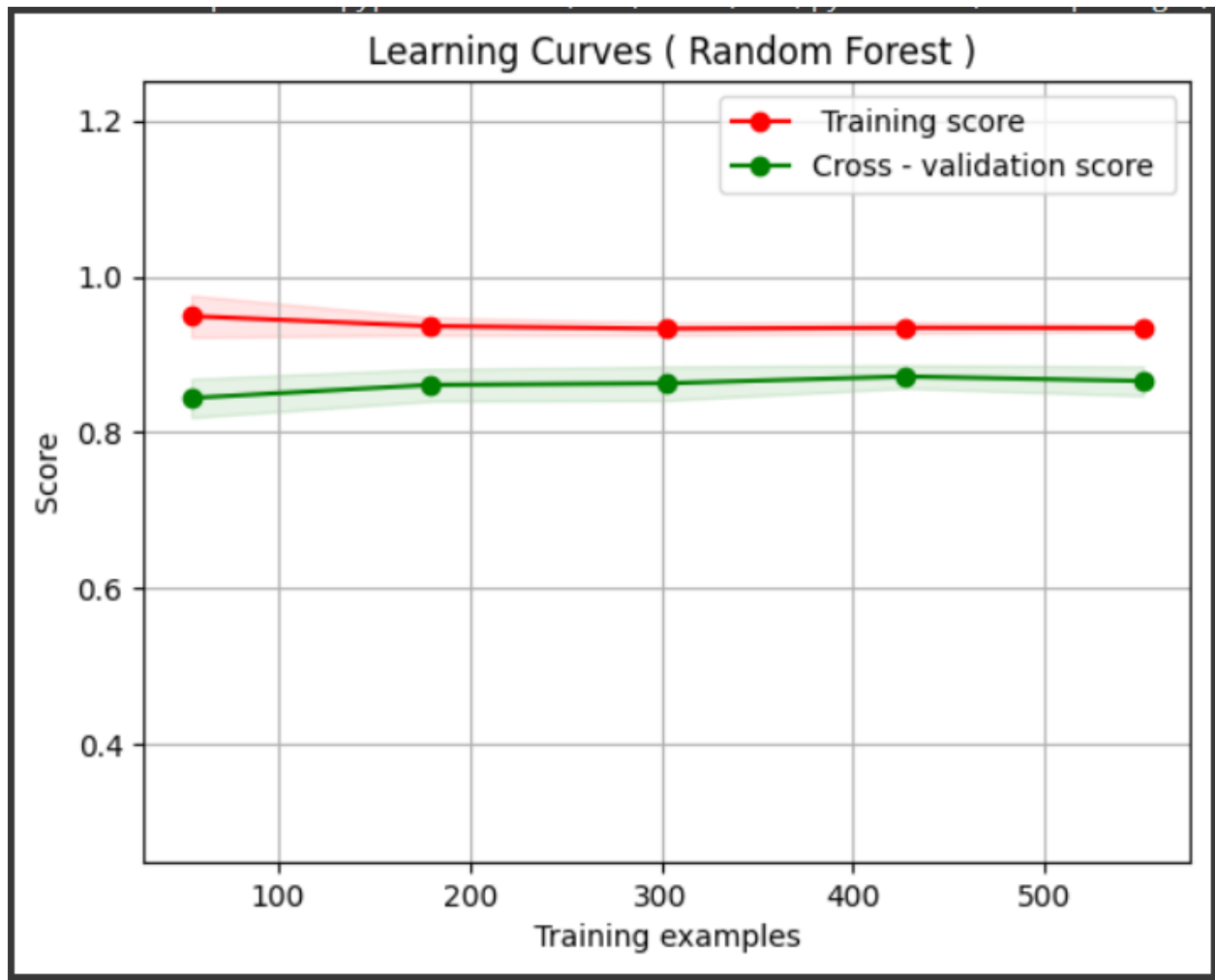
Results(both default and after tuning)-

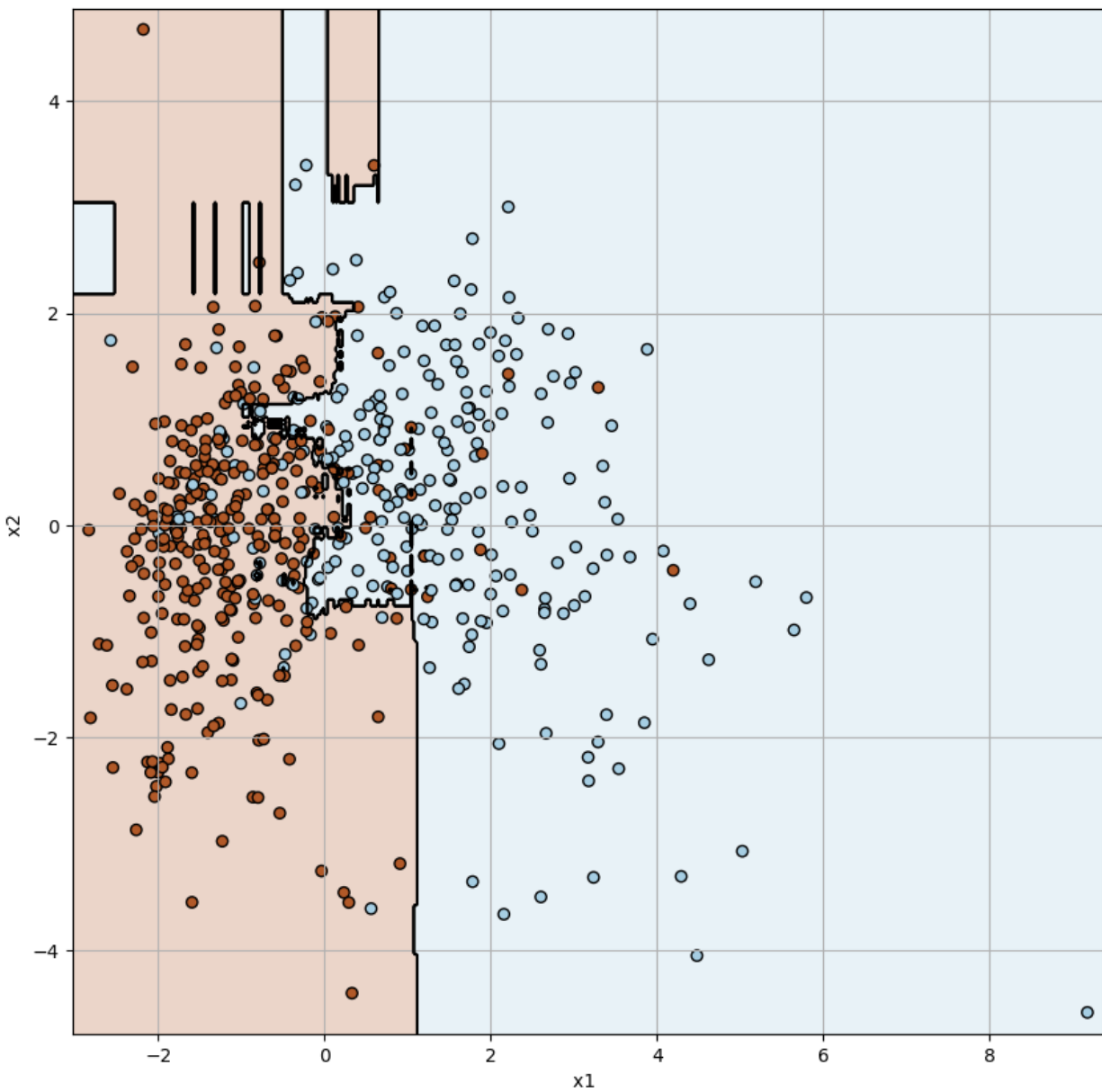
Mean accuracy and confusion matrix without tuning:-

```
Random Forest : Mean Accuracy before tuning 85.2987012987013
```

Confusion Matrix for Random Forest







After observing the results we come to the conclusion that Classification Trees Provide best accuracy to the given dataset and avoid overfitting which is later accompanied by Logistic Regression and atlast SVM comes into the picture.

By doing the Parameter Optimization we have increased the accuracy of SVM manifold but the accuracy of other models have changed a bit.