# ML Code with Feature Selection – Explanation SelectKBest

```
import pandas as pd
from sklearn.model_selection import train_test_split
import time
import numpy as np
from sklearn.preprocessing import StandardScaler
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import chi2
from sklearn.feature_selection import RFE
from sklearn.linear_model import LogisticRegression
import pickle
import matplotlib.pyplot as plt
```

#### 1. Importing required libraries

```
def selectkbest(indep_X,dep_Y,n):
        test = SelectKBest(score_func=chi2, k=n)
        fit1= test.fit(indep_X,dep_Y)
        selectk_features = fit1.transform(indep_X)
       return selectk_features
def split_scalar(indep_X,dep_Y):
        \textbf{X\_train, X\_test, y\_train, y\_test = train\_test\_split(indep\_X, dep\_Y, test\_size = 0.25, random\_state = 0) } 
        sc = StandardScaler()
        X train = sc.fit transform(X train)
       X_test = sc.transform(X_test)
       return X_train, X_test, y_train, y_test
def cm prediction(classifier.X test):
    y_pred = classifier.predict(X_test)
        # Making the Confusion Matrix
     from sklearn.metrics import confusion_matrix
     cm = confusion_matrix(y_test, y_pred)
    from sklearn.metrics import accuracy score
    from sklearn.metrics import classification_report
    Accuracy=accuracy_score(y_test, y_pred )
     report=classification report(y test, y pred)
     return classifier, Accuracy, report, X test, y test, cm
```

- 2. Created a function for Feature selection
- SelectKBest(score func=chi2, k=n) will pick the top n features using the Chi-Square test.
- fit1.fit(indep\_X,dep\_Y) will learn which features are important.
- fit1.transform(indep X) will keep only the best n features then return it.
- 3. Created a function for Split Scalar
- train\_test\_split(...) will Split the data into training (75%) and testing (25%).
- StandardScaler() scales features so they are on the same range.
- fit\_transform on training data learn scaling values (mean, std) and apply.
- transform on test data apply the same scaling.
- 4. Created a function for Prediction using confusion matrix
- Predict(X\_test) The model predict using test data.

- Confusion\_matrix (y\_test, y\_pred) will create the confusion matrix
- Accuracy Correct predictions saved in accuracy
- Classification\_report detailed report created for the predictions.

```
def logistic(X_train,y_train,X_test):
        from sklearn.linear model import LogisticRegression
        classifier = LogisticRegression(random_state = 0)
        classifier.fit(X_train, y_train)
        classifier,Accuracy,report,X_test,y_test,cm=cm_prediction(classifier,X_test)
        return classifier, Accuracy, report, X_test, y_test, cm
def svm_linear(X_train,y_train,X_test):
        from sklearn.svm import SVC
        classifier = SVC(kernel = 'linear', random_state = 0)
        classifier.fit(X_train, y_train)
        classifier,Accuracy,report,X_test,y_test,cm=cm_prediction(classifier,X_test)
        return classifier,Accuracy,report,X_test,y_test,cm
def svm_NL(X_train,y_train,X_test):
        from sklearn.svm import SVC
        classifier = SVC(kernel = 'rbf', random_state = 0)
        classifier.fit(X_train, y_train)
        \verb|classifier,Accuracy,report,X_test,y_test,cm=cm\_prediction(classifier,X_test)|\\
        return classifier,Accuracy,report,X_test,y_test,cm
def Navie(X_train,y_train,X_test):
        # Fitting K-NN to the Training set
        from sklearn.naive_bayes import GaussianNB
        classifier = GaussianNB()
        classifier.fit(X_train, y_train)
        {\tt classifier,} {\tt Accuracy,report,} {\tt X\_test,y\_test,cm=cm\_prediction} ({\tt classifier,X\_test})
        return classifier,Accuracy,report,X_test,y_test,cm
```

- 5. Created a function for logistic regression Algorithm
- LogisticRegression(random\_state = 0) creates an instance of the logistic regression classifier
- classifier.fit(X\_train, y\_train)- Trains the logistic regression model with training data
- cm\_prediction(classifier,X\_test) calls prediction function and make prediction using test data and returns accuracy, report and confusion matrix
- 6. Created a function for SVM Linear Algorithm
- SVC(...) creates SVMLinear classifier with a linear kernel, meaning it tries to find a straight line
- classifier.fit(X\_train, y\_train)- Trains the model with training data
- cm\_prediction(classifier,X\_test) calls prediction function and make prediction using test data and returns accuracy, report and confusion matrix
- 7. Created a function for SVM Non Linear Algorithm
- SVC(...) creates SVM classifier using the RBF kernel. The RBF kernel maps data into a higher-dimensional space to handle cases where data is not linear
- classifier.fit(X\_train, y\_train)- Trains the model with training data
- cm\_prediction(classifier,X\_test) calls prediction function and make prediction using test data and returns accuracy, report and confusion matrix
- 8. Created a function for Navies' bayes Algorithm
- GaussianNB()— Creates an instance of the Gaussian Naive Bayes classifier.
- classifier.fit(X train, y train)- Trains the model with training data
- cm\_prediction(classifier,X\_test) calls prediction function and make prediction using test data and returns accuracy, report and confusion matrix

```
def knn(X_train,y_train,X_test):
         # Fitting K-NN to the Training set
         from sklearn.neighbors import KNeighborsClassifier
         classifier = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski', p = 2)
         classifier.fit(X_train, y_train)
         \verb|classifier,Accuracy,report,X_test,y_test,cm=cm_prediction(classifier,X_test)|\\
         return classifier,Accuracy,report,X_test,y_test,cm
\label{eq:def_Decision} \textbf{def} \ \ \mathsf{Decision}(\mathsf{X}\_\mathsf{train}, \mathsf{y}\_\mathsf{train}, \mathsf{X}\_\mathsf{test}) :
         # Fitting K-NN to the Training set
         from sklearn.tree import DecisionTreeClassifier
         classifier = DecisionTreeClassifier(criterion = 'entropy', random_state = 0)
         classifier.fit(X\_train, y\_train)
         {\tt classifier,} Accuracy, {\tt report,} X\_{\tt test,} y\_{\tt test,} {\tt cm=cm\_prediction} ({\tt classifier,} X\_{\tt test})
         return classifier, Accuracy, report, X test, y test, cm
def random(X_train,y_train,X_test):
         # Fitting K-NN to the Training set
         from sklearn.ensemble import RandomForestClassifier
         classifier = RandomForestClassifier (n\_estimators = 10, criterion = 'entropy', random\_state = 0)
         classifier.fit(X_train, y_train)
         classifier,Accuracy,report,X_test,y_test,cm=cm_prediction(classifier,X_test)
         return classifier,Accuracy,report,X test,y test,cm
```

- 9. Created a function for K-Nearest Neighbors Algorithm
- KNeighborsClassifier(n\_neighbors=5, metric='minkowski', p=2)
  - a. n\_neighbors=5 means the algorithm will look at the 5 closest neighbors to classify a new point.
  - b. metric='minkowski' with p=2 means it uses Euclidean distance to measure closeness between points.
- classifier.fit(X train, y train)- Trains the model with training data
- cm\_prediction(classifier,X\_test) calls prediction function and make prediction using test data and returns accuracy, report and confusion matrix

#### 10. Created a function for Decision Tree Algorithm

- DecisionTreeClassifier(criterion = 'entropy', random\_state = 0) criterion='entropy' tells the tree to use information gain / entropy to decide where to split nodes
- classifier.fit(X train, y train)- Trains the model with training data
- cm\_prediction(classifier,X\_test) calls prediction function and make prediction using test data and returns accuracy, report and confusion matrix

# 11. Created a function for Random Forest Algorithm

- RandomForestClassifier(n\_estimators = 10, criterion = 'entropy', random\_state = 0)
  - a. n\_estimators = 10,-- random forest will create **10 decision trees**, More trees better accuracy
  - b. criterion = 'entropy' -- Each tree will split its nodes using the **entropy (information gain)** method. Measures how pure the data is at each split
- classifier.fit(X train, y train)- Trains the model with training data
- cm\_prediction(classifier,X\_test) calls prediction function and make prediction using test data and returns accuracy, report and confusion matrix

```
def selectk_Classification(acclog,accsvml,accsvmnl,accknn,accnav,accdes,accrf):

dataframe=pd.DataFrame(index=['ChiSquare'],columns=['Logistic','SVMl','SVMnl','KNN','Navie','Decision','Random'])
for number,idex in enumerate(dataframe.index):
    dataframe['togistic'][idex]=acclog[number]
    dataframe['SVMl'][idex]=accsvml[number]
    dataframe['SVMnl'][idex]=accsvml[number]
    dataframe['KNN'][idex]=accknn[number]
    dataframe['Mavie'][idex]=accknn[number]
    dataframe['Decision'][idex]=accdes[number]
    dataframe['Random'][idex]=accrf[number]
    return dataframe
```

#### 12. Created a function for SelectK

- o It accepts inputs that hold accuracy values for different classifiers
- pd.DataFrame(index=['ChiSquare'], columns=['Logistic...] Creates an empty dataframe with one row index - 'ChiSquare' & several col
- o enumerate(dataframe.index): is ['ChiSquare'] --- so this loop runs once
- o enumerate(...) gives two values each iteration:
  - number → the integer position (0 for the first row)
  - index → the actual index label (the string 'ChiSquare')
- o dataframe['Logistic'][index] = acclog[number] Sets the cell at row 'ChiSquare' and column 'Logistic' to the value acclog[0].
- Acclog[0] accuracy value for logistic is set to chiSquare row and logistic col of the dataframe.
- o Like above all the accuracy values will be placed to specific columns.

```
dataset1=pd.read_csv("prep.csv",index_col=None)

df2=dataset1

df2 = pd.get_dummies(df2, drop_first=True)

indep_X=df2.drop('classification_yes', 1)
dep_Y=df2['classification_yes']
```

### 13. Reading Dataset

- dataset1 = pd.read\_csv("prep.csv", index\_col=None) reading the csv file and index\_col=None → means don't treat any column as the row index; just use default numeric indexes (0, 1, 2, ...).
- df2 = pd.get\_dummies(df2, drop\_first=True) to convert the categorical data to numerical we are using get\_dummies
- indep\_X input col in df without classification\_yes col
- dep\_Y output is classification\_yes col

```
kbest=selectkbest(indep_X,dep_Y,6)

acclog=[]
accsvml=[]
accsvmnl=[]
accknn=[]
accnav=[]
accdes=[]
accdes=[]
```

# 14. Calling a function selectkbest

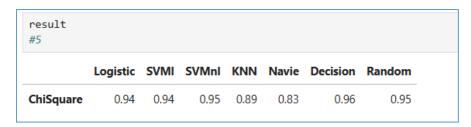
- Parameters are input ,output and n value (6 most important features)
- Empty Python lists created to store accuracy values for different classifiers

```
X_train, X_test, y_train, y_test=split_scalar(kbest,dep_Y)
classifier,Accuracy,report,X_test,y_test,cm=logistic(X_train,y_train,X_test)
acclog.append(Accuracy)
{\tt classifier,Accuracy,report,X\_test,y\_test,cm=svm\_linear(X\_train,y\_train,X\_test)}
accsvml.append(Accuracy)
classifier,Accuracy,report,X_test,y_test,cm=svm_NL(X_train,y_train,X_test)
accsvmnl.append(Accuracy)
classifier,Accuracy,report,X_test,y_test,cm=knn(X_train,y_train,X_test)
accknn.append(Accuracy)
classifier,Accuracy,report,X test,y test,cm=Navie(X train,y train,X test)
accnav.append(Accuracy)
classifier,Accuracy,report,X_test,y_test,cm=Decision(X_train,y_train,X_test)
accdes.append(Accuracy)
classifier,Accuracy,report,X_test,y_test,cm=random(X_train,y_train,X_test)
accrf.append(Accuracy)
result=selectk_Classification(acclog,accsvml,accsvmnl,accknn,accnav,accdes,accrf)
```

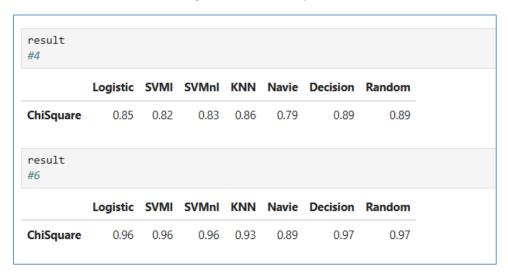
# 15. Training and Evaluation

- X\_train, X\_test, y\_train, y\_test = split\_scalar(kbest, dep\_Y) split the dataset to training and test - kbest(input var selected features)
- classifier, Accuracy, report, X\_test, y\_test, cm = logistic(X\_train, y\_train, X\_test)
- acclog.append(Accuracy)
  - o Trains a **Logistic Regression** model.
  - o Returns:
  - o classifier → trained model
  - O Accuracy → test accuracy
  - report → classification report
  - $\circ$  cm  $\rightarrow$  confusion matrix

- Appends the accuracy to the list acclog.
- Like above all the model will be trained and tested and the accuracy value is appended to the list.
- result = selectk\_Classification(acclog, accsvml, accsvml, accknn, accnav, accdes, accrf) will create the dataframe using the accuracy values .



Printing the result for 5 features. This way we can check by different n features and select the best n number of features which will give more accuracy.



# **Summary:**

- ❖ Model performance improves as we increase the number of selected features from 4 to 6.
- ❖ Decision Tree and Random Forest has the highest accuracy (0.97) with 6 features, making this the best configuration for classification.
- ❖ Naive Bayes performs the weakest across all feature sets.