# Bankruptcy prediction

## of Polish companies

KV Varun S20190020222
Vineesh P S20190020243
Amruth P S20190020242

# 1  Data-set description

The data-set is about bankruptcy prediction of Polish companies.The bankrupt companies were analyzed in the period 2000-2012, while the still operating companies were evaluated from 2007 to 2013.

It contains financial rates from $5^{th}$ year of the forecasting period and corresponding class label that indicates bankruptcy status after 1 year. The data contains 5910 instances (financial statements), 410 represents bankrupted companies, 5500 firms that did not bankrupt in the forecasting period.

Each instance/company has 64 attributes on which the classification is evaluated. This makes the data-set multivariate.

**Class labels**
We have two class labels representing the bankruptcy status of respective company at the end of forecasting period.

- 1 - Indicates that the company went bankrupt.

- 0 - Indicates that the company didn't went bankrupt.

**Attributes (X)**
Few of the 64 attributes are:

1. X1 net profit / total assets

2. X3 working capital / total assets

3. X19 gross profit / sales

**Data-set repository**
https://archive.ics.uci.edu/ml/datasets/Polish+companies+bankruptcy+data

# 2  Problem statement

- Our objective is to classify (binary classification) the given data and to predict corresponding labels with good accuracy for test data containing bank's details.

- Use feature extraction methods to increase accuracy of the model.

# 3  Methodology

1. Importing necessary libraries

2. Importing bankruptcy data-set

3. Processing missing values

4. Splitting data-set into training and test sets

5. Training the logistic regression model, Linear Support Vector Machine, Polynomial Support Vector Machine, Radial basis function SVM, K-NN model

6. Predicting the target labels of test-set

7. Performing Feature extraction

8. Comparing results after and before feature extraction.

**Python code implement the above methodology**
https://github.com/kvv1618/Bankruptcy-prediction

# 4  Explanation

1. Step 1

**Importing necessary libraries**

```
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

2. Step 2

**Importing bankruptcy dataset**

```python
from scipy.io import arff
from io import BytesIO
data = arff.loadarff('5year.arff')
df = pd.DataFrame(data[0])
X = df.iloc[:, :-1].values
Y = df.iloc[:, -1].values
y=[]
for i in range(len(Y)):
    if(Y[i]==b'0'):
        y.append(0)
    else:
        y.append(1)
```

3. Step 3

**Processing missing values**

```python
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(missing_values=np.nan, strategy='mean')
imputer.fit(X)
X = imputer.transform(X)
```
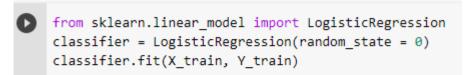
4. Step 4

**Splitting dataset into training and test sets**

```python
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test = train_test_split(X, y,
```

```python
on import train_test_split
Y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
```

5. Logistic regression before feature extraction

## Training the logistic regression model

```
from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(random_state = 0)
classifier.fit(X_train, Y_train)
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/linear
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale
    https://scikit-learn.org/stable/modules/preproces
Please also refer to the documentation for alternativ
    https://scikit-learn.org/stable/modules/linear_mo
  extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG,
LogisticRegression(random_state=0)
```

## Predicting the target labels of testset

```
[6]  from sklearn.metrics import confusion_matrix, accuracy_score
     Y_pred = classifier.predict(X_test)
     cm = confusion_matrix(Y_test, Y_pred)
     print(cm)
     accuracy_score(Y_test, Y_pred)
```

```
[[1382    2]
 [  93    1]]
0.935723951285521
```

6. Feature extraction

# Feacture extraction

```
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
lda = LDA(n_components = 1)
print(len(X_train[0]))
X_train = lda.fit_transform(X_train, Y_train)
X_test = lda.transform(X_test)
print(len(X_train[0]))
```

```
64
1
```

7. Logistic regression after feature extraction

**Training the logistic regression model after feature extraction**

```
[8]  from sklearn.linear_model import LogisticRegression
     classifier = LogisticRegression(random_state = 0)
     classifier.fit(X_train, Y_train)
```

```
LogisticRegression(random_state=0)
```

Predicting the target lables of testset

```
from sklearn.metrics import confusion_matrix, accuracy_score
Y_pred = classifier.predict(X_test)
cm = confusion_matrix(Y_test, Y_pred)
print(cm)
accuracy_score(Y_test, Y_pred)
```

```
[[1378    6]
 [  85    9]]
0.9384303112313938
```

8. Step 8

# ▾ Linear Support Vector Machine

```python
from sklearn.svm import SVC
classifier = SVC(kernel = 'linear', random_state = 0)
classifier.fit(X_train, Y_train)
```

```
SVC(kernel='linear', random_state=0)
```

```python
from sklearn.metrics import confusion_matrix, accuracy_score
Y_pred = classifier.predict(X_test)
cm = confusion_matrix(Y_test, Y_pred)
print(cm)
accuracy_score(Y_test, Y_pred)
```

```
[[1384    0]
 [  94    0]]
0.9364005412719891
```

9. Step 9

# Polynomial Support Vector Machine

```
[12] from sklearn.svm import SVC
     classifier = SVC(kernel = 'poly',degree=3, random_state = 0)
     classifier.fit(X_train, Y_train)

     SVC(kernel='poly', random_state=0)
```

```
    from sklearn.metrics import confusion_matrix, accuracy_score
    Y_pred = classifier.predict(X_test)
    cm = confusion_matrix(Y_test, Y_pred)
    print(cm)
    accuracy_score(Y_test, Y_pred)
```

```
[[1382    2]
 [  93    1]]
0.935723951285521
```

10. Step 10

7

# Radial basis function SVM

```
[14] from sklearn.svm import SVC
     classifier = SVC(kernel = 'rbf', random_state = 0)
     classifier.fit(X_train, Y_train)

     SVC(random_state=0)
```

```
    from sklearn.metrics import confusion_matrix, accuracy_score
    Y_pred = classifier.predict(X_test)
    cm = confusion_matrix(Y_test, Y_pred)
    print(cm)
    accuracy_score(Y_test, Y_pred)
```

```
[[1380    4]
 [  92    2]]
0.9350473612990527
```

11. Step 11

# K-NN model

```
[16] from sklearn.neighbors import KNeighborsClassifier
     classifier = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski',
     classifier.fit(X_train, Y_train)

     KNeighborsClassifier()
```

```
[17] from sklearn.metrics import confusion_matrix, accuracy_score
     Y_pred = classifier.predict(X_test)
     cm = confusion_matrix(Y_test, Y_pred)
     print(cm)
     accuracy_score(Y_test, Y_pred)

     [[1361    23]
      [  87     7]]
     0.925575101488498
```

## 5  Observations

Logistic regression is a supervised learning classification algorithm used to predict the probability of a target variable. We have initially trained the data-set on this model before applying any feature extraction and we got a accuracy score of 93.5723951285521 percent.

LDA is supervised learning dimensionality reduction technique and aims to maximize the distance between the mean of each class and minimize the spreading within the class itself. This is applied on the training data to reduce the dimensionality of the matrix. Since the total number of objective classes are 2, we reduce the training set to a single column for each instance.

After feature extraction results of logistic regression have improved in terms of accuracy. The new accuracy was 93.84303112313938 percent.

## 6  Contributions

- Amruth P
  Selecting data-set from UCI machine learning Repository
  Extraction of required data
  Importing and prepossessing the data.

- Vineesh P
  Training Linear regression classification model on the data.
  Training Linear and Polynomial support vector machine on the data after feature extraction.
  Predicting classification labels for the above mentioned models.

- Varun KV
  Implementing feature extraction.
  Training Radial basis support vector machine and K-NN models on the data after feature extraction.
  Predicting classification labels for the above mentioned models.