

Cancer Prediction using Binary Classification

Omar Kashmar

M.Tech Aerospace Engineering

IIT Bombay

213011004

Self-Project



Abstract:

This report presents the results of a project focused on predicting the occurrence of cancer using binary classification techniques. The primary objective was to develop accurate models that could classify patients into two classes: those with cancer and those without cancer. The study employed various classification algorithms, including Support Vector Machines (SVM), Random Forest, and XGBoost, to achieve this goal. The report discusses the dataset, methodology, model selection, evaluation metrics, and the results obtained.

Introduction

Cancer is a complex and devastating disease that requires early detection and diagnosis for effective treatment. Predictive modeling using machine learning techniques offers the potential to assist medical professionals in identifying potential cases of cancer, enabling earlier interventions and improving patient outcomes. This project aims to explore the applicability of binary classification algorithms to predict the presence of cancer based on relevant medical features.

Dataset

The dataset used for this project consists of medical data collected from patients, including various features such as area_mean, perimeter_mean, concavity_worst status, and more. The target variable indicates whether a patient has cancer (class 1) or does not have cancer (class 0).

```
Project to Explain Classification Last checkpoint: 2 hours ago (autosaved)
Edit View Insert Cell Kernel Widgets Help Not Connected Not Trusted Python 3 (ipykernel)
In [2]: os.chdir('C:\\Noble\\Training\\Top Mentor\\Training\\Presentation\\Project\\Project to Explain Classification\\')
data = pd.read_csv('data.csv')
data

Out[2]:
```

| | id | diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean | smoothness_mean | compactness_mean | concavity_mean | concav points_mean |
|-----|----------|-----------|-------------|--------------|----------------|-----------|-----------------|------------------|----------------|--------------------|
| 0 | 842302 | M | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.30010 | 0.1471 |
| 1 | 842517 | M | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.08690 | 0.0701 |
| 2 | 84300903 | M | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.19740 | 0.1279 |
| 3 | 84348301 | M | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.24140 | 0.1052 |
| 4 | 84358402 | M | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.19800 | 0.1043 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 564 | 926424 | M | 21.56 | 22.39 | 142.00 | 1479.0 | 0.11100 | 0.11590 | 0.24390 | 0.1389 |
| 565 | 926682 | M | 20.13 | 28.25 | 131.20 | 1261.0 | 0.09780 | 0.10340 | 0.14400 | 0.0979 |
| 566 | 926954 | M | 16.60 | 28.08 | 108.30 | 858.1 | 0.08455 | 0.10230 | 0.09251 | 0.0530 |
| 567 | 927241 | M | 20.60 | 29.33 | 140.10 | 1265.0 | 0.11780 | 0.27700 | 0.35140 | 0.1520 |
| 568 | 92751 | B | 7.76 | 24.54 | 47.92 | 181.0 | 0.05263 | 0.04362 | 0.00000 | 0.0000 |

569 rows x 32 columns

Methodology

- Data Preprocessing:** The dataset underwent preprocessing steps, including handling missing values, encoding categorical variables, and scaling numeric features.
- Data Splitting:** The dataset was split into training and testing sets to ensure unbiased evaluation of model performance.
- Resampling:** Given the imbalanced nature of the dataset, resampling techniques, specifically oversampling, were applied to mitigate class imbalance and improve model performance.
- Model Selection:** Three classification algorithms were chosen for this project: Support Vector Machines (SVM), Random Forest, and XGBoost. Hyperparameter tuning was performed using grid search and cross-validation.

- Model Evaluation:** The models were evaluated using various metrics, including accuracy, precision, recall, F1-score, and confusion matrices. Stratified cross-validation was employed to ensure robust evaluation.

Results:

The heat map that combine the relation between the variables:

Normalization

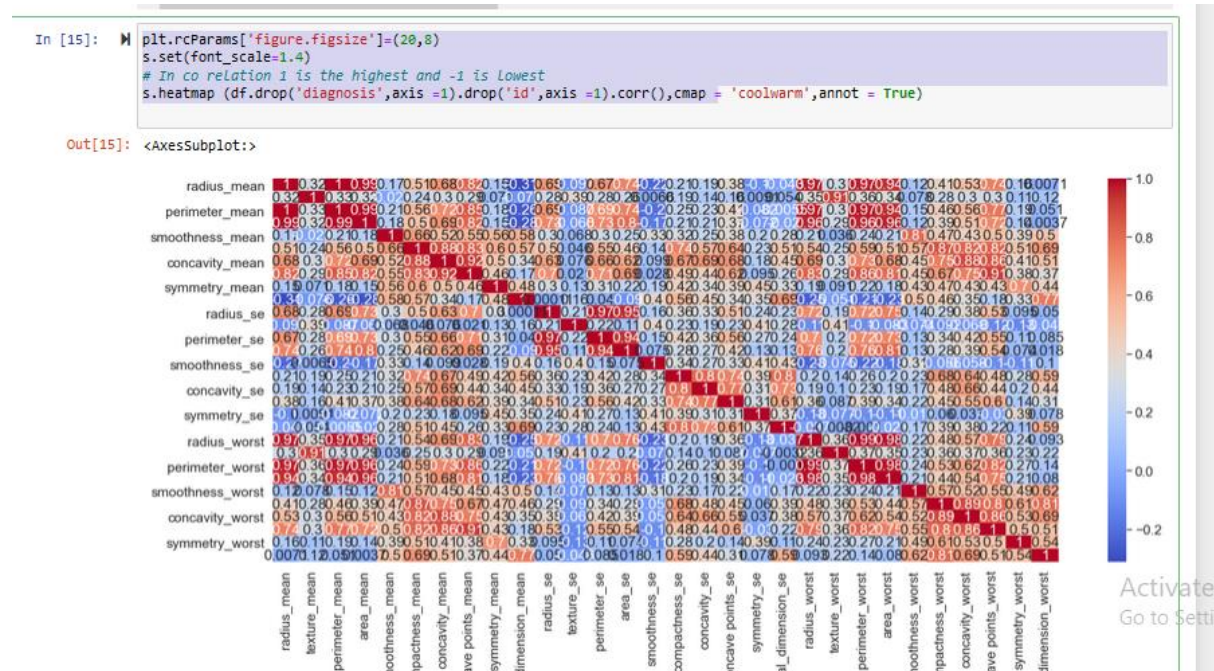
```
df_norm = (x- x.mean()) / (x.max()- x.min())
```

```
df_norm= pd.concat ([df_norm,y], axis =1 )
```

```
df_norm
```

Print Co relation

```
df.drop('diagnosis',axis =1).drop('id',axis =1).corr()
```



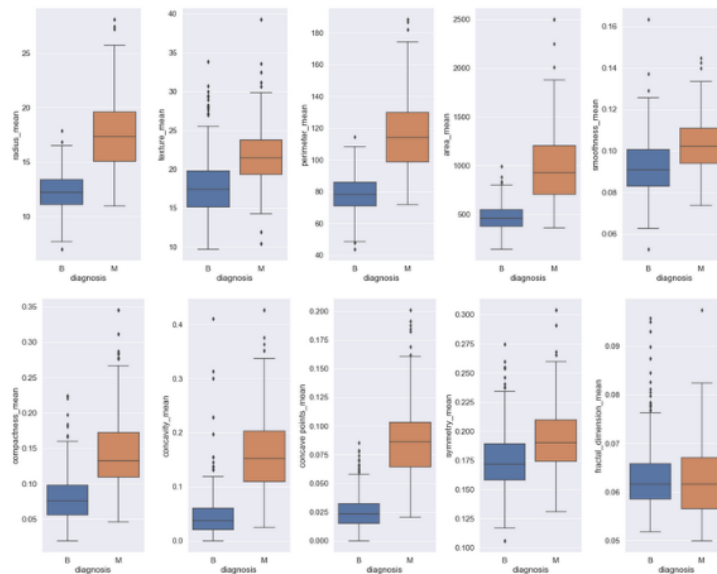
Create Box Plot

#box plot to check the outliers. Not going to remove outliers since this data is important.

```
# Observation, when 'diagnosis' is "B", the values are lower
plt.rcParams['figure.figsize']=(20,8)
f, (ax1,ax2,ax3,ax4,ax5) = plt.subplots (1,5)
s.boxplot ('diagnosis', y = 'radius_mean',data = df , ax = ax1)
s.boxplot ('diagnosis', y = 'texture_mean',data = df , ax = ax2)
s.boxplot ('diagnosis', y = 'perimeter_mean',data = df , ax = ax3)
s.boxplot ('diagnosis', y = 'area_mean',data = df , ax = ax4)
s.boxplot ('diagnosis', y = 'smoothness_mean',data = df , ax = ax5)
f.tight_layout()
```

```
f, (ax1,ax2,ax3,ax4,ax5) = plt.subplots (1,5)
s.boxplot ('diagnosis', y = 'compactness_mean',data = df , ax = ax1)
s.boxplot ('diagnosis', y = 'concavity_mean',data = df , ax = ax2)
s.boxplot ('diagnosis', y = 'concave points_mean',data = df , ax = ax3)
s.boxplot ('diagnosis', y = 'symmetry_mean',data = df , ax = ax4)
s.boxplot ('diagnosis', y = 'fractal_dimension_mean',data = df , ax = ax5)
f.tight_layout()
```

```
s.boxplot('diagnosis', y = 'symmetry_mean', data = df, ax = ax4)
s.boxplot('diagnosis', y = 'fractal_dimension_mean', data = df, ax = ax5)
f.tight_layout()
```



Activate Windows
Go to Settings to activate

Create SVM Model

#cv = cross validation

```
param = {
    'C': [0.1,1,100,1000],
    'gamma':[0.0001,0.001, 0.005, 0.1,1, 3,5,10, 100]
}
```

FitModel (x_norm,y_norm,'SVC',SVC(), param, cv =10)

Create Random Forest

```
param = { 'n_estimators': [100,500,1000,2000] }
```

FitModel (x_norm,y_norm,'Random Forest',RandomForestClassifier(), param, cv =10)

Create Random Forest Normal Way

```
np.random.seed(10)
```

```

x_train,x_test, y_train,y_test = train_test_split(x_norm,y_norm,test_size = 0.2)
forest = RandomForestClassifier (n_estimators = 500)
fit = forest.fit (x_train, y_train)
accuracy = fit.score(x_test,y_test)
predict = fit.predict(x_test)
cmatrix = confusion_matrix (y_test, predict)
print ('Classification Report:',classification_report(y_test,predict))
print ('Accuracy Score', (accuracy_score(y_test,predict)))
print ('Accuracy of Random Forest ', (accuracy))
print ('Confusion Matrix :\n',cmatrix)

```

XG Boost

```

param = { 'n_estimators': [100,500,1000,2000] }
FitModel (x_norm,y_norm,'XGBoost', XGBClassifier(),param, cv = 10)

```

Final Results:

- **Random Forest**

```

In [39]: X_res = pd.DataFrame(X_res)
Y_res = pd.DataFrame(Y_res)
X_res.columns = x_norm.columns
param = { 'n_estimators': [100,500,1000,2000] }
FitModel (X_res [feat_to_keep], Y_res , 'Random Forest', RandomForestClassifier(), param, cv =10)

Fitting 10 folds for each of 4 candidates, totalling 40 fits
[1 1 0 1 0 1 0 0 1 0 1 1 0 1 0 1 1 1 1 0 0 0 0 1 0 1 1 0 0 0 1 1 1 1 1 1
 1 0 0 0 0 1 0 1 1 1 1 0 1 1 0 1 1 0 1 0 1 1 0 1 1 1 1 0 1 1 0 1 0 1 0 0
 1 1 1 0 0 0 1 0 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 0 1 1 0 0 0 0 0 0 0 1
 0 1 0 0 1 0 1 1 1 0 0 1 1 0 0 0 0 0 0 0 1 0 1 1 0 0 1 0 0 0 0 0 0]
Best Params : {'n_estimators': 1000}
Classification Report:

```

| | | precision | recall | f1-score | support |
|--|--------------|-----------|--------|----------|---------|
| | 0 | 0.98 | 0.96 | 0.97 | 68 |
| | 1 | 0.96 | 0.99 | 0.97 | 75 |
| | accuracy | | | 0.97 | 143 |
| | macro avg | 0.97 | 0.97 | 0.97 | 143 |
| | weighted avg | 0.97 | 0.97 | 0.97 | 143 |

```

Accuracy Score 0.972027972027972
Confusion Matrix :
[[65  3]
 [ 1 74]]

In [40]: param = {
          'C': [0.1,1,100,1000],

```

Accuracy for Random forest: from the table :0.97

- SVC

```
[ 1 74]]

In [40]: param = {
          'C': [0.1,1,100,1000],
          'gamma':[0.0001,0.001, 0.005, 0.1,1, 3,5,10, 100]
        }
FitModel (X_res [feat_to_keep], Y_res,'SVC',SVC(), param, cv =5)

Fitting 5 folds for each of 36 candidates, totalling 180 fits
[1 1 0 1 0 1 0 1 1 0 1 1 0 1 1 0 1 1 1 1 1 0 0 0 0 1 0 1 1 0 0 0 1 1 1 1 1 1
 1 0 0 0 0 1 0 1 1 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 1 1 1 0 1 1 0 1 1 1 0 0
 1 1 1 0 0 0 1 0 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 0 1 0 0 0 0 0 0 0 0 0 0 1
 0 1 0 0 1 0 1 1 1 0 0 1 1 0 0 0 0 0 0 1 0 1 1 0 0 1 0 0 0 0 0 1 0]
Best Params : {'C': 100, 'gamma': 10}
Classification Report:

```

| | | precision | recall | f1-score | support |
|--------------|------|-----------|--------|----------|---------|
| | 0 | 1.00 | 0.94 | 0.97 | 68 |
| | 1 | 0.95 | 1.00 | 0.97 | 75 |
| accuracy | | | 0.97 | 143 | |
| macro avg | 0.97 | 0.97 | 0.97 | 143 | |
| weighted avg | 0.97 | 0.97 | 0.97 | 143 | |

```

Accuracy Score 0.972027972027972
Confusion Matrix :
[[64  4]
 [ 0 75]]

```

Accuracy for SVC : from the table :0.97.

- XGBoost:

```

In [41]: param = { 'n_estimators': [100,500,1000,2000] }
FitModel (X_res [feat_to_keep], Y_res,'XGBoost', XGBClassifier(),param, cv = 5)

Fitting 5 folds for each of 4 candidates, totalling 20 fits
[21:33:00] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[1 1 0 1 0 1 0 0 1 0 0 1 0 1 0 1 1 1 1 1 0 0 0 0 1 0 1 1 0 0 0 1 1 1 1 1 1
 1 0 0 0 0 1 0 1 1 1 1 0 1 1 0 1 1 0 1 1 0 1 1 1 1 1 0 1 1 0 1 0 1 0 0
 1 1 1 0 0 0 1 0 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 0 1 1 0 0 0 0 0 0 0 1
 0 1 0 0 1 0 1 1 1 0 0 1 1 0 0 0 0 0 1 0 1 1 0 0 1 0 0 0 0 0 0]
Best Params : {'n_estimators': 100}
Classification Report:

```

| | | precision | recall | f1-score | support |
|--------------|------|-----------|--------|----------|---------|
| | 0 | 0.99 | 0.97 | 0.98 | 68 |
| | 1 | 0.97 | 0.99 | 0.98 | 75 |
| accuracy | | | 0.98 | 143 | |
| macro avg | 0.98 | 0.98 | 0.98 | 143 | |
| weighted avg | 0.98 | 0.98 | 0.98 | 143 | |

```

Accuracy Score 0.9790209790209791
Confusion Matrix :
[[66  2]
 [ 1 74]]

In [42]: param = { 'n_estimators': [100,500,1000,2000] }
FitModel (X_res, Y_res, 'Random Forest', RandomForestClassifier(), param, cv =10)

Fitting 10 folds for each of 4 candidates, totalling 40 fits

```

Accuracy for xGBoost: from the table :0.98.

```
[ 0 75]]

In [44]: param = { 'n_estimators': [100,500,1000,2000] }
FitModel (X_res, Y_res,'XGBoost', XGBClassifier(),param, cv = 5)

Fitting 5 folds for each of 4 candidates, totalling 20 fits
[21:33:44] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[1 1 0 1 0 1 0 0 1 0 0 1 0 1 0 1 1 1 1 1 0 0 0 0 1 0 1 1 0 0 0 0 1 1 1 1 1 1
1 0 0 0 0 1 0 1 1 1 1 0 1 1 0 1 1 0 1 1 0 1 1 1 1 1 0 1 1 0 1 1 0 0 0
1 1 1 0 0 0 1 0 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 0 1 1 0 0 0 0 0 0 0 1
0 1 0 0 1 0 1 1 1 0 0 1 1 0 0 0 0 0 0 1 0 1 1 0 0 1 0 0 0 0 0 0]
Best Params : {'n_estimators': 100}
Classification Report:
              precision    recall  f1-score   support

     0       1.00        0.99        0.99         68
     1       0.99        1.00        0.99         75

 accuracy          0.99        0.99        0.99        143
 macro avg         0.99        0.99        0.99        143
weighted avg         0.99        0.99        0.99        143

Accuracy Score 0.993006993006993
Confusion Matrix :
[[67  1]
 [ 0 75]]

In [45]: load_model = pickle.load(open("XGBoost", "rb"))

In [46]: pred1 = load_model.predict(x_test)
```

Executed Hyperparameter Tuning using **GridsearchCV** resulting in a maximum of **99% =0.99 accuracy** by the XGBoost

Conclusion

The results demonstrate that the chosen classification algorithms, namely **Support Vector Machines, Random Forest, and XGBoost**, all achieved high accuracy and strong performance in **predicting cancer cases**. These models hold promise in aiding medical professionals in early cancer detection, which can significantly impact patient outcomes and treatment strategies. Further research could explore the integration of additional features and **optimization** of model **hyperparameters**.

Import the nesserly librarries:

```
%matplotlib inline

import numpy as np

import pandas as pd

import seaborn as s

# seaborn and pandas for unbalanced data exploration and visualization

from sklearn import model_selection

from sklearn.model_selection import train_test_split

from sklearn.model_selection import KFold

from sklearn.model_selection import cross_val_score

from sklearn.ensemble import RandomForestClassifier

from sklearn import metrics

from sklearn.metrics import confusion_matrix, classification_report, accuracy_score

from sklearn import svm

import matplotlib.pyplot as plt

from sklearn.preprocessing import StandardScaler, LabelEncoder

#optimum parameter choosing

from sklearn.model_selection import GridSearchCV

from sklearn.svm import SVC

from xgboost import XGBClassifier

import pickle

import os

import warnings

warnings.filterwarnings('ignore')
```