

Machine Learning Lab 4 - Predicting Breast Cancer

Submitted By
Name: **Rahad Nishit Shaleeh**
Register Number: **19112014**
Class: **5 BSc Data Science**

Lab Overview

About the Dataset

Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image. n the 3-dimensional space is as described in: [K. P. Bennett and O. L. Mangisarian: "Robust Linear Programming Discrimination of Two Linearly Inseparable Sets", Optimization Methods and Software 1, 1992, 23-34].

This database is also available through the UW CS ftp server:

- ftp.cs.wisc.edu - cd math-prog/tpo-datasets/machine-learn/WDBC/

Also can be found on UCI Machine Learning Repository:

- <https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29>

Attributes: 34

1. 0 = smoothest (local variation in radius lengths)
2. Diagnosis (M = malignant, B = benign)
3. Ten real-valued features are computed for each cell nucleus:
A. radius (mean of distances from center to points on the perimeter)
B. texture (standard deviation of gray-scale values)
C. perimeter
D. area
E. smoothness (local variation in radius lengths)
F. compactness (perimeter²/area - 1.0)
G. concavity (severity of concave portions of the contour)
H. concave points (number of concave portions of the contour)
I. symmetry
J. fractal dimension ("coastline approximation" - 1)
4. The mean, standard error and "worst" or "largest" (one of the three largest values) of these features were computed for each image, resulting in 30 features. For instance, field 3 is Mean Radius, field 13 is Radius SE, field 23 is Worst Radius.
5. All feature values are recorded with four significant digits.
6. Missing attribute values: none
7. Class distribution: 265 benign, 212 malignant

Objective

- Get familiar with the problem statement. Know the dataset thoroughly. Analyse the given dataset by exploring the hidden insights with beautiful visuals and Train & Test the model for accurate classification prediction of Breast Cancer.

Problem Definition

- Understand the Dataset & Features.
- Perform Data Preprocessing Technique to Get Balanced Structured Data.
- Perform Statistical Data Analysis and Derive Valuable Inferences.
- Perform Exploratory Data Analysis and Derive Valuable Insights.
- Train and Test through Different Classification Models for Better Prediction.

Approach

This is an extension to the Problem Definition. Mention the process/approach that you have followed in order to reach out the above problem definition.

- Step 1: Know the dataset thoroughly.
- Step 2: Perform preprocessing on data.
- Step 3: Import needful libraries as an when you try to plot different graphs and evaluate the model.
- Step 4: Perform Statistical Data Analysis and Derive Valuable Inferences.
- Step 5: Perform Exploratory Data Analysis and Derive Valuable Insights.
- Step 6: Train and Test through Different Classification Models for Better Breast Cancer Prediction.
- Step 7: Help the doctors with insights for predicting if a patient is diagnose with breast cancer or not.

Sections

Here, mentioned sections are defined in the below code. For this lab, the sections are -

```
1. Lab Overview
2. Dataset Overview
3. Data Analysis Process
4. About Different Classification Models
5. Implementation and Evaluation of Different Classification Models
6. Conclusion

References
1. https://pandas.pydata.org/
2. https://matplotlib.org/
3. https://seaborn.pydata.org/
4. https://plotly.com/
5. https://rockt-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html
6. https://rockt-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html
7. https://rockt-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html
8. https://www.kaggle.com/ucimlbreast-cancer-wisconsin-data

In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

%matplotlib inline

import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: df = pd.read_csv('BreastCancer.csv')

In [3]: df.shape

Out[3]: (569, 33)
```

```
In [4]: df.columns

Out[4]: Index(['id', 'diagnosis', 'radius_mean', 'texture_mean', 'perimeter_mean', 'area_mean', 'smoothness_mean', 'compactness_mean', 'concavity_mean', 'concave_points_mean', 'symmetry_mean', 'fractal_dimension_mean', 'radius_worst', 'texture_worst', 'perimeter_worst', 'area_worst', 'smoothness_worst', 'compactness_worst', 'concavity_worst', 'concave_points_worst', 'symmetry_worst', 'fractal_dimension_worst', 'Unnamed: 32'],
dtype='object')
```

```
In [5]: df.head()

Out[5]:
```

```
In [6]: df.tail()

Out[6]:
```

```
In [7]: df.drop('Unnamed: 32', axis = 1, inplace = True)
df.drop('id', axis = 1, inplace = True)
```

```
In [8]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 569 entries, 0 to 568
Data columns (total 31 columns):
#   Column              Non-Null Count  Dtype
--  --
0   diagnosis            569 non-null    object
1   radius_mean          569 non-null    float64
2   texture_mean         569 non-null    float64
3   perimeter_mean       569 non-null    float64
4   area_mean            569 non-null    float64
5   smoothness_mean      569 non-null    float64
6   compactness_mean     569 non-null    float64
7   concavity_mean       569 non-null    float64
8   concave_points_mean  569 non-null    float64
9   symmetry_mean        569 non-null    float64
10  fractal_dimension_mean 569 non-null    float64
11  radius_se            569 non-null    float64
12  texture_se           569 non-null    float64
13  perimeter_se         569 non-null    float64
14  area_se              569 non-null    float64
15  smoothness_se       569 non-null    float64
16  compactness_se       569 non-null    float64
17  concavity_se         569 non-null    float64
18  concave_points_se    569 non-null    float64
19  symmetry_se          569 non-null    float64
20  fractal_dimension_se  569 non-null    float64
21  radius_worst         569 non-null    float64
22  texture_worst        569 non-null    float64
23  perimeter_worst      569 non-null    float64
24  area_worst           569 non-null    float64
25  smoothness_worst     569 non-null    float64
26  compactness_worst    569 non-null    float64
27  concavity_worst      569 non-null    float64
28  concave_points_worst 569 non-null    float64
29  symmetry_worst       569 non-null    float64
30  fractal_dimension_worst 569 non-null    float64
dtypes: float64(28), object(3)
memory usage: 137.9+ KB
```

```
In [9]: df.sample(5)

Out[9]:
```

```
In [10]: df.describe()

Out[10]:
```

```
In [11]: df.isnull().sum()

Out[11]:
```

```
diagnosis            0
radius_mean          0
texture_mean         0
perimeter_mean       0
area_mean            0
smoothness_mean      0
compactness_mean     0
concavity_mean       0
concave_points_mean  0
symmetry_mean        0
fractal_dimension_mean 0
radius_se            0
texture_se           0
perimeter_se         0
area_se              0
smoothness_se       0
compactness_se       0
concavity_se         0
concave_points_se    0
symmetry_se          0
fractal_dimension_se 0
radius_worst         0
texture_worst        0
perimeter_worst      0
area_worst           0
smoothness_worst     0
compactness_worst    0
concavity_worst      0
concave_points_worst 0
symmetry_worst       0
fractal_dimension_worst 0
dtype: int64
```

Exploratory Data Analysis

\$ pip install <https://github.com/pandas-profiling/gandas-profilng>
\$ from pandas_profiling import ProfileReport

\$ profile = ProfileReport(df)
\$ profile.to_notebook_iframe()
\$ profile.to_html(output_file = "19112014_NishitRahad_EDA_BreastCancer.html")

Implementing Classification Models

Bifurcating Classification Parameter from the Dataset.

```
In [12]: Y = df['diagnosis']
X = df.drop(['diagnosis'], axis=1)

In [13]: X.sample(5)

Out[13]:
```

```
diagnosis            0
radius_mean          0
texture_mean         0
perimeter_mean       0
area_mean            0
smoothness_mean      0
compactness_mean     0
concavity_mean       0
concave_points_mean  0
symmetry_mean        0
fractal_dimension_mean 0
radius_se            0
texture_se           0
perimeter_se         0
area_se              0
smoothness_se       0
compactness_se       0
concavity_se         0
concave_points_se    0
symmetry_se          0
fractal_dimension_se 0
radius_worst         0
texture_worst        0
perimeter_worst      0
area_worst           0
smoothness_worst     0
compactness_worst    0
concavity_worst      0
concave_points_worst 0
symmetry_worst       0
fractal_dimension_worst 0
dtype: int64
```

```
In [14]: Y.sample(5)

Out[14]:
```

```
In [15]: from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.30, random_state = 8)

In [16]: X_train.sample(5)

Out[16]:
```

```
radius_mean          0
texture_mean         0
perimeter_mean       0
area_mean            0
smoothness_mean      0
compactness_mean     0
concavity_mean       0
concave_points_mean  0
symmetry_mean        0
fractal_dimension_mean 0
radius_se            0
texture_se           0
perimeter_se         0
area_se              0
smoothness_se       0
compactness_se       0
concavity_se         0
concave_points_se    0
symmetry_se          0
fractal_dimension_se 0
radius_worst         0
texture_worst        0
perimeter_worst      0
area_worst           0
smoothness_worst     0
compactness_worst    0
concavity_worst      0
concave_points_worst 0
symmetry_worst       0
fractal_dimension_worst 0
dtype: int64
```

```
In [17]: Y_train.sample(5)

Out[17]:
```

```
In [18]: Y_test.sample(5)

Out[18]:
```

```
In [19]: Y_train.sample(5)

Out[19]:
```

```
In [20]: from sklearn.linear_model import LogisticRegression
logistic_regressor = LogisticRegression()

In [21]: logistic_regressor.fit(X_train, Y_train)

Out[21]: LogisticRegression()
```

```
In [22]: predA = logistic_regressor.predict(X_test)
y_prob = logistic_regressor.predict_proba(X_test)

In [23]: predA[0:5]

Out[23]: array(['B', 'B', 'B', 'B', 'B'], dtype=object)
```

```
In [24]: from sklearn.metrics import accuracy_score
accuracy_score(predA, Y_test)

Out[24]: 0.9532163742690059
```

```
In [25]: y_prob[0:5]

Out[25]: array([[0.99752283, 0.00247717],
[0.98704693, 0.01295307],
[0.9937725, 0.0062275],
[0.9951079, 0.0048921],
[0.9993646, 0.0006354]])
```

K - Nearest Neighbours

Fitting and Predicting

```
In [26]: from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors = 2)
knn.fit(X_train, Y_train)

In [27]: predB = knn.predict(X_test)
accuracy_score(predB, Y_test)

Out[27]: 0.935675146159883
```

Decision Tree

Fitting and Predicting

```
In [27]: from sklearn.tree import DecisionTreeClassifier
dt = DecisionTreeClassifier()
dt.fit(X_train, Y_train)

In [28]: predC = dt.predict(X_test)
accuracy_score(predC, Y_test)

Out[27]: 0.9298245614050888
```

Classification Report

```
In [28]: from sklearn.metrics import classification_report
print(classification_report(Y_test, predA))

precision    recall  f1-score   support

B           0.95         0.96         0.96         105
B           0.94         0.94         0.94          66

accuracy         0.95         0.95         0.95         171
macro avg       0.95         0.95         0.95         171
weighted avg    0.95         0.95         0.95         171
```

Confusion Matrix

```
In [29]: from sklearn.metrics import confusion_matrix

CM = confusion_matrix(Y_test, predA)

Out[30]: array([[181, 4],
[ 4, 62]])
```

```
In [31]: Accuracy = (CM[0][0] + CM[1][1]) / (CM[0][0] + CM[1][0] + CM[0][1] + CM[1][1])
Accuracy

Out[31]: 0.9532163742690059
```

```
In [32]: ErrorRate = (CM[0][1] + CM[1][0]) / (CM[0][0] + CM[1][1] + CM[0][1] + CM[1][0])
ErrorRate

Out[32]: 0.04678362573099415
```

```
In [33]: Sensitivity = CM[0][0] / (CM[0][0] + CM[1][0])
Sensitivity

Out[33]: 0.9619047619047619
```

```
In [34]: Specificity = CM[1][1] / (CM[1][1] + CM[0][1])
Specificity

Out[34]: 0.9393939393939394
```

```
In [35]: Recall = CM[0][0] / (CM[0][0] + CM[1][0])
Recall

Out[35]: 0.9619047619047619
```

```
In [36]: Precision = CM[0][0] / (CM[0][0] + CM[0][1])
Precision

Out[36]: 0.9619047619047619
```

```
In [37]: FScore = (2 * (Precision * Recall)) / (Precision + Recall)
FScore

Out[37]: 0.959047619047619
```

Evaluating the Effect of Parameters For Logistic Regression

```
In [38]: from sklearn.metrics import accuracy_score

def doLogisticRegression(X, Y, test_size = 0.30, random_state = 42, penalty='l2', solvers='lbfgs'):
    X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = test_size, random_state = random_state)
    logistic_regressor = LogisticRegression(penalty = penalty, solver = solver)
    logistic_regressor.fit(X_train, Y_train)
    y_pred = logistic_regressor.predict(X_test)
    acc_score = accuracy_score(Y_test, y_pred)
    return acc_score

In [40]: df2 = pd.DataFrame(columns = ['Test Size', 'Random States', 'Penalty', 'Solvers', 'Accuracy'])
df2

Out[40]:
```

```
In [41]: penalties = ['none', 'l2']
test_size = [0.30, 0.25, 0.20]
random_states = [21, 42, 84]
solvers = ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']

for t_size in test_size:
    for r_state in random_states:
        for penalty in penalties:
            for solver in solvers:
                accuracy = doLogisticRegression(X, Y, t_size, r_state, penalty)
                BreastCancerResults['Test Size'] = t_size
                BreastCancerResults['Random States'] = r_state
                BreastCancerResults['Penalty'] = penalty
                BreastCancerResults['Solvers'] = solver
                BreastCancerResults['Accuracy'] = accuracy
                df2 = df2.append(BreastCancerResults, ignore_index = True)

In [42]: df2.sample(10)

Out[42]:
```

```
Test Size  Random States  Penalty  Solvers  Accuracy
23  0.30                B4         none      sag      0.959047619047619
54  0.25                B4         none      saga     0.9500395
84  0.20                B4         none      saga     0.973684
27  0.30                B4         l2         liblinear 0.959054
77  0.20                42         l2         liblinear 0.964912
9  0.30                21         l2         saga      0.947988
25  0.30                B4         l2         newton-cg 0.959054
43  0.25                42         none      sag      0.9500395
68  0.20                21         l2         sag      0.921053
58  0.25                B4         l2         sag      0.9500395
```

Evaluating the Effect of Parameters For K Nearest Neighbours

```
In [43]: def doKNearestNeighbour(X, Y, test_size = 0.30, randomstate = 8, n = 5):
    X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = test_size, random_state = randomstate)
    clst = KNeighborsClassifier(n_neighbors = n)
    clst.fit(X_train, Y_train)
    predB = clst.predict(X_test)
    acc_score1 = accuracy_score(predY, Y_test)
    return acc_score1

In [44]: test_size = [0.30, 0.25, 0.20, 0.10]
random_states = [8, 27, 42]
n_neighbors = [2, 3, 4, 5]

criteria1 = ['gini', 'entropy']
max_features1 = ['auto', 'sqrt', 'log2']
penalties = ['l1', 'l2', 'l1l2', 'none', 'l2']
solvers = ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']

In [46]: df3 = pd.DataFrame(columns = ['Test Size', 'Random States', 'Number of neighbours', 'K-nearest Neighbour Accuracy'])

for t_size in test_size:
    for r_state in random_states:
        for neigh in n_neighbors:
            k1 = doKNearestNeighbour(X, Y, t_size, r_state, neigh)
            KNearestNeighbours['Test Size'] = t_size
            KNearestNeighbours['Random States'] = r_state
            KNearestNeighbours['Number of neighbours'] = neigh
            KNearestNeighbours['K-nearest Neighbour Accuracy'] = k1
            df3.append(KNearestNeighbours, ignore_index = True)

In [47]: df3.sample(10)

Out[47]:
```

```
Test Size  Random States  Number of neighbours  K-nearest neighbour Accuracy
23  0.25                42.0                5.0                0.951040
6  0.30                27.0                3.0                0.929777
43  0.10                27.0                4.0                0.947768
9  0.30                42.0                3.0                0.941520
2  0.30                8.0                4.0                0.947368
31  0.20                27.0                5.0                0.929825
0  0.30                8.0                2.0                0.935673
28  0.20                27.0                2.0                0.936688
29  0.20                27.0                2.0                0.899685
26  0.20                8.0                4.0                0.946912
```

Evaluating the Effect of Parameters For Decision Tree

```
In [48]: def doDecisionTree(X, Y, test_size = 0.20, randomstate = 8, c = 'gini', m = 'auto'):
    X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = test_size, random_state = randomstate)
    clst = DecisionTreeClassifier(criterion=c, max_features = m)
    clst.fit(X_train, Y_train)
    predB = clst.predict(X_test)
    acc_score2 = accuracy_score(predY, Y_test)
    return acc_score2

In [49]: test_size = [0.30, 0.25, 0.20, 0.10]
random_states = [8, 27, 42]
n_neighbors = [2, 3, 4, 5]

criteria2 = ['gini', 'entropy']
max_features2 = ['auto', 'sqrt', 'log2']
penalties = ['l1', 'l2', 'l1l2', 'none', 'l2']
solvers = ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']

In [50]: df4 = pd.DataFrame(columns = ['Test Size', 'Random States', 'Decision Tree Accuracy', 'Criteria', 'Max Features'])

for t_size in test_size:
    for r_state in random_states:
        for cfs in max_features:
            a2 = doDecisionTree(X, Y, t_size, r_state, cfs, mfs)
            DecisionTree['Test Size'] = t_size
            DecisionTree['Random States'] = r_state
            DecisionTree['Decision Tree Accuracy'] = a2
            DecisionTree['Criteria'] = cfs
            DecisionTree['Max Features'] = mfs
            df4.append(DecisionTree, ignore_index = True)

In [52]: df4.sample(10)

Out[52]:
```

```
Test Size  Random States  Decision Tree Accuracy  Criteria  Max features
23  0.25                8                0.977653      entropy      log2
66  0.10                42                0.912081      gini         auto
14  0.30                42                0.941520      gini         log2
51  0.20                42                0.938696      entropy      auto
29  0.25                27                0.906091      entropy      log2
68  0.10                42                0.947368      gini         log2
65  0.10                42                0.929825      gini         log2
17  0.30                42                0.929825      entropy      log2
59  0.10                8                0.912081      entropy      log2
27  0.25                27                0.910284      entropy      auto
```

Conclusion

In this lab, we have tried to gain the knowledge about data and its variables, further we did some preprocessing to the data in order to bring it into more analyst friendly mode, later we implemented various graphs using various libraries in order to get valuable insights, furthermore, we implemented and evaluated various classification models to get high accuracy in terms of predicting breast cancer which can help the doctors to predict breast cancer for the patients.