

# Machine Learning Lab 5 - Predicting Breast Cancer II

Submitted By  
Name: **Ramkish Chakrabarti**  
Register Number: **3112014**  
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 n-3 dimensional space is as that described in: [K. P. Bennett and O. L. Mangasarian, "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/cpo-dataset/machine-learn/WDBC/

Also can be found on UCI Machine Learning Repository:

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

### Attribute Information:

1. ID number
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 (perimeter<sup>2</sup> / 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" (mean 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: 357 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 medical libraries as in when you try to plot different graphs and evaluate the model.
- Step 4: Perform Statistical Data Analysis and Derive Valuable Insights.
- 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 doctor's diagnosis by predicting if a patient is diagnosed 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

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_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se',
       'compactness_se', 'concavity_se', 'concave points_se', 'symmetry_se',
       'fractal_dimension_se', '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]:
   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
0   842302     M    17.99     10.39     122.80    1001.0     0.11840     0.17760     0.3001     0.14720     0.17917     0.10360     17.32     184.60    2019.0     1866.0     0.1623
1   842517     M    20.57     17.77     132.90    1296.0     0.09674     0.07964     0.0968     0.07017     0.1890     0.10360     18.41     198.50    1965.0     1866.0     0.1228
2   84209093     M    19.69     21.25     130.00    1203.0     0.10960     0.11990     0.1974     0.12790     0.2543     0.12670     12.64     126.70    1709.0     1444.0     0.1444
3   84348031     M    11.42     20.38     77.58     386.1     0.14250     0.28390     0.2414     0.10520     0.2650     0.987     567.7     98.87     567.7     0.2098
4   84354042     M    21.29     14.34     135.10    1297.0     0.10030     0.13260     0.1980     0.10430     0.1607     152.20    1575.0     1575.0     0.1374

5 rows x 33 columns

In [6]:
df.tail()

Out[6]:
   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
564 126424     M    21.56     22.39     142.00    1479.0     0.11100     0.11590     0.14400     0.13990     0.1876     166.10    2027.0     1841.0     0.14100
565 126964     M    20.13     28.25     131.20    1261.0     0.09780     0.10340     0.14400     0.09791     0.3825     155.00    1741.0     1616.0     0.11660
566 126954     M    16.60     26.08     108.30     858.1     0.08465     0.10230     0.06251     0.05302     0.3412     126.70    1244.0     1013.0     0.11390
567 127241     M    20.60     29.33     140.10    1265.0     0.11780     0.27700     0.35140     0.15200     0.3942     184.60    1821.0     1616.0     0.10900
568 127251     B     7.76     24.54     47.92     181.0     0.05033     0.04362     0.00000     0.00000     0.3037     59.16     258.6     258.6     0.08996

5 rows x 33 columns

In [7]:
df.drop('Unnamed: 32', axis = 3, 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(30), object(1)
memory usage: 137.5+ KB

In [9]:
df.sample(5)

Out[9]:
   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
328         B    13.94     13.17     90.31     594.2     0.12480     0.09795     0.10100     0.06615     0.1976     14.62     15.38     94.52     653.
680         B    11.27     12.96     73.16     386.3     0.12370     0.11110     0.07900     0.05550     0.2018     12.84     20.53     84.93     476.
219         M    19.53     32.47     128.00     858.1     0.08420     0.11300     0.11450     0.06637     0.1428     27.90     45.41     180.20     2477.
336         B    12.99     14.23     84.08     514.3     0.09462     0.09995     0.03738     0.02098     0.1652     13.72     16.91     87.38     576.
128         B    15.10     16.39     99.58     674.5     0.11500     0.18070     0.11380     0.08534     0.2001     16.11     18.33     105.90     762.

5 rows x 31 columns

In [10]:
df.describe()

Out[10]:
   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
count  569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000
mean    14.12292     19.09648     101.96023    854.89104     0.096360     0.116041     0.186799     0.048019     0.184162     0.182768     16.249100    25.677223    1077.2
std     3.132409     4.930306     24.280913    201.04219     0.014064     0.020313     0.079720     0.038903     0.027144     0.097600     12.450250     17.650     81.1
min     6.981000     7.010000     43.790000    143.300000     0.062636     0.019380     0.000000     0.000000     0.106000     0.046990     11.930000    12.020000     50.4
25%    11.700000    16.170000    75.170000     420.300000     0.086370     0.064920     0.029560     0.020310     0.161800     0.097700     13.100000    21.980000     84.1
50%    13.370000    18.840000     86.240000    551.100000     0.095870     0.092630     0.061540     0.033500     0.178200     0.061540     14.970000    25.430000     97.6
75%    15.780000    21.900000    104.100000    782.700000     0.105300     0.130400     0.130700     0.074000     0.195700     0.066120     18.790000    29.720000    125.4
max     28.110000    39.280000    188.500000    2501.000000     0.163400     0.345400     0.426800     0.201200     0.336000     0.097440     36.040000    49.540000    251.2

5 rows x 30 columns

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

Exporatory Data Analysis
```

```
$ pip install https://github.com/pandas-profiling/pandas-profiling/archive/master.zip
```

\$ from pandas\_profiling import ProfileReport

\$ profile = ProfileReport(df)

\$ profile.to\_notebook\_iframe()

\$ profile.to\_file(output\_file = "19112014\_NishatRathod\_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]:
   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
32      17.02     23.98     112.80     899.3     0.11970     0.14860     0.24170     0.13200     0.2348     0.06302     20.88     32.09     136.
count      569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000    569.000000
mean      14.12292     19.09648     101.96023    854.89104     0.096360     0.116041     0.186799     0.048019     0.184162     0.182768     16.249100    25.677223    1077.2
std       3.132409     4.930306     24.280913    201.04219     0.014064     0.020313     0.079720     0.038903     0.027144     0.097600     12.450250     17.650     81.1
min       6.981000     7.010000     43.790000    143.300000     0.062636     0.019380     0.000000     0.000000     0.106000     0.046990     11.930000    12.020000     50.4
25%      11.700000    16.170000    75.170000     420.300000     0.086370     0.064920     0.029560     0.020310     0.161800     0.097700     13.100000    21.980000     84.1
50%      13.370000    18.840000     86.240000    551.100000     0.095870     0.092630     0.061540     0.033500     0.178200     0.061540     14.970000    25.430000     97.6
75%      15.780000    21.900000    104.100000    782.700000     0.105300     0.130400     0.130700     0.074000     0.195700     0.066120     18.790000    29.720000    125.4
max      28.110000    39.280000    188.500000    2501.000000     0.163400     0.345400     0.426800     0.201200     0.336000     0.097440     36.040000    49.540000    251.2

5 rows x 30 columns

In [14]:
y.sample(5)

Out[14]:
112     B
254     M
15      B
464     B
132     B
Name: diagnosis, dtype: object

Train Test Split
```

```
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  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
128    15.100     16.39     99.58     674.5     0.11500     0.18070     0.11380     0.08540     0.2001     16.110     18.33     105.90
505     8.676     10.14     64.12     275.5     0.12550     0.12040     0.11800     0.07880     0.2057     0.05475     10.600     18.04     67.
231    11.320     27.48     71.76     395.7     0.06083     0.03813     0.00643     0.00325     0.1869     0.05038     12.080     33.75     73.
306    11.600     24.09     74.23     417.2     0.07474     0.05688     0.01074     0.01330     0.1935     0.05078     12.440     31.62     81.
412     9.387     21.68     59.75     298.8     0.07969     0.06053     0.03735     0.00518     0.1274     0.06724     9.965     27.99     66.

5 rows x 30 columns

In [17]:
X_test.sample(5)

Out[17]:
   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
374    13.840     16.07     87.88     579.1     0.10332     0.10374     0.02556     0.02031     0.1872     0.08689     14.840     20.21     99.
510    11.740     14.69     76.31     436.0     0.08099     0.09661     0.06776     0.02629     0.1480     0.00768     12.450     17.60     81.
126    13.610     24.69     87.76     572.6     0.09268     0.07962     0.05265     0.03085     0.1761     0.06130     16.990     35.94     113.
175     8.671     16.45     54.42     227.2     0.09138     0.04276     0.00000     0.00000     0.1372     0.06724     9.262     17.04     58.
195    12.910     14.33     82.53     516.4     0.07941     0.05366     0.03873     0.02377     0.1829     0.05667     13.880     22.90     90.

5 rows x 30 columns

In [18]:
Y_train.sample(5)

Out[18]:
189     B
238     B
166     B
15      B
586     B
Name: diagnosis, dtype: object

In [19]:
Y_test.sample(5)

Out[19]:
100     M
523     B
457     B
431     B
653     B
Name: diagnosis, dtype: object

In [20]:
from sklearn.metrics import accuracy_score
```

## Decision Tree

### Fitting and Predicting

```
In [21]:
from sklearn.tree import DecisionTreeClassifier
dt = DecisionTreeClassifier()
dt.fit(X_train, Y_train)
predA = dt.predict(X_test)
accuracy_score(predA, Y_test)

Out[21]:
0.945284678362573
```

## Classification Report

```
from sklearn.metrics import classification_report
print(classification_report(Y_test, predA))
```

	precision	recall	f1-score	support
B	0.95	0.95	0.95	105
M	0.92	0.92	0.92	66
accuracy	0.94			
macro avg	0.94	0.94	0.94	171
weighted avg	0.94	0.94	0.94	171

## Confusion Matrix

```
from sklearn.metrics import confusion_matrix
```

```
CM = confusion_matrix(Y_test, predA)
CM
```

```
array([[105,  5],
       [ 5,  61]], dtype=int64)
```

Accuracy =  $\frac{CM[0][0] + CM[1][1]}{CM[0][0] + CM[0][1] + CM[1][1] + CM[1][0]}$

Accuracy

0.945284678362573

ErrorRate =  $\frac{CM[0][1] + CM[1][0]}{CM[0][0] + CM[1][1] + CM[0][1] + CM[1][0]}$

ErrorRate

0.0547153216374269

Sensitivity =  $\frac{CM[0][0]}{CM[0][0] + CM[1][0]}$

Sensitivity

0.95238952389523

Specificity =  $\frac{CM[1][1]}{CM[1][1] + CM[0][1]}$

Specificity

0.9242424242424242

Recall =  $\frac{CM[0][0]}{CM[0][0] + CM[1][0]}$

Recall

0.95238952389523

Precision =  $\frac{CM[0][0]}{CM[0][0] + CM[0][1]}$

Precision

0.95238952389523

F1Score =  $\frac{2 * (Precision * Recall)}{(Precision + Recall)}$

F1Score

0.95238952389523

## Random Forest Classifier

```
from sklearn.ensemble import RandomForestClassifier
rf = RandomForestClassifier(n_estimators=300)
rf.fit(X_train, Y_train)
predB = rf.predict(X_test)
accuracy_score(predB, Y_test)

Out[22]:
0.976768239181286
```

## Classification Report

```
from sklearn.metrics import classification_report
print(classification_report(Y_test, predB))
```

	precision	recall	f1-score	support
B	0.98	1.00	0.99	105
M	1.00	0.92	0.96	66
accuracy	0.97			
macro avg	0.99	0.96	0.97	171
weighted avg	0.97	0.97	0.97	171

## Confusion Matrix

```
CM = confusion_matrix(Y_test, predB)
CM
```

```
array([[105,  0],
       [ 0,  61]], dtype=int64)
```

Accuracy =  $\frac{CM[0][0] + CM[1][1]}{CM[0][0] + CM[0][1] + CM[1][1] + CM[1][0]}$

Accuracy

0.976768239181286

ErrorRate =  $\frac{CM[0][1] + CM[1][0]}{CM[0][0] + CM[1][1] + CM[0][1] + CM[1][0]}$

ErrorRate

0.02323176081871343

Sensitivity =  $\frac{CM[0][0]}{CM[0][0] + CM[1][0]}$

Sensitivity

0.9545454545454546

Specificity =  $\frac{CM[1][1]}{CM[1][1] + CM[0][1]}$

Specificity

1.0

Recall =  $\frac{CM[0][0]}{CM[0][0] + CM[1][0]}$

Recall

0.9545454545454546

Precision =  $\frac{CM[0][0]}{CM[0][0] + CM[0][1]}$

Precision

1.0

F1Score =  $\frac{2 * (Precision * Recall)}{(Precision + Recall)}$

F1Score

0.97674186844517

## Evaluating the Effect of Parameters For Decision Tree

```
In [42]:
def doDecisionTree(X, Y, test_size = 0.30, randomstate = 8, c='gini', mf = 'auto'):
    X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = test_size, random_state = random
```