# APSIT SKILLS INTERNSHIP – PROJECT REPORT

# PO COMPUTER

## Project Batch: B1

## Team Members Names:

Team Leader: - Rahul J. Bhiwande

Member 1: - Saurabh V. Yadav

Member 2: - Urvi V. Aryamane

## Technology Selected:

Python

## Project Topic Name:

Breast Cancer Classification using Machine Learning

## Problem statement:

During their life, among 8% of women are diagnosed with Breast cancer (BC), after lung cancer, BC is the second popular cause of death in both developed and undeveloped worlds. BC is characterized by the mutation of genes, constant pain, changes in the size, color(redness), skin texture of breasts. Today, Machine Learning (ML) techniques are being broadly used in the breast cancer classification problem. They provide high classification accuracy and effective diagnostic capabilities. The early diagnosis of BC can improve the prognosis and chance of survival significantly, as it can promote timely clinical treatment to patients. Further accurate classification of benign tumors can prevent patients undergoing unnecessary treatments. To observe which features are most helpful in predicting malignant or benign cancer and to see general trends that may aid us in model selection. The goal is to classify whether the breast cancer is benign or malignant.

## Detailed Workflow:

- Functions/modules used for the program:
    1. NumPy
    2. Pandas
    3. Matplotlib
    4. Scikit learn

- Flow of the code:

### Step 1: Data Preparation

We will use the UCI Machine Learning Repository for breast cancer dataset.

### Step 2: Data Exploration

We will import the necessary libraries and use the library functions to examine the available data.

### Step 3: Handling the Categorical Data

The data in the dataset is not always in the numeric form. Further processing requires the data to be in numeric form. Hence, in this step, we will convert the non-numeric data into numeric data.

### Step 4: Splitting the dataset

The data we use is usually split into training data and test data. The training set contains a known output and the model learns on this data in order to be generalized to other data later on. We have the test dataset (or subset) in order to test our model's prediction on this subset.

### Step 5: Feature Scaling

Most of the times, the data will be varying over a large range. Hence, this data is needed to be brought in a lower range using normalization.

### Step 6: Model Selection, classification and visualization

We will be using various classification models. In our dataset we have the outcome variable or Dependent variable i.e. Y having only two set of values, either M (Malign) or B(Benign). So, we will use Classification algorithm of supervised learning. Also, the data will be visualized in the form of graphs, using python libraries.
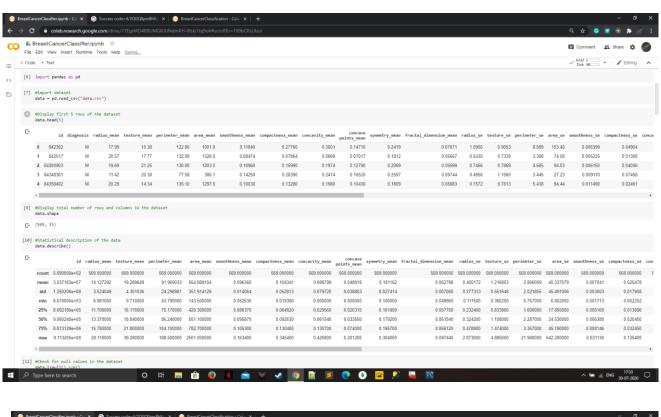
## GitHub / Drive link of project code:

**Data model:**
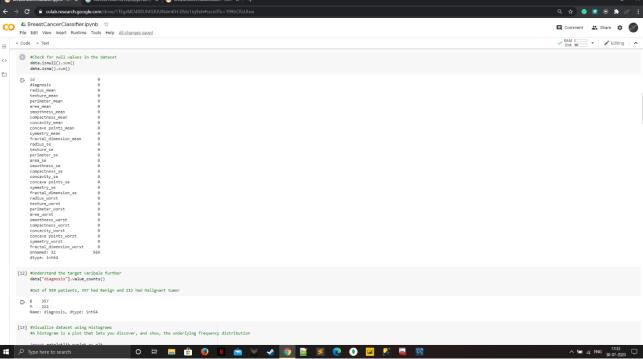https://colab.research.google.com/drive/1TEgrMO48BUMGRJUNdmKH-0fsIs1tq9oh

## Output Screenshots:

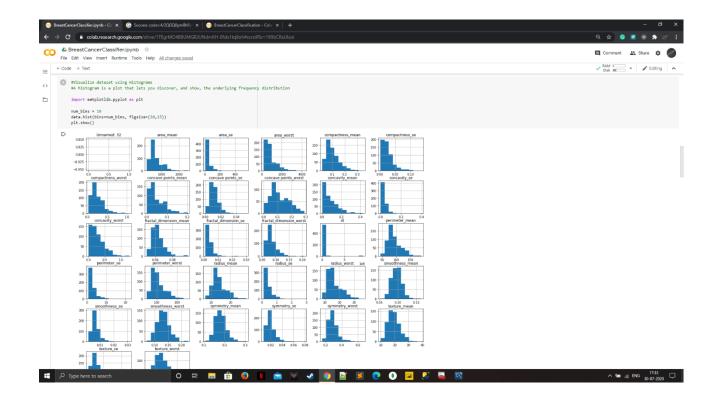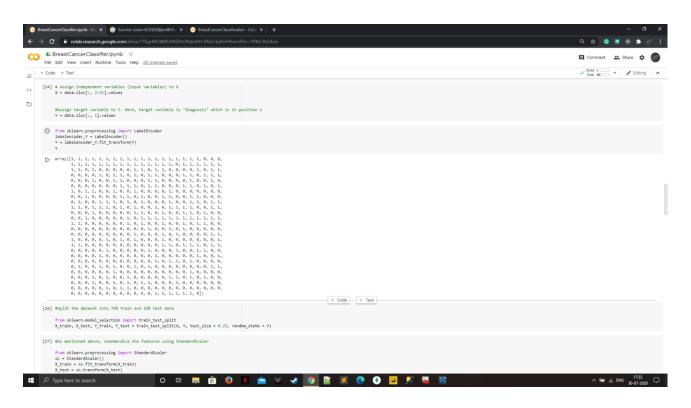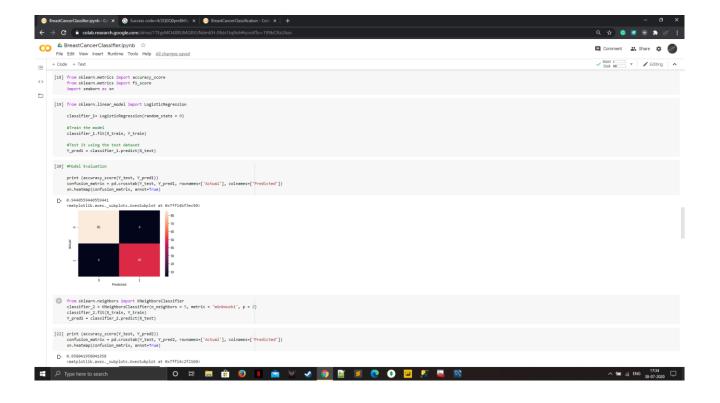## Data Model:

```
[33]  82.57,
      477.1,
      0.1278,
      0.17,
      0.1578,
      0.08089,
      0.2087,
      0.07613,
      0.3345,
      0.8902,
      2.217,
      27.19,
      0.00751,
      0.03345,
      0.03672,
      0.01137,
      0.02165,
      0.005082,
      15.47,
      23.75,
      103.4,
      741.6,
      0.1791,
      0.5249,
      0.5355,
      0.1741,
      0.3985,
      0.1244)
```

```
_input = sc.transform([[18.63,25.11,124.8,1088,0.1064,0.1887,0.2319, 0.1244, 0.2183, 0.06197, 0.8307, 1.466, 5.574, 105, 0.006248, 0.03374, 0.05196, 0.01158, 0.02007, 0.00456, 23.15 ,34.01, 160.5, 1670, 0.1491, 0.4257, 0.6133, 0.1848, 0.3444]])
:lassifier_7.predict(real_input) == [0] :
 int('Benign')
 : :
  int('Malignant')
```

```
Malignant
```

```
[35]  x1 = sc.transform([[12.45, 15.7, 82.57,  477.1, 0.1278, 0.17, 0.1578, 0.08089, 0.2087, 0.07613, 0.3345, 0.8902, 2.217, 27.19, 0.00751, 0.03345, 0.03672, 0.01137, 0.02165, 0.005082, 15.47, 23.75,  103.4 ,741.6, 0.1791, 0.5249, 0.5355, 0.1741, 0.3985
      x2 = sc.transform([[18.63,25.11,124.8,1088,0.1064,0.1887,0.2319, 0.1244, 0.2183, 0.06197, 0.8307, 1.466, 5.574, 105, 0.006248, 0.03374, 0.05196, 0.01158, 0.02007, 0.00456, 23.15 ,34.01, 160.5, 1670, 0.1491, 0.4257, 0.6133, 0.1848, 0.3444]])
```



```
[36]  plt.bar(['M','B'],j,color='cyan',alpha=0.5)
```
<BarContainer object of 2 artists>

```
[37]  from sklearn.preprocessing import LabelEncoder
      labelencoder_Y = LabelEncoder()
      df.iloc[:,1]= labelencoder_Y.fit_transform(df.iloc[:,1].values)
      print(labelencoder_Y.fit_transform(df.iloc[:,1].values))
```



```
[49]  model = models(X_train,Y_train)
```
```
[0]Logistic Regression Training Accuracy: 0.9906103286384976
[1]K Nearest Neighbor Training Accuracy: 0.9765258215962441
[2]Support Vector Machine (Linear Classifier) Training Accuracy: 0.9882629107981221
[3]Support Vector Machine (RBF Classifier) Training Accuracy: 0.9835680751173709
[4]Gaussian Naive Bayes Training Accuracy: 0.9507042253521126
[5]Decision Tree Classifier Training Accuracy: 1.0
[6]Random Forest Classifier Training Accuracy: 0.9953305164319248
```



```
m sklearn.metrics import confusion matrix
```

## Acknowledgment:

We have great pleasure in presenting our topic: Breast Cancer Classification using Machine Learning. We take this opportunity to express our sincere thanks towards our guide Prof. Amol R. Kalugade Department of Computer Engineering, APSIT thane for providing the technical guidelines and suggestions regarding line of work. We would like to express our gratitude towards his constant encouragement, support and guidance through the development of project. We also thank the coordinators of APSIT for their invaluable help rendered during the course of this work.

Student Name1: Urvi Aryamane

Student ID1: 15102001

Student Name2: Saurabh Yadav

Student ID2: 15102021

Student Name3: Rahul Bhiwande

Student ID3: 15102068

## Reference:

1. http://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+%28diagnostic%29
2. https://www.udemy.com/course/machinelearning/learn/lecture/19041246?start=1#overview