# BREAST CANCER PREDICTION
# END TERM REPORT


*by*


# AMAN PANDEY
# &
# ADITYA RAJ JHA


Section: KM118
Roll Numbers: RKM118B51 & RKM118A17

**Department of Intelligent Systems,**

**School of Computer Science Engineering,**

**Lovely Professional University, Jalandhar**

**November, 2022**

# STUDENT DECLARATION

This is to declare that this report has been written by us. No part of the report is copied from other sources. All information included from other sources have been duly acknowledged. We aver that if any part of the report is found to be copied, we are shall take full responsibility for it.

Aman Pandey
RKM118B51

Aditya Raj Jha
RKM118A17

Place: Lovely Professional University, Jalandhar
Date: 08th November 2022

# TABLE OF CONTENTS

|   | TITLE | PAGE NO. |
|---|-------|----------|

# BONAFIDE CERTIFICATE

Certified that this project report "Breast cancer prediction using machine learning" is the bonafide work of Aman Pandey and Aditya Raj verma who carried out the project work under my supervision.

<div align="right">

Dr. Dhanpratap Singh
Associate professor
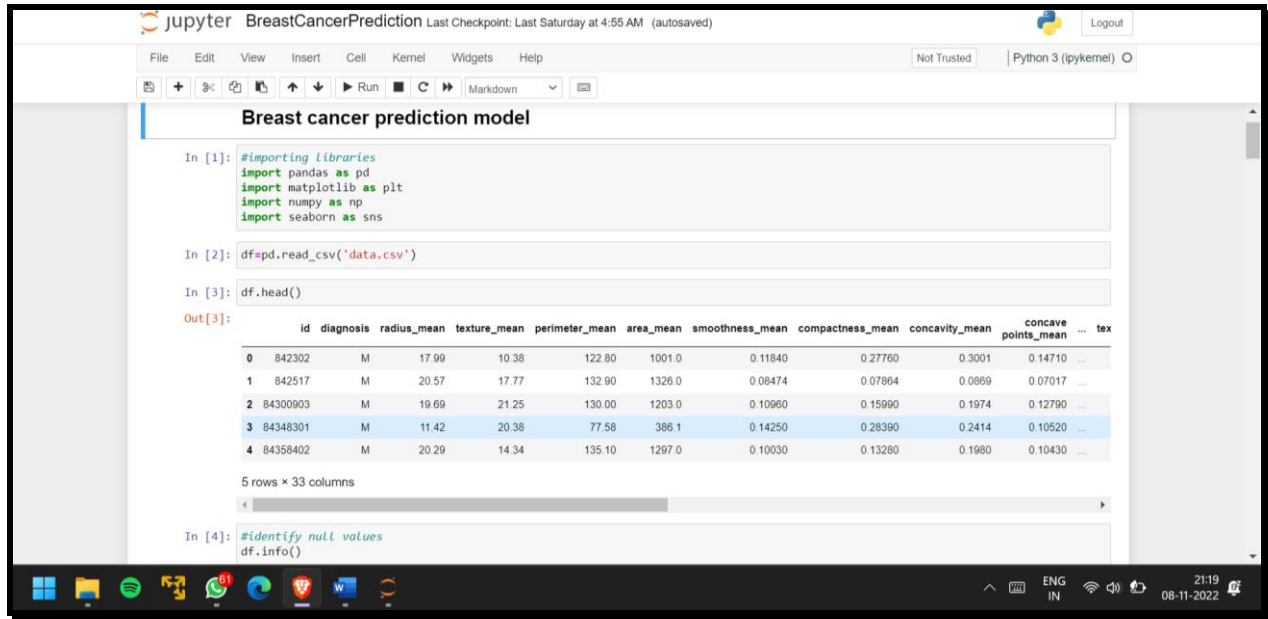25706
Intelligence System 1

</div>

# INTRODUCTION

Breast cancer has become the most recurrent type of health issue among women especially for women in middle age. Early detection of breast cancer can help women cure this disease and death rate can be reduced. In the present-day scenario, to observe breast cancer mammograms are used and they are known be the most effective scanning technique. In this prediction model the detection of cancer cells is done by machine learning technique.

Machine learning is an application of artificial intelligence that provides systems the ability to automatically learn and improve from experience without being explicitly programmed. The basic premise of machine learning is to build algorithms that can receive input data and use statistical analysis to predict an output while updating outputs as new data becomes available. The process of learning begins with observations or data, such as examples, direct experience, or instruction, to look for patterns in data and make better decisions in the future based on the examples that we provide. The primary aim is to allow the computers learn automatically without human intervention or assistance and adjust actions accordingly

Combining multiple risk factors in modelling for breast cancer prediction could help the early diagnosis of the disease with necessary care plans. Collection, storage, and management of different data and intelligent systems based on multiple factors for predicting breast cancer are effective in disease management.

# SYNTAX AND SCREENSHOTS

## Importing Libraries and csv file



## Checking the null values resent in attributes
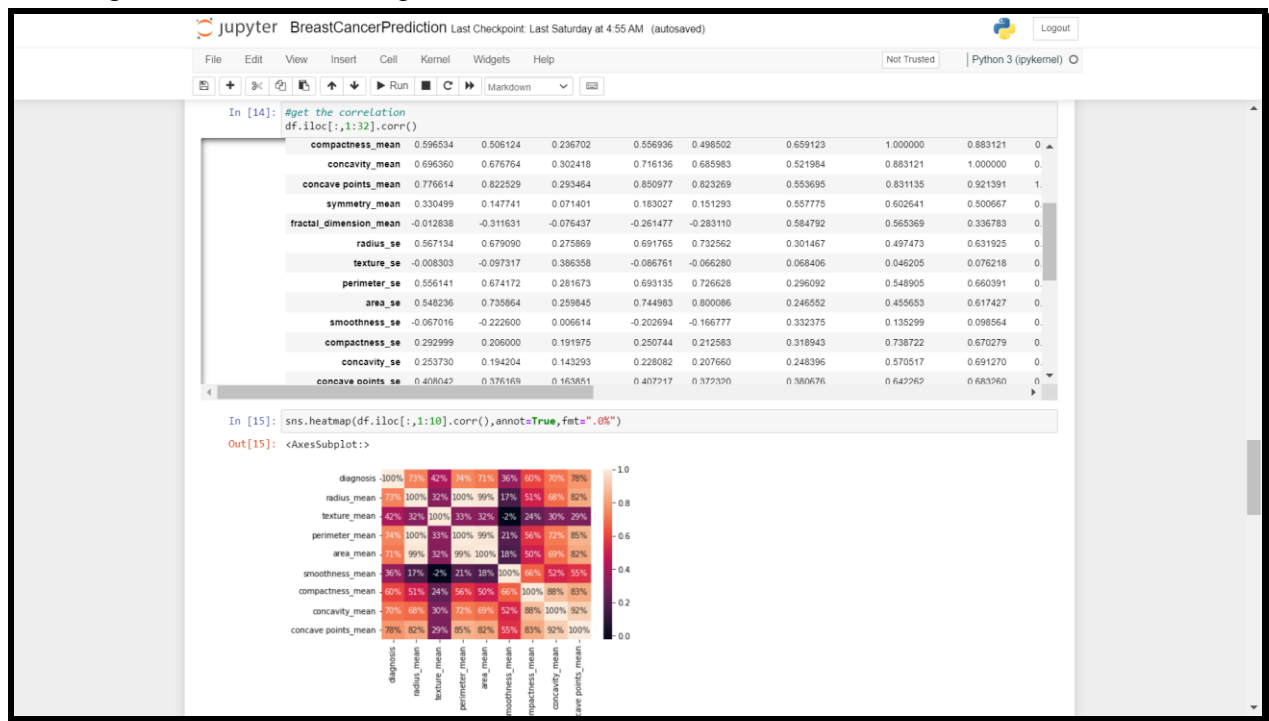
Representing malignant and benignant samples in for of graphs, and changing 'M' and 'B' into 1 and 0 respectively to for a correlation.
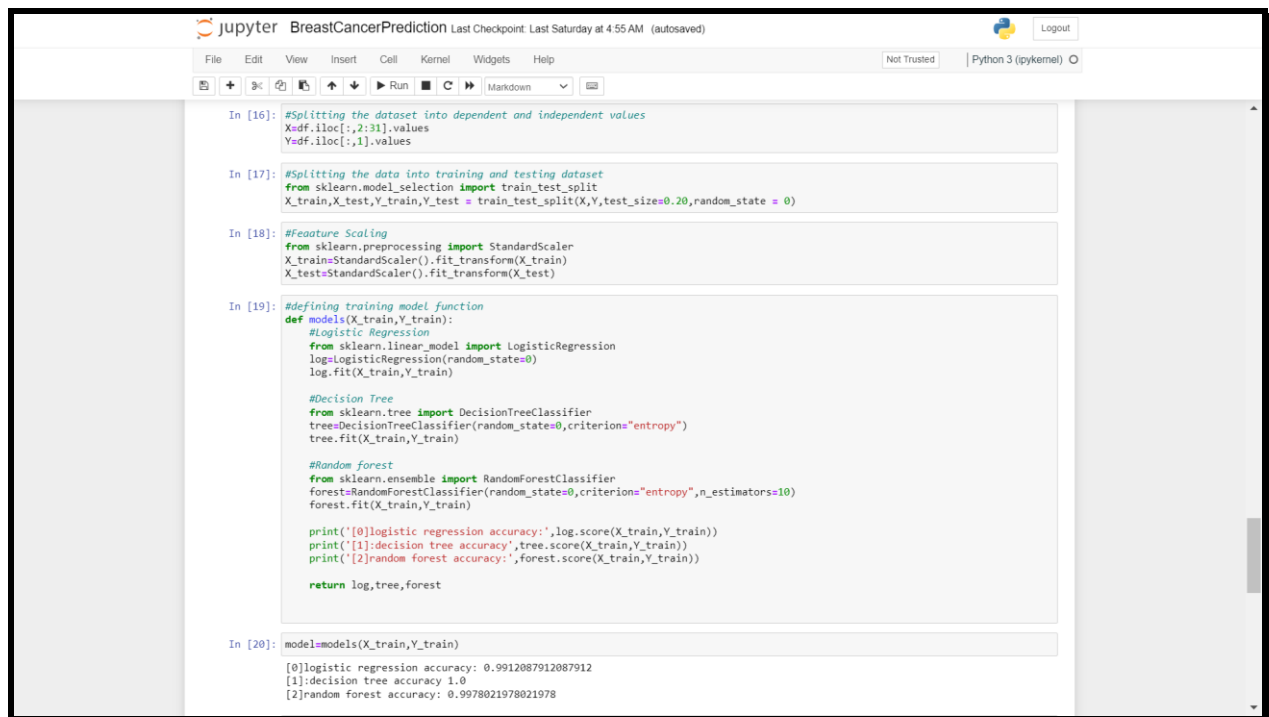


Generating a pairplot

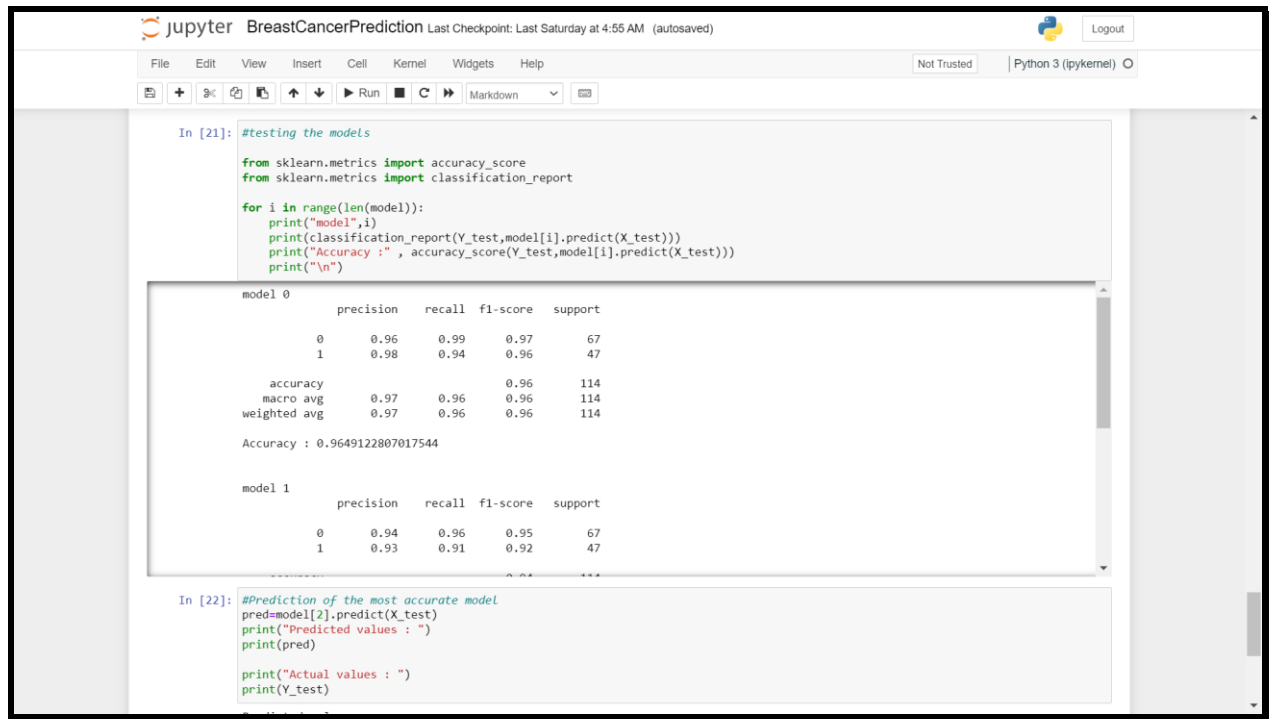## Checking the correlation among all attributes



## Splitting the data for training and testing purposes

## Testing the models



```python
In [21]: #testing the models

from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report

for i in range(len(model)):
    print("model",i)
    print(classification_report(Y_test,model[i].predict(X_test)))
    print("Accuracy :" , accuracy_score(Y_test,model[i].predict(X_test)))
    print("\n")
```

```
model 0
              precision    recall  f1-score   support

           0       0.96      0.99      0.97        67
           1       0.98      0.94      0.96        47

    accuracy                           0.96       114
   macro avg       0.97      0.96      0.96       114
weighted avg       0.97      0.96      0.96       114

Accuracy : 0.9649122807017544


model 1
              precision    recall  f1-score   support

           0       0.94      0.96      0.95        67
           1       0.93      0.91      0.92        47
```
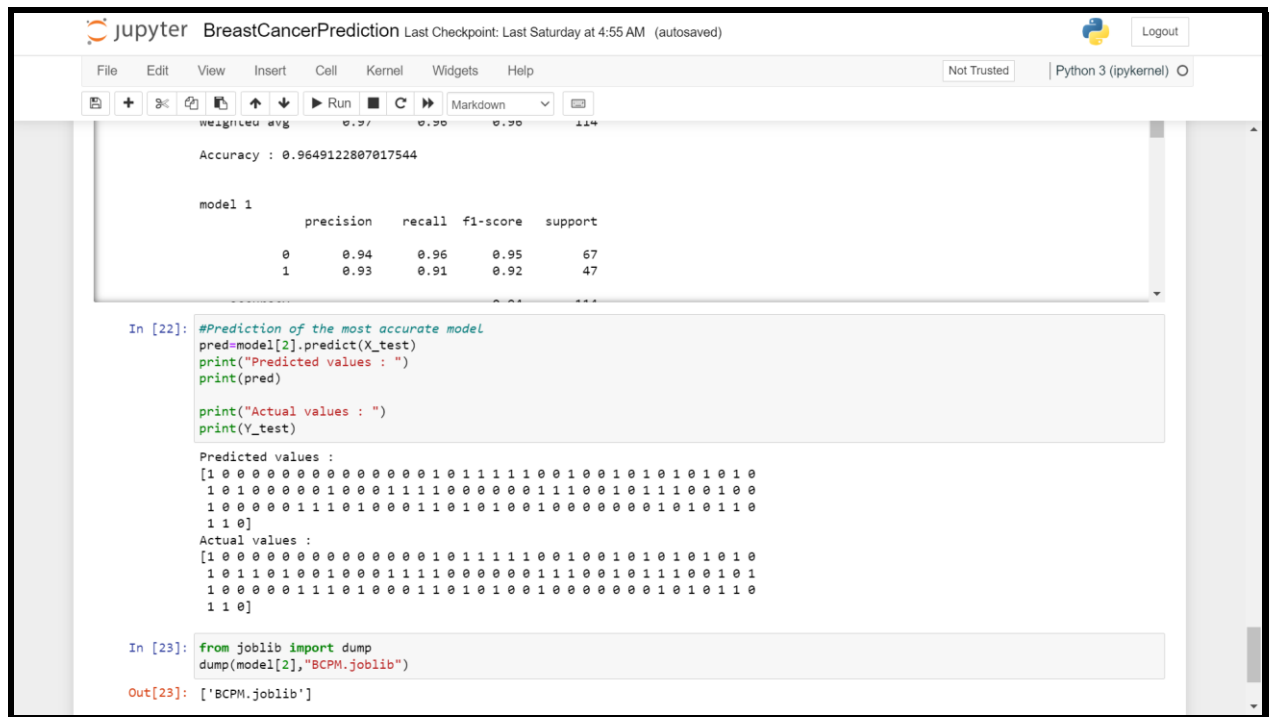
```python
In [22]: #Prediction of the most accurate model
pred=model[2].predict(X_test)
print("Predicted values : ")
print(pred)

print("Actual values : ")
print(Y_test)
```

## Prediction using the test cases



```
weighted avg     0.97      0.96      0.96       114

Accuracy : 0.9649122807017544


model 1
              precision    recall  f1-score   support

           0       0.94      0.96      0.95        67
           1       0.93      0.91      0.92        47
```

```python
In [22]: #Prediction of the most accurate model
pred=model[2].predict(X_test)
print("Predicted values : ")
print(pred)

print("Actual values : ")
print(Y_test)
```

```
Predicted values :
[1 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 1 1 1 0 0 1 0 0 1 0 1 0 1 0 1 0 1 0
 1 0 1 0 0 0 0 0 1 0 0 0 1 1 1 1 0 0 0 0 0 0 0 1 1 1 0 0 1 0 1 1 1 1 0 0 1 0 0
 1 0 0 0 0 1 1 1 0 1 0 0 0 1 1 0 1 0 1 0 0 1 0 0 0 0 0 0 0 0 1 0 1 0 1 1 0
 1 1 0]
Actual values :
[1 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 1 1 1 0 0 1 0 0 1 0 1 0 1 0 1 0 1 0
 1 0 1 1 0 1 0 0 1 0 0 0 1 1 1 1 0 0 0 0 0 0 0 1 1 1 0 0 1 0 1 1 1 0 0 1 0 1 0 1
 1 0 0 0 0 1 1 1 0 1 0 0 0 1 1 0 1 0 1 0 0 1 0 0 0 0 0 0 0 0 1 0 1 0 1 1 0
 1 1 0]
```

```python
In [23]: from joblib import dump
dump(model[2],"BCPM.joblib")

Out[23]: ['BCPM.joblib']
```

# ALGORITHMS USED

## Logistic regression:

It is the appropriate regression analysis to conduct when the dependent variable is dichotomous (binary). Like all regression analyses, the logistic regression is a predictive analysis. Logistic regression is used to describe data and to explain the relationship between one dependent binary variable and one or more nominal, ordinal, interval or ratio-level independent variables.

## Random Forest Classifier:

The Random Forest classifier creates a set of decision trees from a randomly selected subset of the training set. It is basically a set of decision trees (DT) from a randomly selected subset of the training set and then It collects the votes from different decision trees to decide the final prediction.

## Decision tree classifier:

Decision Tree is a Supervised learning technique that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome.

# CONCLUSION

Combining multiple risk factors in modelling for breast cancer prediction could help the early diagnosis of the disease with necessary care plans. Collection, storage, and management of different data and intelligent systems based on multiple factors for predicting breast cancer are effective in disease management.

This particular model considers value from a dataset representing various values and attributes of potential patients. Upon being trained under three different algorithms , the models proves itself to be pretty effective and accurate . The accuracy obtained in the results are as follows

Logistic regression – 96.49%

Decision tree – 9.86%

Random forest classifier – 97.34%

The training model using random forest classifier training model turns out to be the best algorithm for this particular case , coming up with 97.34% accuracy rate.

# BIBLIOGRAPHY

- https://www.javatpoint.com/machine-learning

- https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9175124

- https://www.kaggle.com/datasets/uciml/breast-cancer-wisconsin-data

- https://www.youtube.com/watch?v=HXnDyrraRb0&t=789s&ab_channel=CodeForLife

## PROJECT GITHUB LINK

https://github.com/amanpandey-03/BreastCancerPrediction