

Breast Cancer Risk Prediction System

1.Introduction

Breast cancer (BC) is one of the most common cancers among women worldwide, representing the majority of new cancer cases and cancer-related deaths according to global statistics, making it a significant public health problem in today's society.

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. Thus, the correct diagnosis of BC and classification of patients into malignant or benign groups is the subject of much research. Because of its unique advantages in critical features detection from complex BC datasets, machine learning (ML) is widely recognized as the methodology of choice in BC pattern classification and forecast modelling.

1.1 Overview

In the conventional way of diagnosing breast cancer some tests and procedures are carried out. These tests include Breast exam Mammogram Breast ultrasound Biopsy. As an alternative we can also use Machine Learning techniques for the classification of benign and malignant tumors. The prior diagnosis of Breast Cancer can enhance the prediction and survival rate notably, so that patients can be informed to take clinical treatment at the right time. Classification of benign tumors can help the patients avoid undertaking needless treatments. Thus the research is to be carried for the proper diagnosis of Breast Cancer and categorization of patients into malignant and benign groups. Machine Learning, with its advancements in detection of critical features from the complex datasets is largely acknowledged as the method in the prediction of breast cancer. Application of data mining techniques in the medical field can help in prediction of outcomes, minimizing the cost of medicines, aid people's health, upgrade the healthcare value and to rescue lives of people. This process of classifying benign and malignant tumors can be best done by the application of Classification techniques of machine learning. Lot of research is being conducted in this area by the application of various machine learning and data mining techniques for many different datasets on Breast Cancer. Most of them show that classification techniques give a good accuracy in prediction of the type of tumor.

1.2 PURPOSE

Breast cancer is one of the main causes of cancer death worldwide. Early diagnostics significantly increases the chances of correct treatment and survival, but this process is tedious and often leads to a disagreement between pathologists. Computer-aided diagnosis systems showed the potential for improving diagnostic accuracy. But early detection and prevention can significantly reduce the chances of death. It is important to detect breast cancer as early as possible.

We are building a model in Watson Studio and deploying the model in IBM Watson Machine Learning. To interact with the model we will be using Node-Red and scoring Endpoint.

2.LITERATURE SURVEY

2.1 Existing problem

Breast cancer is the most common malignancy among women, accounting for nearly 1 in 3 cancers diagnosed among women in the United States, and it is the second leading cause of cancer death among women. Breast Cancer occurs because of abnormal growth of cells in the breast tissue, commonly referred to as a Tumor. A tumor does not mean cancer - tumors can be benign (not cancerous), pre-malignant (pre-cancerous), or malignant (cancerous). Tests such as MRI, mammogram, ultrasound, and biopsy are commonly used to diagnose breast cancer performed.

2.1 PROPOSED SOLUTION

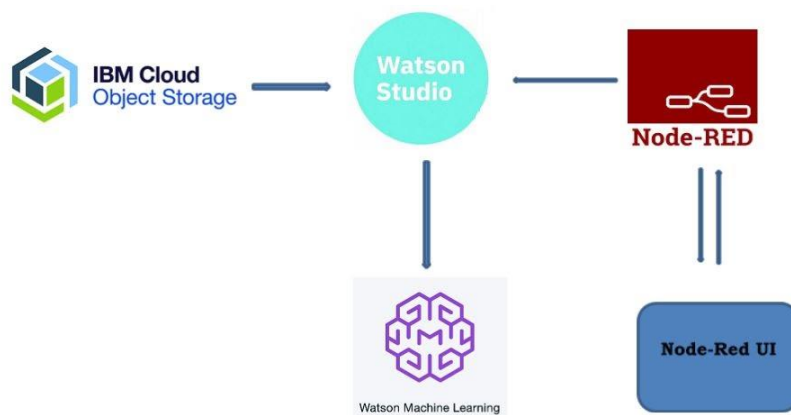
Building a model in Watson Studio and deploying the model in IBM Watson Machine Learning. To interact with the model, I will be using Node-Red and scoring Endpoint.

Solution Requirements:

Develop a model that is capable of detecting the Breast Cancer in early stages. The Machine learning model is trained and deployed on IBM Watson Studio and an endpoint is created. The web application is built using IBM Node-Red.

3.THEORITICAL ANALYSIS

3.1 BLOCK DIAGRAM



3.2 Hardware / Software designing

IBM Node red, IBM Watson Studio, IBM Machine Learning, IBM Cloud Object Storage

3.3 SOLUTION REQUIREMENTS

Develop a model that is capable of detecting the Breast Cancer in early stages. The Machine learning model is trained and deployed on IBM Watson Studio and an endpoint is created. The web application is built using IBM Node-Red.

4.EXPERIMENTAL INVESTIGATION

The dataset in the given problem consist of 32 columns and 569 rows. Totally there are 569 records with 30 featured values.

The prediction was to make whether according to the data from the previous data, is the new patient likely to have breast cancer or not.

Diagnosis column has the output as M and B, where M represents patient with the cancer, whereas B represents no cancer. During analysis of the dataset, and studying the heat map of the correlation between different features, I came to a conclusion that no part of the dataset should not be ignored/deleted.

There is no categorical data in the dataset. Batch of the data is bisected into 80:20 ratio of training test: test data.

5. RESULTS:

The screenshot displays the IBM Watson Studio web interface. The top navigation bar shows the user's account and various service links. The main workspace contains a Jupyter Notebook with the following code and output:

```
X_train_scaled = scaler.transform(X_train)
model = SVC(C=2.0, kernel='rbf')
start = time.time()
model.fit(X_train_scaled, Y_train)
end = time.time()
print("Run Time: %f" % (end-start))
```

Run Time: 0.006097

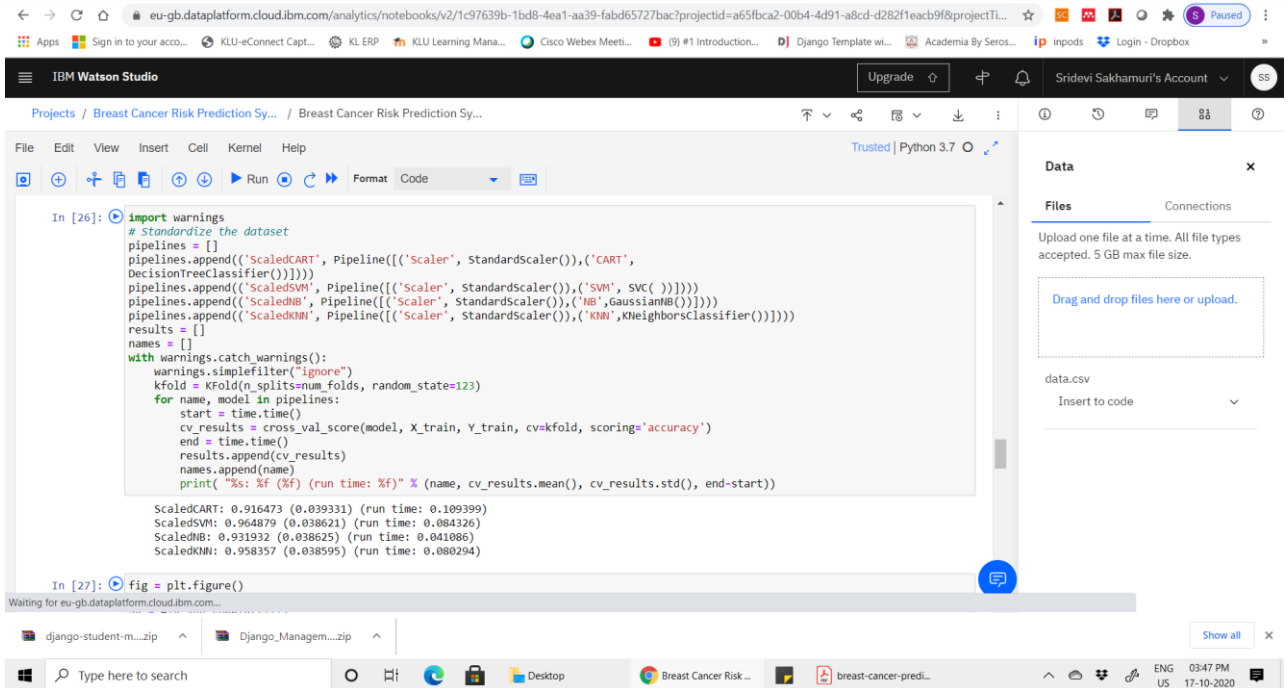
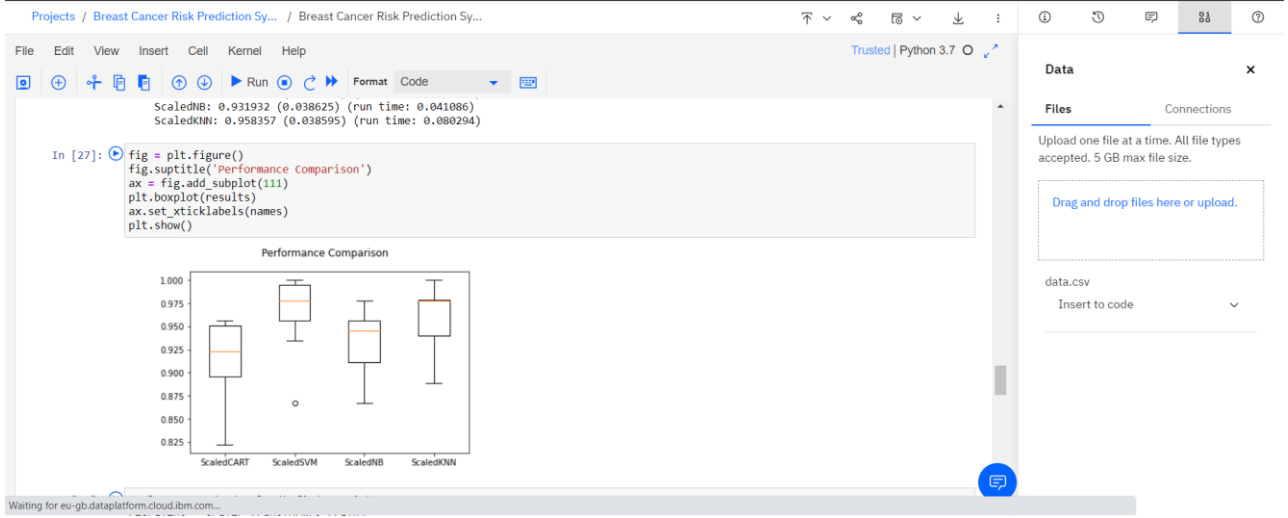
```
In [32]: with warnings.catch_warnings():
         warnings.simplefilter("ignore")
         X_test_scaled = scaler.transform(X_test)
         predictions = model.predict(X_test_scaled)

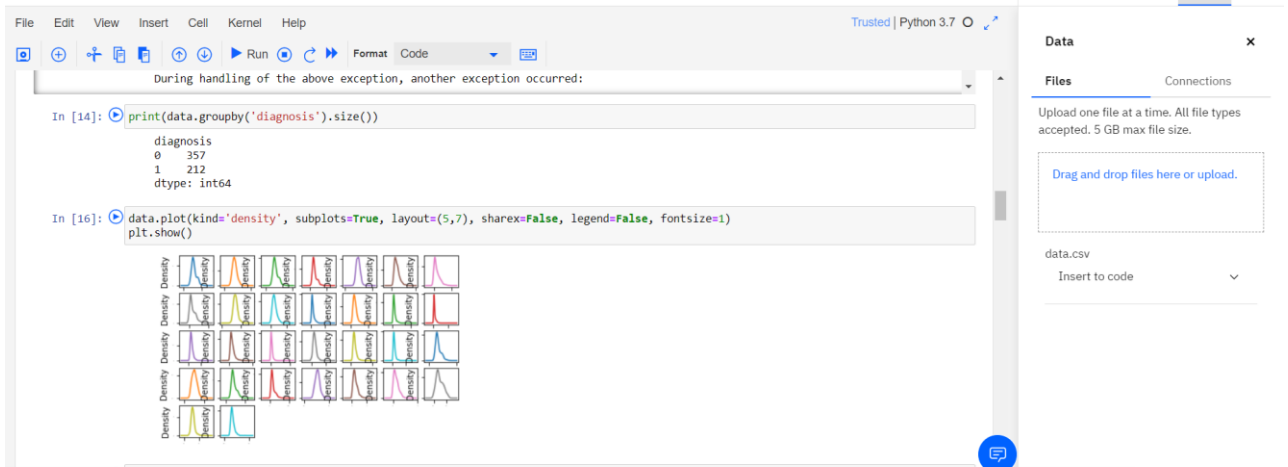
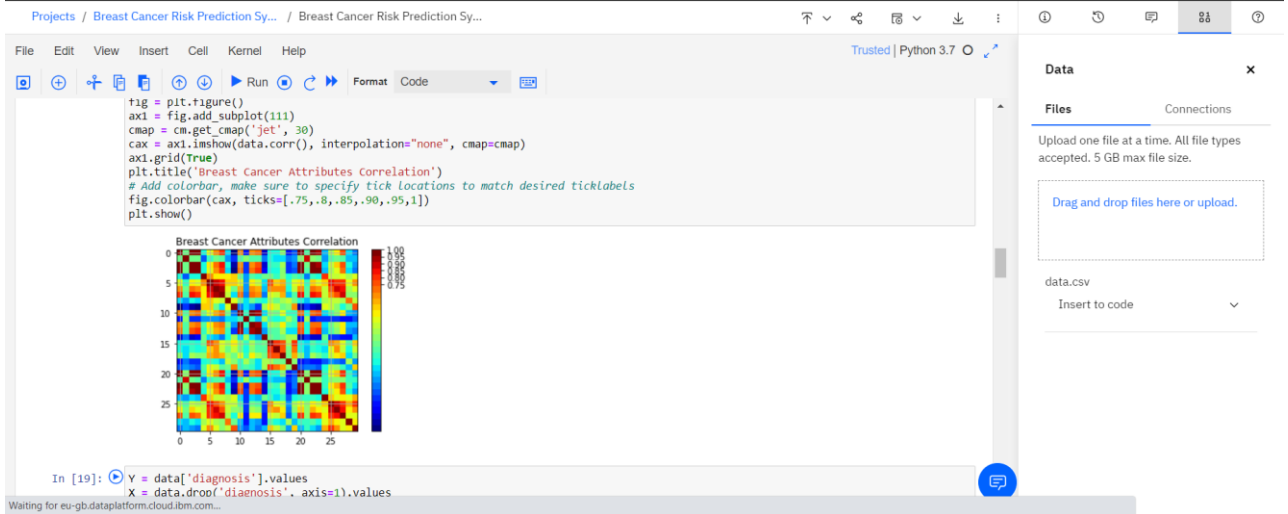
In [33]: print("Accuracy score %f" % accuracy_score(Y_test, predictions))
         print(classification_report(Y_test, predictions))
```

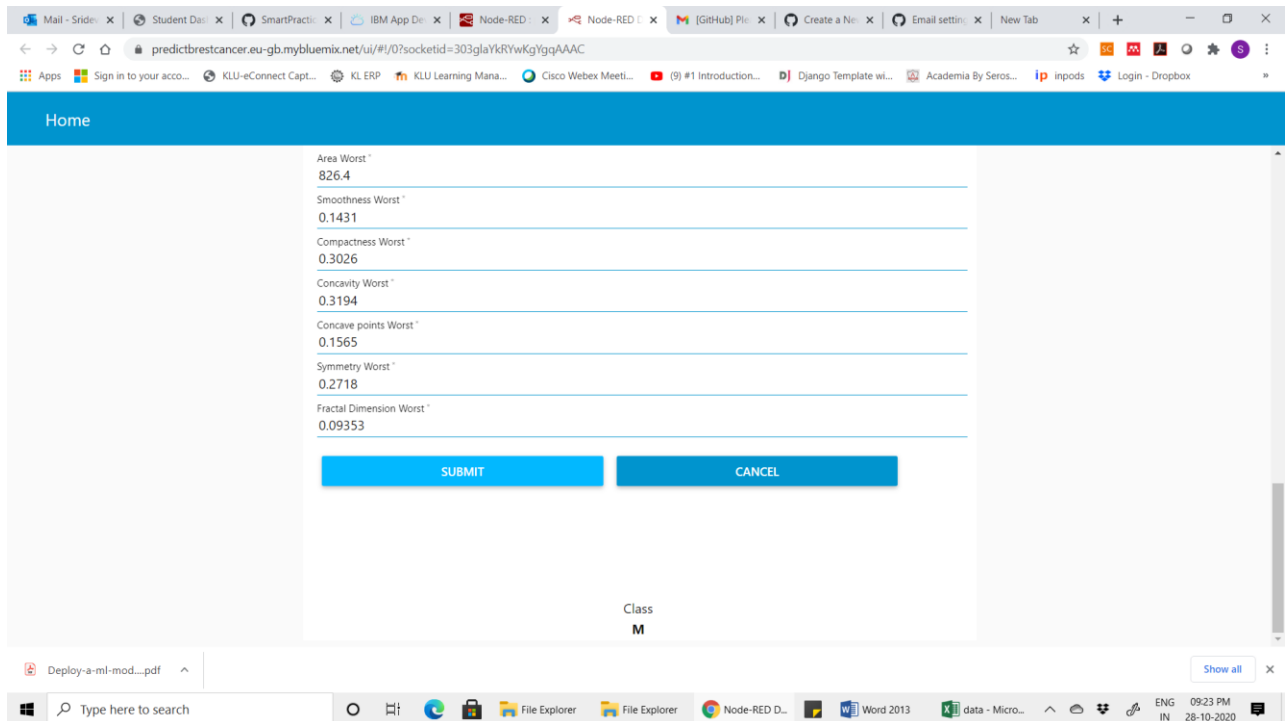
Accuracy score 0.991228

	precision	recall	f1-score	support
0	1.00	0.99	0.99	75
1	0.97	1.00	0.99	39
accuracy			0.99	114
macro avg	0.99	0.99	0.99	114
weighted avg	0.99	0.99	0.99	114

The interface also shows a 'Data' panel on the right with a file upload section and a 'Data.csv' file listed. The bottom status bar indicates the user is waiting for a connection to the IBM Cloud platform.







6. ADVANTAGES & DISADVANTAGES

ADVANTAGES:

1. Both the model and application is lightweight.
2. Prediction speed is high.
3. Server side is authenticated.
4. The prediction is helpful in educational purposes.

DISADVANTAGES:

1. Node-red is not suitable for commercial purposes.
2. Predictions on missing feature can be inaccurate.

7. APPLICATIONS

1. Through this model we can predict, whether a patient is likely to have cancer or not, without even doing medical tests.
2. Medical test can be modified and optimized.
3. We can analyze which starts of population (in women) are more likely to have it in the future.
4. The application can be run on android by SL4A.
5. It would be very useful and handy tool in healthcare.
6. It can run on PC server very fast.
7. It bypasses the first level of manual inspection.

8. CONCLUSION

Because of this ongoing covid-19 pandemic, now people have an urge to have pre knowledge of all medical ailments and advancements. For better working of the model we would be needed actual and large dataset. Since by the level of dataset in the repository, the results are "good". In the source code I have 8 types of models, but on the IBM cloud, Node-red I have deployed it by using only LGBM classification algorithm. The model was trained on a dataset of 569 patients, the total number of features were 30. Feature scaling was very important in this problem set, as the classification algorithm used demands uniform distancing between all features. The model after initial testing was deployed IBM Watson and NODE-RED.

This intermediate level of machine learning is very necessary to understand to leave scope for potential developments.

9. FUTURE SCOPE

As per WHO breast cancer is a deadly cancer, which develops inside the human body without even showing symptoms. Since in the initial days of the disease, no symptoms are witnessed by the patient, the diseases develops into later stages of the deadly cancer .If by regular examination ,we can deploy the dataset into the model to predict ,to very good accuracy we can find about the cancer ,without any medical tests. That is why, this prediction algorithm has great future ahead, if it keeps on learning from bigger dataset. `

