

# Project :- Breast Cancer Risk Prediction using IBM Auto AI

## 1. Introduction -

### 1.1 Overview -

Breast Cancer is the most leading malignancy affecting 2.1 million women each year which leads to greatest number of deaths among women. Early treatment not only helps to cure cancer but also helps in prevention of its recurrence. And hence this system mainly focuses on prediction of breast cancer where it uses different machine learning algorithms for creating models like decision tree, logistic regression, random forest which are applied on pre-processed data which suspects greater accuracy for prediction. Amongst all the models, Random Forest Classification leads to best accuracy with 98.6%. These techniques are coded in python and uses numpy, pandas, seaborn libraries.

### 1.2 Purpose -

According to World health organization, Breast cancer is the most frequent cancer among women and it is the second dangerous cancer after lung cancer. In 2018, from the research it is estimated that total 627,000 women lost their life due to breast cancer that is 15% of all cancer deaths among women. In case of any symptom, people visit to oncologist. Doctors can easily identify breast cancer by using Breast ultrasound, Diagnostic mammogram, Magnetic resonance imaging (MRI), Biopsy. Based on these test results, doctor may recommend further tests or therapy. Early detection is very crucial in breast cancer. If

chances of cancer are predicted at early stage then survivability chances of patient may increase. An alternate way to identify breast cancer is using machine learning algorithms for prediction of abnormal tumor. Thus, the research is carried out for the proper diagnosis and categorization of patients into malignant and benign groups.

## 2. Literature Survey -

### 2.1 Existing Problem -

Over the past decades, a continuous evolution related to cancer research has been performed. Scientists applied different methods, such as screening in early stage, in order to find types of cancer before they cause symptoms. Moreover, they have developed new strategies for the early prediction of cancer treatment outcome. With the advent of new technologies in the field of medicine, large amounts of cancer data have been collected and are available to the medical research community. However, the accurate prediction of a disease outcome is one of the most interesting and challenging tasks for physicians.

As a result, ML methods have become a popular tool for medical researchers. These techniques can discover and identify patterns and relationships between them, from complex datasets, while they are able to effectively predict future outcomes of a cancer type.

### 2.2 Proposed Solution -

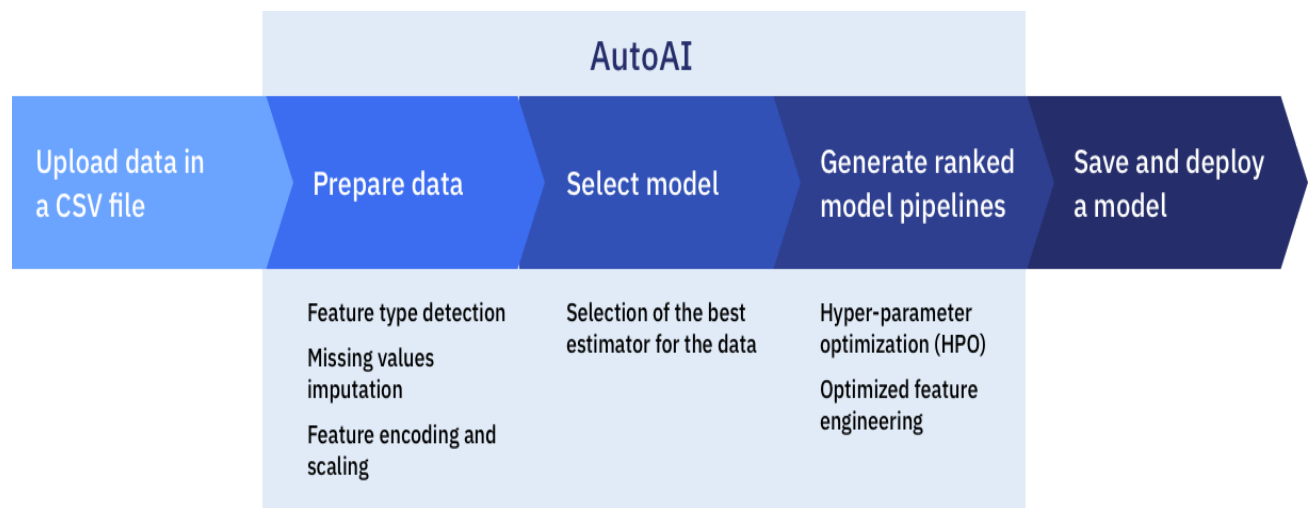
We obtained the breast cancer dataset of Wisconsin Breast Cancer diagnosis dataset and used IBM Auto AI as the platform for the purpose of development of model. The model is been trained using Auto AI service in IBM

watson cloud and that can be deployed in an application such as web or mobile application using Node-RED application.

The AutoAI graphical tool in Watson Studio automatically analyzes your data and generates candidate model pipelines customized for your predictive modeling problem. These model pipelines are created over time as AutoAI analyzes your dataset and discovers data transformations, algorithms, and parameter settings that work best for your problem setting. Results are displayed on a leaderboard, showing the automatically generated model pipelines ranked according to your problem optimization objective.

### 3. THEORITICAL ANALYSIS -

#### 3.1 Block Diagram -



#### 3.2 Hardware/ Software Designing -

The AutoAI process follows this sequence to build candidate pipelines:

- [Data pre-processing](#)
- [Automated model selection](#)
- [Automated feature engineering](#)
- [Hyperparameter optimization](#)

➤ Data pre-processing -

Most data sets contain different data formats and missing values, but standard machine learning algorithms work with numbers and no missing values. AutoAI applies various algorithms, or estimators, to analyze, clean, and prepare your raw data for machine learning. It automatically detects and categorizes features based on data type, such as categorical or numerical. Depending on the categorization, it uses hyper-parameter optimization to determine the best combination of strategies for missing value imputation, feature encoding, and feature scaling for your data.

➤ Automated model selection -

The next step is automated model selection that matches your data. AutoAI uses a novel approach that enables testing and ranking candidate algorithms against small subsets of the data, gradually increasing the size of the subset for the most promising algorithms to arrive at the best match. This approach saves time without sacrificing performance. It enables ranking a large number of candidate algorithms and selecting the best match for the data.

➤ Automated feature engineering -

Feature engineering attempts to transform the raw data into the combination of features that best represents the problem to achieve the most accurate prediction. AutoAI uses a novel approach that explores various feature construction choices in a structured, non-exhaustive manner, while progressively maximizing model accuracy using reinforcement learning. This results in an optimized sequence of transformations for the data that best match the algorithms of the model selection step.

### ➤ Hyperparameter optimization -

Finally, a hyper-parameter optimization step refines the best performing model pipelines. AutoAI uses a novel hyper-parameter optimization algorithm optimized for costly function evaluations such as model training and scoring that are typical in machine learning. This approach enables fast convergence to a good solution despite long evaluation times of each iteration.

### Connecting to Node-Red -

- Create a Node-RED starter application running in IBM Cloud, create machine learning instances of the Watson services, and connect the services to your Node-Red app.
- Launch and configure the Node-RED visual programming editor.
- Install additional Node-RED nodes and create flows that use the Watson services to create the Breast cancer prediction model api.

## 4. Experimental Investigation -

The Breast Cancer (Wisconsin) Diagnosis dataset contains the diagnosis and a set of 30 features describing the characteristics of the cell nuclei present in the digitized image of a of a fine needle aspirate (FNA) of a breast mass.

Eight real-valued features are computed for each cell nucleus:

- radius (mean of distances from center to points on the perimeter);
- texture (standard deviation of gray-scale values);
- perimeter;
- area;
- smoothness (local variation in radius lengths);
- compactness ( $\text{perimeter}^2 / \text{area} - 1.0$ );
- concavity (severity of concave portions of the contour);
- concave points (number of concave portions of the contour);

The mean, standard error (SE) and “worst” or largest (mean of the three largest values) of these features were computed for each image, resulting in

The screenshot shows a Microsoft Excel spreadsheet titled "breast\_cancer\_prediction\_dataset - Microsoft Excel (Product Activation Failed)". The spreadsheet is open to a worksheet with columns labeled A1 through AA and rows numbered 1 through 28. The data is organized into a table with the following columns:

- diagnosis (A1)
- radius\_mean (B1)
- radius\_se (C1)
- radius\_logp (D1)
- perimeter\_mean (E1)
- perimeter\_se (F1)
- perimeter\_logp (G1)
- area\_mean (H1)
- area\_se (I1)
- area\_logp (J1)
- convex\_mean (K1)
- convex\_se (L1)
- convex\_logp (M1)
- concave\_mean (N1)
- concave\_se (O1)
- concave\_logp (P1)
- symmetry\_mean (Q1)
- symmetry\_se (R1)
- symmetry\_logp (S1)
- texture\_mean (T1)
- texture\_se (U1)
- texture\_logp (V1)
- redness\_mean (W1)
- redness\_se (X1)
- redness\_logp (Y1)
- convexity\_mean (Z1)
- convexity\_se (AA1)
- convexity\_logp (AB1)

The data is organized into rows, with the first row (A1) serving as a header for the "diagnosis" column. The subsequent rows (A2 through AA) contain numerical values for each feature, with some cells highlighted in green. The spreadsheet interface includes standard Excel menus (File, Home, Insert, Page Layout, Formulas, Review, View) and toolbars for formatting and data manipulation.

The basic steps for building and training a machine learning model using AutoAI:

- ## 5. FlowChart -

Download the sample training data file

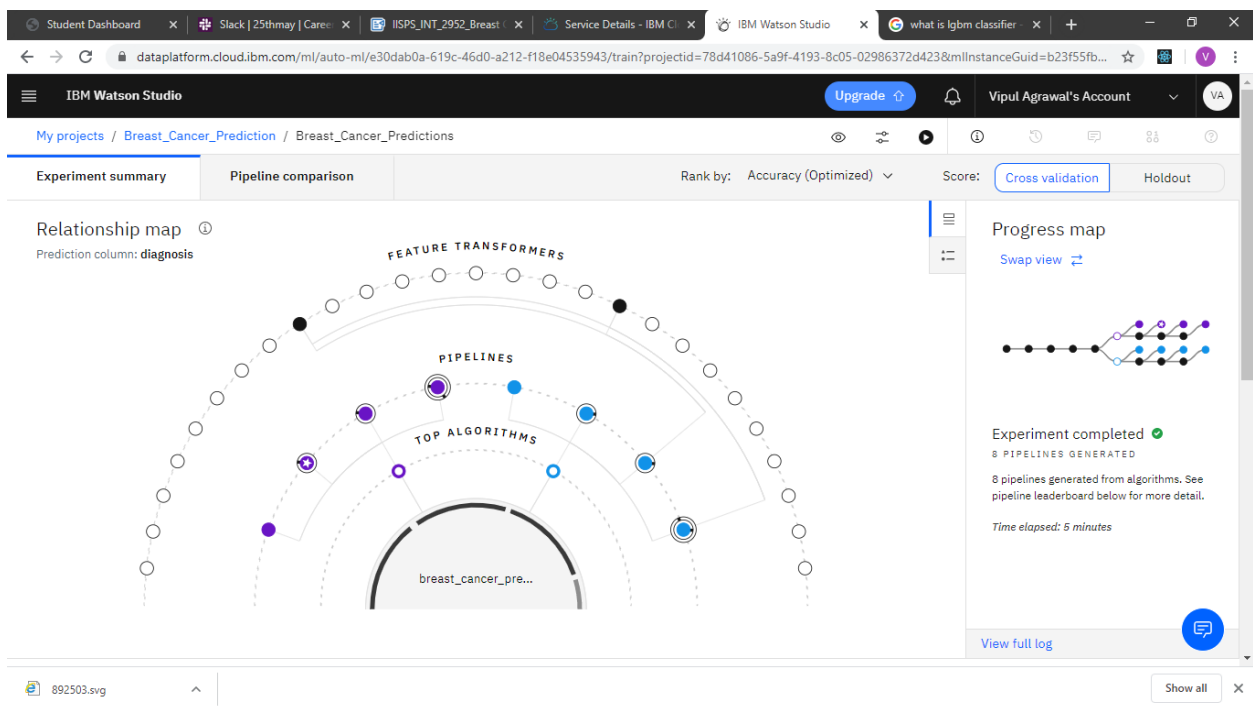
Build and train the experiment

Deploy the trained model

Test the deployed model

## 6. Result -

We Successfully obtained the pipeline with an accuracy score of 0.965. As you can see from images below



Student Dashboard | Slack | 25thmay | Career | IISPS\_INT\_2952\_Breast | Service Details - IBM Cl | IBM Watson Studio | what is lgbm classifier

dataplatform.cloud.ibm.com/ml/auto-ml/e30dab0a-619c-46d0-a212-f18e04535943/train?projectId=78d41086-5a9f-4193-8c05-02986372d423&mlInstanceGuid=b23f55fb...

IBM Watson Studio

My projects / Breast\_Cancer\_Prediction / Breast\_Cancer\_Predictions

Experiment summary | Pipeline comparison | Rank by: Accuracy (Optimized) | Score: Cross validation | Holdout

Pipeline leaderboard

Rank	↑	Name	Algorithm	Accuracy (Optimized)	Enhancements	Build time
>	★ 1	Pipeline 2	LGBM Classifier	0.965	HPO-1	00:00:16
>	2	Pipeline 3	LGBM Classifier	0.965	HPO-1 FE	00:01:05
>	3	Pipeline 4	LGBM Classifier	0.965	HPO-1 FE HPO-2	00:00:41
>	4	Pipeline 8	Gradient Boosting Classifier	0.961	HPO-1 FE HPO-2	00:00:22
>	5	Pipeline 7	Gradient Boosting Classifier	0.955	HPO-1 FE	00:00:45
>	6	Pipeline 1	LGBM Classifier	0.953	None	00:00:01
>	7	Pipeline 6	Gradient Boosting Classifier	0.951	HPO-1	00:00:06

892503.svg

Show all

Output :-

Successfully able to diagnose the patient into benign or malignant.

Student Dashboard | Slack | 25thmay | Career | IISPS\_INT\_2952\_Breast | Service Details - IBM Cl | IBM Watson Studio | what is lgbm classifier

dataplatform.cloud.ibm.com/ml/deployments/56e3fb29-5a5c-4fbb-b655-1eae6b832099/test?projectId=78d41086-5a9f-4193-8c05-02986372d423&mlInstanceGuid=b23f...

IBM Watson Studio

My projects / Breast\_Cancer\_Prediction / Breast\_Cancer\_Predictions - P2 L... / Breast\_Cancer

Breast\_Cancer

Overview | Implementation | Test

Enter input data

0.1622

compactness\_worst

0.6656

concavity\_worst

0.7119

concave points\_worst

0.2654

```

{
  "predictions": [
    {
      "fields": [
        "prediction",
        "probability"
      ],
      "values": [
        "M",
        [
          0.01777776537481579,
          0.9822222346251842
        ]
      ]
    }
  ]
}

```

892503.svg

Show all



## Creation of web application using node-red service from ibm starter kit -

node-red-awnmw.eu-gb.mybluemix.net/ui/#/0?socketid=xp9X4vMdSThaeaWQAAQ

Home

### Breast\_Cancer\_Prediction

radius_mean ^	0.05363
texture_mean ^	0.011587
perimeter_mean ^	23.45
area_mean ^	455.45
smoothness_mean ^	564
compactness_mean ^	65.706
concavity_mean ^	5656.5
concave points_mean ^	655.6
radius_se ^	5605656
texture_se ^	0.56
perimeter_se ^	45540
area_se ^	45540

node-red-awnmw.eu-gb.mybluemix.net/ui/#/0?socketid=xp9X4vMdSThaeaWQAAQ

Home

0.5466566	
radius_worst ^	54656.45
texture_worst ^	4556
perimeter_worst ^	565
area_worst ^	0.54554
smoothness_worst ^	0.1622
compactness_worst ^	0.6656
concavity_worst ^	0.7119
concave points_worst ^	0.2654

Diagnosis **M**

## 7. Advantages And Disadvantages -

### Advantages-

1. Easily identifies trends and patterns
2. No human intervention needed (automation)
3. Continuous Improvement
4. Handling multi-dimensional and multi-variety data
5. Wide Applications

### Disadvantages-

1. Data Acquisition
2. Time and Resources
3. Interpretation of Results
4. High error-susceptibility

## 8. Applications -

### Application of Breast Cancer Risk Prediction Models in Clinical Practice-

Breast cancer risk assessment provides an estimation of disease risk that can be used to guide management for women at all levels of risk. In addition, the likelihood that breast cancer risk is due to specific genetic susceptibility (such as BRCA1 or BRCA2 mutations) can be determined. Recent developments have reinforced the clinical importance of breast cancer risk assessment. Tamoxifen chemoprevention as well as prevention studies such as the Study of Tamoxifen and Raloxifene are available to women at increased risk of developing breast cancer. In addition, specific management strategies are now defined for BRCA1 and BRCA2 mutation carriers. Risk may be assessed as the likelihood of developing breast cancer (using risk assessment models) or as the likelihood of detecting a BRCA1 or BRCA2 mutation (using prior probability models). Each of the models has advantages and disadvantages, and all need to be interpreted in context. We review available risk assessment tools and discuss their application. As illustrated by clinical examples, optimal counseling may require the use of several models, as well as clinical judgment, to provide the most accurate and useful information to women and their families.

## 9. Conclusion -

1. Breast cancer if found at an early stage will help save lives of thousands of women or even men. This project will help the real world patients and doctors to gather as much information as they can. The research on nine papers has helped us gather the data for the project proposed by us. By using machine learning algorithms we will be able to classify and predict the cancer into benign or malignant.
2. Machine learning algorithms can be used for medical oriented research, it advances the system, reduces human errors and lowers manual mistakes.

## 10. Future Scope -

AI is set to change the medical industry in the coming decades — it wouldn't make sense for pathology to not be disrupted too.

Currently, ML models are still in the testing and experimentation phase for cancer prognoses. As datasets are getting larger and of higher quality, researchers are building increasingly accurate models.

Here's what a future cancer biopsy might look like:

You perform clinical tests, either at a clinic or at home. Data is inputted into a pathological ML system. A few minutes later, you receive an email with a detailed report that has an accurate prediction about the development of your cancer.

While you might not see AI doing the job of a pathologist today, you can expect ML to replace your local pathologist in the coming decades, and it's pretty exciting!

ML models still have a long way to go, most models still lack sufficient data and suffer from bias. Yet, something we are certain of is that ML is the next step of pathology, and it will disrupt the industry.

## 10. Bibliography/Appendix -

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

<https://www.ijert.org/breast-cancer-classification-and-prediction-using-machine-learning#:~:text=Breast%20Cancer%20Prediction%20Using%20Genetic,detection%20of%20Breast%20Cancer%20Prediction.>

[https://www.researchgate.net/publication/337275322\\_Breast\\_Cancer\\_Prediction\\_using\\_Supervised\\_Machine\\_Learning\\_Algorithms](https://www.researchgate.net/publication/337275322_Breast_Cancer_Prediction_using_Supervised_Machine_Learning_Algorithms)

[https://dataplatform.cloud.ibm.com/docs/content/wsj/analyze-data/autoai\\_example\\_binary\\_classifier.html](https://dataplatform.cloud.ibm.com/docs/content/wsj/analyze-data/autoai_example_binary_classifier.html)