# 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 g 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 relattionships between them, from complex datasets, while they are able to effectively predict future outcomes of a cancer
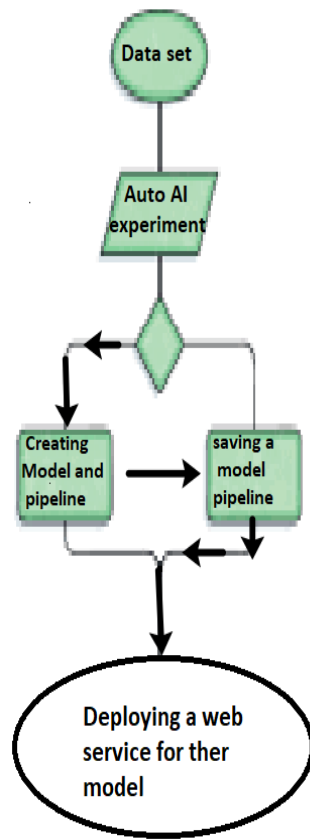
type.

 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 IBMwatson 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.
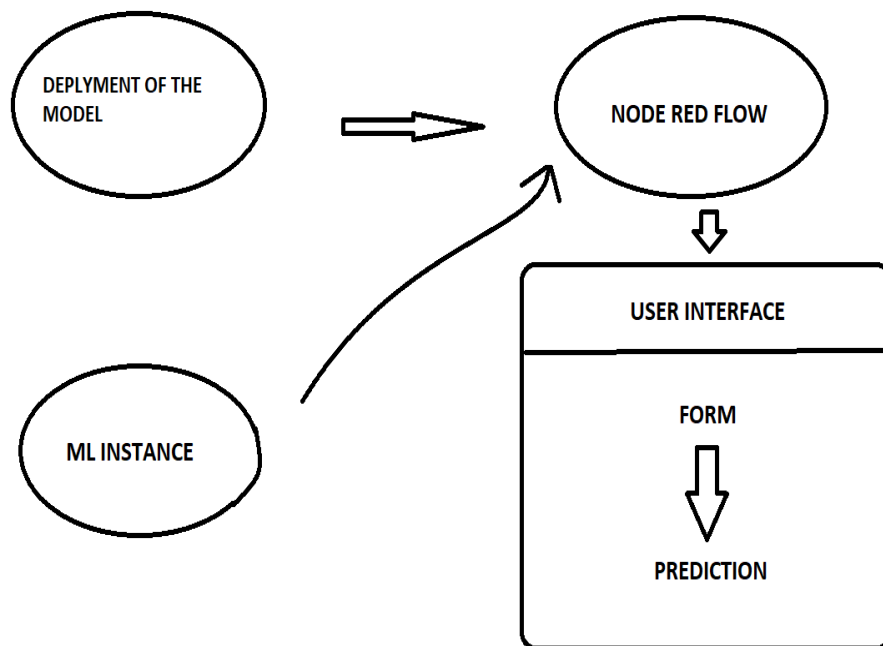
# 3.THEORETICAL ANALYSIS
3.1 Block Diagram

## 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 ● 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 additionalcandidate 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 -  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 (perimeter^2 / 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 24 features. We will analyze the features to understand the predictive value for diagnosis. We will then create models using IBM Auto AI and use the models to predict the diagnosis.
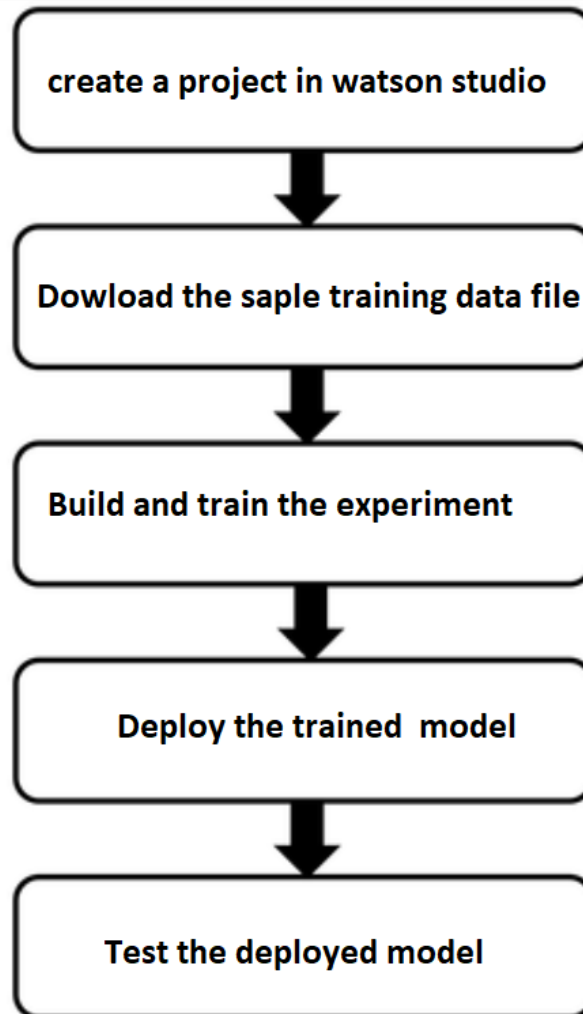


The basic steps for building and training a machine learning

model using AutoAI:

*1. Build and train the model 2. Deploy the trained model 3. Test the deployed model*

## 5. FlowChart -

```
┌─────────────────────────────────┐
│   create a project in watson studio  │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│   Dowload the saple 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

IBM Watson Studio

Upgrade ⇧     Mitali Sakle's Account     MS

My projects / sample1 / breast cancer

Experiment summary     Pipeline comparison

Rank by: Accuracy (Optimized) ⌄     Score: Cross validation     Holdout

## Relationship map ⓘ

Prediction column: **diagnosis**

FEATURE TRANSFORMERS

PIPELINES

TOP ALGORITHMS

- Pipeline 7
Accuracy: ---

100%
of training data
dataset.csv

## Progress map

Swap view ⇄

### Feature engineering

GRADIENT BOOSTING CLASSIFIER

Completed feature engineering for pipeline P7

*Time elapsed: 3 minutes*

View full log

## Pipeline leaderboard

---

IBM Watson Studio

Upgrade ⇧     Mitali Sakle's Account     MS

My projects / sample1 / breast cancer - P3 LGBMClassifie... / deployment

## deployment

Overview     Implementation     **Test**

### Enter input data

56.82

perimeter_worst

123.7

area_worst

2020

smoothness_worst

0.5466

compactness_worst

Predict

```
{
  "predictions": [
    {
      "fields": [
        "prediction",
        "probability"
      ],
      "values": [
        [
          "M",
          [
            0.0051345194708120045,
            0.994865480529188
          ]
        ]
      ]
    }
  ]
}
```

## 7. Advantages And Disadvantages

Advatages

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 PracticeBreast 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 mistake

## 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

## 11. Biblography -

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

https://www.researchgate.net/publication/337275322_Breast_Cancer_Prediction_using_Supervised_Machine_Learning_Algorithms
https://www.ijert.org/breast-cancer-classification-and-prediction-using-machine-learning#:~:text=Breast%20Cancer%20Prediction%20Using%20Genetic,detection%20of%20Breast%20Cancer%20Prediction.

https://dataplatform.cloud.ibm.com/docs/content/wsj/analyze-data/autoai_example_binary_classifier.html

# 12.APPENDIX:-

## A. SOURCE CODE

[{"id":"80fc8dba.aa05a","type":"tab","label":"Flow 2","disabled":false,"info":""},{"id":"3c2e1670.3a554a","type":"ui_form","z":"80fc8dba.aa05a","name":"","label":"","group":"acfa63c0.1b9b4","order":2,"width":0,"height":0,"options":[{"label":"mean_radius","value":"mr","type":"number","required":true,"rows":null},{"label":"mean_texture","value":"mt","type":"number","required":true,"rows":null},{"label":"mean_perimeter","value":"mp","type":"number","required":true,"rows":null},{"label":"mean_area","value":"ma","type":"number","required":true,"rows":null},{"label":"mean_smoothness","value":"ms","type":"number","required":true,"rows":null},{"label":"state1","value":"st1","type":"number","required":true,"rows":null},{"label":"state2","value":"st2","type":"text","required":true,"rows":null}],"formValue":{"mr":"","mt":"","mp":"","ma":"","ms":"","st1":"","st2":""},"payload":"","submit":"submit","cancel":"cancel","topic":"","x":150,"y":1060,"wires":[["217d2375.57153c"]]},{"id":"217d2375.57153c","type":"function","z":"80fc8dba.aa05a","name":"PreToken","func":"global.set(\"mr\",msg.payload.mr)\nglobal.set(\"mt\",msg.payload.mt)\nglobal.set(\"mp\",msg.payload.mp)\nglobal.set(\"ma\",msg.payload.ma)\nglobal.set(\"ms\",msg.payload.ms)\nglobal.set(\"st1\",msg.payload.st1)\nglobal.set(\"st2\",msg.payload.st2)\nvar apikey=\"Q0dqLIDyWxUgR4pFo3lxLPOmtNfmJ-ElCCVoMrBd80__\";\nmsg.headers={\"content-type\":\"application/x-www-form-urlencoded\"}\nmsg.payload={\"grant_type\":\"urn:ibm:params:oauth:grant-type:apikey\",\"apikey\":apikey}\nreturn msg;\n","outputs":1,"noerr":0,"x":220,"y":900,"wires":[["ed5d39e0.2ea778"]]},{"id":"ed5d39e0.2ea778","type":"http request","z":"80fc8dba.aa05a","name":"","method":"POST","ret":"obj","paytoqs":false,"url":"https://iam.cloud.ibm.com/identity/token","tls":"","persist":false,"proxy":"","authType":"","x":390,"y":1100,"wires":[["e80bd4bd.9b0628"]]},{"id":"e80bd4bd.9b0628","type":"function","z":"80fc8dba.aa05a","name":"Pre Prediction","func":"var mr=global.get(\"mr\")\nvar mt=global.get(\"mt\")\nvar mp=global.get(\"mp\")\nvar ma=global.get(\"ma\")\nvar ms=global.get(\"ms\")\nvar st1=global.get(\"st1\")\nvar st2=global.get(\"st2\")\nvar token=msg.payload.access_token\nvar

instance_id=\"2d51206d-d636-4a6d-8348-5c806a2cd38c\"\nmsg.headers={'Content-Type': 'application/json',\"Authorization\":\"Bearer \"+token,\"ML-Instance-ID\":instance_id}\nmsg.payload={\"input_data\": [{\"fields\": [\"state1\",\"state2\",\"mean_radius\", \"mean_texture\", \"mean_perimeter\", \"mean_area\", \"mean_smoothness\"], \"values\": [[st1,st2,mr,mt,mp,ma,ms]]}]}\nreturn msg;","outputs":1,"noerr":0,"x":540,"y":880,"wires":[["55791bce.75da34"]]},{"id":"55791bce.75da34","type":"http request","z":"80fc8dba.aa05a","name":"","method":"POST","ret":"obj","paytoqs":false,"url":"https://eu-gb.ml.cloud.ibm.com/v4/deployments/15844736-fa9b-46a6-8c54-5815c9a3bd4e/predictions","tls":"","persist":false,"proxy":"","authType":"","x":730,"y":1020,"wires":[["b66853a5.2756a","cecf2c30.c5022"]]},{"id":"317fa518.e8e9da","type":"ui_text","z":"80fc8dba.aa05a","group":"acfa63c0.1b9b4","order":1,"width":0,"height":0,"name":"","label":"Dignosis","format":"{{msg.payload}}","layout":"row-spread","x":1040,"y":1140,"wires":[]},{"id":"b66853a5.2756a","type":"function","z":"80fc8dba.aa05a","name":"","func":"msg.payload=msg.payload.predictions[0].values[0][1][1]\nreturn msg;","outputs":1,"noerr":0,"x":770,"y":780,"wires":[["317fa518.e8e9da","7e1c0c6b.6ffd84"]]},{"id":"7e1c0c6b.6ffd84","type":"debug","z":"80fc8dba.aa05a","name":"","active":true,"tosidebar":true,"console":false,"tostatus":false,"complete":"payload","targetType":"msg","x":1010,"y":780,"wires":[]},{"id":"cecf2c30.c5022","type":"debug","z":"80fc8dba.aa05a","name":"","active":true,"tosidebar":true,"console":false,"tostatus":false,"complete":"payload","targetType":"msg","x":1046.892333984375,"y":954.1111450195312,"wires":[]},{"id":"acfa63c0.1b9b4","type":"ui_group","z":"","name":"Breast Cancer Risk Prediction using IBM Auto AI","tab":"a76f3181.b9c51","order":1,"disp":true,"width":"6","collapse":false},{"id":"a76f3181.b9c51","type":"ui_tab","z":"","name":"Home","icon":"dashboard","disabled":false,"hidden":false}]

## B.UI OUTPUT SCREENSHOT

Diagnosis                    M

## form

mean_radius *

13.49

mean_texture *

21.38

mean_perimeter *

86.51

mean_area *

564

mean_smoothness *

0.08752

**SUBMIT**    **CANCEL**