# Breast Cancer Prediction

Group 7

# Group Members

Vrushabh Ajaybhai Desai ( VrushabhAjaybhaiDesai@my.unt.edu ) Document, Code

Varsha Salil ( VarshaSalil@my.unt.edu ) Document, Code

Kareem Baba Shaik ( kareembabashaik@my.unt.edu ) Code, Testing

Manasa Kilaru ( manasakilaru@my.unt.edu ) Dataset, Code

Alekhya Sree Jetti ( AlekhyaSreeJetti@my.unt.edu ) Document, Testing

# Motivation

- In the present world, Breast cancer is one of the major causes of death in women. Over 2 million women get diagnosed with breast cancer every year. The condition is gradually progressing and will only exhibit most of the distinguishing symptoms during more advanced stages. However, the severity of the condition can be curbed by timely prediction and proper treatment.

- The existing diagnosis approaches include physically invasive tests and chemical tests on blood samples. Body cells are then classified into Benign (or non-cancerous) cells or Malignant(cancerous) cells. Cancerous cells can also be detected by analyzing several body parameters like clump thickness, cell size, cell shape etc.

- In this project, we propose a non-invasive study and classification of a patient's body parameters of interest to classify breast cells into two categories i.e, Benign or Malignant, depending on the patient's diagnosis result. This method provides a faster, more economic approach to breast cancer diagnosis.

# Objective

- The main objective of the project is to generate an accurate system which predicts Breast Cancer by analyzing several body parameters such as clump thickness, cell size, cell shape etc.

# Milestones:

- <u>Project proposal</u>: The project proposal is the initial stage where the project idea is presented, and we started outlining the goals, objectives, and requirements of the project.

- <u>Algorithm generation</u>: Once the proposal has been accepted, we started working on the algorithm where we worked on a step-by-step procedure for solving a specific problem.

- <u>Code generation</u>: After the algorithm is created, the next milestone is the implementation of the algorithm in a programming language.

- <u>Testing the code for efficiency</u>: The code generated needs to be tested for efficiency and accuracy to ensure it performs as intended. This milestone involves running the code through different tests to check for errors and make sure it meets the requirements of the project.

- <u>Project report and presentation</u>: summarizing the entire project, including the goals, objectives, requirements, algorithm, and code. The report and presentation also include the testing results and the overall efficiency of the project.
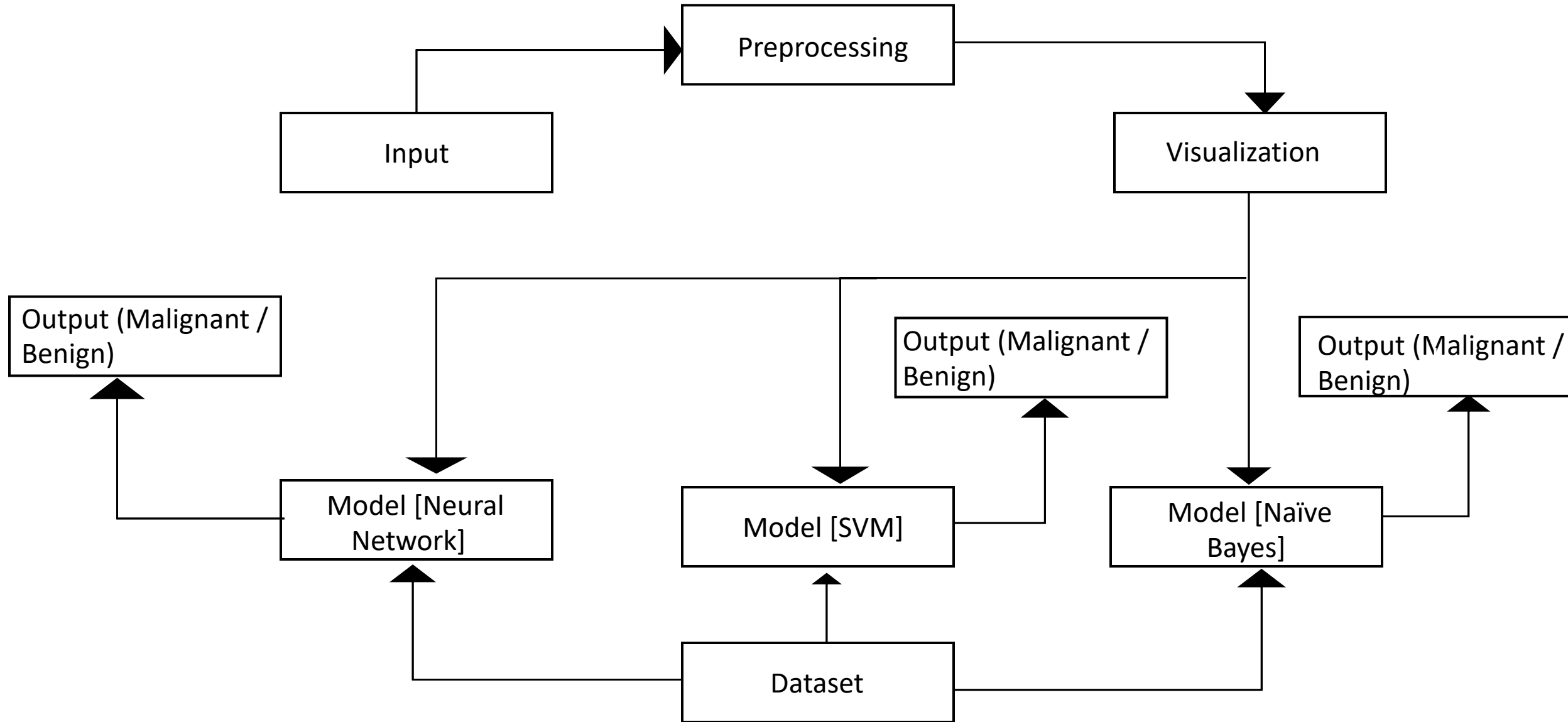
# Project Design

- The first step includes obtaining the relevant data for the diagnosis and data preprocessing for classification, by removing null values, dropping unnecessary columns etc. The following dataset will be used for training and testing purposes: https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/.

- The next step is data visualization.

- A Gaussian Naïve Bayes model is generated next from sklearn.naive_bayes package and is fit with the data for training purpose.

- The model then classifies the data into Benign or malignant and returns the predication along with the accuracy of the prediction.

- The dataset was also classified using a Support Vector Machine (SVM) model and a Neural Network model to compare and improve the accuracy of predictions.

# Data Specification

- This provided dataset is from wisconsin university. It has 31 columns. In which diagnosis column is the outpuut column which we trying to predict.

- This dataset has 357 Benign and 212 Malign data. And it has 569 rows and 33 columns in which two columns id and unnamed: 32 are not useful for the project.

- All provided columns are non-null data and every column are float values except diagnosis which accepts string values.

- We have observed that parameter_se and area_se are corelated with the feature radius_se, the same way as radius_worst with parameter_worst and area_worst and the third correlation is perimeter_mean, area_mean with radius_mean. We observed this using heatmap.

- After this, we can know that breast cancer is majorly depend on three parameters perimeter, area and radius.

- Link for the dataset: https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/

# Workflow

# Models used:

- Gaussian Naïve Bayes Model
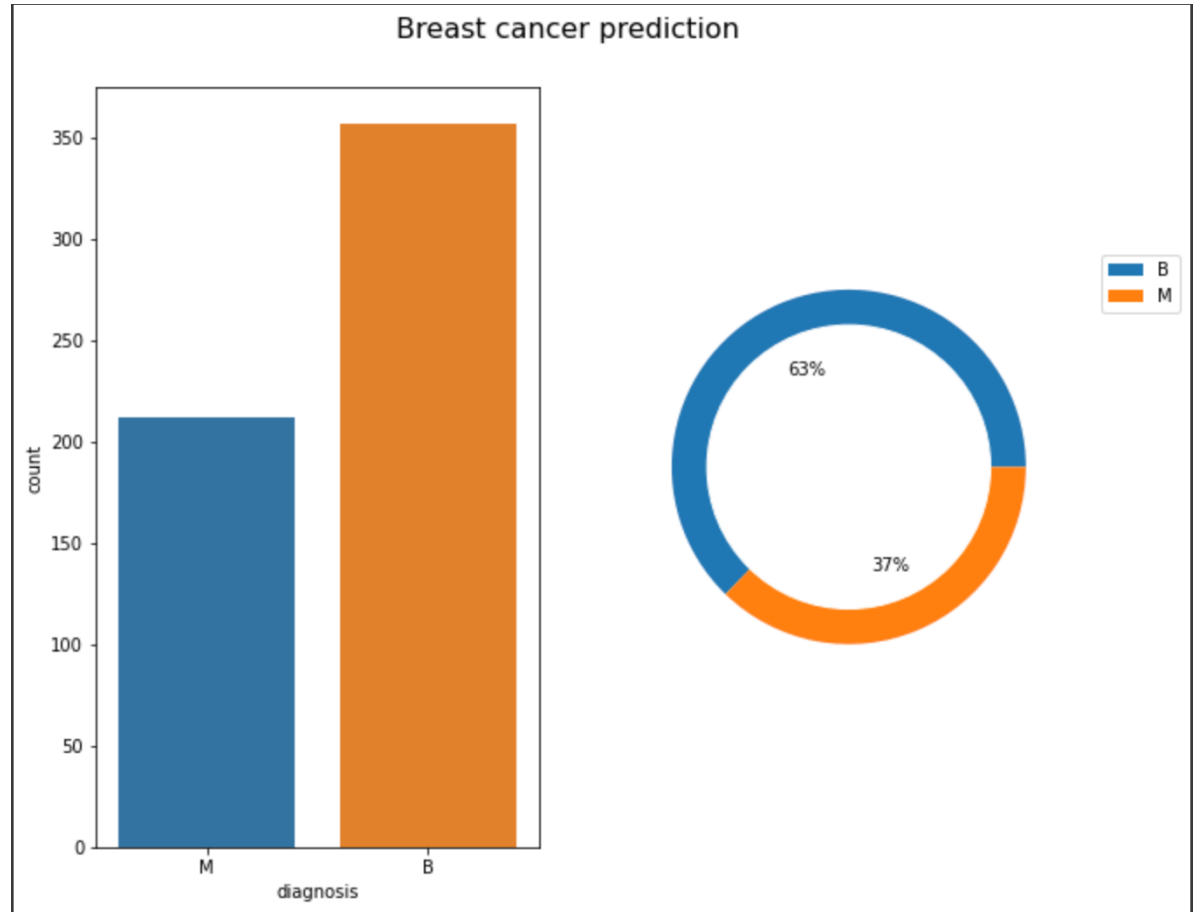- Support Vector Machine
- Neural Network

# Future Scope:

- We have implemented three different models, but we haven't used validation set. So, if for the phase 2, someone can create validation set and tweak the hyper parameters, we can get more accuracy.

- After implementing the models, we want to implement the transfer learning and see how the models work with other datasets.
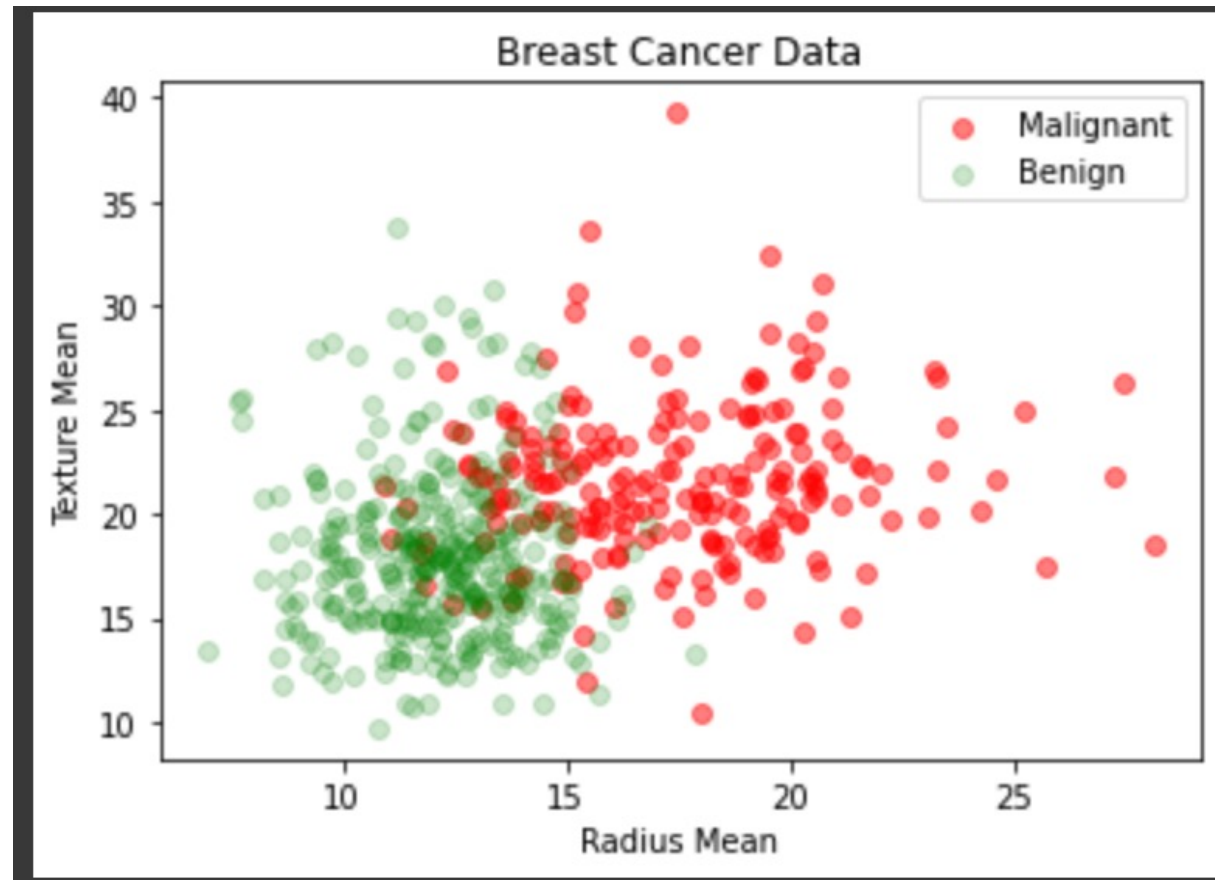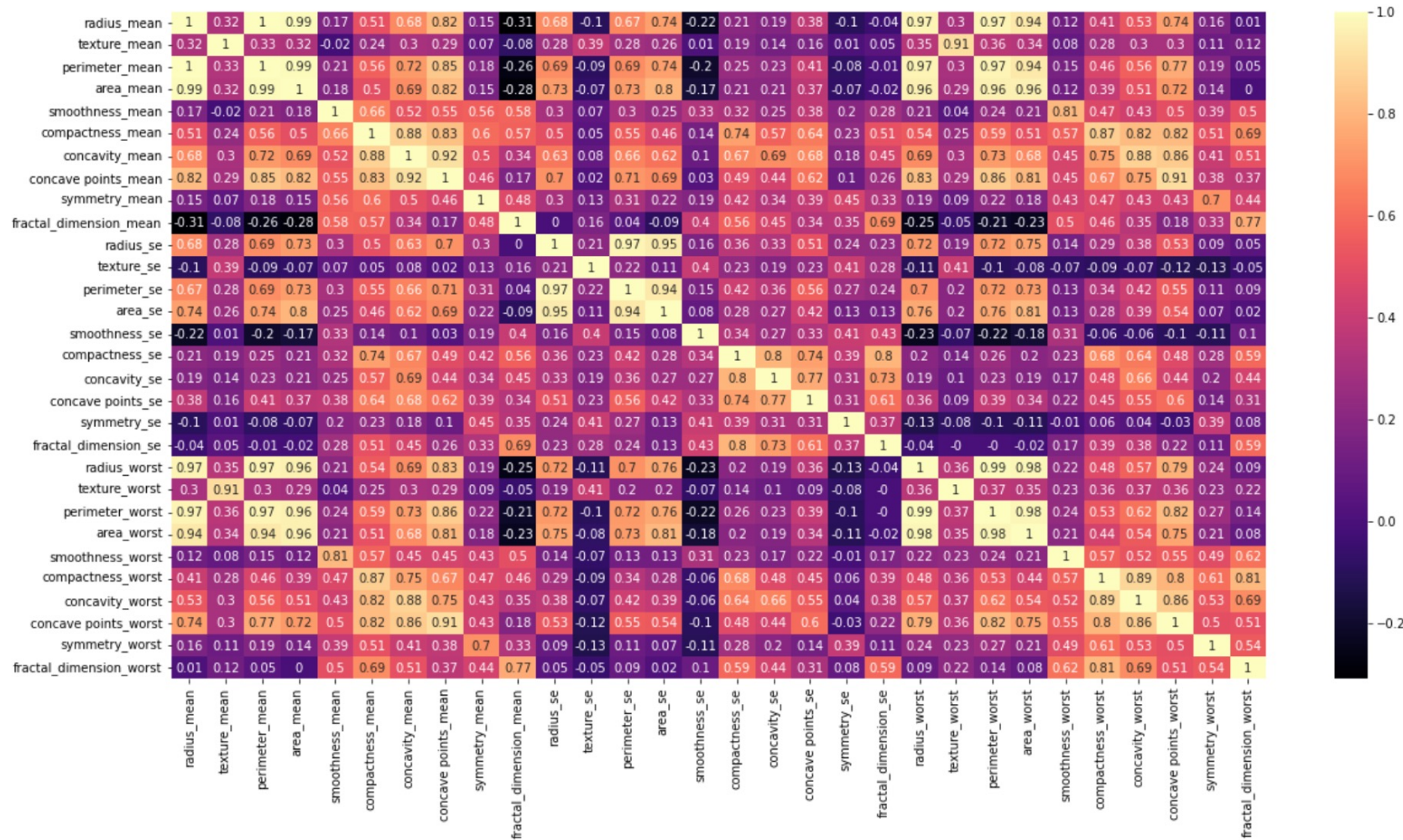
Demonstration

# Data visualization:

# Scatter plot:

# Heat Map:

# Gaussian Naïve Bayes:

```
[ ]  model = GaussianNB()
     model.fit(X_train, y_train)

     GaussianNB()
```

```
▶   predict = model.predict(X_test)
     predict
```

```
array(['B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'M', 'B',
       'M', 'B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'M', 'B', 'B',
       'B', 'M', 'M', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'M',
       'M', 'M', 'B', 'B', 'M', 'B', 'M', 'B', 'B', 'M', 'B', 'M', 'B',
       'B', 'M', 'M', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'M', 'B', 'B',
       'B', 'M', 'B', 'M', 'M', 'B', 'B', 'M', 'M', 'B', 'M', 'B', 'M',
       'M', 'M', 'M', 'B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'M',
       'M', 'M', 'B', 'B', 'M', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B',
       'B', 'M', 'M', 'B', 'M', 'B', 'B', 'M', 'B', 'B', 'M', 'M', 'B',
       'M', 'M', 'M', 'B', 'B', 'M', 'B', 'B', 'M', 'B', 'M', 'B',
       'B', 'M', 'M', 'B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'M',
       'B', 'B', 'M', 'B', 'B', 'B', 'M', 'B', 'B', 'M', 'B', 'B', 'B', 'M',
       'M', 'B', 'M', 'B', 'B', 'M', 'B', 'M', 'M', 'B', 'M', 'B',
       'B', 'B'], dtype='<U1')
```

```
[ ]  model.score(X_test, y_test)

     0.9415204678362573
```

```
[ ]  print(f'Accuracy: {accuracy_score(y_test, predict) * 100}%')

     Accuracy: 94.15204678362574%
```

# Support Vector Machine(SVM):

```
[ ]  model = svm.SVC(kernel='poly', degree=2)
     model.fit(X_train, y_train)

     SVC(degree=2, kernel='poly')

▶    predict = model.predict(X_test)
     predict

👤   array(['B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'M', 'B',
            'B', 'B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
            'B', 'M', 'M', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'M',
            'M', 'M', 'B', 'B', 'M', 'B', 'B', 'M', 'B', 'M', 'B', 'M', 'B',
            'B', 'M', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'M', 'M', 'B',
            'B', 'M', 'B', 'M', 'M', 'B', 'M', 'M', 'M', 'B', 'M', 'B', 'M',
            'M', 'M', 'M', 'B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'M',
            'M', 'M', 'B', 'B', 'M', 'B', 'M', 'B', 'B', 'B', 'M', 'B', 'B',
            'B', 'M', 'M', 'B', 'M', 'B', 'B', 'M', 'B', 'B', 'M', 'M', 'B',
            'M', 'M', 'M', 'B', 'B', 'M', 'B', 'M', 'M', 'B', 'M', 'B', 'M', 'B',
            'B', 'M', 'M', 'B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'M',
            'B', 'B', 'M', 'B', 'B', 'M', 'B', 'B', 'M', 'B', 'B', 'B', 'M',
            'M', 'B', 'M', 'B', 'B', 'B', 'M', 'B', 'B', 'M', 'B', 'M', 'B',
            'B', 'B'], dtype=object)

[ ]  model.score(X_test, y_test)

     0.9824561403508771

[ ]  print(f'Accuracy: {accuracy_score(y_test, predict) * 100}%')

     Accuracy: 98.24561403508771%
```

# Neural Network:

```
[ ]  clf = MLPClassifier(solver='lbfgs', alpha=1e-5, hidden_layer_sizes=(5, 2), random_state=1)
     clf.fit(X_train, y_train)

     MLPClassifier(alpha=1e-05, hidden_layer_sizes=(5, 2), random_state=1,
                   solver='lbfgs')
```

```
▶  predict = clf.predict(X_test)
   predict
```

```
array(['B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'M', 'B',
       'B', 'B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B',
       'B', 'M', 'M', 'B', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'B', 'M',
       'M', 'M', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'M', 'B', 'M', 'B',
       'B', 'M', 'M', 'B', 'B', 'B', 'B', 'B', 'B', 'M', 'M', 'M', 'B',
       'B', 'M', 'B', 'M', 'M', 'B', 'M', 'M', 'M', 'B', 'M', 'B', 'M',
       'B', 'M', 'M', 'B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'M',
       'M', 'M', 'B', 'B', 'M', 'B', 'M', 'B', 'B', 'B', 'M', 'B', 'B',
       'B', 'M', 'M', 'B', 'M', 'B', 'B', 'M', 'B', 'B', 'M', 'M', 'B',
       'M', 'M', 'M', 'B', 'B', 'M', 'B', 'M', 'B', 'M', 'B', 'M', 'B',
       'B', 'M', 'M', 'B', 'B', 'B', 'B', 'M', 'B', 'B', 'B', 'B', 'M',
       'B', 'B', 'M', 'B', 'B', 'M', 'B', 'B', 'M', 'B', 'B', 'B', 'M',
       'M', 'B', 'M', 'B', 'B', 'B', 'M', 'B', 'B', 'M', 'B', 'M', 'B',
       'B', 'B'], dtype='<U1')
```

```
[ ]  clf.score(X_test, y_test)

     0.9649122807017544
```

```
[ ]  print(f'Accuracy: {accuracy_score(y_test, predict) * 100}%')

     Accuracy: 96.49122807017544%
```

# Results

- Using the Naïve Bayes algorithm, a prediction accuracy of 94% was obtained.

- Using SVM, the prediction accuracy was 98%.

- Using Neural Networks, 96% accuracy was obtained.

- By comparison, SVM yields the best prediction accuracy.

# Advantages over existing methods:

- Economic and faster diagnosis of Breast cancer.
- Non-invasive diagnosis using Gaussian Naive Bayes model.

# Incremental Features:

- This project can be further expanded to diagnose other health conditions when the relevant data is available.

- Multiple Machine learning models can be used to compare and improve the diagnosis.

- Implement validation set to tweak the hyper parameters and check correct value for the model.

# References:

- Base paper for the project: https://ieeexplore-ieee-org.libproxy.library.unt.edu/document/9754059

- Breast Cancer Wisconsin (Diagnostic) Data Set - UCI Machine Learning Repository: https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic)
- Breast Cancer Prediction Using Machine Learning - Towards Data Science: https://towardsdatascience.com/breast-cancer-prediction-using-machine-learning-7e83c0b37216
- Breast Cancer Detection Using Machine Learning Techniques: A Review - International Journal of Intelligent Systems and Applications: https://www.researchgate.net/publication/321709047_Breast_Cancer_Detection_Using_Machine_Learning_Techniques_A_Review
- Breast Cancer Diagnosis Using Machine Learning Techniques: A Systematic Review - Applied Sciences: https://www.mdpi.com/2076-3417/10/10/3464
- Breast Cancer Detection Using Machine Learning Algorithms: A Systematic Review - Journal of Medical Systems: https://link.springer.com/article/10.1007/s10916-021-01778-7