

Breast Cancer Prediction Using Classification Model

Name	Id
Md Rifat Khan	20-42704-1
Md Faysal Ahmed	20-42693-1
Al-Ruhul Amin Sabbir	20-42731-1

1. Introduction

1.1 Problem Statement

Breast cancer is a prevalent and potentially life-threatening disease that affects a significant number of women worldwide. Early detection and accurate diagnosis are crucial for effective treatment and increased survival rates. Medical professionals have traditionally relied on various diagnostic methods, including mammography and biopsy. However, with the advent of machine learning, there's an opportunity to enhance the accuracy and efficiency of breast cancer detection. Moreover, the combinations of all the model parameters are considered for simulations. The experimental results of the classifiers show that which classifier model gives better classification accuracy. The existing solution for breast cancer diagnosis relies on manual examination of mammograms, biopsies, and patient history by medical professionals. While effective, this approach can be time-consuming and subject to human error and bias. There is a need for a more automated and objective method that can assist medical practitioners in making accurate and timely diagnoses. This project aims to address this gap by developing a machine learning model that can classify breast tumors as malignant or benign based on quantifiable features. The model will augment the diagnostic process, enabling healthcare professionals to make more informed decisions and potentially leading to earlier detection and improved patient outcomes. In this project NumPy, pandas, matplotlib , seaborn , sklearn libraries are used. Data preprocessing and Exploratory data analysis is present to justify the project as more accurate one. After data analysis unnecessary variables were removed. Then dataset is splitted into training and test set. After that a model is built for each model and test our data. At last it is compared the accuracy score.

1.2 Objective

The primary objective of this project is to design, develop, and evaluate a robust machine learning classification model that enhances the accuracy and efficiency of breast cancer detection. By analyzing quantifiable features extracted from patient data, the model aims to differentiate between

malignant and benign tumors, thereby aiding medical professionals in making informed decisions and potentially facilitating early intervention. The specific objectives of this project include:

i. Dataset Collection and Preprocessing:

Gather a comprehensive dataset comprising various features associated with breast tumor characteristics, patient histories, and relevant medical information. Preprocess the dataset to handle missing values, outliers, and inconsistencies, ensuring the quality and reliability of the data used for model training.

ii. Model Development and Training:

Implement and train machine learning classification algorithms on the preprocessed dataset. Explore a range of algorithms, such as Gaussian Naive Bayes, support vector machines, and K Nearest Neighbours, to identify the most suitable approach for the given task.

iii. Performance Evaluation:

Evaluate the trained models using appropriate evaluation metrics, including accuracy.

By achieving these objectives, we aspire to provide healthcare professionals with a valuable tool that complements their expertise in breast cancer diagnosis. The project's success will be measured by the model's ability to accurately classify tumors, contribute to early detection, reduce false positives/negatives, and enhance overall patient care.

2. Methodology

2.1 Data Collection Procedure

We are loading breast cancer data using a scikit-learn `load_breast_cancer` class.

Number of Instances: 569

Number of Attributes: 30 numeric, predictive attributes and the class

Attribute Information:

- 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)
- symmetry
- fractal dimension ("coastline approximation" - 1)

The mean, standard error, and "worst" or largest (mean of the three worst/largest values) of these features were computed for each image, resulting in 30 features. For instance, field 0 is Mean Radius, field 10 is Radius SE, field 20 is Worst Radius.

- class:
 - WDBC-Malignant
 - WDBC-Benign

Import Libraries

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

Load breast cancer dataset

```
from sklearn.datasets import load_breast_cancer
cancer_dataset = load_breast_cancer()
```

```
1 cancer_dataset
```

Output >>> sklearn {'data': array([[1.799e+01, 1.038e+01, 1.228e+02, ..., 2.654e-01, 4.601e-01, 1.189e-01],
[2.057e+01, 1.777e+01, 1.329e+02, ..., 1.860e-01, 2.750e-01, 8.902e-02],
[1.969e+01, 2.125e+01, 1.300e+02, ..., 2.430e-01, 3.613e-01, 8.758e-02],
...,
[1.660e+01, 2.808e+01, 1.083e+02, ..., 1.418e-01, 2.218e-01,

[illegible]

'target_names': array(['malignant', 'benign'], dtype='<U9'),

'DESCR': '.. _breast_cancer_dataset:\n\nBreast cancer wisconsin (diagnostic) dataset\n-----

-----\n\n**Data Set Characteristics:**\n\n :Number of
Instances: 569\n\n :Number of Attributes: 30 numeric, predictive attributes and the
class\n\n :Attribute Information:\n\n - radius (mean of distances from center to points
on the perimeter)\n\n - texture (standard deviation of gray-scale values)\n\n -
perimeter\n\n - area\n\n - smoothness (local variation in radius lengths)\n\n -
compactness (perimeter² / area - 1.0)\n\n - concavity (severity of concave portions of
the contour)\n\n - concave points (number of concave portions of the contour)\n\n -
symmetry\n\n - fractal dimension ("coastline approximation" - 1)\n\n\n The mean,
standard error, and "worst" or largest (mean of the three\n\n worst/largest values) of
these features were computed for each image,\n\n resulting in 30 features. For
instance, field 0 is Mean Radius, field\n\n 10 is Radius SE, field 20 is Worst Radius.\n\n\n - class:\n\n - WDBC-Malignant\n\n - WDBC-Benign\n\n\n :Summary
Statistics:\n\n =====\n\n =====\n\n\n Min Max\n\n =====\n\n\n =====\n\nradius (mean): 6.981 28.11\n\n texture (mean): 9.71
39.28\n\n perimeter (mean): 43.79 188.5\n\n area (mean):
143.5 2501.0\n\n smoothness (mean): 0.053 0.163\n\n compactness
(mean): 0.019 0.345\n\n concavity (mean): 0.0 0.427\n\n concave points (mean): 0.0 0.201\n\n symmetry (mean): 0.106
0.304\n\n fractal dimension (mean): 0.05 0.097\n\n radius (standard error):
0.112 2.873\n\n texture (standard error): 0.36 4.885\n\n perimeter (standard
error): 0.757 21.98\n\n area (standard error): 6.802 542.2\n\n smoothness (standard error): 0.002 0.031\n\n compactness (standard error):
0.002 0.135\n\n concavity (standard error): 0.0 0.396\n\n concave points
(standard error): 0.0 0.053\n\n symmetry (standard error): 0.008 0.079\n\n fractal dimension (standard error): 0.001 0.03\n\n radius (worst): 7.93
36.04\n\n texture (worst): 12.02 49.54\n\n perimeter (worst):
50.41 251.2\n\n area (worst): 185.2 4254.0\n\n smoothness (worst):
0.071 0.223\n\n compactness (worst): 0.027 1.058\n\n concavity (worst):
0.0 1.252\n\n concave points (worst): 0.0 0.291\n\n symmetry (worst):
0.156 0.664\n\n fractal dimension (worst): 0.055 0.208\n\n\n =====\n\n\n :Missing Attribute Values: None\n\n\n :Class Distribution: 212 - Malignant, 357 - Benign\n\n\n :Creator: Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian\n\n\n :Donor: Nick
Street\n\n\n :Date: November, 1995\n\n\n This is a copy of UCI ML Breast Cancer Wisconsin
(Diagnostic) datasets.\n\n <https://goo.gl/U2Uwz2>\n\n\n Features are computed from a digitized
image of a fine needle\n\n aspirate (FNA) of a breast mass. They describe\n\n characteristics of
the cell nuclei present in the image.\n\n\n Separating plane described above was obtained
using\n\n Multisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree\n\n Construction Via
Linear Programming." Proceedings of the 4th\n\n Midwest Artificial Intelligence and Cognitive
Science Society, npp. 97-101, 1992], a classification method which uses
linear\n\n programming to construct a decision tree. Relevant features\n\n were selected using
an exhaustive search in the space of 1-4\n\n features and 1-3 separating planes.\n\n\n The

actual linear program used to obtain the separating plane\nin the 3-dimensional space is that described in:\n[K. P. Bennett and O. L. Mangasarian: "Robust Linear\nProgramming Discrimination of Two Linearly Inseparable Sets",\nOptimization Methods and Software 1, 1992, 23-34].\n\nThis database is also available through the UW CS ftp server:\n\nftp\nftp.cs.wisc.edu\nncd math-prog/cpo-dataset/machine-learn/WDBC/\n\n.. topic::
References\n\n - W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction \n for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on \n Electronic Imaging: Science and Technology, volume 1905, pages 861-870,\n San Jose, CA, 1993.\n - O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and \n prognosis via linear programming. Operations Research, 43(4), pages 570-577, \n July-August 1995.\n - W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning techniques\n to diagnose breast cancer from fine-needle aspirates. Cancer Letters 77 (1994) \n 163-171.'

```
'feature_names': array(['mean radius', 'mean texture', 'mean perimeter', 'mean area',
                        'mean smoothness', 'mean compactness', 'mean concavity',
                        'mean concave points', 'mean symmetry', 'mean fractal dimension',
                        'radius error', 'texture error', 'perimeter error', 'area error',
                        'smoothness error', 'compactness error', 'concavity error',
                        'concave points error', 'symmetry error',
                        'fractal dimension error', 'worst radius', 'worst texture',
                        'worst perimeter', 'worst area', 'worst smoothness',
                        'worst compactness', 'worst concavity', 'worst concave points',
                        'worst symmetry', 'worst fractal dimension'], dtype='<U23'),
'filename': 'breast_cancer.csv',
'data_module': 'sklearn.datasets.data'}
```

```
1 type(cancer_dataset)
```

Output >>> sklearn.utils.Bunch

The scikit-learn store data in an object bunch like a dictionary.

```
1
2 cancer_dataset.keys()
```

Output >>> dict_keys(['data', 'target', 'target_names', 'DESCR', 'feature_names', 'filename'])

Output >>>

These numeric values are extracted features of each cell.

```
Output >>> array([[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1,
0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1,
1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0,
0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 1,
0, 0, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0,
0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0,
0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1,
0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1,
1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1,
1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1,
1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0,
1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1,
1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0,
0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1])
```

Output >>> array(['malignant', 'benign'], dtype='<U9')

■

Output >>>

```

1 print(cancer_dataset['feature_names'])
2

```

Output >>>

```

1
2 ['mean radius' 'mean texture' 'mean perimeter' 'mean area'
3  'mean smoothness' 'mean compactness' 'mean concavity'
4  'mean concave points' 'mean symmetry' 'mean fractal dimension'
5  'radius error' 'texture error' 'perimeter error' 'area error'
6  'smoothness error' 'compactness error' 'concavity error'
7  'concave points error' 'symmetry error' 'fractal dimension error'
8  'worst radius' 'worst texture' 'worst perimeter' 'worst area'
9  'worst smoothness' 'worst compactness' 'worst concavity'
10 'worst concave points' 'worst symmetry' 'worst fractal dimension']

```

```

1
2 print(cancer_dataset['filename'])

```

Output >>> breast_cancer.csv

2.2. Data Validation Procedure

Data validation is a critical step in ensuring the quality, integrity, and suitability of your dataset for machine learning tasks. In the context of breast cancer prediction, where accurate and trustworthy data is vital, a thorough validation procedure helps identify and address potential issues before training and evaluating classification models.

```
cancer_df.info()
```

Output >>>

```

1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 569 entries, 0 to 568
3 Data columns (total 31 columns):
4 mean radius          569 non-null float64
5 mean texture         569 non-null float64
6 mean perimeter       569 non-null float64
7 mean area           569 non-null float64
8 mean smoothness      569 non-null float64
9 mean compactness     569 non-null float64
10 mean concavity       569 non-null float64
11 mean concave points  569 non-null float64
12 mean symmetry       569 non-null float64
13 mean fractal dimension 569 non-null float64
14 radius error        569 non-null float64

```



```

12 texture error 569 non-null float64
13 perimeter error 569 non-null float64
14 area error 569 non-null float64
15 smoothness error 569 non-null float64
16 compactness error 569 non-null float64
17 concavity error 569 non-null float64
18 concave points error 569 non-null float64
19 symmetry error 569 non-null float64
20 fractal dimension error 569 non-null float64
21 worst radius 569 non-null float64
22 worst texture 569 non-null float64
23 worst perimeter 569 non-null float64
24 worst area 569 non-null float64
25 worst smoothness 569 non-null float64
26 worst compactness 569 non-null float64
27 worst concavity 569 non-null float64
28 worst concave points 569 non-null float64
29 worst symmetry 569 non-null float64
30 worst fractal dimension 569 non-null float64
31 target 569 non-null float64
32 dtypes: float64(31)
33 memory usage: 137.9 KB
34
35
36

```

We have a total of non-null 569 patients' information with 31 features. All feature data types in the float. The size of the DataFrame is 137.9 KB.

Numerical distribution of data. We can know to mean, standard deviation, min, max, 25%,50% and 75% value of each feature.

```

1
2 cancer_df.describe()

```

Output >>>

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst texture	worst perimeter	worst area	worst smoothness	worst compactness	worst concavity	worst concave points	worst symmetry	worst fractal dimension	target
count	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	...	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000
mean	14.127292	19.289649	91.969033	654.889104	0.096360	0.104341	0.088799	0.048919	0.181162	0.062798	...	25.677223	107.261213	880.583128	0.132369	0.254265	0.272188	0.114606	0.290076	0.083946	0.627417
std	3.524049	4.301036	24.298981	351.914129	0.014064	0.052813	0.079720	0.038803	0.027414	0.007060	...	6.146258	33.602542	569.356993	0.022832	0.157336	0.208624	0.065732	0.061867	0.018061	0.483918
min	6.981000	9.710000	43.790000	143.500000	0.052630	0.019380	0.000000	0.000000	0.106000	0.049960	...	12.020000	50.410000	185.200000	0.071170	0.027290	0.000000	0.000000	0.156500	0.055040	0.000000
25%	11.700000	16.170000	75.170000	420.300000	0.086370	0.064920	0.029560	0.020310	0.161900	0.057700	...	21.080000	84.110000	515.300000	0.116600	0.147200	0.114500	0.064930	0.250400	0.071460	0.000000
50%	13.370000	18.840000	86.240000	551.100000	0.095870	0.092630	0.061540	0.033500	0.179200	0.061540	...	25.410000	97.660000	686.500000	0.131300	0.211900	0.226700	0.099930	0.282200	0.080040	1.000000
75%	15.780000	21.800000	104.100000	782.700000	0.105300	0.130400	0.130700	0.074000	0.195700	0.066120	...	29.720000	125.400000	1084.000000	0.146000	0.339100	0.382900	0.161400	0.317900	0.092080	1.000000
max	28.110000	39.280000	188.500000	2501.000000	0.163400	0.345400	0.426800	0.201200	0.304000	0.097440	...	49.540000	251.200000	4254.000000	0.222600	1.058000	1.252000	0.291000	0.663800	0.207500	1.000000

8 rows x 31 columns

```
2 sns.pairplot(cancer_df, hue = 'target')
```

Output >>>

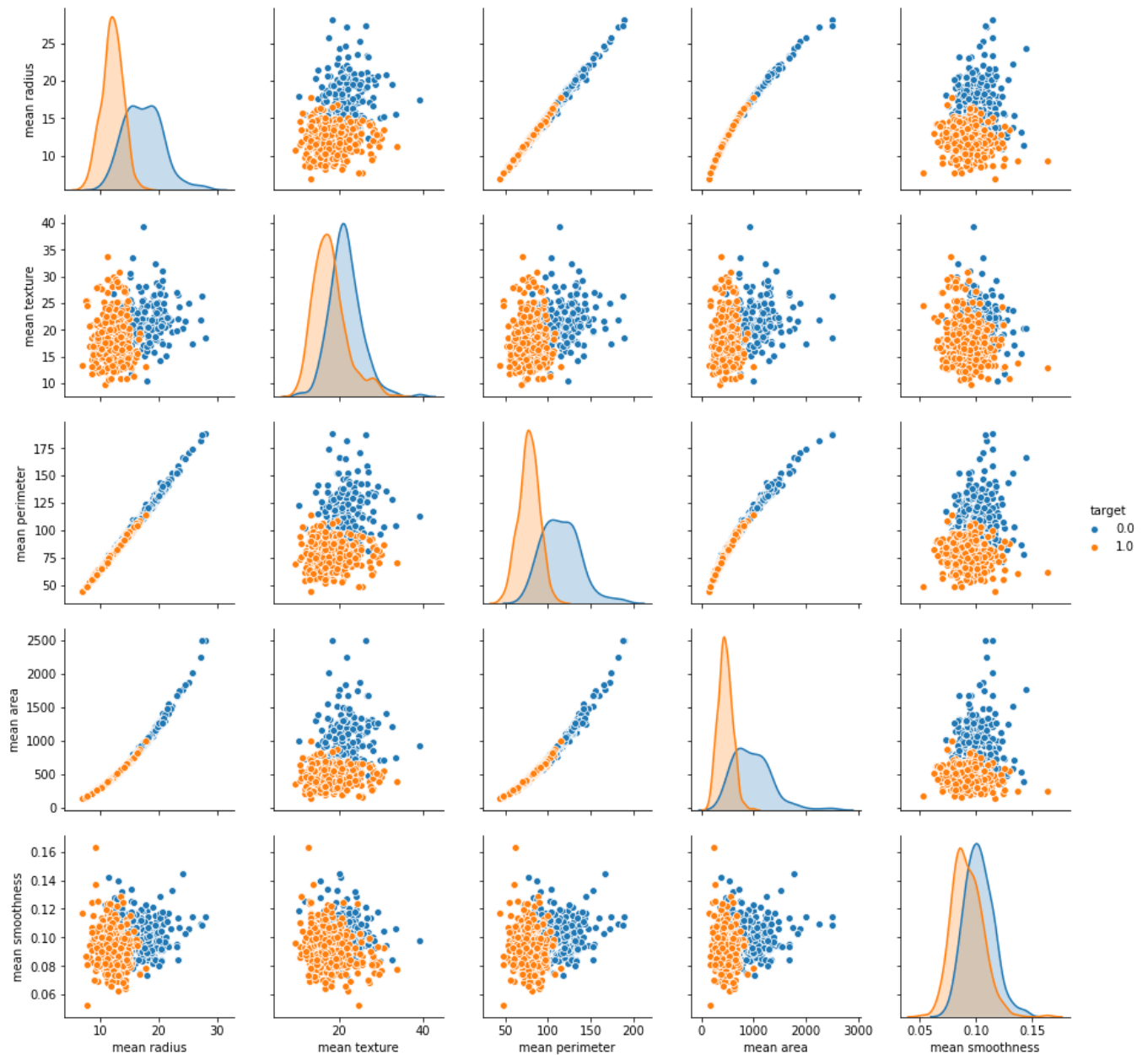


```

1 sns.pairplot(cancer_df, hue = 'target',
2               vars = ['mean radius', 'mean texture', 'mean perimeter', 'mean area', 'mean smoothness'])
3

```

Output >>>



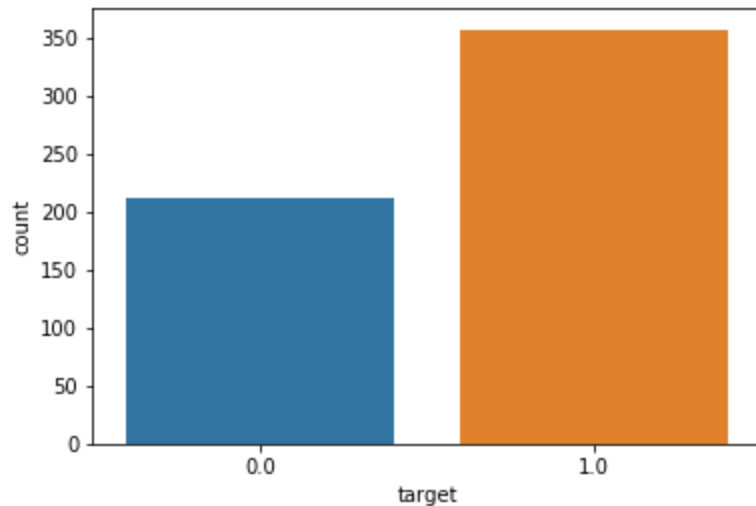
The pair plot showing malignant and benign tumor data distributed in two classes. It is easy to differentiate in the pair plot.

```

1 sns.countplot(cancer_df['target'])
2

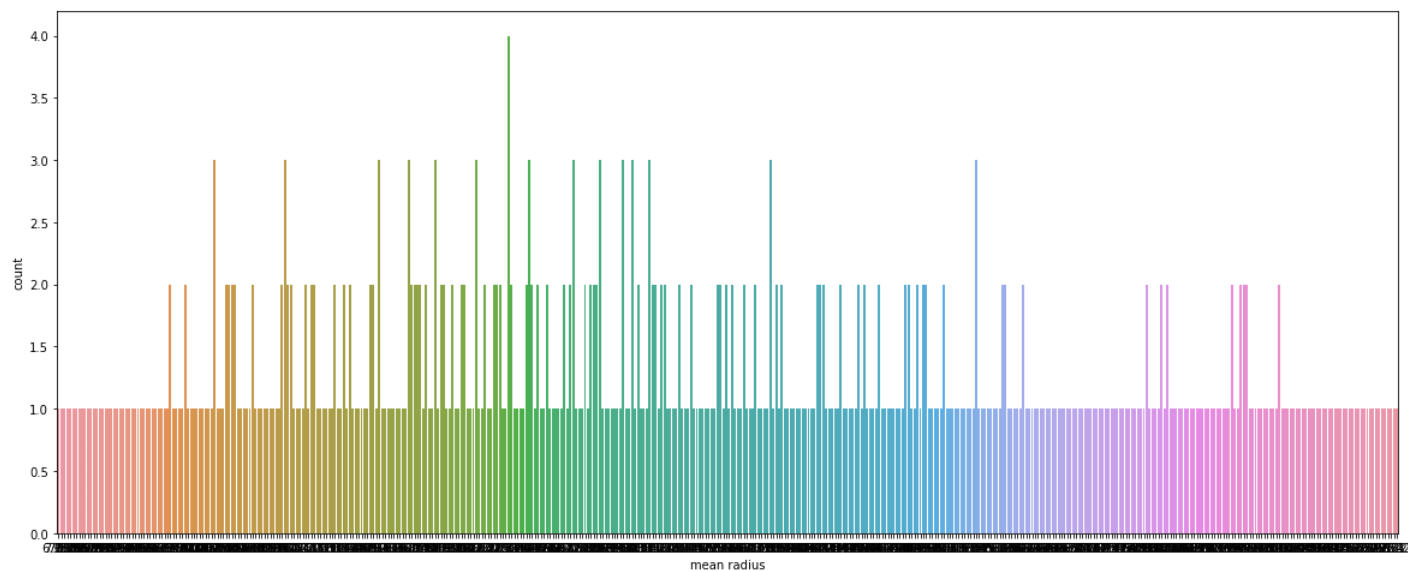
```

Output >>>



In the below counterplot max samples mean radius is equal to 1.

```
1 plt.figure(figsize = (20,8))
2 sns.countplot(cancer_df['mean radius'])
3
```

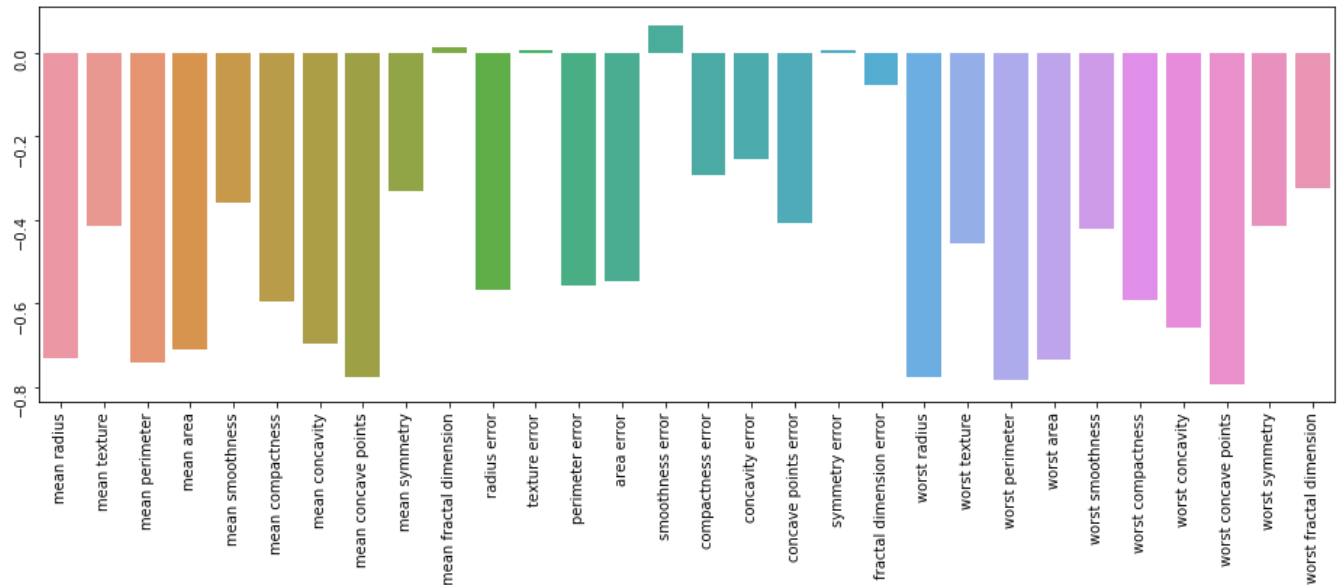


```
1 cancer_df2 = cancer_df.drop(['target'], axis = 1)
2 print("The shape of 'cancer_df2' is : ", cancer_df2.shape)
3
```

Output >>> The shape of 'cancer_df2' is : (569, 30)

```
1 plt.figure(figsize = (16,5))
2 ax = sns.barplot(cancer_df2.corrwith(cancer_df.target).index,
3 cancer_df2.corrwith(cancer_df.target))
4 ax.tick_params(labelrotation = 90)
```

Output >>>



2.3 Data Preprocessing Technique

Data preprocessing is a fundamental step in preparing data for machine learning algorithms. In the context of breast cancer prediction, where accurate and reliable data is crucial, proper preprocessing ensures that the data is clean, consistent, and suitable for training classification models.

```
1
2 X = cancer_df.drop(['target'], axis = 1)
3 X.head(6)
```

Output >>>

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst radius	worst texture	worst perimeter	worst area	worst smoothness	worst compactness	worst concavity	worst concave points	worst symmetry	worst fractal dimension
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	0.07871	...	25.38	17.33	184.60	2019.0	0.1622	0.6656	0.7119	0.2654	0.4601	0.11890
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	0.05667	...	24.99	23.41	158.80	1956.0	0.1238	0.1866	0.2416	0.1860	0.2750	0.08902
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069	0.05999	...	23.57	25.53	152.50	1709.0	0.1444	0.4245	0.4504	0.2430	0.3613	0.08758
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	0.09744	...	14.91	26.50	98.87	567.7	0.2098	0.8663	0.6869	0.2575	0.6638	0.17300
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	0.05883	...	22.54	16.67	152.20	1575.0	0.1374	0.2050	0.4000	0.1625	0.2364	0.07678
5	12.45	15.70	82.57	477.1	0.12780	0.17000	0.1578	0.08089	0.2087	0.07613	...	15.47	23.75	103.40	741.6	0.1791	0.5249	0.5355	0.1741	0.3985	0.12440

6 rows x 30 columns

```
1
2 y = cancer_df['target']
3 y.head(6)
```

Output >>>

```
1 0 0.0
2 1 0.0
3 2 0.0
```

```

3      3      0.0
4      4      0.0
5      5      0.0
6      Name: target, dtype: float64
7
1      from sklearn.model_selection import train_test_split
2      X_train, X_test, y_train, y_test =
3      train_test_split(X, y, test_size = 0.2, random_state= 5)

```

2.4 Feature Extraction Technique

Feature extraction is a crucial step in preparing data for machine learning models. In the context of breast cancer prediction, relevant and informative features can significantly impact the accuracy of the classification model. Converting different units and magnitude data in one unit.

```

1
2      from sklearn.preprocessing import StandardScaler
3      sc = StandardScaler()
4      X_train_sc = sc.fit_transform(X_train)
5      X_test_sc = sc.transform(X_test)

```

2.5 Classification Algorithms

We will use these algorithms to train and predict on your breast cancer dataset. The accuracy results suggest how well these models are performing on our test data.

```

1      from sklearn.metrics import confusion_matrix, classification_report, accuracy_score

```

Support Vector Classifier

```

1      from sklearn.svm import SVC
2      svc_classifier = SVC()
3      svc_classifier.fit(X_train, y_train)
4      y_pred_scv = svc_classifier.predict(X_test)
5      accuracy = accuracy_score(y_test, y_pred_scv).round(4)
6      print("The accuracy of the SVM is:", accuracy)

```

Output >>> The accuracy of the SVM is: 0.9386

K – Nearest Neighbor Classifier

```
1 from sklearn.neighbors import KNeighborsClassifier
2 knn_classifier = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski', p = 2)
3 knn_classifier.fit(X_train, y_train)
4 y_pred_knn = knn_classifier.predict(X_test)
5 accuracy = accuracy_score(y_test, y_pred_knn).round(4)
6 print("The accuracy of the knn is:", accuracy)
```

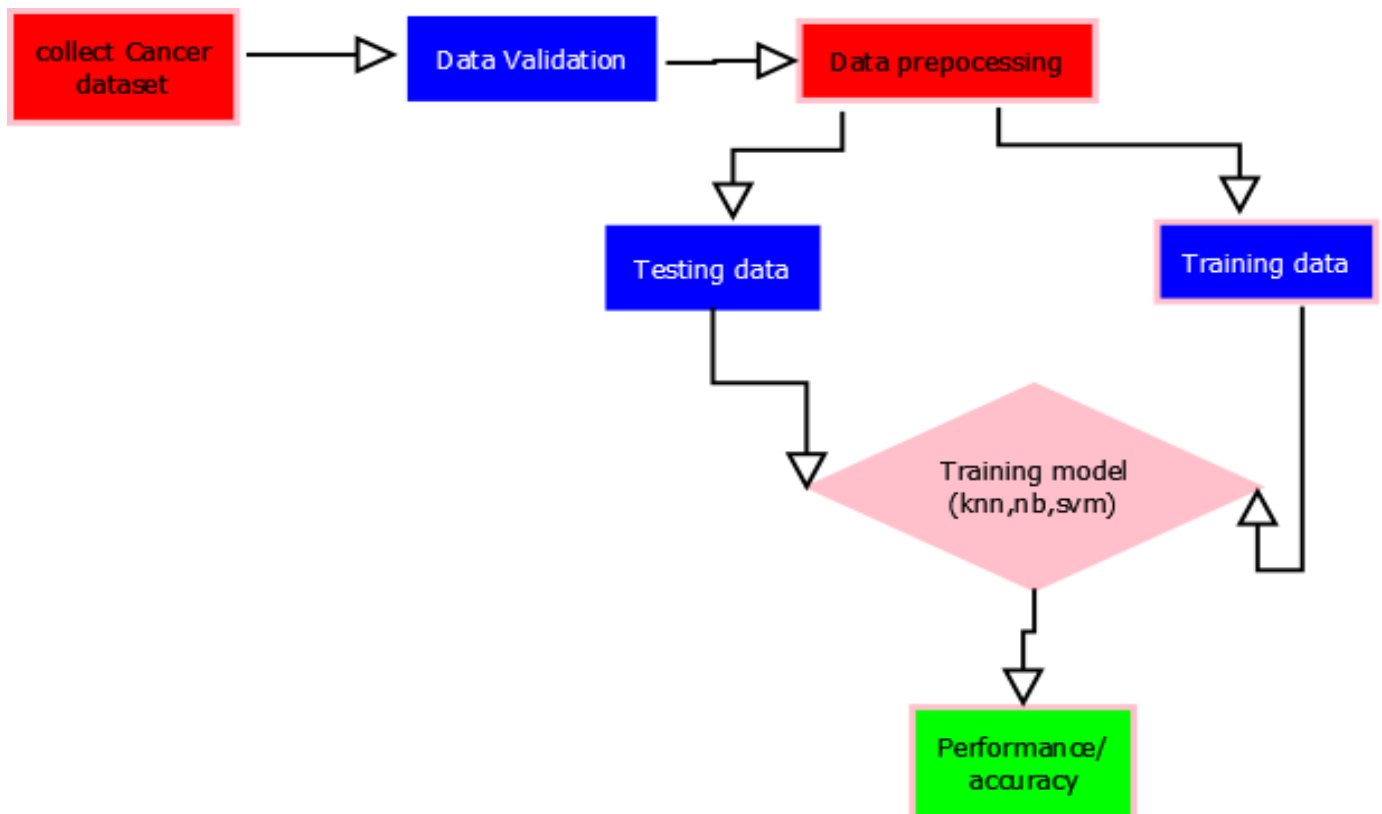
Output >>> The accuracy of the knn is: 0.9386

Naive Bayes Classifier

```
1 from sklearn.naive_bayes import GaussianNB
2 nb_classifier = GaussianNB()
3 nb_classifier.fit(X_train, y_train)
4 y_pred_nb = nb_classifier.predict(X_test)
5 accuracy = accuracy_score(y_test, y_pred_nb).round(4)
6 print("The accuracy of the nb is:", accuracy)
```

Output >>> The accuracy of the nb is: 0.9474

2.6 Block Diagram of Proposed Model



2.7 Data Analysis Techniques

One of the key components of our research on breast cancer was analyzing the data we collected. We used a variety of techniques to gain insights and draw conclusions from the data. One such technique was exploratory data analysis, which allowed us to identify patterns and trends in the data that we could use to inform our models. Another important technique we used was feature selection. This involved identifying the most relevant features in the data set that were most predictive of the outcome variable, in this case, whether a patient had breast cancer or not. By selecting only the most relevant features, we were able to improve the accuracy of our models and reduce the risk of overfitting.

2.8 Experimental Setup

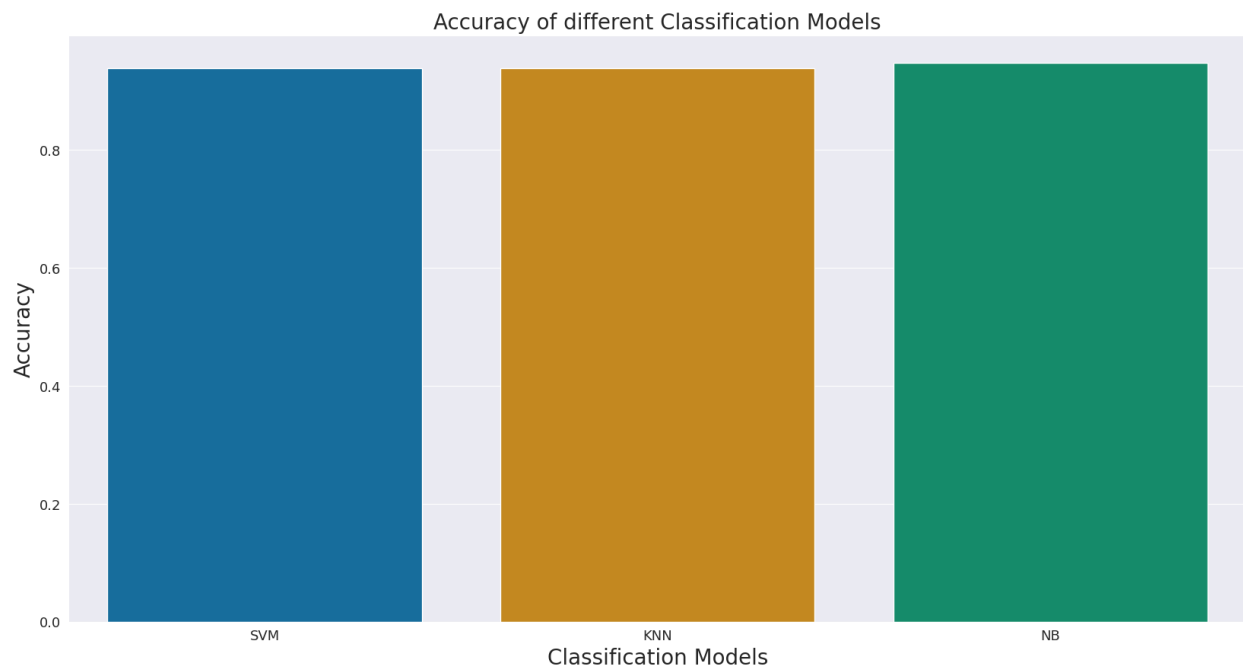
In order to conduct our experiments, we used a dataset of breast cancer patients that had been previously diagnosed. The dataset included various attributes such as age, tumor size, and lymph node involvement. We then preprocessed the data by removing any missing values and normalizing the features. Next, we split the dataset into training and testing sets. We then applied three different classification models - KNN, SVM, and NB - on the training set to predict whether a patient has breast cancer or not. We evaluated the performance of each model using various metrics such as accuracy, precision, and recall.

3. Results and Discussion

3.1 Results Analysis by comparison existing solution

In this report, we did an analysis on a dataset known as the 'load_breast_cancer'. Here, we developed 3 different of classifier model which are Support Vector Machine, k Nearest Neighbours, and Gaussian Naive Bayes. We can see that the accuracy model is K Nearest Neighbours (KNN) which accuracy is 0.9386. The Support Vector Machine (SVM) accuracy is 0.9386, and Gaussian Naive Bayes (NB) accuracy is 0.9474. So, the highest accuracy rate is 0.9474 which is Gaussian Naive Bayes (NB) classifier model. As a result, we can say that the Gaussian Naive Bayes classifier is the best use for this dataset model. The Gaussian Naive Bayes classifier model's accuracy is above 94% because of the dataset. It might perform better if we can train these model on a larger dataset. So in our opinion, depending on this dataset, the Gaussian Naive Bayes classifier is best for predicting the breast cancer. Although, for a larger dataset other model may perform better.

3.2 Results validation by Graphical Representation



4. Conclusion and Future Recommendations

4.1 Significant of outcomes

The outcomes of our study provide valuable insights into the effectiveness of different classification models for breast cancer detection. Our results show that the KNN model had the highest accuracy rate, followed closely by the SVM model. The NB model, while still effective, had a lower accuracy rate compared to the other two models. Furthermore, our data analysis techniques allowed us to identify key features and patterns in the data that were indicative of breast cancer. This information can be used to improve the accuracy of future classification models and ultimately aid in early detection and treatment of breast cancer.

4.2 Recommendation for Future Work

One of the key recommendations for future work is to explore the potential of incorporating other data sources into our analysis. While our current approach has yielded promising results, there may be additional variables that could further improve the accuracy and precision of our classification model. Another area for future research is to investigate the effectiveness of different feature selection techniques. We utilized a standard approach in our analysis, but there may be alternative methods that could yield better performance.