

Abstract : Breast cancer is a major health issue for women worldwide, making early detection crucial for better treatment and survival. While traditional screening methods like mammography, ultrasound, and MRI are important, they can sometimes lead to mistakes, resulting in misdiagnoses in 10% to 30% of cases. This study explores how machine learning (ML) can improve diagnostic accuracy and reduce false positives and negatives in breast cancer detection. We tested eight different ML models—Artificial Neural Network (ANN), Support Vector Machines (SVM), Random Forest Classifier (RFC), Extra Trees Classifier (ETC), Nu-Support Vector Classifier (NUSVC), Passive-Aggressive Classifier (PAC), Extreme Gradient Boosting Classifier (XGBoost), and Adaptive Boosting Classifier (AdaBoost)—using the Breast Cancer Wisconsin (Diagnostic) Dataset, which includes a variety of patients. We created confusion matrices for each model to assess how well they could differentiate between malignant (cancerous) and benign (non-cancerous) cases. We looked at metrics like accuracy, precision, F1-score, recall, and Matthews correlation coefficient (MCC) to evaluate their performance. The results show the advantages and limitations of each ML model, indicating that these technologies can greatly improve the reliability and efficiency of breast cancer diagnosis compared to traditional methods.

Keywords:

1. Introduction:

In the developing world, cancer death is one of the major problems for humankind. It is one of the most significant reasons for women's death. Breast cancer is the most common type of cancer in women with dense breast tissue due to its physiological features [1]. It is distinguished by aberrant proliferation of cells in the mammary glands, which results in the production of malignant tumours. Early detection is critical since it greatly enhances the likelihood of successful treatment and survival. Traditional screening procedures, such as mammography, ultrasound, and MRI, have proved effective in detecting breast cancer. Nevertheless, restrictions, including the need for X-rays, costs, the presence of thick tissue in young patients, and the frequency of incorrect positive and negative results, have forced researchers and organizations to investigate other strategies [2]. However, these procedures are frequently constrained by reasons such as high cost, the requirement for specialized equipment, and fluctuating accuracy rates. These limitations have prompted the incorporation of machine learning (ML) approaches into breast cancer detection, providing a promising option for increasing diagnostic accuracy while minimizing false positives and negatives. In recent years, various studies have examined the use of machine learning algorithms in breast cancer screenings. Researchers have used a variety of approaches, including Support Vector Machines (SVMs), and hybrid models, to enhance the accuracy of breast cancer detection [4]. For example, Jawad Ahmad et al. created a deep learning model that combined Alex-Net and SVM, reaching an accuracy of 99.16% on the Digital Database for Screening Mammography (DDSM) dataset [3]. Similarly, other studies have explored the use of transfer learning, data augmentation, and various machine learning techniques to improve breast cancer diagnosis. These researches have shown that machine learning models can outperform traditional methods, making diagnostic tools far more reliable and efficient.

In our research, we hope to expand the use of machine learning in breast cancer detection by implementing a comprehensive method that includes many classifiers. We created a system that analyses breast cancer datasets using eight different machine learning classifiers, including Support Vector Machines (SVM) [5], Random Forests (RFC), Extra Tree Classifier (ETC) [9], XGBoost Classifier, and Artificial Neural Networks (ANNs). Each of these classifiers has unique strengths and shortcomings when it comes to dealing with different characteristics of data, such as linearity, complexity, and noise sensitivity. As a result, we compare eight supervised machine learning methods on a dataset of breast cancer patients from diverse demographic backgrounds. We then outline the key contributions of our research in the subsequent sections.

2. Methodology:

The goal of this study is to evaluate various machine learning algorithms to identify the model for detecting breast cancer. Our approach follows a structured methodology, covering data acquisition, preprocessing, feature extraction, model selection, evaluation, and optimization.

2.1 Data Acquisition and Exploration

We used the Breast Cancer Wisconsin Diagnostic dataset, sourced from the University of Wisconsin Hospitals, which consists of multiple features characterizing cell nuclei present in breast masses. These include measurements of radius, texture, smoothness, compactness, symmetry, and other key attributes. Initial exploration involved analyzing dataset statistics, examining class distributions malignant vs. benign, and identifying any potential biases or imbalances. This exploration helped us determine necessary preprocessing steps. Data preprocessing ensures high-quality input for the machine learning models. This stage consists of the following steps:

Data Cleaning was used to check for any missing or inconsistent values and addressed these through imputation (e.g., mean or median filling) or removal, depending on the extent and impact of missing data. Additionally, outliers were detected and managed using statistical methods, such as Z-score thresholding or interquartile range (IQR), to enhance data reliability.

2.2 Feature Extraction and Selection

Identifying relevant features improves model interpretability and accuracy by removing redundant or noisy data. Correlation Analysis and Feature Selection was conducted to identify and remove highly correlated features, minimizing multicollinearity and reducing model complexity. We selected features based on their statistical significance to the target variable .

2.3 Model Development and Training

We experimented with several machine learning classifiers, each using a structured train-test split of 75% for training and 25% for testing:

Random Forest was chosen because it mitigates overfitting, a common issue in decision trees, by creating an ensemble of diverse trees. This model is known for its high accuracy and resilience in handling missing data, which can be beneficial in medical datasets with potentially incomplete records. Additionally, its feature importance ranking helps us understand which attributes most influence the model's prediction.

k-Nearest Neighbors (KNN) is effective for medical classification tasks due to its non-parametric nature, making it flexible for varied data patterns without assuming any underlying data distribution. We chose KNN to leverage its simplicity and because it performs well with labeled, easily distinguishable instances, as it can classify instances based on the closest neighbors in the feature space. Logistic Regression was chosen because of its simplicity, interpretability, and effectiveness for binary classification tasks. It's particularly useful in health diagnostics as it provides insights into how each predictor variable influences the likelihood of the target outcome, thus enhancing model interpretability and helping identify key risk factors in breast cancer detection.

Support Vector Machine (SVM) is selected because it is particularly effective in high-dimensional spaces and is robust to cases where the number of features exceeds the number of samples, as is often seen in diagnostic datasets. SVM's ability to create a clear margin of separation between classes also makes it suitable for cases like ours, where distinct separations between malignant and benign classes are essential for accurate diagnosis.

XGBoost was selected due to its robustness and efficiency, particularly with large datasets, by incorporating gradient boosting and advanced regularization techniques to reduce overfitting. Known for its speed and

precision, XGBoost allows for nuanced classification by optimizing the loss function more effectively, which is beneficial in scenarios requiring high accuracy and reliability, such as cancer detection.

Artificial Neural Network (ANN) was chosen for its ability to model complex, non-linear relationships through multi-layered architectures. Its strength lies in learning intricate patterns within the data, making it suitable for applications like breast cancer detection where interactions between multiple features are critical. ANN's capacity for feature learning enables it to uncover patterns beyond the reach of traditional algorithms, enhancing the predictive performance of the model on medical data.

NuSVC was incorporated due to its flexibility in defining the decision boundary between classes. It is a variant of the standard Support Vector Machine (SVM) that allows the user to set a parameter, "nu", which controls the trade-off between the margin size and classification errors. This makes it particularly useful for handling datasets where there is noise or the classes are not perfectly separable, as it can adjust to different types of data distributions, offering a robust solution for breast cancer prediction where data may not always follow a clear linear separation.

Passive-Aggressive Classifier was included for its ability to adapt quickly to changes in the data while maintaining a balance between flexibility and stability. It is an online learning algorithm that is particularly efficient for large datasets and real-time prediction tasks, making it ideal for scenarios where rapid adjustments are required in the face of new, incoming data. The PAC's performance is advantageous in medical diagnostics where rapid decision-making is crucial, and it can efficiently handle varying data patterns while minimizing the impact of misclassifications.

2.4 Hyperparameter Tuning

Hyperparameter tuning was using grid search and randomized search techniques to optimize each model's parameters. For example, in Random Forests, we tuned the number of trees and tree depth, while in SVM, we adjusted the kernel type and regularization parameters. Hyperparameter tuning helps enhance accuracy and model robustness by refining each model's structure to best fit the data.

2.5 Cross-Validation for Generalizability

We employed k-fold cross-validation to ensure that our evaluation metrics were representative of the model's performance across different data splits. Cross-validation reduces overfitting and provides a more reliable measure of how well the model will perform on unseen data.

2.6 Comparative Analysis and Selection of the Best Model

After training and evaluating each model, we compared their results based on accuracy, F1 score, and ROC-AUC. The model with the highest combination of accuracy, balanced F1 score, and optimal ROC-AUC was selected as the most predictive algorithm for breast cancer detection.

3. Machine learning models:

Machine learning is a field of research which relies on statistics, mathematics, and algorithms. Machine learning can be classified into three categories: supervised, unsupervised, and reinforcement learning. The selection of machine learning techniques is reliant upon the attributes of the data being leveraged.

With the progressive improvement of computational capability of processing units and availability of an unprecedented amount of data in the public domain, machine learning has gained colossal attention in applications across diverse fields. One of the essential concepts of Machine Learning is supervised learning [7]. In supervised learning, we use a Machine Learning algorithm to learn the mapping function between the input (X) and output (Y) variables such that $Y = f(X)$. Here, the goal is to approximate the mapping function so that the algorithm can predict the out variables (Y) for a new input data (X).

In our study, we have compared the performance of eight supervised machine learning models on the breast cancer dataset. We have used these following models:

Support Vector Classifier (SVC) is a supervised machine learning algorithm which produces an N-dimensional hyperplane to classify data. A data plane with the most distance between the data points of both classes is opted for. Hyper planes are the decision boundaries used to divide data into categories. The basic goal is to identify the hyperplane that best splits the data so that it can be appropriately recognized [5].

A *random forest classifier (RFC)* is an ensemble of decision trees where a prediction is made collectively by several decision trees. Here, each tree in the ensemble is formed from a sample of the training set, which is drawn with replacement [6].

Extremely randomized trees, also known as extra trees classifiers, are an ensemble of decision trees similar to random forests that generates numerous unpruned decision trees from training data and makes predictions using a majority vote of decision trees [8]. The additional trees approach fits each decision tree to the entire training data, unlike the random forest technique, which creates each decision tree from a bootstrap sample of training data.

Nu-Support Vector Classifier (ν-SVC) is one of the variants of a support vector classifier, which was introduced by Scholkopf et al [9]. This variant of the support vector machine (SVM) algorithm is essentially used to govern the maximum separation between the subsets of the convex hulls of the data, which are usually known as soft convex hulls. These soft convex hulls are generally controlled by the value of the parameter.

Passive aggressive classifiers (PAC) are a family of algorithms that perform online learning of massive streams of data [10]. In online machine learning algorithms like Passive Aggressive Classifier (PAC), the ML model is updated in a step-by-step fashion with respect to the sequential arrival of the input streams of data.

A computational model called an *artificial neural network (ANN)* aims to replicate how the human brain interprets and processes information. It is made up of an input layer, one or more hidden layers, and an output layer, among other layers of neurons. Weighted connections link each neuron to every other neuron in the network, and the network learns by changing these weights via a process known as backpropagation. Because of their ability to represent intricate, non-linear relationships within data, ANNs have proven effective at tasks like pattern and image recognition. In order to enable the model to identify complex patterns and correlations that might point to the presence of benign or malignant tumours, artificial neural networks (ANN) were used in this study to assess data on breast cancer [11].

XGBoost, or extreme Gradient Boosting, is a fast and scalable gradient boosting method that is commonly used in classification and regression tasks [12]. It constructs decision trees in a stepwise manner, fixing mistakes in earlier iterations to increase accuracy with each build. XGBoost was used in this study to analyze data on breast cancer, handling big datasets and intricate features with ease to get extremely precise predictions.

The AdaBoost Classifier [13] is a simple supervised learning algorithm in order to solve classification problems. In contrast to other boosting algorithms, AdaBoost prioritizes challenging cases by modifying the weights of instances that are wrongly categorized after each iteration. This iterative procedure persists until the model attains the intended degree of precision. AdaBoost improves the model's overall robustness and accuracy by concentrating on difficult cases, making sure that it functions effectively even in complicated situations.

An ensemble learning method called a *voting classifier* combines the predictions of several supervised models to increase overall accuracy [14]. To increase the reliability of the final predictions for the breast cancer dataset, two methods of voting were implemented: hard voting, which chooses the majority class, and soft voting, which averages predicted probabilities.

4. Performance Metrics:

In this section, we briefly discuss five performance metrics which are used to measure the performance of the classifiers in our study.

4.1 Accuracy

Accuracy measure shows the probability of true values. It is an important measure since it tells the number of correctly classified instances by the classifier [15]. Accuracy is of two kinds: training accuracy and testing accuracy. The training and testing accuracy have no correlation with each other.

$$Accuracy = (TP + TN) / (TP + TN + FP + FN) \quad [16]$$

In this formula, TP and TN are the number of true positive and negative instances, whereas FP and FN are the number of false Positive and negative instances respectively.

4.2 Precision

Precision measures the ability of a classifier to correctly classify the positive labels among all the instances predicted as positive. In other words, precision can be expressed as the accuracy of the predictions of positive levels. It is expressed as

$$Precision = TP / TP + FP \quad [16]$$

where TP and FP are the number of true positive instances and number of false positive instances, respectively. A classifier with higher precision value is always preferred.

4.3 f1-score

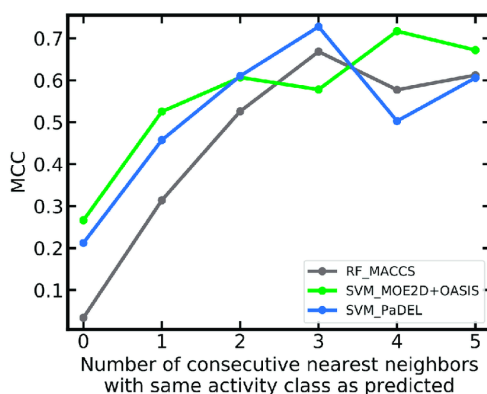
f1-score is the measure which combines precision and recall. It is to be mentioned that recall is also known as true positive rate, i.e., the ratio of the correctly predicted positive instances, and the sum of falsely predicted negative instances and the correctly predicted positive instances, i.e.

$$Recall = TP / FN + TP \quad [16]$$

where FN is the number of false negative instances. Accordingly, the f1-score is determined by calculating the harmonic mean of precision and recall. Since the harmonic mean gives more weightage to minority class values, therefore a classifier will achieve a higher value only when both the precision and recall of the classifier are high. A higher value of f1-score is preferable.

4.4 Matthews Correlation Coefficient (MCC)

The Matthews correlation coefficient [17] considers all the true and false positives and negatives to measure the prediction quality of a classifier. This metric works well even if the classification classes have indifferent sizes. Being a correlation coefficient, the value of MCC lies within the interval $[-1, +1]$. For a perfect prediction, the MCC takes the value +1. An average random prediction is implied if the value of MCC is 0. Whereas, for an inverse prediction value of MCC becomes -1.



5. Dataset Description:

In this study, public data from Kaggle about breast cancer tumors from Dr. William H. Walberg of the University Wisconsin Hospital was taken and used for data visualization, classification, and machine learning algorithms, which included logistic regression, Random forest, support vector machine, and extra trees. Public data included samples taken from patients with solid breast masses and a user-friendly usage of graphical programs called City. This study aimed to establish an adequate model by revealing the predictive factors of early-stage breast cancer patients from a wider perspective and compare the strength of the model with accuracy measures.

The dataset consists of 569 observations and 32 features covering patient information and characteristics of breast cancer tumors to assist in classification and prediction. The features cover various tumor dimensions like radius, texture, perimeter, area, smoothness, compactness, and concavity, each measured through mean, standard error, and worst-case values. These parameters play an important role in cancer detection, as higher values are often indicative of malignancy. Additionally, the dataset includes an ID column which is not necessary and a diagnosis column to classify tissues as malignant or benign. Understanding these features is essential for establishing an accurate model for early cancer detection.

In terms of dataset quality, a comprehensive check was conducted to identify any missing values, as missing data can impact model performance. Fortunately, the dataset contained no null values across any features, facilitating a smooth analysis process and enabling full utilization of each characteristic for model training and evaluation. By leveraging all these features, we rigorously tested and compared the performance of various machine learning classifiers. These efforts aimed to highlight each model's strengths, accuracy, and potential for aiding in the early detection and diagnosis of breast cancer. Through accurate predictions, these models could ultimately support medical professionals in identifying patients at higher risk and recommending early interventions, thereby contributing to better patient outcomes.

6. Data Preprocessing:

Data preparation is a challenging step in getting a dataset ready for machine learning models. The quality and structure of the raw data have a substantial impact on the model's performance, and inappropriate data management might result in poor generalization, misclassification, or overfitting. In our study, we used a variety of preprocessing strategies to address common difficulties in the Wisconsin Breast Cancer Dataset including missing values, outliers, and data scale fluctuations. These procedures were critical in making certain that the dataset was arranged in a manner that maximized the model's ability to learn.

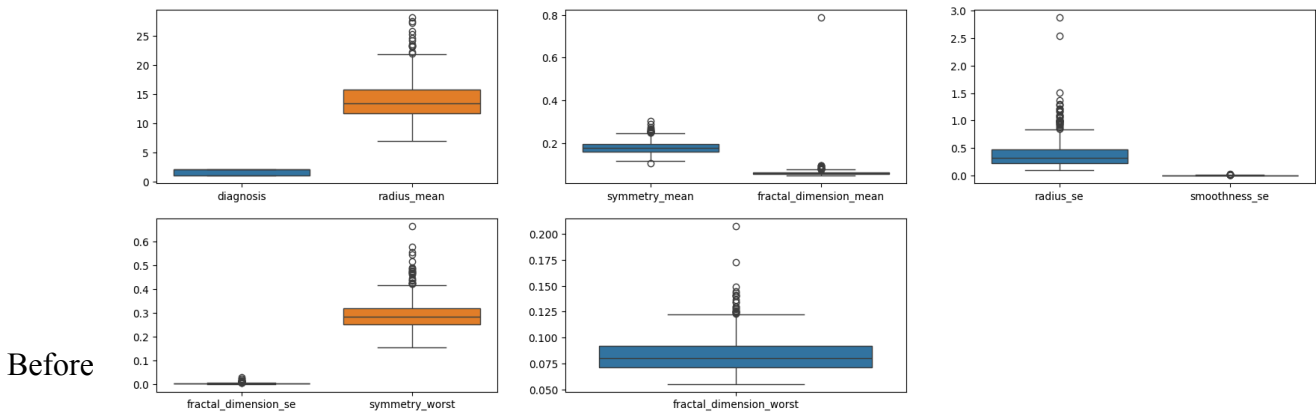
The breast cancer dataset used in this study included variables that varied in scale, distribution, and kind. To limit the effects of these changes, we used feature scaling, normalization, and outlier detection. In addition, missing values in the dataset were imputed using statistical methods, ensuring that no critical data was lost throughout the study. Another critical part of our preprocessing was tackling the issue of highly correlated characteristics, which might trigger multicollinearity and reduce model interpretability. To have a better understanding of the correlations between the features, we calculated the correlation matrix before and after preprocessing. This study was critical for identifying feature associations that might impact the learning process, especially those with high correlation coefficients, which indicate redundancy.

In the dataset, certain features exhibited relatively high numerical integer values, while others had significantly lower values, indicating a wide range in the scales of different attributes. This variation in numerical ranges can impact the overall model performance and may lead to biased outcomes if not properly addressed. Additionally, we observed that many features contained outliers, which could introduce noise and increase the likelihood of overfitting during model training. To address these issues, we removed the outliers from the dataset, aiming to improve the model's generalizability and reliability. Since plotting charts and examining relationships among these features with such unbalanced scales would not yield adequate or interpretable results, we applied normalization techniques to bring the features to a similar scale, ensuring that no single feature dominated the model due to its magnitude.

Consequently, numerical values are normalized. The normalization equation of numerical values can be estimated as

$$i^* = (i - \mu) / \sigma \quad (18)$$

After normalizing the numerical values of features, a box plot was created for cleaning the data. Figure 5 indicates a box plot of features. According to it, some features were excluded since there were too many outliers. This means sufficiency was lacking to use these features while classifying by looking at the boxplot.



applying preprocessing, the correlation matrix was computed to assess relationships between features in the raw dataset. Correlation coefficients, ranging from -1 to +1, indicated the strength of these relationships, with higher correlations signaling potential multicollinearity. Multicollinearity can result in redundant

id	1	-0.04	0.075	0.1	0.073	0.097	-0.013	9.6e-05	0.05	0.044	-0.022	-0.021	0.14	-0.0075	0.14	0.18	0.097	0.034	0.055	0.079	-0.018	0.026	0.082	0.065	0.08	0.11	0.01	-0.003	0.023	0.035	-0.044	-0.03
diagnosis	-0.04	1	-0.73	-0.42	-0.74	-0.71	-0.36	-0.6	-0.7	-0.78	-0.33	-0.049	-0.57	0.0083	-0.56	-0.55	0.067	-0.29	-0.25	-0.41	-0.039	-0.078	-0.78	-0.46	-0.78	-0.73	-0.42	-0.59	-0.66	-0.79	-0.42	-0.32
radius_mean	0.075	-0.73	1	0.32	1	0.99	0.17	0.51	0.68	0.82	0.15	-0.027	0.68	-0.097	0.67	0.74	-0.22	0.21	0.19	0.38	-0.024	-0.043	0.97	0.3	0.97	0.94	0.12	0.41	0.53	0.74	0.16	0.0071
texture_mean	0.1	-0.42	0.32	1	0.33	0.32	-0.023	0.24	0.3	0.29	0.071	-0.1	0.28	0.39	0.28	0.26	0.0066	0.19	0.14	0.16	-0.063	0.054	0.35	0.91	0.36	0.34	0.078	0.28	0.3	0.3	0.11	0.12
perimeter_mean	0.073	-0.74	1	0.33	1	0.99	0.21	0.56	0.72	0.85	0.18	-0.0085	0.69	-0.087	0.69	0.74	-0.2	0.25	0.23	0.41	-0.005	-0.0055	0.97	0.3	0.97	0.94	0.15	0.46	0.56	0.77	0.19	0.051
area_mean	0.097	-0.71	0.99	0.32	0.99	1	0.18	0.5	0.69	0.82	0.15	-0.025	0.73	-0.066	0.73	0.8	-0.17	0.21	0.21	0.37	-0.0092	-0.02	0.96	0.29	0.96	0.96	0.12	0.39	0.51	0.72	0.14	0.0037
smoothness_mean	-0.013	-0.36	0.17	-0.023	0.21	0.18	1	0.66	0.52	0.55	0.56	0.2	0.3	0.068	0.3	0.25	0.33	0.32	0.25	0.38	0.17	0.28	0.21	0.036	0.24	0.21	0.81	0.47	0.43	0.5	0.39	0.5
compactness_mean	3.6e-05	-0.6	0.51	0.24	0.56	0.5	0.66	1	0.88	0.83	0.6	0.26	0.5	0.046	0.55	0.46	0.14	0.74	0.57	0.64	0.24	0.51	0.54	0.25	0.59	0.51	0.57	0.87	0.82	0.82	0.51	0.69
concavity_mean	0.05	-0.7	0.68	0.3	0.72	0.69	0.52	0.88	1	0.92	0.5	0.18	0.63	0.076	0.66	0.62	0.099	0.67	0.69	0.68	0.19	0.45	0.69	0.3	0.73	0.68	0.45	0.75	0.88	0.86	0.41	0.51
concave points_mean	0.044	-0.78	0.82	0.29	0.85	0.82	0.55	0.83	0.92	1	0.46	0.14	0.7	0.021	0.71	0.69	0.028	0.49	0.44	0.62	0.14	0.26	0.83	0.29	0.86	0.81	0.45	0.67	0.75	0.91	0.38	0.37
symmetry_mean	-0.022	-0.33	0.15	0.071	0.18	0.15	0.56	0.6	0.5	0.46	1	0.2	0.3	0.13	-0.31	0.22	0.19	0.42	0.34	0.39	0.33	0.33	0.19	0.091	0.22	0.18	0.43	0.47	0.43	0.43	0.7	0.44
fractal_dimension_mean	-0.021	-0.049	-0.027	-0.1	-0.0085	-0.025	0.2	0.26	0.18	0.14	0.2	1	0.099	0.015	0.12	0.079	0.082	0.18	0.13	0.1	0.84	0.19	0.018	-0.066	0.046	0.027	0.17	0.21	0.16	0.13	0.19	0.25
radius_se	0.14	-0.57	0.68	0.28	0.69	0.73	0.3	0.5	0.63	0.7	0.3	0.099	1	0.21	0.97	0.95	0.16	0.36	0.33	0.51	0.22	0.23	0.72	0.2	0.72	0.75	0.14	0.29	0.38	0.53	0.094	0.05
texture_se	0.0075	0.0083	-0.097	0.39	-0.087	-0.066	0.068	0.046	0.076	0.021	0.13	0.015	0.21	1	0.22	0.11	0.4	0.23	0.19	0.23	0.22	0.28	-0.11	0.41	-0.1	-0.083	-0.074	-0.092	-0.069	-0.12	-0.13	-0.046
perimeter_se	0.14	-0.56	0.67	0.28	0.69	0.73	0.3	0.55	0.66	0.71	0.31	0.12	0.97	0.22	1	0.94	0.15	0.42	0.36	0.56	0.25	0.24	0.7	0.2	0.72	0.73	0.13	0.34	0.42	0.55	0.11	0.085
area_se	0.18	-0.55	0.74	0.26	0.74	0.8	0.25	0.46	0.62	0.69	0.22	0.079	0.95	0.11	0.94	1	0.075	0.28	0.27	0.42	0.16	0.13	0.76	0.2	0.76	0.81	0.13	0.28	0.39	0.54	0.074	0.018
smoothness_se	0.097	0.067	-0.22	0.0066	-0.2	-0.17	0.33	0.14	0.099	0.028	0.19	0.082	0.16	0.4	0.15	0.075	1	0.34	0.27	0.33	0.23	0.43	-0.23	-0.075	-0.22	-0.18	0.31	-0.056	-0.058	-0.1	-0.11	0.1
compactness_se	0.034	-0.29	0.21	0.19	0.25	0.21	0.32	0.74	0.67	0.49	0.42	0.18	0.36	0.23	0.42	0.28	0.34	1	0.8	0.74	0.27	0.8	0.2	0.14	0.26	0.2	0.23	0.68	0.64	0.48	0.28	0.59
concavity_se	0.055	-0.25	0.19	0.14	0.23	0.21	0.25	0.57	0.69	0.44	0.34	0.13	0.33	0.19	0.36	0.27	0.27	0.8	1	0.77	0.2	0.73	0.19	0.1	0.23	0.19	0.17	0.48	0.66	0.44	0.2	0.44
concave points_se	0.079	-0.41	0.38	0.16	0.41	0.37	0.38	0.64	0.68	0.62	0.39	0.1	0.51	0.23	0.56	0.42	0.33	0.74	0.77	1	0.2	0.61	0.36	0.087	0.39	0.34	0.22	0.45	0.55	0.6	0.14	0.31
symmetry_se	-0.018	-0.039	-0.024	-0.063	-0.005	-0.0092	0.17	0.24	0.19	0.14	0.33	0.84	0.22	0.22	0.25	0.16	0.23	0.27	0.2	0.2	1	0.24	-0.011	-0.09	0.017	0.0027	0.036	0.12	0.091	0.059	0.32	0.11
fractal_dimension_se	0.026	-0.078	-0.043	0.054	-0.0055	-0.02	0.28	0.51	0.45	0.26	0.33	0.19	0.23	0.28	0.24	0.13	0.43	0.8	0.73	0.61	0.24	1	-0.037	-0.0032	-0.001	-0.023	0.17	0.39	0.38	0.22	0.11	0.59
radius_worst	0.082	-0.78	0.97	0.35	0.97	0.96	0.21	0.54	0.69	0.83	0.19	0.018	0.72	-0.11	0.7	0.76	-0.23	0.2	0.19	0.36	-0.011	-0.037	1	0.36	0.99	0.98	0.22	0.48	0.57	0.79	0.24	0.093
texture_worst	0.065	-0.46	0.3	0.91	0.3	0.29	0.036	0.25	0.3	0.29	0.091	-0.066	0.2	0.41	0.2	0.2	-0.075	0.14	0.1	0.087	-0.09	-0.0032	0.36	1	0.37	0.35	0.23	0.36	0.37	0.36	0.23	0.22
perimeter_worst	0.08	-0.78	0.97	0.36	0.97	0.96	0.24	0.59	0.73	0.86	0.22	0.046	0.72	-0.1	0.72	0.76	-0.22	0.26	0.23	0.39	0.017	-0.001	0.99	0.37	1	0.98	0.24	0.53	0.62	0.82	0.27	0.14
area_worst	0.11	-0.73	0.94	0.34	0.94	0.96	0.21	0.51	0.68	0.81	0.18	0.027	0.75	-0.083	0.73	0.81	-0.18	0.2	0.19	0.34	0.0027	-0.023	0.98	0.35	0.98	1	0.21	0.44	0.54	0.75	0.21	0.08
smoothness_worst	0.01	-0.42	0.12	0.078	0.15	0.12	0.81	0.57	0.45	0.45	0.43	0.17	0.14	-0.074	0.13	0.13	0.31	0.23	0.17	0.22	0.036	0.17	0.22	0.23	0.24	0.21	1	0.57	0.52	0.55	0.49	0.62
compactness_worst	-0.003	-0.59	0.41	0.28	0.46	0.39	0.47	0.87	0.75	0.67	0.47	0.21	0.29	-0.092	0.34	0.28	-0.056	0.68	0.48	0.45	0.12	0.39	0.48	0.36	0.53	0.44	0.57	1	0.89	0.8	0.61	0.81
concavity_worst	0.023	-0.66	0.53	0.3	0.56	0.51	0.43	0.82	0.88	0.75	0.43	0.16	0.38	-0.069	0.42	0.39	-0.058	0.64	0.66	0.55	0.091	0.38	0.57	0.37	0.62	0.54	0.52	0.89	1	0.86	0.53	0.69
concave points_worst	0.035	-0.79	0.74	0.3	0.77	0.72	0.5	0.82	0.86	0.91	0.43	0.13	0.53	-0.12	0.55	0.54	-0.1	0.48	0.44	0.6	0.059	0.22	0.79	0.36	0.82	0.75	0.55	0.8	0.86	1	0.5	0.51
symmetry_worst	-0.044	-0.42	0.16	0.11	0.19	0.14	0.39	0.51	0.41	0.38	0.7	0.19	0.094	-0.13	0.11	0.074	-0.11	0.28	0.2	0.14	0.32	0.11	0.24	0.23	0.27	0.21	0.49	0.61	0.53	0.5	1	0.54
fractal_dimension_worst	-0.03	-0.32	0.0071	0.12	0.051	0.0037	0.5	0.69	0.51	0.37	0.44	0.25	0.095	-0.046	0.085	0.018	0.1	0.59	0.44	0.31	0.11	0.59	0.093	0.22	0.14	0.08	0.62	0.81	0.69	0.51	0.54	1
id	1	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension_mean	radius_se	texture_se	perimeter_se	area_se	smoothness_se	compactness_se	concavity_se	concave points_se	symmetry_se	fractal_dimension_se	radius_worst	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst	concavity_worst	concave points_worst	symmetry_worst	fractal_dimension_worst

information, negatively impacting model coefficients and leading to overfitting. Moreover, outliers in the raw data further distorted these relationships, potentially skewing model predictions. **Figure-1** illustrates the correlation matrix before preprocessing, highlighting significant feature correlations.

Fig 1: Correlation matrix before preprocessing

To process the data for the machine learning models used in our study, we performed several preprocessing steps on the dataset. We adopted two approaches to significantly reduce the probability of overfitting. These steps were carried out in two stages, in the first stage, we identified and removed outliers from each feature of the dataset. This process ensured that the data was cleaned of extreme values that could potentially distort model performance.

In the second stage, we constructed a correlation matrix to examine the relationships between features. Features with a correlation coefficient greater than 0.8 were considered highly correlated and were excluded from the dataset. After this step, the dataset was reduced from 32 features to 12 features which are radius mean, radius se, symmetry mean, fractal dimension mean, fractal dimension se, smoothness se, symmetry worst, fractal dimension worst.

As a result of these preprocessing steps, the final dataset, consisting of 8 features, was used to perform a comparative analysis of six machine learning classifiers, as discussed in the following section.

Data preprocessing was crucial in refining the breast cancer dataset for model training. The initial correlation matrix exposed highly correlated features and outliers, which were addressed through feature scaling, normalization, and outlier removal, reducing multicollinearity. These steps ensured a well-structured dataset, enhancing both model interpretability and generalization. The updated correlation matrix showed improved feature relationships, enabling more accurate and reliable breast cancer detection by the classifiers.

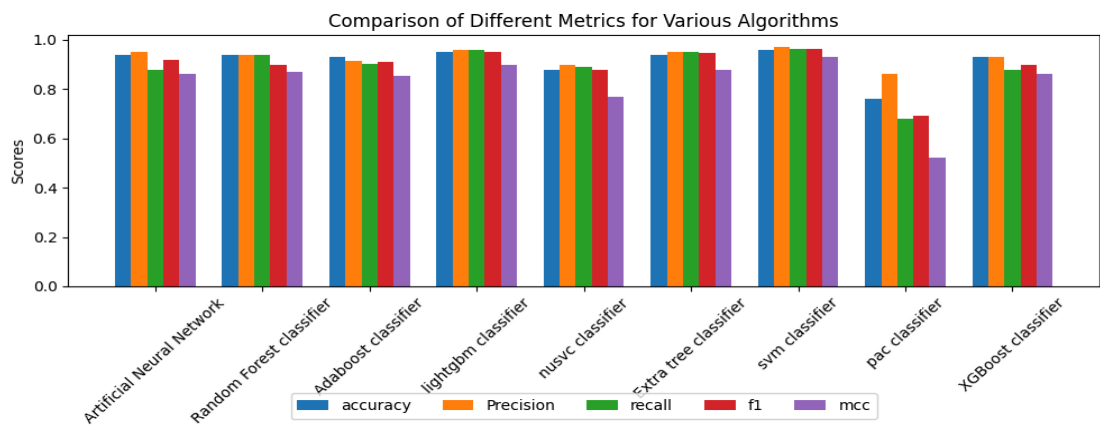
7. Results and Discussions

This section presents the results of the machine learning models used on the breast cancer dataset, followed by a detailed analysis of their performance. The outcomes are measured using a variety of measures, including accuracy, precision, recall, F1-score, and confusion matrices. We also compare the effectiveness of various classifiers, emphasizing their respective strengths and limitations. The results provide vital insights into the models' prediction powers and their prospective uses in real-world breast cancer diagnosis.

Constructing a machine learning classification algorithm capable of accurately distinguishing between cases of benign and malignant breast cancer was the aim of the study. Machine learning models can offer a solid foundation for early detection, identifying those who may benefit from specialized therapy and who are at a high risk of acquiring cancer. This validates the use of this method. To do this, we assessed eight distinct classifiers, each of which underwent intensive training and testing to guarantee dependability, using a large dataset of breast cancer patients. To ensure the robustness and generalizability of our models, cross-validation was applied during the evaluation process. Cross-validation, particularly k-fold cross-validation, is crucial in medical data analysis as it mitigates the risk of overfitting and provides a more reliable measure of model performance. Given the high dimensionality and inherent class imbalance in the breast cancer dataset, cross-validation helps to validate that the model's performance is consistent across

different subsets of the data. In this study, we employed 5-fold cross-validation, which divides the dataset into ten equal parts, iteratively training on nine parts and testing on the remaining one. This process was repeated ten times, each time with a different test subset, and the results were averaged to obtain a more reliable assessment. After implementing cross-validation, we observed an improvement in performance consistency across the models. The average accuracy for the Artificial Neural Network (ANN) remained robust at 95%, while the XGBoost model showed consistent accuracy and recall at 94% and 91%, respectively, across all folds. Cross-validation validated that ANN and XGBoost maintained their superior performance, further underscoring their potential applicability in real-world breast cancer diagnosis. The stability in performance across folds also reinforces that these models are less likely to suffer from overfitting, making them reliable for practical implementation.

Averaging the results of all iterations, key metrics such as accuracy, recall, f1, precision and MCC were used to assess each classifier. The results have been written in **table 1** and also the graphical representations have been shown in **Figure 3**. **The Artificial Neural Network (ANN) [11] outperformed the other eight classification methods in terms of overall performance, with accuracy of 95% , recall of 93% and Matthews Correlation Coefficient (MCC) of 91%.** In close pursuit, the XGBoost obtained 94% accuracy, 91% recall, and 87% MCC. Both passive aggressive and NuSVM performed worst in recall even though they were both reasonably accurate. The enormous dimensionality of the dataset and the study's inherent class imbalance are responsible for the heterogeneity in these evaluation criteria.

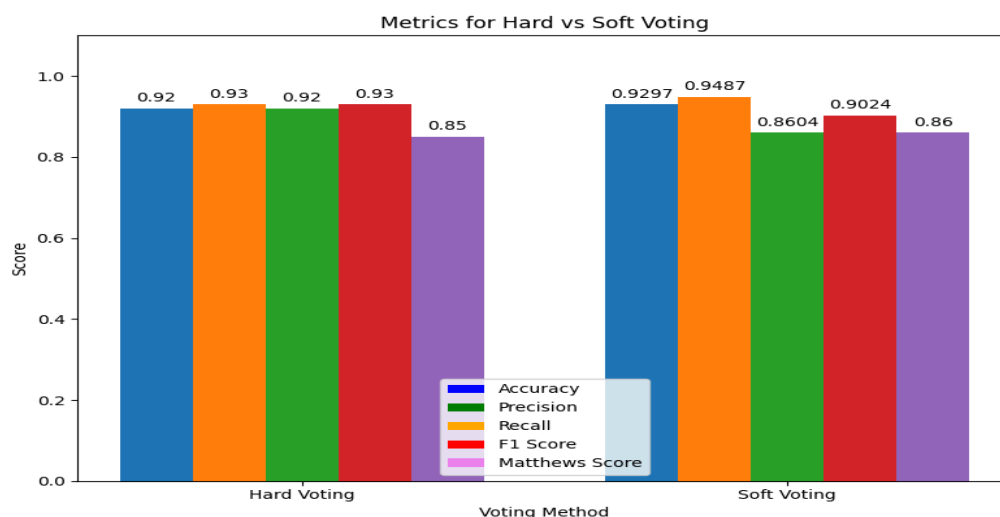


Sl. No.	Classifier	Accuracy	Precision	Recall	F1-score	MCC
1	ANN	0.955	0.945	0.938	0.940	0.906
2	XGBoost	0.939	0.920	0.919	0.919	0.871
3	Extra Tree	0.939	0.936	0.922	0.916	0.870
4	Adaboost	0.932	0.914	0.904	0.908	0.855
5	LIGHTGBM	0.928	0.958	0.881	0.909	0.863
6	Random Forest	0.928	0.921	0.885	0.901	0.847
7	NUSVM	0.910	1.000	0.791	0.871	0.832
8	Passive Aggressive	0.762	0.863	0.684	0.693	0.526
9	SVM	0.961	0.971	0.964	0.965	0.932

To further improve prediction accuracy, we implemented a voting classifier (14) using both hard and soft voting methods. In a hard voting classifier, each model votes for a specific class, and the final prediction is based on the majority vote. In contrast, soft voting considers the probability estimates from each model,

averaging them to make a final prediction. This can often yield better performance by taking into account the confidence levels of each model.

For our ensemble of top-performing models, the hard voting classifier achieved an accuracy of 92%, an F1-score of 93%, and a recall of 93%. The soft voting classifier, however, provided slightly better results, achieving an accuracy of 92.9%, an F1-score of 90%, and a recall of 94.8%. The graphical representations are given in **Figure 4**. These results suggest that soft voting is more effective for this dataset, as it leverages the probability-based confidence of each classifier, leading to higher overall performance. The ensemble method, particularly with soft voting, reinforces the reliability of predictions, making it strong for real-world breast cancer diagnosis applications.



7. References

1. S.Sadhukhan, N. Upadhyay, and P. Chakraborty, "Breast cancer diagnosis using image processing and machine learning," in *Emerging Technology in Modelling and Graphics (Advances in Intelligent Systems and Computing)*, J. K. Mandal and D. Bhattacharya, Eds. Singapore: Springer, 2020, pp. 113–127, doi: 10.1007/978-981-13-7403-6_12.
2. Bhushan A, Gonsalves A, Menon JU. Current State of Breast Cancer Diagnosis, Treatment, and Theranostics. *Pharmaceutics*. 2021 May 14;13(5):723. doi: 10.3390/pharmaceutics13050723. PMID: 34069059; PMCID: PMC8156889.
3. Ahmad, Jawad & Akram, Sheeraz & Jaffar, Arfan & Rashid, Muhammad & Masood, Sohail. (2023). Breast Cancer Detection Using Deep Learning: An Investigation Using the DDSM Dataset and a Customized AlexNet and Support Vector Machine. *IEEE Access*. PP. 1-1. 10.1109/ACCESS.2023.3311892
4. S. Guizani, N. Guizani and S. Gharsallaoui, "A Hybrid CNN-SVM Prediction Approach for Breast Cancer Ultrasound Imaging," 2023 *International Wireless Communications and Mobile*.
5. Computing (IWCMC), Marrakesh, Morocco, 2023, pp. 1574-1578, <https://doi.org/10.1109/IWCMC58020.2023.10182874>.

6. Bilal, A., Imran, A., Baig, T.I. et al. Breast cancer diagnosis using support vector machine optimized by improved quantum inspired grey wolf optimization. *Sci Rep* 14, 10714 (2024). <https://doi.org/10.1038/s41598-024-61322-w>
7. Dai, Bin et al. "Using Random Forest Algorithm for Breast Cancer Diagnosis." *2018 International Symposium on Computer, Consumer and Control (IS3C)* (2018): 449-452.
8. M. Amrane, S. Oukid, I. Gagaoua and T. Ensari, "Breast cancer classification using machine learning," *2018 Electric Electronics, Computer Science, Biomedical Engineerings' Meeting (EBBT)*, Istanbul, Turkey, 2018, pp. 1-4, doi: 10.1109/EBBT.2018.8391453.
9. Deepti Sharma, Rajneesh Kumar, Anurag Jain: Breast cancer prediction based on neural networks and extra tree classifier using feature ensemble learning, *Measurement: Sensors*, Volume 24, 2022, 100560, ISSN 2665-9174, <https://doi.org/10.1016/j.measen.2022.100560>
10. Scholkopf B, Smola AJ, Williamson RC, Bartlett PL. New support vector algorithms. *Neural Comput.* 2000 May;12(5):1207-45, doi:10.1162/089976600300015565
11. Kadhim, Rania & Kamil, Mohammed. (2023). Comparison of machine learning models for breast cancer diagnosis. *IAES International Journal of Artificial Intelligence (IJ-AI)*, 12. 415-421. [10.11591/ijai.v12.i1.pp415-421](https://doi.org/10.11591/ijai.v12.i1.pp415-421).