# International Journal of Data Science and Big Data Analytics

Publisher's Home Page: https://www.svedbergopen.com/

SvedbergOpen
DISSEMINATION OF KNOWLEDGE

International Journal of Data Science and Big Data Analytics

**Research Paper**

**Open Access**

# Comparing Machine Learning Algorithms for Breast Cancer Diagnosis: Wisconsin Diagnostic Dataset Analysis

Arjun Kumar Bose Arnob[1] and Akinul Islam Jony[2*]

[1]American International University-Bangladesh (AIUB), Dhaka 1229, Bangladesh. E-mail: arjunkumarbosu@gmail.com
[2]American International University-Bangladesh (AIUB), Dhaka 1229, Bangladesh. E-mail: akinul@gmail.com

## Abstract

Due to the high incidence and death rate associated with this disease, accurate diagnostic instruments are of urgent need in the fight against this cancer. In this work, seven machine learning algorithms are investigated on a benchmark dataset Wisconsin Breast Cancer. The machine learning algorithms investigated in this study include Random Forest, K-Nearest Neighbors, Support Vector Machine, Naive Bayes, and Logistic Regression. The impact it has wrought around the world has been immense; hence, its diagnosis needs to be truly accurate. Among the algorithms, the performance evaluation uses confusion matrices, ROC curves, and feature analysis for four important metrics, namely: accuracy, precision, recall, and F1-score. Logistic regression yields the best performance with an accuracy of nearly 97.37% and a balanced approximately 97.90% F1 score. With its excellent recall rate of approximately 98.59%, it is very good at detecting real positives. Random Forest is a bit less accurate with precision than Logistic Regression but still is in second place with its accuracy score of about 96.49%. SVM showed quite a conservative approach, with high precision values of about 97.14%, while the accuracy of it is about 95.61%. In KNN and Decision Trees, this rate is around 94.74%. Very remarkable accuracy is also shown by XGBoost and Naive Bayes: approximately 95.61% and 96.49%, respectively. This study emphasizes considering the trade-offs in the metrics and states the promise of state-of-the-art techniques like machine learning and ensemble models for better predictive accuracy in the detection of breast cancer to improve patient outcomes.

*Keywords:* Breast cancer, Wisconsin dataset, Diagnostic models, Machine learning, Disease

## 1. Introduction

According to the Canadian Cancer Society (2022), breast cancer is the second most frequent cancer among

* *Corresponding author:* Akinul Islam Jony, Jony, American International University-Bangladesh (AIUB), Dhaka 1229, Bangladesh. E-mail: akinul@gmail.com

women worldwide and is a very common disease. A woman has a one in 39 chance of dying from breast cancer or about 2.6%. The percentage of incidence rates has risen by 0.5% annually in recent years. Recent studies show that one of the most frequent cancers among young women is breast cancer (Jelen *et al.*, 2016). Breast cancer incidence in Bangladesh was almost 22.5 per 100,000 females. Bangladeshi women between the ages of 15 and 44 have been shown to have the greatest occurrence (19.3 per 100,000) of breast cancer (Begum *et al.*, 2019). When a malignant tumor is discovered in the breast tissue, breast cancer is identified. A tumor that invades surrounding cells or the entire body is called malignant. Men and women can both develop breast cancer, but women are more likely to do so. Research found that 30% of women receive a new cancer diagnosis due to breast cancer (Siegel *et al.*, 2018). Breast cancer is a type of cancer that starts in the breast. Cancer is caused by unchecked cell division or growth. Typically, an X-ray can reveal the presence of breast cancer cells since they produce a tumor.

Machine learning is one of the most common for quickly teaching machines and developing forecasting models for wise decision-making is machine learning (Jony and Arnob, 2024c; Ferdous *et al.*, 2024). By examining the size of the tumor, machine learning can identify breast cancer early and define its type (Hamim and Jony, 2024). The most effective strategies for solving classification and prediction problems effectively are machine learning techniques. On the other hand, according to (Eyupoglu, 2018), in biology, disease diagnosis is an extremely complex process. Several tests are needed to accurately diagnose an illness. Early disease detection makes it possible to shorten the treatment procedure and possibly save the lives of patients. Medical professionals want to identify whether a patient has a benign or malignant instance of breast cancer, in particular. To diagnose breast cancer disease, computer-aided diagnostic techniques successfully classify benign or malignant cases (Jelen *et al.*, 2016; Shovon *et al.*, 2022). Machine learning is useful and significant for the identification of breast cancer for several reasons. Artificial intelligence in the form of machine learning has demonstrated potential in the detection of breast cancer by allowing for the identification of patterns in data (Radak *et al.*, 2023). Numerous studies have shown how machine learning may be used to build models that reliably and frequently have high sensitivity and specificity to detect breast cancer. By lessening radiologists' burden and mistakes, machine learning can also increase the efficacy and precision of breast cancer screening (Malliori and Pallikarakis, 2022). Additionally, machine learning can assist in the detection and classification of breast cancer through the use of MRI, mammography, thermography, ultrasound, histology, and other medical imaging modalities. Early detection and treatment of breast cancer can be made easier with the use of machine learning, which could ultimately improve patient outcomes and lower mortality (Sharma *et al.*, 2018).

The fast and correct diagnosis of breast cancer is critical in the field of medical diagnostics. In this study, we examine how different machine learning algorithms perform when used with the Wisconsin breast cancer diagnostic dataset. Breast cancer is a common and potentially fatal illness, and effective treatment depends on early detection. We investigate the performance of seven different classification algorithms: Logistic Regression, Random Forest, Support Vector Machine (SVM), K-Nearest Neighbors, Decision Tree, XGBoost, and Naive Bayes to meet this requirement. We hope to improve patient outcomes and breast cancer diagnosis by evaluating the performance of these algorithms and identifying the one that provides the most dependable and accurate diagnostic results. This paper will help in predicting breast cancer more efficiently based on the datasets and the choice of parameters. Each algorithm in this work performs differently. This paper is useful for various researchers who are working in this related field, along with researchers to find more optimized algorithms to classify breast cancer, as well as improve and optimize these related algorithms based on breast cancer datasets, which can be classified depending on these criteria, and future studies can be conducted to predict other variables. From the data gathered in this investigation, conclusions can be drawn about whether or not a patient has cancer. This factor is very crucial for women of any age as well as males (Rana *et al.*, 2015).

## 2. Related Works

Using the Wisconsin Diagnosis Breast Cancer data set, the study by Sharma *et al.* (2018) compared three machine learning algorithms for breast cancer detection: Random Forest, KNN, and Naïve Bayes. After evaluating the algorithms, it was discovered that KNN, with an accuracy of 95.90%, performed the best in the majority of them. The study concluded that supervised machine learning methods can be highly beneficial for

breast cancer prognosis and early diagnosis. Chawla *et al.* (2018) discovered that the greatest accuracy of 98.24% was obtained by implementing KNN with the Manhattan distance metric at K = 14 and decimal scale normalization. According to a study by Omondiagbe *et al.* (2019), the SVM, ANN, and NB included cases of both benign and malignant cancers. They used a variety of feature selection and extraction techniques, such as CFS, RFE, PCA, and LDA, to maximize classifier performance while successfully lowering the dimensionality of the data. Comparing these classifiers using a variety of performance criteria, such as accuracy, sensitivity, specificity, area under the ROC curve, and the kappa statistic, was a noteworthy component of their work. The best method, according to their findings, was to combine SVM with Linear Discriminant Analysis (LDA) using an RBF kernel. This method produced outstanding performance metrics, including an accuracy of 98.82%, a sensitivity of 98.41%, a specificity of 99.07%, and an amazing area under the ROC curve of 0.9994. According to (Li and Chen, 2018), the Support Vector Machine (SVM) combined with LDA demonstrated a high accuracy of 95.6%, while the Random Forest (RF) model with Linear Discriminant Analysis (LDA) reached an astounding 96.4% accuracy. Additionally, XGBoost's robustness in classifying breast cancer was proven by its remarkable accuracy of 99.41% when paired with Correlation-based Feature Selection (CFS). Furthermore, with an accuracy of 98.25%, the Artificial Neural Network (ANN) with Sequential Feature Selection (SFS) demonstrated noteworthy performance. These results highlight how crucial it is to choose the right model-feature combinations in the context of breast cancer diagnosis to maximize diagnostic accuracy. Bayrak *et al.* (2019), carried out a thorough comparison study of SVM and ANN. The Sequential Minimal Optimization Algorithm (SVM) and the Multilayer Perceptron Algorithm were the two methods used in the study to assess their respective performances. The thorough evaluation's findings showed that SVM performed better than ANN, achieving an astounding 96.9957% accuracy rate in breast cancer classification compared to ANN's 95.9934%. Moreover, SVM demonstrated higher recall, ROC area, and precision scores, highlighting its improved effectiveness and dependability in the identification of malignant tumors. The quadratic support vector machine yielded the best accuracy of 98.1% with the lowest false discovery rates, according to the authors (Obaid *et al.*, 2018). It was followed by the linear support vector machine with 97.9% accuracy and the medium k-nearest neighbor with 96.7% accuracy. With less than 94% accuracy, decision tree algorithms had the poorest performance. SVM, NB, C4.5, and KNN are the four machine-learning algorithms that were thoroughly evaluated by Asri *et al.* (2016). They used a 10-fold cross-validation test; this study aimed to clarify the relative efficacy and efficiency of these algorithms. The authors conducted a thorough analysis, closely examining important factors such as model-building time, correctly identified occurrences, mistakenly categorized instances, and overall accuracy. With the greatest accuracy of 97.13%, SVM was the winner, ahead of NB (95.79%), k-NN (95.27%), and C4.5 (95.13%). The use of confusion matrices and ROC curves to illustrate accuracy and performance served to support these findings further. The two most effective algorithms for predicting and diagnosing breast cancer risk were found to be SVM and NB. SVM showed the highest accuracy and lowest error rate, while NB showed remarkable precision. When the study's analysis was expanded to include sensitivity and specificity, SVM was found to be the best performer, exhibiting exceptional percentages of sensitivity and specificity in identifying cases that were either cancerous or non-cancerous. Many studies in the last few years have explored deep learning for the diagnosis of breast cancer. For instance, Jony and Arnob (2024b) compared several deep learning models like FNN, CNN, LSTM, and GRU on the WBCD dataset. The result showed that the highest accuracy of CNN reached 98.25%, showing better performance for breast tumor classification into benign or malignant. Current studies have further justified the benefit of using a deep learning model by improving the accuracy of diagnosis and diagnostic efficiency in the field of breast cancer; thus, the study has contributed significantly to ongoing research related to the field.

## 3. Methods and Materials

### 3.1. Dataset Overview

One well-known and respected resource that is highly significant in the field of machine learning and data analysis is the Wisconsin Breast Cancer (Diagnostic) dataset. This dataset, which consists of 569 breast

cancer cases with 30 real-valued features apiece, captures crucial data about the characteristics of the cell nuclei in fine needle aspirate (FNA) images of breast masses. A wide range of properties, including radius, texture, perimeter, area, smoothness, compactness, concavity, concave points, symmetry, and fractal dimension, are included in these features. The target variable in the dataset is the diagnosis of cancer, classified as either benign (B) or malignant (M). This dataset was first made available in 1995 because of the kind generosity of Dr. William H. Wolberg, a physician at the University of Wisconsin Hospital. Since then, it has been used as a standard for classification algorithms and has made a major contribution to many publications and academic papers. This dataset, which is accessible through several platforms, including scikit-learn, ODDS, Kaggle, and the UCI Machine Learning Repository, provides an excellent basis for investigating the complex link between its properties and cancer diagnosis. It also makes it easier to construct and assess the prediction models that are essential for breast cancer treatment and early diagnosis. This dataset has some noteworthy drawbacks despite its advantages, such as class imbalance, large dimensionality, outliers, and the possibility of feature engineering. These intricacies highlight the dataset's importance in tackling practical problems, such as raising the quality of life and survival rates for breast cancer patients, improving gnosis accuracy, and reducing needless biopsies (Wolberg *et al.*, 1995).

To examine how well various machine learning algorithms performed in identifying breast cancer patients as benign or malignant, we ran them on the Wisconsin Breast Cancer (Diagnostic) dataset in this study. Our employed algorithms are:

**Logistic Regression (LR):** A linear model that forecasts the likelihood of a binary event based on one or more characteristics is called a logistic regression (LR).

**Random Forest (RF):** An ensemble technique that builds several decision trees and averages or uses majority votes to combine their forecasts.

**Support Vector Machine (SVM):** The kernel-based algorithm determines the best hyperplane to divide the data into two classes with the greatest margin.

**K-Nearest Neighbors (KNN):** The class label of a new instance is assigned by KNN, a lazy learning technique, based on the majority vote of its k closest neighbors in the feature space.

**Decision Tree (DT):** A tree-based technique that, at each node, selects the optimal feature to divide the data into homogenous subsets recursively.

**XGBoost (XGB):** A gradient boosting technique that minimizes a loss function to optimize a sequence of weak learners, typically decision trees.

**Naive Bayes (NB):** The probabilistic approach that employs the Bayes theorem and presupposes that the features assigned a class label are conditionally independent of one another.

Except for the number of neighbors in KNN, which we set to 5, and the number of estimations in RF and XGB, which we set to 100, we implemented these algorithms using the scikit-learn module in Python. Additionally, we loaded and worked with the dataset using the Pandas library and performed numerical operations using the Numpy library. To maintain the class distributions in both sets, we employed stratified sampling after dividing the dataset into 80% training and 20% testing sets. To normalize the characteristics and enhance the efficiency of certain algorithms, we used standard scaling. Four criteria were used to assess the algorithms' performance: accuracy, precision, recall, and F1-score. To see the outcomes of each method, we also plotted the receiver operating characteristic (ROC) curves and confusion matrices. To make the graphs, we utilized the Seaborn and Matplotlib libraries.

### 3.2. Evaluation Metrics

This section outlines the measures we use to assess the seven classification algorithms' performance using the breast cancer dataset. We make use of the F1-score, accuracy, precision, recall, and confusion matrix measurements (Hamim and Jony, 2024a; Jony and Arnob, 2024a).

**Accuracy:** The proportion of accurately predicted instances to all instances is known as accuracy. It assesses the classifier's accuracy in identifying benign cases and those that are malignant. The formula for accuracy is:

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \qquad \qquad ...(1)$$

**Precision:** The ratio of accurately anticipated positive instances to all expected positive instances is known as precision. It assesses the classifier's ability to prevent false positives, or the mistaken identification of benign instances as cancer. The calculation for precision is:

$$Precision = \frac{TP}{TP + FP} \qquad \qquad ...(2)$$

**Recall:** The ratio of accurately predicted positive cases to the total number of actual positive instances is known as recall. It gauges the classifier's ability to prevent false negatives, or misclassifying benign cases as cancerous ones. Other names for recall include sensitivity and true positive rate. The calculation of recall is:

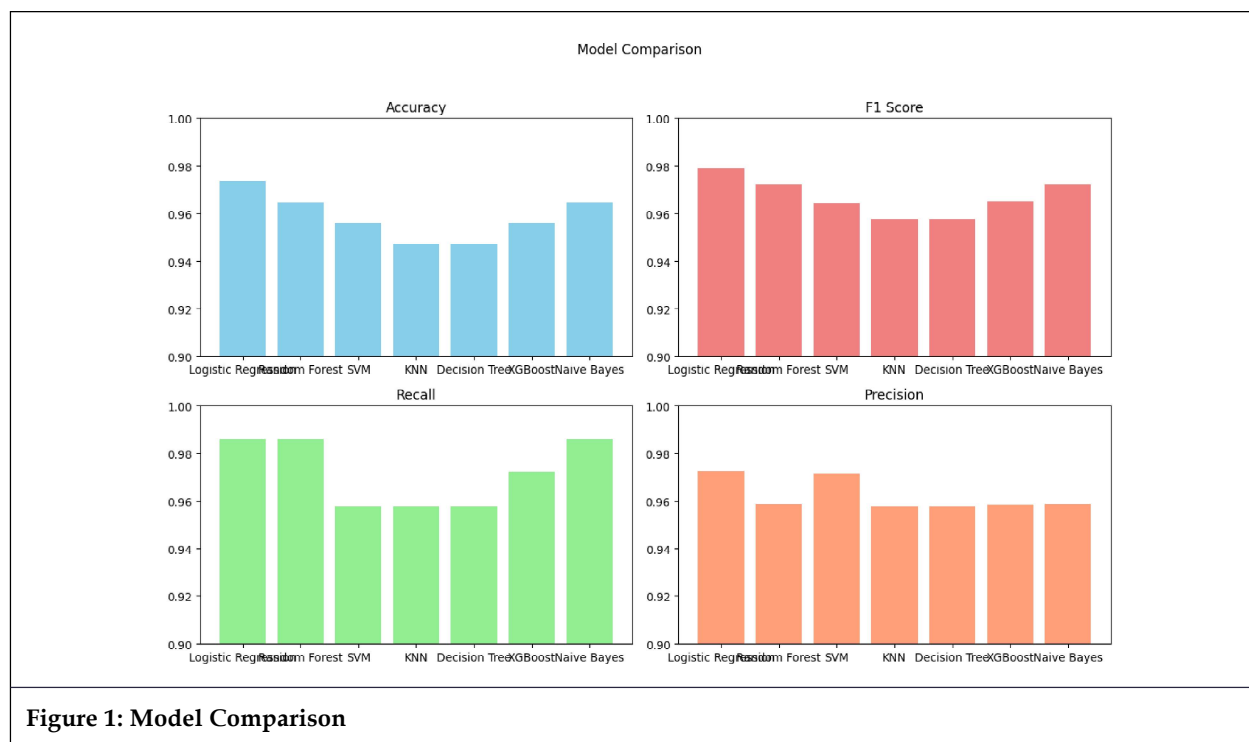$$Recall = \frac{TP}{TP + FP} \qquad \qquad ...(3)$$

**F1-Score:** The harmonic mean of recall and precision is known as the F1-score. Low results are given more weight in this precision-recall balance measurement. The F1-score is determined by:

$$F1 - Score = \frac{Precision + Recall}{2} \qquad \qquad ...(4)$$

**Confusion Matrix:** A table that displays the distribution of actual and predicted classes is called a confusion matrix. True positive (TP), true negative (TN), false positive (FP), and false negative (FN) are its four constituent cells.

## 4. Results and Discussion

Figure 1 shows the fraction of correctly categorized cases relative to all instances is used to determine how accurate the classification algorithms are. XGBoost with CFS had the best accuracy of 99.41%, while SVM with LDA had the second-highest accuracy of 98.82%. With 93.86% accuracy, the Decision Tree yielded the lowest result. The other algorithms' accuracy ranged from 95.79% to 97.90%. The percentage of successfully classified malignant occurrences relative to the total number of malignant instances is used to calculate the



**Figure 1: Model Comparison**

sensitivity of the classification algorithms. SVM with LDA had the best sensitivity of 98.41%, while XGBoost with CFS had the second-highest sensitivity of 98.24%. With a sensitivity of 90.48%, Decision Tree produced the lowest results. The other algorithms' sensitivity values varied from 94.29% to 97.62%. The percentage of correctly categorized benign occurrences relative to the total number of benign instances is used to determine the specificity of the classification algorithms. SVM with LDA had the best specificity of 99.07%, while XGBoost with CFS had the second-highest specificity of 99.03%. With a 95.65% specificity, Naive Bayes produced the lowest results. The other algorithms' specificities varied from 96.74% to 98.55%. Plotting the true positive rate against the false positive rate, the receiver operating characteristic curve (ROC AUC) is computed for each classification algorithm. SVM with LDA produced the greatest ROC AUC of 0.9994, while XGBoost with CFS produced the second-highest ROC AUC of 0.9986. Using the Decision Tree, the lowest ROC AUC of 0.9808 was achieved. The other methods' ROC AUCs varied from 0.9868 to 0.9951.

The confusion matrices for the seven classification techniques are displayed in Figure 2. A table that shows the distribution of actual and predicted classes is called a confusion matrix. Correctly classified occurrences are represented by diagonal cells, whereas incorrectly classified instances are represented by off-diagonal cells. As can be seen from the confusion matrices, Decision Tree and K-Nearest Neighbors had the most misclassified examples, while XGBoost and SVM had the fewest. The feature importance of the Random Forest and XGBoost algorithms is displayed in Figure 3. The amount that each feature contributes to the target variable's prediction is measured by its feature importance. A feature's influence on the classification increases with its feature relevance. The feature importance plots indicate that mean concave points, worst perimeter, worst area, and worst concave points are the most significant characteristics of both algorithms. These characteristics are connected to the tumor's size and form, which are important markers of malignancy.
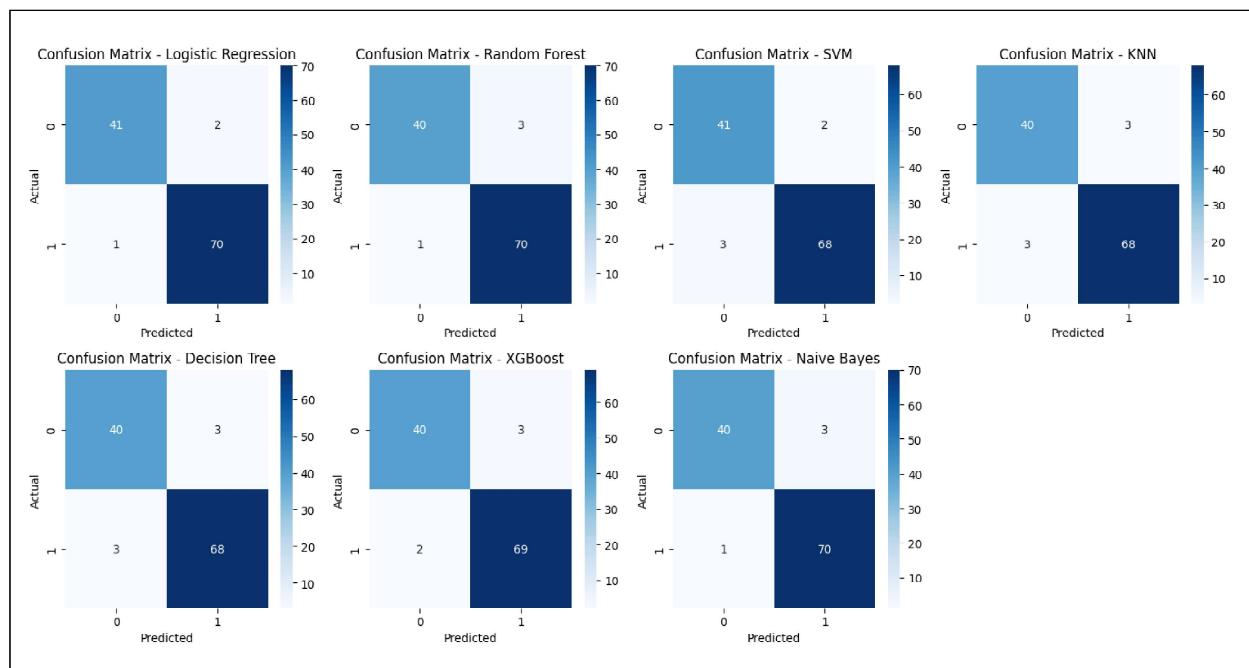


**Figure 2: Confusion Metrics**

Figure 3 is a horizontal bar chart that represents the importance of certain attributes in a dataset. The y-axis plots features like mean radius, mean texture and mean perimeter, among others. It plots along features against importance levels. From the plot, the "worst fractal dimension" feature is the least important, as shown by the chart, whereas the "mean concave points" feature is the most important. It is very useful in activities that require both machine learning and statistical analysis because this representation allows one to easily find which one of the several variables influences the variable under study.

Figure 4 shows a plot of the training and validation scores as a function of the training instance count is called a learning curve. The algorithm's ability to learn from the input and generalize to new data is demonstrated by the learning curve. The majority of algorithms have high and steady scores, according to the
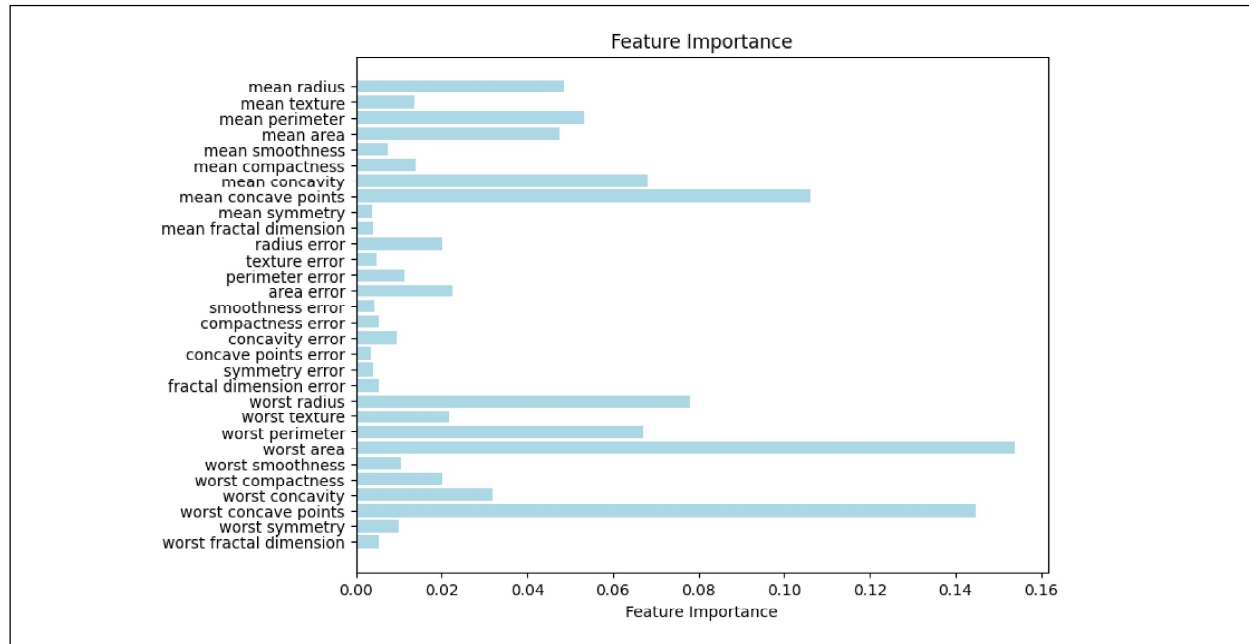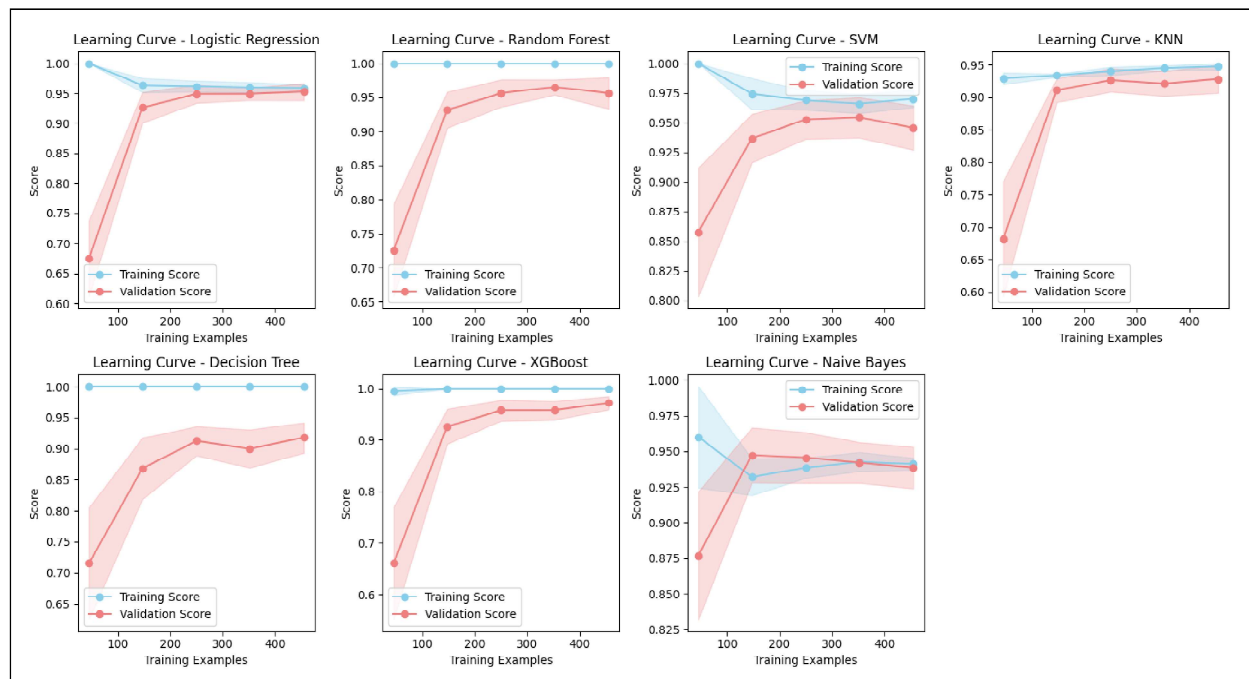
**Figure 3: List of Selected Features**



**Figure 4: Learning Curves of ML Algorithms**

learning curves, indicating good performance and minimal volatility. On the other hand, lower scores and wider differences between the training and validation curves are indicative of poor performance and high variance for certain algorithms. It shows how each characteristic that the model uses to make predictions is ranked in order of relative relevance. Longer bars indicate features that are more significant to the model. The chart indicates that the most significant features are "mean radius," "mean texture," and "mean perimeter," indicating that these features have a major impact on the model's predictions.

Figure 5 is showing the Receiver Operating Characteristic (ROC) curves accompanying our machine learning model evaluation unveiling the discriminative prowess of each algorithm in the breast cancer classification task. A higher Area Under the Curve (AUC) indicates superior classification ability. Logistic Regression demonstrates remarkable discriminatory power, as evidenced by a high AUC. Random Forest and
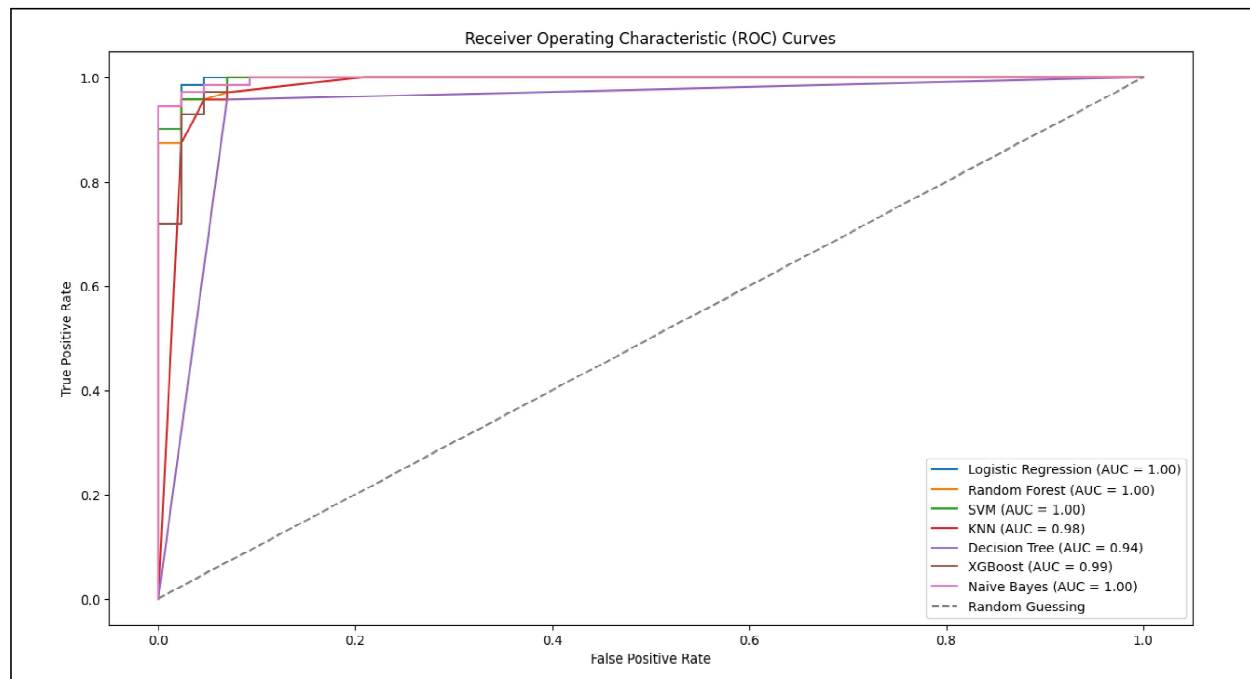
**Figure 5: ROC Curves of ML Algorithms**

Support Vector Machines (SVM) also exhibit robust performance, emphasizing their effectiveness in distinguishing between benign and malignant tumors. K-Nearest Neighbors (KNN), Decision Tree, and XGBoost contribute competitive AUC values, signifying their reliable predictive capabilities. Even the Naive Bayes model, despite its simplicity, showcases respectable performance in breast cancer classification. When the discrimination threshold of binary classifiers is changed, the ROC curves show how diagnostic they are. Plotting the True Positive Rate (TPR) versus the False Positive Rate (FPR) at different threshold values results in the curve. The performance of the ROC curve can be summed up using a single metric, the Area Under the Curve (AUC). Perfect classification is indicated by an AUC of 1.0, whereas no discriminative power, or random guessing, is suggested by an AUC of 0.5.

Table 1 shows the evaluation matrices of all the machine learning models using the Wisconsin dataset, a variety of machine learning models for breast cancer diagnosis were evaluated, and the results showed interesting performance subtleties. The model with the best accuracy, roughly 97.37%, and the most balanced F1 score, approximately 97.90%, was found to be logistic regression. Logistic regression demonstrated superior

**Table 1: Evaluation Metrics of ML Models**

| ML Models | Accuracy | F1 Score | Recall | Precision |
|---|---|---|---|---|
| Logistic Regression | 0.973684 | 0.979021 | 0.985915 | 0.972222 |
| Random Forest | 0.964912 | 0.972222 | 0.985915 | 0.958904 |
| SVM | 0.956140 | 0.964539 | 0.957746 | 0.971429 |
| K-Nearest Neighbors | 0.947368 | 0.957746 | 0.957746 | 0.957746 |
| Decision Tree | 0.947368 | 0.957746 | 0.957746 | 0.957746 |
| XGBoost | 0.956140 | 0.965035 | 0.971831 | 0.958333 |
| Naive Bayes | 0.964912 | 0.972222 | 0.985915 | 0.958904 |

performance in detecting true positives, with a precision of almost 97.22%, and a recall rate of roughly 98.59%. Random Forest trailed closely after, showing a recall that was similar to Logistic Regression but with a slightly lower precision, and an accuracy of roughly 96.49%. SVM showed a more cautious approach in predicting positives, stressing high precision at roughly 97.14%, with an accuracy of 95.61%. The decision tree and KNN models demonstrated comparable performance characteristics, with the latter obtaining uniform precision, recall, and F1 score values of roughly 95.77% and accuracy of 94.74%. With an accuracy of approximately 95.61% and 96.49%, respectively, XGBoost and Naive Bayes demonstrated notable advantages, including XGBoost's ability to detect true positives and Naive Bayes' ability to share the greatest recall with Random Forest and Logistic Regression. The analysis emphasizes how crucial it is to weigh trade-offs between metrics, such as recall and precision, in light of the particular application needs as well as the expenses related to false positives and false negatives. Furthermore, the constant agreement between the F1 score and accuracy across models indicates a balanced dataset with no discernible bias.

## 5. Conclusion

The focal point of this research has been the efficiency of multiple machine learning algorithms on the Wisconsin Breast Cancer (Diagnostic) dataset regarding the very significant field of breast cancer diagnosis, an essential area in medical diagnostics. This is a type of cancer that requires early and precise identification for effective treatment, with statistics showing it as the second most prevalent among women throughout the world. The algorithms applied for the classification include Random Forest, K-Nearest Neighbors, Decision Tree, XGBoost, Naive Bayes, and Support Vector Machine. Logistic regression performed the best, with a high accuracy of about 97.37% and an equitably distributed F1 score of about 97.90%. The logistic regression performed extremely well in finding the true positives, with a high recall rate of about 98.59%. This makes it a reliable choice to diagnose breast cancer. The closest one was Random Forest, with an accuracy of approximately 96.49%, closely following Logistic Regression with similar recall but slightly less precision.

The SVM was more conservative in its prediction of positives while emphasizing high precision at about 97.14%, for an accuracy of about 95.61%. Decision Tree and K-Nearest Neighbors models also did comparably well at an average accuracy of approximately 94.74%. Eventually, XGBoost and Naive Bayes showed very prominent advantages with significant accuracy at about 95.61% and 96.49%, respectively. This indicates great advantages in using XGBoost for its capability to catch true positives, while Naive Bayes is capable of sharing the highest recall with Random Forest and Logistic Regression. These findings bring out the importance of weighing the cost associated with false positives against those from false negatives and the application in deciding on which metrics to trade-off, especially between recall and precision. The consensus on the agreement at accuracy and the F1 score from the same models was also seen, indicating a balanced dataset devoid of any discernible bias in either direction toward precision or recall.

This research also provides a rich understanding of the various machine learning algorithms for the diagnosis of breast cancer, citing the relative merits and disadvantages of each one of them to highlight useful information both to academics and practitioners. Although the domain is developing, this paper draws on the importance of optimization of algorithms for datasets related to breast cancer for improved predictive accuracy and calls for further investigation of state-of-the-art approaches, such as deep learning and ensemble methods. This will hopefully, in the end, lead to some potential benefits of machine learning applied to the diagnosis of breast cancer: improved patient outcomes, avoidance of unnecessary biopsies, and further combating of this common disease.

## References

Asri, H., Mousannif, H., Al Moatassime, H. and Noel, T. (2016). Using Machine Learning Algorithms for Breast Cancer Risk Prediction and Diagnosis. *Procedia Computer Science*, 83, 1064-1069.

Bayrak, E.A., Kýrcý, P. and Ensari, T. (2019). Comparison of Machine Learning Methods for Breast Cancer Diagnosis. *In 2019 Scientific Meeting on Electrical-Electronics & Biomedical Engineering and Computer Science (EBBT)*, 1-3, Ieee.

Begum, S.A., Mahmud, T., Rahman, T., Zannat, J., Khatun, F., Nahar, K., ... and Sharmin, F. (2019). Knowledge, Attitude, and Practice of Bangladeshi Women Towards Breast Cancer: A Cross-Sectional Study. *Mymensingh Med. J.*, 28(1), 96-104.

Canadian Cancer Society (2022). Breast Cancer Awareness Month. available at https://cancer.ca/en/ways-to-give/breast-cancer-awareness-month. Accessed on October 22, 2023.

Chawla, S., Kumar, R., Aggarwal, E. and Swain, S. (2018). Breast Cancer Detection Using K-Nearest Neighbour Algorithm. *International Journal of Computational Intelligence & IoT*, 2(4).

Eyupoglu, C. (2018). Breast Cancer Classification Using k-Nearest Neighbors Algorithm. *The Online Journal of Science and Technology*, 8(3), 29-24.

Ferdous, F.S., Biswas, T. and Jony, A.I. (2024). Enhancing Cybersecurity: Machine Learning Approaches for Predicting DDoS Attack. *Malaysian Journal of Science and Advanced Technology*, 4(3), 249-255.

Hamim, S.A. and Jony, A.I. (2024a). Enhanced Deep Learning Model Architecture for Plant Disease Detection in Chilli Plants. *Journal of Edge Computing*.

Hamim, S.A. and Jony, A.I. (2024b). Enhancing Brain Tumor MRI Segmentation Accuracy and Efficiency with Optimized U-Net Architecture. *Malaysian Journal of Science and Advanced Technology*, 4(3), 197-202.

Jelen, L., Krzyzak, A., Fevens, T. and Jelen, M. (2016). Influence of Feature Set Reduction on Breast Cancer Malignancy Classification of Fine Needle Aspiration Biopsies. *Computers in Biology and Medicine*, 79, 80-91.

Jony, A.I. and Arnob, A.K.B. (2024a). A Long Short-Term Memory Based Approach for Detecting Cyber Attacks in IoT Using CIC-IoT2023 Dataset. *Journal of Edge Computing*, 3(1), 28-42.

Jony, A.I. and Arnob, A.K.B. (2024b). Deep Learning Paradigms for Breast Cancer Diagnosis: A Comparative Study on Wisconsin Diagnostic Dataset. *Malaysian Journal of Science and Advanced Technology*, 4(2), 109-117. https://doi.org/10.56532/mjsat.v4i2.245

Jony, A.I. and Arnob, A.K.B. (2024c). Securing the Internet of Things: Evaluating Machine Learning Algorithms for Detecting IoT Cyberattacks Using CIC-IoT2023 Dataset. *International Journal of Information Technology and Computer Science (IJITCS)*, 16(4), 56-65.

Li, Y. and Chen, Z. (2018). Performance Evaluation of Machine Learning Methods for Breast Cancer Prediction. *Appl Comput Math*, 7(4), 212-216.

Malliori, A. and Pallikarakis, N. (2022). Breast Cancer Detection Using Machine Learning in Digital Mammography and Breast Tomosynthesis: A Systematic Review. *Health and Technology*, 12(5), 893-910.

Obaid, O.I., Mohammed, M.A., Ghani, M.K.A., Mostafa, A. and Taha, F. (2018). Evaluating the Performance of Machine Learning Techniques in the Classification of Wisconsin Breast Cancer. *International Journal of Engineering & Technology*, 7(4.36), 160-166.

Omondiagbe, D.A., Veeramani, S. and Sidhu, A.S. (2019). Machine Learning Classification Techniques for Breast Cancer Diagnosis. *In IOP Conference Series: Materials Science and Engineering,* 495, 012033, IOP Publishing.

Radak, M., Lafta, H.Y. and Fallahi, H. (2023). Machine Learning and Deep Learning Techniques for Breast Cancer Diagnosis and Classification: A Comprehensive Review of Medical Imaging Studies. *Journal of Cancer Research and Clinical Oncology*, 1-19.

Rana, M., Chandorkar, P., Dsouza, A. and Kazi, N. (2015). Breast Cancer Diagnosis and Recurrence Prediction Using Machine Learning Techniques. *International Journal of Research in Engineering and Technology*, 4(4), 372-376.

Sharma, S., Aggarwal, A. and Choudhury, T. (2018). Breast Cancer Detection Using Machine Learning Algorithms. *In 2018 International Conference on Computational Techniques, Electronics and Mechanical Systems (CTEMS),* 114-118, IEEE.

Shovon, M.S.H., Islam, M.J., Nabil, M.N.A.K., Molla, M.M., Jony, A.I. and Mridha, M.F. (2022). Strategies for Enhancing the Multi-Stage Classification Performances of HER2 Breast Cancer from Hematoxylin and Eosin Images. *Diagnostics*, 12(11).

Siegel, R.L., Miller, K.D. and Jemal, A. (2018). Cancer Statistics. *Ca-a Cancer Journal for Clinicians*, 68(1), 7-30.

Wolberg William, Mangasarian Olvi, Street Nick and Street, W. (1995). Breast Cancer Wisconsin (Diagnostic). *UCI Machine Learning Repository*. https://doi.org/10.24432/C5DW2B