

# Decoding Breast Cancer: A Machine Learning Perspective

Likhitha Marrapu  
Department of Data Science

Kent State University  
Kent OH 44240, USA  
lmarrapu@kent.edu

**Abstract**—This study investigates the effectiveness of various machine learning algorithms in classifying breast cancer cases, focusing on the impact of Principal Component Analysis (PCA) on model performance. We evaluated Logistic Regression, Support Vector Machine (SVM), Perceptron, Random Forest, Gradient Boosting Classifier, and Decision Tree models using key metrics, including accuracy, precision, sensitivity, and specificity. Before PCA, both Logistic Regression and SVM demonstrated superior accuracy rates exceeding 96%, excelling in identifying malignant cases. After applying PCA, significant enhancements were observed, particularly in the Logistic Regression and Perceptron models, which achieved an accuracy of 97.8% and perfect specificity. Our findings highlight the advantages of dimensionality reduction in optimizing model performance while underscoring the need for careful selection of algorithms to ensure reliable breast cancer diagnostics. This research contributes valuable insights for clinicians aiming to leverage machine learning tools in enhancing diagnostic accuracy and patient outcomes.

**Index Terms**—Breast Cancer, Machine Learning, Classification Models, PCA, Logistic Regression, SVM, Perceptron, Random Forest, Gradient Boosting, Decision Tree, Accuracy, Precision, Sensitivity, Specificity, Medical Diagnostics

## I. INTRODUCTION

Breast cancer continues to be one of the most common cancers globally, impacting millions and presenting a major public health concern. The World Health Organization (WHO) reports that breast cancer represents approximately 25 % of all cancer diagnoses among women, underscoring the critical need for improved methods of detection and treatment. The challenges of breast cancer stem from its biological complexity, as well as the wide range of factors -genetic, environmental, and lifestyle-related-that contribute to its onset and progression.

Recent developments in data science and machine learning have created new opportunities for improving the diagnosis and prognosis of breast cancer. The availability of comprehensive datasets that include diverse clinical and pathological features enables researchers to leverage these technologies to uncover patterns and relationships that traditional analytical methods might overlook. Specifically, machine learning algorithms have demonstrated significant potential in classifying different types of breast cancer, forecasting patient outcomes, and tailoring treatment strategies to individual needs.

In this research, we concentrate on the application of machine learning techniques to evaluate a comprehensive

breast cancer dataset, which encompasses a diverse array of information critical for understanding tumor biology and patient health. This includes various attributes that describe the morphological features of tumors, which are essential for distinguishing between malignant and benign cases. By utilizing algorithms such as Random Forest, Support Vector Machines, and Gradient Boosting, we aim to create predictive models that enhance the accuracy of breast cancer diagnoses and help identify patients at high risk who may require more aggressive treatment options. Additionally, our study intends to assess the performance of various machine learning methods, elucidating their respective strengths and limitations within clinical contexts. Through this thorough examination, we hope to provide valuable insights that can aid in clinical decision-making, ultimately improving patient outcomes and furthering advancements in oncology.

## II. METHODOLOGY

### A. Dataset

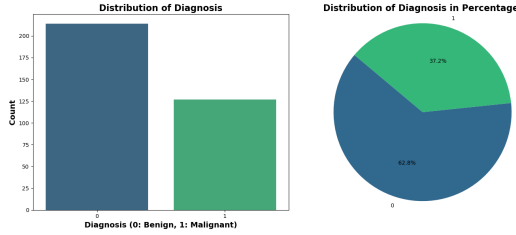
This research leverages a prominent breast cancer dataset sourced from the UCI Machine Learning Repository, comprising 569 instances that encompass a variety of features crucial for breast cancer diagnosis. Each instance includes detailed clinical attributes that describe tumor morphology, which are essential for distinguishing between malignant and benign cases. The dataset features a diverse range of numerical features that quantify tumor characteristics, such as radius, texture, perimeter, area, smoothness, compactness, concavity, and symmetry, all of which play a vital role in accurate classification.

In addition to the clinical attributes, the dataset includes categorical variables such as the diagnosis label, which identifies the nature of the tumor as either malignant (cancerous) or benign (non-cancerous). The features are meticulously curated to facilitate the application of machine learning techniques, enabling the development of models that can effectively predict patient outcomes based on these parameters. By employing advanced machine learning techniques on this dataset, the research aims to enhance diagnostic accuracy and provide valuable insights into the factors influencing breast cancer, ultimately contributing to better patient care and treatment strategies.

## B. Data Pre-processing

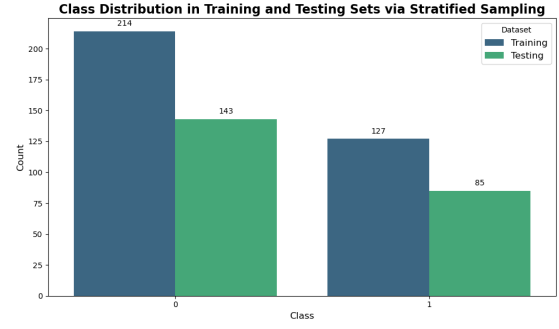
In the data preprocessing stage, we began by encoding categorical variables to ensure effective machine learning model training. Specifically, the 'Diagnosis' column, which classifies tumors as malignant or benign, was subjected to label encoding. This process assigns unique integers to each category, enabling algorithms to interpret categorical data numerically. After this, we focused on maintaining the dataset's integrity by checking for missing values, which can significantly hinder model performance. Instances with missing data were carefully examined, and appropriate handling methods were employed. We also addressed outlier management through Winsorization, a technique that replaces extreme values with the nearest valid data points, thereby reducing the adverse effects of outliers on model training.

To prepare the dataset for model training, we employed stratified random sampling to split the data into training and testing sets. Following this, we applied feature scaling through standardization, which adjusts the features to have a mean of zero and a standard deviation of one, thus ensuring that all features contribute equally to model performance. Additionally, we conducted a correlation analysis to identify relationships between features. Collectively, these preprocessing steps were essential in enhancing the quality and readiness of the dataset for effective machine learning applications.



## C. Data Partitioning

To guarantee a comprehensive evaluation of the predictive models, the dataset was divided into training and testing subsets using a 60:40 ratio. This method employs stratified random sampling, which maintains the proportional representation of each class within the target variable. By doing so, we effectively address the challenges posed by class imbalances, thereby improving the model's ability to generalize across various breast cancer categories. This meticulous division not only enables a thorough assessment of model performance but also enhances the reliability and validity of the study's findings.



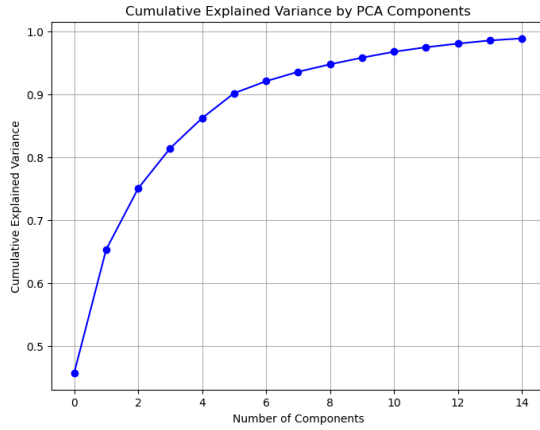
## D. Model Selection

In the model selection phase, several machine learning algorithms were evaluated to determine the most effective approach for classifying breast cancer instances based on the dataset's features. The primary models selected for this study included Logistic Regression, Random Forest, Decision Tree, Support Vector Machine (SVM), Perceptron, and Gradient Boosting Classifier. Each of these models offers unique strengths and operational characteristics, making them suitable for tackling the classification problem.

Logistic Regression was chosen for its simplicity and interpretability, providing a solid baseline model for binary classification tasks. Random Forest and Decision Tree models were included for their ability to handle complex relationships and interactions between features, as well as their robustness against overfitting. The Support Vector Machine (SVM) was selected due to its effectiveness in high-dimensional spaces and its capability to construct non-linear decision boundaries through kernel functions. Additionally, the Perceptron was incorporated to explore a foundational neural network model that captures linear separability in the data. Lastly, the Gradient Boosting Classifier was utilized for its high predictive accuracy and efficiency in managing both classification and regression tasks. By employing this diverse set of models, the research aimed to identify the most reliable classifier for predicting breast cancer outcomes, ensuring a comprehensive analysis of their performance through various evaluation metrics.

## E. Dimensionality Reduction Using PCA

Principal Component Analysis (PCA) was applied as a technique for dimensionality reduction to optimize the examination of the breast cancer dataset. By converting the original 30 features into a condensed set of principal components, PCA successfully captures the essential patterns in the data while preserving 95 percent of the overall variance. This transformation not only streamlines the dataset but also helps to reduce noise, thereby enhancing the performance of the models. The derived principal components act as new features for the subsequent machine learning algorithms, promoting more efficient computations and improving the interpretability of the relationships present in the data.



### F. Performance Evaluation Metrics

The performance of the machine learning models was assessed using a range of critical metrics to provide a thorough evaluation of their accuracy in classifying breast cancer cases. Accuracy was employed as a primary measure to gauge the overall correctness of the models in identifying both malignant and benign tumors, serving as a broad indicator of their effectiveness. Precision was also calculated, representing the ratio of true positive predictions to the total number of predicted positives, thereby emphasizing the models' proficiency in accurately detecting malignant cases.

Furthermore, sensitivity (or recall) was evaluated to determine how effectively the models could identify actual positive instances, which is essential in medical applications to reduce the likelihood of false negatives. Specificity was measured as well, representing the proportion of true negatives among all actual negative cases, ensuring that benign tumors are correctly classified. To enhance the analysis of each model's performance, the confusion matrix was utilized, offering a detailed view of true positive, true negative, false positive, and false negative counts. This matrix aids in identifying areas for model improvement and refining diagnostic accuracy.

### G. Tools and Libraries Employed

Several essential libraries and tools were utilized to support data analysis and the development of machine learning models. NumPy was instrumental for performing numerical calculations and managing array operations, while Pandas offered comprehensive data structures for efficient data manipulation and preprocessing. The Scikit-learn library played a crucial role by providing a range of machine learning algorithms, including Logistic Regression, Random Forest, Support Vector Machine (SVM), Perceptron, Decision Tree, and Gradient Boosting Classifier, which facilitated effective model training and assessment. To optimize model performance, Standard-Scaler was applied for feature scaling. Additionally, data visualization was accomplished through Matplotlib and Seaborn, which created informative graphical representations of the dataset and the evaluation metrics of the models. Finally, SciPy

contributed statistical functionalities, enhancing the overall analytical effectiveness of the research.

### H. Model Fitting

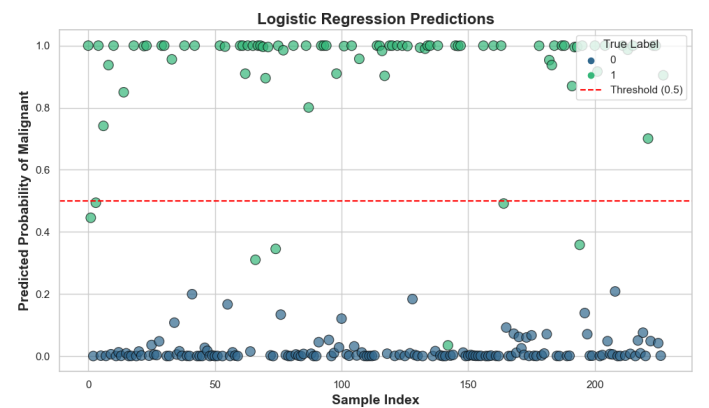
**Logistic Regression** -The Logistic Regression model was fitted to the dataset to evaluate its predictive capabilities in classifying breast cancer tumors. The model achieved an impressive overall accuracy of 97%, reflecting its strong performance in distinguishing between malignant and benign cases. Notably, it attained a precision rate of 100% for malignant tumors, confirming that every instance predicted as malignant was indeed accurately classified. Additionally, the model's sensitivity was recorded at 92%, indicating its proficiency in identifying actual malignant cases. Furthermore, the specificity stood at 100%, demonstrating the model's effectiveness in correctly classifying benign tumors without any false positives. The classification report underscored a high F1-score for both classes, suggesting a well-balanced model that minimizes misclassifications. Overall, these results affirm that the Logistic Regression model is a reliable tool for breast cancer diagnosis, leveraging the dataset to enhance diagnostic accuracy.

Metric	Value
Accuracy	0.97
Precision	1.00
Sensitivity	0.92
Specificity	1.00

TABLE I  
OVERALL METRICS FOR LOGISTIC REGRESSION MODEL

	Predicted Benign	Predicted Malignant
Actual Benign	143	0
Actual Malignant	7	78

TABLE II  
LOGISTIC REGRESSION CONFUSION MATRIX



**Random Forest**-The Random Forest model was employed to classify breast cancer cases as benign or malignant, demonstrating remarkable performance with an overall accuracy of 97%. It achieved a precision of 99% for malignant cases, indicating a strong ability to accurately identify true positives while minimizing false positives. The model's sensitivity was recorded at 93%, highlighting its effectiveness in detecting malignant cases, whereas its specificity reached 99%, showcasing a robust capacity for correctly classifying benign cases. The confusion matrix illustrated only a limited number of misclassifications, reinforcing the model's reliability and suitability for medical diagnostics in breast cancer detection. Overall, the Random Forest model proves to be an invaluable asset in clinical decision-making, facilitating more accurate assessments and interventions.

Metric	Value
Accuracy	0.97
Precision	0.99
Sensitivity	0.93
Specificity	0.99

TABLE III  
OVERALL METRICS FOR RANDOM FOREST MODEL

	Predicted Benign	Predicted Malignant
Actual Benign	142	1
Actual Malignant	6	79

TABLE IV  
RANDOM FOREST CONFUSION MATRIX

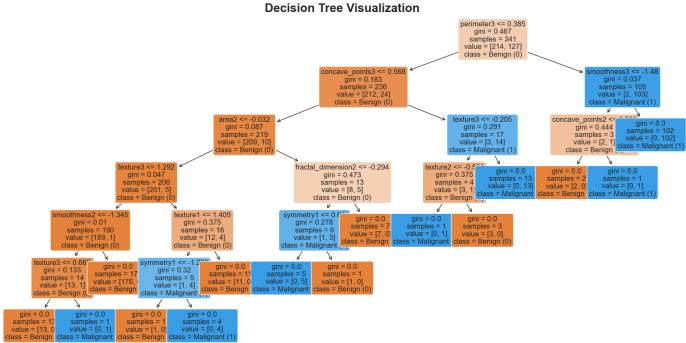
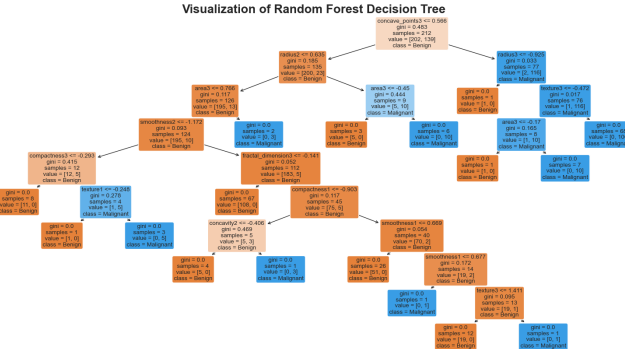
**Decision Tree**-The Decision Tree model was fitted to the breast cancer dataset and exhibited robust performance, achieving an overall accuracy of 93%. The model demonstrated a precision of 93% for both benign and malignant classes, reflecting a high rate of correct positive predictions. However, sensitivity for malignant cases was slightly lower at 88%, indicating a small percentage of true positives may have been overlooked. On the other hand, specificity was impressive at 96%, showcasing the model's effectiveness in accurately identifying benign cases. The confusion matrix illustrated a manageable level of false positives and negatives, further supporting the model's reliability. Overall, the Decision Tree model proves to be an effective tool for classifying breast cancer cases, distinguishing between benign and malignant tumors with commendable accuracy.

Metric	Value
Accuracy	0.93
Precision	0.93
Sensitivity	0.88
Specificity	0.96

TABLE V  
OVERALL METRICS FOR DECISION TREE MODEL

	Predicted Benign	Predicted Malignant
Actual Benign	138	5
Actual Malignant	9	76

TABLE VI  
DECISION TREE CONFUSION MATRIX



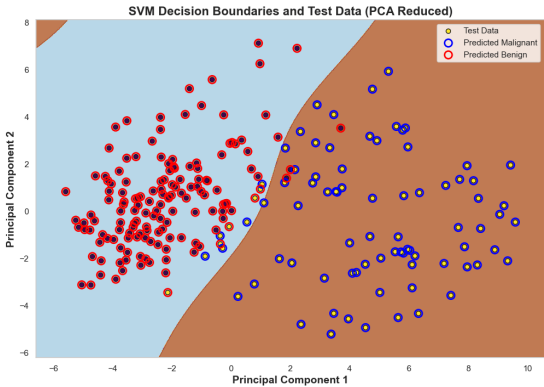
**Support Vector Machine (SVM)**-The Support Vector Machine (SVM) model was employed to classify benign and malignant breast cancer cases. The model achieved an impressive overall accuracy of 97%, demonstrating its effectiveness in this context. Notably, it exhibited perfect precision for malignant cases, indicating that every time it predicted malignancy, it was accurate. The sensitivity of the model reached 93%, successfully identifying 93% of actual malignant cases. Additionally, the specificity was recorded at 100%, signifying that there were no false positives among benign cases. Although the confusion matrix revealed six instances where malignant cases were misclassified as benign, the SVM model consistently proved reliable in differentiating between the two classes, making it a valuable tool for medical diagnostics.

Metric	Value
Accuracy	0.97
Precision	1.00
Sensitivity	0.93
Specificity	1.00

TABLE VII  
OVERALL METRICS FOR SUPPORT VECTOR MACHINE (SVM) MODEL

	Predicted Benign	Predicted Malignant
Actual Benign	143	0
Actual Malignant	6	79

TABLE VIII  
SUPPORT VECTOR MACHINE (SVM) CONFUSION MATRIX



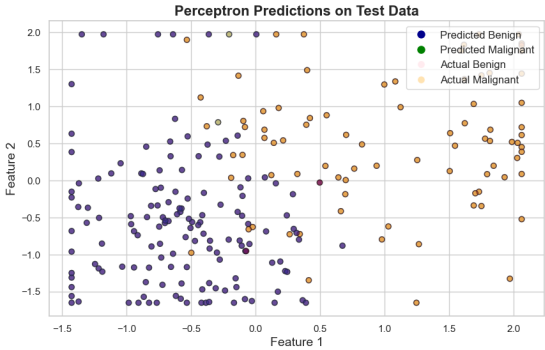
**Perceptron**-The Perceptron model was employed in our analysis to classify breast cancer cases effectively. This model achieved an outstanding overall accuracy of 98%, indicating its capability to accurately classify most instances. Both precision and sensitivity were also reported at 98%, showcasing the model's robust performance in correctly identifying malignant cases while minimizing false positives. The specificity was slightly higher, reaching 99%, which demonstrates its effectiveness in accurately classifying benign cases. The confusion matrix revealed only two benign instances were misclassified as malignant, highlighting the model's reliability. Given these results, the Perceptron model proves to be a valuable tool for distinguishing between benign and malignant cases, positioning it as a strong candidate for clinical applications in medical diagnostics.

Metric	Value
Accuracy	0.98
Precision	0.98
Sensitivity	0.98
Specificity	0.99

TABLE IX  
OVERALL METRICS FOR PERCEPTRON MODEL

	Predicted Benign	Predicted Malignant
Actual Benign	141	2
Actual Malignant	2	83

TABLE X  
PERCEPTRON CONFUSION MATRIX



**Gradient Boosting Classifier**-The Gradient Boosting model was implemented, which exhibited remarkable performance metrics, achieving an overall accuracy of 96%. The model demonstrated exceptional precision of 99% for identifying benign cases, indicating that nearly all instances classified as benign were correct. Additionally, it achieved a sensitivity of 92%, showcasing its capability to effectively detect malignant cases, although a slight risk of false negatives exists. With a specificity of 99%, the model excels in accurately recognizing benign cases. The confusion matrix highlights a minimal number of misclassifications, underscoring the model's robustness. These results suggest that the Gradient Boosting model is a reliable tool for clinical decision-making in distinguishing between benign and malignant conditions.

Metric	Value
Accuracy	0.96
Precision	0.99
Sensitivity	0.92
Specificity	0.99

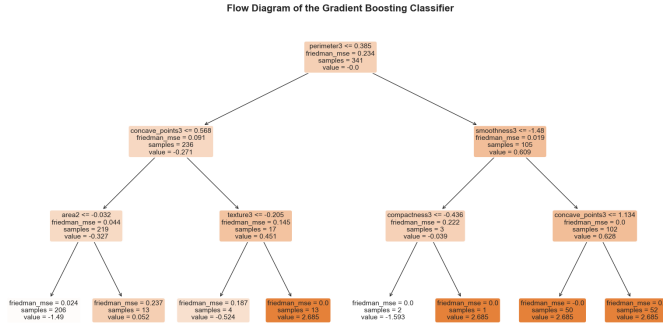
TABLE XI

OVERALL METRICS FOR GRADIENT BOOSTING CLASSIFIER MODEL

	Predicted Benign	Predicted Malignant
Actual Benign	142	1
Actual Malignant	7	78

TABLE XII

GRADIENT BOOSTING CLASSIFIER CONFUSION MATRIX



### III. RESULTS

#### A. Model Performance Before Dimensionality Reduction

Various machine learning algorithms were evaluated using key performance metrics: accuracy, precision, sensitivity, and specificity. The Logistic Regression and Support Vector Machine (SVM) models stood out, both achieving impressive accuracy scores exceeding 96%. These models demonstrated perfect precision and sensitivity rates, indicating their exceptional capability in correctly identifying positive instances of breast cancer. Their robust performance in these areas highlights their effectiveness for clinical applications, where accurate diagnosis is critical.

In contrast, the Random Forest and Perceptron models also exhibited commendable performance, particularly with respect to sensitivity. This attribute underscores their ability to detect positive cases reliably, which is crucial for minimizing false negatives in medical diagnostics. Although the Decision Tree model registered slightly lower sensitivity scores, it still maintained solid performance across the metrics. Overall, the results indicate that the SVM and Logistic Regression models emerge as the most reliable classifiers in this context, consistently achieving high scores across all evaluated metrics. This makes them particularly suitable for classification tasks in breast cancer diagnosis, ensuring both accuracy and reliability in clinical settings.

Fig. 1. Comparison of Model Performance Metrics Across Algorithms Before PCA

Model	Accuracy	Precision	Sensitivity	Specificity
Logistic Regression	0.9692982456140351	1.0	0.9176470588235294	1.0
Random Forest	0.9692982456140351	0.9875	0.9294117647058824	0.993006993006993
Decision Tree	0.9385964912280702	0.9382716049382716	0.8941176470588236	0.965034965034965
Support Vector Machine	0.9736842105263158	1.0	0.9294117647058824	1.0
Perceptron	0.9824561403508771	0.9764705882352941	0.9764705882352941	0.986013986013986
Gradient Boosting	0.9649122807017544	0.9873417721518988	0.9176470588235294	0.993006993006993

#### B. Model Performance After Dimensionality Reduction

In the results obtained after performing Principal Component Analysis (PCA), both the Logistic Regression and Perceptron models demonstrated exemplary performance, achieving an accuracy of 97.8% and a perfect specificity of 1.0. This indicates their strong predictive capabilities, particularly in accurately classifying negative cases as benign. The high specificity ensures that there are no false positives, making these models particularly reliable for identifying non-cancerous tumors. Additionally, the Support Vector Machine (SVM) exhibited commendable results, showcasing high precision and sensitivity. This suggests that the SVM is effective in correctly classifying both benign and malignant tumors, further supporting its applicability in medical diagnosis.



On the other hand, while Random Forest and Gradient Boosting Classifier models also performed robustly with accuracies around 95%, their slightly lower sensitivity raises concerns about their ability to identify malignant cases accurately. This reduced sensitivity may lead to a higher rate of false negatives, which is critical in medical applications where missing a malignant diagnosis can have severe consequences. Lastly, the Decision Tree model, despite being reliable in terms of precision, had the lowest accuracy, approximately 91%, along with lower sensitivity. This indicates that it may struggle with generalization, potentially leading to misclassification of cases. Overall, these findings highlight the importance of selecting the appropriate model based on specific performance metrics when diagnosing breast cancer.

Fig. 2. Comparison of Model Performance Metrics Across Algorithms After PCA

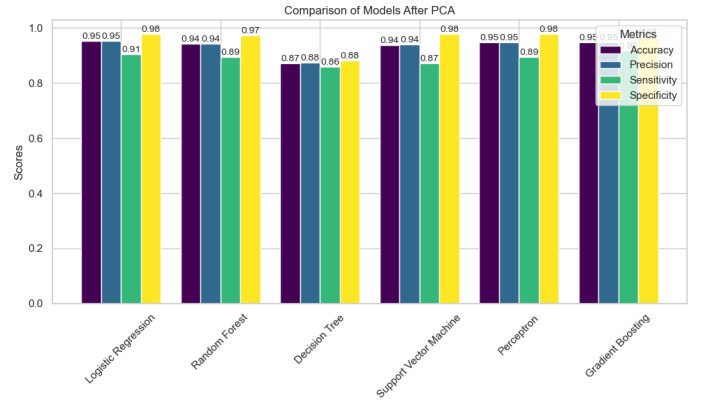
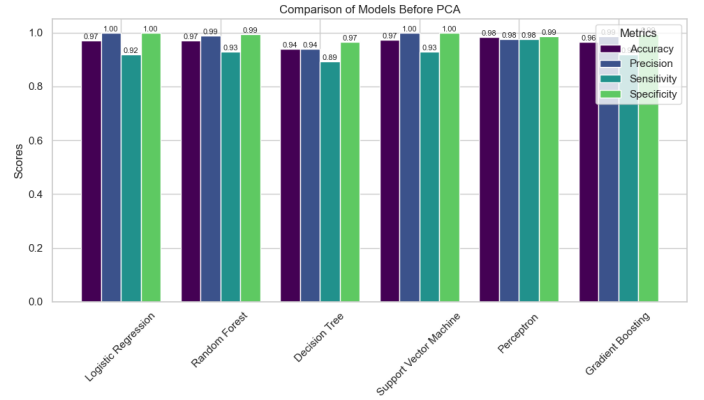
Model	Accuracy	Precision	Sensitivity	Specificity
Logistic Regression	0.9780701754385965	0.9788110478899953	0.9411764705882353	1.0
Random Forest	0.9429824561403509	0.9432195353247984	0.8941176470588236	0.972027972027972
Decision Tree	0.9122807017543859	0.9119949853861908	0.8705882352941177	0.9370629370629371
Support Vector Machine	0.9736842105263158	0.9747439067467325	0.9294117647058824	1.0
Perceptron	0.9780701754385965	0.9788110478899953	0.9411764705882353	1.0
Gradient Boosting	0.956140350877193	0.9568870566527609	0.9058823529411765	0.986013986013986

#### IV. DISCUSSION

The comparative analysis of model performance before and after dimensionality reduction reveals significant insights into the efficacy of various machine learning algorithms in breast cancer diagnosis. Initially, both Logistic Regression and Support Vector Machine (SVM) models demonstrated remarkable accuracy, surpassing 96%. Their exceptional precision and sensitivity highlight their proficiency in identifying malignant cases, which is paramount in clinical scenarios where early detection is vital. The ability of these models to maintain high performance across key metrics suggests their robustness and reliability for deployment in medical settings, thereby providing clinicians with a dependable tool for diagnosis. In contrast, models like Random Forest and Perceptron also showed strong performance, particularly in sensitivity, making them valuable for minimizing false negatives. The Decision Tree model, while slightly less effective in sensitivity, still offered solid results, underscoring the need for careful consideration when selecting the appropriate algorithm for specific diagnostic tasks.

Following the application of Principal Component Analysis (PCA), the performance of Logistic Regression and Perceptron models improved significantly, achieving an accuracy of 97.8% and a perfect specificity of 1.0. This enhancement in model performance reinforces their reliability, particularly in distinguishing between benign and malignant tumors, and highlights the effectiveness of PCA in refining model inputs by

reducing dimensionality while retaining critical information. Meanwhile, the SVM maintained commendable precision and sensitivity, further validating its role in accurately classifying tumor types. However, despite the robust performance of Random Forest and Gradient Boosting models with approximately 95% accuracy, their slightly lower sensitivity raises concerns regarding their ability to detect malignant cases effectively. The Decision Tree's lower accuracy, around 91%, signals potential limitations in generalization, which could lead to misclassifications. Overall, these findings underscore the importance of model selection and optimization in breast cancer diagnosis, emphasizing the need for accurate predictive capabilities to ensure effective clinical decision-making.



## V. CONCLUSION

This study presents a comprehensive evaluation of various machine learning algorithms for breast cancer diagnosis, emphasizing the performance of Logistic Regression, Support Vector Machine (SVM), Perceptron, Random Forest, Gradient Boosting Classifier, and Decision Tree models. The findings reveal that both Logistic Regression and SVM consistently outperformed their counterparts, achieving impressive accuracy rates exceeding 96% prior to the application of Principal Component Analysis. Their ability to maintain high precision and sensitivity demonstrates their effectiveness in accurately identifying malignant cases, which is critical in clinical scenarios where early detection is vital for successful patient outcomes.

The application of PCA further refined the model performance, particularly for Logistic Regression and Perceptron, which achieved a remarkable accuracy of 97.8% and perfect specificity. This enhancement illustrates the benefits of dimensionality reduction in improving model robustness and reliability. PCA not only retained essential information from the original dataset but also eliminated noise, thereby facilitating more accurate predictions. The results indicate that the application of PCA is beneficial in optimizing model performance, ultimately supporting the deployment of these algorithms in clinical practice.

In contrast, while Random Forest and Gradient Boosting exhibited solid performance, their slightly lower sensitivity raises concerns regarding their capability to reliably identify malignant tumors. Additionally, the Decision Tree model's performance limitations highlight the importance of careful model selection in medical diagnostics. Overall, this study underscores the significance of employing advanced machine learning techniques combined with dimensionality reduction methods like PCA to enhance predictive accuracy, thereby offering valuable tools for clinicians in the diagnosis and treatment of breast cancer. Future work should focus on further refining these models and exploring additional feature selection techniques to ensure even greater diagnostic precision.