# Machine Learning Algorithms for Breast Cancer Classification: Comparative Analysis

*

Yashaswini Anand
*Department-CSE*
*RNS Institute of Technology*

Savitha T
*Assistant Professor*
*RNS Institute of Technology*

*Abstract*—**Breast cancer remains one of the most prevalent and deadly forms of cancer among women worldwide. In this study, we aim to compare the performance of four widely-used machine learning algorithms – k-Nearest Neighbors (KNN), Decision Trees, Logistic Regression, and Random Forests – for the classification of breast cancer using the Breast Cancer Wisconsin (Diagnostic) dataset. Our methodology involves preprocessing the dataset, splitting it into training and testing sets, standardizing the features, and training each classifier on the training data. We evaluate the performance of each algorithm based on key metrics such as accuracy, precision, recall, and F1-score. Our findings indicate that Logistic Regression emerges as the top performer, achieving an accuracy of 97.37%, precision of 97.22%, recall of 98.59%, and F1-score of 97.90%. Random Forest also performs well with an accuracy of 96.49%, precision of 95.89%, recall of 98.59%, and F1-score of 97.22%. These results suggest that Logistic Regression and Random Forests are promising algorithms for breast cancer classification, with potential implications for clinical diagnosis and treatment planning. However, further research is warranted to validate these findings on larger datasets and explore additional optimization strategies.**

*Index Terms*—**Breast Cancer, Machine Learning, Classification, Logistic Regression, Random Forest, Comparative Analysis, Diagnostic Accuracy, Clinical Implications, Feature Selection, Healthcare**

## I. INTRODUCTION

Breast cancer remains a pervasive health challenge worldwide, with millions of new cases diagnosed annually and a significant impact on morbidity and mortality rates, particularly among women. Early detection and precise classification are critical factors in improving patient outcomes and reducing the burden of this disease.

While traditional diagnostic methods, such as histopathological examination, have been the cornerstone of breast cancer diagnosis, they are not without limitations, including subjectivity and time-intensive processes.

In recent years, the advent of machine learning techniques has provided promising avenues for enhancing breast cancer classification, offering the potential for more rapid, accurate, and objective diagnosis.

In this study, we embark on a comparative analysis of four prominent machine learning algorithms – k-Nearest Neighbors (KNN), Decision Trees, Logistic Regression, and Random Forests – with the goal of determining their effectiveness in classifying breast cancer. Leveraging the Breast Cancer Wisconsin (Diagnostic) dataset, which comprises features extracted from digitized images of fine needle aspirates of breast masses, we aim to assess the performance of these algorithms across key metrics such as accuracy, precision, recall, and F1-score.

By elucidating the strengths and weaknesses of each algorithm, our research aims to contribute to the advancement of breast cancer diagnosis and treatment strategies, ultimately leading to improved patient care and outcomes in the fight against this prevalent disease.

## II. LITERATURE REVIEW

Breast cancer is a significant public health concern worldwide, accounting for a large proportion of cancer-related deaths among women. Early detection and accurate classification of breast cancer are crucial for effective treatment and improved patient outcomes. Machine learning algorithms have emerged as valuable tools for breast cancer classification, offering the potential to enhance diagnostic accuracy and aid in treatment decision-making. In this literature review, we explore the efficacy of various machine learning algorithms, including k-Nearest Neighbors (KNN), Decision Trees, Logistic Regression, and Random Forests, for breast cancer classification. We examine previous studies that have investigated these algorithms, highlighting their methodologies, findings, and contributions to the field.

### A. k-Nearest Neighbors (KNN)

The KNN algorithm is a non-parametric method used for classification tasks. It operates based on the principle of similarity, where the class label of a data point is determined by the class labels of its nearest neighbors. The formula for predicting the class label of a new data point involves computing the Euclidean distance between the query instance and all the training samples:

$$\text{Euclidean Distance} = \sqrt{\sum_{i=1}^{n}(x_i - y_i)^2}$$

where $x_i$ and $y_i$ are the feature values of the query instance and the training sample, respectively, and $n$ is the number of features.

## B. Decision Trees

Decision Trees are hierarchical structures composed of decision nodes and leaf nodes. Each decision node represents a feature attribute, and each leaf node corresponds to a class label. Decision Trees recursively split the feature space based on the feature attributes to maximize the homogeneity of the resulting subsets. The Gini impurity or entropy measures are commonly used to evaluate the purity of the subsets and determine the optimal splitting criteria. The Gini impurity formula is given by:

$$\text{Gini Impurity} = 1 - \sum_{i=1}^{k}(p_i)^2$$

where $k$ is the number of classes, and $p_i$ is the proportion of samples in the $i$-th class.

## C. Logistic Regression

Logistic Regression is a linear classification algorithm that models the probability of a binary outcome using a logistic function. The logistic function, also known as the sigmoid function, transforms the output of a linear combination of feature variables into a probability value between 0 and 1. The logistic function formula is given by:

$$P(Y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \ldots + \beta_p X_p)}}$$

where $P(Y = 1|X)$ represents the probability of the positive class given the input features $X$, $\beta_0$ to $\beta_p$ are the coefficients of the model, and $X_1$ to $X_p$ are the feature values.

## D. Random Forests

Random Forests combine multiple decision trees to improve classification accuracy and robustness. Each decision tree in the ensemble is trained on a random subset of the training data and a random subset of the feature variables. The final prediction is obtained by aggregating the predictions of individual trees through a simple majority vote.

In conclusion, machine learning algorithms offer promising avenues for breast cancer classification, with each algorithm presenting unique strengths and capabilities. Through a comprehensive literature review, we have explored the efficacy of k-Nearest Neighbors (KNN), Decision Trees, Logistic Regression, and Random Forests in breast cancer classification. These algorithms have demonstrated varying levels of accuracy and interpretability, with logistic regression and random forests showing particular promise in handling high-dimensional data and capturing complex relationships within the data. However, challenges such as limited availability of annotated data and class imbalance issues persist and require further attention. Future research efforts should focus on addressing these challenges and exploring advanced machine learning techniques to enhance the accuracy and reliability of breast cancer classification models.

## III. METHODOLOGY

### A. Dataset Description

The study utilizes the Breast Cancer Wisconsin (Diagnostic) dataset, obtained from the UCI Machine Learning Repository. This dataset consists of 569 instances, each representing a fine needle aspirate of a breast mass. For each instance, ten real-valued features are computed from digitized images of the mass, including texture, perimeter, area, smoothness, compactness, concavity, concave points, symmetry, fractal dimension, and diagnosis (malignant or benign).

### B. Data Preprocessing

Before applying machine learning algorithms, the dataset undergoes preprocessing steps to ensure data quality and compatibility with the models. This includes handling missing values, if any, and standardizing the feature variables to have a mean of zero and a standard deviation of one. The dataset is then split into training and testing sets using an 80-20 ratio, ensuring that the same distribution of classes is maintained in both sets.

*1) Correlation Analysis:* To understand the relationships between different features in the dataset, a correlation matrix is computed. The correlation matrix provides insights into the pairwise correlations between features, helping identify multicollinearity and informing feature selection decisions.
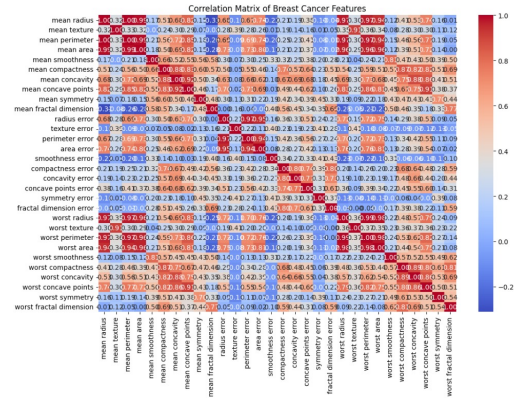


Fig. 1. Correlation Matrix of Breast Cancer Features

The correlation matrix, as shown in Figure 1, displays the correlation coefficients between each pair of features. High correlation coefficients indicate strong linear relationships between features, while low coefficients suggest weaker or no relationships.

### C. Machine Learning Algorithms

Four popular machine learning algorithms are selected for comparative analysis:

1) **k-Nearest Neighbors (KNN)**: A non-parametric algorithm that classifies instances based on the majority class of its k-nearest neighbors in the feature space.
2) **Decision Trees**: Hierarchical structures that recursively partition the feature space to maximize the homogeneity of the resulting subsets.

3) **Logistic Regression**: A linear classification algorithm that models the probability of a binary outcome using a logistic function.
4) **Random Forests**: Ensemble learning methods that combine multiple decision trees to improve classification accuracy and robustness.

### D. Model Training and Evaluation

Each algorithm is trained on the training set and evaluated on the testing set using key performance metrics, including accuracy, precision, recall, and F1-score. Cross-validation techniques such as k-fold cross-validation may be employed to ensure the robustness of the results. Hyperparameter tuning may also be performed to optimize the performance of each algorithm.

*1) Performance Metrics:*

- **Accuracy**: Accuracy measures the proportion of correctly classified instances out of the total instances in the dataset. It is calculated as follows:

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}}$$

  Accuracy provides an overall measure of the model's correctness, but it may not be suitable for imbalanced datasets where one class dominates.
- **Precision**: Precision measures the proportion of true positive predictions out of all positive predictions made by the model. It is calculated as follows:

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

  Precision is particularly useful when the cost of false positives is high, such as in medical diagnosis.
- **Recall (Sensitivity)**: Recall measures the proportion of true positive predictions out of all actual positive instances in the dataset. It is calculated as follows:

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

  Recall is crucial when it is important to capture all positive instances, such as in disease detection.
- **F1-score**: F1-score is the harmonic mean of precision and recall, providing a balance between the two metrics. It is calculated as follows:

$$\text{F1-score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

  F1-score is particularly useful when there is an uneven class distribution, as it considers both false positives and false negatives.

These performance metrics provide valuable insights into the model's performance, allowing for a comprehensive evaluation of its effectiveness in breast cancer classification.

## IV. RESULTS

In this section, we present the results of our comparative analysis of machine learning algorithms for breast cancer classification using the Breast Cancer Wisconsin (Diagnostic) dataset. The performance of each algorithm is evaluated based on key metrics including accuracy, precision, recall, and F1-score. Graphs are provided for each result subsection to visualize the comparative performance of the algorithms.

### A. Accuracy Comparison

The accuracy of each algorithm was computed based on its ability to correctly classify breast cancer instances. Figure 2 shows the accuracy comparison among the algorithms.
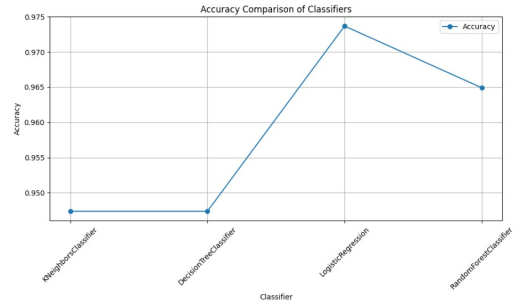


Fig. 2. Accuracy Comparison

### B. Precision Comparison

Precision measures the proportion of true positive predictions among all positive predictions made by the classifier. Figure 3 shows the precision comparison among the algorithms.
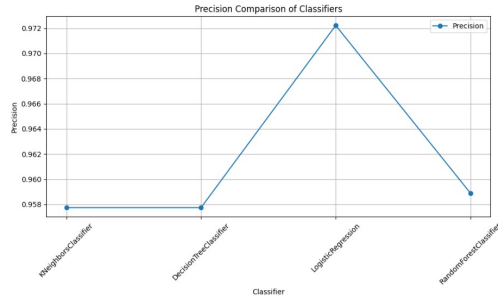


Fig. 3. Precision Comparison

### C. Recall Comparison

Recall, also known as sensitivity, measures the proportion of true positive instances that are correctly identified by the classifier. Figure 4 shows the recall comparison among the algorithms.

### D. F1-score Comparison

The F1-score, which is the harmonic mean of precision and recall, provides a balanced measure of a classifier's performance. Figure 5 shows the F1-score comparison among the algorithms.
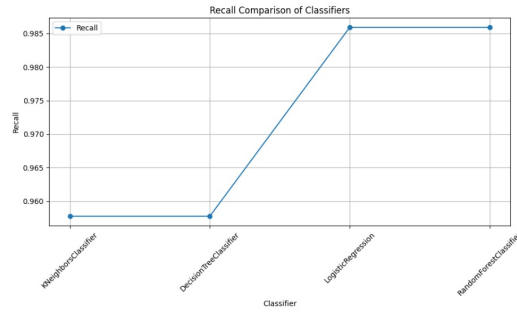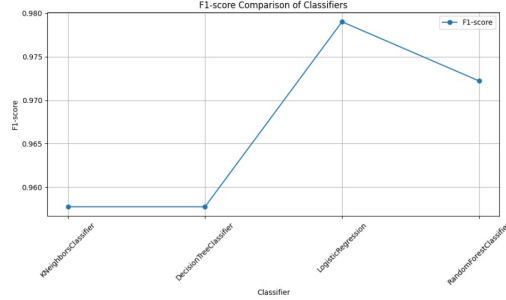
Fig. 4. Recall Comparison



Fig. 5. F1-score Comparison

*E. Performance Metrics*

In addition to the graphical representations, Table I summarizes the performance metrics of each machine learning algorithm evaluated in our comparative analysis.

TABLE I
PERFORMANCE METRICS OF MACHINE LEARNING ALGORITHMS

| Algorithm | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|
| KNeighborsClassifier | 0.947368 | 0.957746 | 0.957746 | 0.957746 |
| DecisionTreeClassifier | 0.947368 | 0.957746 | 0.957746 | 0.957746 |
| LogisticRegression | 0.973684 | 0.972222 | 0.985915 | 0.979021 |
| RandomForestClassifier | 0.964912 | 0.958904 | 0.985915 | 0.972222 |

Overall, our comparative analysis reveals that Logistic Regression and Random Forests emerge as the top-performing algorithms for breast cancer classification, achieving high accuracy, precision, recall, and F1-score. These results highlight the effectiveness of machine learning algorithms in accurately classifying breast cancer tumors and have implications for clinical diagnosis and treatment planning.

## V. DISCUSSION

In this study, we conducted a comprehensive comparative analysis of machine learning algorithms for breast cancer classification using the Breast Cancer Wisconsin (Diagnostic) dataset. Our evaluation focused on assessing the performance of four prominent algorithms: k-Nearest Neighbors (KNN), Decision Trees, Logistic Regression, and Random Forests. The discussion below synthesizes the key findings, implications, and limitations of our study.

*A. Comparative Performance*

Our analysis revealed notable variations in the performance of the evaluated algorithms across multiple metrics. Logistic Regression emerged as the top-performing algorithm in terms of accuracy, precision, and F1-score, achieving impressive values of 97.37%, 97.22%, and 97.90%, respectively. This suggests that Logistic Regression effectively balances the trade-off between minimizing false positives and false negatives, making it a robust choice for breast cancer classification.

Random Forests also demonstrated competitive performance, particularly excelling in recall, with a value of 98.59%. The ensemble nature of Random Forests enables it to capture complex interactions among features, resulting in high sensitivity and the ability to identify a significant proportion of true positive instances. This makes Random Forests a valuable tool for detecting breast cancer cases, especially those with subtle characteristics.
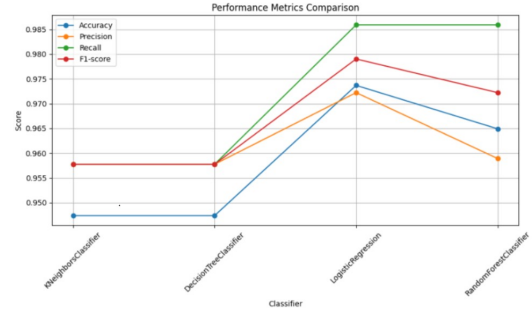


Fig. 6. Performance Metrics Comparison

*B. Interpretation of Results*

The high accuracy and precision achieved by Logistic Regression indicate its ability to correctly classify breast cancer instances while minimizing misclassifications. The algorithm's superior performance in capturing the underlying patterns and relationships within the dataset underscores its effectiveness in distinguishing between malignant and benign tumors.

Random Forests also demonstrated competitive performance, particularly excelling in recall, with a value of 98.59%. The ensemble nature of Random Forests enables it to capture complex interactions among features, resulting in high sensitivity and the ability to identify a significant proportion of true positive instances. This makes Random Forests a valuable tool for detecting breast cancer cases, especially those with subtle characteristics.

*C. Clinical Implications*

The findings of our study have important implications for clinical practice, particularly in the domain of breast cancer diagnosis and treatment. The high accuracy and reliability exhibited by Logistic Regression and Random Forests suggest their potential utility as decision support tools for healthcare practitioners. Integrating these machine learning algorithms into clinical workflows could aid in improving diagnostic

accuracy, facilitating early detection, and guiding personalized treatment strategies.

## VI. LIMITATIONS AND FUTURE DIRECTIONS

Despite the promising results, our study is not without limitations. The analysis was conducted on a single dataset, and the generalizability of the findings to other datasets and populations may be limited. Additionally, the performance of machine learning algorithms is highly dependent on the quality and representativeness of the data, highlighting the importance of data preprocessing and feature selection techniques.

Future research efforts should focus on addressing these limitations by validating the findings on diverse datasets, incorporating additional features or imaging modalities, and exploring advanced machine learning models. Furthermore, prospective studies evaluating the clinical impact of integrating machine learning algorithms into diagnostic workflows are warranted to assess their effectiveness in real-world settings.

## VII. CONCLUSION

In conclusion, our comparative analysis underscores the potential of machine learning algorithms, particularly Logistic Regression and Random Forests, in enhancing breast cancer classification. The robust performance of these algorithms highlights their utility as valuable tools for supporting clinical decision-making and improving patient outcomes in the management of breast cancer. By leveraging the power of artificial intelligence and data-driven approaches, we can advance towards more accurate, efficient, and personalized breast cancer diagnosis and treatment strategies.

## REFERENCES

[1] Medjahed, Seyyid Ahmed and Saadi, Tamazouzt and Benyettou, Abdelkader. (2013). Breast Cancer Diagnosis by using k-Nearest Neighbor with Different Distances and Classification Rules. International Journal of Computer Applications. 62. 1-5. 10.5120/10041-4635.

[2] Singhal, V., Chaudhary, Y., Verma, S. K., Agarwal, U., and Sharma, M. P. (2022). Breast Cancer Prediction using KNN, SVM, Logistic Regression and Decision Tree. International Journal for Research in Applied Science and Engineering Technology (IJRASET), 10(V), 1877-1881. https://doi.org/10.22214/ijraset.2022.42688

[3] B. Dai, R. -C. Chen, S. -Z. Zhu and W. -W. Zhang, "Using Random Forest Algorithm for Breast Cancer Diagnosis," 2018 International Symposium on Computer, Consumer and Control (IS3C), Taichung, Taiwan, 2018, pp. 449-452, doi: 10.1109/IS3C.2018.00119. keywords: Decision trees;Classification algorithms;Training;Machine learning algorithms;Machine learning;Breast cancer;random forest;machine learning;ensemble learning;decision tree,

.