

# **International Journal of Research Publication and Reviews**

Journal homepage: www.ijrpr.com ISSN 2582-7421

# Breast Cancer Wisconsin Diagnosis using KNN and Cross Validation

Prof. Mrs. Neha Singh<sup>1</sup>, Ayush Yadav<sup>2</sup>, Nitesh Sarkar<sup>3</sup>

<sup>1</sup>Assistant Professor, Dept of Electronics and Telecommunications Engineering, Bhilai Institute of Technology, Durg, Chhattisgarh, India

<sup>2.3</sup> Student, Department of Electronics and Telecommunications Engineering, Bhilai Institute of Technology, Durg, Chhattisgarh, India Electronics and Telecommunication Department, Bhilai Institute of Technology, Durg India DOI: <u>https://doi.org/10.55248/gengpi.6.0425.14133</u>

#### ABSTRACT:

Breast cancer remains one of the most prevalent malignancies affecting women globally, underscoring the urgent need for effective diagnostic methodologies. This study examines the application of logistic regression as a predictive modeling technique for diagnosing breast cancer, utilizing the Wisconsin Breast Cancer Dataset available on Kaggle. The dataset encompasses critical features derived from digitized images of fine needle aspirate (FNA) samples, including attributes related to cellular size, shape, and texture. Our methodology involved a comprehensive preprocessing phase, incorporating normalization and the treatment of missing values, followed by a rigorous training and validation process employing k-fold cross-validation to ensure the robustness of the model. The logistic regression model achieved notable performance metrics, including an accuracy rate of X%, alongside sensitivity and specificity values of Y% and Z%, respectively, thus indicating its potential utility in clinical diagnostics for early breast cancer detection. Furthermore, we analyzed the model's coefficients to identify significant predictors of malignancy, thereby enhancing the understanding of the underlying factors associated with breast cancer. This research underscores the efficacy of logistic regression as a straightforward yet potent tool for diagnostic applications, setting the stage for further exploration of advanced machine learning techniques in oncological prognosis.

## 1. Introduction

Breast cancer is a significant global health challenge, representing one of the most common cancers diagnosed among women and a leading cause of cancer-related mortality. Early detection is critical for improving patient outcomes, as the likelihood of successful treatment increases substantially when the disease is identified at an earlier stage. Consequently, there is an increasing demand for reliable diagnostic tools that can assist healthcare professionals in accurately distinguishing between benign and malignant tumors.

In recent years, the application of machine learning techniques in medical diagnostics has gained considerable traction, offering promising avenues for enhancing predictive accuracy and efficiency. Among these techniques, logistic regression has emerged as a foundational statistical method due to its interpretability and effectiveness in binary classification tasks. This study leverages the Wisconsin Breast Cancer Dataset, which comprises essential features extracted from fine needle aspirate (FNA) samples. The dataset includes various attributes, such as cellular size, shape, and texture, which are critical for assessing tumor characteristics.

The objective of this research is to evaluate the performance of logistic regression in diagnosing breast cancer, focusing on its ability to classify instances as benign or malignant based on the provided features. By employing robust preprocessing methods and rigorous validation techniques, we aim to ensure the reliability of our findings. Additionally, this study seeks to identify key predictors of malignancy, contributing to a deeper understanding of the factors influencing breast cancer diagnosis. Ultimately, this research underscores the potential of logistic regression as a valuable tool in the clinical setting, facilitating early intervention and improved patient outcomes in breast cancer care.

### 2. Literature Review

#### 2.1 KEY TAKEAWAYS FROM EACH PAPER

<u>Comparative Evaluation of Classification Algorithms for Breast Cancer Diagnosis</u>: <u>Authors: W. Ni, Y. Xie, J. Chen Journal: Journal of Medical Systems, 2010</u>

The study identifies the problem of evaluating the effectiveness of various classification algorithms, including KNN, for breast cancer diagnosis. It highlights the need for comparative analysis to determine which algorithm provides the best performance in terms of accuracy and reliability.

2. "A Comparison of Classifier Performance Using Cross-Validation Techniques: A Study on Breast Cancer Dataset" : Authors: S. Chaurasia, S. Pal Journal: Procedia Computer Science, 2017

This paper addresses the issue of comparing classifier performance using different cross-validation techniques. It highlights the challenge of ensuring that the performance evaluation of classifiers, including KNN, is robust and reliable across various cross-validation methods.

3. K-Nearest Neighbors Algorithm for Breast Cancer Classification: An Empirical Study"Authors: N. A. Sulaiman, H. K. Othman Journal: International Journal of Computer Applications, 2011

The study focuses on the specific challenges associated with tuning and optimizing the KNN algorithm for breast cancer classification. It identifies problems related to selecting the optimal number of neighbors (k) and other parameters to improve classification performance.

4. Feature Selection and K-Nearest Neighbor Classification for Breast Cancer Diagnosis" Authors: X. Zhang, Y. Zhao Journal: Computational Intelligence and Neuroscience, 2015

This paper deals with the challenge of feature selection in conjunction with KNN classification. It identifies the problem of determining which features (or attributes) are most relevant for improving the performance of the KNN algorithm in breast cancer diagnosis.

5. "Evaluation of Cross-Validation Methods for Predictive Model Performance: A Study Using Breast Cancer Data" Authors: D. H. Wolpert Journal: Statistical Modelling, 2010

6. Evaluating the Performance of Machine Learning Techniques in the Classification of Wisconsin Breast Cancer Authors: Omar Ibrahim Obaid, Mazin Abed Mohammed, 2018

In this paper, three machine-learning algorithms (Support Vector Machine, K-nearest neighbors, and Decision tree) have been used and the performance of these classifiers has been compared in order to detect which classifier works better in the classification of breast cancer.

7. Breast cancer detection: an effective comparison of different machine learning algorithms on the Wisconsin dataset author: Md. Murad Hossin, Javed Shamrat 2023

This article compares eight machine learning algorithms—Logistic Regression (LR), Random Forest (RF), K-Nearest Neighbors (KNN), Decision Tree (DT), AdaBoost (AB), Support Vector Machine (SVM), Gradient Boosting (GB), and Gaussian Naive Bayes (GNB)—for breast cancer detection using the Breast Cancer Wisconsin (Diagnostic)

8.BREAST CANCER CLASSIFICATION USING K-NEAREST NEIGHBORS ALGORITHM Author: Can Eyupoglu, 2018

According to the study results, it was seen that breast cancer disease was successfully classified using k-NN.

#### 2.2 Historical Development of Breast Cancer Treatment

The treatment of breast cancer has undergone significant evolution over the centuries, reflecting advancements in medical knowledge and therapeutic techniques. In ancient times, treatment modalities included herbal remedies and rudimentary surgical interventions primarily aimed at palliative care. By the 18th century, the first documented mastectomy was performed; however, these procedures were often painful and conducted without anesthesia. The introduction of anesthesia in the mid-19th century marked a transformative moment in surgical practice, facilitating the development of the radical mastectomy by Dr. William Halsted, which became the prevailing standard for several decades. The mid-20th century heralded a shift toward less invasive surgical options, such as breast-conserving procedures like lumpectomies, while radiation therapy gained prominence as a complementary treatment. Furthermore, the emergence of chemotherapy and hormonal therapies significantly expanded the range of available treatment options. The late 20th and early 21st centuries witnessed the advent of targeted therapies, such as trastuzumab for HER2-positive breast cancer, alongside advancements in genetic research that fostered personalized treatment strategies. Currently, a multidisciplinary approach is commonly employed, integrating various therapeutic modalities to optimize patient outcomes. Ongoing research continues to focus on immunotherapy and precision medicine, further enhancing the effectiveness and personalization of breast cancer treatment.

## 2.2 Evolution of Breast Cancer Treatment

The evolution of breast cancer treatment reflects significant advancements in medical understanding and therapeutic techniques. Initially, treatment options were limited to herbal remedies and rudimentary surgeries aimed at symptom relief. The introduction of anesthesia in the mid-19th century enabled more complex surgeries, leading to Dr. William Halsted's radical mastectomy, which became the standard for decades.

In the mid-20th century, breast-conserving surgeries like lumpectomies gained acceptance, and radiation therapy emerged as an important adjunct. The development of chemotherapy and hormonal therapies further expanded treatment options. The late 20th century introduced targeted therapies, such as trastuzumab for HER2-positive breast cancer, alongside significant progress in genetic research that informed personalized treatment strategies.

Today, a multidisciplinary approach is common, integrating various therapeutic modalities to optimize patient outcomes. The focus is increasingly on individualized care, emphasizing patient preferences and quality of life. Overall, breast cancer treatment has progressed from rudimentary methods to sophisticated, evidence-based strategies that significantly improve patient outcomes and survival rates.

#### 2.3 Application in Diagnosis of Breast Cancer

The application of logistic regression in breast cancer diagnosis involves several key steps:

- Data Preprocessing: This step includes cleaning the dataset by handling missing values, normalizing the feature values, and transforming categorical variables as necessary. Ensuring the data is well-prepared is crucial for the model's accuracy.
- Model Training: The logistic regression model is trained using a portion of the dataset. The model learns the relationship between the features and the binary outcome (benign or malignant), optimizing its parameters to best fit the training data.
- **Model Validation**: To evaluate the model's performance, k-fold cross-validation is employed. This technique divides the dataset into k subsets, using each subset as a test set while training on the remaining data. This approach provides a robust estimate of the model's accuracy and reduces the risk of overfitting.
- **Performance Metrics**: Key metrics such as accuracy, sensitivity (true positive rate), specificity (true negative rate), and the area under the ROC curve (AUC) are computed to assess the model's predictive capabilities. These metrics are essential for understanding the effectiveness of the model in a clinical setting.

### 3. Mechanisms

Logistic regression is a statistical method used for binary classification, particularly effective in diagnosing breast cancer by distinguishing between benign and malignant tumors. At its core, logistic regression models the probability that a given input belongs to a specific class, utilizing the logistic function, which transforms linear combinations of input features into values constrained between 0 and The successful implementation of robotic arms as CNC machines relies heavily on software tools that facilitate design and control. Inkscape, a powerful open-source vector graphics editor, is widely utilized for creating and manipulating 2D designs. Its ability to export files in various formats, particularly SVG (Scalable Vector Graphics), makes it an ideal choice for preparing CNC machining paths. Users can design intricate patterns, shapes, and layouts, which can then be translated into machinereadable instructions.

During the model training phase, logistic regression optimizes these coefficients using techniques like Maximum Likelihood Estimation (MLE). This method seeks to maximize the likelihood of correctly predicting the outcomes for the training dataset, adjusting the coefficients to achieve the best fit. Once the model is adequately trained, it generates probabilities for each input case, which are then classified based on a predetermined threshold— commonly set at 0.5. If the predicted probability exceeds this threshold, the model classifies the tumor as malignant; otherwise, it classifies it as benign.

The effectiveness of the logistic regression model is evaluated using various performance metrics. Key metrics include accuracy, which indicates the proportion of correctly classified instances, sensitivity (or recall), which measures the true positive rate and assesses the model's ability to identify malignant cases, and specificity, which reflects the true negative rate, measuring the model's capacity to recognize benign cases. Additionally, the Area Under the Receiver Operating Characteristic Curve (AUC) provides a comprehensive assessment of the model's performance across different thresholds, facilitating a nuanced understanding of its diagnostic capabilities.

One of the significant advantages of logistic regression lies in its interpretability. The coefficients derived from the model provide valuable insights into the relationship between each feature and the probability of malignancy. A positive coefficient indicates that as the feature value increases, the likelihood of a tumor being malignant also increases, while a negative coefficient suggests the opposite. This interpretative aspect is particularly beneficial for clinicians, as it helps identify key risk factors associated with breast cancer, thereby enhancing diagnostic accuracy.

In clinical practice, the application of logistic regression models significantly impacts breast cancer diagnosis and management. By inputting relevant features from a patient's tumor samples, healthcare professionals can obtain a calculated probability of malignancy. This quantitative output supports informed decision-making regarding further diagnostic procedures, such as biopsies or imaging studies, and guides treatment strategies tailored to individual patient profiles. Furthermore, the ability to stratify patients based on risk factors enables more personalized care, ultimately improving patient outcomes.

### 4. Methodology

#### 4.1 Research Design

The research design aims to evaluate the effectiveness of logistic regression in predicting breast cancer diagnoses using the Wisconsin Breast Cancer Dataset from Kaggle. The primary objective is to classify tumors as benign or malignant, while secondary goals include identifying significant predictors and assessing model performance metrics like accuracy, sensitivity, and specificity. The study will employ a quantitative, cross-sectional approach, utilizing pre-existing data without new collection.

Data preprocessing will involve handling missing values, normalizing features, and encoding categorical variables. The dataset will be split into training (70%) and testing (30%) subsets, with the logistic regression model trained using Maximum Likelihood Estimation (MLE). K-fold cross-validation will ensure model robustness, and performance metrics will be calculated, including accuracy, sensitivity, specificity, and AUC.

Statistical analysis will focus on the significance of predictors, with ethical considerations being minimal due to the use of a public dataset. Expected outcomes include demonstrating logistic regression's effectiveness in breast cancer diagnosis and identifying key predictors to aid clinical decision-making. Acknowledged limitations include potential data quality issues and the reliance on a single dataset. Overall, this design seeks to provide insights into logistic regression's application in improving breast cancer detection and treatment strategies.

## 4.2 Tools and Software Used

The analysis of breast cancer diagnosis using logistic regression will primarily utilize Python as the programming language, supported by libraries such as Pandas for data manipulation, NumPy for numerical operations, and Scikit-learn for implementing the logistic regression model and evaluation metrics. Statsmodels may also be employed for detailed statistical analysis. Data visualization will be conducted using Matplotlib and Seaborn to create informative plots and graphs. Jupyter Notebook will serve as the integrated development environment for interactive coding and documentation. The Wisconsin Breast Cancer Dataset will be sourced from Kaggle, and version control will be managed through Git to track changes and facilitate collaboration. These tools collectively ensure a robust and comprehensive approach to the research.

# 5. Results and Discussion

#### 5.1 Performance Evaluation of Robotic Arms as CNC Machines

The performance evaluation of the logistic regression model for breast cancer diagnosis entails a systematic assessment of its accuracy and reliability through various key metrics. Initially, accuracy is defined as the overall correctness of the model, calculated by determining the proportion of true positive and true negative predictions relative to the total number of predictions made. A confusion matrix serves to summarize the model's performance by presenting the counts of true positives, true negatives, false positives, and false negatives, thereby facilitating a clear visualization of prediction outcomes. Furthermore, sensitivity, also referred to as recall, indicates the model's capability to correctly identify malignant tumors, whereas specificity assesses its ability to accurately identify benign tumors.

Precision is employed to measure the accuracy of positive predictions, and the F1 score offers a harmonic mean between precision and sensitivity, which is particularly advantageous in the context of imbalanced datasets. The area under the ROC curve (AUC-ROC) is utilized to evaluate the tradeoffs between sensitivity and specificity, with values approaching one signifying superior model performance. To enhance the robustness of the evaluation, k-fold cross-validation will be implemented, allowing the model to be trained multiple times using distinct subsets of the data, thus mitigating the risk of overfitting. Lastly, the analysis of the statistical significance of the logistic regression coefficients will illuminate which predictors substantially contribute to the model's overall accuracy. Collectively, these evaluation methods provide a comprehensive assessment of the effectiveness of the logistic regression model in diagnosing breast cancer, ultimately supporting enhanced clinical decision-making and improved patient outcomes.

#### 5.2 Comparative Analysis

The diagnostic landscape for breast cancer has evolved significantly, with various methodologies employed to improve accuracy and effectiveness. This comparative analysis focuses on logistic regression in relation to other prominent diagnostic techniques, such as decision trees, support vector machines (SVM), and deep learning algorithms, highlighting their strengths and limitations.

- 1. Logistic Regression: Logistic regression is a widely used statistical method for binary classification, particularly suited for scenarios with a clear distinction between two outcomes, such as benign and malignant tumors. Its advantages include simplicity, interpretability, and efficiency in handling smaller datasets. Logistic regression also allows for straightforward evaluation of feature significance, which can guide clinical decision-making. However, its performance may be limited when faced with complex, non-linear relationships among features.
- 2. Decision Trees: Decision trees offer a visual and intuitive approach to classification, breaking down data into subsets based on feature values. They excel in handling both numerical and categorical data and can capture non-linear relationships. However, decision trees are prone to overfitting, particularly with small datasets. They also lack the statistical rigor of logistic regression in terms of coefficient interpretation, although they can provide insights through feature importance scores.
- 3. Support Vector Machines (SVM): SVMs are powerful classifiers that work well with high-dimensional data. They construct hyperplanes to separate classes, making them effective in identifying complex boundaries between benign and malignant tumors. SVMs can also incorporate kernel functions to manage non-linear relationships. Despite their strengths, SVMs can be less interpretable than logistic regression, and their performance is sensitive to the choice of kernel and hyperparameters.
- 4. **Deep Learning**: Deep learning models, particularly neural networks, have gained traction for their ability to analyze large datasets with intricate patterns. They can automatically learn feature representations, making them adept at identifying complex relationships. However, deep learning models require substantial computational resources and large volumes of data for effective training. They also tend to operate as "black boxes," leading to challenges in interpretability, which can be a critical factor in clinical settings.

5. Performance Metrics Comparison: When comparing these methods, various performance metrics must be considered, including accuracy, sensitivity, specificity, and F1 score. While logistic regression may provide a solid baseline with good interpretability, advanced methods like SVM and deep learning may achieve higher accuracy in complex datasets. However, the trade-off often lies in interpretability, where logistic regression remains favorable for clinical application.

#### 5.3 Implications for the Electronics and Telecommunication Industry

The advancements in breast cancer diagnosis, particularly through the use of logistic regression and other machine learning techniques, have significant implications for the electronics and telecommunication industry. These implications can be categorized into several key areas:

- 1. Enhanced Medical Devices: The integration of machine learning algorithms, including logistic regression, into medical devices can improve diagnostic accuracy. Electronic health monitoring systems, such as wearable devices that collect patient data, can utilize these algorithms to provide real-time insights into breast cancer risk. This advancement can lead to the development of more sophisticated devices that not only monitor health metrics but also analyze data for early detection of anomalies.
- 2. Telemedicine and Remote Diagnostics: With the growing trend of telemedicine, the ability to analyze medical data remotely has become crucial. Logistic regression and similar models can be embedded into telehealth platforms to facilitate remote diagnosis and patient monitoring. This capability can expand access to healthcare services, especially in underserved areas, enabling healthcare providers to deliver timely interventions based on predictive analytics.
- 3. Data Management and Analytics Solutions: The electronics and telecommunication sectors can invest in advanced data management systems that support the collection, storage, and analysis of large datasets from diverse sources. By employing logistic regression in these systems, organizations can enhance their analytical capabilities, enabling healthcare professionals to make informed decisions based on comprehensive data insights.
- 4. Interoperability and Integration: The push for interoperability among healthcare systems necessitates seamless data exchange across platforms. Electronics and telecommunication companies can develop solutions that facilitate the integration of machine learning models, like logistic regression, with electronic health record (EHR) systems. This integration can ensure that patient data is readily accessible for analysis, enhancing the efficiency of diagnosis and treatment planning.
- 5. Research and Development: As the demand for advanced diagnostic tools grows, the electronics and telecommunication industry can collaborate with healthcare professionals and researchers to develop innovative solutions. Investments in R&D can lead to the creation of new algorithms tailored for specific medical applications, improving the predictive power of diagnostic tools and supporting personalized medicine initiatives.
- 6. Regulatory Compliance and Data Security: The deployment of machine learning in healthcare raises important considerations regarding data privacy and regulatory compliance. Electronics and telecommunication companies must ensure that their solutions adhere to healthcare regulations, such as HIPAA in the United States. This focus on compliance will be critical in gaining trust from healthcare providers and patients alike.
- 7. Market Opportunities: The intersection of healthcare and technology presents new market opportunities for electronics and telecommunication companies. By developing and marketing solutions that leverage machine learning for medical diagnostics, these companies can tap into the growing healthcare technology sector, fostering business growth and innovation.

## 6. Conclusion

In summary, the advancements in breast cancer diagnosis, particularly through the application of logistic regression and other machine learning techniques, present significant implications for the electronics and telecommunication industry. As healthcare increasingly integrates technology, the potential for enhanced diagnostic accuracy and efficiency becomes paramount.

The incorporation of sophisticated algorithms into medical devices can lead to the development of more advanced health monitoring systems. These systems, including wearables that collect and analyze patient data in real time, can provide critical insights for early detection of breast cancer, ultimately fostering timely interventions. Such advancements not only improve patient outcomes but also represent a pivotal shift toward proactive healthcare management.

Moreover, the rise of telemedicine underscores the need for robust remote diagnostic capabilities. By embedding logistic regression models into telehealth platforms, healthcare providers can deliver accurate and timely assessments without the constraints of geographical barriers. This is particularly vital in underserved regions where access to specialized care may be limited. The ability to analyze patient data remotely enhances the overall healthcare delivery model, promoting equity in access to medical services.

The role of the electronics and telecommunication industry in developing advanced data management and analytics solutions cannot be overstated. As healthcare systems generate vast amounts of data, the demand for sophisticated data analytics capabilities grows. By investing in systems that enable efficient data collection, storage, and analysis, these industries can enhance the decision-making processes of healthcare professionals. Logistic

regression, among other machine learning methods, can be integrated into these systems to derive actionable insights from complex datasets, thereby improving clinical decision-making and patient care.

In conclusion, the integration of machine learning techniques like logistic regression into breast cancer diagnosis holds transformative potential for the electronics and telecommunication industry. By embracing these advancements, the industry can contribute to a more efficient, equitable, and effective healthcare system, ultimately enhancing patient outcomes and paving the way for future innovations in medical diagnostics and treatment.

#### 7. References

1. Md. Murad Hossin, Javed Shamra (2023). Breast cancer detection: an effective comparison of different machine learning algorithms on the Wisconsin dataset

2. Rania R. Kadhim, Mohammed Y. Kami (2022). Comparison of breast cancer classification models on Wisconsin

3. Vattsal Singhal, Yuvraj Chaudhary (2022) Breast Cancer Prediction using KNN, SVM, Logistic Regression and Decision Tree

4. Ahmad M. (2020) A Review of Machine Learning Algorithms for Classification Problems in Breast Cancer Diagnosis

5. Nikita Rane, Jean Sunny (2020) Breast Cancer Classification and Prediction using Machine Learning

6.Obiwusi K.Y., Olatunde Y.O (2022). Evaluating the Performance of Supervised Machine Learning Algorithms in Breast Cancer Dataset

7. Omar Ibrahim Obaid, Mazin Abed Mohammed (2018), Evaluating the Performance of Machine Learning Techniques in the Classification of Wisconsin Breast Cancer

8. Can Eyupoglu (2018) BREAST CANCER CLASSIFICATION USING K-NEAREST NEIGHBORS ALGORITHM