

# Breast Cancer Diagnosis Prediction

Md. Imamul Mursalin sujoy<sup>†</sup>

<sup>†</sup>Department of Computer Science and Engineering, BRAC University

May 08, 2024

## Abstract

Breast cancer diagnosis is a critical area in oncology, demanding accurate predictive models to aid in timely interventions and patient care. This project presents a comprehensive investigation into the development of such models using machine learning techniques. Leveraging a dataset comprising diverse clinical features, including demographic information and tumor characteristics, this study explores the efficacy of K-Nearest Neighbors (KNN) and Support Vector Machine (SVM) algorithms in breast cancer prediction. Through extensive exploratory data analysis (EDA) and feature engineering, the dataset is prepared for model implementation, ensuring robustness and generalizability. The performance of the KNN and SVM models is rigorously evaluated using various metrics, including accuracy, precision, recall, and area under the curve (AUC). The results demonstrate the remarkable predictive capabilities of both models, highlighting their potential utility in clinical settings. This research contributes to advancing the field of medical decision support systems by showcasing the feasibility and effectiveness of machine learning approaches in breast cancer diagnosis prediction. Further refinement and validation of these models hold promise for enhancing healthcare outcomes and patient management strategies in oncology.

**Keywords:** SVM, KNN, recall, f1, AUC

**Corresponding author:** Mr. Annajiat Alim Rasel *E-mail address:* annajiat@bracu.ac.bd

## ■ INTRODUCTION

Breast cancer remains one of the most prevalent and life-threatening diseases affecting women worldwide. Early detection and accurate diagnosis are paramount for successful treatment and improved patient outcomes. In this context, the integration of machine learning techniques into the realm of medical decision-making has emerged as a promising approach to augment traditional diagnostic methods. This project focuses on the development and evaluation of predictive models for breast cancer diagnosis using a dataset encompassing a wide array of clinical attributes. By leveraging advanced algorithms such as K-Nearest Neighbors (KNN) and Support Vector Machine (SVM), this study aims to harness the power of computational analysis to enhance diagnostic accuracy and efficiency. Through a systematic exploration of the dataset, including thorough exploratory data analysis (EDA) and meticulous feature engineering, we seek to uncover meaningful patterns and relationships that can inform the predictive modeling process.

The implementation of KNN and SVM models allows for the creation of robust classifiers capable of distinguishing between malignant and benign breast tumors with high accuracy. By rigorously evaluating the performance of these models using established metrics and validation techniques, we endeavor to demonstrate their potential as valuable tools in clinical decision support. This research not only contributes to the growing body of knowledge in machine learning applications in healthcare but also holds significant implications for improving breast cancer management strategies and patient care protocols.

## ■ DATASET DESCRIPTION

The dataset utilized in this project comprises clinical features extracted from breast cancer patients, including demographic information, tumor characteristics, and diagnostic measures. It contains both malignant and benign instances, enabling the development of a classification model to differentiate between the two classes.

The dataset contains 5 classes for the breast cancer prediction. These are mean radius, mean texture, mean perimeter, mean area and mean smoothness. A comprehensive EDA was conducted to gain insights into the dataset's characteristics and underlying patterns. This involved visualizing distributions of individual features, identifying correlations between variables, and exploring potential associations with the target variable (i.e., cancer diagnosis). Prior to

model development, feature engineering techniques were employed to preprocess and enhance the dataset. This step included handling missing values, scaling numerical features, encoding categorical variables, and selecting relevant attributes based on domain knowledge and statistical analyses.

## ■ MODEL

Two machine learning algorithms, namely KNN and SVM, were implemented for breast cancer prediction. KNN is a non-parametric method that classifies instances based on similarity measures, while SVM constructs a hyperplane to maximize the margin between different classes. Both models were trained on the preprocessed dataset to learn the underlying patterns and relationships between features and target labels. The performance of each model was evaluated using appropriate metrics such as accuracy, precision, recall, and F1-score. Additionally, receiver operating characteristic (ROC) curves and area under the curve (AUC) values were utilized to assess the models' discriminative capabilities and robustness.

## ■ RESULT ANALYSIS

The experimental results demonstrated the efficacy of both KNN and SVM algorithms in accurately predicting breast cancer diagnosis. The models achieved remarkable accuracy rates, indicating their potential utility in clinical practice for assisting healthcare professionals in decision-making processes. The accuracy from SVM I get 88.89% and from KNN I get 90.06%. The successful implementation of predictive models for breast cancer diagnosis signifies the importance of machine learning techniques in augmenting medical decision support systems. Further research could focus on refining model architectures, incorporating additional features, and evaluating performance across diverse patient populations to enhance generalizability and clinical relevance.

## ■ CONCLUSION

In conclusion, this project has demonstrated the efficacy of machine learning algorithms, specifically K-Nearest Neighbors (KNN) and Support Vector Machine (SVM), in breast cancer diagnosis prediction. By leveraging a comprehensive dataset and employing rigorous analytical techniques, we have successfully developed robust predictive models capable of accurately distinguishing between malignant and

benign tumors. The achieved results underscore the potential of computational approaches to augment traditional diagnostic methods and enhance clinical decision-making in oncology.

Furthermore, extending this work to incorporate additional datasets from different sources and populations could provide valuable opportunities for validation and validation of the proposed models. By integrating data from various sources, including genomic, proteomic, and imaging data, a more comprehensive understanding of breast cancer biology and progression can be achieved, leading to more accurate and personalized diagnostic and treatment strategies. Moreover, the application of machine learning techniques to other areas of oncology and healthcare, such as prognosis prediction, treatment response modeling, and drug discovery, holds promise for further advancements in precision medicine and patient care. By harnessing the power of data-driven approaches, we can continue to push the boundaries of medical research and innovation, ultimately leading to improved outcomes for patients battling breast cancer and other complex diseases.