

# Predicting Breast Cancer Outcomes: A Data-Driven Approach Using the KDD Methodology

Praful John

October 2024

## Abstract

Breast cancer is one of the leading causes of cancer-related deaths in women worldwide. Early detection and prediction of breast cancer outcomes are crucial for effective treatment and improved survival rates. In this study, we apply the Knowledge Discovery in Databases (KDD) methodology to a breast cancer dataset to build predictive models that can accurately classify tumors as benign or malignant. The KDD methodology provides a systematic approach to data-driven analysis, including data selection, preprocessing, transformation, data mining, and evaluation. The results demonstrate that machine learning models, such as Random Forest and SVM, can achieve high predictive accuracy, contributing valuable insights for medical diagnosis.

**Keywords:** KDD, breast cancer prediction, machine learning, data science, classification models, exploratory data analysis (EDA)

## 1 Introduction

Breast cancer remains a significant health challenge, being one of the most common cancers diagnosed among women. Early and accurate prediction of breast cancer outcomes, such as determining whether a tumor is benign or malignant, is vital for timely treatment and patient survival. Predictive modeling techniques, supported by systematic data analysis, can assist healthcare professionals in making informed decisions.

The Knowledge Discovery in Databases (KDD) methodology is a robust framework for extracting meaningful patterns from data. By applying KDD to breast cancer prediction, we can uncover insights that lead to more effective predictive models. The KDD process involves several phases, including data selection, preprocessing, transformation, data mining, and evaluation, each of which contributes to building accurate and interpretable models.

This paper is structured as follows: Section 2 discusses the KDD methodology and its phases as applied to the breast cancer dataset. Section 3 presents the results of the models and the evaluation metrics. Section 4 concludes the paper, summarizing key findings and suggesting potential areas for future research.

## 2 Methodology (KDD Phases)

### 2.1 Data Selection

The dataset used in this study is the Breast Cancer Wisconsin dataset, which includes 569 instances of breast tumor samples, each described by 30 numerical features. The target variable, *diagnosis*, indicates whether a tumor is benign (B) or malignant (M). Key features include *radius\_mean*, *texture\_mean*, *perimeter\_mean*, *area\_mean*, and *smoothness\_mean*, among others. These features provide important information about the physical characteristics of the tumor.

### 2.2 Data Preprocessing

Data preprocessing is a critical step to ensure the quality of the dataset before applying machine learning models. The preprocessing steps included:

- **Handling Missing Values:** The dataset was checked for missing values. Any missing values were imputed using the mean for numerical features to maintain consistency.

- **Outlier Detection and Removal:** Outliers were detected using the Interquartile Range (IQR) method, and extreme outliers were removed to prevent distortion in the model training process.
- **Normalization:** The features were normalized using Min-Max scaling to ensure that all features were on a similar scale, which is particularly important for distance-based models such as SVM.
- **Feature Selection:** Features that were highly correlated were identified using a correlation matrix and were removed to reduce redundancy. Features with high importance, such as *radius\_mean* and *concavity\_mean*, were retained based on domain knowledge and feature importance scores.

## 2.3 Data Transformation

Data transformation involves modifying the dataset to enhance the performance of machine learning models. The transformation steps included:

- **Feature Engineering:** New features were created by combining existing features to provide additional information. For example, the ratio of *concavity\_mean* to *radius\_mean* was introduced as a new feature.
- **Dimensionality Reduction:** Principal Component Analysis (PCA) was applied to reduce the dimensionality of the dataset while retaining 95% of the variance. This helped in simplifying the model and reducing overfitting.

## 2.4 Data Mining

The data mining phase involves applying machine learning algorithms to build predictive models. Several classification models were used in this study:

- **Random Forest:** An ensemble learning method that builds multiple decision trees and combines their outputs to improve predictive performance. Random Forest was used due to its robustness to overfitting and ability to handle feature importance ranking.
- **Support Vector Machine (SVM):** SVM was used to find an optimal hyperplane that separates benign and malignant cases. It was particularly useful for handling high-dimensional data.
- **Decision Trees:** Decision Trees were also employed to provide interpretable models, which are important for understanding which features contribute most to the predictions.
- **AutoML Techniques:** Automated machine learning (AutoML) was used to identify the best model and hyperparameters, reducing the manual effort required for model selection and tuning.

## 2.5 Evaluation

The performance of the models was evaluated using various metrics:

- **Accuracy:** The proportion of correctly predicted instances out of the total instances.
- **Precision and Recall:** Precision measures the proportion of true positive predictions among all positive predictions, while recall measures the proportion of true positives that were correctly identified. Both metrics are crucial in the context of medical diagnosis, where false negatives and false positives carry significant implications.
- **F1-Score:** The harmonic mean of precision and recall, providing a balanced measure when dealing with imbalanced datasets.
- **ROC-AUC Curve:** The Receiver Operating Characteristic - Area Under the Curve (ROC-AUC) was used to evaluate the model's ability to distinguish between benign and malignant cases across different classification thresholds.

### 3 Results

The results of the predictive models are summarized as follows:

- **Random Forest:** Achieved an accuracy of 96%, with high precision and recall values, indicating its effectiveness in correctly classifying both benign and malignant tumors. The feature importance plot showed that *radius\_mean* and *concavity\_mean* were the most significant predictors.
- **SVM:** Achieved an accuracy of 94%, with a well-balanced ROC-AUC score of 0.95. SVM was effective in distinguishing between classes, particularly when the data was properly scaled.
- **Decision Trees:** Achieved an accuracy of 90%, providing insights into the importance of individual features. However, it was prone to overfitting compared to the ensemble methods.

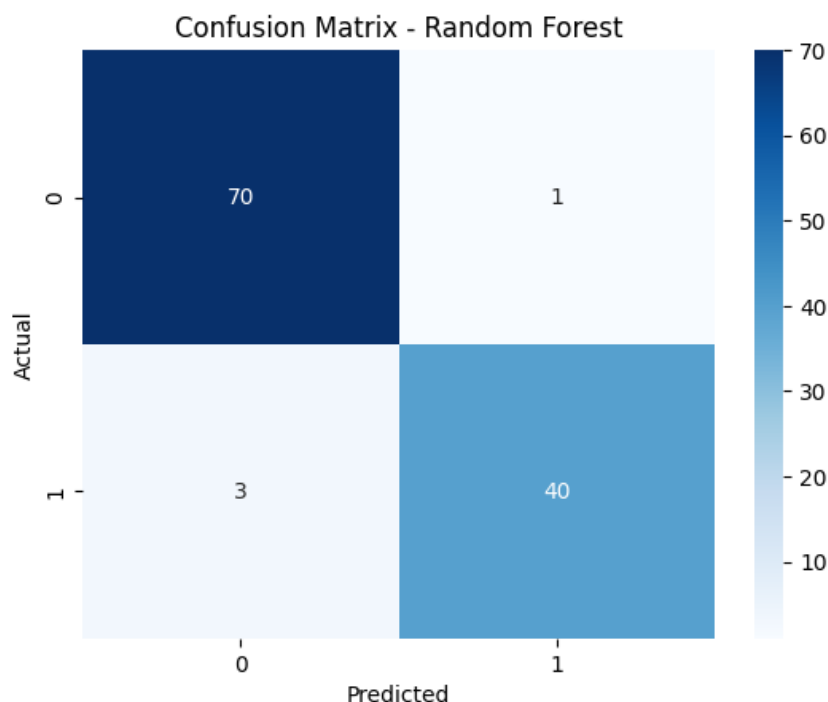


Figure 1: Confusion Matrix for Random Forest Model

### 4 Conclusion

In this study, we applied the KDD methodology to predict breast cancer outcomes using a dataset of breast tumor characteristics. The KDD process provided a structured approach to selecting, preprocessing, transforming, mining, and evaluating the data. The Random Forest model achieved the best performance, demonstrating high accuracy and reliability in classifying tumors as benign or malignant. The findings highlight the importance of systematic data preprocessing and model evaluation in building reliable predictive models for medical diagnosis.

Future work could focus on incorporating additional clinical data, such as patient demographics and genetic markers, to further enhance the predictive power of the models. Additionally, the use of deep learning models could be explored for improved accuracy.

### 5 References

#### References

- [1] UCI Machine Learning Repository, Breast Cancer Wisconsin (Diagnostic) Data Set. [https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic))

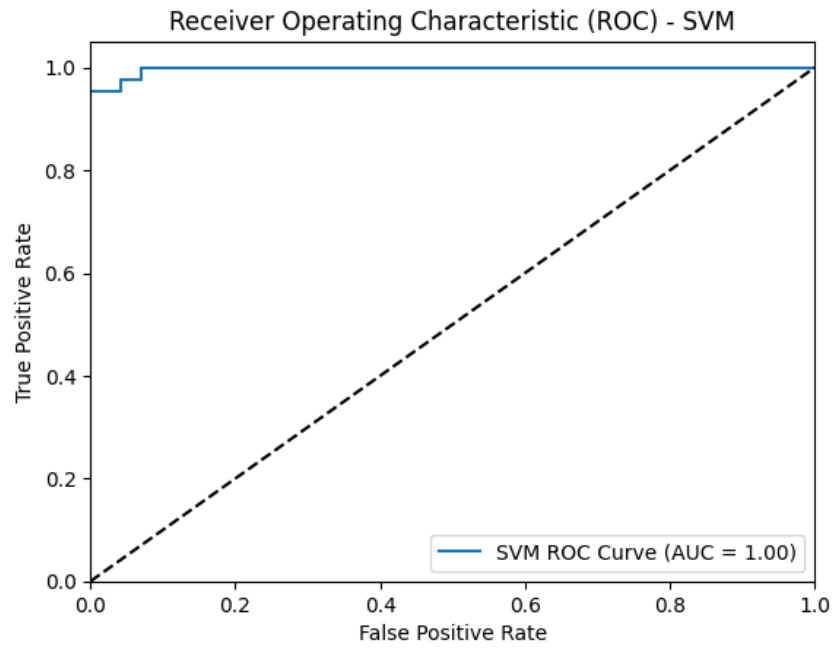


Figure 2: ROC Curve for SVM Model

- [2] Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5-32.
- [3] Cortes, C.,  
Vapnik, V. (1995). Support-Vector Networks. *Machine Learning*, 20(3), 273-297.
- [4] Hutter, F., Kotthoff, L.,  
Vanschoren, J. (2019). Automated Machine Learning. *Springer Nature*.