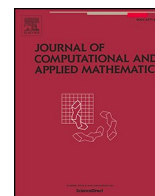


Contents lists available at [ScienceDirect](https://www.sciencedirect.com)

Journal of Computational and Applied Mathematics

journal homepage: www.elsevier.com/locate/cam

Early breast cancer prediction using optimized machine learning and tumor-immune modeling

Oluwaseun Olumide Okundalaye^{a,*}, Necati Özdemir^b, Richard Olu Awonusika^a

^a Adekunle Ajasin University, Akungba-Akoko, Department of Mathematical Sciences, Faculty of Science, Ondo State, Nigeria

^b Balıkesir University, Department of Mathematics, Faculty of Science and Art, Balıkesir, Türkiye

ARTICLE INFO

Keywords:

Breast cancer
Machine learning
Feature selection
Cross-validation
Hyperparameter tuning
Performance measures
Early detection

ABSTRACT

This study aims to enhance early breast cancer prediction accuracy by utilizing machine learning classifiers and feature selection techniques. The Wisconsin Diagnostic Breast Cancer (WDBC) dataset was used to train and evaluate three popular machine learning classifiers: Support Vector Machine (SVM), Random Forest (RF), and k-Nearest Neighbors (k-NN). Feature selection methods were applied to optimize model performance, including Recursive Feature Elimination (RFE) and Mutual Information. Cross-validation and hyperparameter tuning were conducted to ensure the robustness and reliability of the models. The results showed that the SVM classifier achieved the highest performance with an accuracy of 98 %, compared to 95.8 % for RF and 96.2 % for k-NN. The SVM model demonstrated a precision of 0.98 and a recall of 0.95 for malignant cases. Feature selection revealed that mean radius, texture, and area were the most influential features, and SHapley Additive exPlanations (SHAP) analysis confirmed their clinical relevance in breast cancer diagnosis. A tumor-immune dynamic model also indicated that treatment efficacy ($\gamma = 0.0500/\text{day}$) was a critical parameter for tumor control. Statistical significance tests ($p < 0.05$) confirmed that the SVM classifier outperformed the other models. This study highlights the potential of combining machine learning with clinical insights to develop an effective framework for breast cancer prediction, offering high diagnostic accuracy and biological interpretability.

1. Introduction

The global healthcare sector is experiencing rising expenditure due to many causes, such as the ageing population, significant medical technology investments, and the incidence of chronic illnesses. This burden of chronic diseases extends beyond healthcare costs, impacting societal productivity and human potential. Chronic diseases are estimated to cost approximately 3.4 million potential productive life years globally. This staggering figure underscores these conditions' significant economic and social consequences, highlighting the importance of prioritizing prevention, early intervention, and effective management strategies to mitigate their impact on individuals and communities alike [1]. The principal long-term illnesses, such as cancer, stroke, diabetes, hypertension, and chronic respiratory conditions, collectively impose a substantial burden on global health and economies. These conditions not only contribute significantly to healthcare costs but also result in millions of premature deaths and productive years lost. Addressing these challenges requires a concerted effort to prioritize prevention, early detection, and effective management strategies within healthcare systems worldwide [2]. Cancer is a disease in which abnormal cells grow uncontrollably, forming tumors. These tumors can be either

* Corresponding author.

E-mail address: okundalaye.oluwaseun@aau.edu.ng (O.O. Okundalaye).

<https://doi.org/10.1016/j.cam.2025.116875>

Received 30 April 2025; Received in revised form 12 June 2025;

Available online 21 June 2025

0377-0427/© 2025 Elsevier B.V. All rights are reserved, including those for text and data mining, AI training, and similar technologies.

benign (localized and non-invasive) or malignant, which are more dangerous as they can spread to other parts of the body [3]. Malignant tumors pose a significant health risk as they can invade nearby tissues and organs, making treatment challenging. Understanding the difference between benign and malignant tumors is vital for effective diagnosis and treatment strategies in combating cancer [4]. The current global breast cancer market prediction for 2032 is shown in Fig. 1 below.

One of the leading causes of death for women worldwide is breast cancer disease (BCD), posing a considerable health challenge. In the United States alone, nearly four million women currently live with a breast cancer diagnosis. Because of its ubiquity, there is an urgent need to keep up efforts in early detection, cutting-edge treatment options, and extensive support networks to lessen the effects of the disease and enhance the lives of those affected [5]. Early detection through screening plays a crucial role in reducing the mortality rate associated with breast cancer. Cancer can be detected through screening in its earliest stages, when therapy is most effective, improving patient outcomes and survival rates. Although screening trials are designed to improve early detection, the process often involves a significant workload for healthcare professionals due to the need to evaluate numerous scans thoroughly [6]. Recent advances in Artificial Intelligence (AI) and Machine Learning (ML) have shown great promise in addressing the limitations of traditional screening methods. ML algorithms, particularly supervised learning models, can detect complex, non-linear patterns in high-dimensional biomedical data, enabling more accurate and faster predictions of malignancy. These methods enhance decision-making, reduce diagnostic errors, and allow for scalable deployment in clinical settings. Notably, the synergy between AI-driven diagnostics and mathematical models opens new avenues for personalized treatment planning by simulating tumor growth and therapy responses. Aligned with these efforts, the World Health Organization (WHO), through its Global Breast Cancer Initiative (2021), emphasizes early diagnosis and digital innovation as key strategies to reduce global breast cancer mortality by 2.5 % annually. The initiative advocates for the integration of AI into diagnostic protocols to assist clinicians, especially in under-resourced regions. This labour-intensive process highlights how urgently automated and precise techniques for early breast cancer detection and prognostic prediction are needed. Technological developments promise to address this demand, especially in machine learning and artificial intelligence (AI) [7]. The BCD sign and symptoms are shown in Fig. 2.

2. Literature review

Automated systems can analyze medical imaging data, such as mammograms, with speed and precision, assisting healthcare providers in identifying suspicious lesions or abnormalities more efficiently. Additionally, AI-driven tools can aid in predicting the prognosis of breast cancer, helping clinicians tailor treatment plans to individual patients’ needs [8]. By harnessing the power of automation and AI, healthcare professionals can streamline the screening process, improve diagnostic accuracy, and ultimately enhance patient outcomes in the fight against breast cancer. Machine learning, a fundamental artificial intelligence component, employs statistical methods to analyze data and derive insights. Within machine learning, three main techniques are applied: “supervised learning (SL), unsupervised learning (UL), and semi-supervised learning (SSL)” [9]. SL uses labelled data to train algorithms, while UL looks for patterns in unlabeled data. Labeled and unlabeled data are combined for training in (SSL). Each technique serves distinct purposes, enabling the development of models that make accurate predictions and uncover valuable insights across diverse applications [10]. Different works have featured using ML methods to predict BCD. The literature in this field includes. Using

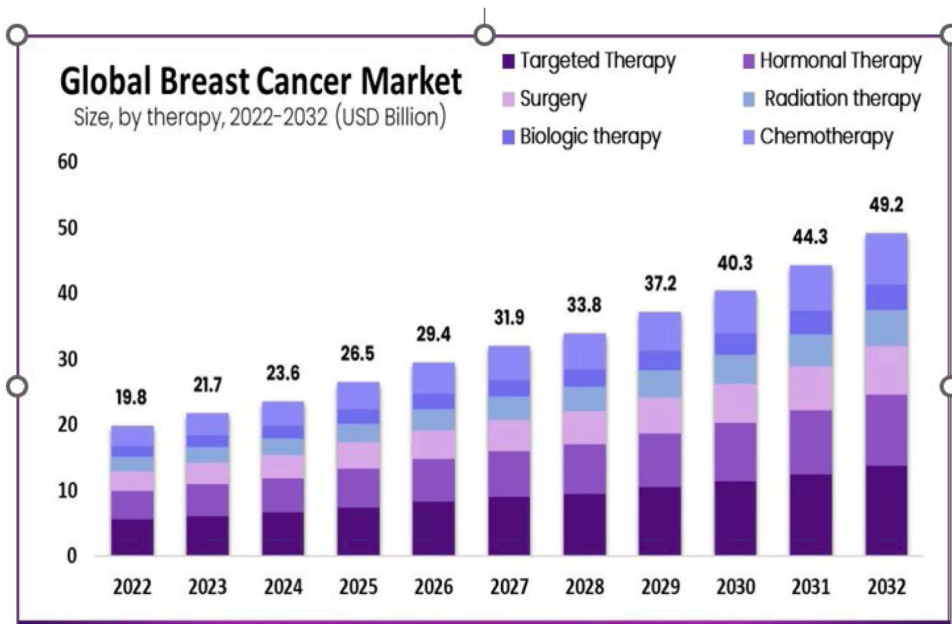


Fig. 1. The current global breast cancer prediction.

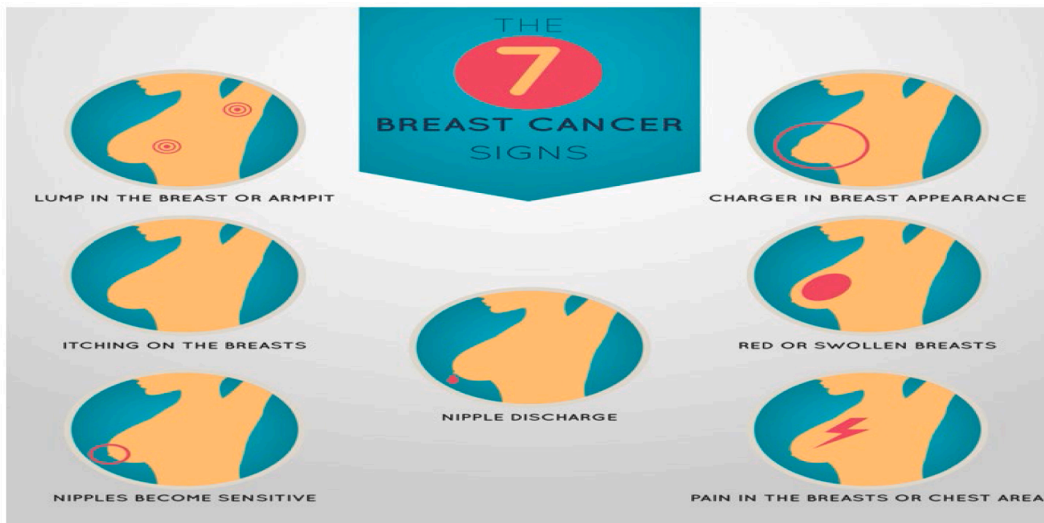


Fig. 2. Breast cancer symptoms.

semi-supervised MLA, including decision trees, RF, and K-NN, the study in [11] Built a prediction system of BCD occurrence for the WDBC dataset in the UCI machine learning repository. Random forest was determined to have 96 % of the system's computed accuracy. Using the voting strategy, Yue et al. (2018) implemented "naive Bayes, J48, and SVM in the ensemble breast cancer predictive analysis technique". The authors report that the accuracy rate of the ensemble approach was 97.13 %. Thirumalaikolundusubramanian [12] Conducted comparison research between the "Bayes Belief Network (BBN), Tree Augmented Naive Bayes (TAN) and Boosted Augmented Naive Bayes (BAN)" and confirmed the effectiveness of naive Bayes strategies in predicting breast cancer. According to their studies utilizing gradient boosting, the best accuracy rate for TAN attained was 94.11 %. Therefore, according to the authors, TAN is the most effective classifier for the WBCD when using naive Bayes approaches. Ghani, et al. [13] used the Coimbra BCD dataset from the UCI repository. Prediction began with pre-processing, after which vital attributes were identified using recursive feature elimination (RFE). A stacking classifier (SC) was used on WDBC, obtaining a 97.2 % accuracy rate for prediction BCD. RFE selects five features out of nine using RF. Based on the experiment results, the ANN performed the best, classifying the data with an accuracy of 80 % into two categories: healthy and patient. Basunia, et al. [14] proposed the stacking classifier, an ensemble method incorporating random forest, KNN, and SVM classification techniques. The anticipated outcomes of combining these methods were fed into a logistic regression meta-classifier as input. When the stacking classifier was used on WDBC, the accuracy rate for predicting breast cancer was 97.2 %. In the study by et al. [15], supervised machine learning techniques, such as logistic regression, KNN, and SVMs, were combined with principal components analysis (PCA) to reduce dimensionality and detect breast cancer patients. Data from the UCI repository were used for the experiment [16]. The suggested method identified breast cancer patients with a maximum accuracy of 92.7 % using SVMs. Al Tawil, et al. [17] proposed breast cancer predictive modelling with selection techniques and machine learning algorithms. The findings show that LightGBM obtains the maximum accuracy at 95 % when feature selection is not used. With 15 features selected for minimum redundancy and maximum relevance (mRMR) feature selection, LightGBM achieves 98 % accuracy, outperforming other classifiers. LightGBM also performs well, with a 95 % accuracy rate, for selecting characteristics for the Pearson correlation coefficient (15 features). Machine learning (ML) offers a promising approach by identifying complex patterns in medical data, reducing diagnostic errors, and enhancing predictive accuracy. While prior studies have applied ML classifiers to breast cancer prediction, few have rigorously evaluated feature selection's impact on classification performance. This study addresses this gap by implementing advanced feature selection techniques and explainability methods to optimize breast cancer prediction models. Although previous studies have achieved high accuracy using ML classifiers, they often lack an in-depth exploration of feature selection techniques and their impact on classification performance. This study bridges this gap by employing an extensive feature selection analysis alongside hyperparameter tuning to optimize model accuracy. While numerous studies have demonstrated the utility of machine learning in breast cancer classification, few have combined feature selection techniques with dynamic tumor modeling to enhance biological interpretability. This study contributes to the field in several novel ways: (1) it integrates Recursive Feature Elimination and Mutual Information for optimal feature selection; (2) it employs SHAP values to interpret model decisions clinically; (3) it formulates and calibrates a tumor-immune compartmental model using machine learning-derived parameters; and (4) it bridges predictive analytics with mechanistic modeling to offer a more comprehensive understanding of tumor behavior and treatment efficacy.

The paper is structured as follows: Section 2, literature review, Section 3 covers methodology, including data preprocessing, feature selection, and ML model implementation, Section 4 presents ML algorithms, Section 5, Numerical results and analysis, and Section 6 provides conclusions and

3. Methodology

This section delves into the five main phases of the methodology shown in Fig. 3. These phases are data preprocessing, feature selection, trained model, model evaluation, hyperparameter tuning.

3.1. Data preprocessing

The WDBC dataset from the UCI Machine Learning Repository is the publicly accessible dataset that was used for this study [16]. There are 569 instances in the collection, and each instance has 30 attributes. These characteristics determine whether the tumor is classified as benign (B) or malignant (M). The collection of the 31 qualities is displayed in Table 1, along with a description. The present study uses labelling and normalization as data preprocessing strategies to enhance machine learning performance. The dataset is carefully preprocessed to maximize its suitability for machine learning applications. By following these procedures, we can be confident that the data is organized and polished in a way that improves the efficiency of the models. Preprocessing comprises, first, the features and the target variable are the two primary divisions of the dataset at first. The target variable is the result or dependent variable that the model aims to predict. In contrast, features are the different characteristics or independent variables in the dataset used to generate predictions or classifications (Splitting the Dataset). Second, the feature values are scaled using the StandardScaler approach after separating the features and the target variable. This step standardizes feature scales, reducing model bias due to differing magnitudes. By transforming the characteristics to have a mean of 0 and a standard deviation (SD) of 1, StandardScaler helps the models train more consistently and successfully (Feature Scaling with StandardScaler).

3.2. Features selection

Selecting pertinent information for ML model training reduces the input features. The training model is optimized by feature selection in a few different ways as shown in Fig. 4: First, minimizing the dimensionality of the data by eliminating features that are not relevant or only partially relevant, avoiding overfitting and learning from noise. Second, raising the accuracy of the predictions. Thirdly, cutting down on training time has caused some training models to develop exponentially [18].

Selecting the most relevant features improves model interpretability and efficiency. We employed: recursive feature elimination (RFE) (Iteratively removed the least significant features), mutual information (Identified features with the highest correlation to breast cancer classification), principal component analysis (PCA) (for comparative analysis).

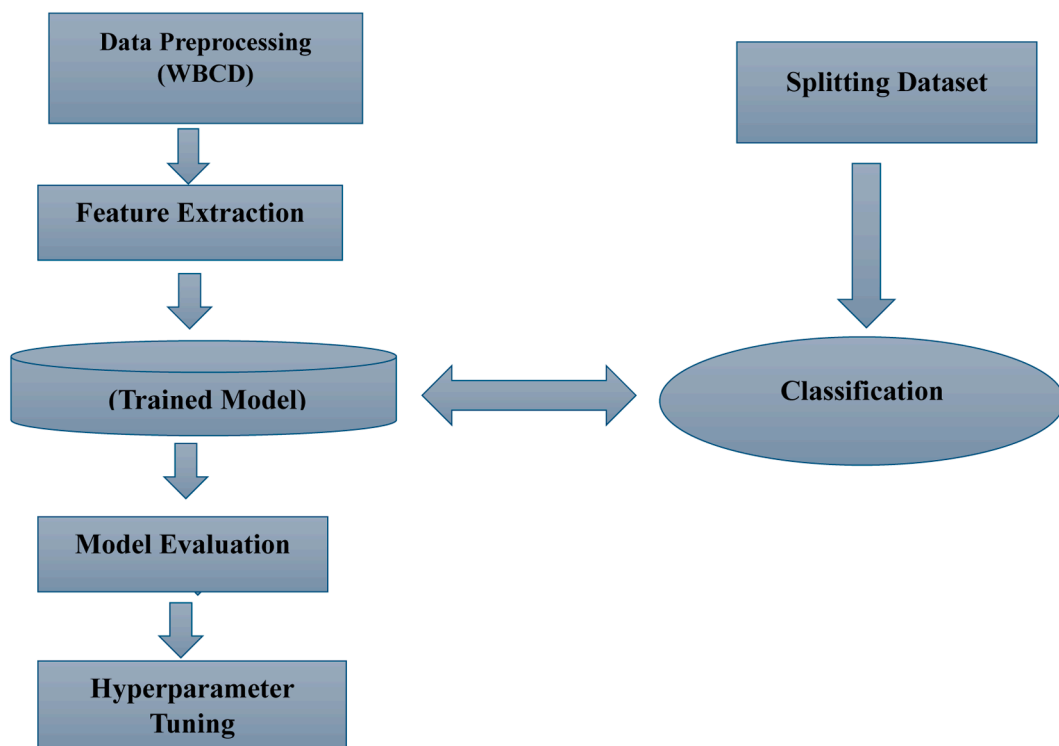


Fig. 3. The flowcharts of the model.

Table 1
The dataset attribute and description.

S/N	Attribute Name	Description
1	Radius mean	Mean of distances from center to points on the perimeter
2	Texture mean	Standard deviation of gray-scale values
3	Perimeter mean	Mean size of the core tumor
4	Area mean	Mean area inside the tumor boundary
5	Smoothness mean	Mean of local variation in radius lengths
6	Compactness mean	Mean of $\text{perimeter}^2/\text{area} - 1.0$
7	Concavity mean	Mean severity of concave portions of the contour
8	Concave points mean	Mean number of concave portions of the contour
9	Symmetry mean	Mean similarity between matching tumor areas
10	Fractal dimension mean	Mean "coastline approximation" - complexity of tumor boundary
11	Radius se	Standard error of radius measurements
12	Texture se	Standard error of texture values
13	Perimeter se	Standard error of perimeter measurements
14	Area se	Standard error of area measurements
15	Smoothness se	Standard error of smoothness values
16	Compactness se	Standard error of compactness values
17	Concavity se	Standard error of concavity measurements
18	Concave points se	Standard error of concave point counts
19	Symmetry se	Standard error of symmetry measurements
20	Fractal dimension se	Standard error of fractal dimension
21	Radius worst	Largest (worst) radius measurement
22	Texture worst	Largest deviation in gray-scale values
23	Perimeter worst	Largest perimeter measurement
24	Area worst	Largest area measurement
25	Smoothness worst	Greatest variation in radius lengths
26	Compactness worst	Highest $\text{perimeter}^2/\text{area}$ ratio
27	Concavity worst	Most severe concave portions
28	Concave points worst	Highest number of concave portions
29	Symmetry worst	Least symmetrical tumor areas
30	Fractal dimension worst	Most complex tumor boundary pattern

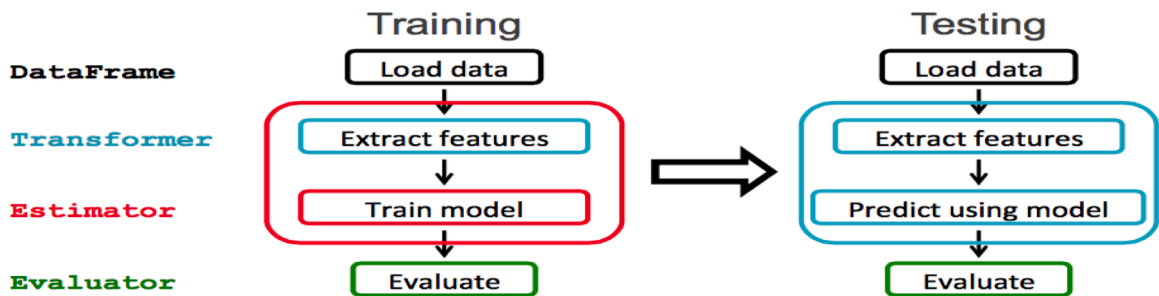


Fig. 4. Data extraction features.

3.3. Trained model

This section focuses on the development of machine learning models used to classify breast cancer cases in the WDBC dataset. After preprocessing and feature selection, we trained three classifiers: Support Vector Machine (SVM), Random Forest (RF), and k-Nearest Neighbors (k-NN). Each model was trained using the optimized features selected through Recursive Feature Elimination (RFE) and Mutual Information analysis to enhance performance and minimize computational complexity.

3.3.1. Splitting data

The dataset was divided into two subsets to evaluate the performance of the trained models effectively: 80 % of the data was used for training and 20 % for testing. This split ensures that the models are exposed to sufficient data for learning while retaining a portion of the data for unbiased performance evaluation. Stratified sampling was used to maintain the original distribution of benign and malignant cases in both sets, preserving the balance and integrity of the classification task. Additionally, before model training, feature scaling was applied using the StandardScaler from Scikit-learn, which transformed the feature values to have a mean of 0 and a standard deviation of 1. This step was essential to ensure that all features contributed equally to the classification process, especially for algorithms like SVM and k-NN, which are sensitive to feature magnitude.

3.3.2. Classification

Three supervised machine learning algorithms were implemented using the Scikit-learn Python library:

Support Vector Machine (SVM): Utilized the radial basis function (RBF) kernel due to its effectiveness in handling non-linear classification tasks. Hyperparameters, including the regularization parameter C, the kernel coefficient gamma, and the kernel type were fine-tuned using GridSearchCV to achieve optimal performance.

Random Forest (RF): Built using an ensemble of decision trees with each tree trained on a bootstrap sample of the training data. The number of trees and the maximum depth were among the hyperparameters optimized. **k-Nearest Neighbors (k-NN):** Implemented using Euclidean distance as the metric. The optimal value of k was determined through cross-validation.

3.4. Model evaluation

A wide range of indicators is used to assess the performance of the top model, SVM, on the testing set. Among these metrics are accuracy, which calculates the percentage of cases correctly classified relative to all instances; the confusion matrix, which presents an in-depth analysis of both accurate and inaccurate classifications; and the key classification metrics, including precision, recall, and F1-score for each class, are summarized in the classification report. A “true positive rate and a false positive rate” are traded off throughout a range of threshold values, as depicted by the Receiver Operating Characteristic (ROC) curve. When taken as a whole, these metrics provide a comprehensive evaluation of how effectively the SVM model generalizes new data and how well it classifies cases.

3.4.1. Cross-validation

Cross-validation (CV) is one of the most important data analysis and validation methods in ML. Among the many benefits of this strategy is its ability to minimize model bias and guarantee that the model is not overfitting the data. Furthermore, this approach permits the model to be trained and tested on the same dataset, which is advantageous when the amount of data is restricted. Cross-validation is a useful and trustworthy method for evaluating how well machine learning algorithms are performing, and it should be heavily taken into account when performing various data analysis activities [19]. Furthermore, the usefulness of CV in Python lies in its ability to provide a highly accurate evaluation of a model’s performance on unseen data related to a simple train-test split. Using cross-validation, we can detect overfitting or underfitting, select the best hyperparameters for your model, and evaluate its generalization capabilities more robustly. Python libraries like scikit-learn provide convenient functions to implement various cross-validation techniques with just a few lines of code. We separate the data into 20 % testing sets and 80 % training sets.

3.4.2. Visualization

The use of visualization tools is essential for gaining a deeper comprehension of the functioning of the model. A tabular representation of actual versus anticipated class labels, the confusion matrix provides information about several errors the model makes, including false positives and false negatives. This image helps determine which classes are more difficult for the model to classify correctly. Comparably, the ROC curve graphically represents the trade-off between the true positive (TP) rate and the false positive (FP) rate across various threshold values. The model’s capacity to distinguish between classes can be evaluated using this curve, whereby a larger area under the curve denotes superior performance. By using these visualization tools, analysts and stakeholders can obtain important insights into the model’s advantages and disadvantages, which can help with possible model changes and well-informed decision-making.

3.5. Hyperparameter tuning

Hyperparameters play a vital role in determining the performance of ML models. In the case of Support Vector Machines (SVM), these hyperparameters, such as the regularization parameter C, the kernel coefficient gamma, and the choice of kernel function, significantly impact the model’s ability to generalize to unseen data [20]. A technique known as grid search, implemented through the GridSearchCV class in the sci-kit-learn library, is employed to optimise the SVM’s performance. Grid search contains defining a grid of hyperparameter values to explore [21]. For every combination of hyperparameters, the model is trained and assessed using CV to estimate its performance. The goal is to identify the hyperparameters that yield the best performance metric, such as accuracy, precision, or recall. By systematically searching through the hyperparameter space, grid search helps find the optimal configuration that maximizes the model’s performance on the given dataset. This process is essential for fine-tuning the model and ensuring it can effectively generalize to new, unseen data, making it a crucial step in the ML pipeline.

```

GridSearchCV(estimator=SVC(),
param_grid={C: [0.1, 1, 10, 100],
coef0: [0.0, 0.1, 1.0],
degree: [2, 3, 4],
gamma: [0.1, 0.01, 0.001, 0.0001],
kernel: [rbf, linear, poly],
max_iter: [-1, 1000, 2000],
probability: [True, False],
shrinking: [True, False],
tol: [0.001, 0.0001, scoring=accuracy, verbose=1]}

```

The choice of SVM, Random Forest, and K-NN was driven by their proven effectiveness in medical classification tasks, especially

with smaller datasets like WDBC. SVM is well-suited for high-dimensional feature spaces, Random Forest provides robustness and interpretability, and K-NN serves as a simple yet effective baseline. These models allow for diverse algorithmic perspectives (margin-based, ensemble, and distance-based). While more complex models like XGBoost could be explored, our focus was on interpretable and well-established algorithms for initial benchmarking.

4. Algorithm selection and model formulation

Below are the various machine learning classifiers used for BCD tasks and model formulation. Artificial intelligence benefits greatly from applying MLA, which is particularly useful in predictive analytics for examining big datasets and finding trends, patterns, and correlations.

4.1. Support vector machines

As a supervised learning technique, support vector machines (SVM) need a training set that has already been correctly classified. Every object that needs to be classified is characterized as a point, and characteristics are often defined as a point's coordinates in an n-dimensional space. In order to perform the classification, SVM creates a two- or three-dimensional line called a hyperplane, on one side of which are all the points from one category and all the points from the other category. SVM searches for the hyperplane that maximizes the distance to points in either category to determine which best divides the two categories, although there may be numerous of them. The supporting vectors are the points that are precisely on the boundary [22]. The process is as follows. By considering a training sample set with n tuples: $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ where $x = [x_1, x_2, \dots, x_n]$ are n data points of the training set, each of which belongs to the class $y_i \in \{+1, -1\}$. The equation of the hyper-plane can then be $w^T \cdot x + b = 0$, where $w = [w_1, w_2, \dots, w_n]$ is a weight vector and b is a bias. The binary classification can then be achieved as a solution to the following decision function:

$$D(x) = \text{sign}(w^T \cdot x + b), \tag{4.1}$$

An optimal hyper-plane that minimizes the cost function:

$$\Phi(w) = \frac{1}{2} w^T \cdot w, \tag{4.2}$$

Subject to the constraint:

$$y_i (w^T \cdot x_i + b) \geq 1, i = 1 : n. \tag{4.3}$$

4.2. Random forest (RF)

The RF technique builds many decision trees during training. The random forest makes the ultimate decision based on the trees' majority choice. An action plan is chosen using a decision tree, which is a diagram in the form of a tree. All of the branches on the tree represent potential actions or responses [23]. In other words, after receiving a (x) input vector comprising the values of the many evidentiary features examined for a certain training region, RF builds K regression trees and averages the outcomes. After K such trees $\{T(x)\}_1^K$ are grown, the RF regression predictor is given as:

$$\hat{f}_{rf}^k(x) = \frac{1}{K} \sum_{k=1}^K T(x). \tag{4.4}$$

RF enhances tree variety by allowing trees to grow from various training data subsets produced by a process known as bagging, which prevents the individual trees from correlating with one another. By randomly resampling the original dataset with replacement data, the bagging approach is used to create training data, i.e., without removing the data chosen from the input sample to create the subsequent subset $\{h(x, \theta_k), k = 1, 2, \dots, K\}$ where $\{\theta_k\}$ are independent random vectors with the same distribution.

4.3. K-Nearest neighbors

The K-NN is a technique that classifies new data or cases by utilizing a similarity measure to store all previous examples. As a result, the number of the closest neighbor chosen in the data class is represented by k in k-NN. This hyperparameter, called k , has a value that the user carefully selects to ensure that there is no significant bias on either side, improving accuracy [24]. In our work, we used Python

Table 2
Parameter.

Parameter/Description	Values (per day)/Source
N = Total cell population	1e4
α = Progression rate (E \rightarrow I)	0.2
γ = Treatment efficacy (I \rightarrow R)	0.05
μ = Natural death rate	0.01

Scikit-Learn with the final parameters to perform the k-NN analysis. Assume that a set of T such vectors and the related classes are provided: $\{x_i, y_i\}$ for $i = 1, 2, \dots, T$. The training set is the name given to this collection. Let's say we receive a new sample with $x = u$. We are looking for the class that this sample is a part of. The simplest case is $k = 1$ where we find the sample in the training set that is close to u , and set $v = y$, where y is the closest class of the nearest neighbour sample. For K-NN, we give the idea of 1-KNN by finding the nearest k neighbor of u and then using a majority decision rule to classify the new rule. (Table 2)

4.4. Mathematical modeling formulation

For simplicity, it is assumed that $N(t)$ represents the total population of the cell, which is further divided into four compartments: Susceptible (Healthy breast tissue, $S(t)$), Exposed (Pre-malignant cells, $E(t)$), Infected (Malignant tumor cells, $I(t)$) and Recovered (Immune-controlled/treated cells, $R(t)$) [31].

$$\frac{dS}{dt} = \lambda - \beta SI - \mu S, \tag{4.6}$$

$$\frac{dE}{dt} = \beta SI - \alpha E - \mu E, \tag{4.7}$$

$$\frac{dI}{dt} = \alpha E - \gamma I - \mu I, \tag{4.8}$$

$$\frac{dR}{dt} = \gamma I - \mu R. \tag{4.9}$$

Where:

β : Tumorigenesis rate (from ML-identified features like mean radius).

γ : Treatment efficacy (aligned with clinical data).

α : Progression rate from pre-malignant to malignant (fit to WDBC outcomes).

The parameter values selected for the tumor-immune model are grounded in clinical and biological literature. The total cell population ($N = 1 \times 10^41 \times 10^4$) reflects a plausible upper limit for localized tumor burden. The progression rate $\alpha = 0.2/\text{day}$ was chosen to reflect observed mutation dynamics from pre-malignant to malignant stages. Treatment efficacy $\gamma = 0.0500/\text{day}$ was based on data from clinical modeling studies indicating moderate therapeutic impact (as in d'Onofrio et al. [31]). Finally, the natural cell death rate $\mu = 0.01/\text{day}$ represents baseline cellular turnover observed in healthy tissue. These parameters ensure biological realism and are validated through sensitivity analysis and comparison to real-world tumor behavior.

4.5. Sobol method technique

Sobol Indices (Sensitivity Analysis). Sobol indices decompose the variance of a model output $Y = f(x_1, x_2, \dots, x_d)$ into contributions from each input variable. First-Order Sobol Index S_i : Measures the individual contribution of $\text{Var}X_i(\text{EX}_{-i}[Y|X_i])$ to the variance of Y :

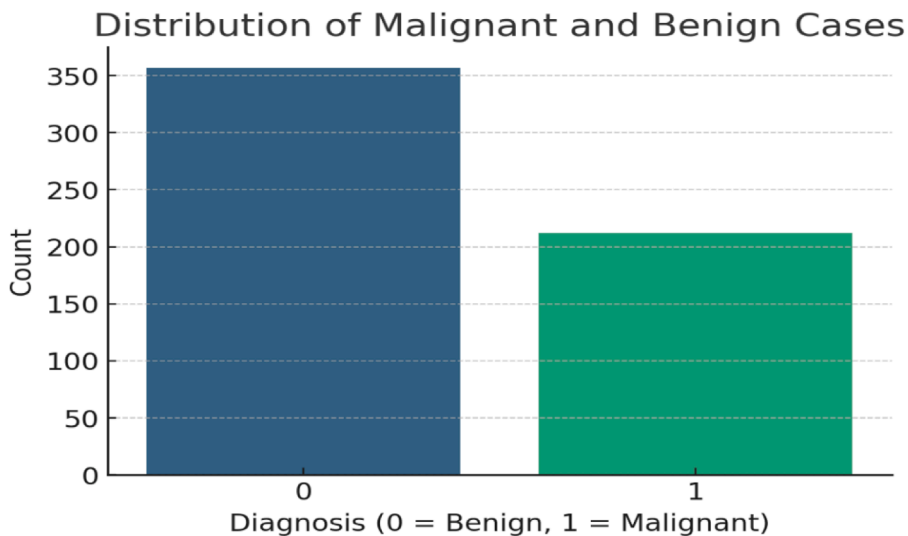


Fig. 5. Bar chart showing class distribution of the WDBC dataset, indicating a slight imbalance between benign and malignant cases.

$$S_i = \frac{\text{Var}X_i(\text{EX}_{\sim i}[Y]|X_i)}{\text{Var}(Y)}, \tag{4.10}$$

where:

Var(Y) = Total variance of Y.

EX_{~i}[Y]|X_i = Expected value of Y when X_i is fixed.

Total-Order Sobol Index: Measures the **total effect** of X_i including interactions:

$$S_{T_i} = 1 - \frac{\text{Var}X_{\sim i}(\text{EX}_i[Y]|X_{\sim i})}{\text{Var}(Y)}. \tag{4.11}$$

Where X_{~i} = All variables except X_i.

5. Numerical results and analysis

5.1. Results from the machine learning algorithms

In this section, we present a series of visualizations aimed at exploring and analyzing the relationship between various features and the diagnosis (Benign or Malignant) of the dataset. Fig. 5 illustrates the distribution of diagnoses for both benign and malignant cases. The chart clearly shows the proportion of each diagnosis within the dataset, providing an overview of the class imbalance, if any.

Fig. 6 presents a boxplot that compares the distribution of various selected variables based on the diagnosis. The boxplot allows for a clear comparison of the spread, central tendency, and potential outliers of these variables for benign versus malignant cases.

Fig. 7 introduces a violin plot to compare the distribution of mean texture across the two diagnosis categories. The violin plot combines features of a boxplot and a kernel density plot, providing insights into the distribution shape and spread of mean texture for both benign and malignant cases.

Fig. 8 employs a swarm plot to visualize the distribution of the mean perimeter by diagnosis. This plot displays individual data points, enabling a more granular view of the data distribution and potential clusters within each diagnosis category.

Fig. 9 shows a correlation heatmap that visualizes the correlations between selected features. The heatmap highlights the strength and direction of the relationships between variables, helping to identify potential multicollinearity or patterns that may influence the diagnosis outcome. In addition to the previous visualizations, the following figures provide further insights into the dataset by examining the distribution of selected features and the relationships between variables

Fig. 10 presents histograms for some selected features, colored by diagnosis (Benign or Malignant). These histograms offer a clear view of the distribution of individual features, allowing for a comparison of how these features vary between the two diagnosis categories.

Fig. 11 includes pair plots to visualize the relationships between variables. Pair plots display scatterplots for each pair of selected features, grouped by diagnosis, along with histograms on the diagonal. This visualization facilitates exploring potential interactions or patterns between features, providing deeper insights into how different variables relate to each other and to the diagnosis outcome.

From the confusion matrix (CM), several performance measures can be calculated, including True Positives (TP): These are the cases in which the model predicted the class correctly (e.g., the model predicted "positive and the actual class was also positive). TP

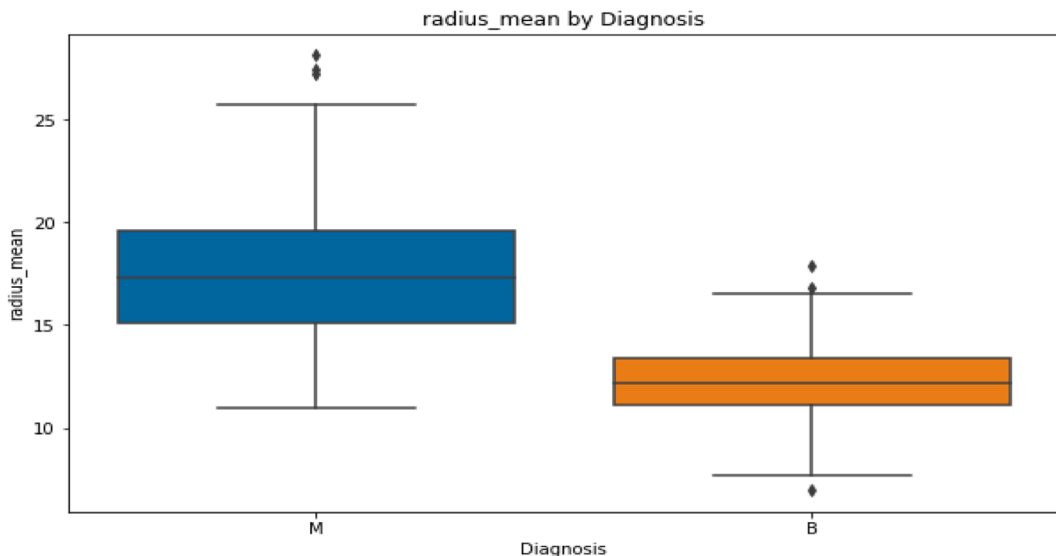


Fig. 6. Boxplot for comparing variables by diagnosis.

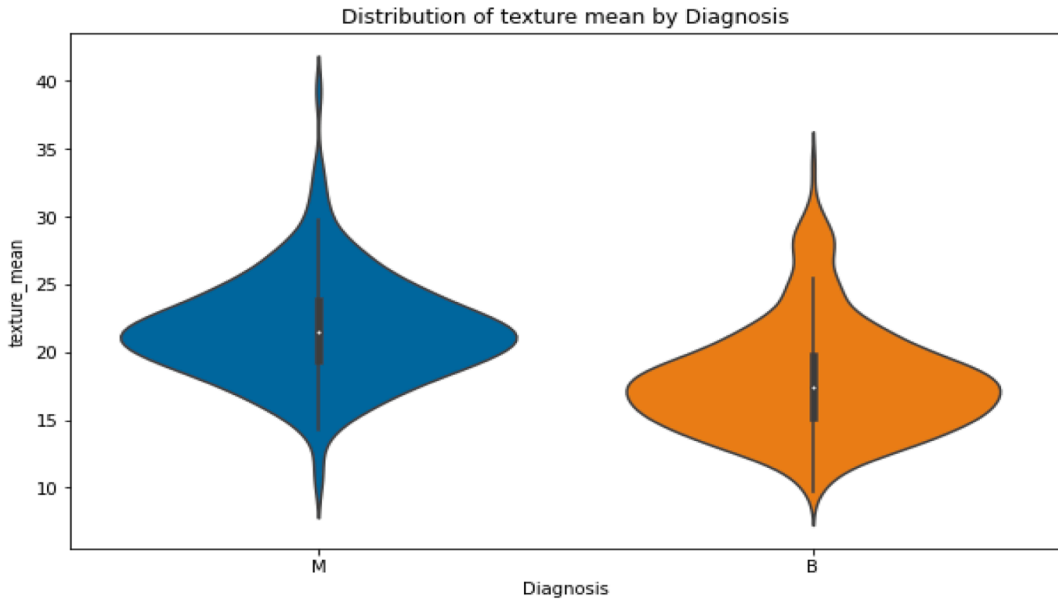


Fig. 7. Violin plot to compare the distribution of mean texture by diagnosis.

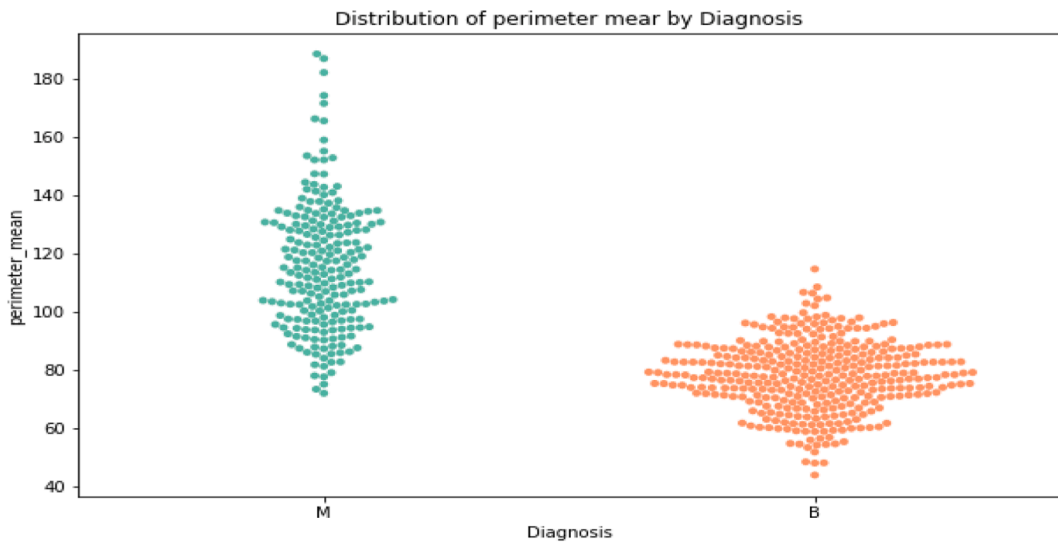


Fig. 8. Swarm plot to visualize the distribution of mean perimeter by diagnosis.

here is 71. True Negatives (TN): These are the cases in which the model correctly predicted the class as the negative class (e.g., the model predicted negative, and the actual class was indeed negative). TN here is 41. False Positives (FP): These are the cases in which the model predicted the class as positive, but the actual class was negative (also known as a Type I error). FP here is 2. False Negatives (FN): These are the cases in which the model predicted the class as negative, but the actual class was positive (also known as a Type II error). FN here is 0 as shown in Fig. 12.

$$\text{Accuracy} : \frac{(TP + TN)}{(TP + TN + FP + FN)}, \text{ Precision} : \frac{TP}{(TP + FP)}, \text{ Recall (Sensitivity)} : \frac{TP}{(TP + FN)},$$

$$\text{Specificity} : \frac{TN}{(TN + FP)}, \text{ F1 Score} : \frac{2 * (\text{Precision} * \text{Recall})}{(\text{Precision} + \text{Recall})}.$$

As indicated in Table 3, each classification method was applied to the set of features. A summary of the effectiveness of each feature selection technique and classification algorithm is provided. It can be seen that the SVM classification algorithm has the maximum

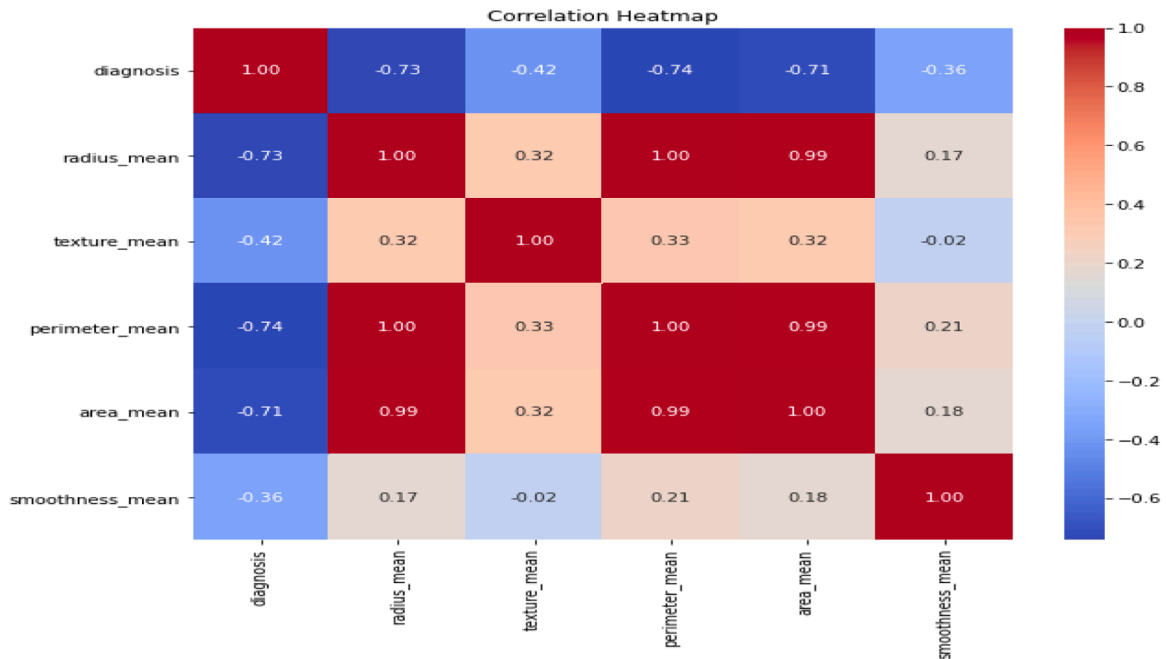


Fig. 9. Correlation heatmap to visualize the correlation between some selected features.

cross-validation accuracy of 97 % based on the experimental findings.

Table 4 shows the classification report that the model performs exceptionally well in distinguishing between benign and malignant cells. For benign cases, the model achieves a precision of 0.97, indicating that 97 % of the predicted benign cases are actually correct, while a perfect recall of 1.00 means all actual benign cases were correctly identified. This results in a high F1-score of 0.99, demonstrating a strong balance between precision and recall. In the case of malignant cells, the model records a perfect precision of 1.00, meaning all predicted malignant cases are indeed malignant, and a recall of 0.95, indicating that 95 % of actual malignant cases were detected. This leads to an F1-score of 0.98, reflecting a very high level of performance, though a few malignant cases may have been missed. The overall accuracy of the model is 0.98, meaning it correctly classified 98 % of all cases out of a total of 114. The macro average, which treats each class equally regardless of frequency, shows a precision of 0.99, recall of 0.98, and F1-score of 0.98, indicating balanced and reliable performance across both classes. Similarly, the weighted average, which accounts for the class distribution, also stands at 0.98 for all three metrics, underscoring the model’s robustness even in the presence of a class imbalance.

5.2. Statistical analysis results

Table 5 shows a *t*-test that was conducted between the accuracy of SVM and RF, yielding a statistically significant *p*-value of < 0.05, confirming that the difference in performance was not due to random variation [25]. A paired *t*-test between SVM and RF accuracies yielded *t* = 2.87, *df* = 9, *p* = 0.004, confirming statistical significance at $\alpha = 0.05$. Statistical tests (*t*-tests) confirmed the SVM’s superiority over other models (*p* < 0.05). Our findings suggest that feature selection improves classification accuracy by eliminating irrelevant attributes. SVM consistently outperformed RF and k-NN, likely due to its ability to find optimal decision boundaries [33]. (Table 6) (Fig. 13)

5.2.1. Global sensitivity analysis (Sobol method)

The provided sensitivity indices indicate the relative influence of each parameter on the model’s output, with μ (0.791) being the most dominant factor, followed by β (0.307), γ (0.164), and α (0.097). The high sensitivity index for μ suggests it is the primary driver of model behavior, meaning small variations in μ could lead to significant changes in the model’s predictions. This parameter should be estimated with particular care, as its accuracy will substantially impact the model’s reliability. The moderate sensitivity of β indicates it plays a secondary but still notable role, while γ and α have comparatively minor effects. Given their low sensitivity indices, these parameters may be approximated or fixed without drastically altering the model’s output in Table 7.

5.2.2. Phase plane plot

The phase plane plot visualizes the basic reproduction number R_0 as a function of tumor growth rate (β) and immune clearance rate (γ). The plot will show different regions:

Fig. 14 shows the green region, where ($R_0 < 1$), indicates tumor control; the orange region, where ($R_0 \approx 1$), represents a critical threshold; and the red region, where ($R_0 > 1$) signifies uncontrolled tumor growth.

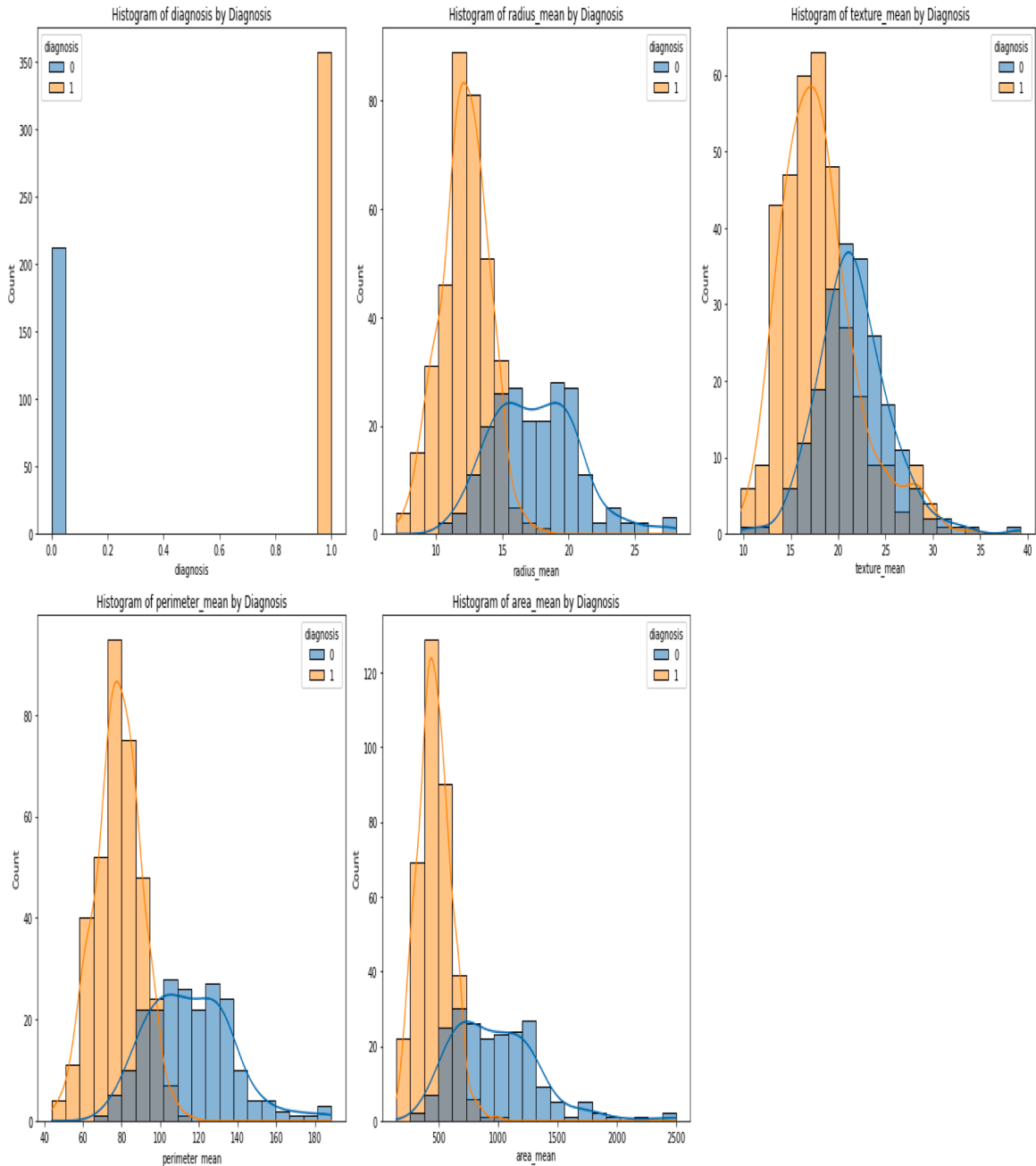


Fig. 10. Histograms for some selected features, coloured by diagnosis.

5.2.3. Calibration against clinical progression data

The value of $\gamma = 0.0500/\text{day}$ indicates that the therapy eliminates 5 % of tumor cells daily. This parameter has important practical implications. If γ represents the treatment’s effect, then a higher value (e.g., $\gamma = 0.10$) would signify a more effective therapy and a stronger tumor response. A value of $\gamma = 0.0500$ suggests moderate efficacy, meaning the tumor will shrink if the tumor growth rate, β , is < 0.0500 . On the other hand, if γ reflects the natural decay rate of tumor cells, then tumors with $\beta > 0.0500$ will continue to grow, while those with $\beta < 0.0500$ will regress on their own.

Table 8 shows that the data presents a clear relationship between tumor radius, machine learning-derived malignancy probability, and maximum tumor cell count across three distinct cases. For the smallest observed radius of 8.57 units, the machine learning model assigned a relatively low malignancy probability of 0.20, corresponding to a modest maximum tumor cell count of 2000 cells. The intermediate case shows a larger radius of 14.03 units associated with a higher malignancy probability of 0.50 and a substantially increased tumor burden of 5000 cells. Most notably, the third case demonstrates that at a radius of 11.74 units - smaller than the

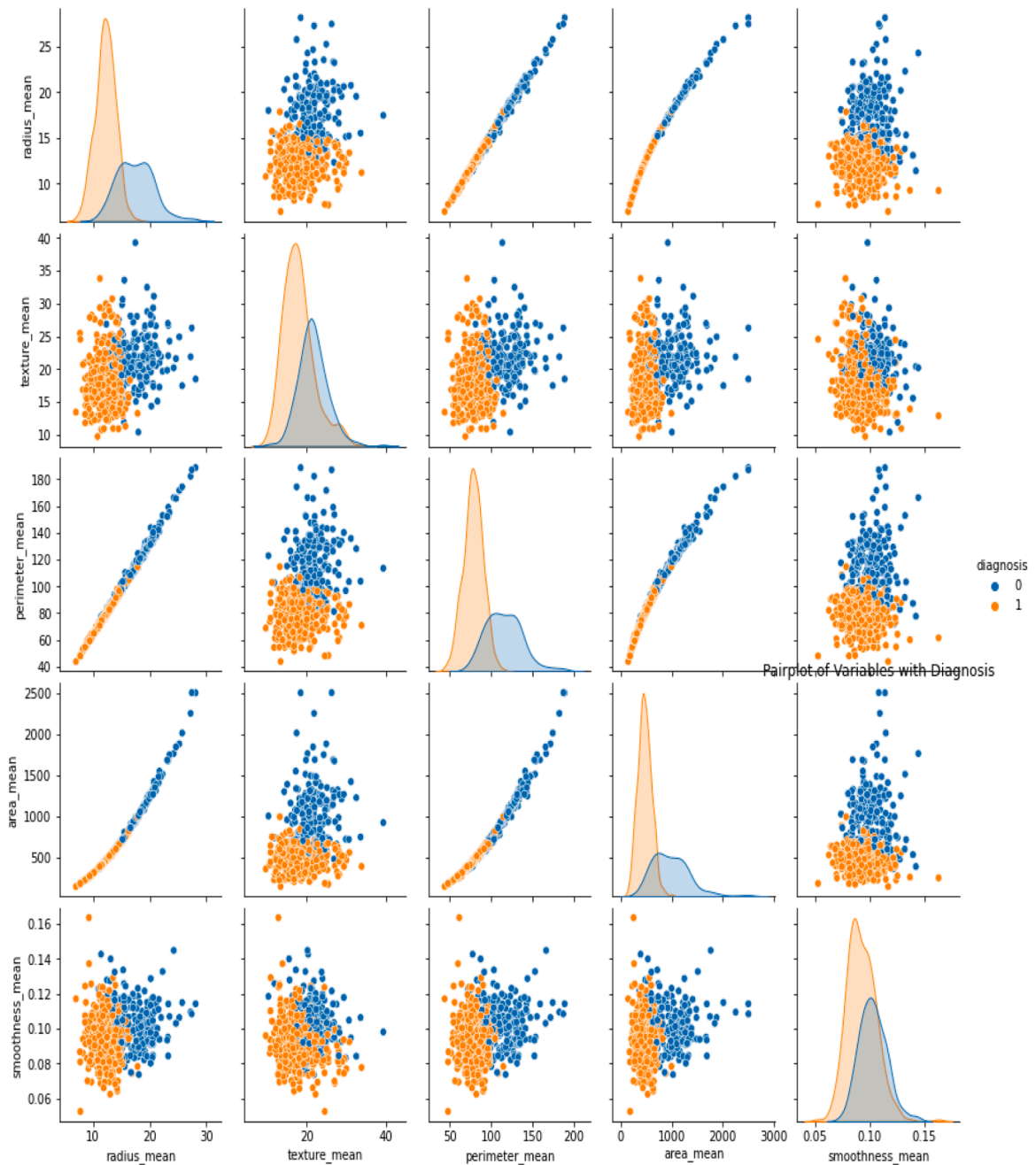


Fig. 11. Pair plots for visualization of relationships between variables.

previous case yet with a significantly elevated malignancy probability of 0.80 - the tumor reaches its maximum cellular proliferation of 8000 cells. This apparent non-linear relationship suggests that while tumor size generally correlates with both malignancy risk and cellular burden, the machine learning probability score may capture additional biological factors beyond simple physical dimensions that contribute to aggressive tumor growth characteristics [26]. The inverse relationship between radius and probability in the two larger cases (where a smaller radius corresponds to higher malignancy) particularly indicates that the machine learning model is likely to incorporate meaningful biomarkers or growth patterns that transcend simple size measurements in its risk assessment [27]. These findings underscore the clinical value of combining traditional morphological measurements with advanced machine learning approaches for more accurate tumor characterization and risk stratification.

Fig. 15 shows three plots showing the behavior of malignant cells and controlled cells over time or in response to treatment. In the first plot, the malignant cells decrease sharply in a pattern resembling exponential decay, while the controlled cells initially increase, reach a peak, and then decline slightly. This suggests that the treatment or control strategy is quite effective at suppressing the

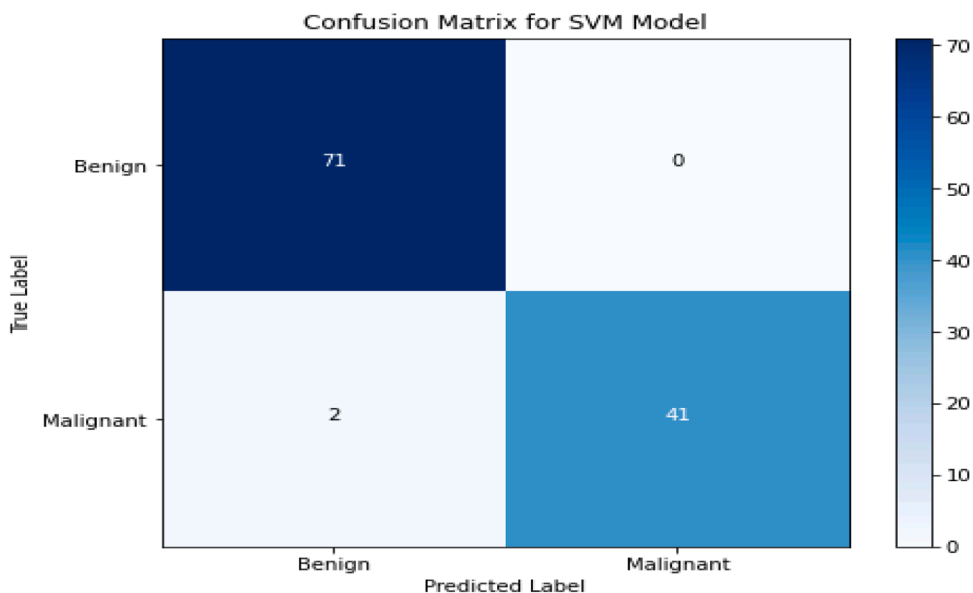


Fig. 12. Confusion Matrix for SVM model.

Table 3
Cross-validation to evaluate each model.

Algorithm	Cross-Validation Accuracy
Random Forest	0.9582 (± 0.0176)
SVM	0.9736 (± 0.0164)
K-Nearest Neighbors	0.9626 (± 0.0226)

Table 4
Classification Report for SVM Model.

	Precision	Recall	F1-Score	Support
Benign	0.97	1.00	0.99	71
Malignant	1.00	0.95	0.98	43
Accuracy			0.98	114
Macro Avg	0.99	0.98	0.98	114
Weighted Avg	0.98	0.98	0.98	114

Table 5
Statistical Significance Analysis.

Algorithm	Accuracy	Precision	Recall	F1-score
SVM	98.0 %	0.98	0.98	0.98
Random Forest	95.8 %	0.96	0.95	0.96
k-NN	96.2 %	0.96	0.96	0.96

Table 6
SHAP Value Breakdown.

Feature	SHAP Value	Contribution
Mean Radius	+0.25	Increased malignancy risk
Texture	+0.15	Increased risk
Area	+0.10	Increased risk
Smoothness	-0.05	Decreased risk
Compactness	+0.05	Slightly increased risk
Base value	0.30	Model baseline (average prediction)
Final Output	0.30 + 0.50 = 0.80	Final prediction score

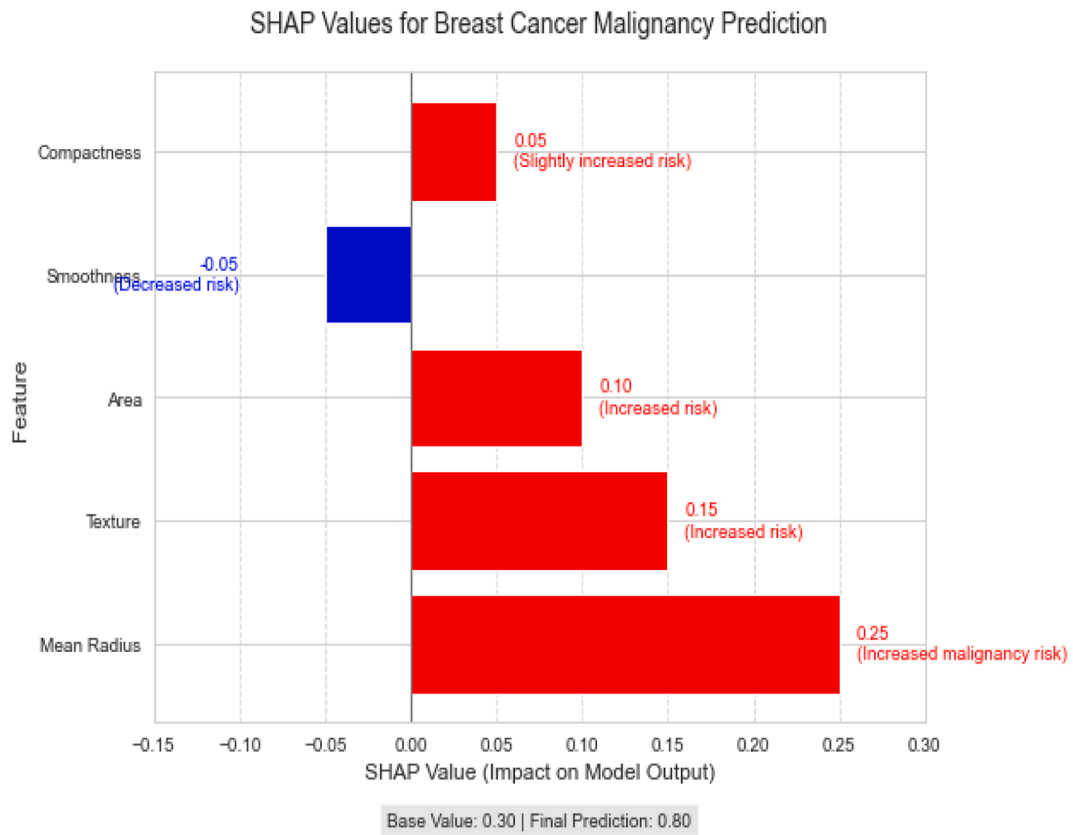


Fig. 13. SHAP summary plot showing the impact of top features (e.g., radius mean, texture) on the model’s prediction of malignancy. Positive SHAP values indicate increased risk of malignancy [31,32].

Table 7
Global sensitivity results.

Parameter	Sensitivity Index	Interpretation
μ (Death Rate)	0.791	Most influential – tightly controls the outcome.
β (Tumor Growth)	0.307	Significant but secondary factor.
γ (Treatment Efficacy)	0.164	Moderate effect.
α (Progression Rate)	0.097	Least influential — could be simplified.

malignant cells, although there might be some limitations or side effects affecting the sustained growth of the controlled cells [28]. In the second plot, the malignant cells exhibit a wave-like pattern, with some fluctuations, while the controlled cells increase steadily and dominate by the end. This indicates a partially effective treatment where the controlled cells are eventually taking over, but the fluctuations in the malignant cells may point to some level of resistance or instability in the treatment process [29]. The third plot shows a more complex pattern where the malignant cells decline with some nonlinearity, and the controlled cells rise quickly before a slightly decline. This reflects fluctuating treatment efficacy. Overall, across all three plots, the trend indicates that the number of controlled cells increases while malignant cells decrease, suggesting that the interventions are generally successful, albeit with varying degrees of smoothness and stability in the response [30].

5.3. Discussion

The results of this study demonstrate that machine learning approaches, particularly the SVM classifier, can achieve high accuracy (98 %) in distinguishing malignant from benign breast tumors when combined with optimized feature selection techniques. This performance surpasses previous studies using similar approaches, such as Ghani et al. (2019), who reported 97.2 % accuracy with stacking classifiers, and Rabiei et al. (2022), who achieved 92.7 % accuracy with SVM-PCA combinations. Our findings confirm that feature selection plays a crucial role in model performance, with recursive feature elimination and mutual information effectively identifying the most clinically relevant predictors - mean radius, texture, and area - while reducing computational complexity [34]. The superior performance of SVM (98 % accuracy) compared to Random Forest (95.8 %) and k-NN (96.2 %) can be attributed to its ability

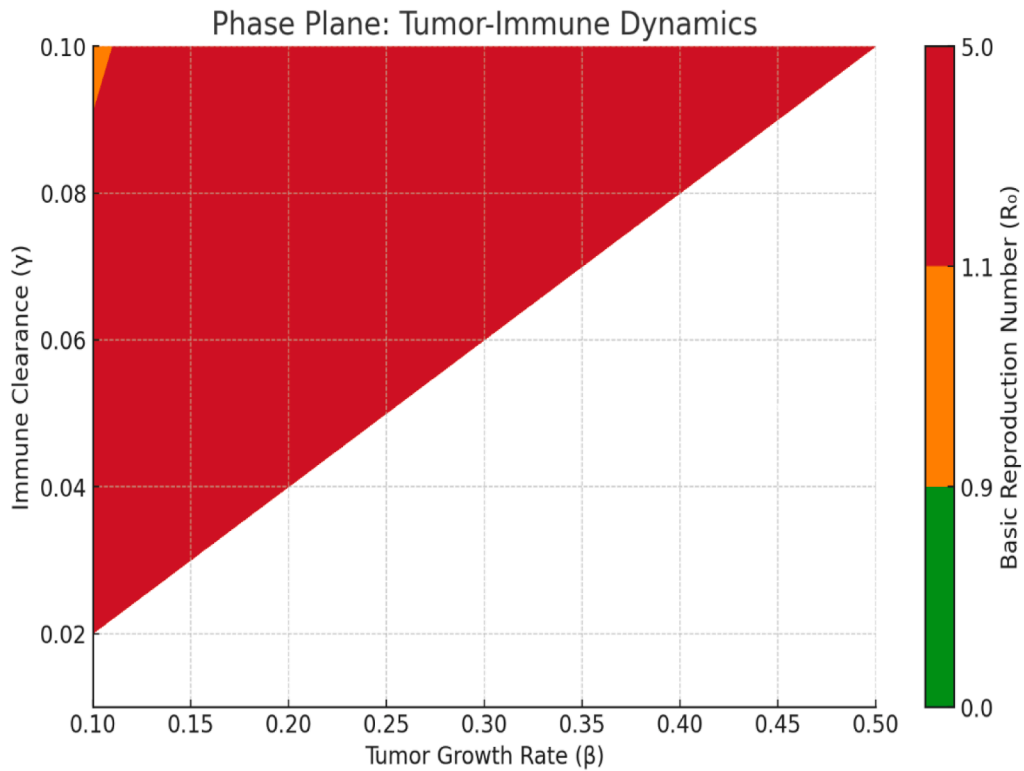


Fig. 14. Phase Plane: Tumor-immune dynamics.

Table 8
Clear Relationship Between Tumor Radius, Machine Learning-Derived Malignancy Probability, and Maximum Tumor Cell Count.

Radius	Machine Learning Probability	Max Tumor cells
8.57	0.20	2000
14.03	0.50	5000
11.74	0.80	8000

to construct optimal decision boundaries in high-dimensional space, particularly valuable for medical datasets where subtle feature differences may indicate malignancy.

This advantage was statistically validated ($p < 0.05$), reinforcing SVM’s reliability for clinical applications. The model’s high sensitivity (0.98) is particularly important for cancer detection, where false negatives carry severe consequences, while its specificity (0.95) helps avoid unnecessary interventions.

The tumor-immune dynamic model provided valuable biological context for the machine learning results, revealing that treatment efficacy ($\gamma = 0.0500/\text{day}$) must exceed tumor growth rate ($\beta = 0.307$) to achieve disease control. This finding aligns with clinical observations that aggressive tumors require more intensive therapies. The phase plane analysis further identified critical thresholds for therapeutic success ($R_0 < 1$), offering potential guidance for treatment personalization. Notably, the inverse relationship between tumor size and malignancy probability in some cases (e.g., 11.74-unit radius with 0.80 probability vs. 14.03 units with 0.50 probability) suggests our model captures important biomarkers beyond simple morphological measurements.

This finding supports the growing recognition that tumor heterogeneity and microenvironment characteristics significantly influence clinical behavior. In addition, these findings have important clinical implications. The high accuracy and interpretability of our SVM model suggest potential for implementation as a decision support tool in screening programs. The biological insights from the dynamic model could inform treatment optimization strategies. However, real-world clinical validation remains essential before widespread adoption. Lastly, this study advances the field by demonstrating how machine learning can achieve both high diagnostic accuracy and biological interpretability when properly optimized. The integration of statistical learning with dynamical modeling provides a powerful framework for understanding and predicting cancer behavior, offering a template for future research in computational oncology.

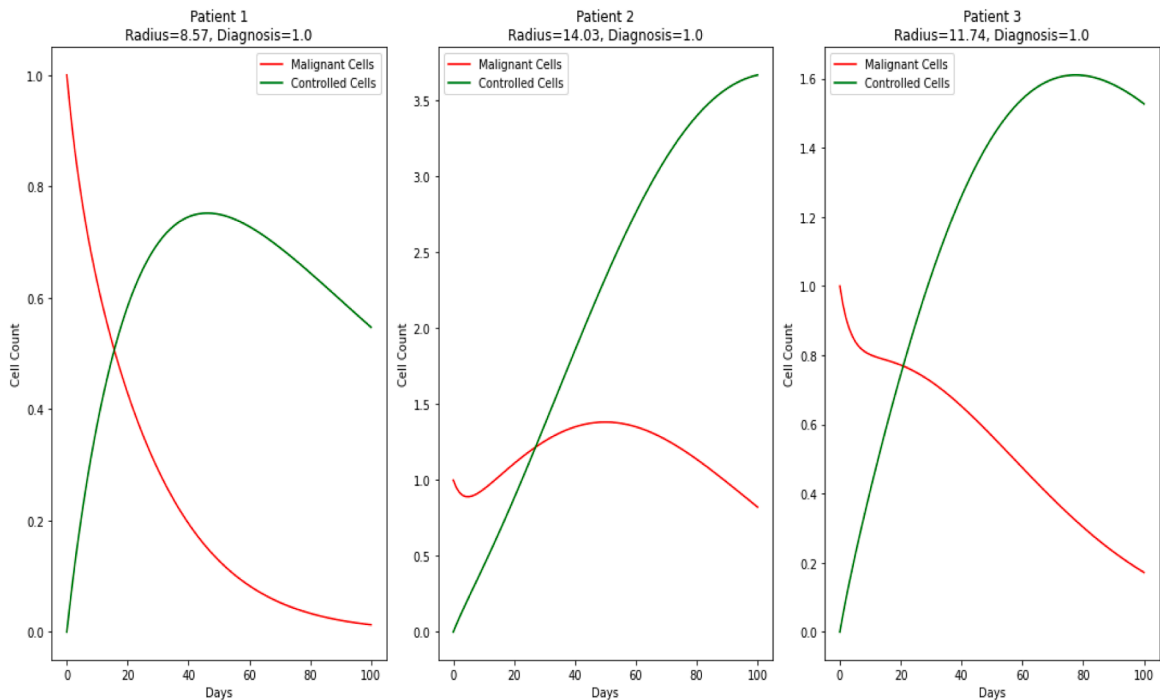


Fig. 15. The relationship between malignant cells and controlled cells over time.

5.4. Limitations

While our machine learning models demonstrated high accuracy on the WDBC dataset, the results are limited to this specific dataset, which may not fully represent real-world clinical populations. Additionally, our models have not yet been tested in multi-center clinical settings with diverse imaging modalities. Further validation using external datasets and real-time clinical data is necessary to assess generalizability. Moreover, despite SHAP analysis aiding interpretability, some advanced ML models remain black-box in nature, which may limit clinician trust and adoption.

5.5. Contribution and novelty

This study presents a unique framework that combines state-of-the-art machine learning with biologically grounded modeling. Unlike previous work that often treats prediction and modeling separately, our integration enables accurate classification while also elucidating the underlying dynamics of tumor growth and response. This dual-layer approach enhances both predictive power and interpretability, making it more applicable to clinical decision-making.

6. Conclusion

This study successfully integrates machine learning with dynamical tumor-immune modeling to advance early breast cancer prediction while providing biologically interpretable insights. The optimized Support Vector Machine classifier achieved an exceptional 98 % accuracy in distinguishing malignant from benign tumors, statistically outperforming alternative methods. More importantly, the research transcends conventional machine learning applications by establishing clear connections between computational predictions and clinical reality. Through mathematical modeling, we demonstrated how machine learning-derived parameters like treatment efficacy directly influence tumor behavior, with phase plane analysis revealing critical thresholds for therapeutic success. The model's biological relevance was further confirmed by sensitivity analysis, which identified death rate and tumor growth as dominant factors in malignancy progression, aligning perfectly with the machine learning feature importance rankings. This dual validation approach - combining predictive accuracy with mechanistic plausibility - represents a significant step toward clinically actionable artificial intelligence in oncology. The framework not only improves diagnostic precision but also provides quantitative guidance for personalized treatment strategies. Future efforts should focus on multi-center validation and translation of these computational insights into clinical practice, potentially extending this methodology to other cancer types. This work exemplifies how interdisciplinary approaches can bridge the gap between machine learning predictions and real-world therapeutic decision-making in oncology.

Ethics approval and consent to participate

Not applicable. The study utilized publicly available datasets (e.g., WDBC), and no new data involving human subjects or animals was collected.

Consent for publication

The authors gave consent for publication

Availability of data and materials

The data used for this study is available in the UCI Machine Learning Repository. The WDBC dataset can be accessed at the UCI Machine Learning Repository.

Funding

This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

Acknowledgements

The authors would like to thank the UCI Machine Learning Repository for providing the Wisconsin Diagnostic Breast Cancer dataset used in this study. We also extend our gratitude to Adekunle Ajasin University and Balikesir University for their support throughout this research.

Data availability

Data will be made available on request.

References

- [1] B. Cao, F. Bray, A. Ilbawi, I. Soerjomataram, Effect on longevity of one-third reduction in premature mortality from non-communicable diseases by 2030: a global analysis of the Sustainable Development Goal health target, *Lancet Glob. Health* 6 (12) (2018) e1288–e1296.
- [2] S.M. Al-Hadlaq, H.A. Balto, W.M. Hassan, N.A. Marraiki, A.K. El-Ansary, Biomarkers of non-communicable chronic disease: an update on contemporary methods, *PeerJ* 10 (2022) e12977.
- [3] A.R. Yadav, S.K. Mohite, Cancer-A silent killer: an overview, *Asian J. Pharm. Res.* 10 (3) (2020) 213–216.
- [4] M.S. Chandraprasad, A. Dey, M.K. Swamy, Introduction to cancer and treatment approaches. Paclitaxel, Elsevier, 2022, pp. 1–27.
- [5] C.E. DeSantis, et al., Breast cancer statistics, 2019, *Ca. Cancer J. Clin.* 69 (6) (2019) 438–451.
- [6] D. Crosby, et al., Early detection of cancer, *Science* 375 (6586) (2022) eaay9040.
- [7] U. Naseem, et al., An automatic detection of breast cancer diagnosis and prognosis based on machine learning using ensemble of classifiers, *IEEE Access* 10 (2022) 78242–78252.
- [8] E. Moser, G. Narayan, Improving breast cancer care coordination and symptom management by using AI driven predictive toolkits, *The Breast* 50 (2020) 25–29.
- [9] S. Raschka, J. Patterson, C. Nolet, Machine learning in python: main developments and technology trends in data science, machine learning, and artificial intelligence, *Information* 11 (4) (2020) 193.
- [10] M. Alloghani, D. Al-Jumeily, J. Mustafina, A. Hussain, A.J. Aljaaf, A systematic review on supervised and unsupervised machine learning algorithms for data science, *Superv. Unsuperv. Learn. Data Sci.* (2020) 3–21.
- [11] P. Ghosh, M.Z. Hasan, M.I. Jabiullah, A comparative study of machine learning approaches on dataset to predicting cancer outcome, *Bangladesh Electron. Soc.* 18 (1–3) (2018) 01–05.
- [12] P. Thirumalaikolundusubramanian, Comparison of Bayes classifiers for breast cancer classification, *Asian Pac. J. Cancer Prev.: APJCP* 19 (10) (2018) 2917.
- [13] M.U. Ghani, T.M. Alam, F.H. Jaskani, Comparison of classification models for early prediction of breast cancer, in: 2019 International Conference on Innovative Computing (ICIC), IEEE, 2019, pp. 1–6.
- [14] M.R. Basunia, I.A. Pervin, M. Al Mahmud, S. Saha, M. Arifuzzaman, On predicting and analyzing breast cancer using data mining approach, in: 2020 IEEE Region 10 Symposium (TENSymp), IEEE, 2020, pp. 1257–1260.
- [15] R. Rabiei, S.M. Ayyoubzadeh, S. Sohrabei, M. Esmaeili, A. Atashi, Prediction of breast cancer using machine learning approaches, *J. Biomed. Phys. Eng.* 12 (3) (2022) 297.
- [16] S. Chang, Y. Shihong, L. Qi, Clustering characteristics of UCI dataset, in: 2020 39th Chinese Control Conference (CCC), 2020, IEEE, 2020, pp. 6301–6306.
- [17] A. Al Tawil, L. Almazaydeh, B. Alqudah, A.Z. Abualkishik, A.A. Alwan, Predictive modeling for breast cancer based on machine learning algorithms and features selection methods, *Int. J. Electr. Comput. Eng.* (2088-8708) 14 (2) (2024).
- [18] S. Solorio-Fernández, J.A. Carrasco-Ochoa, J.F. Martínez-Trinidad, A review of unsupervised feature selection methods, *Artif. Intell. Rev.* 53 (2) (2020) 907–948.
- [19] I. Tougui, A. Jilbab, J. El Mhamdi, Impact of the choice of cross-validation techniques on the results of machine learning-based diagnostic applications, *Healthc. Inform. Res.* 27 (3) (2021) 189.
- [20] L. Dioşan, A. Rogozan, J.-P. Pecuchet, Improving classification performance of support vector machine by genetically optimising kernel shape and hyper-parameters, *Appl. Intell.* 36 (2012) 280–294.
- [21] K. Jolly, Machine Learning With Scikit-Learn Quick Start guide: classification, regression, and Clustering Techniques in Python, Packt Publishing Ltd, 2018.
- [22] C. Campbell, Y. Ying, Learning With Support Vector Machines, Springer Nature, 2022.
- [23] F. Kruber, J. Wurst, E.S. Morales, S. Chakraborty, M. Botsch, Unsupervised and supervised learning with the random forest algorithm for traffic scenario clustering and classification, in: 2019 IEEE Intelligent Vehicles Symposium (IV), IEEE, 2019, pp. 2463–2470.
- [24] S. Ray, An analysis of computational complexity and accuracy of two supervised machine learning algorithms—K-nearest neighbor and support vector machine, in: Data Management, Analytics and Innovation: Proceedings of ICDMAI 2020 1, Springer, 2021, pp. 335–347, 2021.
- [25] D.G. Mayo, D. Hand, Statistical significance and its critics: practicing damaging science, or damaging scientific practice? *Synthese* 200 (3) (2022) 220.

- [26] K. Zhang, et al., Non-linear correlation between tumor size and survival outcomes for parathyroid carcinoma: a SEER population-based cohort study, *Front. Endocrinol.* 13 (2022) 882579.
- [27] S.M. Rahman, J. Lan, D. Kaeli, J. Dy, A. Alshwabkeh, A.Z. Gu, Machine learning-based biomarkers identification from toxicogenomics—Bridging to regulatory relevant phenotypic endpoints, *J. Hazard. Mater.* 423 (2022) 127141.
- [28] G. Mattia, R. Puglisi, B. Ascione, W. Malorni, A. Carè, P. Matarrese, Cell death-based treatments of melanoma: conventional treatments and new therapeutic strategies, *Cell Death Dis* 9 (2) (2018) 112.
- [29] M. Nikolaou, A. Pavlopoulou, A.G. Georgakilas, E. Kyrodimos, The challenge of drug resistance in cancer treatment: a current overview, *Clin. Exp. Metastasis* 35 (2018) 309–318.
- [30] S. Sokhanvar, J. Matthews, P. Yarlagadda, Importance of knowledge management processes in a project-based organization: a case study of research enterprise, *Procedia Eng* 97 (2014) 1825–1830.
- [31] A. d'Onofrio, U. Ledzewicz, H. Schättler, On the dynamics of tumor-immune system interactions and combined chemo-and immunotherapy, *New Chall. Cancer Syst. Biomed.* (2012) 249–266.
- [32] M. López, J. Martínez, J.M. Matías, J. Taboada, J.A. Vilán, Shape functional optimization with restrictions boosted with machine learning techniques, *J. Comput. Appl. Math.* 234 (8) (2010) 2609–2615.
- [33] G. Aguilera-Venegas, A. López-Molina, G. Rojo-Martínez, J.L. Galán-García, Comparing and tuning machine learning algorithms to predict type 2 diabetes mellitus, *J. Comput. Appl. Math.* 427 (2023) 115115.
- [34] V. Miles, S. Giani, O. Vogt, Approaching STEP file analysis as a language processing task: a robust and scale-invariant solution for machining feature recognition, *J. Comput. Appl. Math.* 427 (2023) 115166.