

Project Report For Machine Learning-SVM Model

Enhancing Breast Tumor Diagnosis: Leveraging Support Vector Machine Models

Acquah, Theophilus B. K.

tbak16@gmail.com

November 24, 2023



Contents

1	Introduction	2
2	Methodology	3
3	The Data	3
4	Exploratory Data Analysis	5
5	Support Vector Machine (SVM)	10
5.1	Step 1: Selecting Important Parameters	10
5.2	Step 2: Fitting the SVM Model with Key Parameters	11
5.3	Step 3: Model Performance	12
5.4	Step 4: Model Accuracy - Receiver Operating Characteristic (ROC) curve	14
5.5	Step 5: Variable Importance	14
5.6	Step 6: Decision Boundaries	15
	Conclusion	17
6	Limitations of Study	17

List of Figures

1	Plot of Histograms of mean variables	7
2	Plot of Histograms of mean variables	7
3	Plot of Histograms of se variables	8
4	Plot of Histograms of se variables	8
5	Plot of Histograms of worst variables	8
6	Plot of Histograms of worst variables	8
7	Correlation matrix	9
8	Mean Variables Correlation Matrix	9
9	Performance of Train dataset	12
10	Confusion Matrix for Test Data	13
11	ROC Plot for Test Data	14
12	Decision boundaries for Perimeter Worst and Concave points worst	16
13	SVM Decision Boundary and Support Vectors	16
14	SVM Decision Boundary and Margin	16

List of Tables

1	Breast Cancer Dataset Variables Description	4
2	Summary Statistics for Selected Variables	5
3	Summary Statistics for Selected Variables	6
4	Summary Statistics for Variables: Fractal Dimension Mean to Radius SE	6
5	Summary Statistics for Variables: Area SE to Perimeter Worst	6
6	Summary Statistics for Variables: Concave Points SE to Perimeter Worst	6
7	Summary Statistics for Variables: Texture Worst to Perimeter Worst	6
8	Summary Statistics for Selected Variables	7
9	Summary Statistics for Selected Variables	7
10	Best Parameters for Different SVM Kernels	11
11	Summary of SVM Model Fitting with Key Parameters	12
12	Model Evaluation Indicators	13
13	Variable Importance Scores for Breast Cancer Diagnosis	15

Abstract

This project presents a comprehensive analysis of Support Vector Machines (SVM) in the classification of breast cancer tumors from fine-needle aspiration (FNA) test results. We explore the effectiveness of SVM in differentiating between malignant and benign tumors and identify the most significant features that contribute to classification accuracy. Our methodology encompasses a multi-phase process, including exploratory data analysis, model development, and optimization of the SVM classifier. Our results indicate that the SVM model with a linear kernel, achieved a classification accuracy of 97.06% on the test dataset, with a sensitivity of 95.24% and a specificity of 98.13%. The model's precision is evidenced by positive and negative predictive values of 96.77% and 97.22%, respectively. Furthermore, the model's ability to discriminate between tumor classes is demonstrated by an Area Under the Curve (AUC) of 96.68%, as determined by Receiver Operating Characteristic (ROC) curve analysis. The study's findings highlight the SVM model's potential as a diagnostic tool in the medical field, with the capacity to provide early and accurate diagnosis of breast cancer, leading to improved patient outcomes. The incorporation of this machine learning approach in clinical settings could revolutionize breast cancer diagnostics, offering a reliable, non-invasive diagnostic alternative.

1 Introduction

Breast cancer represents one of the most significant health challenges in modern medicine, being the most commonly diagnosed cancer among women worldwide and the second leading cause of cancer death among women. It arises from the abnormal growth of cells in the breast tissue, manifesting as a tumor, which can be benign (non-cancerous), pre-malignant (pre-cancerous), or malignant (cancerous). Diagnosing breast cancer typically involves various tests, including MRI, mammogram, ultrasound, and biopsy. The present study focuses on the use of Predictive Analytics to classify breast cancer tumors based on fine-needle aspiration (FNA) tests. FNA is a minimally invasive procedure, similar to drawing a blood sample but used to extract cells from a breast lesion or cyst for examination.

1.1 Expected Outcome

The goal of this study is to build a predictive model capable of classifying breast cancer tumors as either malignant (cancerous) or benign (non-cancerous) based on the results from FNA tests. The classification is binary: '1' indicating the presence of malignant (cancerous) cells and '0' indicating benign (non-cancerous) cells.

1.2 Objective

Given the discrete nature of the labels in our dataset (malignant or benign), this problem falls into the category of a classification problem in machine learning. The objective is to accurately classify the nature of breast cancer (benign or malignant) and to predict the recurrence and non-recurrence of malignant cases after a certain period. This involves employing machine learning classification methods to fit a function capable of predicting the discrete class of new inputs.

1.3 Research Questions

Given the challenges in diagnosing breast cancer and the potential of machine learning in predictive analytics, the following research questions frame this study:

- **RQ1:** How effectively can machine learning algorithms, specifically Support Vector Machines (SVM) classify breast cancer tumors as malignant or benign based on fine-needle aspiration (FNA) test results?
- **RQ2:** What are the key features in the FNA test results that most significantly contribute to the accurate classification of breast cancer tumors?

These research questions aim to explore the comprehensive capabilities of machine learning in the context of breast cancer diagnosis, from classification accuracy to feature importance and

prognostic prediction.

2 Methodology

The analysis encompasses several key phases:

- **Part 1: Identifying the Problem and Data Sources** - This phase involves understanding the types of information contained in the dataset utilizing R to familiarize with the data.
- **Part 2: Exploratory Data Analysis** - This phase involves data exploration and visualization. The aim is to thoroughly explore and visualize the data to understand the relationships between the variables and the response variable.
- **Part 3: Predictive Model Using SVM** - Developing a predictive model using the Support Vector Machine (SVM) algorithm to classify the breast tumor diagnosis, followed by model evaluation using tools like confusion matrix and ROC curves.
- **Part 4: Optimizing the Support Vector Classifier** - Tuning the parameters of the SVM classifier to improve its performance.

This study aims to leverage the capabilities of machine learning in predictive analytics to contribute to the early and accurate diagnosis of breast cancer, potentially leading to improved patient outcomes.

3 The Data

The Breast Cancer dataset is a multivariate dataset from the University of California, Irvine's machine learning repository. It consists of 569 instances, each with 30 features, and is used primarily for classification tasks within the health and medicine field. The features are real-valued, computed from digitized images of fine needle aspirates (FNAs) of breast masses.

The dataset's goal is to classify tumors into malignant (M) or benign (B) categories based on these features, which represent various characteristics of the cell nuclei present in the images. The first two columns of the dataset contain unique ID numbers and the diagnosis, while columns 3 to 32 contain the computed features.

These features include details about the radius, texture, perimeter, area, smoothness, compactness, concavity, number of concave points, symmetry, and fractal dimension of the cell nuclei. They are recorded under different conditions, as indicated by the numerical suffixes 1, 2, 3, etc. There are no missing values in the dataset, indicating it is well-maintained and suitable for analysis without the need for preliminary data cleaning for missing data.

The dataset is associated with a diagnostic challenge for Wisconsin Breast Cancer and can be utilized to develop predictive models. The features were selected through an exhaustive search and the separating planes used in the original analysis were derived using the Multisurface Method-Tree (MSM-T), a decision tree construction method through linear programming.

Table 1: **Breast Cancer Dataset Variables Description**

Variable Name	Role	Type	Description
ID	ID	Categorical	Unique identifier
Diagnosis	Target	Categorical	M = malignant, B = benign
radius1	Feature	Continuous	Mean of distances from center to points on the perimeter
texture1	Feature	Continuous	Standard deviation of gray-scale values
perimeter1	Feature	Continuous	Perimeter size
area1	Feature	Continuous	Area of the tumor
smoothness1	Feature	Continuous	Local variation in radius lengths
compactness1	Feature	Continuous	Perimeter squared over area minus 1.0
concavity1	Feature	Continuous	Severity of concave portions of the contour
concave_points1	Feature	Continuous	Number of concave portions of the contour
symmetry1	Feature	Continuous	Symmetry of the tumor
fractal_dimension1	Feature	Continuous	"Coastline approximation" - 1
radius2	Feature	Continuous	Mean of distances from center to points on the perimeter (SE)
texture2	Feature	Continuous	Standard deviation of gray-scale values (SE)
perimeter2	Feature	Continuous	Perimeter size (SE)
area2	Feature	Continuous	Area of the tumor (SE)
smoothness2	Feature	Continuous	Local variation in radius lengths (SE)
compactness2	Feature	Continuous	Perimeter squared over area minus 1.0 (SE)
concavity2	Feature	Continuous	Severity of concave portions of the contour (SE)
concave_points2	Feature	Continuous	Number of concave portions of the contour (SE)
symmetry2	Feature	Continuous	Symmetry of the tumor (SE)
fractal_dimension2	Feature	Continuous	"Coastline approximation" - 1 (SE)
radius3	Feature	Continuous	"Worst" or largest mean value for mean of distances from center to points on the perimeter
texture3	Feature	Continuous	"Worst" or largest mean value for standard deviation of gray-scale values
perimeter3	Feature	Continuous	"Worst" or largest mean value for perimeter size
area3	Feature	Continuous	"Worst" or largest mean value for area of the tumor
smoothness3	Feature	Continuous	"Worst" or largest mean value for local variation in radius lengths
compactness3	Feature	Continuous	"Worst" or largest mean value for perimeter squared over area minus 1.0
concavity3	Feature	Continuous	"Worst" or largest mean value for severity of concave portions of the contour
concave_points3	Feature	Continuous	"Worst" or largest mean value for number of concave portions of the contour
symmetry3	Feature	Continuous	"Worst" or largest mean value for symmetry of the tumor
fractal_dimension3	Feature	Continuous	"Worst" or largest mean value for "Coastline approximation" - 1

4 Exploratory Data Analysis

Exploratory Data Analysis (EDA) is an essential precursor to modeling, providing critical insights into data characteristics without assumptions. It helps understand data structure, distribution, outliers, and relationships. The aim is to reveal data trends, assess quality, and hypothesize using summary statistics and visualization. Descriptive statistics quantify data features, while visualization uncovers patterns, both aiding in preprocessing and guiding subsequent analysis stages.

The summary statistics of the dataset reveal key insights across multiple variables. The minimum and maximum values highlight the range of each feature, with 'Radius Mean' varying from 6.981 to 28.110, indicating diverse tumor sizes. The first and third quartiles, along with the median, outline the central tendency and dispersion, reflecting that the majority of the data points lie within these bounds. The mean values closely follow the medians, suggesting a relatively symmetrical distribution for most features. However, the significant difference between the mean and median in certain variables like 'Area Worst' might indicate skewness, warranting further investigation. These statistics are essential for understanding the data's underlying structure, guiding subsequent data preprocessing and feature selection for model building.

The correlation plots for "mean", "se", and "worst" variables in a breast cancer dataset reveal key insights. High positive correlations exist among measurements related to tumor size (radius, perimeter, area) across all variable types, indicating that these features tend to increase together. Shape-related features (compactness, concavity, concave points) are also strongly interrelated, particularly in their worst measurements, which may reflect more severe tumor characteristics.

Texture measurements show variable degrees of correlation with other features, suggesting a more complex relationship. The correlation strength is indicated by a color scale from -1 to 1, with higher values showing stronger positive relationships. These patterns are crucial for understanding feature interdependencies and can inform the feature selection process for predictive modeling, as overlapping information might be pruned to simplify models.

Diagnosis & Summary	Statistic	Radius Mean	Texture Mean	Perimeter Mean	Area Mean
Benign (n=357)	Min.	6.981	9.71	43.79	143.5
	1st Qu.	11.700	16.17	75.17	420.3
Malignant (n=212)	Median	13.370	18.84	86.24	551.1
	Mean	14.127	19.29	91.97	654.9
	3rd Qu.	15.780	21.80	104.10	782.7
	Max.	28.110	39.28	188.50	2501.0

Table 2: Summary Statistics for Selected Variables

Statistic	Smoothness Mean	Compactness Mean	Concavity Mean	Concave Points Mean	Symmetry Mean
Min.	0.05263	0.01938	0.00000	0.00000	0.1060
1st Qu.	0.08637	0.06492	0.02956	0.02031	0.1619
Median	0.09587	0.09263	0.06154	0.03350	0.1792
Mean	0.09636	0.10434	0.08880	0.04892	0.1812
3rd Qu.	0.10530	0.13040	0.13070	0.07400	0.1957
Max.	0.16340	0.34540	0.42680	0.20120	0.3040

Table 3: Summary Statistics for Selected Variables

Statistic	Fractal Dimension Mean	Radius SE	Texture SE	Perimeter SE
Min.	0.04996	0.1115	0.3602	0.757
1st Qu.	0.05770	0.2324	0.8339	1.606
Median	0.06154	0.3242	1.1080	2.287
Mean	0.06280	0.4052	1.2169	2.866
3rd Qu.	0.06612	0.4789	1.4740	3.357
Max.	0.09744	2.8730	4.8850	21.980

Table 4: Summary Statistics for Variables: Fractal Dimension Mean to Radius SE

Statistic	Area SE	Smoothness SE	Compactness SE	Concavity SE
Min.	6.802	0.001713	0.002252	0.00000
1st Qu.	17.850	0.005169	0.013080	0.01509
Median	24.530	0.006380	0.020450	0.02589
Mean	40.337	0.007041	0.025478	0.03189
3rd Qu.	45.190	0.008146	0.032450	0.04205
Max.	542.200	0.031130	0.135400	0.39600

Table 5: Summary Statistics for Variables: Area SE to Perimeter Worst

Statistic	Concave Points SE	Symmetry SE	Fractal Dimension SE	Radius Worst
Min.	0.000000	0.007882	0.0008948	7.93
1st Qu.	0.007638	0.015160	0.0022480	13.01
Median	0.010930	0.018730	0.0031870	14.97
Mean	0.011796	0.020542	0.0037949	16.27
3rd Qu.	0.014710	0.023480	0.0045580	18.79
Max.	0.052790	0.078950	0.0298400	36.04

Table 6: Summary Statistics for Variables: Concave Points SE to Perimeter Worst

Statistic	Texture Worst	Perimeter Worst	Area Worst
Min.	12.02	50.41	185.2
1st Qu.	21.08	84.11	515.3
Median	25.41	97.66	686.5
Mean	25.68	107.26	880.6
3rd Qu.	29.72	125.40	1084.0
Max.	49.54	251.20	4254.0

Table 7: Summary Statistics for Variables: Texture Worst to Perimeter Worst

Statistic	Smoothness Worst	Compactness Worst	Concavity Worst
Min.	0.07117	0.02729	0.0000
1st Qu.	0.11660	0.14720	0.1145
Median	0.13130	0.21190	0.2267
Mean	0.13237	0.25427	0.2722
3rd Qu.	0.14600	0.33910	0.3829
Max.	0.22260	1.05800	1.2520

Table 8: Summary Statistics for Selected Variables

Statistic	Concave Points Worst	Symmetry Worst	Fractal Dimension Worst
Min.	0.00000	0.1565	0.05504
1st Qu.	0.06493	0.2504	0.07146
Median	0.09993	0.2822	0.08004
Mean	0.11461	0.2901	0.08395
3rd Qu.	0.16140	0.3179	0.09208
Max.	0.29100	0.6638	0.20750

Table 9: Summary Statistics for Selected Variables

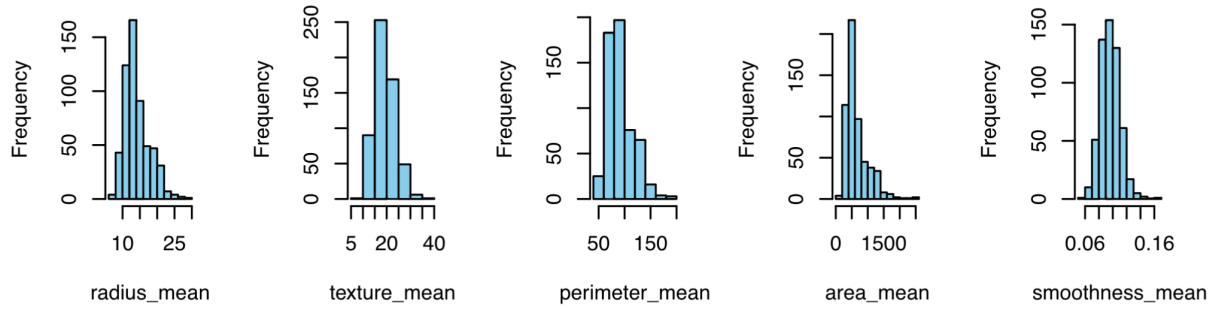


Figure 1: Plot of Histograms of mean variables

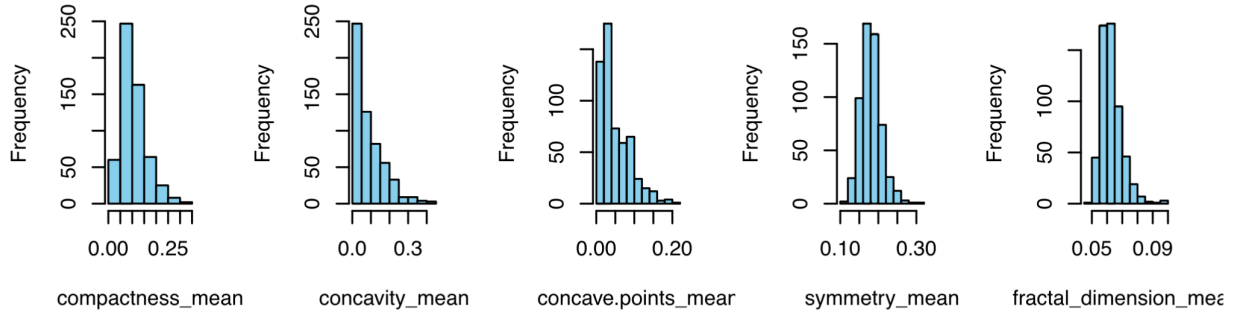


Figure 2: Plot of Histograms of mean variables

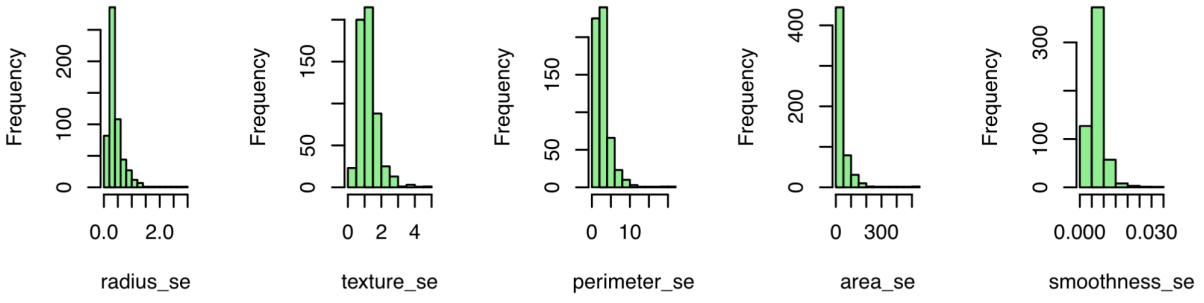


Figure 3: Plot of Histograms of se variables

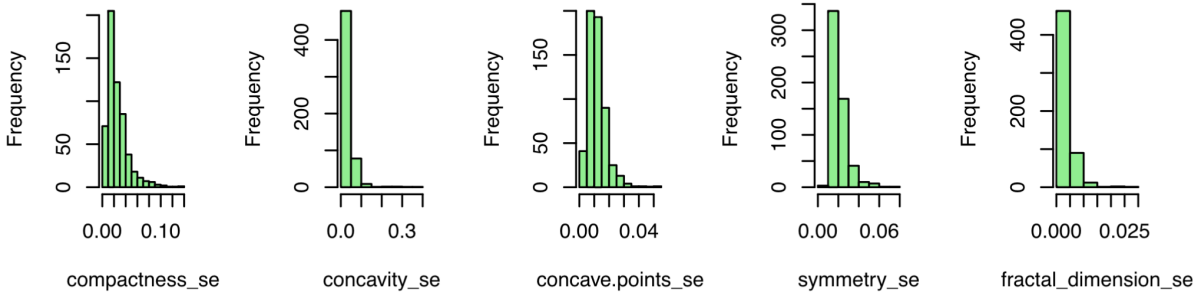


Figure 4: Plot of Histograms of se variables

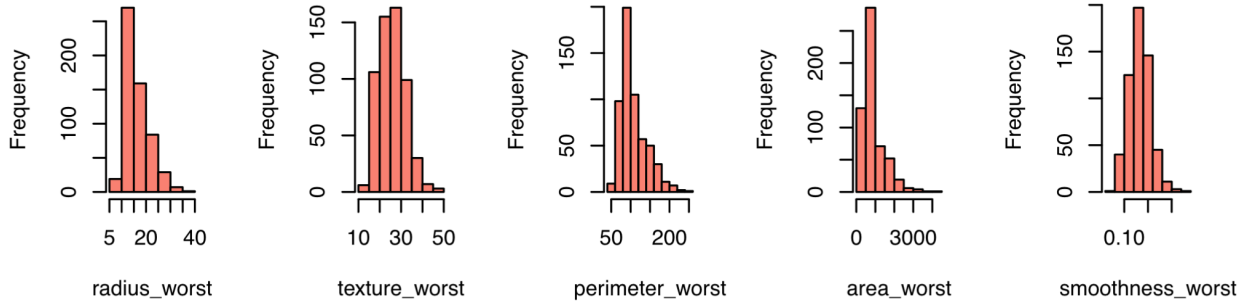


Figure 5: Plot of Histograms of worst variables

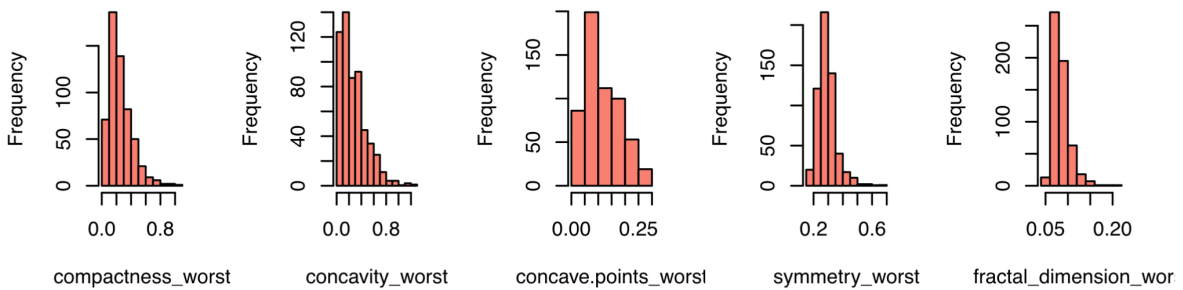
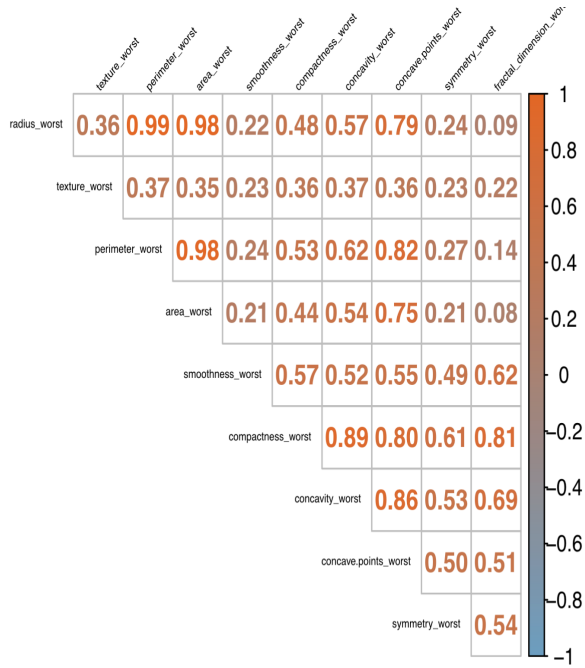
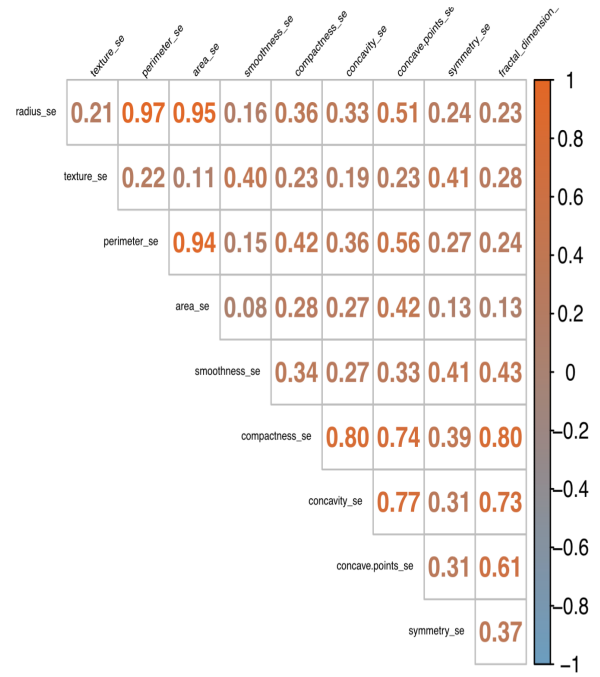


Figure 6: Plot of Histograms of worst variables



(a) Worst Variables Correlation Matrix



(b) Se Variables Correlation Matrix

Figure 7: Correlation matrix

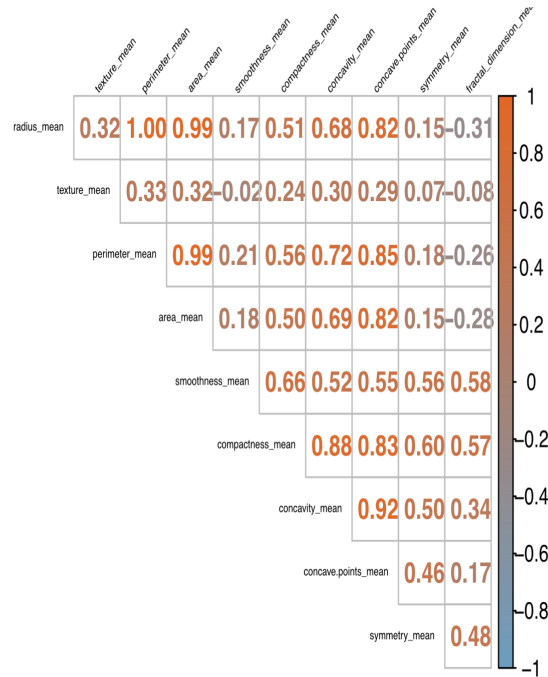


Figure 8: Mean Variables Correlation Matrix

5 Support Vector Machine (SVM)

In this segment of our study, we will develop a predictive model utilizing the Support Vector Machine (SVM) learning algorithm. Recognized as a prominent classification technique, SVMs notably transform nonlinear data, enabling the application of a linear algorithm for linear model fitting, as outlined by Cortes and Vapnik in 1995. SVMs, particularly when kernelized, are robust and versatile, demonstrating effective performance across diverse datasets. They are adept at creating intricate decision boundaries, even with datasets characterized by a limited number of features. SVMs exhibit proficiency in handling both low-dimensional and high-dimensional data. However, their scalability is less efficient with an increasing number of samples, making them more suitable for datasets with up to 10,000 samples. For larger datasets, SVMs may encounter challenges related to runtime and memory usage. One critical aspect of employing SVMs is the necessity for meticulous data preprocessing and parameter tuning. It's important to note that SVM models are not particularly transparent, making it somewhat challenging to interpret specific predictions or to explain the model's workings to someone without a background in the field. Despite these considerations, the strength of SVMs in handling various types of data makes them a valuable tool in our predictive modeling arsenal.

5.1 Step 1: Selecting Important Parameters

The efficacy of kernel SVMs is significantly influenced by a few key parameters. These include:

- **Regularization parameter (C):** The Regularization parameter (C) in SVM plays a critical role in managing the trade-off between minimizing error on the training data and reducing model complexity to enhance generalization. Essentially, a higher value of C targets lower training error, potentially increasing the risk of overfitting. Conversely, a lower C value favors a simpler, more generalized model. The “cost” option specifies the cost of violating the margin. We tried costs 0.001, 0.01, 0.1, 1, 10, 100, 1000. In the specific context of our study, the optimal cost identified is 0.1. This suggests a preference for a simpler decision boundary, which is likely to generalize more effectively to new, unseen data. This value of C reflects a strategic choice to prioritize model generalizability over fitting to every nuance in the training dataset.
- **Choice of Kernel:** The choice of the kernel in SVMs is crucial as it dictates how the data is transformed for classification. Common kernel options include linear, radial basis function (RBF), and polynomial, each significantly impacting the decision boundary's shape. For the dataset, the linear kernel emerged as the optimal choice, indicating that the classes in the feature space might be effectively separable by linear boundaries. The linear kernel works by identifying the best straight line (or hyperplane in higher dimensions) to differentiate be-

tween classes. This selection suggests that the data’s underlying structure can be adequately captured with linear relationships, favoring simplicity and interpretability in the model.

- **Kernel-specific Parameters:** For different kernels, there are specific parameters to tune. For instance, the gamma parameter in the RBF kernel determines the influence radius of a single training example.
- **Gamma:** In SVM models, the parameters gamma and C are critical in determining model complexity. Large values for either result in more intricate models. For optimal performance, these parameters should be adjusted together due to their interdependence. Gamma, specifically, is crucial for non-linear hyperplanes in the kernel function. A lower gamma value indicates a larger similarity radius, leading to a more generalized decision boundary. In our model, with gamma set to 1, this suggests a broader influence of individual training samples on the decision boundary, pointing towards a balance between model complexity and generalization.

Kernel	Cost	Gamma	Error Rate	Dispersion
Linear	0.1	1	0.0276	0.0322
Radial	10	1	0.3711	0.1098
Polynomial	0.001	1	0.0425	0.0206

Table 10: Best Parameters for Different SVM Kernels

Considering the error rates, the linear kernel seems to perform the best on the dataset, showcasing the lowest error rate among the tested kernels. The selected parameters for the SVM model, focusing on low cost and gamma with a linear kernel, demonstrate a preference for simplicity and generalizability, aiming to create a less complex model that effectively generalizes to new data.

5.2 Step 2: Fitting the SVM Model with Key Parameters

The model identified 46 support vectors from the training data. These are the crucial data points that define the decision boundary and contribute to making classifications. The numbers (23 23) possibly represent the count of support vectors for each class ('B' and 'M'). Thus, for the 46 points violating the separating hyperplane or the margin, 23 are in each class.

In the process of fitting the SVM model with the chosen key parameters, the initial evaluation using a confusion matrix reveals promising results. The model correctly classified 62.4% (249 cases) as true positives and 36.1% (144 cases) as true negatives, indicating strong predictive accuracy for both classes. Only a minimal fraction of 0.3% (1 case) was misclassified as false negatives, and 1.3% (5 cases) as false positives, showcasing the model’s ability to discern between classes effectively with few errors. The recall rates, represented by 99.6% for true positives and 96.6%

Parameter	Details
SVM-Type	C-classification
SVM-Kernel	linear
Cost	0.1
Number of Support Vectors	46 (23 23)
Levels	"B" "M"

Table 11: Summary of SVM Model Fitting with Key Parameters

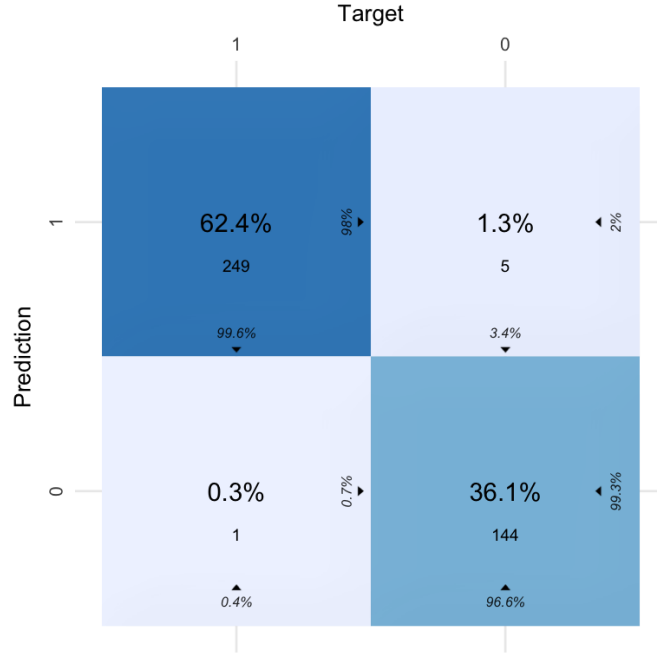


Figure 9: Performance of Train dataset

for true negatives, emphasize the model’s reliability in identifying correct class labels. However, these figures are based on the training data, and to mitigate the risk of overfitting, a thorough validation on an independent test set is imperative to confirm the model’s ability to generalize to new, unseen data.

5.3 Step 3: Model Performance

In the third step of our analysis, the SVM model’s performance was rigorously evaluated using a test dataset. The confusion matrix showed a high percentage of correct classifications: 61.8% true negatives and 35.3% true positives, indicating a proficient identification of both benign and malignant instances. The model achieved an impressive 97.06% accuracy, suggesting that it can correctly classify the data with high reliability.

Sensitivity and specificity rates stood at 95.24% and 98.13%, respectively, which shows the model’s strong capability in correctly detecting malignant cases and its precision in identifying benign cases. The precision of the model, measured by the positive predictive value, was 96.77%, and

the negative predictive value was 97.22%, further underscoring the model's consistent performance. An excellent Kappa statistic of 0.9367 pointed to a substantial agreement beyond chance.

These test data-based indicators are critical for assessing the model's generalizability to new data. While the model's numerical performance is robust, it's crucial to contextualize these results within the clinical setting, particularly considering the implications of any misclassifications. The model's high accuracy and discriminative power suggest it could be a valuable tool for medical diagnostics, provided it is used in conjunction with expert clinical judgment.

Table 12: Model Evaluation Indicators

Metric	Value
Accuracy	0.9706
95% CI	(0.9327, 0.9904)
No Information Rate	0.6294
p-Value [Acc > NIR]	< 0.0001
Kappa	0.9367
Sensitivity	0.9524
Specificity	0.9813
Pos Pred Value	0.9677
Neg Pred Value	0.9722
Prevalence	0.3706
Detection Rate	0.3529
Detection Prevalence	0.3647
Balanced Accuracy	0.9668
Area under the curve	0.9668
Root Mean Squared Error (RMSE)	0.15056

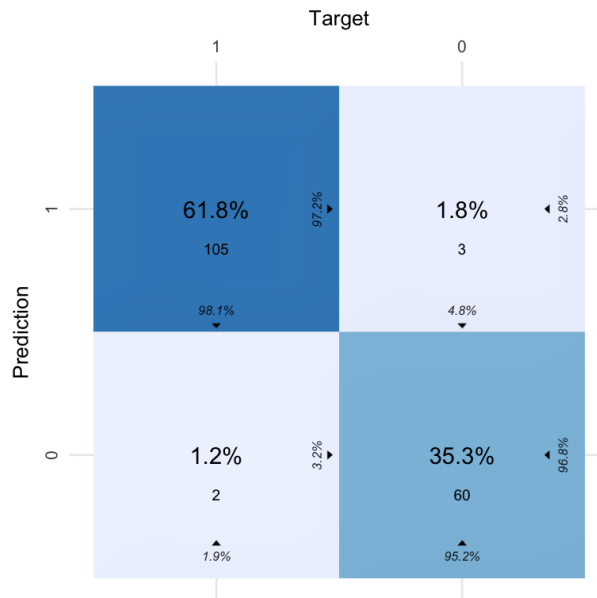


Figure 10: Confusion Matrix for Test Data

5.4 Step 4: Model Accuracy - Receiver Operating Characteristic (ROC) curve

In statistical modeling and machine learning, a commonly-reported performance measure of model accuracy for binary classification problems is the Area Under the Curve (AUC). In an ROC curve, the "True Positive Rate" (TPR) is plotted on the Y-axis, and the "False Positive Rate" (FPR) is plotted on the X-axis. TPR represents the model's ability to correctly identify positive instances among all actual positives, while FPR indicates the ratio of incorrect identifications of negative instances among all actual negatives.

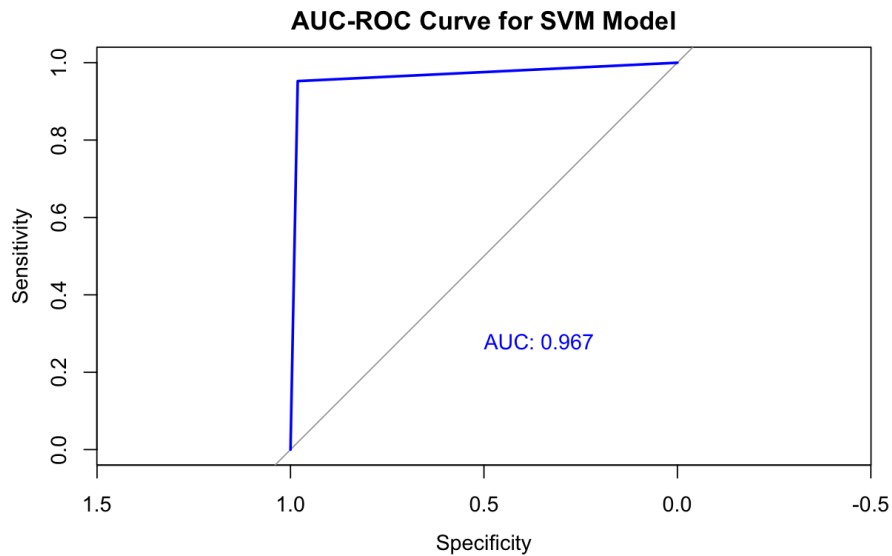


Figure 11: ROC Plot for Test Data

The ROC curve illustrates the trade-off between TPR and FPR at various decision thresholds set by the model. The AUC, the area under this curve, quantifies the model's overall performance. A higher AUC signifies better discriminative ability and overall model performance. The AUC for the ROC curve was 96.68%, demonstrating the model's exceptional ability to distinguish between the classes. A high AUC is indicative of the model's effectiveness in differentiating between benign and malignant tumors.

5.5 Step 5: Variable Importance

The reported information pertains to the variable importance scores calculated using ROC curve analysis. The 20 most important variables (out of a total of 30) are presented based on their importance scores. Table 13 exhibits the variables and their respective importance scores. Each variable is ranked based on its discriminative power for the diagnosis prediction task, as quantified by their importance values.

Variable	Importance Score
perimeter_worst	0.9755
concave.points_worst	0.9708
radius_worst	0.9703
area_worst	0.9689
concave.points_mean	0.9630
perimeter_mean	0.9465
concavity_mean	0.9453
area_mean	0.9360
radius_mean	0.9358
concavity_worst	0.9337
area_se	0.9229
compactness_worst	0.8820
compactness_mean	0.8793
perimeter_se	0.8699
radius_se	0.8666
concave.points_se	0.7987
concavity_se	0.7966
texture_worst	0.7898
texture_mean	0.7867
smoothness_worst	0.7568

Table 13: Variable Importance Scores for Breast Cancer Diagnosis

The higher the importance score, the greater the variable’s contribution to the model’s discriminatory ability. Variables like *perimeter_worst*, *concave.points_worst*, *radius_worst*, and *area_worst* demonstrate notably high importance scores, indicating their significant role in differentiating between tumor classes.

5.6 Step 6: Decision Boundaries

The classification was finalized using a linear kernel, resulting in a linear boundary. Support vectors are denoted as "x", while other points are represented as "o". This linear boundary visualizes the decision space created by the SVM model to separate different classes. The line dividing the two areas represents the decision boundary determined by the SVM. This is the line where the model calculates a 50% chance of being either class. The two different colors or shades in the plot show the areas where new data points would be classified as one class or the other. Data points falling in the red region would be classified as M (malignant), while those in the yellow region would be classified as B (benign).

The SVM classification plot in figure 14 provides a visual representation of the model’s decision boundary and support vectors in a two-dimensional feature space. Data points are color-coded to indicate their respective classifications, with stars denoting the support vectors that are critical in defining the boundary. In Figure 14, the plot suggests a clear separation between classes, reflecting

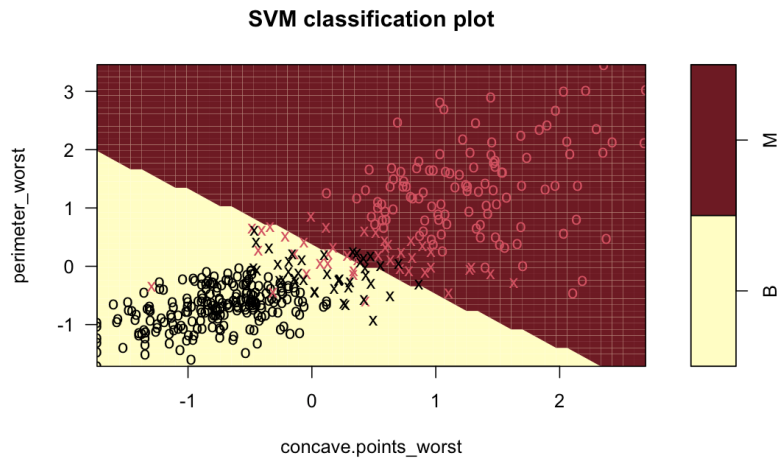


Figure 12: Decision boundaries for Perimeter Worst and Concave points worst

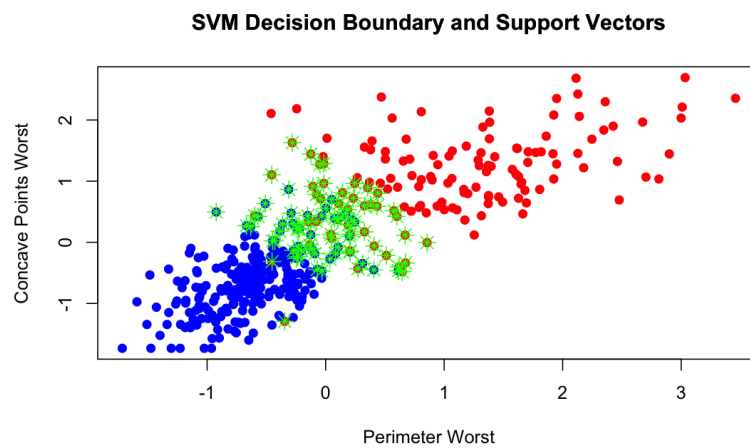


Figure 13: SVM Decision Boundary and Support Vectors

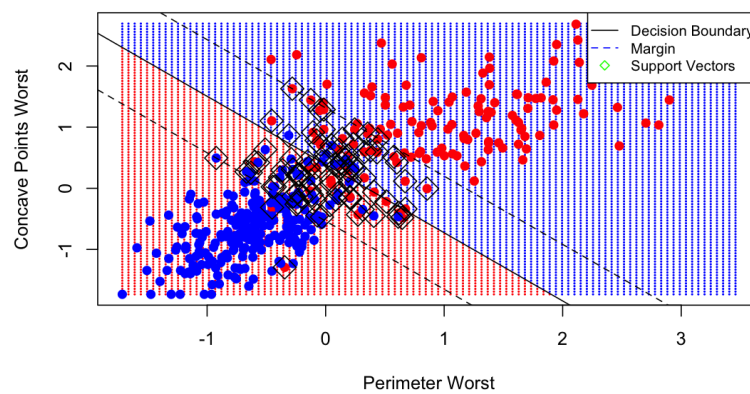


Figure 14: SVM Decision Boundary and Margin

the SVM’s ability to differentiate between categories such as benign and malignant. It should be noted, however, that the actual SVM operates in a multidimensional space, and this plot simplifies its complex decision-making to two dimensions for interpretability.

Conclusion

This study meticulously examined the application of Support Vector Machines (SVM) for classifying breast cancer tumors using fine-needle aspiration (FNA) test results. Our research questions were centered on the SVM’s classification effectiveness and the identification of significant features within FNA test results. Our methodical approach included comprehensive data analysis and predictive modeling with SVM, where the model demonstrated a high accuracy rate of 97.06% on the test dataset. This precision is indicative of the model’s robustness, further substantiated by a Kappa statistic of 0.9367, reflecting substantial agreement beyond chance.

The model’s sensitivity and specificity rates stood impressively at 95.24% and 98.13%, respectively. These rates affirm the model’s adeptness in correctly detecting malignant cases and its precision in identifying benign cases. Moreover, the high precision of the model, indicated by a positive predictive value of 96.77%, and the negative predictive value of 97.22%, underscore the model’s consistent performance.

The ROC curve analysis culminated in an AUC score of 96.68%, showcasing the model’s exceptional capability to distinguish between benign and malignant tumor classes. The variable importance analysis yielded significant predictors for tumor classification, with variables such as *perimeter_worst* and *concave.points_worst* exhibiting particularly high importance scores. The decision boundary visualizations elucidated the SVM’s strategy for class differentiation, with support vectors markedly influencing the linear boundary. These vectors were instrumental in defining the decision space, ensuring accurate class segregation.

In conclusion, our SVM model’s outcomes, supported by quantitative evidence, confirm its high efficacy in classifying breast cancer tumors. The model’s ability to generalize was affirmed by its exceptional performance metrics, suggesting its potential as a valuable asset in the realm of medical diagnostics. Further validation is advocated to solidify the model’s applicability in clinical settings.

6 Limitations of Study

While our study aims to provide insights into predictive modeling using machine learning techniques, it is essential to acknowledge several limitations that could impact the interpretation and generalization of our findings. These limitations encompass various aspects, including the dataset used, the modeling approach employed, and the scope of analysis.

1. **Dataset Size:** The study's performance might be influenced by the size of the dataset. A larger dataset often provides better generalization, and your findings could change with more extensive data.
2. **Absence of Feature Selection:** The study lacks the implementation of feature selection techniques, a critical aspect in refining predictive models. This omission may result in a model that includes redundant or less impactful features, potentially compromising predictive accuracy, interpretability, and computational efficiency. The absence of feature selection undermines the opportunity to identify the most informative variables, raising concerns about the model's robustness and its potential performance on unseen data.
3. **Model Evaluation:** The study could benefit from utilizing various evaluation metrics beyond error rates, such as precision, recall, or F1-score, for a more comprehensive assessment.
4. **Model Parameters:** While the chosen parameters show promise, further optimization or exploring a wider range of hyperparameters might yield even better-performing models.
5. **Sensitivity to Scaling:** Certain machine learning models, including SVMs, can be sensitive to data scaling. Ensure that the scaling technique used doesn't introduce biases or affect model performance disproportionately.
6. **Algorithm Selection:** Limiting the study to SVMs might overlook the potential advantages of other algorithms or ensemble methods that could perform better on the given data.

Recognizing these constraints/limitations is crucial in understanding the boundaries within which our study operates and in highlighting areas for future research and improvement.