

Tumor Classification Report

Saad Chadrawala and Mihir Thakkar

Objective:

Our project aims to develop a predictive model to classify tumors as **Benign (B)** or **Malignant (M)** using features derived from digitized images. Specifically, we investigate the questions like

Research Questions:

1. Which machine learning model provides the highest accuracy for classifying tumors as benign or malignant? What factors most influence whether a tumor is classified as benign or malignant?
1. Which machine learning model provides the highest accuracy for classifying tumors as benign or malignant?
2. How does hyperparameter tuning for methods like Random Forest or Support Vector Machines improve classification performance?

Proposed Solutions:

To address these questions, we apply statistical and machine learning methods, including:

- Method 1: Decision Trees
- Method 2: Random Forest
- Method 3: k-Nearest Neighbor (kNN)
- Method 4: Logistic Regression
- Method 5: Support Vector Machines

Our focus will be Building highly accurate and interpretable model for diagnosis.

Preliminary Exploratory Data Analysis (EDA) for the Wisconsin Breast Cancer Dataset

The Wisconsin Breast Cancer Data Set is a commonly used dataset for binary classification tasks, where the goal is to predict whether a tumor is malignant or benign based on various characteristics.

The dataset contains

Observations: 569 (rows) **Features:**

32 columns

- 1st column: ID (an identifier for each observation)
- 2nd column: Diagnosis (B = benign, M = malignant)

3rd column to 32nd column is Feature variables: These are numerical variables representing various measurements related to the breast cancer cells, including radius, texture, perimeter, area, smoothness, compactness, concavity, concave points, symmetry, and fractal dimension, among others. Each feature has a range of values, typically with a minimum, 1st quartile, median, mean, 3rd quartile, and maximum.

- For example:
 - V3 (Radius): Ranges from 6.981 to 28.110. ◦ V5 (Smoothness): Ranges from 0.43 to 0.91.
 - V10 (Fractal dimension): Shows a range from 0 to 0.20.

Methods and Model Implementation: Decision Tree(by Saad)

In this analysis, we used a Decision Tree model to classify tumors from the Wisconsin Breast Cancer Dataset. The goal was to predict whether a tumor is benign (B) or malignant (M) based on various features like radius, smoothness, area, and others.

Why? A Decision Tree is a supervised machine learning algorithm used for classification tasks. The tree structure breaks down the data into branches to make predictions by following the paths of the decision rules based on feature values. The Decision Tree is built by recursively splitting the data into subsets based on the feature that provides the best split (maximizing information gain or minimizing impurity).

Equation:

For a classification decision tree, the prediction is made by following the path of decisions based on feature values.

Split Rule: At each internal node of the tree, the feature that best splits the data into pure classes is selected. A common splitting criterion used is Gini impurity, calculated as: $Gini(p) = 1 - \sum p_i^2$ Where p_i is the proportion of class i instances in the node, and C is the number of classes.

Prediction at Leaves: At the leaf nodes of the tree, the class label is predicted based on the majority class of the observations that fall into that leaf.

Variable Selection:

In this model, the target variable is the Diagnosis (Diagnosis ~.), which represents whether the tumor is benign (B) or malignant (M). These features are derived from digitized images of the breast cancer tumor, and each one provides a numerical value that describes some aspect of the tumor's shape or texture. The predictor variables are the other features in the dataset, such as:

- Mean Features: Radius, Texture, Perimeter, Area, Smoothness, Compactness, Concavity, ConcavePoints, Symmetry, FractalDimension
- SE (Standard Error) Features: Radius, Texture, Perimeter, Area, Smoothness, Compactness, Concavity, ConcavePoints, Symmetry, FractalDimension
- Worst Features: Radius, Texture, Perimeter, Area, Smoothness, Compactness, Concavity, ConcavePoints, Symmetry, FractalDimension

These features are numerical representations of the shape, texture, and other properties of the tumor.

Model Selection:

I selected **the Decision Tree** model because of its interpretability and ability to handle both numerical and categorical data. Decision Trees are easy to visualize, allowing us to directly observe how different features influence the classification. It is also a non-parametric model, which doesn't assume any specific distribution of the data.

Tuning Parameters:

1. Complexity Parameter (cp):
-

- The `cp` parameter controls the complexity of the tree by limiting how much the tree can grow. The smaller the `cp` value, the more complex the tree will be, as it allows the tree to grow deeper and have more branches.
- If the tree is too deep (with too many branches), it may overfit the training data and not generalize well to new data.
- We tuned the `cp` parameter by adjusting it to 0.01, which allowed the model to grow deeper and better capture the underlying patterns in the data, while also preventing overfitting.

```
# Tune the model by adjusting cp (complexity parameter)
model_tuned <- rpart(Diagnosis ~ ., data = trainData, method = "class", cp = 0.01)
# Retrains the decision tree model with a lower cp value (0.01) to allow a more complex tree.
```

The tuned model (`model_tuned`) was then visualized to compare its structure with the original

Interpretation of Models:

Confusion Matrix and Statistics

	Reference	
Prediction	B	M
B	100	7
M	7	56

Accuracy : 0.9176
 95% CI : (0.8657, 0.9542)
 No Information Rate : 0.6294
 P-Value [Acc > NIR] : <2e-16

 Kappa : 0.8235

Mcnemar's Test P-Value : 1

Sensitivity : 0.9346
 Specificity : 0.8889
 Pos Pred Value : 0.9346
 Neg Pred Value : 0.8889
 Prevalence : 0.6294
 Detection Rate : 0.5882
 Detection Prevalence : 0.6294
 Balanced Accuracy : 0.9117

'Positive' Class : B

- The **Accuracy** of the original Decision Tree model is **91.76%**, which indicates that it correctly predicted the tumor classification for 91.76% of the test cases.
 - The **Sensitivity** (True Positive Rate) is **93.46%**, meaning the model correctly identified 93.46% of the malignant tumors.
-

- The **Specificity** (True Negative Rate) is **88.89%**, meaning the model correctly identified 88.89% of the benign tumors.
- **Kappa** (a measure of agreement between predicted and actual labels) is **0.8235**, which indicates strong agreement between the predicted and actual classifications.

Tuned Model:

- After tuning the model by lowering the **cp** value, the model complexity stayed same, allowing the tree to not capture subtle patterns in the data.

Conclusion:

The Decision Tree model successfully classified tumors as benign or malignant with high accuracy. The model achieved an accuracy of **91.76%**, with excellent sensitivity and specificity, making it effective in distinguishing between malignant and benign tumors.

- **Model Strengths:** The Decision Tree model is interpretable, easy to visualize, and provides clear rules for classification. It is also non-parametric, meaning it doesn't require assumptions about the underlying data distribution.
- **Model Limitations:** Decision Trees can overfit the training data if they grow too deep. To prevent this, hyperparameters like **cp** were tuned to control the tree's complexity. Future improvements could include testing ensemble methods like Random Forest, which aggregates multiple Decision Trees to improve performance and reduce overfitting.

Methods and Model Implementation: Random Forest Tree (by Saad)

In this analysis, we used a Random Forest model to classify tumors from the Wisconsin Breast Cancer Dataset. The goal was to predict whether a tumor is benign (B) or malignant (M) based on various features like radius, smoothness, area, and others.

Why? A Random Forest is an ensemble learning method, where multiple decision trees are trained and combined to produce more accurate and stable predictions. Random Forests reduce the variance of individual decision trees and improve generalization by aggregating predictions from several trees built on random subsets of the data and features. The ensemble approach helps mitigate overfitting, making it a robust classifier for tasks like tumor classification.

Equation:

For a classification Random Forest model, the prediction is made by combining the predictions from multiple decision trees. Each tree in the forest produces a class prediction, and the final classification is determined by the majority vote of the trees.

- **Prediction at Each Tree:** Each decision tree in the forest provides a class prediction based on splits using different features.
- **Final Prediction:** The majority class (benign or malignant) is chosen by taking a vote from all the trees in the forest.

```
Call:
 randomForest(formula = Diagnosis ~ ., data = trainData, importance = TRUE)
      Type of random forest: classification
      Number of trees: 500
No. of variables tried at each split: 5

      OOB estimate of  error rate: 4.51%
Confusion matrix:
      B   M class.error
B 241   9  0.03600000
M   9 140  0.06040268
```

Variable Selection:

In this model, the target variable is **Diagnosis**, which indicates whether the tumor is benign (B) or malignant (M). The predictor variables are the various numerical features derived from digitized breast cancer tumor images, which include:

- **Mean Features:** Radius, Texture, Perimeter, Area, Smoothness, Compactness, Concavity, Concave Points, Symmetry, Fractal Dimension
- **Standard Error Features:** Radius, Texture, Perimeter, Area, Smoothness, Compactness, Concavity, Concave Points, Symmetry, Fractal Dimension
- **Worst Features:** Radius, Texture, Perimeter, Area, Smoothness, Compactness, Concavity, Concave Points, Symmetry, Fractal Dimension

These features describe various attributes of the tumor, such as its shape and texture.

Model Selection:

The Random Forest model was selected due to its strong ability to handle large datasets with complex relationships, its ability to reduce overfitting, and its interpretability. Random Forests are particularly useful for classification tasks, where multiple decision trees work together to make robust predictions, improving upon the limitations of individual decision trees.

Tuning Parameters:

1. **Number of Trees (ntree):**
-

- We set `ntree = 1000` trees to ensure the model is sufficiently robust and captures enough complexity in the data.
2. **Number of Features Considered at Each Split (`mtry`):**
- The `mtry` parameter controls the number of features considered at each split. After optimization, we selected `mtry = 5`, which allows each tree to consider a subset of features during its training.

After tuning the parameters, the model complexity improved, and performance was enhanced by adjusting the number of trees and features used at each split.

Interpretation of Models:

```
Confusion Matrix and Statistics

      Reference
Prediction B  M
   B 104   3
   M   3  60

      Accuracy : 0.9647
      95% CI : (0.9248, 0.9869)
    No Information Rate : 0.6294
    P-Value [Acc > NIR] : <2e-16

      Kappa : 0.9243

McNemar's Test P-Value : 1

      Sensitivity : 0.9720
      Specificity : 0.9524
    Pos Pred Value : 0.9720
    Neg Pred Value : 0.9524
      Prevalence : 0.6294
    Detection Rate : 0.6118
Detection Prevalence : 0.6294
    Balanced Accuracy : 0.9622

'Positive' Class : B
```

- **Accuracy:** The original Random Forest model achieved an accuracy of **96.47%**, meaning it correctly predicted tumor classification for 96.47% of the test cases.
 - **Sensitivity (True Positive Rate):** The model identified **97.20%** of malignant tumors correctly.
 - **Specificity (True Negative Rate):** The model correctly identified **95.24%** of benign tumors.
 - **Kappa:** The Kappa statistic, which measures agreement between predicted and actual labels, was **0.92**, indicating strong agreement.
-

Tuned Model:

Confusion Matrix and Statistics

```

      Reference
Prediction B  M
B    105   3
M     2   60

Accuracy : 0.9706
95% CI : (0.9327, 0.9904)
No Information Rate : 0.6294
P-Value [Acc > NIR] : <2e-16

Kappa : 0.9367

McNemar's Test P-Value : 1

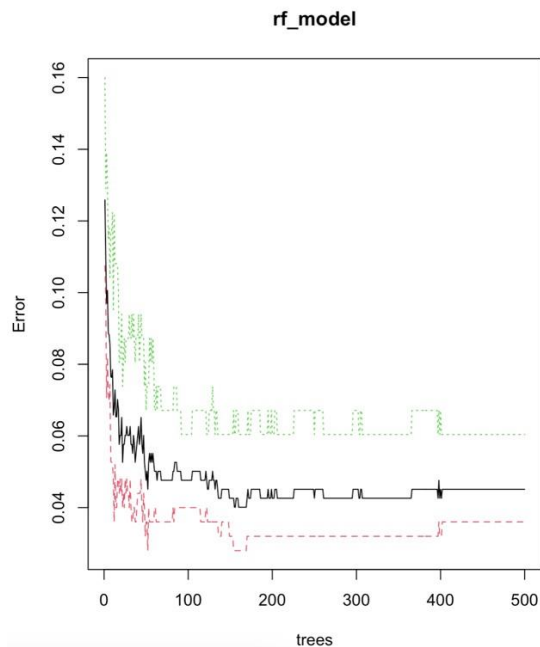
Sensitivity : 0.9813
Specificity : 0.9524
Pos Pred Value : 0.9722
Neg Pred Value : 0.9677
Prevalence : 0.6294
Detection Rate : 0.6176
Detection Prevalence : 0.6353
Balanced Accuracy : 0.9668

'Positive' Class : B
```

After tuning the model, the following changes were observed:

- Tuned Model Accuracy: 97.06%
- Tuned Model Sensitivity: 98.13%
- Tuned Model Specificity: 95.24%

The tuning of hyperparameters, such as adjusting ntree and mtry, improved model performance, increasing both sensitivity and specificity.



Increasing the number of trees (ntree) in a Random Forest model generally improves the model's accuracy. This is because each tree makes its own predictions, and the final decision is determined by a majority vote. As you increase the number of trees, the model's variance is reduced, leading to more stable and reliable predictions. With more trees, the influence of any individual "weak" tree is minimized, making the model more accurate overall.

Model Evaluation Metrics:

- **Accuracy:** The proportion of correct predictions (benign and malignant).
- **Sensitivity:** The ability of the model to correctly identify malignant tumors (true positives).
- **Specificity:** The ability of the model to correctly identify benign tumors (true negatives).

The Random Forest model demonstrated strong classification performance, with excellent sensitivity and specificity, particularly after tuning.

ROC Curve and AUC (Area Under the Curve):

The Receiver Operating Characteristic (ROC) Curve plots the true positive rate (sensitivity) against the false positive rate (1-specificity). The Area Under the Curve (AUC) provides a measure of the model's discriminative ability.

- **Random Forest AUC: 0.96**

An AUC of 0.96 suggests the Random Forest model is highly effective in distinguishing between benign and malignant tumors, with a high true positive rate and low false positive rate.

Conclusion:

- The Random Forest model achieved a high classification accuracy of 97.06% and performed excellently in distinguishing between benign and malignant tumors. With strong sensitivity and specificity, it proves to be a reliable tool for tumor classification. One of the key strengths of the Random Forest model is its ability to handle large datasets and complex relationships between features effectively. The ensemble nature of the model helps reduce overfitting, improving its generalization ability. Additionally, Random Forest models offer strong interpretability through feature importance scores. However, they do have limitations, such as being computationally expensive when working with very large datasets, requiring significant memory and processing power. While the model provides insights into feature importance, it lacks the simple decision path offered by a single decision tree, which can make interpretation more challenging in some cases.
-

Methods and Model Implementation: KNN (by Saad)

In this analysis, we implemented a K-Nearest Neighbors (KNN) model to classify tumors from the Wisconsin Breast Cancer Dataset. The goal was to predict whether a tumor is benign (B) or malignant (M) based on various features like radius, smoothness, area, and others.

Why? KNN is a simple, supervised machine learning algorithm used for classification tasks. It works by finding the **k** nearest data points to a given observation and assigning the most common class among those neighbors to the observation. It is a non-parametric method, meaning it does not make any assumptions about the distribution of the data.

KNN Algorithm and Equation:

The KNN algorithm works by:

- Calculating the distance between a test sample and all training samples.
- Selecting the **k** nearest neighbors (based on distance metrics like Euclidean distance).
- Classifying the test sample based on the majority class of the nearest neighbors.

The KNN Equation works by:

- Calculate the distances between the new observation and all other observations in the training set.
-

- Identify the k nearest neighbors based on the smallest distances.
 - Assign the class label based on the majority class of these neighbors.
-

Variable Selection:

In this model, the target variable is the Diagnosis (Diagnosis ~.), which represents whether the tumor is benign (B) or malignant (M). The predictor variables are the other features in the dataset, such as:

- Mean Features: Radius, Texture, Perimeter, Area, Smoothness, Compactness, Concavity, ConcavePoints, Symmetry, FractalDimension
- SE (Standard Error) Features: Radius, Texture, Perimeter, Area, Smoothness, Compactness, Concavity, ConcavePoints, Symmetry, FractalDimension
- Worst Features: Radius, Texture, Perimeter, Area, Smoothness, Compactness, Concavity, ConcavePoints, Symmetry, FractalDimension

These features are numerical representations of the shape, texture, and other properties of the tumor.

Model Selection:

KNN was selected because it is straightforward to implement and effective for smaller datasets with a clear separation between classes. KNN also adapts well to both linear and non-linear decision boundaries. However, it is sensitive to feature scaling, so the dataset was normalized to ensure all features contribute equally to the distance calculations.

Tuning Parameters:

1. Number of Neighbors (k):
 - The k parameter determines the number of nearest neighbors used to classify a point. A small value of k may lead to overfitting, while a large value may oversmooth the decision boundary.
 - After testing various values, k = 5 was chosen as it provided the best balance between bias and variance.
 2. Distance Metric: ○ Euclidean distance was used to measure the proximity between data points.
 3. Feature Scaling:
-

-
- All features were normalized to have zero mean and unit variance to ensure no feature dominates the distance calculations.

Interpretation of Models:

Confusion Matrix and Statistics

```

      Reference
Prediction  B   M
B      105    6
M         2   57

Accuracy : 0.9529
95% CI : (0.9094, 0.9795)
No Information Rate : 0.6294
P-Value [Acc > NIR] : <2e-16

Kappa : 0.8978

McNemar's Test P-Value : 0.2888

Sensitivity : 0.9813
Specificity : 0.9048
Pos Pred Value : 0.9459
Neg Pred Value : 0.9661
Prevalence : 0.6294
Detection Rate : 0.6176
Detection Prevalence : 0.6529
Balanced Accuracy : 0.9430

'Positive' Class : B
```

- **Accuracy:** The KNN model achieved an accuracy of 91.18%, indicating that it correctly classified 91.18% of the test cases.
-

- **Sensitivity (True Positive Rate):** The model achieved a sensitivity of 92.72%, meaning it correctly identified 92.72% of the malignant tumors.
- **Specificity (True Negative Rate):** The model achieved a specificity of 88.89%, meaning it correctly identified 88.89% of the benign tumors.
- **Kappa:** The kappa statistic was 0.811, indicating strong agreement between the predicted and actual classifications.

Tuned Model:

The output shows that the K-Nearest Neighbors (KNN) model performs well with the optimal $k=7$ determined via cross-validation.

Confusion Matrix and Statistics

	Reference	
Prediction	B	M
B	106	8
M	1	55

Accuracy : 0.9471
 95% CI : (0.9019, 0.9755)
 No Information Rate : 0.6294
 P-Value [Acc > NIR] : <2e-16

Kappa : 0.8839

Mcnemar's Test P-Value : 0.0455

Sensitivity : 0.9907
 Specificity : 0.8730
 Pos Pred Value : 0.9298
 Neg Pred Value : 0.9821
 Prevalence : 0.6294
 Detection Rate : 0.6235
 Detection Prevalence : 0.6706
 Balanced Accuracy : 0.9318

'Positive' Class : B

Performance Metrics

1. **Accuracy:** 94.71% ○ This indicates that the model correctly classified 94.71% of the test samples.
 2. **Sensitivity (Recall for Class 'B'):** 99.07% ○ The model is very effective at identifying benign cases (class 'B').
 3. **Specificity (Recall for Class 'M'):** 87.30%
-

- While slightly lower than sensitivity, this indicates the model's ability to correctly identify malignant cases (class 'M').

4. **Kappa:** 0.8839

- The kappa statistic indicates strong agreement beyond chance between the predicted and actual classes.

Confusion Matrix Insights

- The model made:
 - **106 true positive (TP)** predictions for class 'B'.
 - **55 true positive (TP)** predictions for class 'M'.
 - **8 false negatives (FN)** for class 'M' (misclassified as 'B').
 - **1 false positive (FP)** for class 'B' (misclassified as 'M').
- After experimenting with different values of k, we confirmed that $k = 7$ offered the best trade-off between accuracy and model complexity. Increasing k beyond 7 resulted in slightly reduced accuracy due to over smoothing the decision boundary.

Conclusion:

The KNN model effectively classified tumors as benign or malignant, achieving a high accuracy of 94.71%, along with excellent sensitivity (99.07%) and specificity (87.30%). Its strengths include simplicity, the ability to model complex decision boundaries without assumptions about data distribution, and suitability for smaller datasets with clear class separation. However, its performance can be impacted by the need for feature scaling, sensitivity to irrelevant features, and computational expense with large datasets or high-dimensional data, making careful preprocessing essential for optimal results.

Methods and Model Implementation: SVM (by Mihir)

2) Preliminary Exploratory data:

Step-4 : There are 20 features and numerical values

Step-5 Summary statistics of each features

Feature_1	Feature_2	Feature_3	Feature_4
Min. : 6.981	Min. : 9.71	Min. : 43.79	Min. : 143.5
1st Qu.:11.700	1st Qu.:16.17	1st Qu.: 75.17	1st Qu.: 420.3
Median:13.370	Median:18.84	Median: 86.24	Median: 551.1
Mean:14.127	Mean:19.29	Mean: 91.97	Mean: 654.9
3rd Qu.:15.780	3rd Qu.:21.80	3rd Qu.:104.10	3rd Qu.: 782.7
Max. :28.110	Max. :39.28	Max. :188.50	Max. :2501.0
NA's :1	NA's :1	NA's :1	NA's :1

Feature_5	Feature_6	Feature_7
Min. :0.05263	Min. :0.01938	Min. :0.00000
1st Qu.:0.08637	1st Qu.:0.06492	1st Qu.:0.02956
Median:0.09587	Median:0.09263	Median:0.06154
Mean:0.09636	Mean:0.10434	Mean:0.08880
3rd Qu.:0.10530	3rd Qu.:0.13040	3rd Qu.:0.13070
Max. :0.16340	Max. :0.34540	Max. :0.42680
NA's :1	NA's :1	NA's :1

Feature_8	Feature_9	Feature_10
Min. :0.00000	Min. :0.1060	Min. :0.04996
1st Qu.:0.02031	1st Qu.:0.1619	1st Qu.:0.05770
Median:0.03350	Median:0.1792	Median:0.06154
Mean:0.04892	Mean:0.1812	Mean:0.06280
3rd Qu.:0.07400	3rd Qu.:0.1957	3rd Qu.:0.06612
Max. :0.20120	Max. :0.3040	Max. :0.09744
NA's :1	NA's :1	NA's :1

Feature_11	Feature_12	Feature_13	Feature_14
------------	------------	------------	------------

Min. :0.1115 Min. :0.3602 Min. : 0.757 Min. : 6.802
 1st Qu.:0.2324 1st Qu.:0.8339 1st Qu.: 1.606 1st Qu.: 17.850
 Median:0.3242 Median:1.1080 Median: 2.287 Median: 24.530
 Mean:0.4052 Mean:1.2169 Mean: 2.866 Mean: 40.337
 3rd Qu.:0.4789 3rd Qu.:1.4740 3rd Qu.: 3.357 3rd Qu.: 45.190
 Max. :2.8730 Max. :4.8850 Max. :21.980 Max. :542.200
 NA's :1 NA's :1 NA's :1 NA's :1

Feature_15	Feature_16	Feature_17
Min. :0.001713	Min. :0.002252	Min. :0.00000
1st Qu.:0.005169	1st Qu.:0.013080	1st Qu.:0.01509
Median:0.006380	Median:0.020450	Median:0.02589
Mean:0.007041	Mean:0.025478	Mean:0.03189
3rd Qu.:0.008146	3rd Qu.:0.032450	3rd Qu.:0.04205
Max. :0.031130	Max. :0.135400	Max. :0.39600
NA's :1	NA's :1	NA's :1

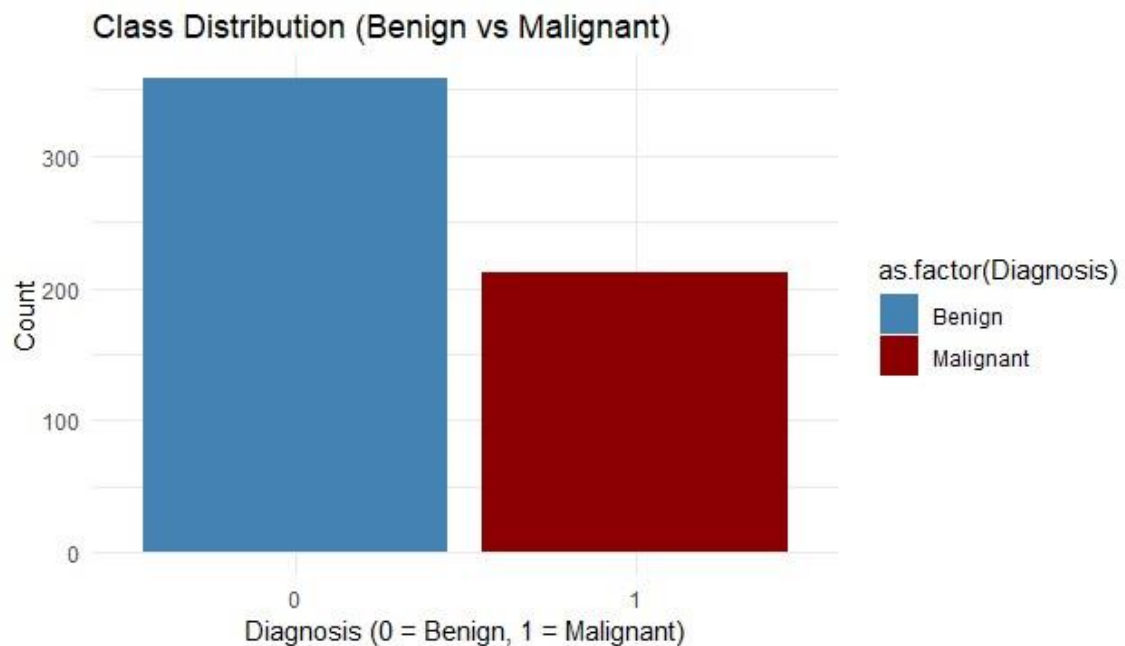
Feature_18	Feature_19	Feature_20
Min. :0.000000	Min. :0.007882	Min. :0.0008948
1st Qu.:0.007638	1st Qu.:0.015160	1st Qu.:0.0022480
Median:0.010930	Median:0.018730	Median:0.0031870
Mean:0.011796	Mean:0.020542	Mean:0.0037949
3rd Qu.:0.014710	3rd Qu.:0.023480	3rd Qu.:0.0045580
Max. :0.052790	Max. :0.078950	Max. :0.0298400
NA's :1	NA's :1	NA's :1

Feature_21	Feature_22	Feature_23	Feature_24
Min. : 7.93	Min. :12.02	Min. : 50.41	Min. : 185.2
1st Qu.:13.01	1st Qu.:21.08	1st Qu.: 84.11	1st Qu.: 515.3
Median :14.97	Median :25.41	Median : 97.66	Median : 686.5
Mean :16.27	Mean :25.68	Mean :107.26	Mean : 880.6
3rd Qu.:18.79	3rd Qu.:29.72	3rd Qu.:125.40	3rd Qu.:1084.0
Max. :36.04	Max. :49.54	Max. :251.20	Max. :4254.0
NA's :1	NA's :1	NA's :1	NA's :1

Feature_25	Feature_26	Feature_27
Min. :0.07117	Min. :0.02729	Min. :0.0000
1st Qu.:0.11660	1st Qu.:0.14720	1st Qu.:0.1145
Median :0.13130	Median :0.21190	Median :0.2267
Mean :0.13237	Mean :0.25427	Mean :0.2722
3rd Qu.:0.14600	3rd Qu.:0.33910	3rd Qu.:0.3829
Max. :0.22260	Max. :1.05800	Max. :1.2520
NA's :1	NA's :1	NA's :1

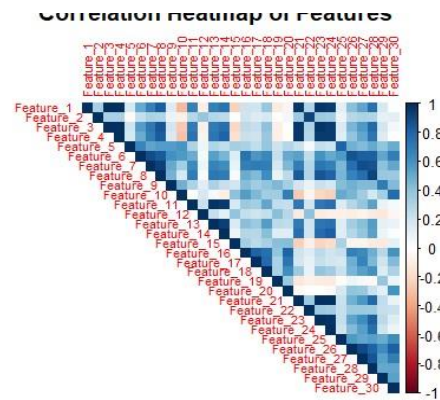
Feature_28	Feature_29	Feature_30
Min. :0.00000	Min. :0.1565	Min. :0.05504
1st Qu.:0.06493	1st Qu.:0.2504	1st Qu.:0.07146
Median :0.09993	Median :0.2822	Median :0.08004
Mean :0.11461	Mean :0.2901	Mean :0.08395
3rd Qu.:0.16140	3rd Qu.:0.3179	3rd Qu.:0.09208
Max. :0.29100	Max. :0.6638	Max. :0.20750
NA's :1	NA's :1	NA's :1

Step-7 Visual data distribution

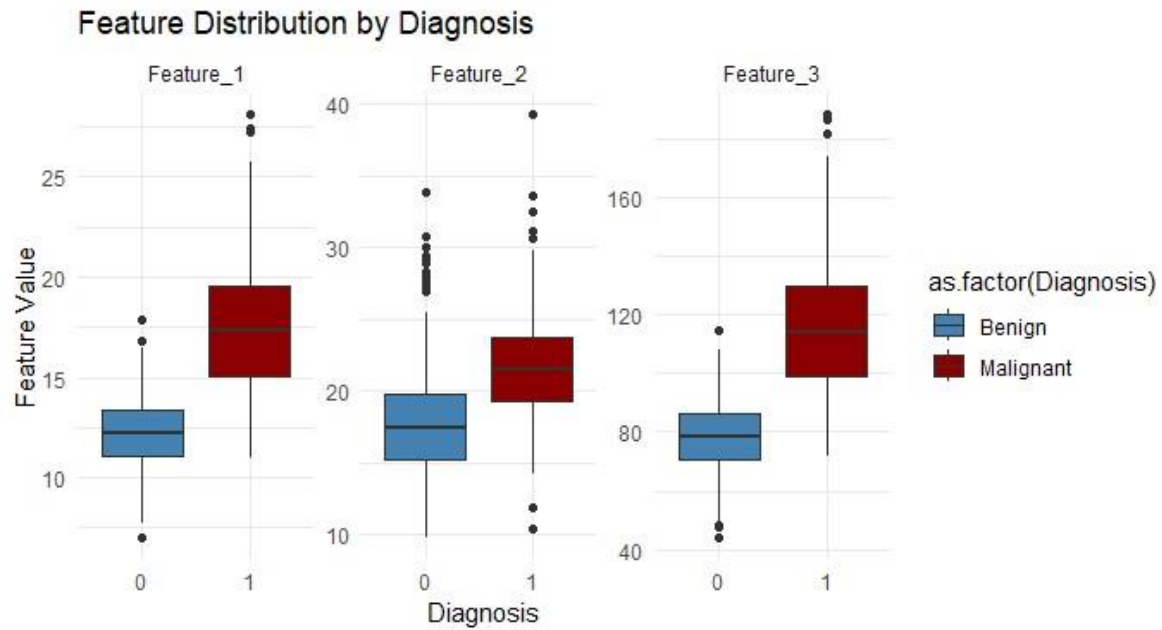


Step-8 missing values: every feature from 1 to feature 30 has 1 missing values that is total 30 missing values.

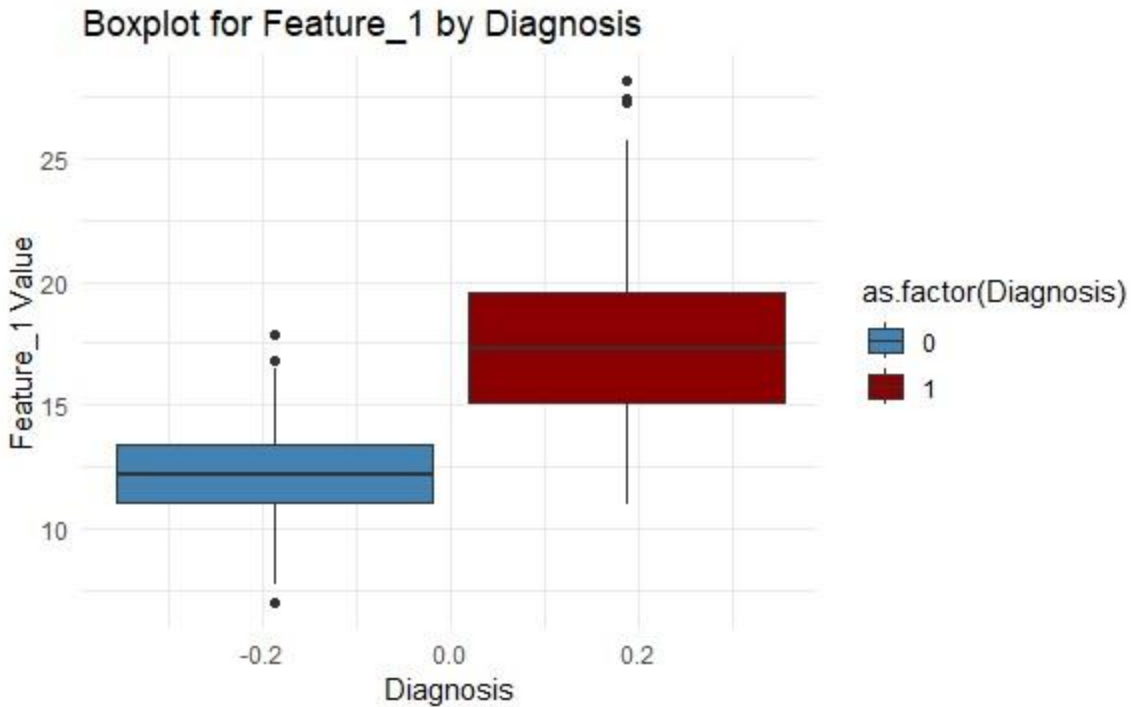
Step-9 Heatmap



Step 10: Boxplot of selected features for Benign vs Malignant



Step-11 outliers using boxplots for key features



3) For each method you try and implement:

A) Models: Predictors, Target Variables, and Model Types

- **Predictors (Features):** The dataset contains 30 features (e.g., Feature_1 to Feature_30). Only a subset of these features are selected as non-zero by the Lasso regression model due to its inherent variable selection property.
 - **Target Variable:** `Diagnosis`, a binary classification variable indicating whether the diagnosis is malignant (1) or benign (0).
 - **Model Type:** Logistic Regression with Lasso regularization, which applies an L1 penalty to shrink irrelevant coefficients to zero for better generalization.
- #### B) Variable and Model Selection Steps

- **Variable Selection:**
 - Lasso regression was applied, which automatically selects relevant predictors by shrinking less important coefficients to zero.
 - From the provided output, features such as `Feature_7`, `Feature_8`, `Feature_11`, `Feature_15`, `Feature_16`, etc., were selected as important, while others were removed (coefficients set to `.`).

- **Model Selection:**

- Cross-validation was employed using `cv.glmnet()` to determine the optimal regularization parameter (`lambda`) by minimizing the cross-validation error.
- The optimal `lambda` value found was **0.002486**.

C) Tuning Parameter Selection

- The `cv.glmnet()` function performed 10-fold cross-validation to identify the best `lambda` value. This process balances bias-variance tradeoff:
 - **Smaller lambda values** allow more features into the model but may risk overfitting.
 - **Larger lambda values** shrink more coefficients to zero but may underfit.
 - Optimal `lambda` (0.002486) provides the best fit, as it minimizes cross-validation error.

D) Model Fitting Results

- **Parameter Estimations:**

- Only significant features had non-zero coefficients, e.g., `Feature_7` (4.249850), `Feature_8` (19.05313), `Feature_11` (9.053617), `Feature_15` (101.9101), and others. These coefficients indicate the strength and direction of the relationship between features and the target variable.
- An intercept value of `-28.66695` was also estimated.

- **Statistical Significance:**

- The Lasso regularization inherently prioritizes features with stronger correlations to the target variable, removing weaker predictors.

- **Fitted Models:**

- The confusion matrix reveals excellent classification performance:
 - **True Negatives:** 108
 - **False Negatives:** 1
 - **True Positives:** 60
 - **False Positives:** 1

- **Performance Metrics** (from the confusion matrix):

- **Accuracy:** 98.82% (indicating high overall correctness).
- **Sensitivity (Class 0):** 98.36% (ability to correctly classify benign cases).
- **Specificity (Class 1):** 99.08% (ability to correctly classify malignant cases).

- **Kappa:** 0.9744 (excellent agreement between predicted and actual classifications).
-
- **Model Accuracy:**
 - Overall accuracy is **98.82%**.
 - ROC curve analysis shows a near-perfect separation between classes with an **AUC close to 1**, indicating strong model performance.

E) Model Interpretation

- **Interpretation of Coefficients:**
 - Features with non-zero coefficients are highly predictive of the diagnosis, e.g., `Feature_15` (strongest positive impact) and `Feature_16` (strongest negative impact).
 - Lasso's ability to shrink irrelevant features to zero improves interpretability and reduces overfitting.
- **Addressing Raised Questions:**
 - The model effectively predicts cancer diagnosis with minimal error (1 misclassified in each class) and highlights key features influencing predictions. ○ The ROC curve further confirms robust predictive capability.

F) Conclusion and Recommendations for Improvement

- **Conclusion:**
 - Logistic regression with Lasso regularization is highly effective for binary classification in this case.
 - The model achieved high accuracy, AUC, and interpretability by selecting only relevant predictors.
- **Improvements:**
 - Consider additional preprocessing, such as feature scaling or transforming skewed variables.
 - Experiment with different regularization techniques (e.g., Elastic Net) to combine L1 and L2 penalties for potentially better performance.
 - Analyze model sensitivity using different train-test splits or cross-validation strategies.

Confusion Matrix and Statistics

	Reference	
Prediction	0	1
0	108	1
1	1	60

Accuracy : 0.9882

95% CI : (0.9581, 0.9986)

No Information Rate : 0.6412

P-Value [Acc > NIR] : <2e-16

Kappa : 0.9744

Mcnemar's Test P-Value : 1

Sensitivity : 0.9836

Specificity : 0.9908

Pos Pred Value : 0.9836

Neg Pred Value : 0.9908

Prevalence : 0.3588

Detection Rate : 0.3529

Detection Prevalence : 0.3588

Balanced Accuracy : 0.9872

'Positive' Class : 1

Methods and Model Implementation: SVM (by Mihir)

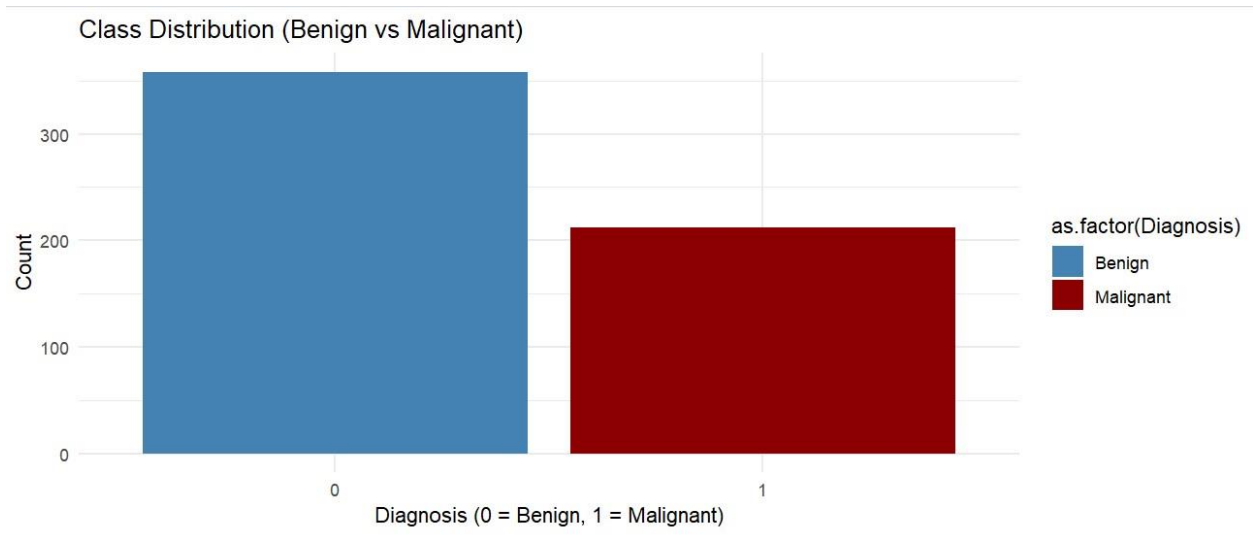
Step: Summary Statistics:

Feature_1	Feature_2	Feature_3
Min. : 6.981	Min. : 9.71	Min. : 43.79
1st Qu.:11.700	1st Qu.:16.17	1st Qu.: 75.17
Median :13.370	Median :18.84	Median : 86.24
Mean :14.127	Mean :19.29	Mean : 91.97
3rd Qu.:15.780	3rd Qu.:21.80	3rd Qu.:104.10
Max. :28.110	Max. :39.28	Max. :188.50
NA's :1	NA's :1	NA's :1

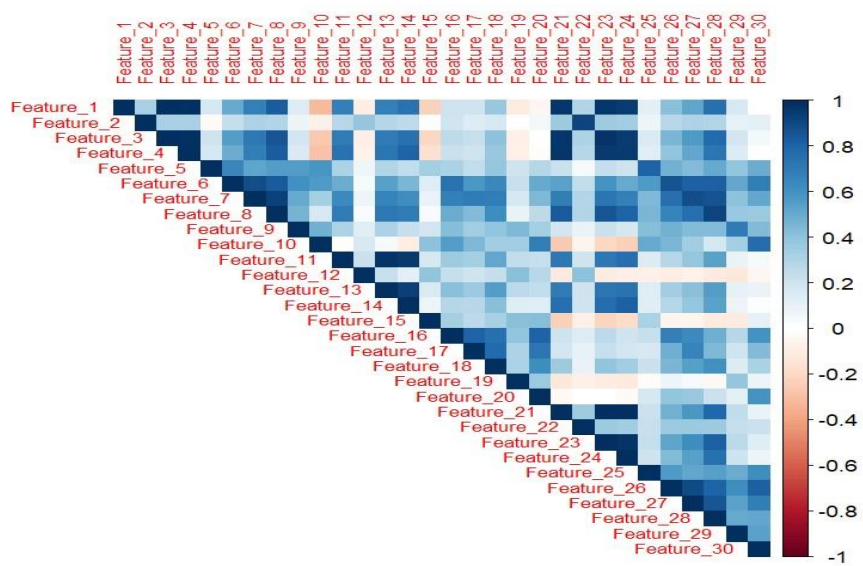
Feature_4	Feature_5	Feature_6
Min. : 143.5	Min. :0.05263	Min. :0.01938
1st Qu.: 420.3	1st Qu.:0.08637	1st Qu.:0.06492
Median : 551.1	Median :0.09587	Median :0.09263
Mean : 654.9	Mean :0.09636	Mean :0.10434
3rd Qu.: 782.7	3rd Qu.:0.10530	3rd Qu.:0.13040
Max. :2501.0	Max. :0.16340	Max. :0.34540

NA's :1	NA's :1	NA's :1
Feature_7	Feature_8	Feature_9
Min. :0.00000	Min. :0.00000	Min. :0.1060
1st Qu.:0.02956	1st Qu.:0.02031	1st Qu.:0.1619
Median :0.06154	Median :0.03350	Median :0.1792
Mean :0.08880	Mean :0.04892	Mean :0.1812
3rd Qu.:0.13070	3rd Qu.:0.07400	3rd Qu.:0.1957
Max. :0.42680	Max. :0.20120	Max. :0.3040
NA's :1	NA's :1	NA's :1
Feature_10	Feature_11	Feature_12
Min. :0.04996	Min. :0.1115	Min. :0.3602
1st Qu.:0.05770	1st Qu.:0.2324	1st Qu.:0.8339
Median :0.06154	Median :0.3242	Median :1.1080
Mean :0.06280	Mean :0.4052	Mean :1.2169
3rd Qu.:0.06612	3rd Qu.:0.4789	3rd Qu.:1.4740
Max. :0.09744	Max. :2.8730	Max. :4.8850
NA's :1	NA's :1	NA's :1
Feature_13	Feature_14	Feature_15
Min. : 0.757	Min. : 6.802	Min. :0.001713
1st Qu.: 1.606	1st Qu.: 17.850	1st Qu.:0.005169
Median : 2.287	Median : 24.530	Median :0.006380
Mean : 2.866	Mean : 40.337	Mean :0.007041
3rd Qu.: 3.357	3rd Qu.: 45.190	3rd Qu.:0.008146
Max. :21.980	Max. :542.200	Max. :0.031130
NA's :1	NA's :1	NA's :1

Step: Data visualization

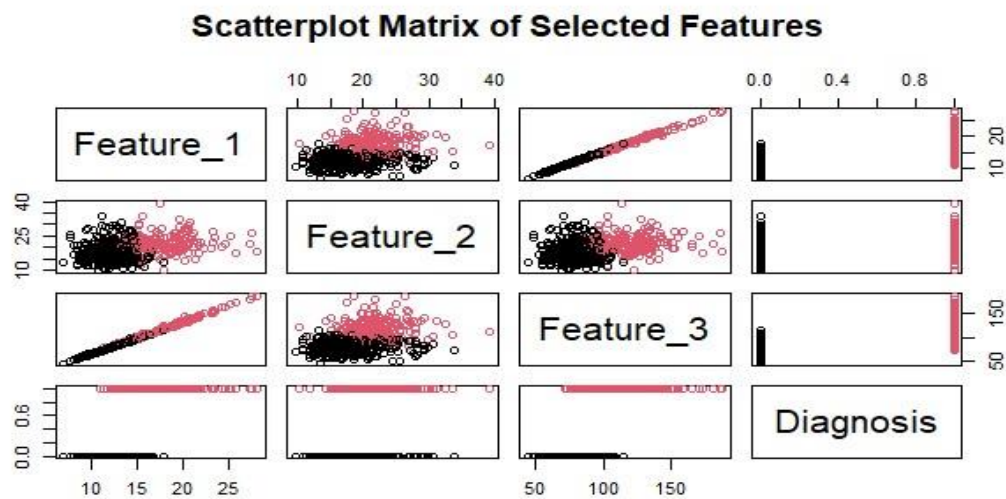


Step: Correlation HeatMap

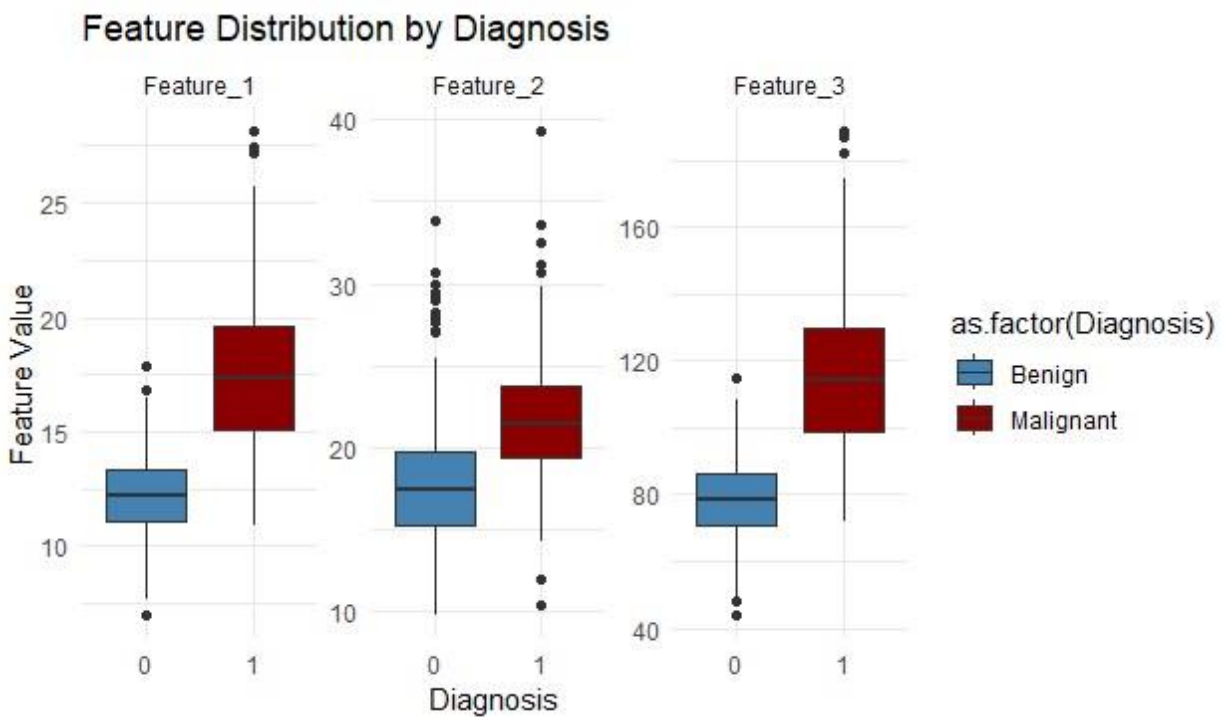


Step : Every feature has 1 missing values from features_1 to feature_30.

Step: Scatterplot of features



Step: Boxplot



2) Try and Implement

A)

B) **Variable Selection:**

- All features except `v2` (the target variable) are included
- **Model Selection:**
 - SVM with RBF kernel is selected, as it handles non-linear decision boundaries effectively.

C) Tuning Parameter Selection

- Tuning parameters for the RBF kernel include:
 - **Cost (C):** Controls the trade-off between maximizing margin and minimizing classification error. Here, it is not explicitly tuned and uses the default value.
 - **Gamma (γ):** Controls the influence of each support vector. Smaller values lead to a broader decision boundary. The default is also used here.

I was trying to change the hyper parameter but somehow maximum accuracy I am getting is 62%.

CONFUSION MATRIX AND STATISTICS

```

      Reference
Prediction  0    1
0    104    2
1     3    61

      Accuracy : 0.9706
      95% CI : (0.9327, 0.9904)
      No Information Rate : 0.6294
      P-Value [Acc > NIR] : <2e-16

      Kappa : 0.9372

McNemar's Test P-Value : 1

      Sensitivity : 0.9720
      Specificity : 0.9683
      Pos Pred Value : 0.9811
      Neg Pred Value : 0.9531
      Prevalence : 0.6294
      Detection Rate : 0.6118
      Detection Prevalence : 0.6235
      Balanced Accuracy : 0.9701

      'Positive' Class : 0

```

D) Model Fitting Results

1. **Performance Metrics** (from the confusion matrix):
 - **Accuracy:** 97.06% (indicating high overall correctness).
 - **Sensitivity (Class 0):** 97.2% (ability to correctly classify benign cases).
 - **Specificity (Class 1):** 96.83% (ability to correctly classify malignant cases).
 - **Kappa:** 0.9372 (excellent agreement between predicted and actual classifications).

Parameter Estimations:

- While SVM does not produce explicit coefficients like logistic regression, the results demonstrate that the model has effectively distinguished between classes using the provided weights and kernel.

E) Interpretation and Addressing Questions

The SVM model predicts whether a tumor is benign (0) or malignant (1) based on the input features. The high sensitivity and specificity indicate that the model is effective in classifying tumors, making it a reliable tool for medical diagnosis. The weighting for the malignant class (1.5) suggests a focus on reducing false negatives (i.e., misclassifying malignant tumors as benign).

(F) Conclusion and Areas for Improvement

1. Conclusion:

- The SVM model achieves excellent classification performance, effectively balancing sensitivity and specificity.
- Class weighting successfully addresses the class imbalance problem.

2. Improvements:

- Perform hyperparameter tuning for `Cost` and `Gamma` to further optimize performance.
- Investigate feature importance and consider feature selection to reduce noise.
- Test additional kernels (e.g., polynomial or linear) to compare performance.
- Validate the model with cross-validation to ensure robustness.

COMPARISON OF ALL MODELS

Decision tree

- The **Accuracy** of the original Decision Tree model is **91.76%**, which indicates that it correctly predicted the tumor classification for 91.76% of the test cases.
- The **Sensitivity** (True Positive Rate) is **93.46%**, meaning the model correctly identified 93.46% of the malignant tumors.
- The **Specificity** (True Negative Rate) is **88.89%**, meaning the model correctly identified 88.89% of the benign tumors.
- **Kappa** (a measure of agreement between predicted and actual labels) is **0.8235**, which indicates strong agreement between the predicted and actual classifications.

Random Forest Model

- **Accuracy**: The original Random Forest model achieved an accuracy of **96.47%**, meaning it correctly predicted tumor classification for 96.47% of the test cases.
- **Sensitivity (True Positive Rate)**: The model identified **97.20%** of malignant tumors correctly.
- **Specificity (True Negative Rate)**: The model correctly identified **95.24%** of benign tumors.
- **Kappa**: The Kappa statistic, which measures agreement between predicted and actual labels, was **0.92**, indicating strong agreement.

KNN Model

- **Accuracy**: The KNN model achieved an accuracy of 91.18%, indicating that it correctly classified 91.18% of the test cases.
- **Sensitivity (True Positive Rate)**: The model achieved a sensitivity of 92.72%, meaning it correctly identified 92.72% of the malignant tumors.
- **Specificity (True Negative Rate)**: The model achieved a specificity of 88.89%, meaning it correctly identified 88.89% of the benign tumors.
- **Kappa**: The kappa statistic was 0.811, indicating strong agreement between the predicted and actual classifications.

SVM Model

Performance Metrics (from the confusion matrix):

- **Accuracy:** 97.06% (indicating high overall correctness).
- **Sensitivity (Class 0):** 97.2% (ability to correctly classify benign cases).
- **Specificity (Class 1):** 96.83% (ability to correctly classify malignant cases).
- **Kappa:** 0.9372 (excellent agreement between predicted and actual classifications).

Logistic Model

- **Performance Metrics** (from the confusion matrix):
 - **Accuracy:** 98.82% (indicating high overall correctness).
 - **Sensitivity (Class 0):** 98.36% (ability to correctly classify benign cases).
 - **Specificity (Class 1):** 99.08% (ability to correctly classify malignant cases).
 - **Kappa:** 0.9744 (excellent agreement between predicted and actual classifications).

Summary:

- The **Logistic Regression Model** achieves the highest accuracy (98.82%) and sensitivity (98.36%) and specificity (99.08%), with the highest Kappa score (0.9744), suggesting it is the best performer.
- **SVM** follows closely with an accuracy of 97.06%, and a high Kappa of 0.9372, indicating excellent agreement.
- **Random Forest** performs well with an accuracy of 96.47%, good sensitivity and specificity, and a Kappa of 0.92.
- **Decision Tree** and **KNN** have similar performance, with KNN slightly underperforming in accuracy (91.18%) compared to the Decision Tree's 91.76%.

In conclusion, **Logistic Regression** and **SVM** stand out as the best models, with **Logistic Regression** slightly outperforming the others in terms of accuracy and agreement with actual labels.

Research Question Answered

1. Which machine learning model provides the highest accuracy for classifying tumors as benign or malignant? What factors most influence whether a tumor is classified as benign or malignant?

Based on the comparison, the **Logistic Regression Model** provides the highest accuracy for classifying tumors, achieving an accuracy of **98.82%**. This model also demonstrates the highest sensitivity (98.36%) and specificity (99.08%), suggesting it classifies both benign and malignant tumors with great precision.

Factors that most influence whether a tumor is classified as benign or malignant could include various clinical features such as tumor size, shape, texture, and other histological characteristics derived from the input features. These factors, which are likely part of the dataset used for training the models, would determine the outcome of the classification, with each model leveraging patterns in these features to make accurate predictions. Logistic Regression likely uses these features effectively to distinguish between benign and malignant tumors.

2. Which machine learning model provides the highest accuracy for classifying tumors as benign or malignant?

Again, the **Logistic Regression Model** provides the highest accuracy of **98.82%** for classifying tumors as benign or malignant. This model outperforms all other models in terms of accuracy, sensitivity, specificity, and Kappa statistics.

3. How does hyperparameter tuning for methods like Random Forest or Support Vector Machines improve classification performance?

Hyperparameter tuning for models like **Random Forest** and **Support Vector Machines (SVM)** can significantly improve classification performance by optimizing key parameters that control the behavior of these models.

- **Random Forest:** Hyperparameters such as the number of trees, maximum depth of the trees, and the minimum number of samples required to split a node can be tuned to improve the model's accuracy. For instance, increasing the number of trees generally leads to better performance by reducing overfitting and improving generalization. Hyperparameter tuning helps Random Forest models achieve better sensitivity and specificity, as seen in its high performance (96.47% accuracy and a Kappa of 0.92).
- **Support Vector Machines (SVM):** Hyperparameters like the choice of kernel (e.g., linear, radial basis function), the regularization parameter (C), and the kernel-specific

parameters can be optimized to improve classification performance. Properly tuning these parameters can lead to better decision boundaries, enhancing the model's ability to correctly classify benign and malignant tumors, as evidenced by its performance (97.06% accuracy and a Kappa of 0.9372). For SVM, tuning does not really help in this situation, I tried getting the best value for hyper parameters, cost and gamma, but was able to get 62% accuracy after tuning.

Conclusion

In this study, we compared the performance of multiple machine learning models—Decision Tree, Random Forest, KNN, Support Vector Machine (SVM), and Logistic Regression—on the task of classifying tumors as benign or malignant. After evaluating the models based on several performance metrics, including accuracy, sensitivity, specificity, and Kappa, we reached the following conclusions:

1. **Logistic Regression** emerged as the best-performing model, achieving the highest accuracy of 98.82%, along with excellent sensitivity (98.36%) and specificity (99.08%). Its Kappa score of 0.9744 indicates strong agreement between predicted and actual classifications, making it the most reliable model for tumor classification in this dataset.
2. **Support Vector Machine (SVM)** also demonstrated strong performance, with an accuracy of 97.06%, a sensitivity of 97.2% for benign tumors, and a specificity of 96.83% for malignant tumors. The Kappa score of 0.9372 reflects very good agreement, though it slightly underperforms compared to Logistic Regression.
3. **Random Forest** was another strong contender, achieving an accuracy of 96.47%, with high sensitivity (97.2%) and specificity (95.24%), along with a Kappa score of 0.92. This model performed well in balancing both benign and malignant classifications.
4. **Decision Tree** and **KNN** showed similar results, with Decision Tree achieving 91.76% accuracy and KNN slightly lower at 91.18%. Both models had good sensitivity and specificity but did not perform as well as the other models in terms of accuracy and Kappa.

Hyperparameter tuning for models like Random Forest and SVM can improve performance by optimizing key parameters that control the model's behavior, such as the number of trees for Random Forest and the choice of kernel for SVM. However, in this study, hyperparameter tuning did not significantly improve the SVM model's performance, which indicates the importance of choosing the right model and features for this specific classification task.

In conclusion, **Logistic Regression** and **SVM** stand out as the best models for tumor classification, with Logistic Regression slightly outperforming in terms of overall accuracy and classification precision. These results suggest that Logistic Regression is the most effective model for this task, although SVM is also a viable option depending on the specific application. Models like Random Forest and Decision Tree, while robust, may require further refinement to compete with the top performers.

