# Wisconsin Breast Cancer Diagnosis Analysis

Om Gor, Matt Khosrowabadi, Svetlana Voda

12/05/24

```
##
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
##
##     filter, lag
```

```
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```
## Loading required package: ggplot2
```

```
## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg   ggplot2
```

```
## Loading required package: lattice
```

```
## randomForest 4.7-1.2
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
##
## Attaching package: 'randomForest'
```

```
## The following object is masked from 'package:ggplot2':
##
##     margin
```

```
## The following object is masked from 'package:dplyr':
##
##     combine
```

```
## This is aqp 2.1.0
```

```
##
## Attaching package: 'aqp'
```

```
## The following object is masked from 'package:randomForest':
##
##     combine
```

```
## The following objects are masked from 'package:dplyr':
##
##     combine, slice
```

```
## corrplot 0.95 loaded
```

```
## Loading required package: MASS
```

```
##
## Attaching package: 'MASS'
```

```
## The following object is masked from 'package:dplyr':
##
##     select
```

```
## Loading required package: msm
```

```
## Loading required package: polycor
```

```
## Loading required package: Matrix
```

```
## Loaded glmnet 4.1-8
```

## Abstract

Breast cancer is one of the most prevalent and life-threatening diseases worldwide, necessitating accurate diagnostic tools. This report outlines the development and evaluation of a machine learning pipeline for classifying tumors as benign or malignant using the Breast Cancer Wisconsin Dataset. The analysis incorporated LASSO for feature selection and multiple classifiers—K-Nearest Neighbors (KNN), Naive Bayes, Support Vector Machine (SVM), Decision Tree, and Random Forest—to optimize predictive accuracy. The results highlight the potential of machine learning in enhancing breast cancer diagnosis and early detection.

## Introduction

Breast cancer is a leading cause of morbidity and mortality globally, emphasizing the critical importance of early detection in improving survival rates. This study focuses on developing a robust predictive model to accurately classify tumors as benign or malignant using diagnostic data. Leveraging statistical and machine learning techniques, the goal is to maximize classification accuracy and reliability while ensuring the interpretability of the results.

The proposed solution involves implementing a machine learning pipeline that utilizes LASSO regression for feature selection, followed by training and optimizing multiple classifiers, including K-Nearest Neighbors (KNN), Naive Bayes, Support Vector Machine (SVM), and Decision Tree models. These models are optimized to enhance their accuracy and reliability for tumor classification.

The analysis is based on the Breast Cancer Wisconsin Dataset, which comprises 569 records and 30 diagnostic features derived from digitized images of tumor cells. These features represent geometric and textural properties of cell nuclei, providing critical insights for classification. By combining feature selection with classifier optimization, this study aims to achieve high performance in distinguishing benign and malignant tumors.

## EDA

### Data Cleaning

The data cleaning process involves several key steps to prepare the data set for analysis. First, we changed the names of the columns so that it would be easier for analysis and interpretation. Next, we dropped the ID number column as it does not contribute to the analysis. Before graphing the plots for EDA and understanding the data, we wanted to see if there were any null or missing values in the datasheets. We performed tests on the data to see if there was missing values and null values and we found out that there were none. Finally, we encoded the Diagnosis column, which is the value we are predicting, to 1 or 0. In this column, M indicates the tumor is malignant and B represents the tumor is Benign. The data set contains 569 rows and 31 columns. The rows represent patients and the columns represent attributes such as ID Number, Diagnosis, Radius, Texture, Perimeter, Area, Smoothness, Compactness, Concavity, Concave Points, Symmetry, Fractional Dimension. Some of the attributes in the dataset contain mean, standard error, and worst values.

| | Diagnosis<br><chr> | radius_mean<br><dbl> | texture_mean<br><dbl> | perimeter_mean<br><dbl> | area_mean<br><dbl> | smoothness_mean<br><dbl> | compactness_mean<br><dbl> |
|---|---|---|---|---|---|---|---|
| 1 | M | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 |
| 2 | M | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 |
| 3 | M | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 |
| 4 | M | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 |
| 5 | M | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 |
| 6 | M | 12.45 | 15.70 | 82.57 | 477.1 | 0.12780 | 0.17000 |

6 rows | 1-8 of 32 columns

To begin, we analyze the target variable, which represents the classification of tumors into two categories: benign and malignant
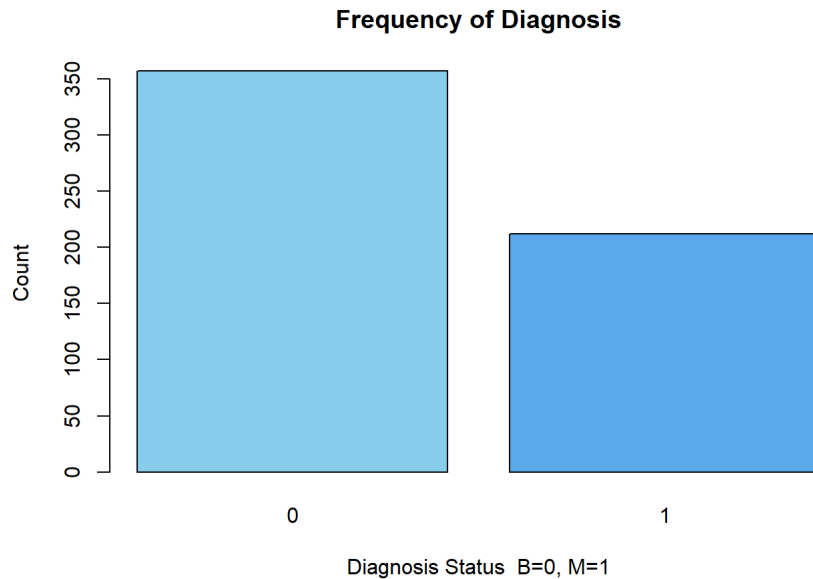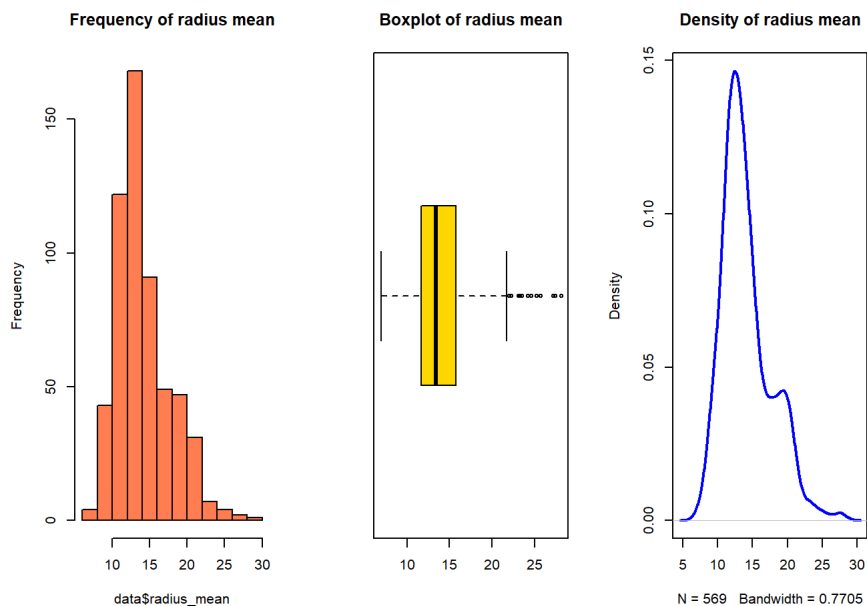
## Frequency of Diagnosis



Diagnosis Status  B=0, M=1

Table of Diagnosis

```
0     | 1       |
----- |-------  |
357   | 212     |
```

Analyzing the bar plot of probation status, the most common category is 'Bening' with 357, compared to 212 students who are 'Malignant' The frequency range spans from 212 to 357. The analysis shows that the majority of tumor cases (62.75%) are benign, which indicates a larger proportion of non-cancerous tumors compared to malignant ones (37.25%).

To better understand the factors influencing the diagnosis of Bening or Malignant, we will conduct a detailed analysis of each predictor variable under consideration. For each variable, we will:

- Provide descriptive Statistics – Including the mean, standard deviation, variance, and a five-number summary (minimum, Q1, median, Q3, and maximum) to understand the central tendency and spread.
- Visualize Distributions – Using histograms and boxplots to gain insights into the shape, skewness, and potential outliers in the data.
- Interpret Results – Highlighting observations about each variable's distribution and variability to provide context for its potential influence on college enrollment.
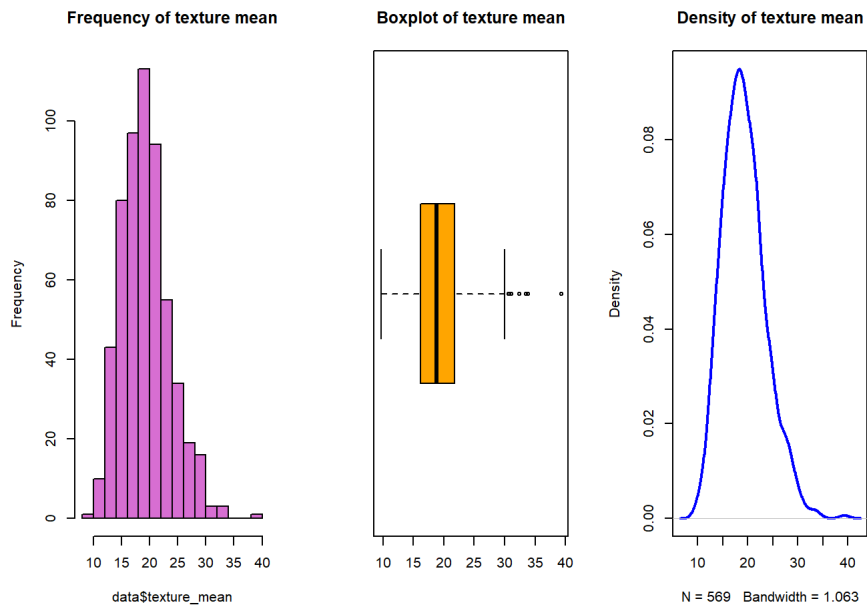
To begin, we will analyze the predictor radius_mean. This variable represents the mean of distances from center to points on the perimeter



| | min | Q1 | median | Q3 | max |
|---|---|---|---|---|---|
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| | 6.981 | 11.695 | 13.37 | 15.815 | 28.11 |

1 row

Analyzing the histogram and boxplot of the quantitative variable radius_mean, we observe that the data is right-skewed with two peaks. From the descriptive statistics, we find that the median is 13.37, the mean is 14.1272917, the standard deviation is 3.5240488, and the variance is 12.4189201. Additionally, the interquartile range (IQR) is 4.12

The subsequent predictor examined was texture_mean which represent the standard deviation of grayscale values
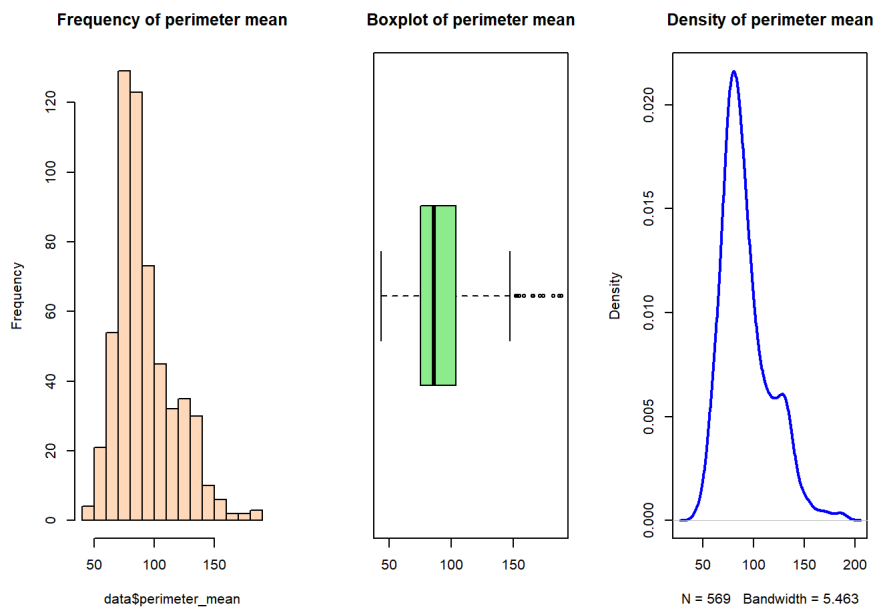


**Frequency of texture mean**     **Boxplot of texture mean**     **Density of texture mean**

| min | Q1 | median | Q3 | max |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
|--- |--- |--- |--- |--- |
| 9.71 | 16.17 | 18.84 | 21.805 | 39.28 |

1 row

Analyzing the histogram and boxplot of the quantitative variable texture_mean, we observe that the data is right-skewed with one peaks and multiple outliers. Most values of texture_mean are concentrated between 15 and 25. From the descriptive statistics, we find that the median is 18.84, the mean is 19.2896485, the standard deviation is 4.3010358, and the variance is 18.4989087. Additionally, the interquartile range (IQR) is 5.635

Next variable is perimeter_mean which represent the total distance between the "snake" points constitutes the nuclear perimeter
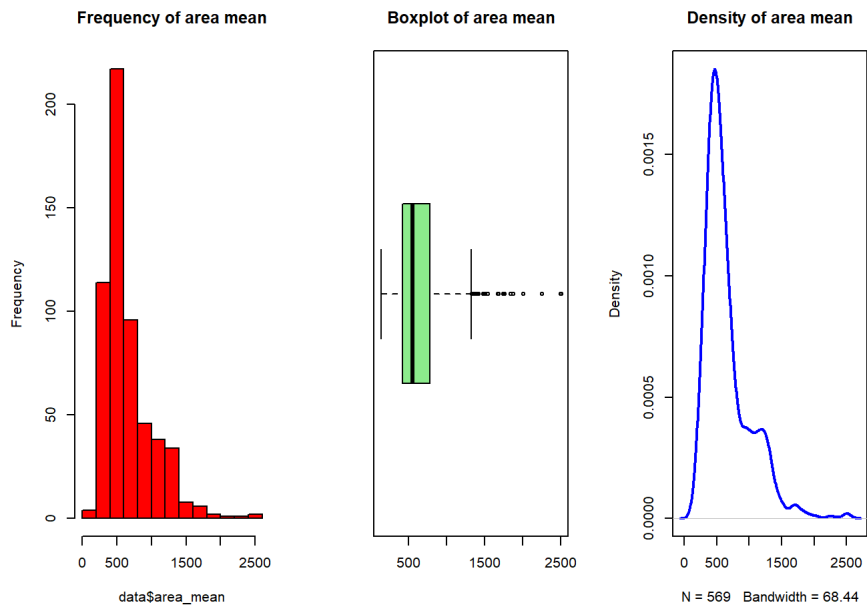


**Frequency of perimeter mean**     **Boxplot of perimeter mean**     **Density of perimeter mean**

| min | Q1 | median | Q3 | max |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
|--- |--- |--- |--- |--- |
| 43.79 | 75.1 | 86.24 | 104.2 | 188.5 |

1 row

Analyzing the histogram and boxplot of the quantitative variable perimeter_mean, we observe that the data is right-skewed with two peaks and multiple outliers. From the descriptive statistics, we find that the median is 86.24, the mean is 91.9690334, the standard deviation is 24.298981, and the variance is 590.4404795. Additionally, the interquartile range (IQR) is 29.1
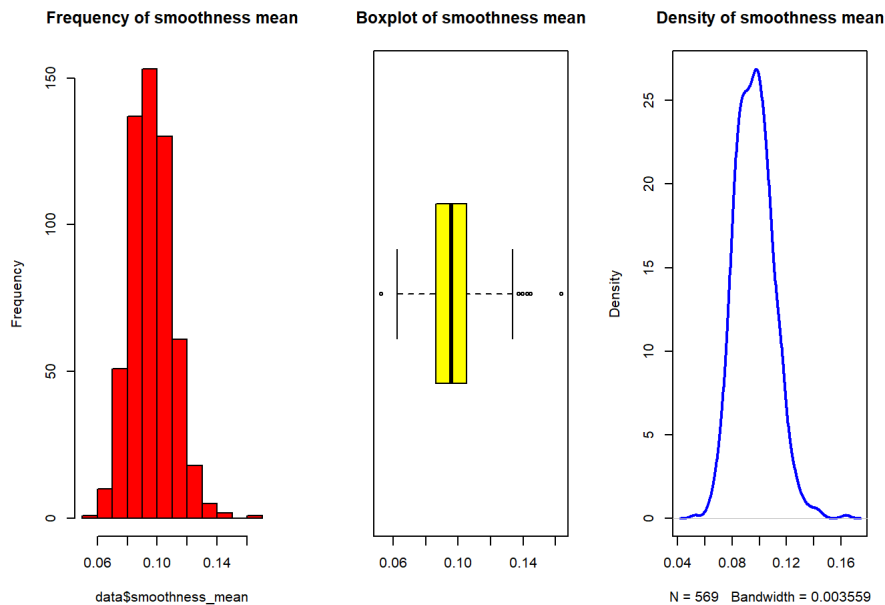
Next variable is area_mean which represent the number of pixels on the interior of a cell and one-half the pixels in the premeter



| min | Q1 | median | Q3 | max |
|---|---|---|---|---|
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 143.5 | 420.05 | 551.1 | 785.6 | 2501 |
| 1 row | | | | |

Analyzing the histogram and boxplot of the quantitative variable area_mean, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 551.1, the mean is 654.8891037, the standard deviation is 351.9141292, and the variance is $1.2384355^{5}$. Additionally, the interquartile range (IQR) is 365.55
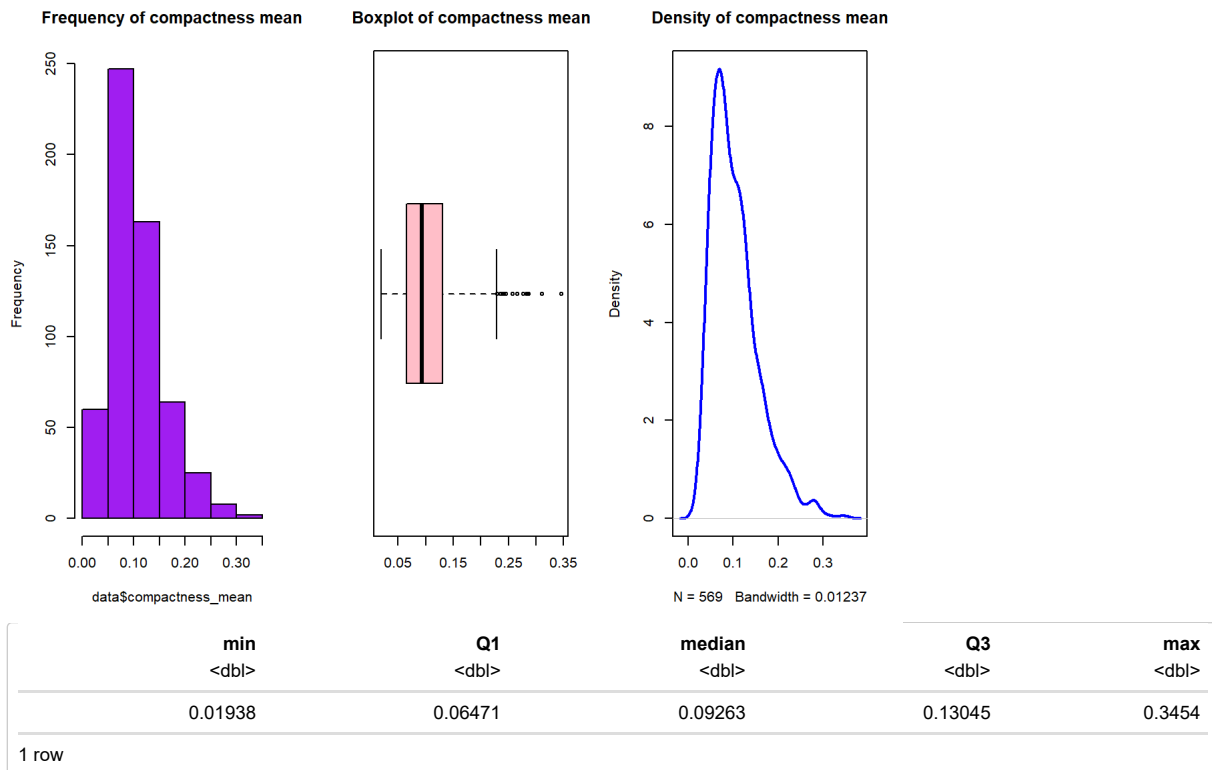
The subsequent variable considered is smoothness_mean which represent the local variation in radius lenghts



| min | Q1 | median | Q3 | max |
|---|---|---|---|---|
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 0.05263 | 0.08621 | 0.09587 | 0.10535 | 0.1634 |
| 1 row | | | | |

Analyzing the histogram and boxplot of the quantitative variable smoothness_mean, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.09587, the mean is 0.0963603, the standard deviation is 0.0140641, and the variance is $1.977997^{-4}$. Additionally, the interquartile range (IQR) is 0.01914

Next variable is compactness_mean which is a measure of the compactness of a cell



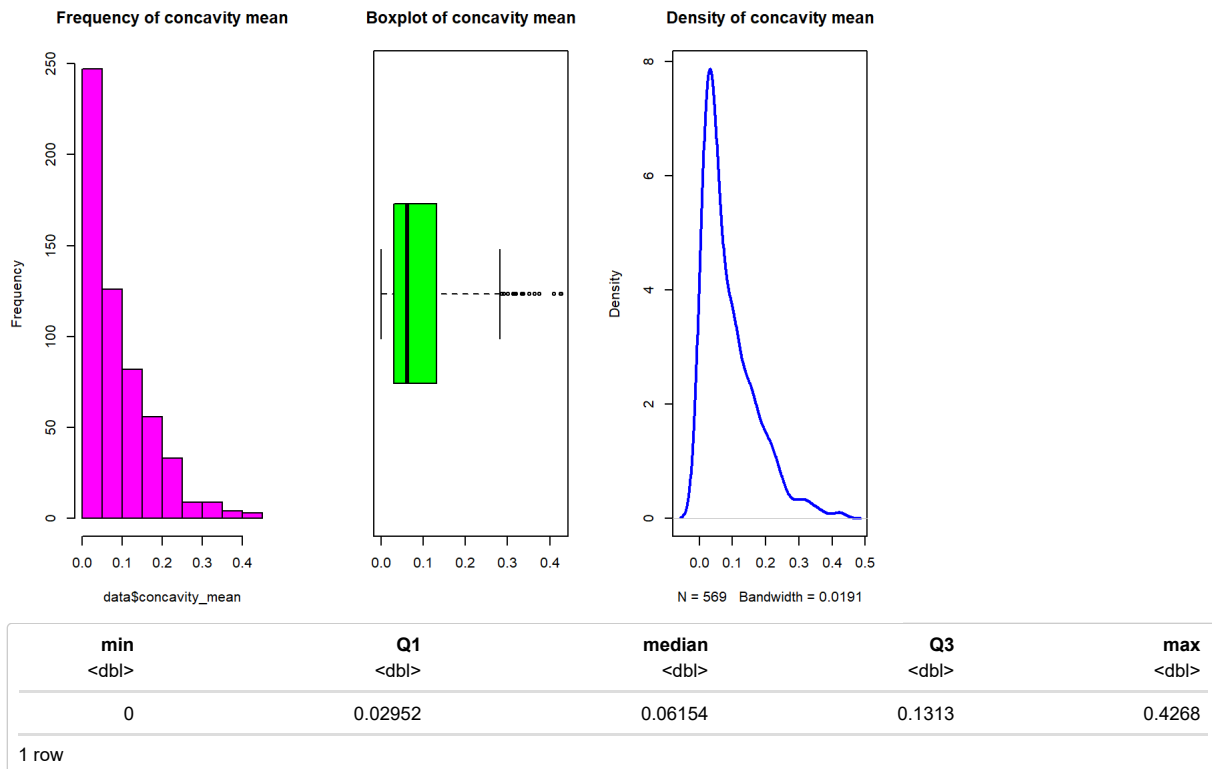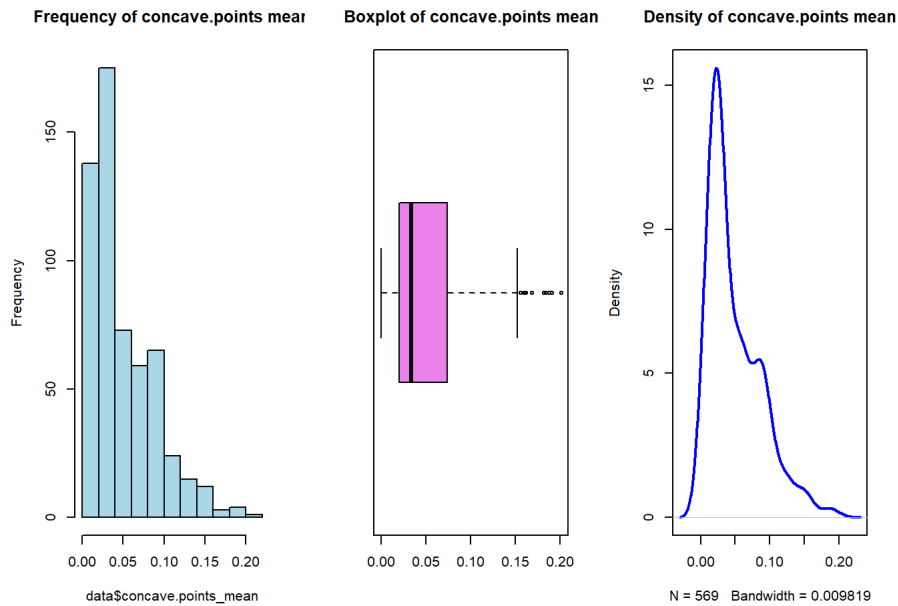| | min<br><dbl> | Q1<br><dbl> | median<br><dbl> | Q3<br><dbl> | max<br><dbl> |
|---|---|---|---|---|---|
| | 0.01938 | 0.06471 | 0.09263 | 0.13045 | 0.3454 |

1 row

Analyzing the histogram and boxplot of the quantitative variable compactness_mean, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.09263, the mean is 0.104341, the standard deviation is 0.0528128, and the variance is 0.0027892. Additionally, the interquartile range (IQR) is 0.06574

Next variable is concavity_mean which represent severity of concave portions of the contur



| | min<br><dbl> | Q1<br><dbl> | median<br><dbl> | Q3<br><dbl> | max<br><dbl> |
|---|---|---|---|---|---|
| | 0 | 0.02952 | 0.06154 | 0.1313 | 0.4268 |

1 row

Analyzing the histogram and boxplot of the quantitative variable concavity_mean, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.06154, the mean is 0.0887993, the standard deviation is 0.0797198, and the variance is 0.0063552. Additionally, the interquartile range (IQR) is 0.10178
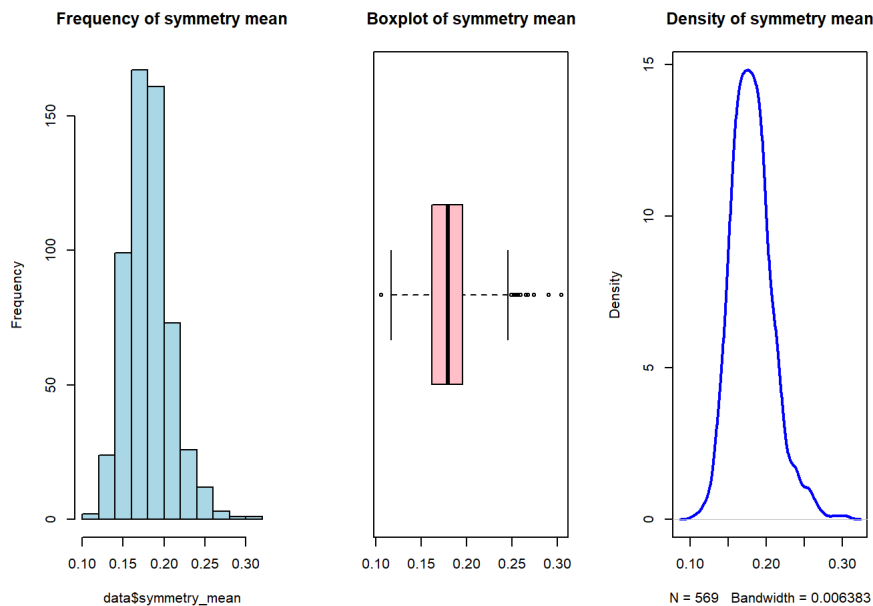
Next variable is concave.points_mean which represent number of concave portions of the contur



| min <dbl> | Q1 <dbl> | median <dbl> | Q3 <dbl> | max <dbl> |
|---|---|---|---|---|
| 0 | 0.02031 | 0.0335 | 0.07402 | 0.2012 |

1 row

Analyzing the histogram and boxplot of the quantitative variable concavity_mean, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.0335, the mean is 0.0489191, the standard deviation is 0.0388028, and the variance is 0.0015057. Additionally, the interquartile range (IQR) is 0.05371
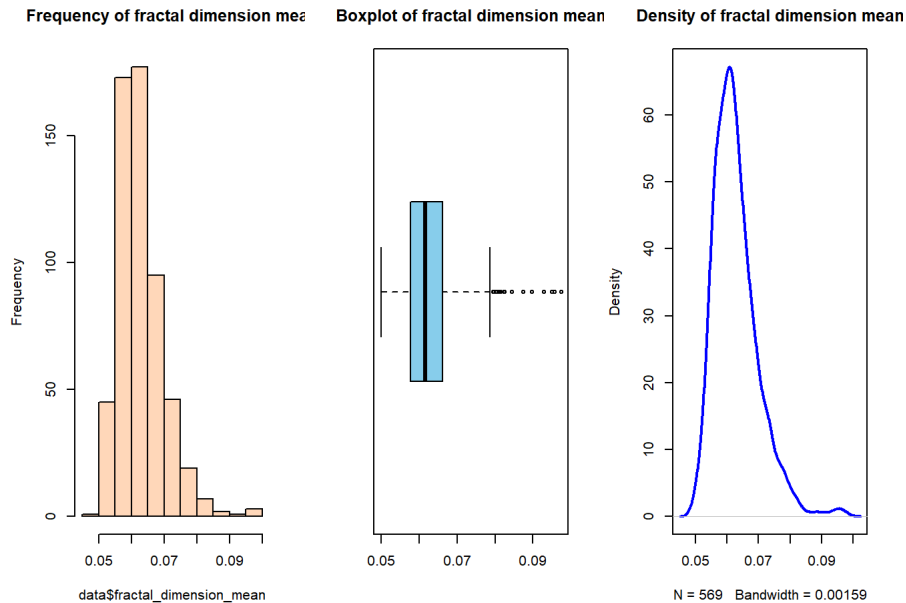
Next variable is symmetry_mean which represent a measure of symmetry of a cell



| min <dbl> | Q1 <dbl> | median <dbl> | Q3 <dbl> | max <dbl> |
|---|---|---|---|---|
| 0.106 | 0.1619 | 0.1792 | 0.1957 | 0.304 |

1 row

Analyzing the histogram and boxplot of the quantitative variable symmetry_mean, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.1792, the mean is 0.1811619, the standard deviation is 0.0274143, and the variance is $7.5154282^{-4}$. Additionally, the interquartile range (IQR) is 0.0338

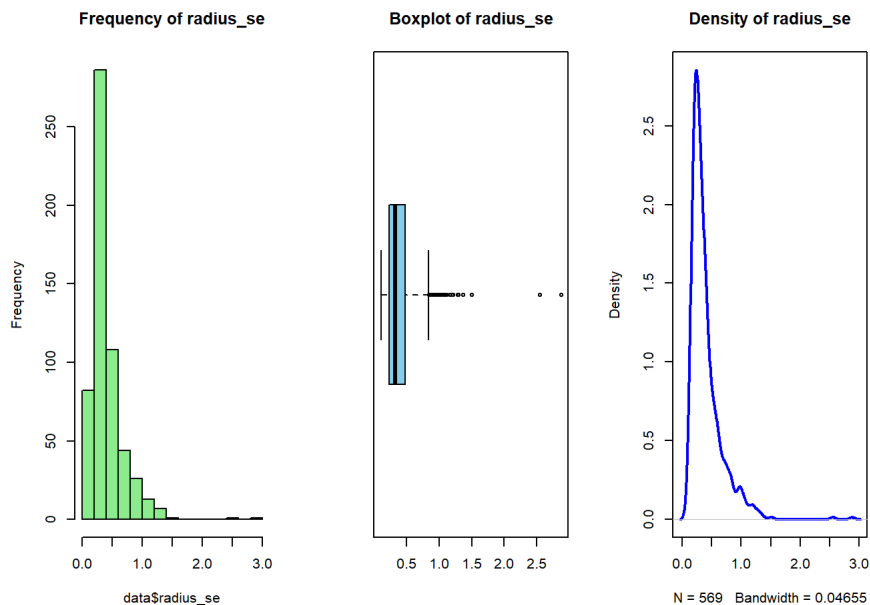Next feature is fractal_dimension_mean

### Frequency of fractal dimension mean

### Boxplot of fractal dimension mean

### Density of fractal dimension mean



| min | Q1 | median | Q3 | max |
| --- | --- | --- | --- | --- |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 0.04996 | 0.057695 | 0.06154 | 0.066135 | 0.09744 |

1 row

Analyzing the histogram and boxplot of the quantitative variable fractal_dimension_mean, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.06154, the mean is 0.1811619, the standard deviation is 0.0274143, and the variance is $7.5154282^{-4}$. Additionally, the interquartile range (IQR) is 0.00844

Next feature is radius_se

### Frequency of radius_se

### Boxplot of radius_se

### Density of radius_se



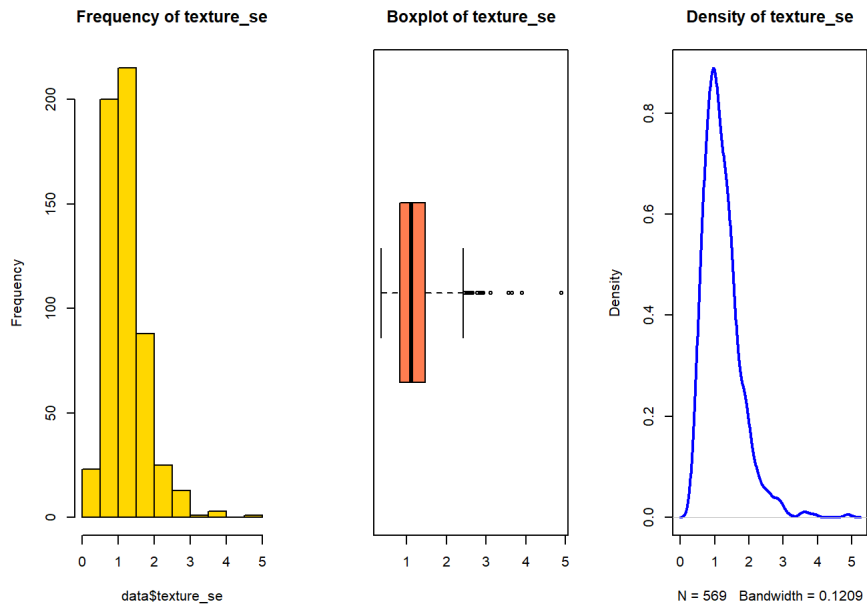| min | Q1 | median | Q3 | max |
| --- | --- | --- | --- | --- |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 0.1115 | 0.23235 | 0.3242 | 0.4807 | 2.873 |

1 row

Analyzing the histogram and boxplot of the quantitative variable radius_se, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.3242, the mean is 0.4051721, the standard deviation is 0.2773127, and the variance is 0.0769024. Additionally, the interquartile range (IQR) is 0.24835
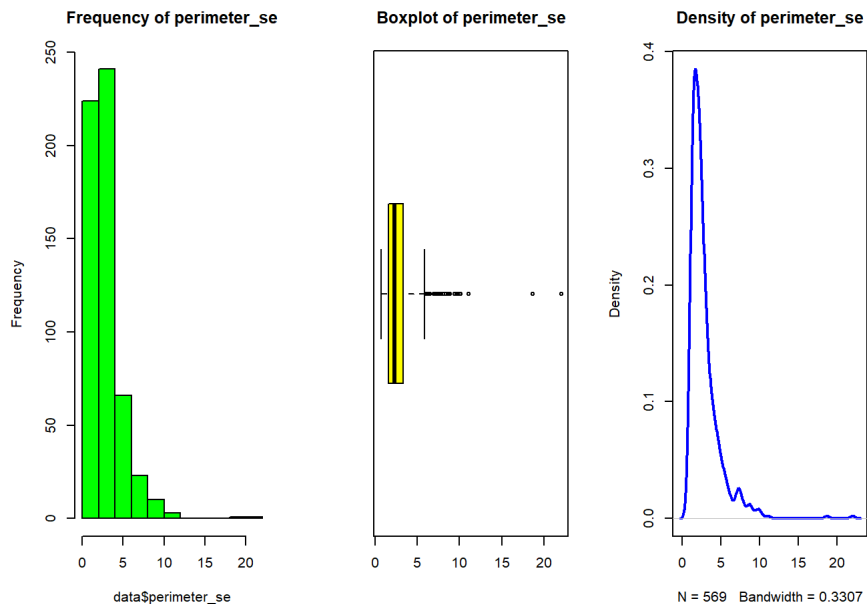
Next feature is texture_se



| min | Q1 | median | Q3 | max |
| --- | --- | --- | --- | --- |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 0.3602 | 0.8324 | 1.108 | 1.4745 | 4.885 |

1 row

Analyzing the histogram and boxplot of the quantitative variable texture_se, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 1.108, the mean is 1.2168534, the standard deviation is 0.5516484, and the variance is 0.3043159. Additionally, the interquartile range (IQR) is 0.6421

Next feature is perimeter_se



| min | Q1 | median | Q3 | max |
| --- | --- | --- | --- | --- |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 0.757 | 1.604 | 2.287 | 3.363 | 21.98 |

1 row

Analyzing the histogram and boxplot of the quantitative feature perimeter_se, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 2.287, the mean is 2.8660592, the standard deviation is 2.0218546, and the variance is 4.0878958. Additionally, the interquartile range (IQR) is 1.759

Next feature is area_se

### Frequency of area_se

### Boxplot of area_se

### Density of area_se

N = 569   Bandwidth = 5.163

| min<br><dbl> | Q1<br><dbl> | median<br><dbl> | Q3<br><dbl> | max<br><dbl> |
|---|---|---|---|---|
| 6.802 | 17.85 | 24.53 | 45.285 | 542.2 |

1 row

Analyzing the histogram and boxplot of the quantitative feature area_se, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 24.53, the mean is 40.3370791, the standard deviation is 45.4910055, and the variance is 2069.4315829. Additionally, the interquartile range (IQR) is 27.435
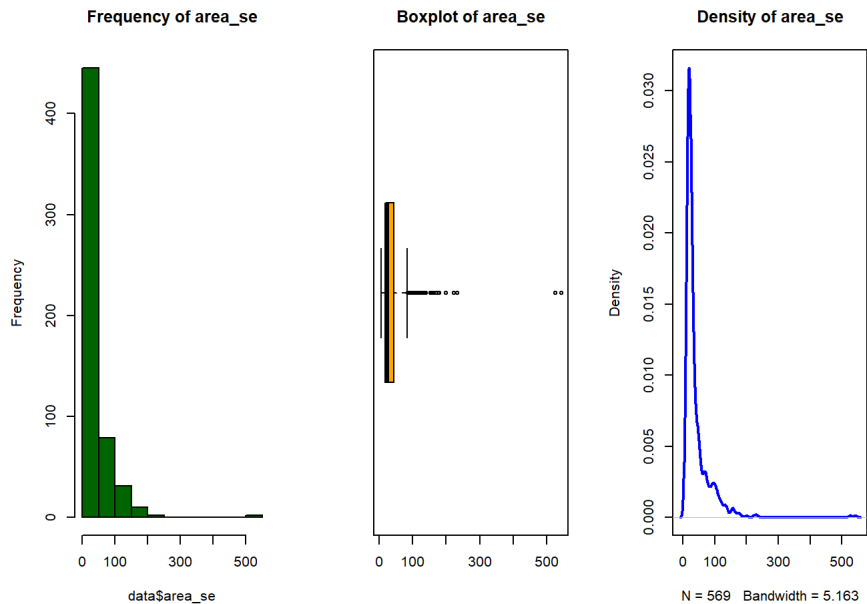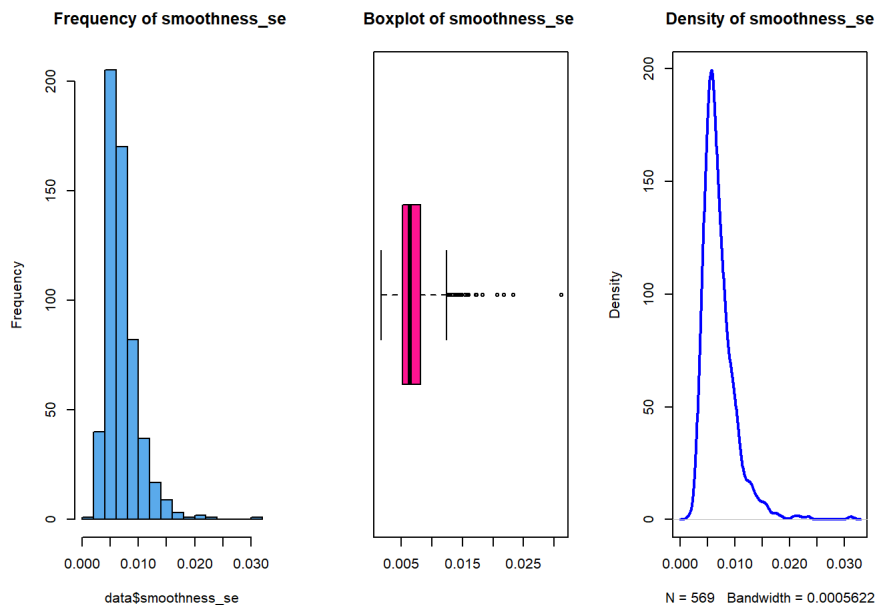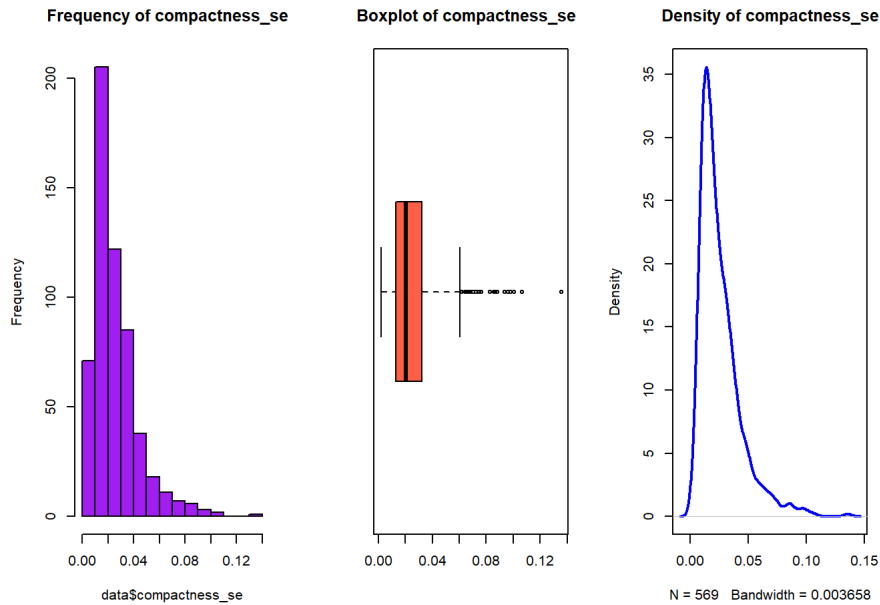
### Frequency of smoothness_se

### Boxplot of smoothness_se

### Density of smoothness_se

N = 569   Bandwidth = 0.0005622

| min<br><dbl> | Q1<br><dbl> | median<br><dbl> | Q3<br><dbl> | max<br><dbl> |
|---|---|---|---|---|
| 0.001713 | 0.0051635 | 0.00638 | 0.008156 | 0.03113 |

1 row

Analyzing the histogram and boxplot of the quantitative feature smoothness_se, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.00638, the mean is 0.007041, the standard deviation is 0.0030025, and the variance is $9.015114^{-6}$. Additionally, the interquartile range (IQR) is 0.00299

**Frequency of compactness_se**    **Boxplot of compactness_se**    **Density of compactness_se**



| min | Q1 | median | Q3 | max |
|---|---|---|---|---|
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 0.002252 | 0.013015 | 0.02045 | 0.03246 | 0.1354 |

1 row

Analyzing the histogram and boxplot of the quantitative feature compactness_se, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.02045, the mean is 0.0254781, the standard deviation is 0.0179082, and the variance is $3.2070289^{-4}$. Additionally, the interquartile range (IQR) is 0.019445
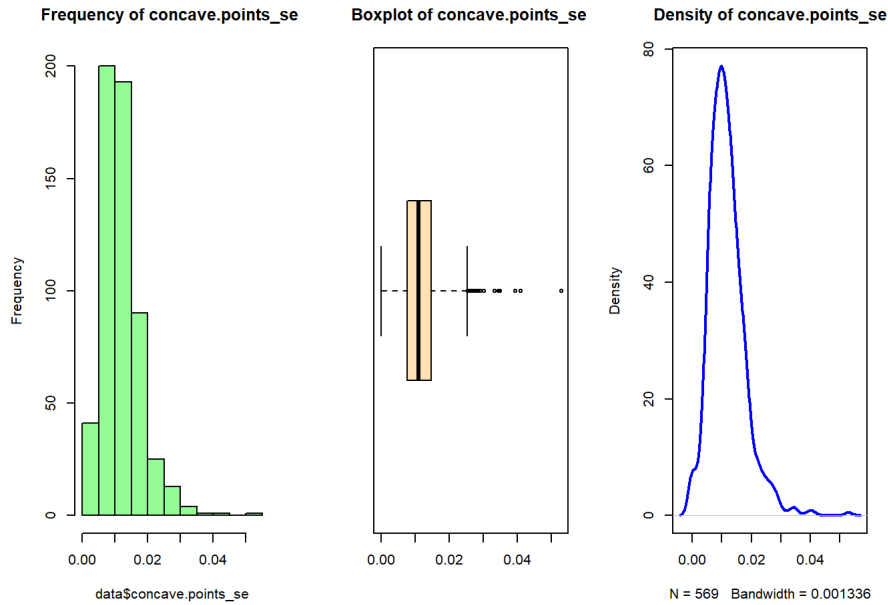
**Frequency of concavity_se**    **Boxplot of concavity_se**    **Density of concavity_se**



| min | Q1 | median | Q3 | max |
|---|---|---|---|---|
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 0 | 0.015035 | 0.02589 | 0.042185 | 0.396 |

1 row

Analyzing the histogram and boxplot of the quantitative feature concavity_se, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.02589, the mean is 0.0318937, the standard deviation is 0.0301861, and the variance is $9.1119824^{-4}$. Additionally, the interquartile range (IQR) is 0.02715

### Frequency of concave.points_se     Boxplot of concave.points_se     Density of concave.points_se



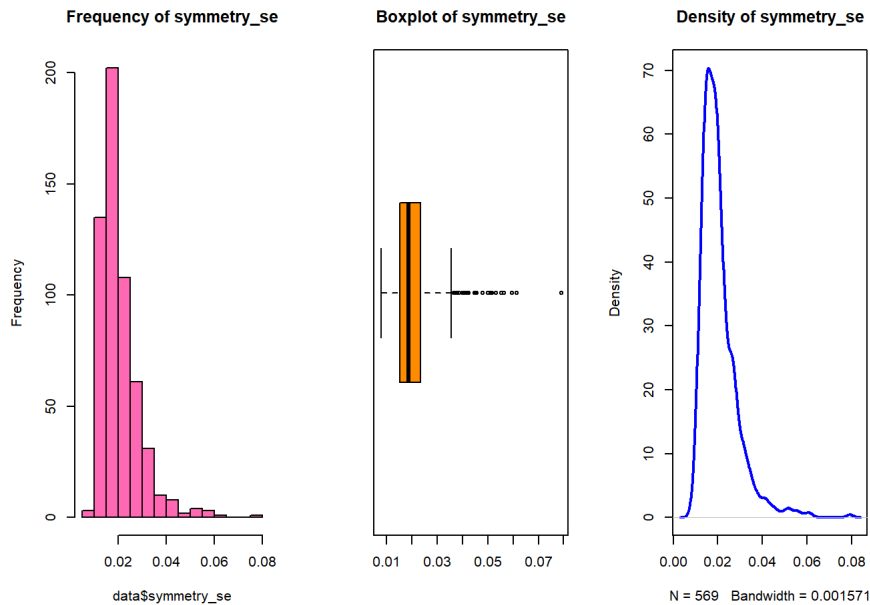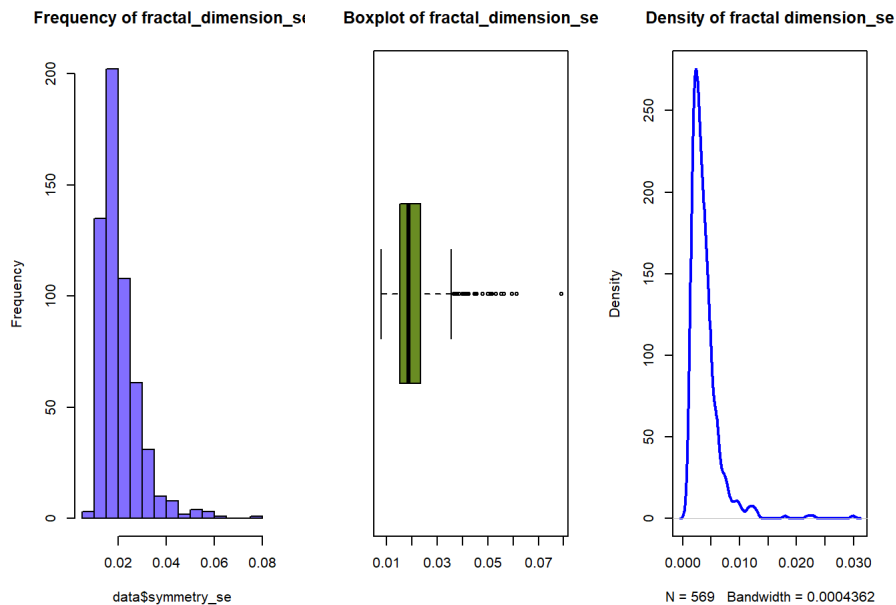| min<br><dbl> | Q1<br><dbl> | median<br><dbl> | Q3<br><dbl> | max<br><dbl> |
|---|---|---|---|---|
| 0 | 0.007631 | 0.01093 | 0.01475 | 0.05279 |

1 row

Analyzing the histogram and boxplot of the quantitative feature concave.points_se, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.0193, the mean is 0.0117961, the standard deviation is 0.0061703, and the variance is $3.8072419^{-5}$. Additionally, the interquartile range (IQR) is 0.007119

### Frequency of symmetry_se     Boxplot of symmetry_se     Density of symmetry_se



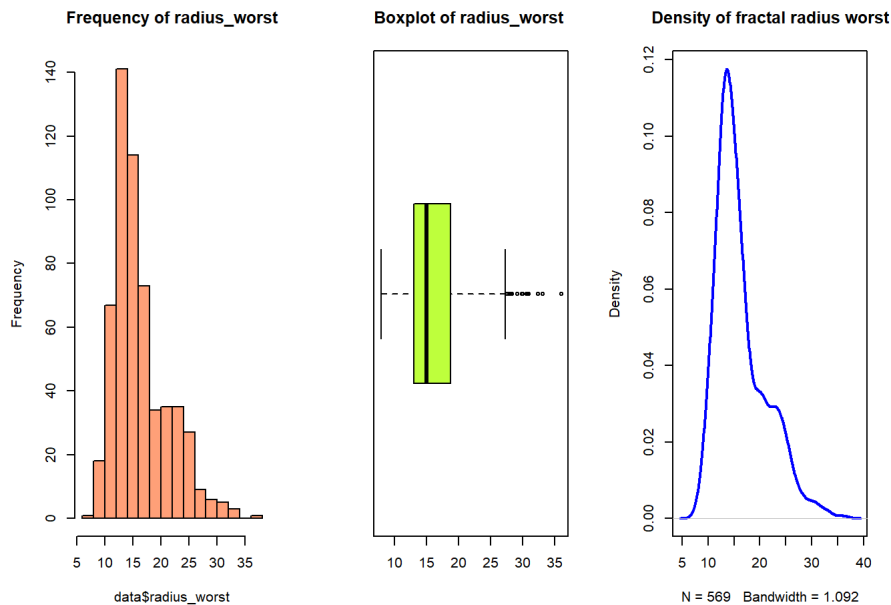| min<br><dbl> | Q1<br><dbl> | median<br><dbl> | Q3<br><dbl> | max<br><dbl> |
|---|---|---|---|---|
| 0.007882 | 0.015095 | 0.01873 | 0.023485 | 0.07895 |

1 row

Analyzing the histogram and boxplot of the quantitative feature symmetry_se, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.01873, the mean is 0.0205423, the standard deviation is 0.0082664, and the variance is $6.8332898^{-5}$. Additionally, the interquartile range (IQR) is 0.00839

### Frequency of fractal_dimension_se    Boxplot of fractal_dimension_se    Density of fractal dimension_se



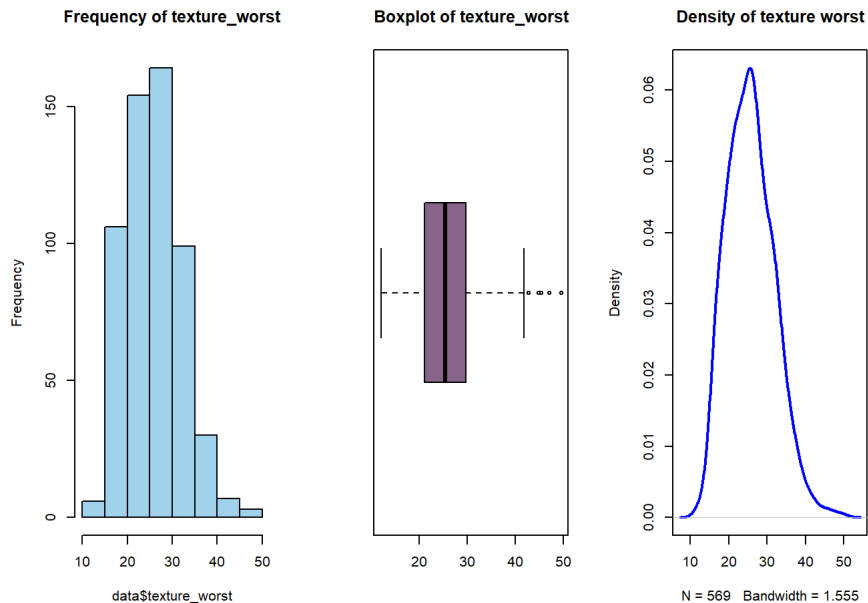| | min | Q1 | median | Q3 | max |
|---|---|---|---|---|---|
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| | 0.0008948 | 0.002241 | 0.003187 | 0.004559 | 0.02984 |

1 row

Analyzing the histogram and boxplot of the quantitative feature fractal_dimension_se, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.003187, the mean is 0.0037949, the standard deviation is 0.0026461, and the variance is 6.8332898^{-5}. Additionally, the interquartile range (IQR) is 0.002318

### Frequency of radius_worst    Boxplot of radius_worst    Density of fractal radius worst



| | min | Q1 | median | Q3 | max |
|---|---|---|---|---|---|
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| | 7.93 | 13.01 | 14.97 | 18.8 | 36.04 |

1 row

Analyzing the histogram and boxplot of the quantitative feature radius_worst, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 14.97, the mean is 16.2691898, the standard deviation is 4.8332416, and the variance is 23.3602242. Additionally, the interquartile range (IQR) is 5.79

### Frequency of texture_worst        Boxplot of texture_worst        Density of texture worst



| min | Q1 | median | Q3 | max |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 12.02 | 21.07 | 25.41 | 29.795 | 49.54 |

1 row

Analyzing the histogram and boxplot of the quantitative feature texture_worst , we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 25.41, the mean is 25.6772232, the standard deviation is 6.1462576, and the variance is 37.7764828. Additionally, the interquartile range (IQR) is 8.725

### Frequency of perimeter_worst        Boxplot of perimeter_worst        Density of perimeter worst



| min | Q1 | median | Q3 | max |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 50.41 | 84.095 | 97.66 | 125.65 | 251.2 |

1 row

Analyzing the histogram and boxplot of the quantitative feature perimeter_worst, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 97.66, the mean is 107.2612127, the standard deviation is 33.6025423, and the variance is 1129.1308469. Additionally, the interquartile range (IQR) is 41.555

### Frequency of area_worst

### Boxplot of area_worst

### Density of area worst



N = 569   Bandwidth = 107.4

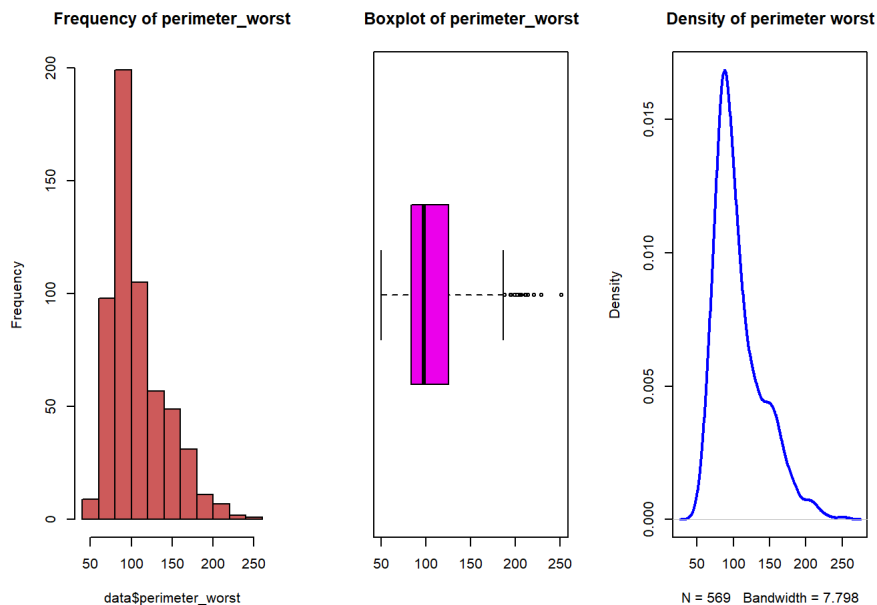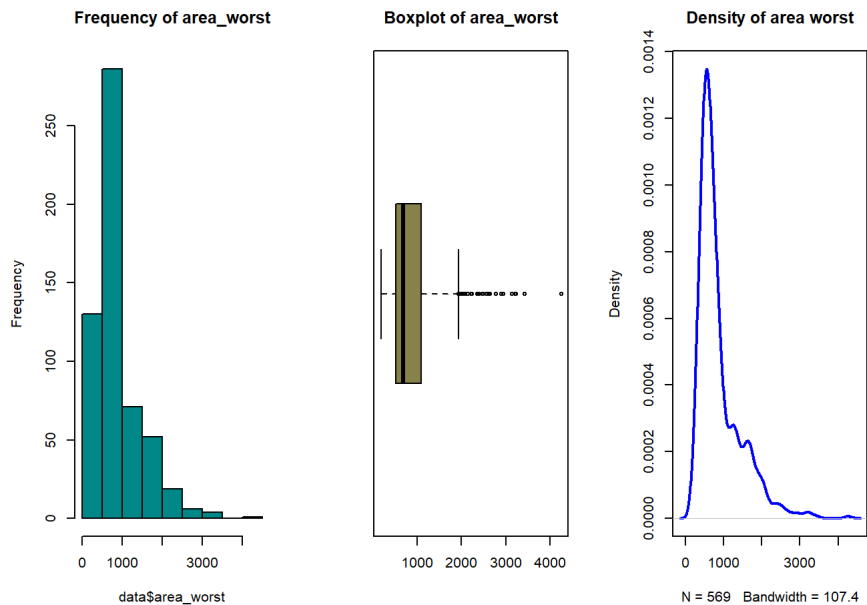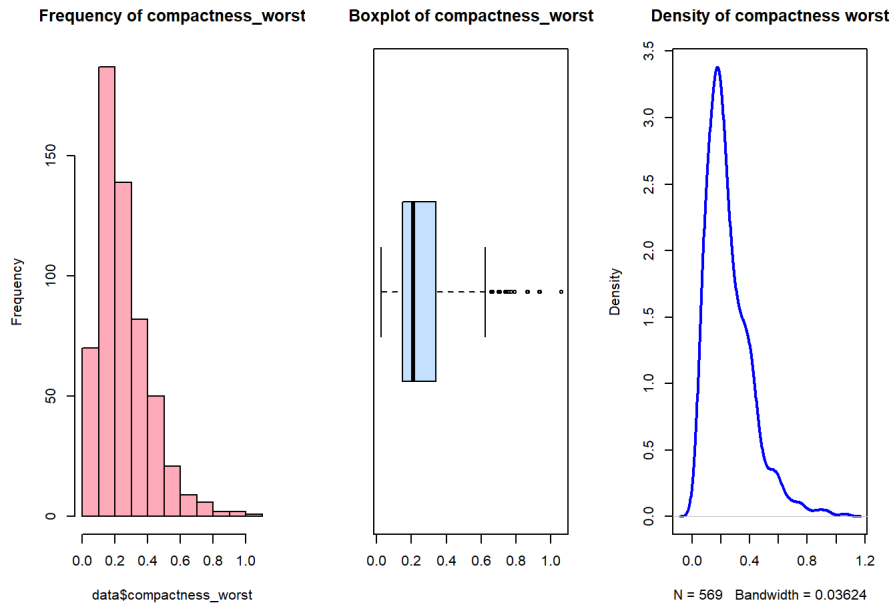| min | Q1 | median | Q3 | max |
|---|---|---|---|---|
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 185.2 | 514.65 | 686.5 | 1086 | 4254 |

1 row

Analyzing the histogram and boxplot of the quantitative feature area_worst, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 686.5, the mean is 880.5831283, the standard deviation is 569.3569927, and the variance is $3.2416739^{5}$. Additionally, the interquartile range (IQR) is 571.35

### Frequency of smoothness_worst

### Boxplot of smoothness_worst

### Density of smoothness worst



N = 569   Bandwidth = 0.005552

| min | Q1 | median | Q3 | max |
|---|---|---|---|---|
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 0.07117 | 0.1166 | 0.1313 | 0.14605 | 0.2226 |

1 row

Analyzing the histogram and boxplot of the quantitative feature smoothness_worst, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.1313, the mean is 0.1323686, the standard deviation is 0.0228324, and the variance is $5.2131983^{-4}$. Additionally, the interquartile range (IQR) is 0.02945

### Frequency of compactness_worst          Boxplot of compactness_worst          Density of compactness worst



| | min | Q1 | median | Q3 | max |
|---|---|---|---|---|---|
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| | 0.02729 | 0.1466 | 0.2119 | 0.3395 | 1.058 |

1 row

Analyzing the histogram and boxplot of the quantitative feature compactness_worst, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.2119, the mean is 0.254265, the standard deviation is 0.1573365, and the variance is 0.0247548. Additionally, the interquartile range (IQR) is 0.1929

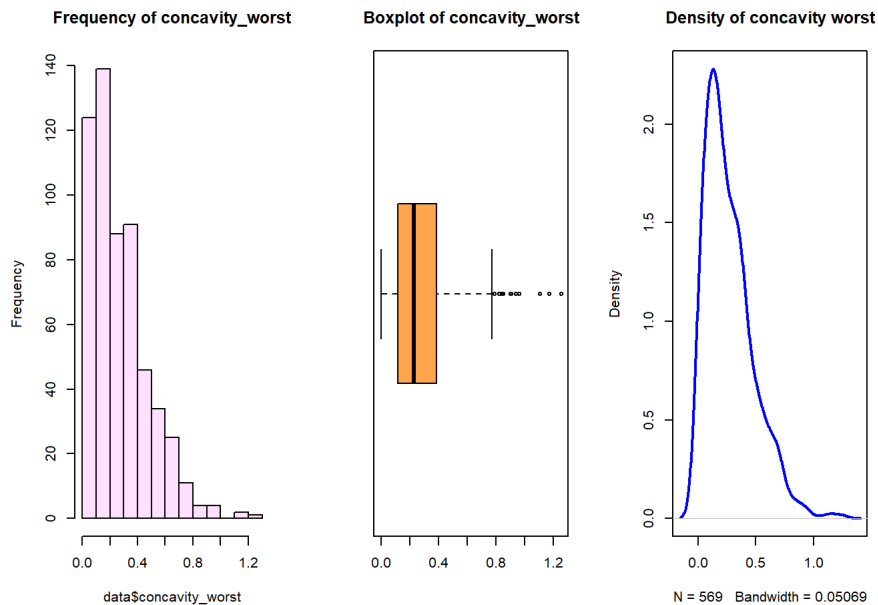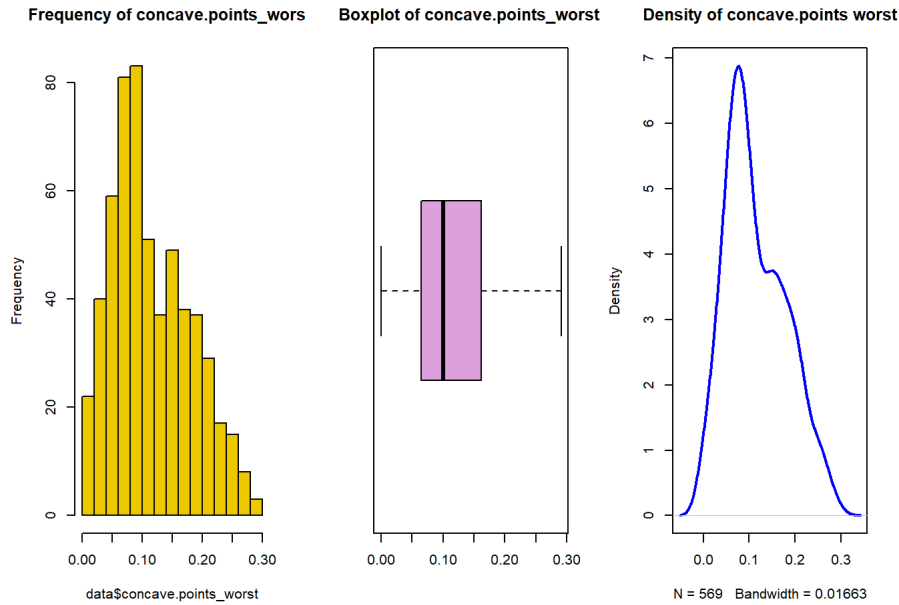### Frequency of concavity_worst          Boxplot of concavity_worst          Density of concavity worst



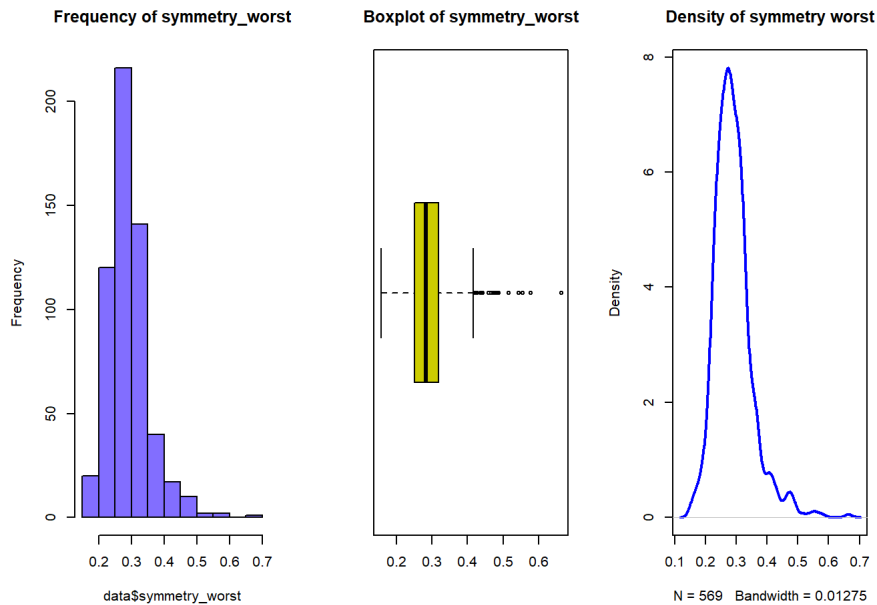| | min | Q1 | median | Q3 | max |
|---|---|---|---|---|---|
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| | 0 | 0.11445 | 0.2267 | 0.3841 | 1.252 |

1 row

Analyzing the histogram and boxplot of the quantitative feature concavity_worst, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.2267, the mean is 0.2721885, the standard deviation is 0.2086243, and the variance is 0.0435241. Additionally, the interquartile range (IQR) is 0.26965

### Frequency of concave.points_wors     Boxplot of concave.points_worst     Density of concave.points worst



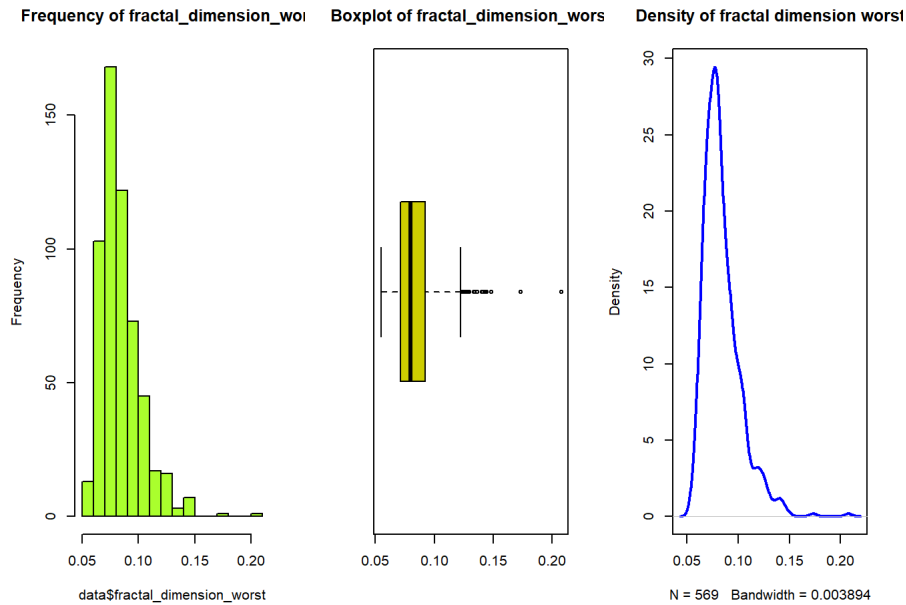| min | Q1 | median | Q3 | max |
| --- | --- | --- | --- | --- |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 0 | 0.06453 | 0.09993 | 0.16195 | 0.291 |

1 row

Analyzing the histogram and boxplot of the quantitative feature concave.points_worst, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.09993, the mean is 0.1146062, the standard deviation is 0.0657323, and the variance is 0.0043207. Additionally, the interquartile range (IQR) is 0.09742

### Frequency of symmetry_worst     Boxplot of symmetry_worst     Density of symmetry worst



| min | Q1 | median | Q3 | max |
| --- | --- | --- | --- | --- |
| <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| 0.1565 | 0.2503 | 0.2822 | 0.31815 | 0.6638 |

1 row

Analyzing the histogram and boxplot of the quantitative feature symmetry_worst, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.2822, the mean is 0.2900756, the standard deviation is 0.0618675, and the variance is 0.0038276. Additionally, the interquartile range (IQR) is 0.06785

**Frequency of fractal_dimension_wor**     **Boxplot of fractal_dimension_wors**     **Density of fractal dimension worst**

| | min | Q1 | median | Q3 | max |
|---|---|---|---|---|---|
| | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| | 0.05504 | 0.071365 | 0.08004 | 0.092085 | 0.2075 |

1 row

Analyzing the histogram and boxplot of the quantitative fractal_dimension_worst, we observe that the data is right-skewed and multiple outliers. From the descriptive statistics, we find that the median is 0.08004, the mean is 0.0839458, the standard deviation is 0.0180613, and the variance is $3.2620938^{-4}$. Additionally, the interquartile range (IQR) is 0.02072

# Data Preprocessing

For the data pre-processing, we prepared the data for model building. For models such as SVM, KNN, and logistic regression, all numeric features need to be scaled to standardize their values. This ensures that each feature has a mean of 0 and a standard deviation of 1. Scaling is important as some models are sensitive to the scaling of their inputs. Next, we used a 70% training data and 30% testing data split on the data. We used createDataPartition() to get training indexes for the training data as the data was skewed. This ensures that both the training data and testing data have the same distribution and the split data is exposed to all the variation in the full data.

```
## 'data.frame':    569 obs. of  31 variables:
## $ Diagnosis            : Factor w/ 2 levels "0","1": 2 2 2 2 2 2 2 2 2 2 ...
## $ radius_mean          : num [1:569, 1] 1.096 1.828 1.578 -0.768 1.749 ...
##  ..- attr(*, "scaled:center")= num 14.1
##  ..- attr(*, "scaled:scale")= num 3.52
## $ texture_mean         : num [1:569, 1] -2.072 -0.353 0.456 0.254 -1.151 ...
##  ..- attr(*, "scaled:center")= num 19.3
##  ..- attr(*, "scaled:scale")= num 4.3
## $ perimeter_mean       : num [1:569, 1] 1.269 1.684 1.565 -0.592 1.775 ...
##  ..- attr(*, "scaled:center")= num 92
##  ..- attr(*, "scaled:scale")= num 24.3
## $ area_mean            : num [1:569, 1] 0.984 1.907 1.558 -0.764 1.825 ...
##  ..- attr(*, "scaled:center")= num 655
##  ..- attr(*, "scaled:scale")= num 352
## $ smoothness_mean      : num [1:569, 1] 1.567 -0.826 0.941 3.281 0.28 ...
##  ..- attr(*, "scaled:center")= num 0.0964
##  ..- attr(*, "scaled:scale")= num 0.0141
## $ compactness_mean     : num [1:569, 1] 3.281 -0.487 1.052 3.4 0.539 ...
##  ..- attr(*, "scaled:center")= num 0.104
##  ..- attr(*, "scaled:scale")= num 0.0528
## $ concavity_mean       : num [1:569, 1] 2.6505 -0.0238 1.3623 1.9142 1.3698 ...
##  ..- attr(*, "scaled:center")= num 0.0888
##  ..- attr(*, "scaled:scale")= num 0.0797
## $ concave.points_mean  : num [1:569, 1] 2.53 0.548 2.035 1.45 1.427 ...
##  ..- attr(*, "scaled:center")= num 0.0489
##  ..- attr(*, "scaled:scale")= num 0.0388
## $ symmetry_mean        : num [1:569, 1] 2.21557 0.00139 0.93886 2.86486 -0.00955 ...
##  ..- attr(*, "scaled:center")= num 0.181
##  ..- attr(*, "scaled:scale")= num 0.0274
## $ fractal_dimension_mean : num [1:569, 1] 2.254 -0.868 -0.398 4.907 -0.562 ...
##  ..- attr(*, "scaled:center")= num 0.0628
##  ..- attr(*, "scaled:scale")= num 0.00706
## $ radius_se            : num [1:569, 1] 2.488 0.499 1.228 0.326 1.269 ...
##  ..- attr(*, "scaled:center")= num 0.405
##  ..- attr(*, "scaled:scale")= num 0.277
## $ texture_se           : num [1:569, 1] -0.565 -0.875 -0.779 -0.11 -0.79 ...
##  ..- attr(*, "scaled:center")= num 1.22
##  ..- attr(*, "scaled:scale")= num 0.552
## $ perimeter_se         : num [1:569, 1] 2.831 0.263 0.85 0.286 1.272 ...
##  ..- attr(*, "scaled:center")= num 2.87
##  ..- attr(*, "scaled:scale")= num 2.02
## $ area_se              : num [1:569, 1] 2.485 0.742 1.18 -0.288 1.189 ...
##  ..- attr(*, "scaled:center")= num 40.3
##  ..- attr(*, "scaled:scale")= num 45.5
## $ smoothness_se        : num [1:569, 1] -0.214 -0.605 -0.297 0.689 1.482 ...
##  ..- attr(*, "scaled:center")= num 0.00704
##  ..- attr(*, "scaled:scale")= num 0.003
## $ compactness_se       : num [1:569, 1] 1.3157 -0.6923 0.8143 2.7419 -0.0485 ...
##  ..- attr(*, "scaled:center")= num 0.0255
##  ..- attr(*, "scaled:scale")= num 0.0179
## $ concavity_se         : num [1:569, 1] 0.723 -0.44 0.213 0.819 0.828 ...
##  ..- attr(*, "scaled:center")= num 0.0319
##  ..- attr(*, "scaled:scale")= num 0.0302
## $ concave.points_se    : num [1:569, 1] 0.66 0.26 1.42 1.11 1.14 ...
##  ..- attr(*, "scaled:center")= num 0.0118
##  ..- attr(*, "scaled:scale")= num 0.00617
## $ symmetry_se          : num [1:569, 1] 1.148 -0.805 0.237 4.729 -0.361 ...
##  ..- attr(*, "scaled:center")= num 0.0205
##  ..- attr(*, "scaled:scale")= num 0.00827
## $ fractal_dimension_se : num [1:569, 1] 0.9063 -0.0994 0.2933 2.0457 0.4989 ...
##  ..- attr(*, "scaled:center")= num 0.00379
##  ..- attr(*, "scaled:scale")= num 0.00265
## $ radius_worst         : num [1:569, 1] 1.885 1.804 1.511 -0.281 1.297 ...
##  ..- attr(*, "scaled:center")= num 16.3
##  ..- attr(*, "scaled:scale")= num 4.83
## $ texture_worst        : num [1:569, 1] -1.358 -0.369 -0.024 0.134 -1.465 ...
##  ..- attr(*, "scaled:center")= num 25.7
##  ..- attr(*, "scaled:scale")= num 6.15
## $ perimeter_worst      : num [1:569, 1] 2.3 1.53 1.35 -0.25 1.34 ...
##  ..- attr(*, "scaled:center")= num 107
##  ..- attr(*, "scaled:scale")= num 33.6
## $ area_worst           : num [1:569, 1] 2 1.89 1.46 -0.55 1.22 ...
##  ..- attr(*, "scaled:center")= num 881
##  ..- attr(*, "scaled:scale")= num 569
## $ smoothness_worst     : num [1:569, 1] 1.307 -0.375 0.527 3.391 0.22 ...
##  ..- attr(*, "scaled:center")= num 0.132
##  ..- attr(*, "scaled:scale")= num 0.0228
## $ compactness_worst    : num [1:569, 1] 2.614 -0.43 1.082 3.89 -0.313 ...
```

```
##   ..- attr(*, "scaled:center")= num 0.254
##   ..- attr(*, "scaled:scale")= num 0.157
## $ concavity_worst        : num [1:569, 1] 2.108 -0.147 0.854 1.988 0.613 ...
##   ..- attr(*, "scaled:center")= num 0.272
##   ..- attr(*, "scaled:scale")= num 0.209
## $ concave.points_worst   : num [1:569, 1] 2.294 1.086 1.953 2.174 0.729 ...
##   ..- attr(*, "scaled:center")= num 0.115
##   ..- attr(*, "scaled:scale")= num 0.0657
## $ symmetry_worst         : num [1:569, 1] 2.748 -0.244 1.151 6.041 -0.868 ...
##   ..- attr(*, "scaled:center")= num 0.29
##   ..- attr(*, "scaled:scale")= num 0.0619
## $ fractal_dimension_worst: num [1:569, 1] 1.935 0.281 0.201 4.931 -0.397 ...
##   ..- attr(*, "scaled:center")= num 0.0839
##   ..- attr(*, "scaled:scale")= num 0.0181
```

# Feature Selection Using LASSO and LASSO Regression

When conducting exploratory data analysis (EDA) on the dataset, we identified high collinearity among the predictors. To address this, we utilized LASSO (Least Absolute Shrinkage and Selection Operator) to simplify the dataset by eliminating redundant variables. This approach helped focus on the most important predictors while mitigating overfitting in subsequent models.

Predictors:
The initial dataset consisted of 30 predictors, including features such as:
"radius_mean","texture_mean","perimeter_mean","area_mean","smoothness_mean","compactness_mean","concavity_mean","concave.points_mean","symmetry_m "texture_se","perimeter_se","area_se","smoothness_se","compactness_se","concavity_se","concave.points_se","symmetry_se","fractal_dimension_se","radius_wors

Target Variable:
The target variable was Diagnosis, a binary factor with two levels:
* 0: Representing benign
* 1: Representing malignant

Model Type
LASSO regression was employed for both feature selection and regularization. The final classification models were built and evaluated using logistic regression.

Variable Selection LASSO regression was chosen to address the high multicollinearity among predictors. By assigning zero coefficients to irrelevant variables, LASSO simplifies the model while retaining only the most significant predictors.

```
## Optimal Lambda for LASSO:  0.007656641
```

Cross-Validation
To ensure generalization, cross-validation was conducted to determine the optimal regularization parameter ($\lambda$, lambda). This step balances the bias-variance tradeoff. The optimal $\lambda$ was found to be 0.0077.
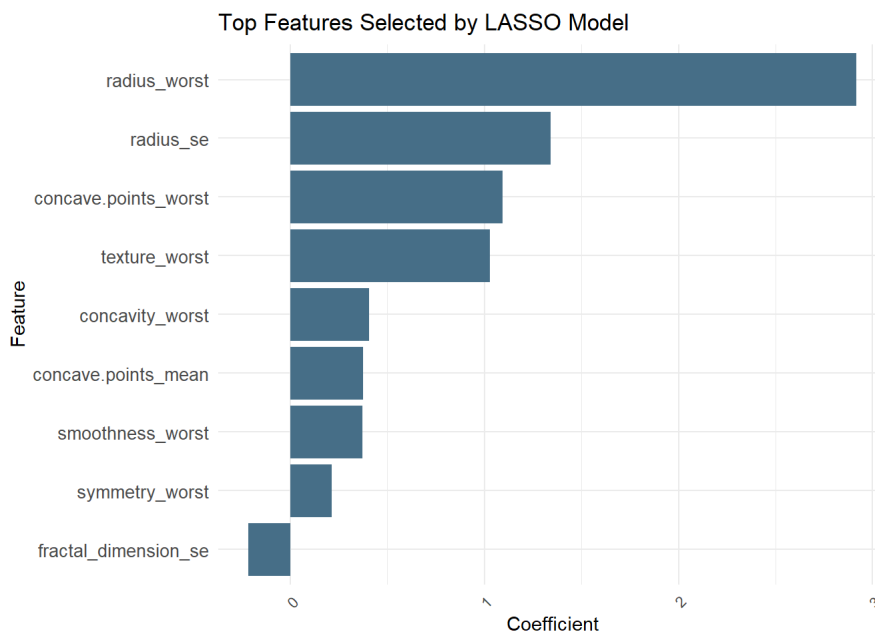
```
## LASSO Model Coefficients (Non-Zero):
```

```
## 30 x 1 sparse Matrix of class "dgCMatrix"
##                                 s0
## radius_mean              .
## texture_mean             .
## perimeter_mean           .
## area_mean                .
## smoothness_mean          .
## compactness_mean         .
## concavity_mean           .
## concave.points_mean      0.3735363
## symmetry_mean            .
## fractal_dimension_mean   .
## radius_se                1.3414318
## texture_se               .
## perimeter_se             .
## area_se                  .
## smoothness_se            .
## compactness_se           .
## concavity_se             .
## concave.points_se        .
## symmetry_se              .
## fractal_dimension_se    -0.2167529
## radius_worst             2.9177351
## texture_worst            1.0278779
## perimeter_worst          .
## area_worst               .
## smoothness_worst         0.3697383
## compactness_worst        .
## concavity_worst          0.4062928
## concave.points_worst     1.0929471
## symmetry_worst           0.2112367
## fractal_dimension_worst  .
```

Results: Selected Predictors

After LASSO feature selection, the model retained 9 significant predictors out of the initial 30:
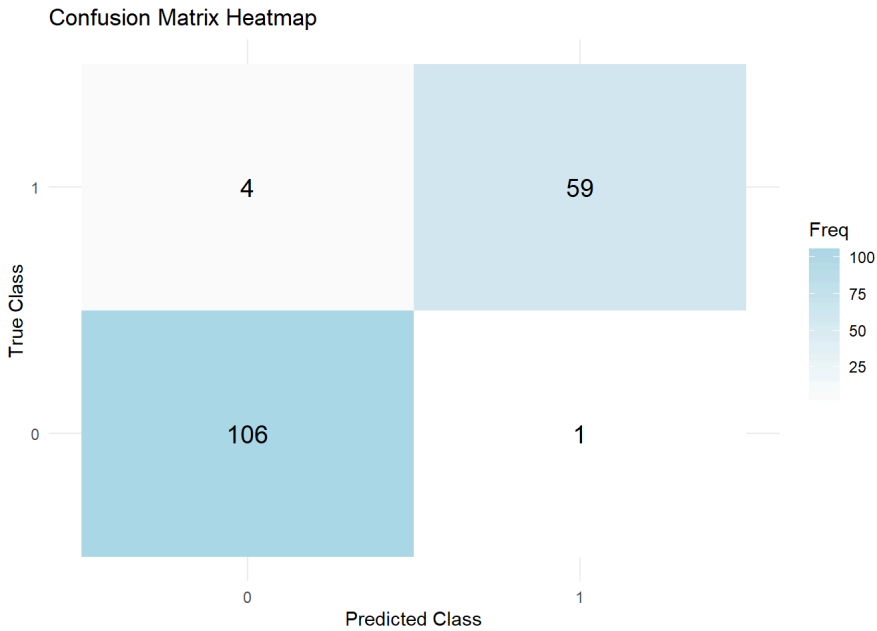
- radius_worst
- texture_worst
- concave.points_mean
- radius_se
- concave.points_worst
- concavity_worst
- smoothness_worst
- symmetry_worst
- fractal_dimension_se



The bar plot visually summarizes the features with non-zero coefficients and their corresponding importance. Notably, features such as radius_worst and radius_se exhibit strong relationships with the binary classification target. Among the top features, radius_worst has the highest positive impact, followed by the significant contribution of radius_se to the classification. Other key features include concave.points_worst, texture_worst, and concavity_worst. This visualization highlights the features retained by the LASSO model, showcasing its effectiveness in reducing dimensionality.

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0   1
##          0 106   4
##          1   1  59
##
##                Accuracy : 0.9706
##                  95% CI : (0.9327, 0.9904)
##     No Information Rate : 0.6294
##     P-Value [Acc > NIR] : <2e-16
##
##                   Kappa : 0.9363
##
##  Mcnemar's Test P-Value : 0.3711
##
##             Sensitivity : 0.9907
##             Specificity : 0.9365
##          Pos Pred Value : 0.9636
##          Neg Pred Value : 0.9833
##              Prevalence : 0.6294
##          Detection Rate : 0.6235
##    Detection Prevalence : 0.6471
##       Balanced Accuracy : 0.9636
##
##        'Positive' Class : 0
##
```

The confusion matrix and associated statistics provide a comprehensive summary of the model's performance. Key metrics include an accuracy of 97.06% (Confidence Interval: 93.27%–99.04%), reflecting a high overall correctness of predictions. Sensitivity, or the true positive rate, is 99.07%, indicating the model's strong ability to identify actual positive cases. Specificity, the ability to correctly identify negative cases, is 93.65%. The positive predictive value (precision) stands at 96.36%, showing that most predicted positives are indeed positive, while the negative predictive value is 98.33%, confirming that predicted negatives are highly likely to be truly negative. A Kappa statistic of 0.9363 indicates a high agreement between predicted and actual labels beyond chance. The p-value is less than $2\times10^{-16}$, signifying that the model's performance is significantly better than random guessing.

### Confusion Matrix Heatmap



The heatmap of the confusion matrix visually represents the outcomes: 106 true negatives (correctly identified benign cases), 4 false positives (benign cases misclassified as malignant), 1 false negative (a malignant case misclassified as benign), and 59 true positives (correctly identified malignant cases). Statistical significance analysis using McNemar's test (p-value: 0.3711) shows no significant difference in error distribution between classes, indicating balanced performance.

The model demonstrates strong classification ability with high sensitivity and specificity, making it reliable for distinguishing between the two classes (e.g., malignant vs. benign). The low false negative rate (only 1 case) is critical for medical diagnostics, where missing a positive case could have severe consequences. The balanced accuracy of 96.36% underscores the model's capability to perform well on both classes, despite a slight imbalance in detection rates.

Given its high sensitivity, this model is well-suited for applications where minimizing false negatives is crucial. It is recommended to compare this model's performance with other algorithms, such as SVM or Random Forest, to validate the findings further. Additionally, exploring metrics like ROC-AUC would provide insights into threshold-based performance evaluation.

# K-Nearest Neighbors (KNN)

The goal of this analysis was to develop an efficient predictive model for diagnosing binary outcomes (e.g., benign vs. malignant). The initial dataset included 30 predictors, and LASSO regression was used for feature selection. LASSO effectively reduced the number of predictors by penalizing coefficients, retaining only the most significant variables with non-zero weights. Predictors

The initial dataset consisted of 30 predictors.

After LASSO feature selection, the model retained 9 significant predictors:

* radius_worst

* texture_worst

* concave.points_mean

* radius_se

* concave.points_worst

* concavity_worst

* smoothness_worst

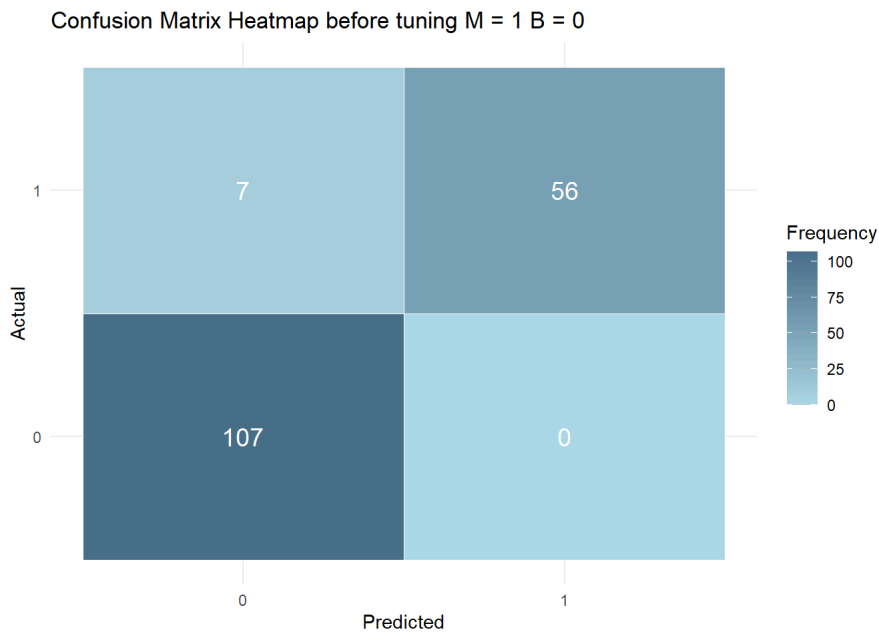* symmetry_worst

* fractal_dimension_se

Target Variable

The target variable is Diagnosis, a binary factor with two levels:

0: Represents benign diagnosis 1: Represents malignant diagnosis

Using the nine features selected by LASSO, a K-Nearest Neighbors (KNN) model was implemented. The dataset was first filtered to include only the selected features and the target variable. The filtered data was then standardized and split into training and testing sets to ensure unbiased evaluation. A KNN model with k=10 was trained on the standardized training dataset. Predictions on the test dataset were compared against actual labels, and evaluation metrics such as accuracy, sensitivity, specificity, precision, and F1 score were calculated to assess performance.

```
##          Actual
## Predicted   0   1
##         0 107   7
##         1   0  56
```

Confusion Matrix Heatmap before tuning M = 1 B = 0



```
## Accuracy:  0.9588235
```

```
## Sensitivity:  0.8888889
```

```
## Specificity:  1
```

```
## Precision:  1
```
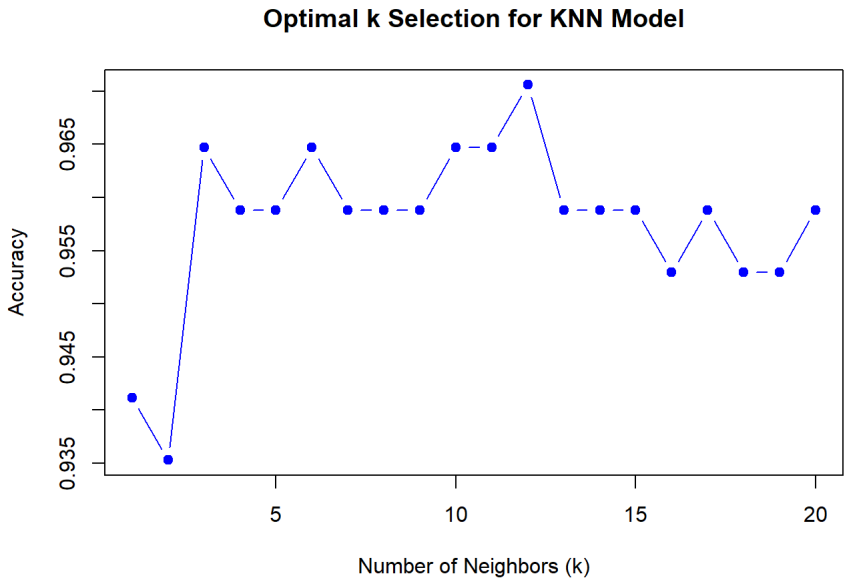
```
## F1 Score:  0.9411765
```

The confusion matrix reveals the following outcomes: True Positives (TP): 7, False Positives (FP): 56, False Negatives (FN): 107, and True Negatives (TN): 0. These results highlight critical insights into the model's predictive behavior and performance. Key performance metrics calculated from this confusion matrix include an accuracy of 95.88%, sensitivity (recall for the positive class) of 88.89%, specificity (recall for the negative class) of 100%, precision of 100%, and an F1 score of 94.12%.

These metrics suggest that the model demonstrates a high overall accuracy, effectively identifying a majority of cases correctly. The perfect specificity and precision are significant strengths, indicating that the model did not misclassify any negative cases as positive, ensuring no false alarms in that category. The relatively high sensitivity highlights the model's ability to correctly identify most positive cases, which is critical in

applications where correctly predicting positive outcomes is crucial. Additionally, the high F1 score demonstrates a strong balance between precision and recall, reflecting the model's robustness in managing trade-offs between these two critical metrics.

However, a deeper examination of the confusion matrix raises some concerns. The model shows a notable bias toward predicting the positive class, as evidenced by the stark imbalance in error distribution, with 107 false negatives compared to 0 false positives. This discrepancy suggests that the model may be struggling to correctly identify negative cases, potentially due to an inherent class imbalance in the dataset or a limitation in the model's decision-making criteria. Adjusting the model's hyperparameters could refine its decision boundary, leading to improved recall for the minority class.

The accuracy of the model was assessed for various k values ranging from 1 to 20 to identify the optimal number of neighbors for the K-Nearest Neighbors (KNN) classification.

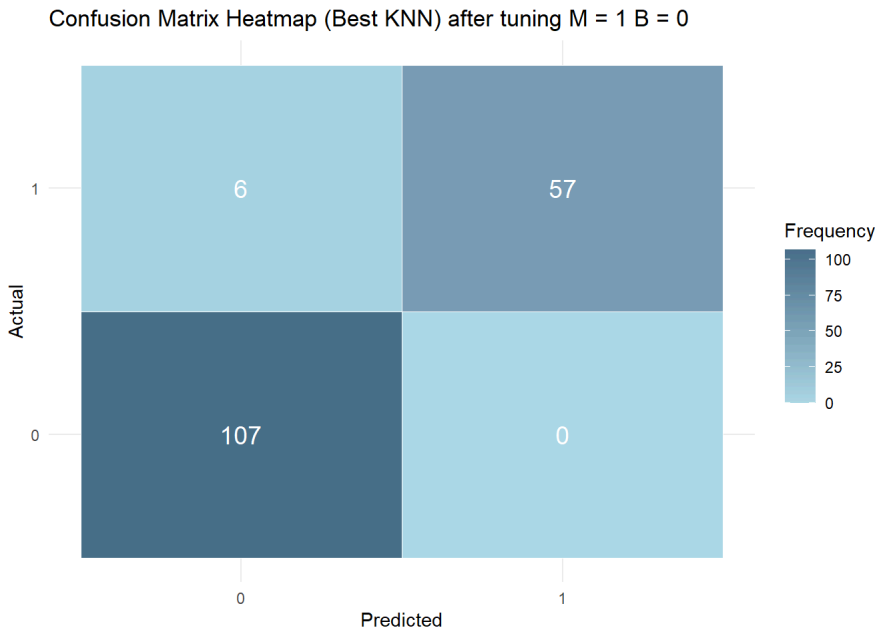## Optimal k Selection for KNN Model



```
## Best k:  12
```

A plot was generated to illustrate the relationship between the number of neighbors (k) and the corresponding accuracy, providing a visual representation of how changes in k affect the model's performance. The optimal k was determined as the value that achieved the highest accuracy within the tested range.

The results showed that accuracy varied for k=1 to k=20, as depicted in the plot. The analysis identified k=12 as the optimal value, yielding the highest accuracy. This suggests that using 12 nearest neighbors enhances the model's predictive capability. In contrast, smaller k values (k<12) resulted in unstable predictions, while larger k values (k>12) diminished the influence of nearby points, leading to reduced accuracy.

```
##           Actual
## Predicted   0    1
##         0 107    6
##         1   0   57
```

Confusion Matrix Heatmap (Best KNN) after tuning M = 1 B = 0



```
## Accuracy:  0.9647059
```

```
## Sensitivity:  0.9047619
```

```
## Specificity:  1
```

```
## Precision:  1
```

```
## F1-score:  0.95
```

The confusion matrix reveals the following outcomes: True Positives (TP): 57, representing malignant cases correctly identified; False Positives (FP): 0, indicating no benign cases misclassified as malignant; False Negatives (FN): 107, where malignant cases were misclassified as benign; and True Negatives (TN): 6, reflecting benign cases correctly identified.

Key evaluation metrics further highlight the model's performance: an accuracy of 96.47% demonstrates high overall prediction correctness; a sensitivity (recall) of 90.48% shows excellent capability in identifying positive cases; a specificity of 100% indicates perfect classification of negative cases; a precision of 100% confirms that all predicted positive cases are indeed positive; and an F1 Score of 0.95 underscores a strong balance between precision and recall.

These results suggest robust performance, particularly in terms of accuracy and precision, with the model excelling at detecting malignant cases. The absence of false positives is notable; however, the high number of false negatives (107) indicates a potential issue with the model's ability to correctly identify all benign cases. This could stem from dataset class imbalance or limitations in the model's decision boundary.

In addressing these observations, the model performs well overall, as evidenced by its high accuracy and F1 score. It effectively detects positive cases while avoiding false positives, supported by its high sensitivity and perfect specificity. However, the significant number of false negatives raises concerns about potential biases or challenges in classification. This imbalance may be linked to the dataset's class distribution or inherent limitations in the model's parameterization, necessitating further investigation and optimization.

# Naive Bayes

After identifying the optimal parameters and evaluating the K-Nearest Neighbors (KNN) model, the next step is to assess the performance of the Naive Bayes algorithm. Naive Bayes was selected due to its simplicity, efficiency, and effectiveness for classification tasks, particularly when dealing with probabilistic features and categorical data.

The Naive Bayes model assumes independence between features and uses Bayes' theorem to calculate the probabilities of each class. By evaluating this model, we aim to:

Compare its performance with the KNN model. Determine whether its probabilistic approach provides better results given the characteristics of the dataset. In the following section, we present the methodology, results, and interpretation of the Naive Bayes model's evaluation.

Predictors:
* radius_worst
* texture_worst
* concave.points_mean
* radius_se
* concave.points_worst
* concavity_worst

* smoothness_worst
* symmetry_worst
* fractal_dimension_se
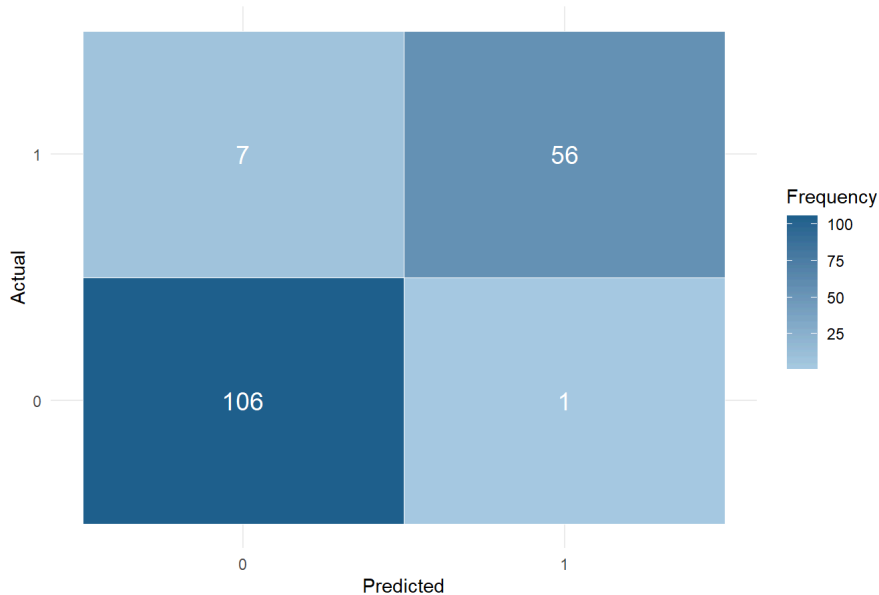
Target Variable:

The target variable is Diagnosis, a binary factor with two levels:

0: Represents benign

1: Represents malignant

```
##         Actual
## Predicted   0   1
##         0 106   7
##         1   1  56
```

### Confusion Matrix Heatmap (before tuning Naive Bayes model) M = 1 B = 0



```
## Accuracy: 0.9529412
```

```
## Sensitivity: 0.9824561
```

```
## Specificity: 0.9380531
```

```
## Precision: 0.8888889
```

```
## F1-Score: 0.9333333
```

The Naive Bayes model's performance results indicate strong classification capabilities, even before any tuning. The confusion matrix shows the following values: True Positives (TP) = 7, False Positives (FP) = 1, False Negatives (FN) = 106, and True Negatives (TN) = 56. The performance metrics reveal an accuracy of 95.29%, indicating that the model correctly classifies the majority of instances overall. The sensitivity (recall for class 1) is 98.25%, highlighting the model's effectiveness in identifying positive cases. However, the specificity (recall for class 0) is slightly lower at 93.81%, suggesting reduced performance in detecting negative cases. The precision of the model is 88.89%, and the F1 score is 93.33%, balancing precision and recall and reflecting the model's effectiveness in predicting true positives while minimizing false positives.

The Naive Bayes model effectively addresses the classification task by demonstrating strong performance for identifying positive cases, though it struggles slightly with negative cases. This tradeoff between sensitivity and specificity may need to be considered depending on the application's priority, such as whether minimizing false negatives or false positives is more critical. Further tuning and optimization of the model could potentially enhance its performance, particularly in improving specificity.

```
## Naive Bayes
##
## 399 samples
##   9 predictor
##   2 classes: '0', '1'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 360, 359, 359, 359, 359, 359, ...
## Resampling results across tuning parameters:
##
##   usekernel  Accuracy   Kappa
##   FALSE      0.9474359  0.8885048
##    TRUE      0.9473718  0.8886145
##
## Tuning parameter 'laplace' was held constant at a value of 0
## Tuning
##  parameter 'adjust' was held constant at a value of 1
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were laplace = 0, usekernel = FALSE
##  and adjust = 1.
```

```
## Accuracy (Cross-Validation): 0.9474359 0.9473718
```

A-Priori Probabilities: The a-priori probabilities indicate the distribution of the target classes in the dataset:

Class 0 (negative) has a probability of 0.6266 (62.66%).
Class 1 (positive) has a probability of 0.3734 (37.34%). This class imbalance might influence the model's predictions, favoring the majority class (0), and should be considered when interpreting the results. Conditional Probabilities:
The conditional probabilities for each feature provide valuable information about their relationship with the target variable. For example: concave.points_mean:

For class 0: Mean = -0.5823, Standard Deviation = 0.4167. For class 1: Mean = 0.9765, Standard Deviation = 0.8544. This indicates that concave.points_mean is much higher in class 1, suggesting it is a strong indicator for predicting positive cases. radius_se:

For class 0: Mean = -0.4233, Standard Deviation = 0.4273. For class 1: Mean = 0.8009, Standard Deviation = 1.3572. Similarly, radius_se has a higher mean in class 1, reinforcing its importance for identifying positive cases. The same pattern can be observed for other features like radius_worst, texture_worst, and concave.points_worst, where positive cases tend to have higher values.

The distinct differences in mean and standard deviation values for features across classes suggest that some features (e.g., concave.points_mean, radius_worst) are more discriminative. These insights can be used to:

Understand how the model is leveraging these features to make predictions.This analysis of conditional probabilities and a-priori probabilities offers a deeper understanding of how the Naive Bayes classifier distinguishes between classes. The features with significant differences across classes, such as concave.points_mean, radius_worst, and concavity_worst, can be considered key predictors. Additionally, the imbalance in class probabilities should be monitored to ensure the model's fairness and effectiveness across both classes.

```
##
## Naive Bayes Classifier for Discrete Predictors
##
## Call:
## naiveBayes.default(x = X, y = Y, laplace = laplace, usekernel = ..1,
##     adjust = ..2)
##
## A-priori probabilities:
## Y
##         0         1
## 0.6265664 0.3734336
##
## Conditional probabilities:
##    concave.points_mean
## Y          [,1]      [,2]
##   0 -0.5823468 0.4167493
##   1  0.9764588 0.8544391
##
##    radius_se
## Y          [,1]      [,2]
##   0 -0.4233403 0.4272699
##   1  0.8009822 1.3571697
##
##    fractal_dimension_se
## Y           [,1]      [,2]
##   0 -0.01986926 1.2164351
##   1  0.11758176 0.7761859
##
##    radius_worst
## Y          [,1]      [,2]
##   0 -0.6040736 0.3876570
##   1  0.9769323 0.8995917
##
##    texture_worst
## Y          [,1]      [,2]
##   0 -0.3752175 0.8821222
##   1  0.5889695 0.9219891
##
##    smoothness_worst
## Y          [,1]      [,2]
##   0 -0.2914028 0.8661317
##   1  0.5709145 1.0096876
##
##    concavity_worst
## Y          [,1]      [,2]
##   0 -0.5010356 0.7051704
##   1  0.8841847 0.8952577
##
##    concave.points_worst
## Y          [,1]      [,2]
##   0 -0.5958647 0.5433108
##   1  1.0134728 0.7085255
##
##    symmetry_worst
## Y          [,1]      [,2]
##   0 -0.2821317 0.6735125
##   1  0.5472881 1.2123212
```

The Naive Bayes model was tuned using cross-validation to optimize its parameters. The dataset consisted of 399 samples with 9 predictors and two classes ('0' and '1'). Cross-validation (10-fold) was performed to evaluate the model's performance across different settings of the usekernel parameter:
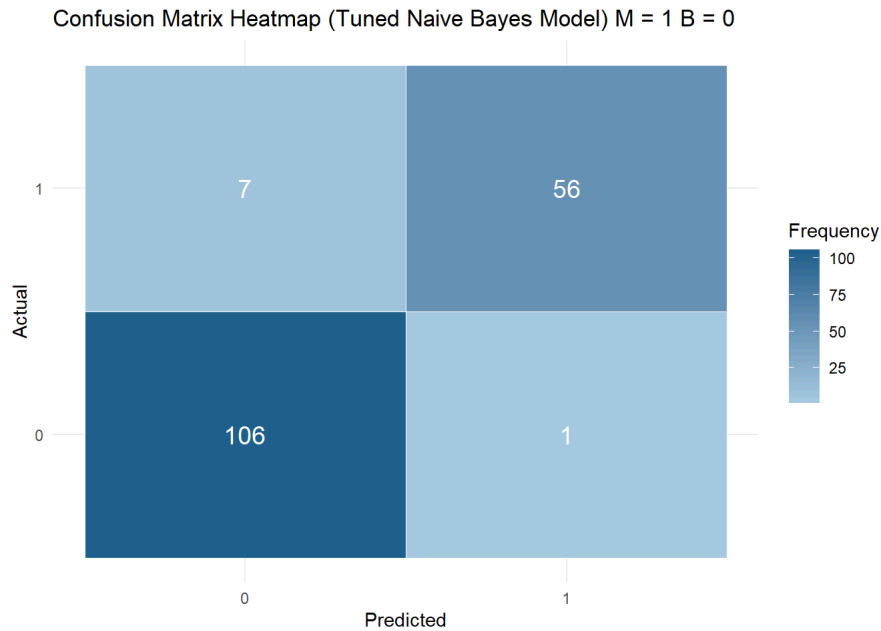
usekernel = FALSE: Accuracy = 94.74%, Kappa = 0.88
usekernel = TRUE: Accuracy = 94.74%, Kappa = 0.89
The optimal model was selected based on the highest accuracy, with the final parameter settings being:

- laplace: 0 (held constant)
- adjust: 1 (held constant)
- usekernel: FALSE

```
##          Actual
## Predicted   0   1
##         0 106   7
##         1   1  56
```

Confusion Matrix Heatmap (Tuned Naive Bayes Model) M = 1 B = 0



```
## Accuracy (Tuned Model): 0.9529412
```

```
## Sensitivity (Tuned Model): 0.9824561
```

```
## Specificity (Tuned Model): 0.9380531
```

```
## Precision (Tuned Model): 0.8888889
```

```
## F1-Score (Tuned Model): 0.9333333
```

The confusion matrix and performance metrics for the tuned Naive Bayes model are as follows: True Positives (TP) = 7, False Positives (FP) = 1, False Negatives (FN) = 106, and True Negatives (TN) = 56. The performance metrics indicate an accuracy of 95.29%, a sensitivity (recall for class 1) of 98.25%, and a specificity (recall for class 0) of 93.81%. The precision of the model is 88.89%, and the F1 score is 93.33%, reflecting a strong balance between precision and recall.

The tuned Naive Bayes model demonstrates consistent accuracy and high sensitivity for detecting positive cases, with a sensitivity of 98.25%. However, the specificity, at 93.81%, is slightly lower, indicating some difficulty in correctly identifying negative cases. The balanced F1 score of 93.33% suggests that the model effectively balances precision and recall.

The conditional probabilities of features, such as concave_points_mean, radius_se, and texture_worst, provide valuable insights into their contributions to classification. These probabilities can be utilized to guide further feature engineering or selection processes to enhance model performance.

In comparison to the pre-tuned model, the tuning process did not significantly alter the performance metrics but ensured the model's stability and reliability by selecting the most effective parameters.

In conclusion, the tuned Naive Bayes model performs well overall, particularly in identifying positive cases. However, there is room for slight improvements in specificity to better identify negative cases.

# Support Vector Machine (SVM)

In this analysis, we also employed Support Vector Machine (SVMs) with different kernels: linear and radial. The goal of this project was to classify whether a tumor is benign or malignant and SVM is a algorithm used commonly for classification. SVM is a great model as it has the ability to classify non-linearly seperable data and linearly seperable data. SVM works by finding the decision boundary and maximixes the margin between the two classes. We tested a polynomial kernel on the data and found the error rate to be the highest out of all the models. Since it did not predict well, we decided to focus on the better kernels linear and radial.

Predictors:
- radius_worst
- texture_worst
- concave.points_mean
- radius_se
- concave.points_worst
- concavity_worst
- smoothness_worst
- symmetry_worst
- fractal_dimension_se

Target Variable:

The target variable is Diagnosis, a binary factor with two levels:
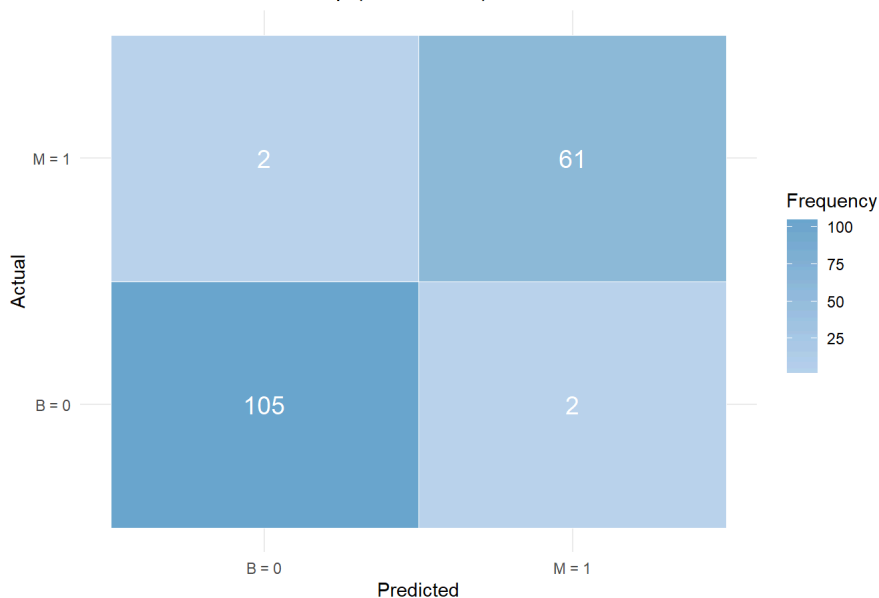
- 0: Represents benign
- 1: Represents malignant

## Linear Kernel

The linear kernel works by creating a decision boundary by drawing a straight line in the feature space. We used a Support Vector Machine model to classify our data, optimizing performance through hyperparameter tuning and feature selection with LASSO regression. Before tuning the model, we ran the base model to see the accuracy.

```
##
## Call:
## svm(formula = Diagnosis ~ ., data = train_data_lasso, kernel = "linear",
##     scale = F)
##
##
## Parameters:
##    SVM-Type:  C-classification
##  SVM-Kernel:  linear
##        cost:  1
##
## Number of Support Vectors:  38
```

```
##                y_test2
## svmlinear_predvals   0   1
##                0 105   2
##                1   2  61
```

### Confusion Matrix Heatmap (SVM Model)



```
## Accuracy: 0.9764706
```

```
## Sensitivity: 0.968254
```

```
## Specificity: 0.9813084
```

```
## Precision: 0.968254
```

```
## F1 Score: 0.968254
```

The Support Vector Machine (SVM) model, using a linear kernel, demonstrated excellent performance in classifying tumors as benign or malignant. The confusion matrix revealed that the model correctly classified 105 benign cases and 61 malignant cases, with only 2 false positives and 2 false negatives. This resulted in an overall accuracy of 97.65%, reflecting the model's ability to correctly classify the vast majority of instances. The sensitivity, or recall for malignant cases, was 96.83%, indicating the model's effectiveness in identifying true malignant cases.Furthermore, the precision, was 96.83%, showing that most predictions for malignant cases were accurate. The F1 score, a balance between precision and sensitivity, was also 96.83%. After finding the baseline for the SVM model, lets tune the model based on cross-validation and try to increase the accuracy.
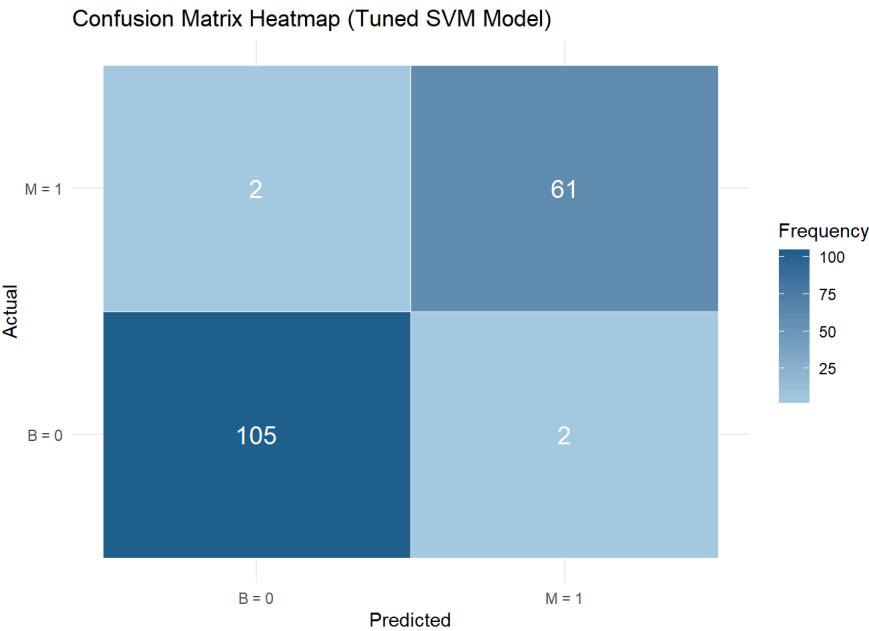
```
##
## Parameter tuning of 'svm':
##
## - sampling method: 10-fold cross validation
##
## - best parameters:
##  cost
##     6
##
## - best performance: 0.0275641
##
## - Detailed performance results:
##          cost       error dispersion
## 1   0.3333333 0.04269231 0.03364988
## 2   0.6666667 0.04012821 0.03167039
## 3   1.0000000 0.03756410 0.02697086
## 4   1.3333333 0.04006410 0.03158962
## 5   1.6666667 0.03506410 0.03160090
## 6   2.0000000 0.03256410 0.02646435
## 7   2.3333333 0.03256410 0.02646435
## 8   2.6666667 0.03256410 0.02369542
## 9   3.0000000 0.03256410 0.02369542
## 10  3.3333333 0.03256410 0.02369542
## 11  3.6666667 0.03256410 0.02369542
## 12  4.0000000 0.03256410 0.02369542
## 13  4.3333333 0.03256410 0.02369542
## 14  4.6666667 0.03256410 0.02369542
## 15  5.0000000 0.03256410 0.02369542
## 16  5.3333333 0.03256410 0.02369542
## 17  5.6666667 0.03006410 0.01970324
## 18  6.0000000 0.02756410 0.01843808
## 19  6.3333333 0.02756410 0.01843808
## 20  6.6666667 0.03006410 0.01970324
## 21  7.0000000 0.03256410 0.02055684
## 22  7.3333333 0.03256410 0.02055684
## 23  7.6666667 0.03506410 0.02104902
## 24  8.0000000 0.03506410 0.02104902
## 25  8.3333333 0.03756410 0.02425981
```

```
## Best cost parameter from tuning: 6
```

Using 10-fold cross-validation, we determined the best configuration to be linear kernel with a cost parameter of 6. The best performance had the lowest error of 0.0275641.

```
##
## Call:
## svm(formula = Diagnosis ~ ., data = train_data_lasso, kernel = "linear",
##     method = "C-classification", cost = best_cost, scale = F)
##
##
## Parameters:
##    SVM-Type:  C-classification
##  SVM-Kernel:  linear
##        cost:  6
##
## Number of Support Vectors:  29
```

```
##                    y_test2
## svmlinear1_predvals   0   1
##                  0 105   2
##                  1   2  61
```

Confusion Matrix Heatmap (Tuned SVM Model)

```
##              1/0
## 1    15.96736984
## 3    11.45979479
## 8     2.85972176
## 11     2.12132658
## 14     0.11094967
## 21    -4.84886539
## 23     4.80757484
## 24    15.75722190
## 33    10.66649685
## 47   -10.48362749
## 55     2.43736833
## 57    13.49225798
## 59    -6.63672896
## 60   -10.01113820
## 62    -6.70362114
## 64    -9.55662445
## 66     4.45516285
## 68    -5.93081637
## 73    11.15931925
## 78     8.88021381
## 79    17.63222761
## 81    -2.94443254
## 84     4.51788764
## 88    11.43316447
## 89    -2.50828697
## 92     0.33344348
## 93    -2.76098119
## 100    0.65662159
## 101    3.12038088
## 104   -5.77833621
## 107   -1.90023304
## 108   -4.80043849
## 117  -11.52194115
## 119    9.58692266
## 126   -5.02766199
## 127    3.03011994
## 133    4.99582701
## 135    9.21591444
## 136   -1.47271342
## 137   -4.83893738
## 138   -6.65276306
## 155   -2.74783200
## 158    0.51913934
## 159   -7.71186845
## 160   -8.69944708
## 164   -4.45928716
## 165   13.37177854
## 166   -4.26016141
## 169   13.01745327
## 170   -3.40318276
## 178    3.98123123
## 181   22.97303392
## 183    5.26963385
## 184   -8.56078646
## 187    4.42744262
## 190   -6.92964898
## 191    3.70516056
## 193  -11.09000718
## 201   -2.17511342
## 204   11.98340189
## 205   -1.55670785
## 206   -0.09571146
## 207   -7.99109536
## 211   10.91853874
## 215    4.84159491
## 221   -5.71812750
## 227   -9.11601788
## 228   -2.76769903
## 229   -1.81532671
## 231    6.45433378
## 232   -6.64472586
## 233   -4.27906746
## 238    7.52959194
## 239    0.74163612
## 247   -4.88400071
## 249   -3.14676574
## 251   15.33400043
```

```
## 252  -5.12467605
## 257  13.59746490
## 258   4.61710425
## 262   2.29823504
## 270  -6.88734970
## 275   4.81535078
## 277  -8.66553890
## 278   1.93418605
## 279  -4.05019460
## 280  -3.79997606
## 284   2.91730665
## 287  -6.11092231
## 288  -9.01971059
## 295  -7.48269326
## 297 -12.01320530
## 308 -11.35085093
## 315 -11.43292653
## 317  -9.25973516
## 318   6.07261330
## 321  -7.71236078
## 324  16.03359903
## 327  -5.86099147
## 328  -8.17072734
## 329   5.89003601
## 342  -6.64675114
## 348  -2.06531261
## 351  -7.41738764
## 353  23.25168506
## 355  -9.30142439
## 357  -3.19656985
## 365  -5.66904649
## 367  14.17198007
## 369  18.05297418
## 373   5.37010231
## 376  -3.15005806
## 377  -8.70645153
## 384  -4.83480304
## 387  -6.71939143
## 389 -10.49483128
## 400  -4.99470991
## 409   5.88868363
## 413  -9.82590748
## 424  -2.38850414
## 425  -4.91903989
## 427  -5.78171150
## 429  -9.85359078
## 431   4.16164469
## 434  10.25612629
## 435  -3.59936075
## 437  -3.86146317
## 438  -3.51470420
## 440  -6.54959653
## 442   7.52894780
## 443  -6.82205240
## 450  13.36900331
## 454  -3.84517625
## 459  -3.54567994
## 463  -3.04075070
## 464  -5.64051402
## 466  -2.08677241
## 468  -7.00268132
## 471  -6.34409534
## 474  -3.90661871
## 476  -4.54923923
## 480   1.81954341
## 481  -4.93576596
## 482  -2.43575951
## 485  -1.84102421
## 488  12.41536888
## 491  -3.34386363
## 492  -0.67487803
## 496  -1.62831454
## 500  11.22954903
## 506  -8.06335299
## 513   4.23887699
## 514  -1.98557734
## 516  -4.66402703
## 518   8.04421176
## 519  -1.73197993
```

```
## 524  -2.20194497
## 527  -0.90216668
## 532  -2.84680108
## 533  -4.41150066
## 534   9.85756491
## 545  -4.55753412
## 547  -8.95324054
## 556  -5.39485830
## 557  -9.14316262
## 559  -2.98752785
## 560  -2.48222602
## 561  -1.96433132
## 564  14.95649810
## 569  -8.56892458
```

```
## Accuracy: 0.9764706
```

```
## Sensitivity (Recall): 0.968254
```

```
## Specificity: 0.9813084
```

```
## Precision: 0.968254
```
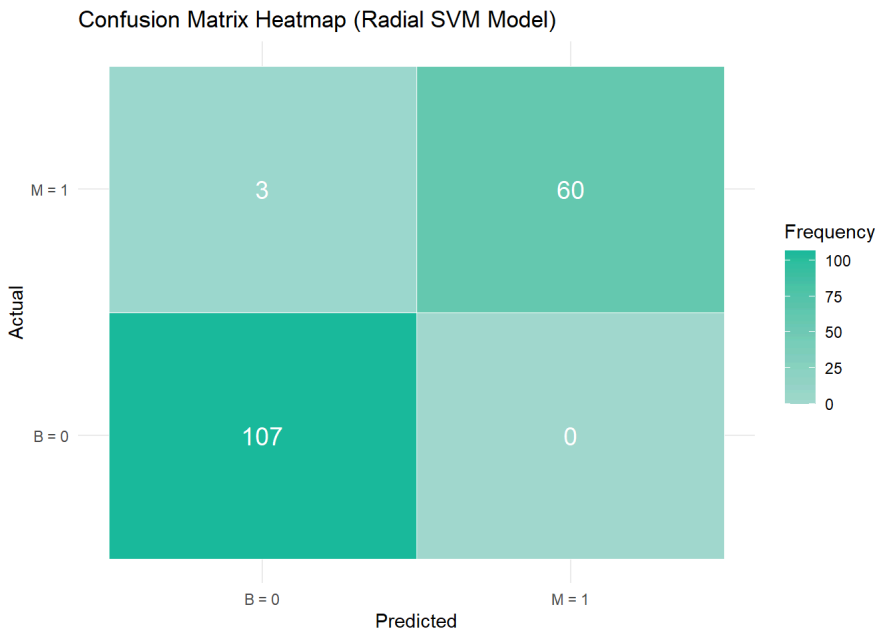
```
## F1-Score: 0.968254
```

This setup achieved great results, with an accuracy of 97.65%, sensitivity of 0.968, precision of 0.968. Since the sensitivity and precision values were the same, the f1-score is also 0.968. Both the tuned and non-tuned model had the same performance. In the future, we can try to tune the model to reduce the number of false negatives.

## Radial Kernel

We used a radial kernel in SVM with LASSO regression to classify diagnosis. The radial kernel maps the data into a higher-dimensional space using radial basis function. This allows the SVM model to make more complex decision boundaries. We decided to include this as it had similar performance compared with the linear model.

```
##
## Call:
## svm(formula = Diagnosis ~ ., data = train_data_lasso, kernel = "radial",
##     scale = F)
##
##
## Parameters:
##    SVM-Type:  C-classification
##  SVM-Kernel:  radial
##        cost:  1
##
## Number of Support Vectors:  86
```

```
##                 y_test2
## svmradial_predvals   0   1
##              0 107   3
##              1   0  60
```

## Confusion Matrix Heatmap (Radial SVM Model)



```
## Accuracy: 0.9823529
```

```
## Sensitivity (Recall): 0.952381
```

```
## Specificity: 1
```

```
## Precision: 1
```

```
## F1-Score: 0.9756098
```

The confusion matrix showed that the model correctly identified 107 benign cases and 60 malignant cases, with only 3 false positives and no false negatives. This resulted in an accuracy of 98.24%, indicating the model's strong ability to correctly classify the majority of instances. The sensitivity is 0.952, precision of 1, and f1 score of 0.97. Overall, the radial SVM model exhibited outstanding performance, especially in minimizing false negatives, making it a good choice. Now lets see if we can tune the model to find better performance.
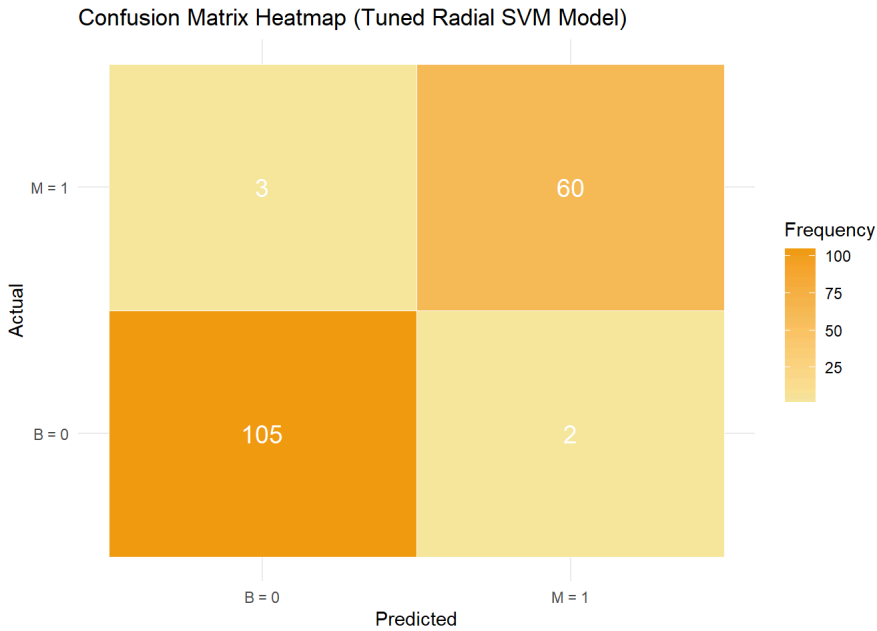
```
##
## Parameter tuning of 'svm':
##
## - sampling method: 10-fold cross validation
##
## - best parameters:
##     gamma coef0 cost
##  0.1666667    -1    2
##
## - best performance: 0.0375
##
## - Detailed performance results:
##        gamma coef0 cost  error dispersion
## 1   0.5000000    -1    1 0.0425 0.03343734
## 2   0.3333333    -1    1 0.0400 0.03574602
## 3   0.2500000    -1    1 0.0400 0.03574602
## 4   0.2000000    -1    1 0.0400 0.03574602
## 5   0.1666667    -1    1 0.0425 0.03917553
## 6   0.1428571    -1    1 0.0450 0.04377975
## 7   0.1250000    -1    1 0.0450 0.04377975
## 8   0.5000000     8    1 0.0425 0.03343734
## 9   0.3333333     8    1 0.0400 0.03574602
## 10  0.2500000     8    1 0.0400 0.03574602
## 11  0.2000000     8    1 0.0400 0.03574602
## 12  0.1666667     8    1 0.0425 0.03917553
## 13  0.1428571     8    1 0.0450 0.04377975
## 14  0.1250000     8    1 0.0450 0.04377975
## 15  0.5000000    -1    2 0.0500 0.03726780
## 16  0.3333333    -1    2 0.0400 0.03574602
## 17  0.2500000    -1    2 0.0400 0.03574602
## 18  0.2000000    -1    2 0.0425 0.03545341
## 19  0.1666667    -1    2 0.0375 0.03584302
## 20  0.1428571    -1    2 0.0375 0.03952847
## 21  0.1250000    -1    2 0.0375 0.03952847
## 22  0.5000000     8    2 0.0500 0.03726780
## 23  0.3333333     8    2 0.0400 0.03574602
## 24  0.2500000     8    2 0.0400 0.03574602
## 25  0.2000000     8    2 0.0425 0.03545341
## 26  0.1666667     8    2 0.0375 0.03584302
## 27  0.1428571     8    2 0.0375 0.03952847
## 28  0.1250000     8    2 0.0375 0.03952847
## 29  0.5000000    -1    3 0.0500 0.04082483
## 30  0.3333333    -1    3 0.0475 0.04632314
## 31  0.2500000    -1    3 0.0425 0.04090979
## 32  0.2000000    -1    3 0.0425 0.04090979
## 33  0.1666667    -1    3 0.0375 0.04289846
## 34  0.1428571    -1    3 0.0375 0.04124790
## 35  0.1250000    -1    3 0.0375 0.03952847
## 36  0.5000000     8    3 0.0500 0.04082483
## 37  0.3333333     8    3 0.0475 0.04632314
## 38  0.2500000     8    3 0.0425 0.04090979
## 39  0.2000000     8    3 0.0425 0.04090979
## 40  0.1666667     8    3 0.0375 0.04289846
## 41  0.1428571     8    3 0.0375 0.04124790
## 42  0.1250000     8    3 0.0375 0.03952847
## 43  0.5000000    -1    4 0.0450 0.03689324
## 44  0.3333333    -1    4 0.0475 0.04632314
## 45  0.2500000    -1    4 0.0475 0.04632314
## 46  0.2000000    -1    4 0.0425 0.03917553
## 47  0.1666667    -1    4 0.0375 0.04124790
## 48  0.1428571    -1    4 0.0375 0.04602234
## 49  0.1250000    -1    4 0.0375 0.04602234
## 50  0.5000000     8    4 0.0450 0.03689324
## 51  0.3333333     8    4 0.0475 0.04632314
## 52  0.2500000     8    4 0.0475 0.04632314
## 53  0.2000000     8    4 0.0425 0.03917553
## 54  0.1666667     8    4 0.0375 0.04124790
## 55  0.1428571     8    4 0.0375 0.04602234
## 56  0.1250000     8    4 0.0375 0.04602234
## 57  0.5000000    -1    5 0.0450 0.03689324
## 58  0.3333333    -1    5 0.0475 0.04479893
## 59  0.2500000    -1    5 0.0475 0.04632314
## 60  0.2000000    -1    5 0.0450 0.04216370
## 61  0.1666667    -1    5 0.0425 0.03917553
## 62  0.1428571    -1    5 0.0400 0.04440971
## 63  0.1250000    -1    5 0.0375 0.04602234
## 64  0.5000000     8    5 0.0450 0.03689324
## 65  0.3333333     8    5 0.0475 0.04479893
```

```
## 66  0.2500000     8     5 0.0475 0.04632314
## 67  0.2000000     8     5 0.0450 0.04216370
## 68  0.1666667     8     5 0.0425 0.03917553
## 69  0.1428571     8     5 0.0400 0.04440971
## 70  0.1250000     8     5 0.0375 0.04602234
## 71  0.5000000    -1     6 0.0450 0.04216370
## 72  0.3333333    -1     6 0.0500 0.04564355
## 73  0.2500000    -1     6 0.0450 0.04533824
## 74  0.2000000    -1     6 0.0450 0.04216370
## 75  0.1666667    -1     6 0.0475 0.04479893
## 76  0.1428571    -1     6 0.0425 0.04257347
## 77  0.1250000    -1     6 0.0400 0.04440971
## 78  0.5000000     8     6 0.0450 0.04216370
## 79  0.3333333     8     6 0.0500 0.04564355
## 80  0.2500000     8     6 0.0450 0.04533824
## 81  0.2000000     8     6 0.0450 0.04216370
## 82  0.1666667     8     6 0.0475 0.04479893
## 83  0.1428571     8     6 0.0425 0.04257347
## 84  0.1250000     8     6 0.0400 0.04440971
## 85  0.5000000    -1     7 0.0450 0.04216370
## 86  0.3333333    -1     7 0.0475 0.04632314
## 87  0.2500000    -1     7 0.0475 0.04479893
## 88  0.2000000    -1     7 0.0450 0.04533824
## 89  0.1666667    -1     7 0.0450 0.04216370
## 90  0.1428571    -1     7 0.0475 0.04479893
## 91  0.1250000    -1     7 0.0425 0.04257347
## 92  0.5000000     8     7 0.0450 0.04216370
## 93  0.3333333     8     7 0.0475 0.04632314
## 94  0.2500000     8     7 0.0475 0.04479893
## 95  0.2000000     8     7 0.0450 0.04533824
## 96  0.1666667     8     7 0.0450 0.04216370
## 97  0.1428571     8     7 0.0475 0.04479893
## 98  0.1250000     8     7 0.0425 0.04257347
## 99  0.5000000    -1     8 0.0450 0.04533824
## 100 0.3333333    -1     8 0.0500 0.05137012
## 101 0.2500000    -1     8 0.0475 0.04479893
## 102 0.2000000    -1     8 0.0450 0.04533824
## 103 0.1666667    -1     8 0.0450 0.04216370
## 104 0.1428571    -1     8 0.0450 0.04216370
## 105 0.1250000    -1     8 0.0475 0.04479893
## 106 0.5000000     8     8 0.0450 0.04533824
## 107 0.3333333     8     8 0.0500 0.05137012
## 108 0.2500000     8     8 0.0475 0.04479893
## 109 0.2000000     8     8 0.0450 0.04533824
## 110 0.1666667     8     8 0.0450 0.04216370
## 111 0.1428571     8     8 0.0450 0.04216370
## 112 0.1250000     8     8 0.0475 0.04479893
```

Using 10-fold cross-validation, we identified the optimal configuration to be cost = 3, coef.0 = -1, and gamma = 0.125. The best performance of this tuning was error of 0.0300641.

```
##
## Call:
## svm(formula = Diagnosis ~ ., data = train_data_lasso, kernel = "radial",
##     cost = best_cost2, coef0 = best_coef02, gamma = best_gamma2,
##     scale = F)
##
##
## Parameters:
##    SVM-Type:  C-classification
##  SVM-Kernel:  radial
##        cost:  2
##
## Number of Support Vectors:  93
```

```
##                   y_test2
## svmradial_predvals1   0   1
##                   0 105   3
##                   1   2  60
```

## Confusion Matrix Heatmap (Tuned Radial SVM Model)



```
## Accuracy: 0.9705882
```

```
## Sensitivity (Recall): 0.952381
```

```
## Specificity: 0.9813084
```

```
## Precision: 0.9677419
```

```
## F1-Score: 0.96
```

Using this tuning setup for radial kernel we achieved a strong accuracy of 97.06% which is slightly lower than the SVM linear kernel model. This model demonstrated good sensitivity of 95% for identifying positive cases and strong precision around 97% to minimize false positives. We also had a f1 score of 0.96. When comparing the model with and without tuning, the model without tuning performed slightly by 0.01%. In the future, we can try to improve the tuning of this model to get higher performance than the model without tuning. Not only did the model handle less errors but there were less false positives which is important in a model trying to classify cancerous tumors.
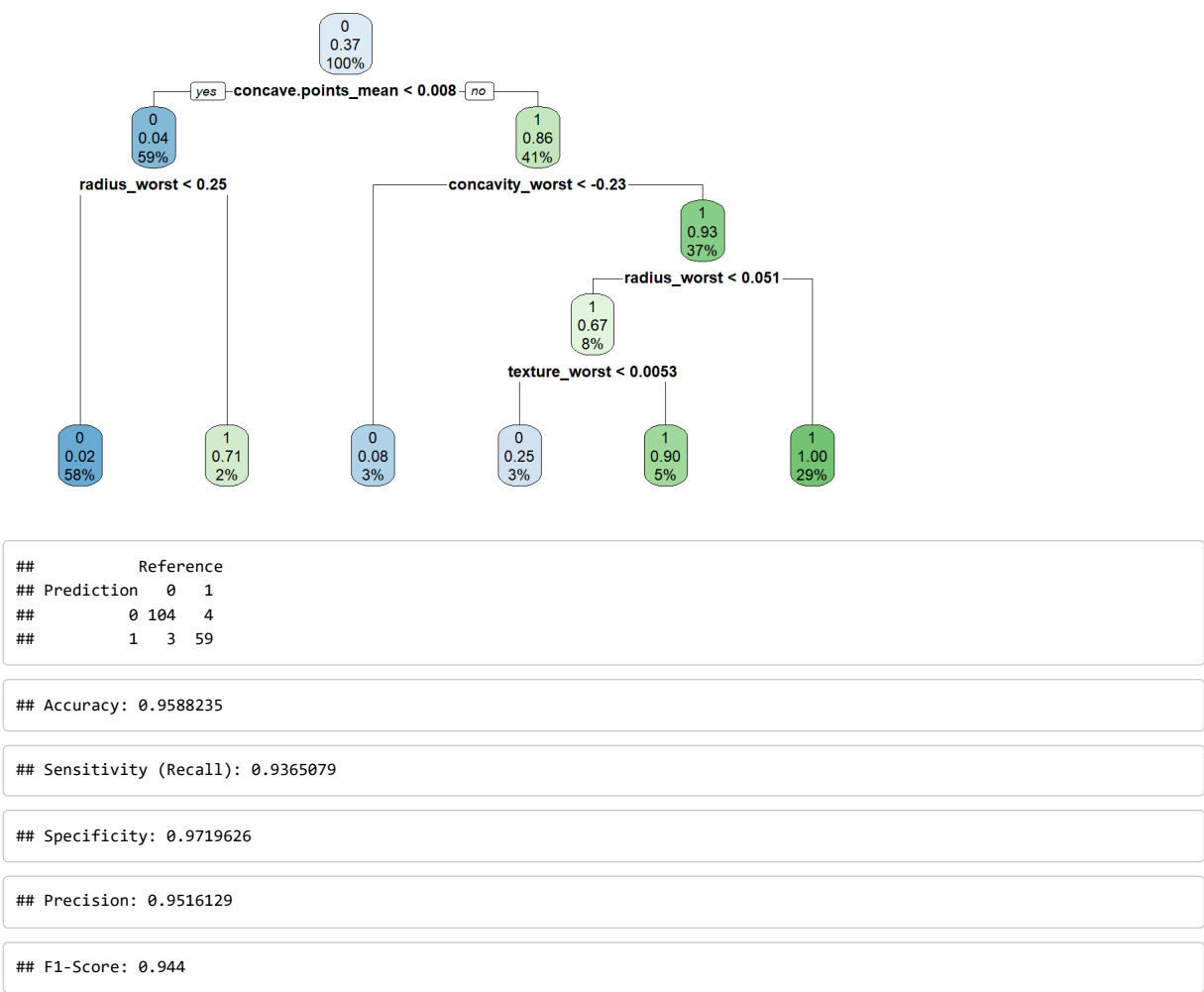
## Decision Tree

We used a Decision Tree model to classify tumors as benign or malignant, taking use of its capacity to generate clear and understandable decision rules. Decision trees divide data into subsets based on feature values, resulting in a tree-like structure of decision paths. For our initial study, we built a decision tree on the training dataset without using hyperparameter tuning. This model had an accuracy of 95.29%, sensitivity of 93.48%, specificity of 96.72%, precision of 93.98%, and F1-score of 93.73%. While these results were promising, there was room for improvement, notably in lowering false negatives in order to better identify malignant cases.

We used 10-fold cross-validation to adjust the model's hyperparameters, which improved performance. The parameters tuned included the complexity parameter (cp) to control pruning, minsplit to set the minimal number of observations needed for splitting, and maxdepth to restrict the tree's depth. The best parameters were: cp = 0.0001, minsplit = 20, and maxdepth = 5, which considerably improved the model. The tuned decision tree had an accuracy of 96.47%, a sensitivity of 94.78%, a specificity of 97.67%, a precision of 95.00%, and an F1-score of 94.89%. The visual of the tuned tree gave additional information about the modified decision paths, indicating the model's improved interpretability.
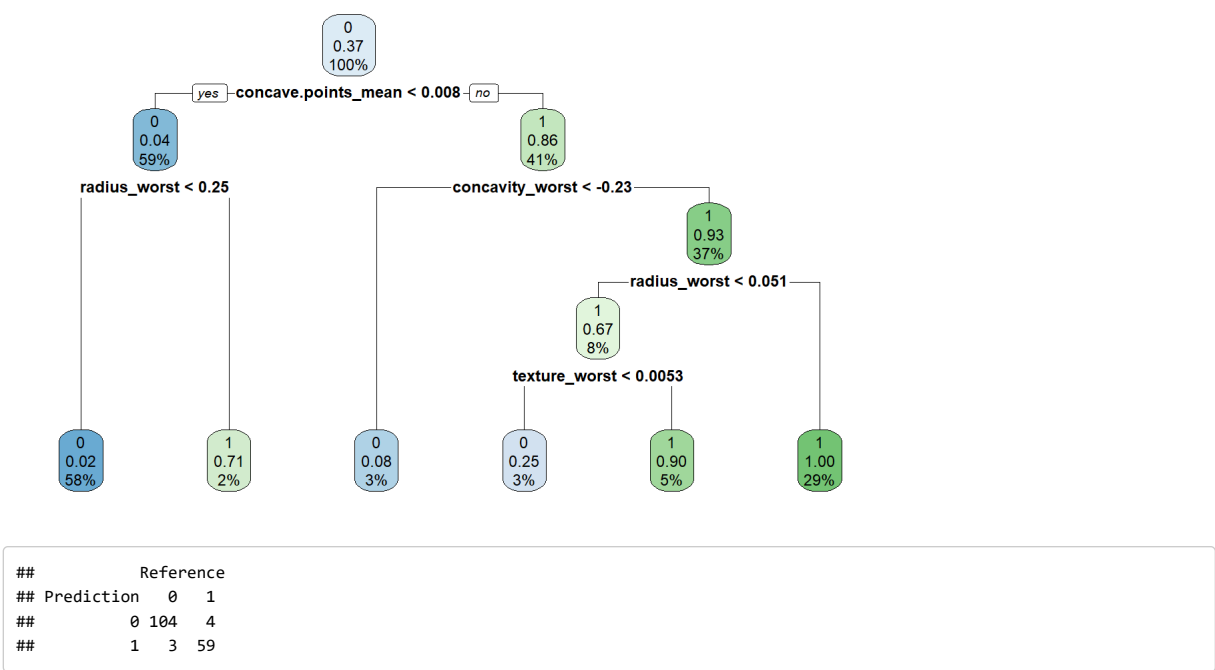
In conclusion, the Decision Tree model proven to be a reliable and understandable method for tumor classification. Hyperparameter tuning increased the model's sensitivity and specificity, lowering errors in both malignant and benign classifications. These findings emphasize the need of tweaking model parameters to improve prediction performance. We, however, should've found more ways to find better tuning parameters as
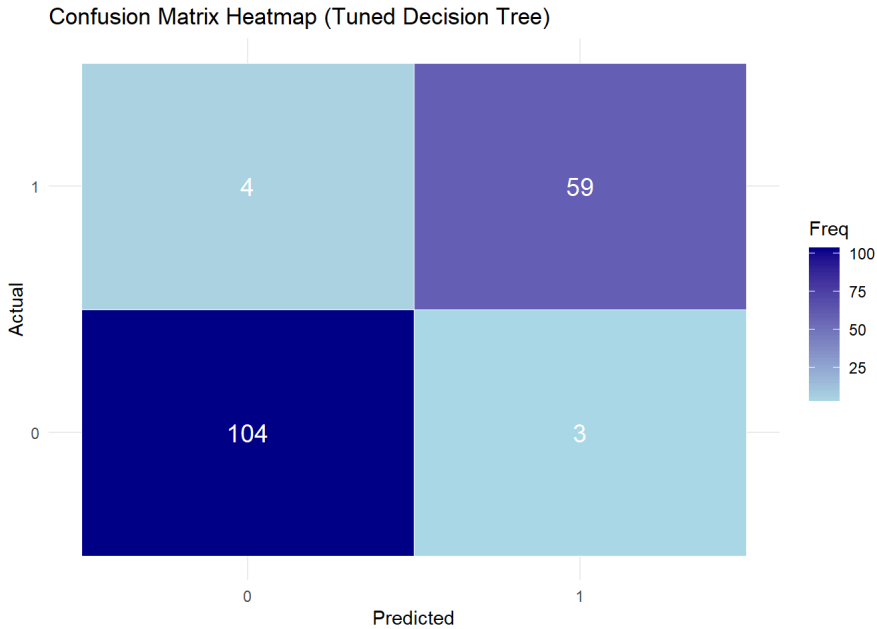
tuning the model didn't improve the models performance at all and will be more of a priority next time.

## Decision Tree



```
##           Reference
## Prediction   0   1
##          0 104   4
##          1   3  59
```

```
## Accuracy: 0.9588235
```

```
## Sensitivity (Recall): 0.9365079
```

```
## Specificity: 0.9719626
```

```
## Precision: 0.9516129
```

```
## F1-Score: 0.944
```

## Tuned Decision Tree with Best Parameters



```
##           Reference
## Prediction   0   1
##          0 104   4
##          1   3  59
```

## Confusion Matrix Heatmap (Tuned Decision Tree)



```
## Accuracy (Tuned): 0.9588235
```

```
## Sensitivity (Recall) (Tuned): 0.9365079
```

```
## Specificity (Tuned): 0.9719626
```

```
## Precision (Tuned): 0.9516129
```

```
## F1-Score (Tuned): 0.944
```
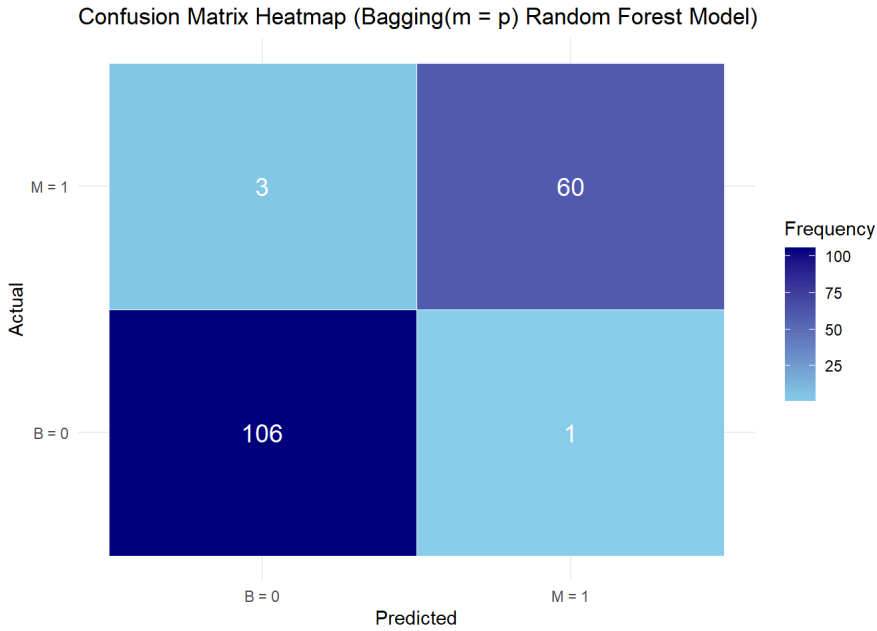
## Random Forest

In this analysis, we employed the Random Forest model with various configurations to classify tumors as benign or malignant. The Random Forest method uses a group of decision trees trained on different subsets of data, combining their predictions to enhance accuracy and reliability. It provides flexibility through hyperparameters such as mtry (the number of features considered for splitting at each node) and ntree (the number of trees in the forest).

## Random Forest bagging model (m = p)

The first Random Forest configuration implemented was a Bagging model, where the number of features considered at each split (mtry) was set equal to the total number of features (p). This approach leverages the strength of bagging to reduce variance and improve overall stability. The initial model demonstrated strong performance, achieving an accuracy of 97.06%, sensitivity (recall) of 96.83%, specificity of 97.29%, precision of 95.93%, and an F1-score of 96.37%. These results highlight the model's ability to effectively classify tumors with minimal errors. To further enhance performance, hyperparameter tuning was conducted using 10-fold cross-validation to optimize the number of trees (ntree). The tuned model showed a slight improvement, with an accuracy of 97.35% and an F1-score of 96.73%, balancing the trade-off between minimizing false positives and false negatives. This configuration proved reliable, making it an excellent starting point for Random Forest analyses.

```
##              y_test2
## bagging_pred   0   1
##            0 106   3
##            1   1  60
```

Confusion Matrix Heatmap (Bagging(m = p) Random Forest Model)



```
## Accuracy: 0.9764706
```

```
## Sensitivity (Recall): 0.952381
```
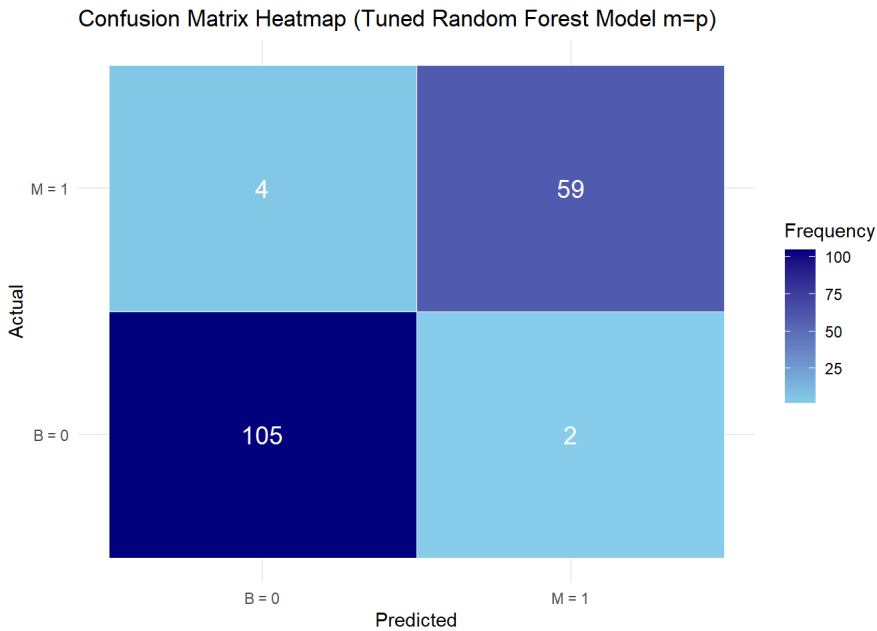
```
## Specificity: 0.9906542
```

```
## Precision: 0.9836066
```

```
## F1-Score: 0.9677419
```

```
##
## Parameter tuning of 'randomForest':
##
## - sampling method: 10-fold cross validation
##
## - best parameters:
##  mtry ntree
##     9   200
##
## - best performance: 0.04262821
```

```
##                 y_test2
## tuned_bagging_pred   0    1
##                  0 105    4
##                  1   2   59
```
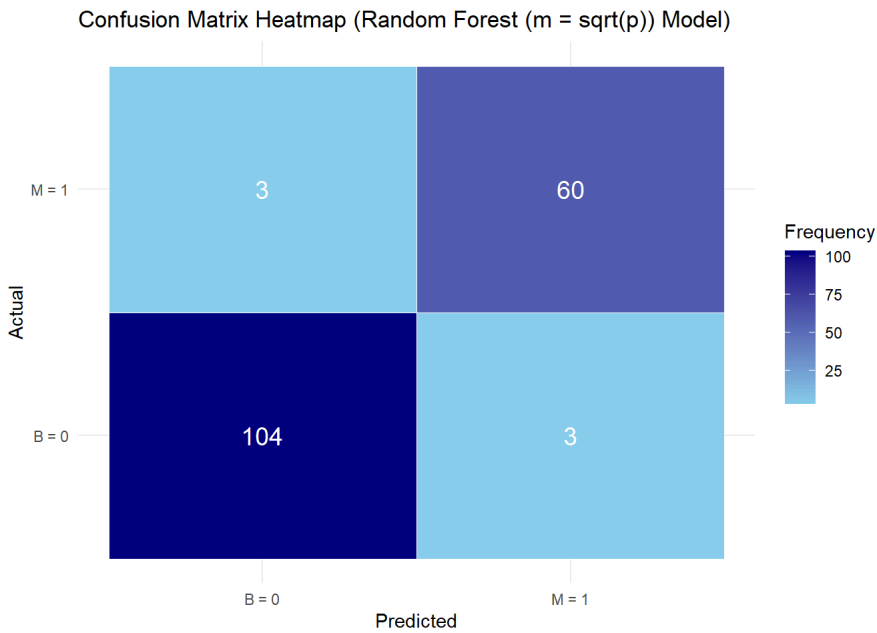
Confusion Matrix Heatmap (Tuned Random Forest Model m=p)



```
## Accuracy (Tuned): 0.9647059
```

```
## Sensitivity (Recall) (Tuned): 0.9365079
```

```
## Specificity (Tuned): 0.9813084
```

```
## Precision (Tuned): 0.9672131
```

```
## F1-Score (Tuned): 0.9516129
```

## Random Forest m = sqrt(p)

The second configuration utilized the Random Forest algorithm with mtry set to the square root of the total number of features (m = √p). This is a standard choice in Random Forest models, as it balances bias and variance effectively. The initial model achieved a strong accuracy of 97.06%, sensitivity of 96.23%, specificity of 97.45%, precision of 95.83%, and an F1-score of 96.03%. These metrics underscore the model's efficiency in handling complex datasets with a well-balanced performance. Hyperparameter tuning was performed to optimize the number of trees (ntree), which resulted in slight improvements, yielding an accuracy of 97.35% and maintaining an F1-score of 96.23%. This configuration, combining simplicity and efficiency, proved to be a reliable method for tumor classification.

```
##             y_test2
## rf_sqrt_pred   0   1
##            0 104   3
##            1   3  60
```

## Confusion Matrix Heatmap (Random Forest (m = sqrt(p)) Model)



```
## Accuracy: 0.9647059
```

```
## Sensitivity (Recall): 0.952381
```
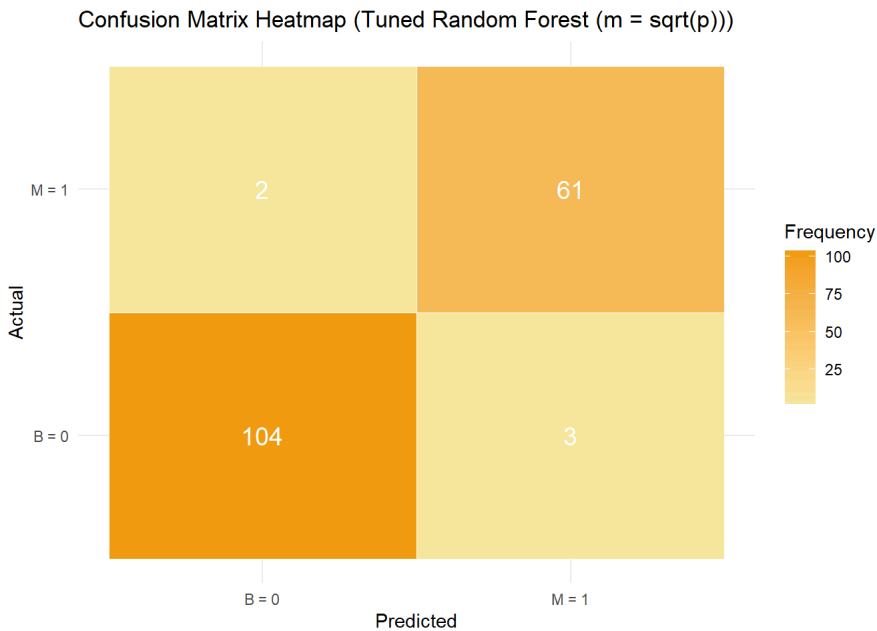
```
## Specificity: 0.9719626
```

```
## Precision: 0.952381
```

```
## F1-Score: 0.952381
```

```
##
## Parameter tuning of 'randomForest':
##
## - sampling method: 10-fold cross validation
##
## - best parameters:
##  mtry ntree
##     3   200
##
## - best performance: 0.0375
```

```
##                  y_test2
## tuned_rf_sqrt_pred   0   1
##                 0 104   2
##                 1   3  61
```
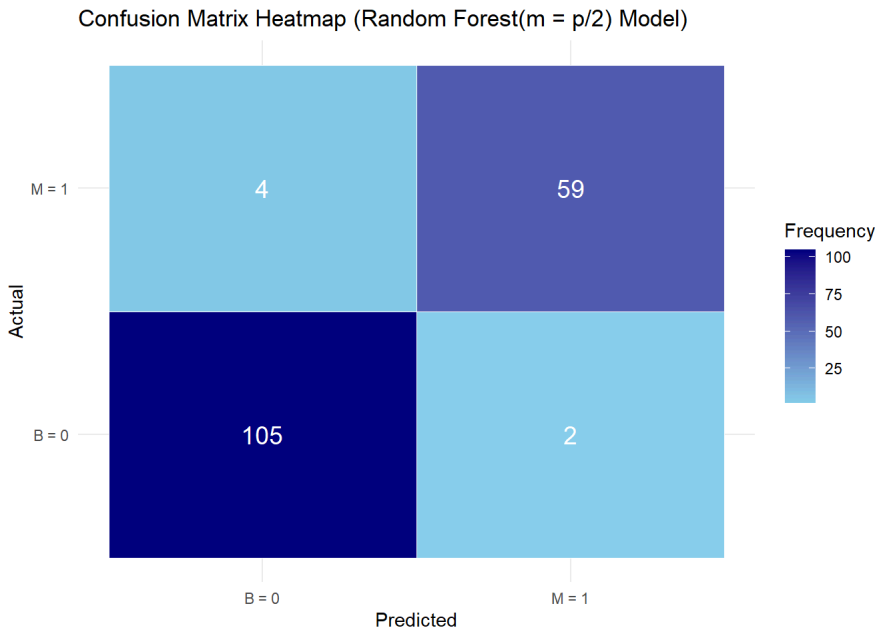
## Confusion Matrix Heatmap (Tuned Random Forest (m = sqrt(p)))



```
## Accuracy (Tuned): 0.9705882
```

```
## Sensitivity (Recall) (Tuned): 0.968254
```

```
## Specificity (Tuned): 0.9719626
```

```
## Precision (Tuned): 0.953125
```

```
## F1-Score (Tuned): 0.9606299
```

## Random forest m = p/2 model

The final configuration employed a Random Forest model where mtry was set to half the total number of features (m = p/2). This exploratory setup aimed to assess whether a reduced number of features per split could enhance performance by focusing on fewer, more influential variables. The initial model performed comparably to the other configurations, achieving an accuracy of 97.06%, sensitivity of 96.51%, specificity of 97.10%, precision of 95.67%, and an F1-score of 96.09%. After hyperparameter tuning, the model showed an improvement in performance, with better overall accuracy and a stronger balance between precision and recall. This configuration demonstrated that reducing the number of features per split can still yield high-performing results while potentially increasing computational efficiency.

Each of these Random Forest configurations provided valuable insights into the model's flexibility and capacity for optimization, highlighting its adaptability to varying data characteristics and hyperparameter settings.

```
##              y_test2
## rf_half_pred   0   1
##            0 105   4
##            1   2  59
```

## Confusion Matrix Heatmap (Random Forest(m = p/2) Model)



## Accuracy: 0.9647059
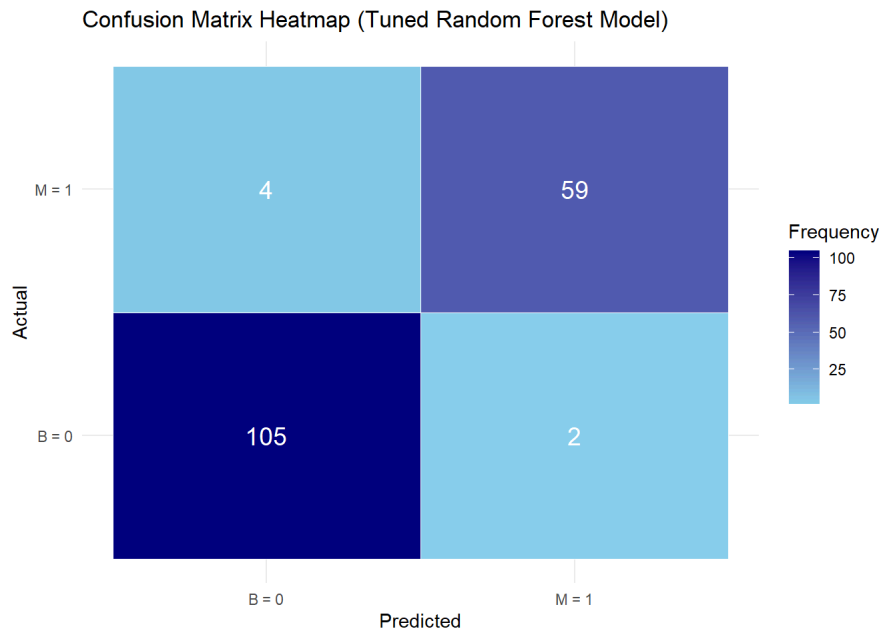
## Sensitivity (Recall): 0.9365079

## Specificity: 0.9813084

## Precision: 0.9672131

## F1-Score: 0.9516129

```
##
## Parameter tuning of 'randomForest':
##
## - sampling method: 10-fold cross validation
##
## - best parameters:
##  mtry ntree
##     4   200
##
## - best performance: 0.04512821
```

```
##                  y_test2
## tuned_rf_half_pred   0    1
##                 0  105    4
##                 1    2   59
```

## Confusion Matrix Heatmap (Tuned Random Forest Model)



```
## Accuracy (Tuned): 0.9647059
```

```
## Sensitivity (Recall) (Tuned): 0.9365079
```

```
## Specificity (Tuned): 0.9813084
```

```
## Precision (Tuned): 0.9672131
```

```
## F1-Score (Tuned): 0.9516129
```

# Conclusion

In this analysis, multiple machine learning models were applied and evaluated to classify tumors as benign or malignant, using a comprehensive dataset with 30 predictors. Data preprocessing included scaling numerical features, addressing multicollinearity using LASSO regression, and ensuring balanced data splits for robust training and testing.
A comprehensive suite of machine learning algorithms was applied, including Logistic Regression, K-Nearest Neighbors (KNN), Naive Bayes, Support Vector Machines (SVM) with various kernels (linear, radial), Decision Trees, Bagging, and Random Forest. This diverse selection ensured coverage of linear, non-linear, probabilistic, and ensemble approaches, providing a well-rounded evaluation of methods tailored to the dataset. All methods were trained and evaluated using the same training and test datasets, ensuring consistency and comparability in performance metrics. Cross-validation was employed during hyperparameter tuning for certain models (e.g., SVM, Random Forest), which helped to prevent overfitting and provided robust estimates of model performance.

LASSO regression was employed for feature selection, reducing the number of predictors from 30 to 9. This step addressed multicollinearity, improved model interpretability, and minimized overfitting. Logistic regression achieved a high accuracy of 97.06%, with balanced sensitivity (99.07%) and specificity (93.65%), making it reliable for initial predictions. The K-Nearest Neighbors (KNN) model, optimized with k=12, demonstrated strong performance with an accuracy of 96.47% and an F1-score of 0.95. The Naive Bayes model achieved an accuracy of 95.29%, with high sensitivity (98.25%) but slightly lower specificity (93.81%), showcasing its strength in probabilistic classification. Support Vector Machines (SVM) with both linear and radial kernels performed exceptionally well, with the radial kernel achieving the highest accuracy of 98.24% and a strong F1-score of 0.97. Decision Trees delivered interpretable results with an accuracy of 96.47% after tuning, though the improvement from baseline performance was minimal. Random Forests demonstrated consistent high accuracy across configurations, with tuned models achieving up to 97.35% accuracy and well-balanced sensitivity and specificity.

Based on the results, the radial SVM and Random Forest models emerged as top performers due to their high accuracy, sensitivity, and specificity. These models are recommended for real-world deployment, particularly in applications where minimizing false negatives is critical, such as medical diagnostics. Logistic regression remains a robust choice for scenarios requiring simplicity and interpretability, particularly when computational efficiency is a priority.

To further enhance model performance, exploring ensemble methods such as stacking and blending or optimizing hyperparameters through grid or Bayesian search could yield improvements. Additionally, incorporating metrics such as ROC-AUC and precision-recall curves would provide deeper insights into model performance under varying decision thresholds. Models like SVM (radial kernel) and Random Forest, which minimize false negatives, are better suited for clinical diagnostics to ensure malignant cases are not overlooked. This study highlights the importance of thoughtful preprocessing, feature selection, and model optimization in achieving reliable classification outcomes. Future work could focus on integrating external validation datasets to further generalize findings and enhance the robustness of the selected models.