

# Minor Project

## By Siddhartha Sinha

### 1) What is the Problem Statement?

: - **Breast Cancer Detection.**

Breast Cancer is a tumour in the chest region which turns Malignant. A tumour can be Malignant (cancerous) or Benign (non- cancerous). A tumour can be analysed using its different features and then labelled as Malignant or Benign. A doctor uses his years of knowledge assisted by other medical tests to label the tumour. After analysing and collecting breast tumour attributes and its class we see a distinct pattern that can be used to classify tumour as Cancerous or Non-cancerous. With help of Machine Learning we can train the algorithm to classify tumours.

### 2) Which dataset is chosen?

: - Dataset for this project is breast\_cancer.csv available at Breast Cancer Wisconsin [https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic))  
It has patient id along with 32 features and its respective diagnosis.

### 3) Libraries chosen?

**caret**: Classification and Regression Training  
**caTools**: Moving window statistics  
**corrplot**: Visualisation of a correlation matrix  
**dplyr**: Rules for data manipulation  
**ggplot2**: Data Visualisation such as heatmap ,etc  
**gridExtra**: Miscellaneous Function for grid graphics  
**pROC**: Display and analyse ROC curves  
**readr**: Read Rectangular Text data  
**MASS**: Support Functions and Datasets for venerable and Ripley's MASS

### 4) How cleaning/EDA was performed?

Firstly the dataset was converted into a dataframe. Then the information of dataframe was shown using  
**str(data)**

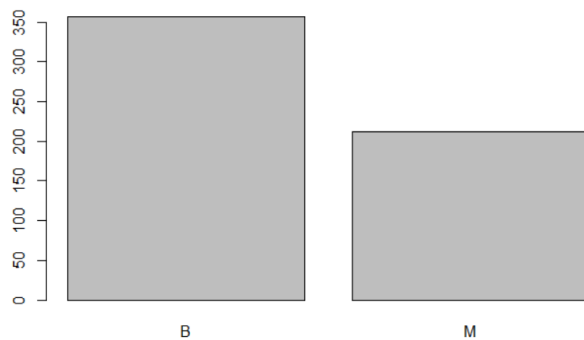
With this we selected the relevant dependent and independent features and also removed the null valued Column using

```
data$diagnosis <- as.factor(data$diagnosis)  
data[,33] <- NULL
```

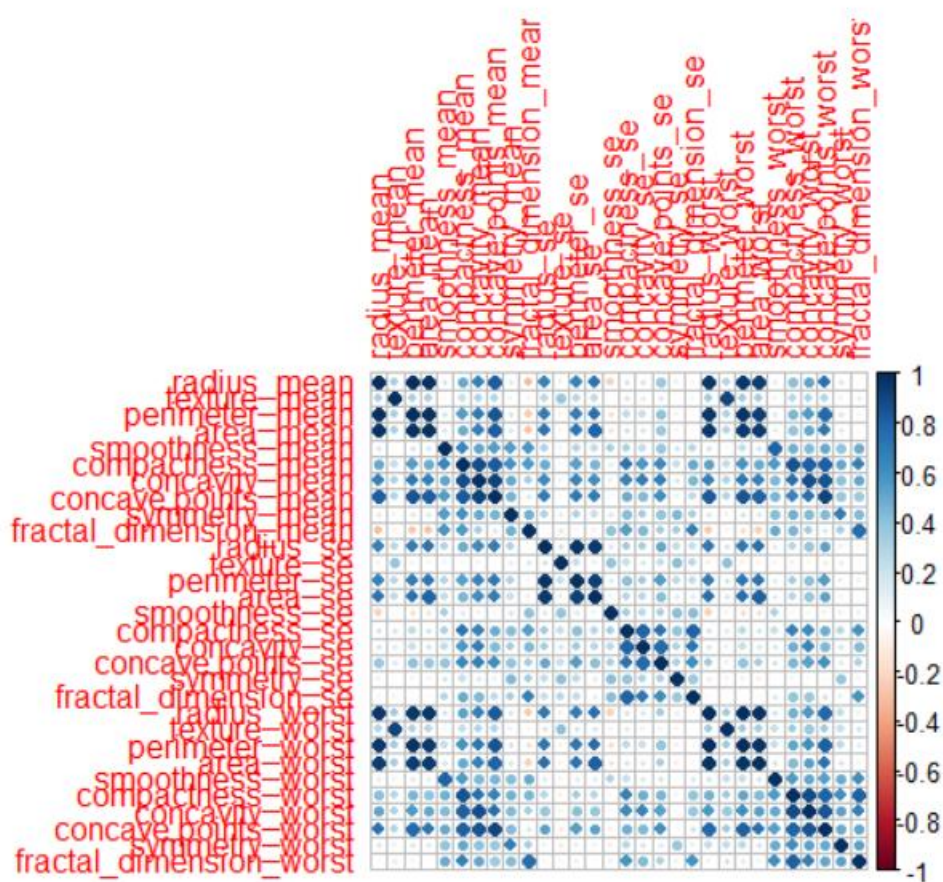
Then to find the amount of Malignant and Benign cases in the dataset

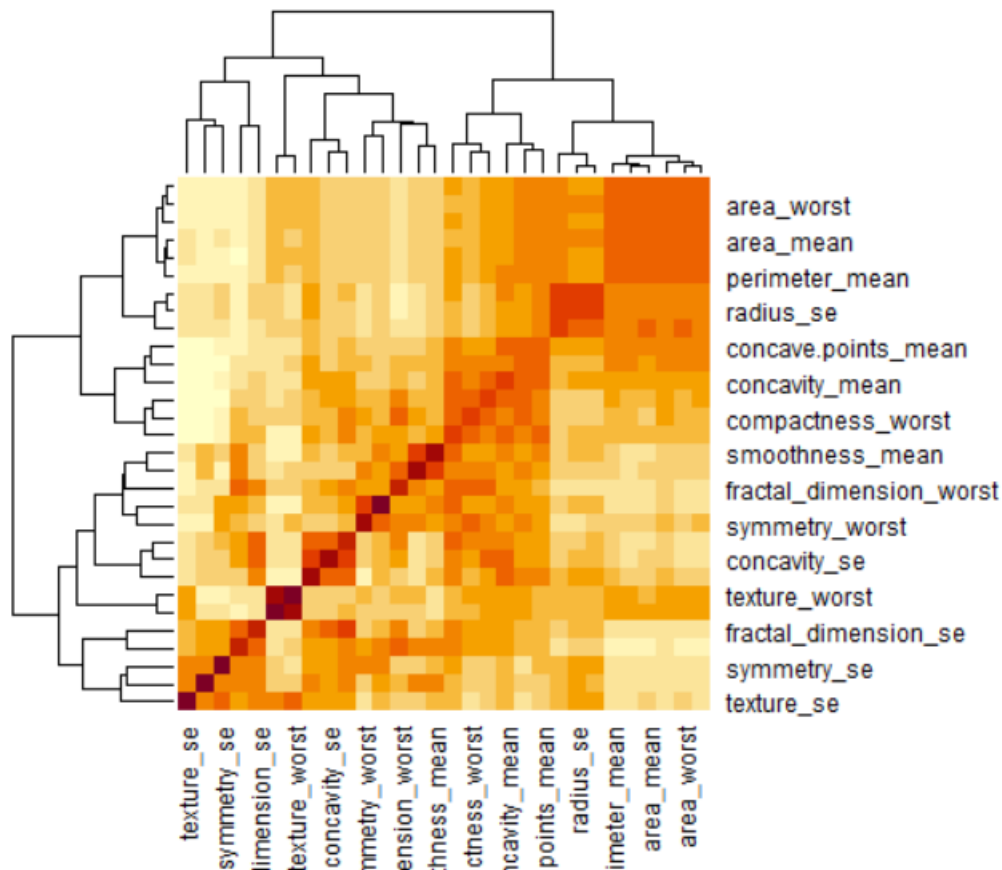
`prop.table(table(data$diagnosis))`

`plot(data$diagnosis)` was used, from that we deduced that the dataset had 62.74% Benign cases and 37.25% Malignant cases.

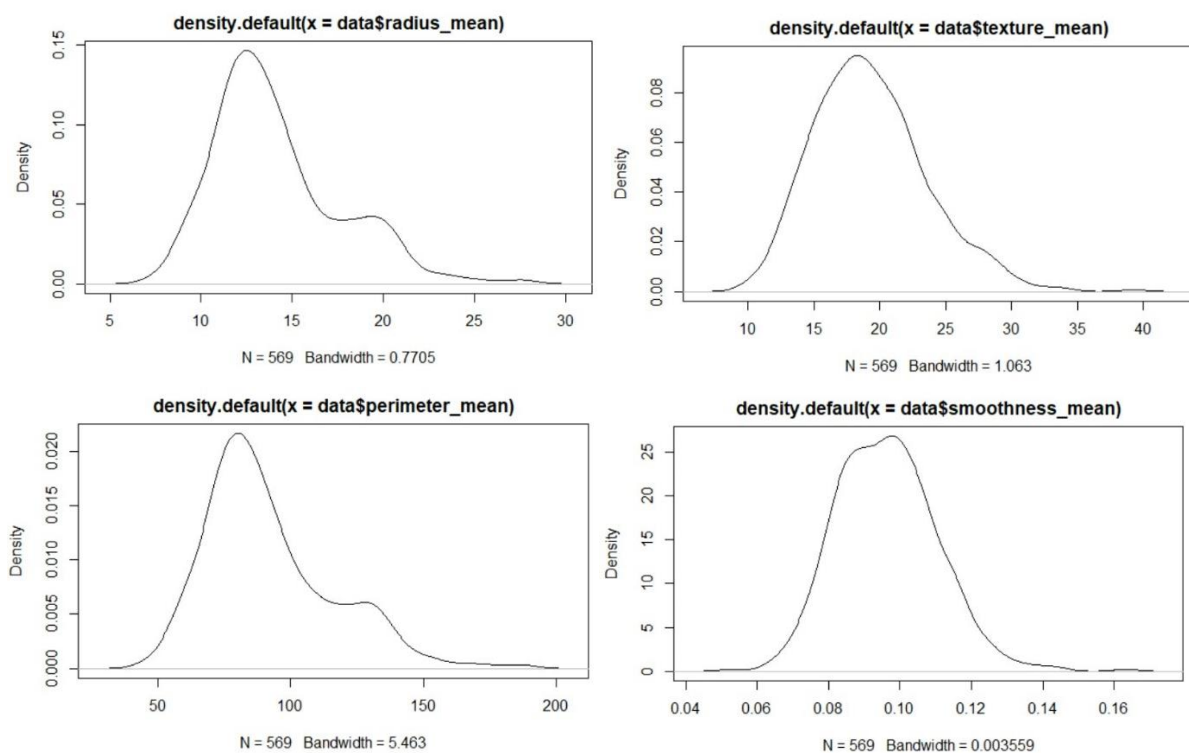


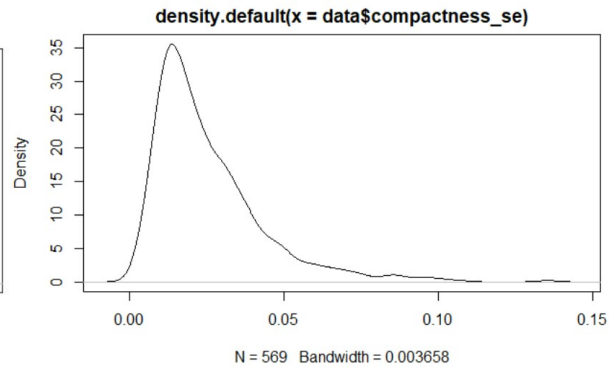
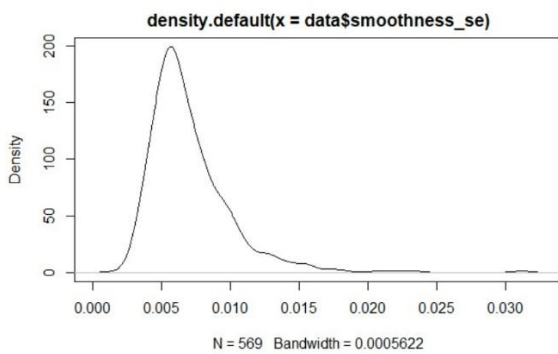
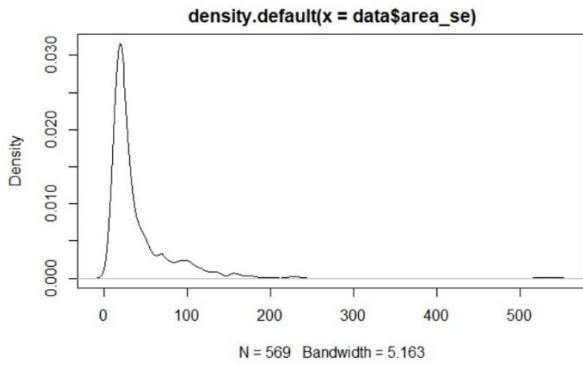
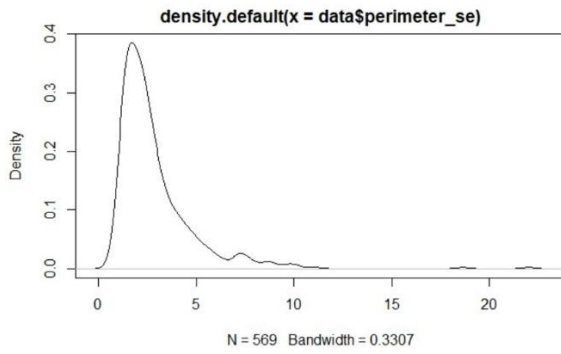
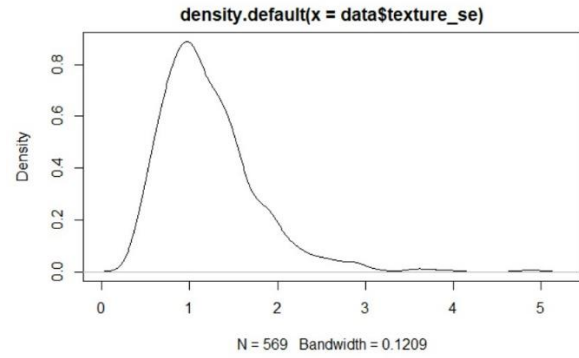
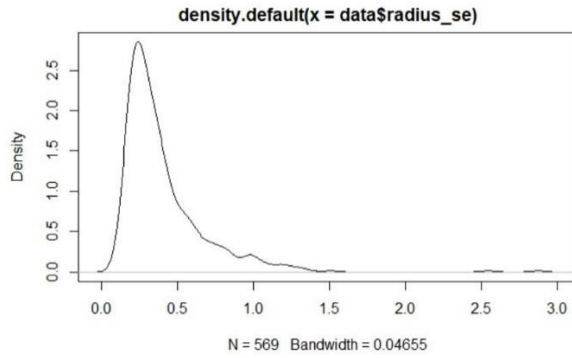
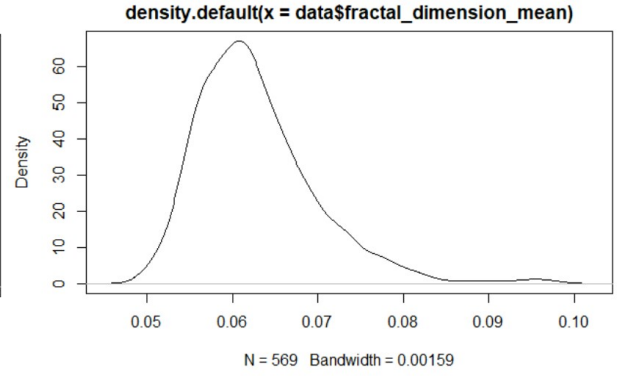
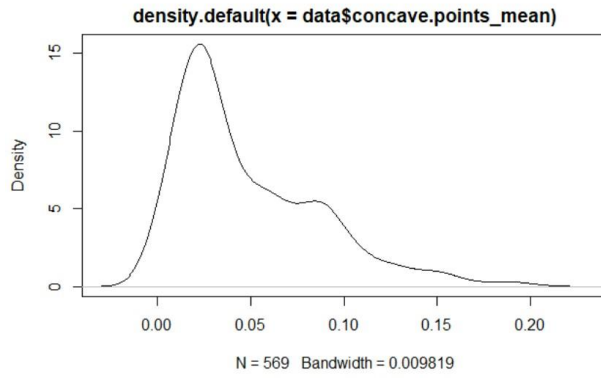
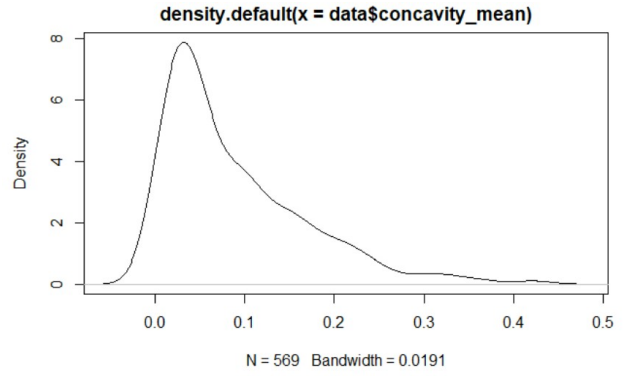
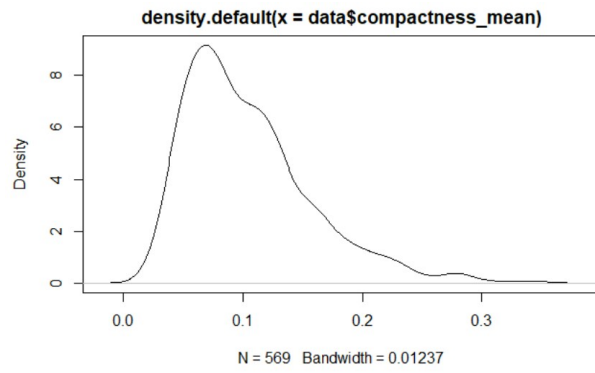
To see the correlation between the features we plotted the correlation matrix and heatmap.

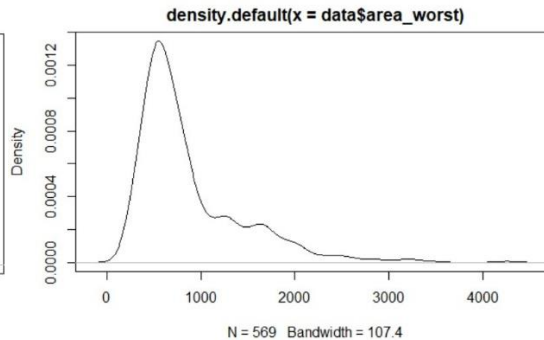
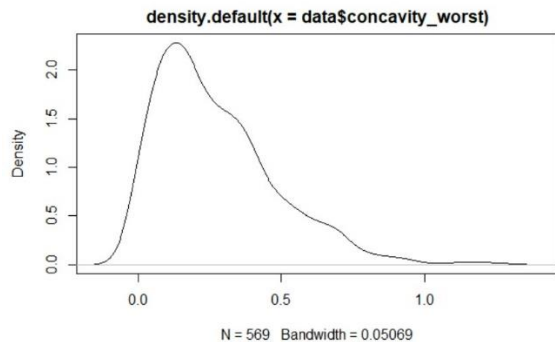
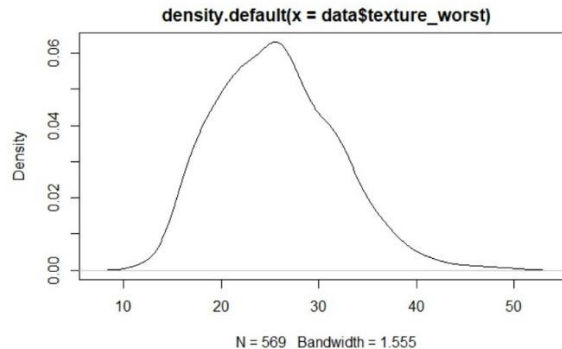
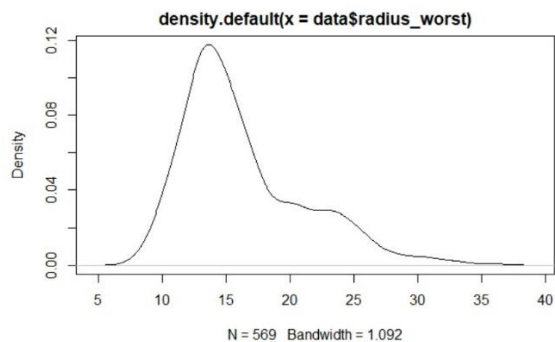
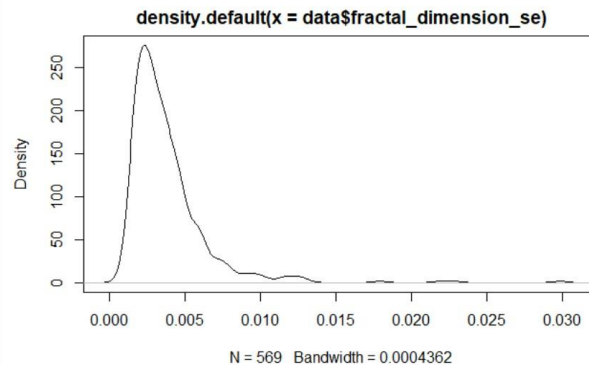
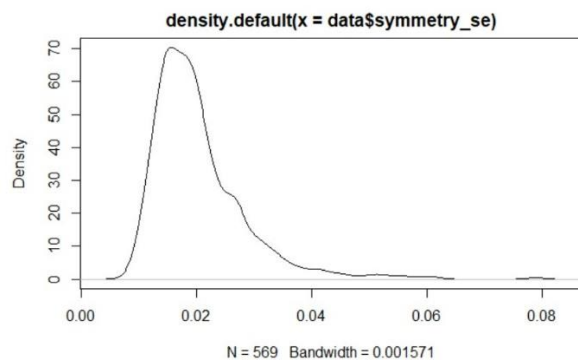
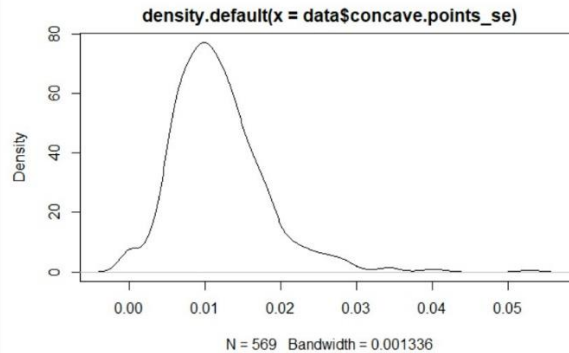
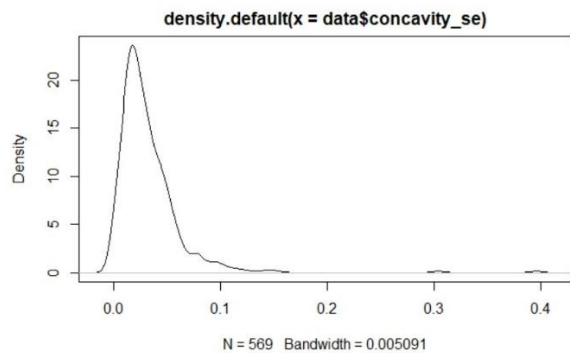


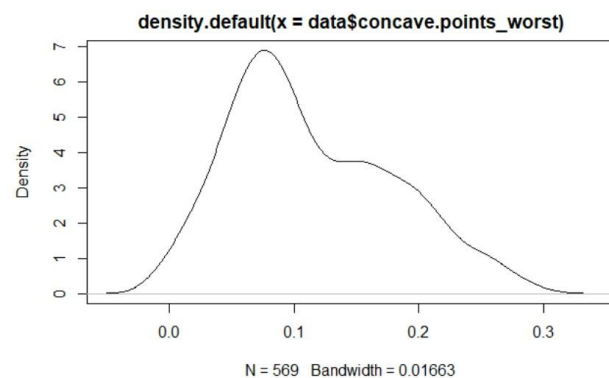
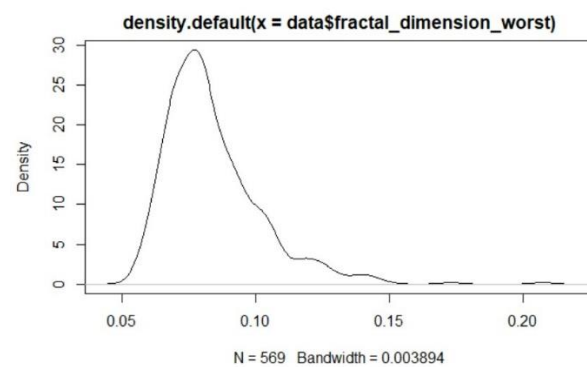
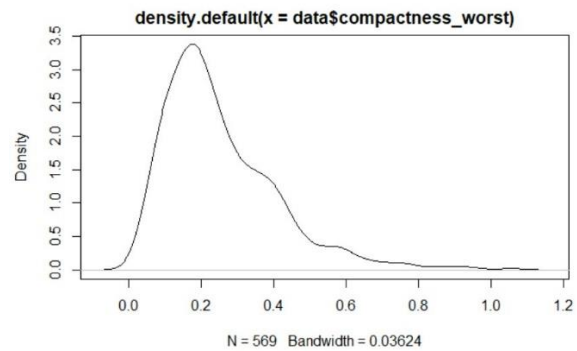
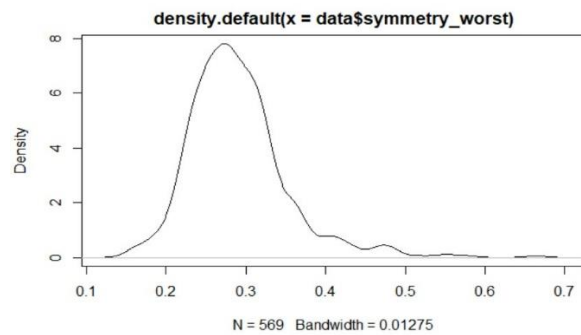
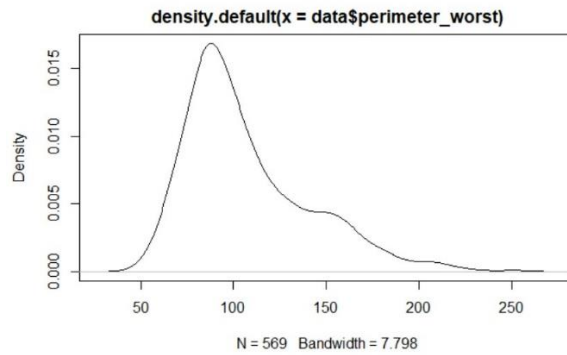
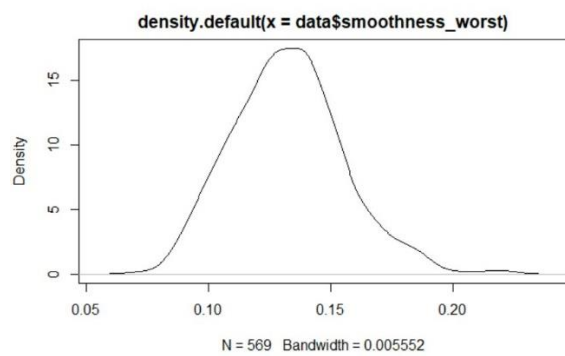


Then in order to properly fit the features we needed to know the density of the individual features.









## 5) Independent and Dependent feature?

: - Independent Feature: diagnosis

Dependent Feature:

radius\_mean  
 texture\_mean  
 perimeter\_mean  
 area\_mean  
 smoothness\_mean  
 compactness\_mean  
 concavity\_mean  
 concave points\_mean  
 symmetry\_mean  
 fractal\_dimension\_mean  
 radius\_se  
 texture\_se  
 perimeter\_se



```
area_se
smoothness_se
compactness_se
concavity_se
concave points_se
symmetry_se
fractal_dimension_se
radius_worst
texture_worst
perimeter_worst
area_worst
smoothness_worst
compactness_worst
concavity_worst
concave points_worst
symmetry_worst
fractal_dimension_worst
```

## 6) Why and how selection/engineering/scaling were performed?

∴ Feature scaling had to be done as the features ranged from magnitude of  $10^1$  to  $10^3$  so using `set.seed(1234)`

```
data_index <- createDataPartition(data$diagnosis, p=0.7, list = FALSE)
train_data <- data[data_index, -1]
test_data <- data[-data_index, -1]
```

```
fitControl <- trainControl(method="cv",
  number = 5,
  preProcOptions = list(thresh = 0.9),
  classProbs = TRUE,
  summaryFunction = twoClassSummary)
```

Feature scaling was done.

## 7) Which Classifier was chosen and why?

Random Forest Classifier was chosen on this dataset as its accuracy was found to be better than Naïve Bayes (94.71% which is greater than NB 91.3%). However KNN and boosted tree Accuracy (97%) was more than Random Forest, after hyper parameter tuning we obtained the accuracy close to boosted tree method (96.47%).

