



Project Cover Sheet

Assignment Title:	Healthcare Dataset Stroke Prediction		
Assignment No:	Final Project 1	Date of Submission:	25 December 2023
Course Title:	Introduction To Data Science		
Course Code:	CSC4180	Section:	B
Semester:	Fall	2023-24	Course Teacher: Tohedul Islam

Declaration and Statement of Authorship:

1. I/we hold a copy of this Assignment/Case-Study, which can be produced if the original is lost/damaged.
2. This Assignment/Case-Study is my/our original work and no part of it has been copied from any other student's work or from any other source except where due acknowledgement is made.
3. No part of this Assignment/Case-Study has been written for me/us by any other person except where such collaboration has been authorized by the concerned teacher and is clearly acknowledged in the assignment.
4. I/we have not previously submitted or currently submitting this work for any other course/unit.
5. This work may be reproduced, communicated, compared and archived for the purpose of detecting plagiarism.
6. I/we give permission for a copy of my/our marked work to be retained by the Faculty for review and comparison, including review by external examiners.
7. I/we understand that Plagiarism is the presentation of the work, idea or creation of another person as though it is your own. It is a form of cheating and is a very serious academic offence that may lead to expulsion from the University. Plagiarized material can be drawn from, and presented in, written, graphic and visual form, including electronic data, and oral presentations. Plagiarism occurs when the origin of them arterial used is not appropriately cited.
8. I/we also understand that enabling plagiarism is the act of assisting or allowing another person to plagiarize or to copy my/our work.

* Student(s) must complete all details except the faculty use part.

** Please submit all assignments to your course teacher or the office of the concerned teacher.

Group Name/No.: 19

No	Name	ID	Program	Signature
1	MOHAMMAD BIN HARUN	21-44583-1	BSc [CSE]	
2	MD. HARUN OR RASHID	21-44586-1	BSc [CSE]	

Faculty use only

FACULTY COMMENTS	Marks Obtained	
	Total Marks	

Description of Dataset:

The Healthcare Dataset on stroke the dataset is used analysis from the specific method. On this dataset predicted stroke occurrence based on naïve bayes classification. Before using model classification first prepare the dataset to handle missing value and outliers. Using the Pearson's Chi-squared test find the significant relationship between attribute and class attribute (stroke). If there is no significant relationship found for the attribute remove the attribute from the dataset because there is no relationship with the class attribute (stroke). After all that using the naïve bayes classification model for the train data and predict the test data. Find the predictive accuracy of the naïve bayes classifier. And generate the confusion matrix calculate the recall, precision and f-measure value of classifier.

Attributes:

The dataset contains several attributes that are various aspects related to stroke. The key attributes:

1. id: This id is unique identifier for each instance.
2. gender: gender of the individuals in the dataset.
3. age: age of individuals dataset
4. hypertension: Binary value of an individuals has hypertension (1) or not (0).
5. heart disease: Binary value of an individuals has a heart disease (1) or not (0).
6. ever married: ever married of individuals has ever been married (Yes) or not (No).
7. work type: The type of work individual is engaged in categorizing the occupation.
8. Residence type: Binary indicator of whether an individual resides in urban or rural area.
9. smoking status: The smoking habits of individuals categorizing as formerly smoked, never smoked or smokes.
10. BMI: Body Mass Index. Measure of body fat based on height and weight.
11. Average glucose level: The average glucose level in blood.
12. Stroke: Binary indicator of an individuals had a stroke (1) or not (0). And also, this is a class attribute.

Data Preprocessing:

For Dataset are preprocessing clean data involved handle missing values, outliers. Imputation of missing values such as fill the mean for BMI. Outlier's removal based on the interquartile range (IQR) method.

Library called:

`library(readxl), library(dplyr), library(e1071), library(caret)`

Description:

The library "readxl" is called for reading the excel file. The library "dplyr" is called for manipulation the data. The library "e1071" is called that provide statistical learning, naïve bayes classifier. And the library "caret" called for provide training and evaluating model.

Read dataset:

Read the dataset from the excel file using the `read_excel()`. And show the dataset on the terminal.

Code	
<pre>stroke<- read_excel("E:/AIUB/9th semester/Data Science/Final/Dataset/healthcare-dataset-stroke-data.xlsx") print(stroke)</pre>	
Output	Description
<pre># A tibble: 5,110 x 12 id gender age hypertension heart_disease ever_married work_type Residence_type avg_glucose_level bmi <dbl> <chr> <dbl> <dbl> <dbl> <chr> <chr> <dbl> <chr> 1 9046 Male 67 0 1 Yes Private Urban 229. 36.6 2 51676 Female 61 0 0 Yes Self-employ... Rural 202. N/A 3 31112 Male 80 0 1 Yes Private Rural 106. 32.5 4 60182 Female 49 0 0 Yes Private Urban 171. 34.4 5 1665 Female 79 1 0 Yes Self-employ... Rural 174. 24 6 56669 Male 81 0 0 Yes Private Urban 186. 29 7 53882 Male 74 1 1 Yes Private Rural 70.1 27.4 8 10434 Female 69 0 0 No Private Urban 94.4 22.8 9 27419 Female 59 0 0 Yes Private Rural 76.2 N/A 10 60491 Female 78 0 0 Yes Private Urban 58.6 24.2 # i 5,100 more rows # i 2 more variables: smoking_status <chr>, stroke <dbl> # i Use `print(n = ...)` to see more rows</pre>	Declare a stroke variable to the read dataset file from the directory folder and using the print function show the dataset.

Attribute names:

To get the dataset attribute using the names () function which also be set the object.

Code	
<pre>names(stroke)</pre>	
Output	Description
<pre>[1] "id" "gender" "age" "hypertension" "heart_disease" [6] "ever_married" "work_type" "Residence_type" "avg_glucose_level" "bmi" [11] "smoking_status" "stroke"</pre>	This returns the all-attributes names for the dataset.

Summary

To obtain different types of attributes.

Code					
summary(stroke)					
Output					Description
<pre> id gender age hypertension heart_disease ever_married Min. : 67 Length:5110 Min. : 0.08 Min. :0.00000 Min. :0.00000 Length:5110 1st Qu.:17741 Class :character 1st Qu.:25.00 1st Qu.:0.00000 1st Qu.:0.00000 Class :character Median :36932 Mode :character Median :45.00 Median :0.00000 Median :0.00000 Mode :character Mean :36518 Mean :43.23 Mean :0.09746 Mean :0.05401 3rd Qu.:54682 3rd Qu.:61.00 3rd Qu.:0.00000 3rd Qu.:0.00000 Max. :72940 Max. :82.00 Max. :1.00000 Max. :1.00000 work_type Residence_type avg_glucose_level bmi smoking_status stroke Length:5110 Length:5110 Min. : 55.12 Length:5110 Length:5110 Min. :0.00000 Class :character Class :character 1st Qu.: 77.25 Class :character Class :character 1st Qu.:0.00000 Mode :character Mode :character Median : 91.89 Mode :character Mode :character Median :0.00000 Mean :106.15 Mean :0.04873 3rd Qu.:114.09 3rd Qu.:0.00000 Max. :271.74 Max. :1.00000 </pre>					This will provide a summary of each attribute of dataset stroke including measures minimum, 1 st quantile, 3 rd quantile, median, mean and maximum for numerical attributes as count for factors.

Dataset structure:

To display the structure of an R objects, use the str(stroke).

Code					
str(stroke)					
Output					Description
<pre> tibble [5,110 × 12] (S3: tbl_df/tbl/data.frame) \$ id : num [1:5110] 9046 51676 31112 60182 1665 ... \$ gender : chr [1:5110] "Male" "Female" "Male" "Female" ... \$ age : num [1:5110] 67 61 80 49 79 81 74 69 59 78 ... \$ hypertension : num [1:5110] 0 0 0 0 1 0 1 0 0 0 ... \$ heart_disease : num [1:5110] 1 0 1 0 0 0 1 0 0 0 ... \$ ever_married : chr [1:5110] "Yes" "Yes" "Yes" "Yes" ... \$ work_type : chr [1:5110] "Private" "Self-employed" "Private" "Private" ... \$ Residence_type : chr [1:5110] "Urban" "Rural" "Rural" "Urban" ... \$ avg_glucose_level: num [1:5110] 229 202 106 171 174 ... \$ bmi : chr [1:5110] "36.6" "N/A" "32.5" "34.4" ... \$ smoking_status : chr [1:5110] "formerly smoked" "never smoked" "never smoked" "smokes" ... \$ stroke : num [1:5110] 1 1 1 1 1 1 1 1 1 1 ... > </pre>					This str will be show the data types and some values of the dataset attributes.

Conversion:

BMI has some N/A which are contain in categorized. BMI is numeric value for Body Mass Index. Measure of body fat based on height and weight.

Code:

```
stroke$bmi<- as.integer(stroke$bmi) stroke$hypertension<-as.factor(stroke$hypertension) stroke$heart_disease<-
as.factor(stroke$heart_disease) stroke$stroke<-as.factor(stroke$stroke)
```

Description:

Convert the BMI attribute char to integer. Convert the hypertension, heart_disease and stroke attribute numerical data to factor level.

Missing value:

Check is the dataset has any missing instance or not.

Code						
colSums(is.na(stroke))						
Output						Description
id	gender	age	hypertension	heart_disease	ever_married	On this dataset bmi attribute has 201 missing values found.
0	0	0	0	0	0	
work_type	Residence_type	avg_glucose_level	bmi	smoking_status	stroke	
0	0	0	201	0	0	

Handle missing values:

For the bmi attribute missing values are handle using measure of central tendency: mean

Code						
stroke\$bmi[is.na(stroke\$bmi)] <- mean(stroke\$bmi, na.rm = TRUE)						
colSums(is.na(stroke))						
Output						Description
id	gender	age	hypertension	heart_disease	ever_married	Replace all missing value using the mean.
0	0	0	0	0	0	
work_type	Residence_type	avg_glucose_level	bmi	smoking_status	stroke	
0	0	0	0	0	0	

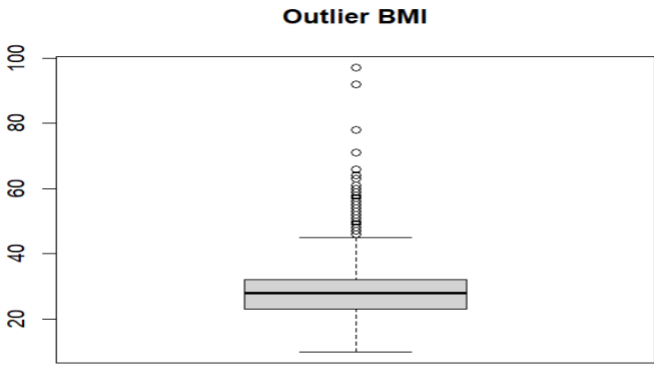
Noisy value:

Now since outliers are an issue terms of getting accurate results from this dataset, we remove them and replace them with the mean values of those attribute columns.

Box plot:

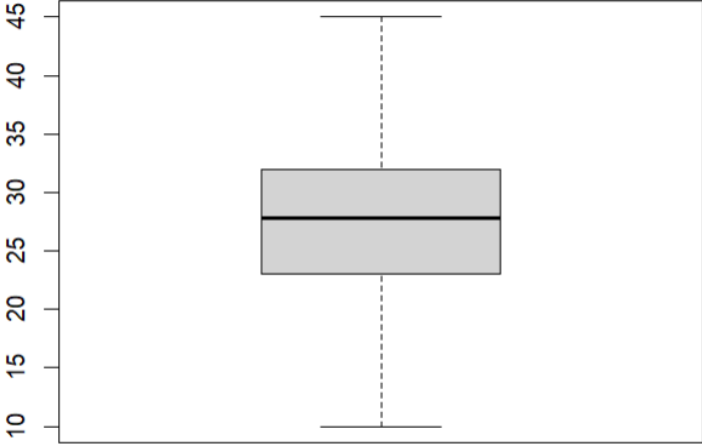
To Finding the outliers first using the box-plot graphical representation. After that remove the outliers.

BMI

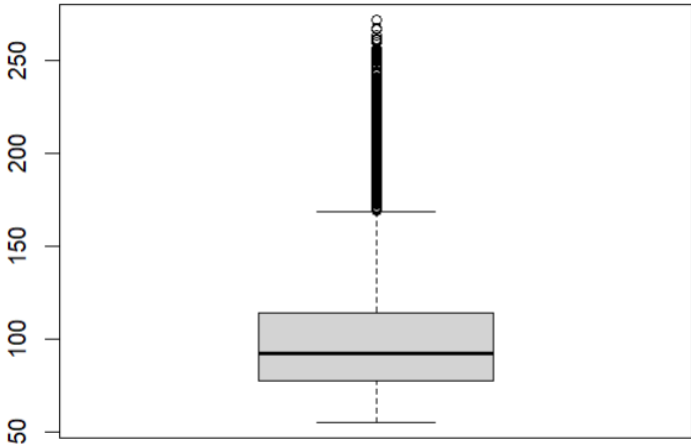
Code	
<pre>boxplot(stroke\$bmi,main="Outlier BMI")</pre>	
Output	Description
	Generates a box plot for the 'bmi' attribute in the dataset, providing a visual representation of its distribution and identifying potential outliers.

Handle BMI outlier

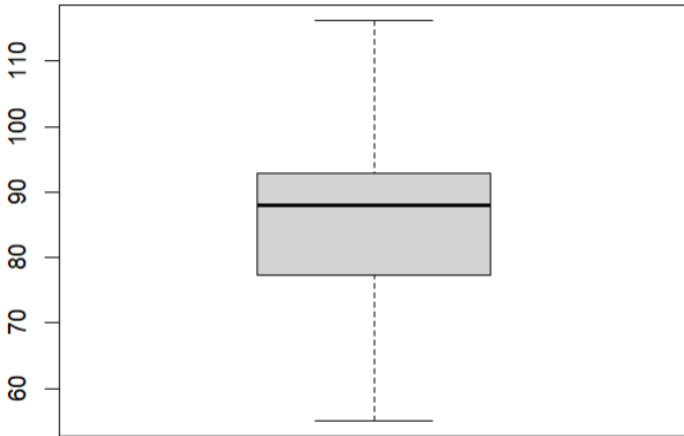
Code
<pre>q1<-quantile(stroke\$bmi,0.25) q3<-quantile(stroke\$bmi,0.75) iqr<-q3-q1 lower<- q1-1.5*iqr upper<- q3+1.5*iqr outliers_bmi<- stroke\$bmi < lower stroke\$bmi > upper stroke\$bmi <- ifelse(outliers_bmi,NA,stroke\$bmi) colSums(is.na(stroke)) stroke\$bmi[is.na(stroke\$bmi)]<-mean(stroke\$bmi,na.rm = TRUE) boxplot(stroke\$bmi, main='Remove Outlier BMI')</pre>

Output	Description
<p style="text-align: center;">Remove Outlier BMI</p> 	<p>To handle the outliers, computes the first quartile (q1), third quartile (q3), and interquartile range (iqr) for 'bmi', identifies outliers based on a iqr criterion, and replaces those outliers with NA, then replace the NA values with the mean of the 'bmi' attribute.</p>

Average glucose level

Code	
<pre>boxplot(stroke\$avg_glucose_level,main="Outlier of Avg Glucose level")</pre>	
Output	Description
<p style="text-align: center;">Outlier of Avg Glucose level</p> 	<p>Generates a box plot for the ' avg_glucose_level ' attribute in the dataset, providing a visual representation of its distribution and identifying potential outliers.</p>

Handle average glucose level

Code	
<pre>q1 <- quantile(stroke\$avg_glucose_level,0.25) q3 <- quantile(stroke\$avg_glucose_level,0.75) iqr <- q3 - q1 lower <- q1-1.5*iqr upper<- q3+1.5*iqr outliers_avg_glucose_level<- stroke\$avg_glucose_level < lower stroke\$avg_glucose_level > upper stroke\$avg_glucose_level <- ifelse(outliers_avg_glucose_level,NA,stroke\$avg_glucose_level) colSums(is.na(stroke)) stroke\$avg_glucose_level[is.na(stroke\$avg_glucose_level)]<-mean(stroke\$avg_glucose_level,na.rm = TRUE) colSums(is.na(stroke)) boxplot(stroke\$avg_glucose_level,main="Remove Outlier Avg Glucose level")</pre>	
Output	Description
<p>Remove Outlier Avg Glucose level</p> 	<p>To handle the outliers, computes the first quartile (q1), third quartile (q3), and interquartile range (iqr) for 'avg_glucose_level', identifies outliers based on a iqr criterion, and replaces those outliers with NA, then replace the NA values with the mean of the 'avg_glucose_level' attribute.</p>

Analysis Technique

Statistical correlation:

Using correlation technique Pearson correlation test for find the significant attribute along with target attribute. If it didn't find any significant there will be removed the attribute from the dataset.

Code

```
Pearsons_Chi_squared<- function(attribute){  
  count_instances<- table(stroke[[attribute]],stroke$stroke)  
  p_value<- chisq.test(count_instances)  
  return(p_value)  
}  
  
pearson_correlation_id <- Pearsons_Chi_squared('id')  
pearson_correlation_gender <- Pearsons_Chi_squared('gender')  
pearson_correlation_age <- Pearsons_Chi_squared('age')  
pearson_correlation_hypertension <- Pearsons_Chi_squared('hypertension')  
pearson_correlation_heart_disease <- Pearsons_Chi_squared('heart_disease')  
pearson_correlation_ever_married <- Pearsons_Chi_squared('ever_married')  
pearson_correlation_work_type <- Pearsons_Chi_squared('work_type')  
pearson_correlation_Residence_type <- Pearsons_Chi_squared('Residence_type')  
pearson_correlation_smoking_status <- Pearsons_Chi_squared('smoking_status')  
pearson_correlation_bmi <- Pearsons_Chi_squared('bmi')  
pearson_correlation_avg_glucose_level <- Pearsons_Chi_squared('avg_glucose_level')  
pearson_correlation_stroke <- Pearsons_Chi_squared('stroke')  
  
if(pearson_correlation_id$p.value<0.05){  
  print("There is a significant relationship between id and stroke.")  
}else{  
  stroke<- subset(stroke,select = -id)  
  print("No significant relationship found.")  
}  
  
if(pearson_correlation_gender$p.value<0.05){  
  print("There is a significant relationship between gender and stroke.")
```

```
}else{

stroke<- subset(stroke,select = -gender)
print("No significant relationship found.")
}

if(pearson_correlation_age$p.value<0.05){
print("There is a significant relationship between age and stroke.")
}else{
stroke<- subset(stroke,select = -age)
print("No significant relationship found.")
}

if(pearson_correlation_hypertension$p.value<0.05){
print("There is a significant relationship between hypertension and stroke.")
}else{
stroke<- subset(stroke,select = -hypertension)
print("No significant relationship found.")
}

if(pearson_correlation_heart_disease$p.value<0.05){
print("There is a significant relationship between heart_disease and stroke.")
}else{
stroke<- subset(stroke,select = -heart_disease)
print("No significant relationship found.")
}

if(pearson_correlation_ever_married$p.value<0.05){
print("There is a significant relationship between ever married and stroke.")
}else{
stroke<- subset(stroke,select = -ever_married)
print("No significant relationship found.")
}
```

```

}
if(pearson_correlation_work_type$p.value<0.05){
  print("There is a significant relationship between work type and stroke.")
}else{
  stroke<- subset(stroke,select = -work_type)
  print("No significant relationship found.")
}
if(pearson_correlation_Residence_type$p.value<0.05){
  print("There is a significant relationship between residence type and stroke.")
}else{
  stroke<- subset(stroke,select = -Residence_type)
  print("No significant relationship found.")
}
if(pearson_correlation_smoking_status$p.value<0.05){
  print("There is a significant relationship between smoking status and stroke.")
}else{
  stroke<- subset(stroke,select = -smoking_status)
  print("No significant relationship found.")
}
if(pearson_correlation_bmi$p.value<0.05){
  print("There is a significant relationship between bmi and stroke.")
}else{
  stroke<- subset(stroke,select = -bmi)
  print("No significant relationship found.")
}
if(pearson_correlation_avg_glucose_level$p.value<0.05){
  print("There is a significant relationship between avg glucose level and stroke.")
}else{
  stroke<- subset(stroke,select = -avg_glucose_level)
  print("No significant relationship found.")
}

```

```

if(pearson_correlation_stroke$p.value<0.05){
  print("There is a significant relationship between stroke and stroke.")
}
else{
  stroke<- subset(stroke,select = -stroke)
  print("No significant relationship found.")
}

```

Output	Description
<pre> > pearson_correlation_stroke <- Pearsons_Chi_squared('stroke') > if(pearson_correlation_id\$p.value<0.05){ + print("There is a significant relationship between id and stroke.") + }else{ + stroke<- subset(stroke,select = -id) + print("No significant relationship found.") + } [1] "No significant relationship found." > if(pearson_correlation_gender\$p.value<0.05){ + print("There is a significant relationship between gender and stroke.") + }else{ + stroke<- subset(stroke,select = -gender) + print("No significant relationship found.") + } [1] "No significant relationship found." > if(pearson_correlation_age\$p.value<0.05){ + print("There is a significant relationship between age and stroke.") + }else{ + stroke<- subset(stroke,select = -age) + print("No significant relationship found.") + } [1] "There is a significant relationship between age and stroke." > if(pearson_correlation_hypertension\$p.value<0.05){ + print("There is a significant relationship between hypertension and stroke.") + }else{ + stroke<- subset(stroke,select = -hypertension) + print("No significant relationship found.") + } [1] "There is a significant relationship between hypertension and stroke." > if(pearson_correlation_heart_disease\$p.value<0.05){ + print("There is a significant relationship between heart_disease and stroke.") + }else{ + stroke<- subset(stroke,select = -heart_disease) + print("No significant relationship found.") + } [1] "There is a significant relationship between heart_disease and stroke." </pre>	<p>Here create a function <code>pearson_correlation_stroke</code> function using the Pearson correlation for significant relationship. Check the condition with p value if the p value less 0.05 then there will be significant relationship otherwise no significant found.</p>

Data splitting:

Code
<pre> set.seed(123) split_Index<- createDataPartition(stroke\$stroke, p = 0.8, list = FALSE) train_data<- stroke[split_Index,] test_data <- stroke[-split_Index,] </pre>

Output	Description
<p>train data:</p> <pre># A tibble: 4,089 × 8 age hypertension heart_disease ever_married work_type bmi smoking_status stroke <dbl> <fct> <fct> <chr> <chr> <dbl> <chr> <fct> 1 67 0 1 Yes Private 36 formerly smoked 1 2 80 0 1 Yes Private 32 never smoked 1 3 49 0 0 Yes Private 34 smokes 1 4 79 1 0 Yes Self-employed 24 never smoked 1 5 81 0 0 Yes Private 29 formerly smoked 1 6 74 1 1 Yes Private 27 never smoked 1 7 69 0 0 No Private 22 never smoked 1 8 59 0 0 Yes Private 28.4 Unknown 1 9 78 0 0 Yes Private 24 Unknown 1 10 81 1 0 Yes Private 29 never smoked 1 # i 4,079 more rows # i Use `print(n = ...)` to see more rows</pre> <p>test data:</p> <pre># A tibble: 1,021 × 8 age hypertension heart_disease ever_married work_type bmi smoking_status stroke <dbl> <fct> <fct> <chr> <chr> <dbl> <chr> <fct> 1 61 0 0 Yes Self-employed 28.4 never smoked 1 2 79 0 1 Yes Private 28 never smoked 1 3 50 1 0 Yes Self-employed 30 never smoked 1 4 60 0 0 No Private 37 never smoked 1 5 57 0 1 No Govt_job 28.4 Unknown 1 6 52 1 0 Yes Self-employed 27.8 never smoked 1 7 82 0 1 Yes Private 32 Unknown 1 8 69 0 1 Yes Self-employed 28 smokes 1 9 72 1 0 Yes Private 23 formerly smoked 1 10 76 1 0 Yes Private 33 never smoked 1 # i 1,011 more rows # i Use `print(n = ...)` to see more rows</pre>	<p>This process helps in creating a random and representative split of the dataset for training and testing a predictive model. The list = FALSE argument ensures that the output is a vector of indices rather than a list.</p>

Naïve bayes model:

Naïve bayes model using for applied a predicting target data from the test data.

Code	
<pre>nb_model <- naiveBayes(stroke ~ ., data = train_data) predictions <- predict(nb_model, test_data) table(predictions)</pre>	
Output	Description
<pre>predictions 0 1 974 47</pre>	<p>Here predict the test data using the naïve bayes classification among the train data.</p>

Accuracy: calculate the naïve bayes predictive accuracy using the divided the data into training and test set.

Code
<pre>accuracy <- sum(predictions == test_data\$stroke) / nrow(test_data) cat("Naive Bayes Accuracy:", round(accuracy,2), "\n")</pre>

Output	Description
Naïve Bayes Accuracy: 0.93	The naïve bayes predictive accuracy of training and test set is 0.93. the round operation can configure the value in 2 digits.

Confusion Matrix:

To generate a confusion matrix for dataset using naïve bayes classification and calculate the recall, precision and f-measure of the classifier.

Code	
<pre> confusion_matrix<- table(predictions,test_data\$stroke) recall_matrix<- confusion_matrix[2,2]/sum(confusion_matrix[2,]) precision_matrix<- confusion_matrix[2,2]/sum(confusion_matrix[,2]) f_measure <- 2 * (precision_matrix * recall_matrix) / (precision_matrix + recall_matrix) cat("Confusion Matrix: ",confusion_matrix,"\n") cat("Recall: ",recall_matrix_second_class,"\n") cat("Precision: ",precision_matrix_second_class,"\n") cat("f-measure: ",f_measure,"\n") </pre>	
Output	Description
<pre> > cat("Confusion Matrix: ",confusion_matrix,"\n") Confusion Matrix: 936 36 38 11 > cat("Recall: ",recall_matrix_second_class,"\n") Recall: 0.2340426 > cat("Precision: ",precision_matrix_second_class,"\n") Precision: 0.962963 > cat("f-measure: ",f_measure,"\n") f-measure: 0.2291667 </pre>	The output shown the calculated confusion matrix, recall, precision and f-measure value of classifier.