

**A
PROJECT REPORT**

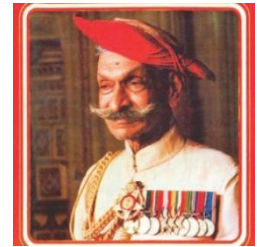
ON

**“The Statistical Analysis Of Comfortable Seat And
Check The Health Status Of Driver”**

Submitted to

DEPARTMENT OF STATISTICS,

**SHREE SHAHAJI CHH. MAHAVIDYALAYA,
KOLHAPUR**



For the partial fulfillment of degree

B.Sc. IN STATISTICS

By

- **MR. JADHAV RAM DEVIDAS**
- **MISS. JADHAV SHARWARI BHIKAJI**
- **MISS. KUMBHAR PRANALEE VILAS**
- **MISS. REDDY PRATYUSHA VENKATSRINIVAS**
- **MISS. SUTAR ASHWINI NAMDEV**

Under the Guidance of

Mr. BHOSALE .A. B.

Mrs. PANWALKAR .S.M.

Mr. MALI .K. M.

DEPARTMENT OF STATISTICS,

**SHREE SHAHAJI CHH. MAHAVIDYALAYA,
KOLHAPUR.**

CERTIFICATE



This is to certify that this Project report on,
“The Statistical Analysis of Comfortable Seat and Check the
Health Status of Driver” is

Submitted by,

- MR. JADHAV RAM DEVIDAS
- MISS. JADHAV SHARWARI BHIKAJI
- MISS. KUMBHAR PRANALEE VILAS
- MISS. REDDY PRATYUSH VENKATSRINIVAS
- MISS. SUTAR ASHWINI NAMDEV

Under the guidance of **Mr. BHOSALE A. B. , Mrs. PANWALKAR S. M., Mr. MALI K. M.** to the B. Sc. Part III, Department of Statistics, Shree Shahaji Chh. Mahavidyalaya, Kolhapur has satisfactorily completed the project work during academic year 2018 -2019.

Date:-

Project Guide

Examiner

Head of Department

-: Acknowledgement:-

We have great pleasure while submitting this project report on “The Statistical analysis of comfortable seat and check the health status of driver” in partial fulfillment of B. Sc. III.

We thank Mr. Shinde sir (D. M. Of MSRTC) and Mr. Dhudhakar sir (Head of KMT Office) for giving permission for data collection. We also thank Dr. R. K. Shanediwan sir (Principal) for giving permission for doing the project.

We wish to thank all teachers, non-teaching staff and students of third year from Shree Shahaji Chh. Mahavidyalaya, Kolhapur for their co-operation in collecting the necessary data. We would like to express our profound gratitude to Mrs. Panwalkar S.M. madam for valuable guidance in completion of this project

We would like to thank Mr. Bhosale. A. B sir. , Head of Department of Statistics, Shree Shahaji Chh. Mahavidyalaya, Kolhapur for providing us necessary facilities .

We are also thankful to Mr. Mali K.M. for his timely suggestions and encouragement.

CONTENTS :-

Sr.No.	Contents	Page No.
1.	Introduction	5
2.	Objective	6
3.	Data Collection Method	7
4.	Statistical Tools	8
5.	Statistical Analysis	9-21
6.	Conclusion	22
7.	Bibliography	23
8.	Appendix	24-34

➤ Introduction:-

In our daily life we use car, travels ,bus, KMT for travelling and tractor ,truck for transporting goods. But the main work of this travelling or transporting is done by the driver of the respective vehicle.

Our aim is to find the adjusted level of comfortable seat while driving and to check their health status. Because he is the important factor of travelling and transporting. For this study we have taken the sample of 202 drivers who filled the questionnaire consisting 22 questions from various vehicles, And the questionnaire considered the various variable value and health related information. The number of forms collected from various vehicle are given below Car=43, Travels=23, Bus=39, KMT=41, Tractor= 20, Truck=36.

Using various statistical techniques, we analyzed the collected data and verified our objectives.

Objectives

- To study adjustment of comfortable seat for drivers.
- To estimate measurements of comfortable seat.
- To check the health status of drivers.

❖ **Data collection method**

The main objective of our project is to study the position of comfortable seat and to check the health status of driver in the Kolhapur City. So, for this project we have collected data from drivers of various types of vehicles like Cars, Buses, Trucks, Travels and Tractors in the Kolhapur City.

For this project we have prepared a questionnaire satisfying some characteristics of good questionnaire. In overall city nearly about 202 different types of vehicles (Car,Bus,Truck,Travel, Tractor) data has been collected. In this project we have collected some of the information of driver and some data about the seating level or measurements for better adjustment of seat for comfort of the driver. From the overall collection of the data we have been collected data of 41 Cars, 23 Travels, 39 Buses, 43 KMT, 20 Tractors, 36 Trucks.

The main objective of this project is to define seating comforts, consequently the response of the vehicle drivers was required. Hence a driver study was planned, to obtain general information of the drivers and their preferred seating adjustment of different types of vehicles.

❖ **Statistical Tools Used For Analysis.**

- **Graphical Presentation**

- 1) Pie-Chart
- 2) Bar Diagram

- **Statistical Test**

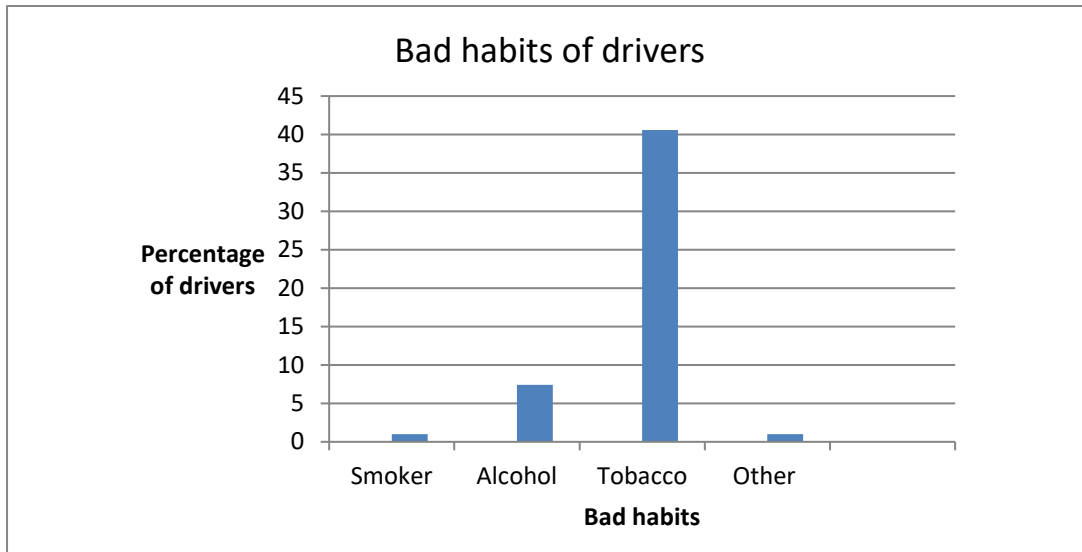
- 1) Z-test for testing single population mean and Proportion.
- 2) K-nearest neighbour classifier.
- 3) Naive Bayesian Classifier.
- 4) Logistic Regression model.

- **Softwares used**

1. MS-Excel
2. R-Software
3. MS-Word

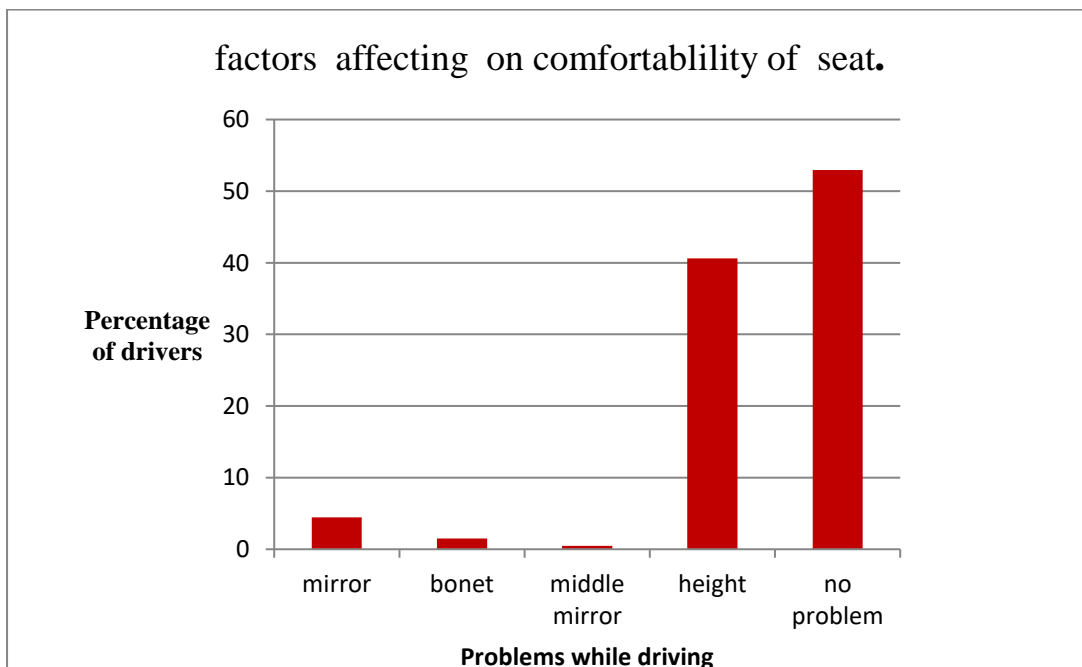
- **Statistical Analysis**

Bar diagram for percentage of bad habits of drivers



Conclusion:- 40.59% drivers have habit of tobacco chwibing and 7.43% drivers have habit of drinking alcohol.

- **Bar diagram for factors affecting comfortability of seat.**



Conclusion:- 1) 52.97% of drivers feel comfortable on seat.

2) In driving , there is maximum problem due to height of driver.

❖ Summary of vehicles:-

Here the study of average given by various vehicles (in km/liter), is as given below

▪ Car :

Min.	Q1.	Median	Mean	Q3.	Max.
10.0	12.5	14.0	15.82	18.00	28.0

- Here it is seen that, the average of cars is 15.82 km/liter, 25% of cars have average below 12.5 km/liter and 75 % cars have average below 18 km/liter.

▪ KMT:-

Min.	Q1.	Median	Mean	Q3.	Max.
3.0	3.60	4.0	4.332	4.500	8.0

- Here it is seen that, the average of KMT is 4.332 km /liter, 25% of KMT have average below 3.6 km/liter and 75 % KMT have average below 4.5 km/liter.

▪ ST. Bus:-

Min.	Q1.	Median	Mean	Q3.	Max.
3.50	4.0	4.5	4.345	4.775	5.0

- Here it is seen that, the average of ST.Bus is 4.345 km /liter, 25% of ST.Bus have average below 4 km/liter and 75 % ST.Bus have average below 4.775km/liter.

▪ **Tractor:**

Min.	Q1.	Median	Mean	Q3.	Max.
3.00	6.75	7.0	7.60	9.0	14.0

- Here it is seen that, the average of Tractor is 7.60 km /liter, 25% of Tractor have average below 6.75 km/liter and 75 % Tractor have average below 9 km/liter.

▪ **Travels:-**

Min.	Q3.	Median	Mean	Q3.	Max.
2.500	3.0	4.0	5.239	4.75	12.

- Here it is seen that, the average of Travels is 5.239 km /liter, 25% of Travels have average below 3km/liter and 75 % Travels have average below 4.750km/liter.

▪ **Truck:-**

Min.	Q1.	Median	Mean	Q3.	Max.
3.0	3.5	4.0	4.306	5.0	7.5

- Here it is seen that, the average of Truck is 4.306 km /liter, 25% of Truck have average below 3.50 km/liter and 75 % Truck have average below 5km/liter.

➤ **Testing Proportion of drivers having no disease is 0.95 or not.**

Hypothesis

We have to test the hypothesis,

H_0 : $P = 0.95$ i.e. 95% drivers have no disease.

H_1 : $P \neq 0.95$ i.e. 95% drivers have disease.

p = Sample proportion of drivers having no disease.

P_0 = Specified value of population proportion of drivers have no disease = 0.95

Observation table:

Suffered from disease	Observed Sample	Observed Proportion (p)	Test Proportion
Yes	7	0.04	0.95
No	195	0.96	
Total	202	1	

Calculation:

Sample size	Z-cal	P-value	95% C.I. for P
202	0.70459	0.4013	(0.9269,0.9847)

Let l.o.s is $\alpha = 0.05$,

As p value is greater than α (0.05), Hence we accept H_0 at 5 % l.o.s.

Conclusion:

95% drivers do not suffer from any disease.

➤ **Testing Proportion of drivers having no pain is 0.8 or more.**

Hypothesis ,

We have to test the hypothesis,

$H_0: P = 0.80$ i.e 80% drivers have no pain.

$H_1: P > 0.80$ i.e more than 80% drivers have no pain.

p = Sample proportion of drivers having no pain.

P_0 = Specified value of population proportion of drivers having no pain = 0.8

Observation table:

Suffered from pain	Observed Sample	Observed Proportion (p)	Test Proportion
Yes	28	0.14	0.80
No	174	0.86	
Total	202	1	

Calculation:

Sample size	Z-cal	P-value	95% C.I. for P
202	4.3815	0.03633	(0.8041 , 0.9044)

Let l.o.s is $\alpha = 0.05$,

As p value is less than alpha (0.05) reject H_0 at 5 % l.o.s.,

Conclusion:

More than 80% drivers have no pain.

➤ **Testing Proportion of drivers having no surgical history is 0.75 or more .**

Hypothesis,

We have to test the hypothesis,

$H_0: P = 0.75$ i.e. 75% drivers have no surgical history.

$H_1: P > 0.75$ i.e. more than 75% drivers have no surgical history.

p = sample proportion of drivers having no surgical history.

P_0 = Specified value of population proportion of drivers having no surgical history = 0.75

Observation table:-

Suffered from surgical history	Observed Sample	Observed Proportion (p)	Test Proportion
Yes	27	0.14	0.75
No	175	0.86	
Total	202	1	

Calculation:-

Sample size	Z-cal	P-value	95% C.I. for P
202	13.967	0.000186	(0.8097 , 0.9086)

Let l.o.s is $\alpha = 0.05$.

As p value is less than alpha (0.05) we reject H_0 at 5 % l.o.s.

Conclusion:

More than 75% drivers have no surgical history.

➤ **Testing Proportion of drivers having bad habits is 0.5 or more.**

Hypothesis :-

We have to test the hypothesis,

$H_0: P = 0.50$ i.e. 50% drivers have bad habits .

$H_1: P > 0.50$ i.e. more than 50% drivers have bad habits.

p = sample proportion of drivers having bad habits.

P_0 = Specified value of population proportion of drivers having bad habits = 0.5

Observation table:-

Suffered from bad habits	Observed Sample	Observed Proportion (p)	Test Proportion
Yes	101	0.50	0.50
No	101	0.50	
Total	202	1	

Calculation:-

Sample size	Z-cal	P-value	95% C.I. for P
202	0	1	(0.4316 0.5683)

Let l.o.s is $\alpha = 0.05$.

As p value is greater than alpha (0.05) we accept H_0 . at 5 % l.o.s..

Conclusion:

more than 50% drivers have bad habits.

➤ **Testing Proportion of drivers don't feeling comfortable on seat is 0.5 or less .**

P_0 = Specified value of population proportion of drivers don't feeling comfortable on seat = 0.5

Hypothesis :-

We have to test the hypothesis,

$H_0: P = 0.50$ i.e 50% drivers don't feel comfortable on seat.

$H_1: P < 0.50$ i.e less than 50% drivers don't feel comfortable on seat.

p = sample proportion of drivers feeling comfortable on seat.

P_0 = Specified value of population proportion of drivers don't feeling comfortable on seat = 0.5

Observation table:-

Suffered from comfortable seat	Observed Sample	Observed Proportion (p)	Test Proportion
Yes	108	0.53	0.50
No	94	0.47	
Total	202	1	

Calculation:-

Sample size	Z-cal	P-value	95% C.I. for P
202	0.8366	0.3604	(0.3954, 0.5366)

Let l.o.s is $\alpha = 0.05$.

As p value is greater than alpha (0.05). Hence we accept H_0 at 5 % l.o.s.

Conclusion:-

50% drivers don't feel comfortable on seat.

➤ **Testing daily average distance driven by drivers is 250 or not.**

μ = daily average distance driven by drivers

Hypothesis :-

We have to test the hypothesis,

$H_0: \mu = 250$ i.e. daily average distance driven by drivers is 250 km.

$H_1: \mu \neq 250$ i.e. daily average distance driven by drivers is not 250 km.

X: The daily distance driven by the drivers

μ_0 = Specified value of population mean = 250

$n = 202$

$\bar{X} = 278.302, \quad s = 192.2513$

Under H_0 test statistic is

$$Z = \frac{\bar{X} - \mu_0}{\frac{s}{\sqrt{n}}} \sim N(0,1)$$

$$Z = \frac{278.302 - 250}{\frac{192.2513}{\sqrt{202}}} = 0.00729$$

$$|Z| = 0.00729$$

Critical value at 5% l.o.s $Z_{\alpha/2} = 1.96$

$$|Z| > Z_{\alpha/2}$$

Conclusion:-

Daily average distance driven by drivers is 250 km.

➤ **Classification Model Building**

Our data contains 202 observations and 37 variables for these classification model purpose we select comparability with seat (comfort “Yes” or “No”) is a response variable and predictors as (Y= response variable, X1=distance between seat to bottom, X2=Breadth of seat, X3=Horizontal length of seat, X4=Vertical length of seat, X5=distance between steering to seat, X6=distance between clutch to seat, X7= distance between gear to seat, X8= distance between brake to seat).

we used KNN, Naive Bays and logistic regression classifier algorithms. These models have been selected for this study because of their popularity in the recent literature. We first give a short description of these classification models.

k-nearest neighbour classifier

The k-nearest-neighbour method was first described in 1950s. Nearest-neighbour classifiers are based on learning by analogy, that is, by comparing a given test tuple with training tuples that are similar to it. The training tuples are described by n attributes. Each tuple represents a point in an n-dimensional space. In this way, all the training tuples are stored in an n-dimensional pattern space. When given an unknown tuple, a k-nearest-neighbour classifier searches the pattern space for the k training tuples that closest to the unknown tuple. These k training tuples are the k “nearest neighbours” of the unknown tuple.

Naive Bayesian Classifier:-

The Naive Bayes algorithm is a simple probabilistic classifier that calculates a set of probabilities by counting the frequency and combinations of values in a given dataset. The algorithm uses Bayes theorem and assumes all attributes to be independent given the value of class variable. This conditional independence assumption rarely holds true in real world applications, hence the characterization as Naive yet the algorithm tends to perform well and learn rapidly in various supervised classification problems. Naive Bayes classifier is based on Bayes theorem and the theorem of total probability.

In this classifier we compute the conditional probability
 $P(C_j/X)$ and

assign X to those class C_i having large probability i.e. $X \in C_j$ if $P(C_j/X) > P(C_i/X)$ for all $i \neq j = 1, 2, \dots, m$.

Logistic Regression model:-

Logistic regression is the generalization of linear regression. It is used primarily for predicting binary or multiclass dependent variables. Because the response variable is discrete, it cannot be modeled directly by linear regression. Therefore, rather than predicting point estimate of the event itself, it builds the model to predict the odds of its occurrence. In two class problem odds greater than 50% would mean that the case is assigned to the class designed as “1” and “0” otherwise. While logistic regression is a powerful modeling tool, it assumes that the response variable is linear in the coefficients of the predictor variable.

Measures for performance evaluation:-

In this study we used three performance measures: accuracy, sensitivity, specificity.

- i) Accuracy = It is a rate of true classified instances.
- ii) Sensitivity = It is a rate of true positive classified instances.
- iii) Specificity = It is a ratio of true negative instances and total observed negative instances.

To develop a model, we partitioned the data into two parts training dataset and testing dataset. First we build the model on training dataset and then check the performance of model on testing dataset. The methodology and performance of every model is as follows -

KNN Classifier:

We used R software to build KNN classifier with starting point is 17. The confusion matrix by using testing data is,

obs.\pred.	No.	Yes	Total
No.	18	7	25
Yes	13	23	36
Total	31	30	61

Where,

Observed response No & predicted response No (True negative) = 18

Observed response No & predicted response Yes (False negative) = 7

Observed response Yes & predicted response No (False positive) =13
 Observed response Yes & predicted response Yes (True positive) =23

Naive Bayesian Classifier:-

We used R-software to build Naive Bayesian classifier. The confusion matrix by using testing data is,

obs.\pred.	No	Yes	Total
No	17	7	24
Yes	8	29	37
Total	25	36	61

Logistic Regression model:-

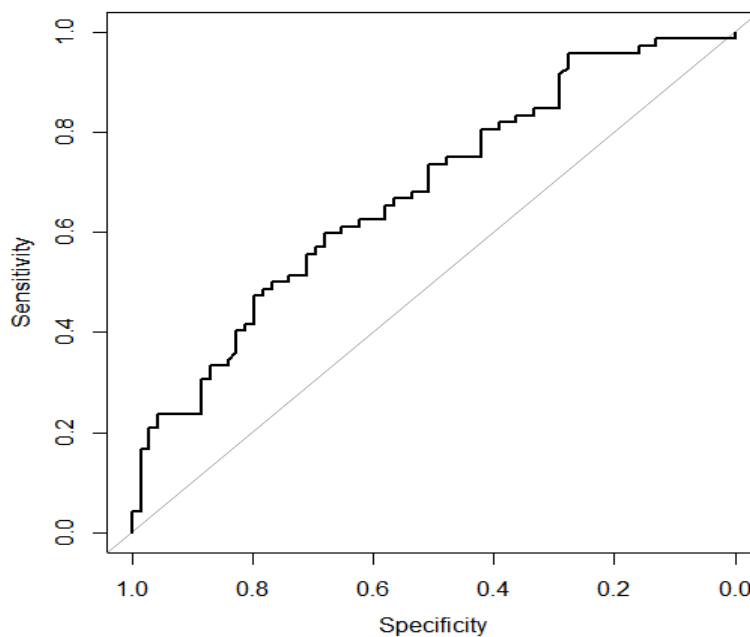
We developed a model on training dataset by using R-software. Stepwise variable selection procedure is used for variable selection and result is as follows,

Step	Effect	Estimate	Pr(> z)
1	(Intercept)	-0.787316	0.8403
2	BMI	0.012683	0.7789
3	Ht.of.S.down	-0.009190	0.7442
4	Breadth.of.S	0.115105	0.0163
5	Lgt.of.S..horz..	- 0.061566	0.1931
6	Ht.of.S.top.	0.040477	0.0720
7	S.Stering..dist..	0.040396	0.3161
8	S.Gear.dist..	-0.020479	0.4822
9	S.Clutch.dist.	0.004658	0.9204
10	S.Brake.dist.	-0.004550	0.9223

Utilization and number of credit lines variables removed in stepwise procedure.

coefficients	Breadth.of.S	Ht.of.S.top.
estimate	0.115105	0.040477

Roc :



Since area under the curve is 0.6787, hence our model is appropriate model for given data. From above ROC curve cut off point is 0.5. The confusion matrix of testing data is,

Observed\predicted	No	Yes
No	21	4
Yes	24	12

Oversampling Results:-

	Accuracy	Sensitivity	Specificity
KNN	0.6721	0.6389	0.72
Logistic regression	0.5409	0.3333	0.84
Naive Bayesian	0.7541	0.68	0.8056

From above table it is clear that accuracy of Naive Bayesian is larger than all other classifiers. Accuracy achieved by Naïve Bayesian is 75.41% with sensitivity 33.33% and specificity 80.56%.

From above table it is clear that all the classifiers have low area under curve (AUC) among them accuracy and sensitivity achieved by Naive Bayesian is higher than all the classifier. Hence it is good model for given data among all classifier models.

➤ Conclusion

In these project we conclude that,

- 1] Disease, pain, driving time & dinner time are positively associated.
- 2] 95% driver have no disease & 80% driver have no pain.
- 3] 75% drivers are not operated by surgical procedure& 50 % driver have bad habits.
- 4] Less than 60% driver have no comfortable seats.
- 5] Many of the drivers drive daily distance about 250km.
- 6] Accuracy of Naïve Bayesian is larger than all other classifiers. Accuracy achieved by Naïve Bayesian is 75.41% with sensitivity 33.33% and specificity 80.56%.

Bibliography:-

- 1) B. L. Agarwal Program Statistics, second edition, 2005.
- 2) <https://www.statistics.com>
- 3) www.statisticsworldwide.com/data
- 4) Data mining concept & techniques.

Appendix

R-code

```
rm=(list=ls())
```

```
data=read.csv("C:/Users/user/Desktop/Ram.csv",header=TRUE)
```

```
names(data)
```

```
x=data$X.1
```

```
y=data$Avg0
```

```
car=subset(y,x==1)
```

```
summary(car)
```

```
KMT=subset(y,x==2)
```

```
summary(KMT)
```

```
ST=subset(y,x==3)
```

```
summary(ST)
```

```
tractor=subset(y,x==4)
```

```
summary(tractor)
```

```
travels=subset(y,x==5)
```

```
summary(travels)
```



```
truck=subset(y,x==6)
```

```
summary(truck)
```

```
##proportion test##
```

```
##H0:proportion of disease
```

```
x=data$X.6
```

```
n=202
```

```
y=table(x);y
```

```
prop.test(y,n,0.95)
```

```
####H0:proportion of pain
```

```
x=data$X.8
```

```
n=202
```

```
y=table(x);y
```

```
prop.test(y,n,0.8)
```

```
##H0:proportion of opretion
```

```
x=data$X.12
```

```
n=202
```

```
y=table(x);y
```

```
prop.test(y,n,0.75)
```

```
##H0:proportion of bad habbitess
```

```
x=data$X.14
```

```
n=202
```

```
y=table(x);y
```

```
prop.test(y,n,0.5)
```

```
##H0:proportion of no conforttable seat
```

```
x=data$X.16
```

```
n=202
```

```
y=table(x);y
```

```
prop.test(y,n,0.50)
```

```
##data mining ##
```

```
data=read.csv("C:/Users/user/Desktop/Ram.csv",header=TRUE)
```

```
names(data)
```

```
w1=data[,c(6,30,32:39)]
```

```
qq=ifelse(w1[,2]==0,"N","Y")
```

```
w2=data[,c(6,32:39)]
```

```
d=cbind(qq,w2)
```

```
head(d)
```

```
##### KNN#####
```

```

#install.packages("ROCR")

library(class);library(e1071);library(lattice);library(ggplot2);library(Rcpp)

library(ROCR)

set.seed(3)

## 70% of the sample size

sep_size <- floor(0.70 * nrow(d))

train_ind <- sample(seq_len(nrow(d)), size = sep_size)

dataTr=d[train_ind, ]

dataTe=d[-train_ind, ]

#head(data)

row=nrow(dataTr);col=ncol(dataTr);

d1=dataTr[,c(2:10)]

#head(d1)

d2=dataTe[,c(2:10)]

y=dataTr[,1]

y1=dataTe[,1]

#K Choose

accuracy1=0

for(i in 1:30)

```

```

{

s1=knn(d1,d2,y,k=i)

q=table(y1,s1)

accuracy1[i]=(sum(diag(q))/length(y1))

}

plot(accuracy1,type="o")

s=knn(d1,d2,y,k=6)

q=table(y1,s)

accuracy=(sum(diag(q))/length(y1));

Sensitivity=(q[2,2])/sum(q[2,])

Specificity=(q[1,1])/sum(q[1,])

accuracy;Sensitivity;Specificity

##ROC####

s=class::knn(d1, d2, y, k=5, prob=TRUE)

prob1=attr(s, "prob")

prob=2*ifelse(s == "Y", 1-prob1, prob1) - 1

pred_knn1=prediction(prob, y1)

pred_knn <- performance(pred_knn1, "tpr", "fpr")

plot(pred_knn, avg= "threshold", colorize=T, lwd=3, main="ROC curve")

```

```
####naive Bays###
```

```
library(caret) ;library(Rcpp)
```

```
#install.packages("e1071");
```

```
N=naiveBayes(d1,y,laplace=0)
```

```
prd=predict(N,newdata=d2)
```

```
z=confusionMatrix(prd,y1) #install.packages("Rcpp");library(caret)
```

```
z
```

```
##ROC for Naive Bays
```

```
predvec <- ifelse(prd=="Y", 1, 0)
```

```
realvec <- ifelse(y1=="Y", 1, 0)
```

```
pr <- prediction(predvec, realvec)
```

```
prf <- performance(pr, "tpr", "fpr")
```

```
a=plot(prf)
```

```
auc(a)
```

```
###logistic regression classifire##
```

```
model=glm(qq ~.,family=binomial(link='logit'),data=dataTr)
```

```
summary(model)
```

```
dd=dataTr[,c(1,4,6)]
```

```
modelfit <- glm(formula=qq ~ ., family=binomial(), data=dd, na.action=na.omit)
```

```

##ROC logistic

prob=predict(modelfit,type=c("response"))

modelfit$prob=prob

library(pROC)

g <- roc(qq ~ prob, data =dd)

plot(g)

auc(g)

dataTe2=dataTe[,c(1,4,6)]

fit2=predict(modelfit,dataTe2)

predt=ifelse(fit2 > 0.7057,1,0) #AUC=0.928

q2=table(y1,predt)

accuracy2=(sum(diag(q2))/length(y1));

Sensitivity2=(q2[2,2])/sum(q2[2,])

Specificity2=(q2[1,1])/sum(q2[1,])

accuracy2;Sensitivity2;Specificity2

```

