

ML Assignemnt 3

Ananth Kumar

11/10/2021

R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
.libPaths("C:\\Users\\Ananth\\OneDrive\\Desktop\\MSBA Kent\\Fall 2021\\Fundamentals of Machine Learning\\Assignment\\Ass 2")
```

```
library(reshape) library(caret) library(e1071)
```

readin the excel data into dataframe

```
rm(list=ls())
NB3 <- read.csv("UniversalBank.csv")
head(NB3)
```

```
##   ID Age Experience  Income ZIP.Code Family CCAvg Education Mortgage
## 1  1  25           1     49   91107      4   1.6           1         0
## 2  2  45          19     34   90089      3   1.5           1         0
## 3  3  39          15     11   94720      1   1.0           1         0
## 4  4  35           9    100   94112      1   2.7           2         0
## 5  5  35           8     45   91330      4   1.0           2         0
## 6  6  37          13     29   92121      4   0.4           2        155
##   Personal.Loan Securities.Account CD.Account Online CreditCard
## 1           0                   1           0         0         0
## 2           0                   1           0         0         0
## 3           0                   0           0         0         0
## 4           0                   0           0         0         0
## 5           0                   0           0         0         1
## 6           0                   0           0         1         0
```

Converting data into factors(categorical) mainly the one which are important to this.

```
NB3$Personal.Loan = as.factor(NB3$Personal.Loan) # converting Personal Loan into categorical data
NB3$Online = as.factor(NB3$Online) # converting Online into categorical data
NB3$CreditCard = as.factor(NB3$CreditCard) # converting CreditCard into categorical data
```

#Data partition 60 % training and 40 % into validation

```
set.seed(1)
train.index <- sample(row.names(NB3), 0.6*dim(NB3)[1]) # 60 % of data into training set
valid.index <- setdiff(row.names(NB3), train.index) # 40 % into validation set
train.df <- NB3[train.index, ] # assigning the train.index into data frame
valid.df <- NB3[valid.index, ] # assigning the validation index into data frame
train <- NB3[train.index, ] # Making a copy of the data frame train.df
valid = NB3[valid.index,] # Making a copy of the data frame valid.df
```

A. Create a pivot table for the training data with Online as a column variable, CC as a row variable, and Loan as a secondary row variable. The values inside the table should convey the count. In R use functions melt() and cast(), or function table().

Pivot table For CreditCard , Personal loan as row variables and Online in column.

```
library(reshape2)
```

```
## Warning: package 'reshape2' was built under R version 4.0.5
```

```
melt = melt(train,id=c("CreditCard","Personal.Loan"),variable= "Online")
```

```
## Warning: attributes are not identical across measure variables; they will be
## dropped
```

```
cast=dcast(melt,CreditCard+Personal.Loan~Online) # dcast is to convert the data in CC , Personal loan
```

```
## Aggregation function missing: defaulting to length
```

```
cast[,c(1,2,3,14)] # casting column no 14 which credit card and 1 , 2 , 3 column is , personal loan,
```

```
##   CreditCard Personal.Loan   ID Online
## 1          0             0 1924    1924
## 2          0             1  198     198
## 3          1             0  801     801
## 4          1             1   77      77
```

B. Consider the task of classifying a customer who owns a bank credit card and is actively using online banking services. Looking at the pivot table, what is the probability that this customer will accept the loan offer? [This is the probability of loan acceptance (Loan = 1) conditional on having a bank credit card (CC = 1) and being an active user of online banking services (Online = 1)].

```
LoanCC1 <- 77/3000 # 77 is the value for Loan and CC =1 as per pivot table. and 3000 is the total count
LoanCC1 # which is 2.6 %.
```

```
## [1] 0.02566667
```

C. Create two separate pivot tables for the training data. One will have Loan (rows) as a function of Online (columns) and the other will have Loan (rows) as a function of CC.

```
melt1 = melt(train,id=c("Personal.Loan"),variable = "Online") # Melting Personal loan and Online data
```

```
## Warning: attributes are not identical across measure variables; they will be
## dropped
```

```
melt2 = melt(train,id=c("CreditCard"),variable = "Online") # Melting Credicard data with reference to o
```

```
## Warning: attributes are not identical across measure variables; they will be
## dropped
```

```
cast1 =dcast(melt1,Personal.Loan~Online) # Casting Personal loan and online values
```

```
## Aggregation function missing: defaulting to length
```

```
cast2=dcast(melt2,CreditCard~Online) # Casting Personal loan and online values
```

```
## Aggregation function missing: defaulting to length
```

```
Loanline=cast1[,c(1,13)]
LoanCC = cast2[,c(1,14)]
```

```
Loanline # indicates personal loan count in reference with online
```

```
##   Personal.Loan Online
## 1             0    2725
## 2             1     275
```

```
LoanCC # Indicates Credit Card count in reference with online.
```

```
##   CreditCard Online
## 1          0    2122
## 2          1     878
```

D. Compute the following quantities [P (A | B) means “the probability of A given B”]: P (CC = 1 | Loan = 1) (the proportion of credit card holders among the loan acceptors) P(Online=1|Loan=1) P (Loan = 1) (the proportion of loan acceptors) P(CC=1|Loan=0) P(Online=1|Loan=0) P(Loan=0)

```
table(train[,c(14,10)]) # Creating a pivot table for column 14 and 10 which is credit card and persona
```

```
##           Personal.Loan
## CreditCard    0      1
##           0 1924  198
##           1  801   77
```

```
table(train[,c(13,10)]) # Creating a pivot table for column 13 and 10 which is online and personal loan
```

```
##      Personal.Loan
## Online    0     1
##      0 1137  109
##      1 1588  166
```

```
table(train[,c(10)]) # Pivot table for Personal loan. There are 2725 and 275.
```

```
##
##      0     1
## 2725  275
```

$P(CC = 1 | Loan = 1)$

```
CCLoan1 = 77/(77+198) # by referring the above pivot table we can get the CC= 1 and Loan = 1 values, wh
CCLoan1
```

```
## [1] 0.28
```

$P(Online=1|Loan=1)$

```
ONLoan1 =166/(166+109) # by referring the above pivot table we can get the online = 1 and Loan = 1 valu
ONLoan1
```

```
## [1] 0.6036364
```

$P(Loan = 1)$

```
Loan1 =275/(275+2725) # by referring the above pivot table we can get the Loan = 1
Loan1
```

```
## [1] 0.09166667
```

$P(CC=1|Loan=0)$

```
CCLoan01= 801/(801+1924) # by referring the above pivot table we can get the CC = 1 and Loan = 0 values
CCLoan01
```

```
## [1] 0.293945
```

$P(Online=1|Loan=0)$

```
O1LO= 1588/(1588+1137) # by referring the above pivot table we can get the online = 1 and Loan = 0 val
O1LO
```

```
## [1] 0.5827523
```

$P(Loan=0)$

```
Loan0= 2725/(2725+275) # by referring the above pivot table we can get the Loan = 0 values
Loan0
```

```
## [1] 0.9083333
```

E. Use the quantities computed above to compute the naive Bayes probability $P(\text{Loan} = 1 \mid \text{CC} = 1, \text{Online} = 1)$.

```
Naivebayes = ((77/(77+198))*(166/(166+109))*(275/(275+2725)))/(((77/(77+198))*(166/(166+109))*(275/(275+2725))))
Naivebayes # 90 % is the probability
```

```
## [1] 0.09055758
```

F. Compare this value with the one obtained from the pivot table in (b). Which is a more accurate estimate? 9.05% are very similar to the 9.7% the difference between the exact method and the naive-bayes method is the exact method would need the the exact same independent variable classifications to predict, where the naive bayes method does not.

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 4.0.5
```

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

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

```
library(e1071)
naive.train = train.df[,c(10,13,14)] # training data is from Personal loan, Creditcard and online. columns
naive.test = valid.df[,c(10,13,14)] # testing set data from the same columns of data
naivebayes = naiveBayes(Personal.Loan~.,data=naive.train) # applying naive bayes to personal loan and test data
naivebayes
```

```
##
## Naive Bayes Classifier for Discrete Predictors
##
## Call:
## naiveBayes.default(x = X, y = Y, laplace = laplace)
##
## A-priori probabilities:
## Y
##      0      1
## 0.9083333 0.0916667
##
## Conditional probabilities:
##      Online
## Y      0      1
## 0 0.4172477 0.5827523
## 1 0.3963636 0.6036364
##
```

```
##      CreditCard
## Y          0          1
##  0 0.706055 0.293945
##  1 0.720000 0.280000
```

G. Which of the entries in this table are needed for computing $P(\text{Loan} = 1 \mid \text{CC} = 1, \text{Online} = 1)$? In R, run naive Bayes on the data. Examine the model output on training data, and find the entry that corresponds to $P(\text{Loan} = 1 \mid \text{CC} = 1, \text{Online} = 1)$. Compare this to the number you obtained in (E).

Answer:

the naive bayes is the exact same output we recieved in the previous methods. $(.280)(.603)(.09)/(.280.603.09+.29.58.908) = .09$ which is the same response provided as above.