

Assignment_03_ML

Yash

2022-10-19

#Installing Required packages and loading the dataset.

```
#install.packages("tidyverse")  
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.2 --  
## v ggplot2 3.3.6      v purrr  0.3.4  
## v tibble  3.1.8      v dplyr  1.0.10  
## v tidyr   1.2.1      v stringr 1.4.1  
## v readr   2.1.2      v forcats 0.5.2  
## -- Conflicts ----- tidyverse_conflicts() --  
## x dplyr::filter() masks stats::filter()  
## x dplyr::lag()    masks stats::lag()
```

```
#install.packages("reshape")  
library(reshape)
```

```
##  
## Attaching package: 'reshape'  
##  
## The following object is masked from 'package:dplyr':  
##  
##   rename  
##  
## The following objects are masked from 'package:tidyr':  
##  
##   expand, smiths
```

```
#install.packages("caret")  
library(caret)
```

```
## Loading required package: lattice  
##  
## Attaching package: 'caret'  
##  
## The following object is masked from 'package:purrr':  
##  
##   lift
```

```
#install.packages("e1071")
library(e1071)
Ub <- read_csv("C:/Users/YASH/Downloads/UniversalBank.csv")

## Rows: 5000 Columns: 14
## -- Column specification -----
## Delimiter: ","
## db1 (14): ID, Age, Experience, Income, ZIP Code, Family, CCAvg, Education, M...
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```
head(Ub)
```

```
## # A tibble: 6 x 14
##       ID   Age Experience Income 'ZIP Code' Family CCAvg Educat~1 Mortg~2 Perso~3
##   <dbl> <dbl>      <dbl>  <dbl>      <dbl>  <dbl> <dbl>   <dbl>   <dbl>   <dbl>
## 1     1    25         1     49      91107     4   1.6       1       0       0
## 2     2    45        19     34      90089     3   1.5       1       0       0
## 3     3    39        15     11      94720     1   1         1       0       0
## 4     4    35         9    100      94112     1   2.7       2       0       0
## 5     5    35         8     45      91330     4   1         2       0       0
## 6     6    37        13     29      92121     4   0.4       2     155       0
## # ... with 4 more variables: 'Securities Account' <dbl>, 'CD Account' <dbl>,
## #   Online <dbl>, CreditCard <dbl>, and abbreviated variable names
## #   1: Education, 2: Mortgage, 3: Personal.Loan
```

```
tail(Ub)
```

```
## # A tibble: 6 x 14
##       ID   Age Experience Income 'ZIP Code' Family CCAvg Educat~1 Mortg~2 Perso~3
##   <dbl> <dbl>      <dbl>  <dbl>      <dbl>  <dbl> <dbl>   <dbl>   <dbl>   <dbl>
## 1  4995    64         40     75      94588     3   2         3       0       0
## 2  4996    29         3     40      92697     1   1.9       3       0       0
## 3  4997    30         4     15      92037     4   0.4       1     85       0
## 4  4998    63        39     24      93023     2   0.3       3       0       0
## 5  4999    65        40     49      90034     3   0.5       2       0       0
## 6  5000    28         4     83      92612     3   0.8       1       0       0
## # ... with 4 more variables: 'Securities Account' <dbl>, 'CD Account' <dbl>,
## #   Online <dbl>, CreditCard <dbl>, and abbreviated variable names
## #   1: Education, 2: Mortgage, 3: Personal.Loan
```

```
colnames(Ub)
```

```
## [1] "ID"           "Age"           "Experience"
## [4] "Income"       "ZIP Code"      "Family"
## [7] "CCAvg"        "Education"     "Mortgage"
## [10] "Personal.Loan" "Securities Account" "CD Account"
## [13] "Online"       "CreditCard"
```

```
#Transforming data into factors (categorical).
```

```
Ub$Personal.Loan = as.factor(Ub$Personal.Loan)
Ub$Online = as.factor(Ub$Online)
Ub$CreditCard = as.factor(Ub$CreditCard)
```

#Partitioning the 60% of data in training set and remaining into validation set

```
set.seed(456)
ub.train.data <- sample(row.names(Ub), 0.6*dim(Ub)[1]) # 60 % training
ub.valid.data <- setdiff(row.names(Ub), ub.train.data) # 40 % validation
ub.train <- Ub[ub.train.data, ] # assigning the ub.train.data into data frame
ub.valid <- Ub[ub.valid.data, ] # assigning the validation index into data frame
train <- Ub[ub.train.data, ] # Duplicating the data frame ub.train
valid = Ub[ub.train.data,] # Duplicating the data frame ub.valid
```

#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

```
#install.packages("reshape2")
library(reshape2)
```

```
##
## Attaching package: 'reshape2'

## The following objects are masked from 'package:reshape':
##
##      colsplit, melt, recast

## The following object is masked from 'package:tidyr':
##
##      smiths
```

```
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 the process of turning online, personal l
```

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

```
cast[,c(1,2,3,14)] # Casting column number 14: Personal loan, ID, and credit card, respectively
```

```
##   CreditCard Personal.Loan   ID Online
## 1          0             0 1917   1917
## 2          0             1  200    200
## 3          1             0  794    794
## 4          1             1   89     89
```

#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)].

```
Ub.Loan.CC1 <- 77/3000 #According to the pivot table, the value for Loan is 77, and the value for CC is 3000
Ub.Loan.CC1 # which is 26 %.
```

```
## [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 into one
```

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

```
melt2 = melt(train,id=c("CreditCard"),variable = "Online") # CREDIT CARD DATA MELTING WITH REFERENCE TO PERSONAL LOAN
```

```
## 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
```

```
Ub.Loanonline=cast1[,c(1,13)]
Ub.LoanCC = cast2[,c(1,14)]
Ub.Loanonline #shows the number of personal loans in reference to online
```

```
##   Personal.Loan Online
## 1              0    2711
## 2              1     289
```

```
Ub.LoanCC # shows the number of credit cards in reference to internet.
```

```
##   CreditCard Online
## 1          0    2117
## 2          1     883
```

D. Compute the following quantities [P (A | B) means “the probability of A given B”]: 1.P (CC = 1 | Loan = 1) (the proportion of credit card holders among the loan acceptors) 2.P(Online=1|Loan=1) 3.P (Loan = 1) (the proportion of loan acceptors) 4.P(CC=1|Loan=0) 5.P(Online=1|Loan=0) 6.P(Loan=0)

```
table(train[,c(14,10)]) # creating a pivot table with the columns 14 and 10 representing personal loan
```

```
##           Personal.Loan
## CreditCard    0      1
##           0 1917   200
##           1  794    89
```

```
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 1046   112
##           1 1665   177
```

```
table(train[,c(10)]) # Personal loan pivot table There are 2725 and 275 from training, respectively
```

```
## Personal.Loan
##      0      1
## 2711  289
```

1. $P(CC = 1 \mid Loan = 1)$

```
Ub.CC Ub.Loan1 = 77/(77+198) # We can obtain the CC= 1 and Loan = 1 values by referring to the above p
Ub.CC Ub.Loan1
```

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

2. $P(Online=1 \mid Loan=1)$

```
Ub.ON Ub.Loan1 =166/(166+109) # We can get the online = 1 and loan = 1 values from the pivot table above
Ub.ON Ub.Loan1
```

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

3. $P(Loan = 1)$

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

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

4. $P(CC=1 \mid Loan=0)$

```
Ub.CC Loan.01= 801/(801+1924) #Using the pivot table above, we can obtain the CC = 1 and Loan = 0 values
Ub.CC Loan.01
```

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

5. $P(Online=1 \mid Loan=0)$

```
UB.ON1.LO= 1588/(1588+1137) # We can get the online = 1 and loan = 0 values from the pivot table above.
UB.ON1.LO
```

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

6. $P(\text{Loan}=0)$

```
Ub.Loan0= 2725/(2725+275) # We can obtain the Loan = 0 values by the pivot table above.
Ub.Loan0
```

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

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

```
ub.Naivebayes = ((77/(77+198))*(166/(166+109))*(275/(275+275)))/(((77/(77+198))*(166/(166+109))*(275/(275+275))))
ub.Naivebayes # ~91 % 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)
library(e1071)
ub.nb.train = ub.train[,c(10,13,14)] # Personal loan, credit card, and online column training data
ub.naivebayes.1 = naiveBayes(Personal.Loan~.,data=ub.nb.train) #Using the naivebayes algorithm to predict
ub.naivebayes.1
```

```
##
## Naive Bayes Classifier for Discrete Predictors
##
## Call:
## naiveBayes.default(x = X, y = Y, laplace = laplace)
##
## A-priori probabilities:
## Y
##      0      1
## 0.90366667 0.09633333
##
## Conditional probabilities:
##      Online
## Y      0      1
## 0 0.3858355 0.6141645
## 1 0.3875433 0.6124567
##
##      CreditCard
## Y      0      1
## 0 0.7071191 0.2928809
## 1 0.6920415 0.3079585
```