

# Naive Bayes Classifier

Truong Thi An Hai

12/12/2020

## \*\*\* Part I \*\*\*

In this assignment you will train a Naive Bayes classifier on categorical data and predict individuals' incomes. Import the nbtrain.csv file. Use the first 9010 records as training data and the remaining 1000 records as testing data.

```
#Import nbtrain.csv
data = read.csv("C:/Users/hp/Desktop/nbtrain.csv")
#Split training data vs testing data
library(caret)

## Loading required package: lattice

## Loading required package: ggplot2

data_train <- head(data, 9010)
data_test <- tail(data, 1000)
```

a. Construct the Naïve Bayes classifier from the training data, according to the formula "income ~ age + sex + educ". To do this, use the "naiveBayes" function from the "e1071" package. Provide the model's a priori and conditional probabilities.

```
library(e1071)
NBClassifier <- naiveBayes(as.factor(income) ~ age + sex + educ, data=data_train)
NBClassifier

##
## Naive Bayes Classifier for Discrete Predictors
##
## Call:
## naiveBayes.default(x = X, y = Y, laplace = laplace)
##
## A-priori probabilities:
## Y
##      10-50K      50-80K      GT 80K
## 0.80266371 0.12563818 0.07169611
##
## Conditional probabilities:
##      age
##      20-30      31-45      GT 45
## 10-50K 0.20796460 0.34457965 0.44745575
## 50-80K 0.08393887 0.39752650 0.51943463
## GT 80K 0.06811146 0.34055728 0.59133127
##
##      sex
##      F      M
## 10-50K 0.4798119 0.5201881
## 50-80K 0.2871025 0.7128975
## GT 80K 0.2058824 0.7941176
##
##      educ
##      College      Others      Prof/Phd
## 10-50K 0.24585177 0.73976770 0.01438053
## 50-80K 0.49558304 0.44257951 0.06183746
## GT 80K 0.53869969 0.29566563 0.16563467
```

A-priori probabilities:

- Income is in the range 10-50K: 0.803
- Income is in the range 50-80K: 0.126
- Income is in the range GT 80K: 0.072

Conditional probabilities

Age

```
NBClassifier$Tables$age

##      age
##      20-30      31-45      GT 45
## 10-50K 0.20796460 0.34457965 0.44745575
## 50-80K 0.08393887 0.39752650 0.51943463
## GT 80K 0.06811146 0.34055728 0.59133127

SEX

NBClassifier$Tables$sex

##      sex
##      F      M
## 10-50K 0.4798119 0.5201881
## 50-80K 0.2871025 0.7128975
## GT 80K 0.2058824 0.7941176

Education

NBClassifier$Tables$educ

##      educ
##      College      Others      Prof/Phd
## 10-50K 0.24585177 0.73976770 0.01438053
## 50-80K 0.49558304 0.44257951 0.06183746
## GT 80K 0.53869969 0.29566563 0.16563467
```

b. Score the model with the testing data and create the model's confusion matrix. Also, calculate the overall, 10-50K, 50-80K, and GT 80K misclassification rates. Explain the variation in the model's predictive power across income classes.

```
testPred <- predict(NBClassifier, data_test, type="class")
message("Confusion Matrix for Test Data")

## Confusion Matrix for Test Data

Matrix <- confusionMatrix(testPred, as.factor(data_test$income))
Matrix

## Confusion Matrix and Statistics
##
##      Reference
## Prediction 10-50K 50-80K GT 80K
## 10-50K      787      127      67
## 50-80K       0       0       0
## GT 80K       6       5       8
##
## Overall Statistics
##
##      Accuracy : 0.795
##      95% CI : 0.7606, 0.8196
##      No Information Rate : 0.793
##      P-Value [Acc > NIR] : 0.4564
##
##      Kappa : 0.0709
##
##      Mcnemar's Test P-Value : <2e-16
##
## Statistics by Class:
##
##      Class: 10-50K Class: 50-80K Class: GT 80K
## Sensitivity      0.9924      0.000      0.1067
## Specificity      0.0628      1.000      0.9881
## Pos Pred Value   0.8022      NaN      0.4211
## Neg Pred Value   0.6842      0.868      0.9317
## Prevalence       0.7930      0.132      0.0750
## Detection Rate   0.7870      0.000      0.0690
## Detection Prevalence 0.9810      0.000      0.0100
## Balanced Accuracy 0.5276      0.500      0.5474

The overall misclassification rate: 1 - Accuracy = 0.205

library(shipunov)

## package 'shipunov', version 1.12

Misclass(testPred, as.factor(data_test$income))

## Classification table:
##      obs
## pred  10-50K 50-80K GT 80K
## 10-50K      787      127      67
## 50-80K       0       0       0
## GT 80K       6       5       8
## Misclassification errors (%):
## 10-50K 50-80K GT 80K
## 0.8 100.0 69.3 0.000
## Mean misclassification error: 63.4%

• The 10-50K misclassification rate: 0.8%
• The 50-80K misclassification rate: 100%
• The GT 80K misclassification rate: 89.3%
```

In this model variation is explained mostly by confusion matrix

## \*\*\* Part II \*\*\*

As in assignment I, import the nbtrain.csv file. Use the first 9010 records as training data and the remaining 1000 records as testing data.

```
#Import nbtrain.csv
data = read.csv("C:/Users/hp/Desktop/nbtrain.csv")
#Split training data vs testing data

data_train <- head(data, 9010)
data_test <- tail(data, 1000)
```

a. Construct the classifier according to the formula "sex ~ age + educ + income", and calculate the overall, female, and male misclassification rates. Explain the misclassification rates?

```
NBClassifier <- naiveBayes(as.factor(sex) ~ age + income + educ, data=data_train)
NBClassifier

##
## Naive Bayes Classifier for Discrete Predictors
##
## Call:
## naiveBayes.default(x = X, y = Y, laplace = laplace)
##
## A-priori probabilities:
## Y
##      F      M
## 0.43596 0.56404
##
## Conditional probabilities:
##      age
##      20-30      31-45      GT 45
## F 0.1802444 0.3475951 0.4722595
## M 0.1837859 0.3536009 0.4626131
##
##      income
##      10-50K      50-80K      GT 80K
## F 0.80340122 0.00273951 0.8338947
## M 0.74025974 0.15879575 0.10094451
##
##      educ
##      College      Others      Prof/Phd
## F 0.32128310 0.65797739 0.02163951
## M 0.28040142 0.68103189 0.03856749

testPred <- predict(NBClassifier, data_test, type="class")
Matrix <- confusionMatrix(testPred, as.factor(data_test$sex))

The overall misclassification rate: 1 - Accuracy = 0.418

Misclass(testPred, as.factor(data_test$sex))

## Classification table:
##      obs
## pred  F      M
## F 106 97
## M 321 476
## Misclassification errors (%):
## F      M
## 75.2 16.9
## Mean misclassification error: 46.1%

• The female misclassification rate: 75.2%
• The male misclassification rate: 16.9%
```

b. Divide the training data into two partitions, according to sex, and randomly select 3500 records from each partition. Reconstruct the model from part (a) from these 7000 records. Provide the model's a priori and conditional probabilities.

```
library(dplyr)

#Divide the training data into two partitions
data_female = subset(data_train, data_train$sex == 'F')
data_male = subset(data_train, data_train$sex == 'M')
#Randomly select 3500 records from each partition
data_female = sample_n(data_female, 3500)
data_male = sample_n(data_male, 3500)
new_data = rbind(data_female, data_male)
model <- naiveBayes(as.factor(sex) ~ age + income + educ, data=new_data)
message("Model Naive Bayes Classifier")
model

##
## Naive Bayes Classifier for Discrete Predictors
##
## Call:
## naiveBayes.default(x = X, y = Y, laplace = laplace)
##
## A-priori probabilities:
## Y
##      F      M
## 0.5 0.5
##
## Conditional probabilities:
##      age
##      20-30      31-45      GT 45
## F 0.1802487 0.3477143 0.47280900
## M 0.1805714 0.3508571 0.4626714
##
##      income
##      10-50K      50-80K      GT 80K
## F 0.80114286 0.00514286 0.83371429
## M 0.744000000 0.15714286 0.09085714
##
##      educ
##      College      Others      Prof/Phd
## F 0.322000000 0.65571429 0.02228571
## M 0.28028571 0.68257143 0.03714286
```

The a priori probabilities are equal and the conditional probabilities are very similar.

c. How well does the model classify the testing data?

```
testPred <- predict(model, data_test, type="class")
Matrix <- confusionMatrix(testPred, as.factor(data_test$sex))
Matrix$table

##      Reference
## Prediction F      M
##      F 106 97
##      M 321 476
## Misclassification errors (%):
## F      M
## 75.2 16.9
## Mean misclassification error: 46.1%

message("Accuracy")

## Accuracy

Matrix$overall[1]

## Accuracy
##      0.53
```

d. Repeat step (b) 4 several times. What effect does the random selection of records have on the model's performance?

```
1.

#Divide the training data into two partitions
data_female = subset(data_train, data_train$sex == 'F')
data_male = subset(data_train, data_train$sex == 'M')
#Randomly select 3500 records from each partition
data_female = sample_n(data_female, 3500)
data_male = sample_n(data_male, 3500)
new_data = rbind(data_female, data_male)
model <- naiveBayes(as.factor(sex) ~ age + income + educ, data=new_data)
model

##
## Naive Bayes Classifier for Discrete Predictors
##
## Call:
## naiveBayes.default(x = X, y = Y, laplace = laplace)
##
## A-priori probabilities:
## Y
##      F      M
## 0.5 0.5
##
## Conditional probabilities:
##      age
##      20-30      31-45      GT 45
## F 0.1774286 0.3408000 0.4745714
## M 0.1831429 0.3480857 0.4680000
##
##      income
##      10-50K      50-80K      GT 80K
## F 0.80200000 0.00485714 0.83314286
## M 0.736000000 0.16000000 0.10400000
##
##      educ
##      College      Others      Prof/Phd
## F 0.31800000 0.66057143 0.02142857
## M 0.28005714 0.68028571 0.03805714

testPred <- predict(model, data_test, type="class")
Matrix <- confusionMatrix(testPred, as.factor(data_test$sex))
message("Accuracy")

## Accuracy

Matrix$overall[1]

## Accuracy
##      0.53

2.

#Divide the training data into two partitions
data_female = subset(data_train, data_train$sex == 'F')
data_male = subset(data_train, data_train$sex == 'M')
#Randomly select 3500 records from each partition
data_female = sample_n(data_female, 3500)
data_male = sample_n(data_male, 3500)
new_data = rbind(data_female, data_male)
model <- naiveBayes(as.factor(sex) ~ age + income + educ, data=new_data)
model

##
## Naive Bayes Classifier for Discrete Predictors
##
## Call:
## naiveBayes.default(x = X, y = Y, laplace = laplace)
##
## A-priori probabilities:
## Y
##      F      M
## 0.5 0.5
##
## Conditional probabilities:
##      age
##      20-30      31-45      GT 45
## F 0.1817143 0.3442857 0.4740000
## M 0.1837143 0.3517143 0.4645714
##
##      income
##      10-50K      50-80K      GT 80K
## F 0.80428571 0.00257143 0.83314286
## M 0.74028571 0.16028571 0.09942857
##
##      educ
##      College      Others      Prof/Phd
## F 0.32028571 0.65057143 0.02114286
## M 0.28371429 0.6785714 0.03742857

testPred <- predict(model, data_test, type="class")
Matrix <- confusionMatrix(testPred, as.factor(data_test$sex))
message("Accuracy")

## Accuracy

Matrix$overall[1]

## Accuracy
##      0.53

3.

#Divide the training data into two partitions
data_female = subset(data_train, data_train$sex == 'F')
data_male = subset(data_train, data_train$sex == 'M')
#Randomly select 3500 records from each partition
data_female = sample_n(data_female, 3500)
data_male = sample_n(data_male, 3500)
new_data = rbind(data_female, data_male)
model <- naiveBayes(as.factor(sex) ~ age + income + educ, data=new_data)
model

##
## Naive Bayes Classifier for Discrete Predictors
##
## Call:
## naiveBayes.default(x = X, y = Y, laplace = laplace)
##
## A-priori probabilities:
## Y
##      F      M
## 0.5 0.5
##
## Conditional probabilities:
##      age
##      20-30      31-45      GT 45
## F 0.1817143 0.3442857 0.4740000
## M 0.1848571 0.3551429 0.4600000
##
##      income
##      10-50K      50-80K      GT 80K
## F 0.80285714 0.00257143 0.83457143
## M 0.740000000 0.15742857 0.10257143
##
##      educ
##      College      Others      Prof/Phd
## F 0.32114286 0.65771429 0.02114286
## M 0.27542857 0.68371429 0.04085714

testPred <- predict(model, data_test, type="class")
Matrix <- confusionMatrix(testPred, as.factor(data_test$sex))
message("Accuracy")

## Accuracy

Matrix$overall[1]

## Accuracy
##      0.53

Conditional probabilities are very close over the entire sample.
```