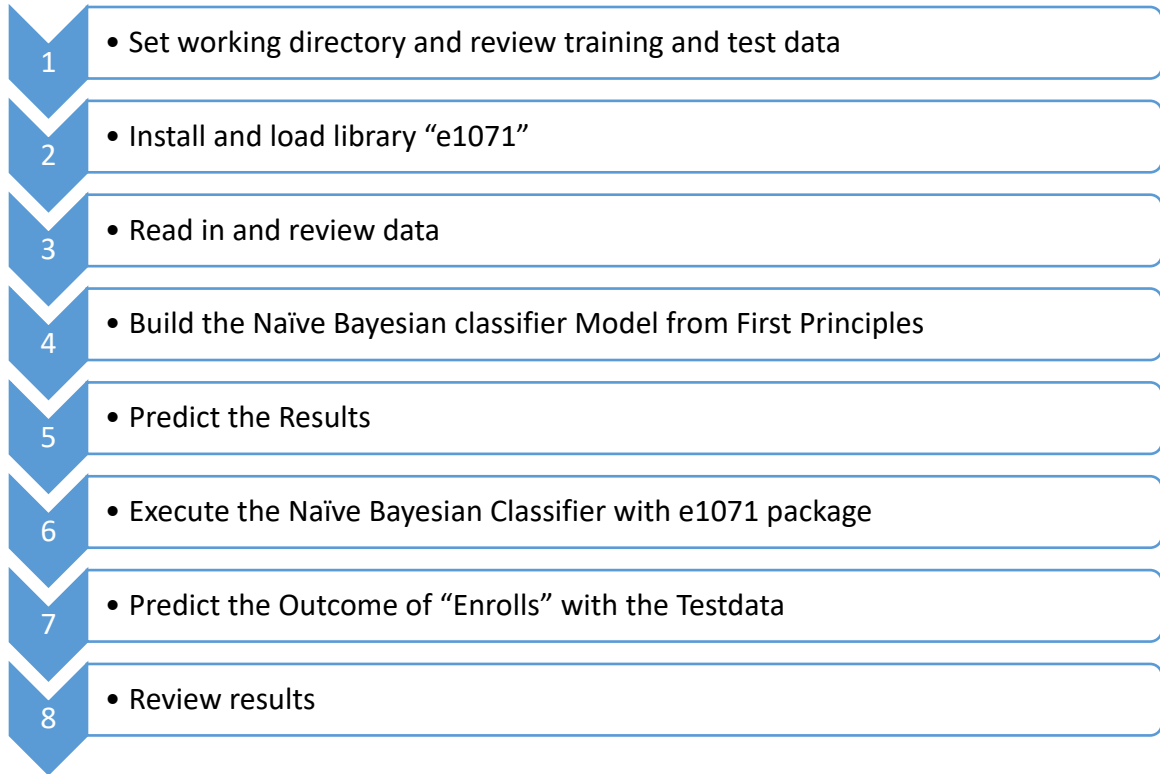**Hannah Roach**
**5-12-2019**

# Lab Exercise 8: Naïve Bayesian Classifier

| Purpose: | This lab is designed to investigate and practice the Naïve Bayesian Classifier analytic technique. After completing the tasks in this lab you should be able to: |
|---|---|
| | • Use R functions for Naïve Bayesian Classification<br>• Apply the requirements for generating appropriate training data<br>• Validate the effectiveness of the Naïve Bayesian Classifier with the big data |

| Tasks: | Tasks you will complete in this lab include: |
|---|---|
| | • Use R –Studio environment to code the Naïve Bayesian Classifier<br>• Use the ODBC connection to the "census" database to create a training data set for Naïve Bayesian Classifier from the big data<br>• Use the Naïve Bayesian Classifier program and evaluate how well it predicts the results using the training data and then compare the results with original data |

# Part 1 – Building Naïve Bayesian Classifier

# Workflow Overview

| | |
|---|---|
| **1** | • Set working directory and review training and test data |
| **2** | • Install and load library "e1071" |
| **3** | • Read in and review data |
| **4** | • Build the Naïve Bayesian classifier Model from First Principles |
| **5** | • Predict the Results |
| **6** | • Execute the Naïve Bayesian Classifier with e1071 package |
| **7** | • Predict the Outcome of "Enrolls" with the Testdata |
| **8** | • Review results |

# LAB Instructions

| Step | Action |
|------|--------|
| 1 | Log in with GPADMIN credentials on to R-Studio. |
| 2 | **Set working directory and review training and test data**<br>1. Set the working directory using the following command:<br>`> setwd("~/LAB08")`<br><br>• The **"sample1.csv"** file in this directory represents the data worked with in the instructor led training session. The file has a header row, followed by 14 rows of training data.<br>• The **testing data** on which you will predict the results should be appended after the **training data**. The data set should read:<br><br>`Age,Income,Jobstaisfaction,Desire,Enrolls` ←------- Header<br>`<=30,High,No,Fair,No`<br>`<=30,High,No,Excellent,No`<br>`31 to 40,High,No,Fair,Yes`<br>`>40,Medium,No,Fair,Yes`<br>`>40,Low,Yes,Fair,Yes`<br>`>40,Low,Yes,Excellent,No`<br>`31 to 40,Low,Yes,Excellent,Yes`<br>`<=30,Medium,No,Fair,No`<br>`<=30,Low,Yes,Fair,Yes`<br>`>40,Medium,Yes,Fair,Yes`<br>`<=30,Medium,Yes,Excellent,Yes`<br>`31 to 40,Medium,No,Excellent,Yes`<br>`31 to 40,High,Yes,Fair,Yes`<br>`>40,Medium,No,Excellent,No`<br>`<=30,Medium,Yes,Fair,` ←--------- testing data |
| 3 | **Install and load library "e1071"**<br>Execute the following command to install the required packages and load the libraries:<br>`> install.packages("e1071")`<br>`> library("e1071")` |

| Step | Action |
|------|--------|
| 4 | **Read in and review data**<br>1. Execute the following to read in the data.<br><br>```<br>> # read the data into a table from the file<br>> sample <- read.table("sample1.csv",header=TRUE,sep=",")<br>> # we will now define the data frames to use the NB<br>classifier<br>> # we will now define the data frames to use the NB<br>classifier<br>> traindata <- as.data.frame(sample[1:14,])<br>> testdata <- as.data.frame(sample[15,])<br>```<br><br>You now have two data frame objects "**traindata**" and "**testdata**" for running the NB Classifier.<br><br>2. Execute the following command to display the data frames, to ensure they are loaded properly.<br><br>```<br>> #Display data frames<br>> traindata<br>> testdata<br>```<br><br>3. Screenshot the output for traindata and test data<br><br> |

| 5 | **Build the Naïve Bayesian classifier Model from First Principles:** |
|---|---|

1. The first step in building the model is the computation of prior probabilities. The independent variables here are the "Age", "Income", "Jobsatisfaction" and "Desire". The dependent variable is "Enrolls"
Compute the prior probabilities of enrollment, P(no), P(yes) first, the counts :

```
> tprior <- table(traindata$Enrolls)
```
then, normalize over the total number of instances to get the probabilities
```
> tprior <- tprior/sum(tprior)
> tprior
```

```
> tprior

           No   Yes
0.000 0.357 0.643
```

2. Compute the summaries that you need to create a Bayes model: P(A|b), b={no, yes}
First, count up "no" and "yes" by Age:
```
> ageCounts <-table(traindata[,c("Enrolls", "Age")])
```

3. Then, normalize by the total number of "no" and "yes" each to get the conditional probabilities
```
> ageCounts <- ageCounts/rowSums(ageCounts)
```

Display the results on the console and review the conditional probabilities
```
> ageCounts
```
```
> ageCounts
          Age
Enrolls   <=30 31 to 40    >40

    No  0.600     0.000 0.400
    Yes 0.222     0.444 0.333
```

4. Do the same for the other variables.

```
> incomeCounts <- table(traindata[,c("Enrolls",
"Income")])
> incomeCounts <- incomeCounts/rowSums(incomeCounts)
>incomeCounts
```
**Screenshot**

| Step | Action |
|------|--------|
| | ``` > incomeCounts           Income Enrolls  High    Low Medium      No  0.400 0.200  0.400     Yes 0.222 0.333  0.444 ```<br><br>```> jsCounts <- table(traindata[,c("Enrolls", "Jobsatisfaction")])```<br>```> jsCounts<-jsCounts/rowSums(jsCounts)```<br>```>jsCounts```<br>==Screenshot==<br>``` > jsCounts         Jobsatisfaction Enrolls     No    Yes      No  0.800 0.200     Yes 0.333 0.667 ```<br><br>```> desireCounts <- table(traindata[,c("Enrolls", "Desire")])```<br>```> desireCounts <- desireCounts/rowSums(desireCounts)```<br>```>desireCounts```<br>==Screenshot==<br>``` > desireCounts         Desire Enrolls Excellent  Fair      No        0.600 0.400     Yes        0.333 0.667 ``` |

| Step | Action |
|---|---|
| 6 | **Predict the Results:**<br><br>1. Use the Naïve Bayesian Classifier formula to compute product of P(A\|b), for b={no, yes}. The maximum of the two is the "predicted" result of the dependent variable. In the test data we need to predict the "Enrolls" given the for Age<=30, Income = Medium, Jobsatisfaction = yes and Desire = Fair<br><br>```<br>> pyes <-<br>        ageCounts["Yes","<=30"]*<br>        incomeCounts["Yes","Medium"]*<br>        jsCounts["Yes","Yes"]*<br>        desireCounts["Yes","Fair"]*<br>        tprior["Yes"]<br>```<br>followed by<br>```<br>> pno <-<br>        ageCounts["No","<=30"]*<br>        incomeCounts["No","Medium"]*<br>        jsCounts["No","Yes"]*<br>        desireCounts["No","Fair"]*<br>        tprior["No"]<br>```<br>2. The prediction will be max(pyes,pno).<br>```<br>> pyes<br>> pno<br>> max(pyes,pno)<br>```<br><br>3. What is the predicted result for "Enrolls" for someone's age less than 30, income medium ,JobSatisfaction yes, and desidre Fiar?<br><br>```<br>> print (pyes)<br>    Yes<br>0.0282<br>> print (pno)<br>     No<br>0.00686<br>> print(max(pyes,pno))<br>[1] 0.0282<br>``` |

| 7 | **Execute the Naïve Bayesian Classifier with e1071 package:** |
|---|---|
| | The Naïve Bayes function computes the conditional a-posterior probabilities of a categorical class variable given independent categorical predictor variables using the Bayes rule. The usage takes the form of naiveBayes(formula, data,…) where the arguments are defined as follows:<br><br>    &bull;  **formula** A formula of the form class ~ x1 + x2 + .... Interactions are not allowed.<br>    &bull;  **data**    Either a data frame of factors or a contingency table.<br><br>&bull;  You are modeling for attribute " Enrolls".<br><br>1.  Use the following commands to execute the model and display the results.<br>```> # use the NB classifier``` <br>```> model <- naiveBayes(Enrolls ~.,traindata)``` <br>```> # display model``` <br>```> model``` <br><br>2.  <mark>SCREENSHOT THE RESULTS</mark> and compare these results to the **apriori probabilities** you manually computed earlier in step 5. <mark>Are they the same or different ?</mark><br><mark>They are different</mark> |

```
> model

Naive Bayes Classifier for Discrete Predictors

Call:
naiveBayes.default(x = X, y = Y, laplace = laplace)

A-priori probabilities:
Y
          No    Yes
0.000 0.357 0.643

Conditional probabilities:
     Age
Y       <=30 31 to 40    >40

  No  0.600     0.000 0.400
  Yes 0.222     0.444 0.333

     Income
Y       High    Low Medium

  No  0.400 0.200  0.400
  Yes 0.222 0.333  0.444

     Jobsatisfaction
Y          No    Yes

  No  0.800 0.200
  Yes 0.333 0.667

     Desire
Y     Excellent  Fair

  No       0.600 0.400
  Yes      0.333 0.667
```

| 8 | **Predict the Outcome of "Enrolls" with the Testdata:** |
|---|---|

1.  To use the predict function, type in the following:

```
> # predict with testdata
> results <- predict (model,testdata)
> # display results
> results
```

2.  Review the results (Prediction for "Enrolls") on the console. What is the prediction Yes or No?

```
> results
[1] Yes
Levels:  No Yes
```

| 9 | **Review results** |
|---|---|
| | 1. |
| | |
| | 2.  Build another NB model, with Laplace smoothing  model2 = naiveBayes(Enrolls ~.,traindata, laplace=0.01) |

```
> model1

Naive Bayes Classifier for Discrete Predictors

Call:
naiveBayes.default(x = X, y = Y, laplace = laplace)

A-priori probabilities:
Y
          No   Yes
0.000 0.357 0.643

Conditional probabilities:
     Age
Y         <=30 31 to 40      >40
       0.33333  0.33333 0.33333
  No   0.59841  0.00199 0.39960
  Yes 0.22259  0.44408 0.33333

     Income
Y      High   Low Medium
       0.333 0.333  0.333
  No   0.400 0.201  0.400
  Yes 0.223 0.333  0.444

     Jobsatisfaction
Y        No   Yes
       0.500 0.500
  No   0.799 0.201
  Yes 0.334 0.666

     Desire
Y      Excellent  Fair
           0.500 0.500
  No       0.600 0.400
  Yes      0.334 0.666

>
```

3. Compare the probabilities here with those of the first model

Note down your observations in the space provided below:

The probabilities are very similar. The probability between 31 and 40 for the second model is slightly higher than the first model. In the second model there probability is measured at higher levels of precision with five numbers trailing the decimal point for age.

*End of Lab Exercise*