

Naive Bayes - An example using R language

The purpose of this code snippet is to demonstrate to the reader how a Naive Bayes (NB) classifier can be trained on given dataset using the R language. Note that building and deploying NB-enabled systems in industry settings would be more complicated than this simplified example (toy example) and such a discussion is out of the scope of this work.

The dataset: We use the well-known Iris dataset (included with R) for this purpose. It contains four features (measurements) for 150 flowers representing three types of iris plants: Setosa, Versicolor and Virginica.

Machine learning (ML) task: We will train a NB model to identify the plant type based on the four measurements of a given flower. So, our ML task in this problem would be a classification.

```
data(iris) # Attach the Iris dataset to the R environment
mydata <- iris # Let's rename the dataset as mydata
dim(mydata) # check dimensions of mydata

## [1] 150    5

sapply(mydata, class) # check the data types of each feature

## Sepal.Length Sepal.Width Petal.Length  Petal.Width   Species
##      "numeric"      "numeric"      "numeric"      "numeric"    "factor"

levels(mydata$Species) # check different levels (values) for each class

## [1] "setosa"      "versicolor" "virginica"

head(mydata) # have a look at top data points in mydata

##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1         5.1         3.5         1.4         0.2  setosa
## 2         4.9         3.0         1.4         0.2  setosa
## 3         4.7         3.2         1.3         0.2  setosa
## 4         4.6         3.1         1.5         0.2  setosa
## 5         5.0         3.6         1.4         0.2  setosa
## 6         5.4         3.9         1.7         0.4  setosa

summary(mydata) # a summary of class distributions

##   Sepal.Length   Sepal.Width   Petal.Length   Petal.Width
## Min.    :4.300   Min.    :2.000   Min.    :1.000   Min.    :0.100
## 1st Qu.:5.100   1st Qu.:2.800   1st Qu.:1.600   1st Qu.:0.300
## Median :5.800   Median :3.000   Median :4.350   Median :1.300
## Mean    :5.843   Mean    :3.057   Mean    :3.758   Mean    :1.199
## 3rd Qu.:6.400   3rd Qu.:3.300   3rd Qu.:5.100   3rd Qu.:1.800
## Max.    :7.900   Max.    :4.400   Max.    :6.900   Max.    :2.500
##      Species
## setosa    :50
## versicolor:50
## virginica :50
##
##
##
```

Creating training and validation datasets: We're going to do an 80/20 partitioning for the training and validation sets. We use the createDataPartition function from the caret package for this purpose.

```
#install.packages("caret") #If the caret package is not installed on your system, uncomment this line t
library(caret) #Loading the library
```

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

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

```
tr_index <- createDataPartition(mydata$Species, p=0.80, list=FALSE) # List of 80% of the rows
trainSet <- mydata[tr_index,] # select 80% of the data for the trainSet
testSet <- mydata[-tr_index,] # Select the remaining 20% of data for testSet
```

Building a NB classifier: Now we will train our NB classifier using the above trainSet. We will utilize the e1071 package in R for this purpose. Note the prior probabilities when outputting the following lines of code. We get the same prior probabilities for all classes.

```
#install.packages("e1071") #If the e1071 package is not installed on your system, uncomment this line t
library(e1071)
NBclassifier=naiveBayes(Species~., data=trainSet) # Once you call this line, R fits the NB model using t
print(NBclassifier) # Check the newly fitted model to see if everything is OK.
```

```
##
## Naive Bayes Classifier for Discrete Predictors
##
## Call:
## naiveBayes.default(x = X, y = Y, laplace = laplace)
##
## A-priori probabilities:
## Y
##      setosa versicolor  virginica
## 0.3333333 0.3333333 0.3333333
##
## Conditional probabilities:
##      Sepal.Length
## Y      [,1]      [,2]
## setosa  5.0325 0.2938581
## versicolor 5.9225 0.5005062
## virginica 6.5650 0.6447222
##
##      Sepal.Width
## Y      [,1]      [,2]
## setosa  3.4575 0.3280771
## versicolor 2.7900 0.2959556
## virginica 2.9650 0.3324732
##
##      Petal.Length
## Y      [,1]      [,2]
## setosa  1.4975 0.1656185
## versicolor 4.2725 0.4701268
## virginica 5.5300 0.5459407
##
##      Petal.Width
## Y      [,1]      [,2]
## setosa  0.2525 0.1131994
## versicolor 1.3275 0.1974030
## virginica 2.0325 0.2644661
```

Make predictions: Now let's apply the above model to assign labels for test cases in testSet. Then we create the confusion matrix, a table that is often used to describe the performance of a classifier.

```
testPrediction=predict(NBclassifier, newdata=testSet, type="class") # Assign labels for each test case
confusionMatrix(testPrediction, testSet$Species) # Print confusion matrix
```

```
## Confusion Matrix and Statistics
##
##               Reference
## Prediction  setosa versicolor virginica
##   setosa      10           0           0
##   versicolor   0           9           1
##   virginica    0           1           9
##
## Overall Statistics
##
##               Accuracy : 0.9333
##               95% CI : (0.7793, 0.9918)
##   No Information Rate : 0.3333
##   P-Value [Acc > NIR] : 8.747e-12
##
##               Kappa : 0.9
##
##   McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##               Class: setosa Class: versicolor Class: virginica
## Sensitivity           1.0000           0.9000           0.9000
## Specificity           1.0000           0.9500           0.9500
## Pos Pred Value        1.0000           0.9000           0.9000
## Neg Pred Value        1.0000           0.9500           0.9500
## Prevalence            0.3333           0.3333           0.3333
## Detection Rate        0.3333           0.3000           0.3000
## Detection Prevalence  0.3333           0.3333           0.3333
## Balanced Accuracy      1.0000           0.9250           0.9250
```

We can see that the accuracy is 91.3%. It was a small test set, but whether or not that accuracy is sufficient depends on the problem context and is based on many other factors such as the cost of misclassification to the business.