

# Naive Bayes - An example using R language

The main purpose of this document is to educate the reader with the practical aspect of the application of the Naive Bayes (NB) algorithm. In this direction, here we provide the necessary code snippet to show the reader how to apply the NB algorithm in a data classification problem. We use a publicly available dataset and R language for this demonstration. All the important steps of the implementation of the algorithm are explained below.

**The dataset:** In order to implement the above code, we use the well-known Iris dataset (included with R). It consists of four features (measurements), namely sepal length, sepal width, petal length, and petal width for 150 flowers. The dataset contains information about three types of iris plants: Setosa, Versicolor and Virginica. In the R package, likelihoods in the Naive Bayes formula for numerical features (e.g. measurements) are calculated using probability distributions.

**Machine learning (ML) task:** In this application, we aim to construct an ML model using NB algorithm to identify (classify) the specie of an Irish flower, based on the values of the four features ( sepal length, sepal width, petal length, and petal width ) of a given Irish flower. So, we're tackling a classification task here.

```
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) # some descriptive statistics of the data and 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 follow the convention of 80/20 samples ratio to partition the dataset to 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
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. For this purpose, we will utilize e1071 package in R. Note that following lines of code will fit the NB model to the dataset and will also display the model details. In this case, the priori probabilities for all classes are the same.

```
#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   4.9800 0.3631769
## versicolor 5.8900 0.5042384
## virginica 6.6375 0.6019999
##
##      Sepal.Width
## Y      [,1]      [,2]
## setosa   3.4200 0.3943186
## versicolor 2.7800 0.3022989
## virginica 3.0025 0.3158119
##
##      Petal.Length
## Y      [,1]      [,2]
## setosa   1.4475 0.1839558
## versicolor 4.2650 0.4907085
## virginica 5.5825 0.5439115
##
##      Petal.Width
## Y      [,1]      [,2]
## setosa   0.2350 0.09753369
## versicolor 1.3250 0.20223114
```

```
## virginica 2.0825 0.24166888
```

**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          10         4
```

```
## virginica   0          0          6
```

```
##
```

```
## Overall Statistics
```

```
##
```

```
##           Accuracy : 0.8667
```

```
##           95% CI : (0.6928, 0.9624)
```

```
## No Information Rate : 0.3333
```

```
## P-Value [Acc > NIR] : 2.296e-09
```

```
##
```

```
##           Kappa : 0.8
```

```
##
```

```
## McNemar's Test P-Value : NA
```

```
##
```

```
## Statistics by Class:
```

```
##
```

```
##           Class: setosa Class: versicolor Class: virginica
```

```
## Sensitivity          1.0000          1.0000          0.6000
```

```
## Specificity          1.0000          0.8000          1.0000
```

```
## Pos Pred Value       1.0000          0.7143          1.0000
```

```
## Neg Pred Value       1.0000          1.0000          0.8333
```

```
## Prevalence           0.3333          0.3333          0.3333
```

```
## Detection Rate       0.3333          0.3333          0.2000
```

```
## Detection Prevalence 0.3333          0.4667          0.2000
```

```
## Balanced Accuracy     1.0000          0.9000          0.8000
```