# Tutorial 7

## DSA1101
## Introduction to Data Science

### October 19, 2018

**Exercise 1.** The Naïve Bayes Classifier

This week, we will look at the CSV dataset "Titanic.csv" which provides information on the fate of passengers on the fatal maiden voyage of the ocean liner *Titanic*, and includes the variables economic status (class), sex, age and survival. We will train a naïve Bayes classifier using this dataset, and predict survival.

(a) Load the dataset "Titanic.csv" which has been posted under the folder for Tutorial 7.

```
1 Titanic_dataset= read.csv("Titanic.csv")
2 dim(Titanic_dataset)
3 head(Titanic_dataset)
```

(b) Compute the probabilities $P(Y = 1)$ (survived) and $P(Y = 0)$ (did not survive).

```
1 tprior <- table(Titanic_dataset$Survived)
2 tprior
3 tprior <- tprior/sum(tprior)
4 tprior
```

(c) Compute the conditional probabilities $P(X_i = x_i | Y = 1)$ and $P(X_i = x_i | Y = 0)$ , where $i = 1, 2, 3, 4$ for the feature variables $X = \{class, sex, age\}$.

```
1 classCounts <- table(Titanic_dataset[,c("Survived", "Class")])
2 classCounts <- classCounts/rowSums(classCounts)
3 classCounts
4
5 genderCounts <- table(Titanic_dataset[,c("Survived", "Sex")])
6 genderCounts <- genderCounts/rowSums(genderCounts)
7 genderCounts
8
9 ageCounts <- table(Titanic_dataset[,c("Survived", "Age")])
10 ageCounts <- ageCounts/rowSums(ageCounts)
11 ageCounts
```

(d) Predict survival for an adult female passenger in $2^{nd}$ class cabin.

```
1 prob_survived <-
2 classCounts["Yes","2nd"]*
3 genderCounts["Yes","Female"]*
4 ageCounts["Yes","Adult"]*
5 tprior["Yes"]
6
7 prob_not_survived <-
8 classCounts["No","2nd"]*
9 genderCounts["No","Female"]*
10 ageCounts["No","Adult"]*
11 tprior["No"]
12
13 prob_survived
14 prob_not_survived
```

(e) Compare your prediction in (d) with the one performed by the `naiveBayes` function in package 'e1071'

```r
library(e1071)

model <- naiveBayes(Survived ~.,
Titanic_dataset)

test <- data.frame(Class="2nd", Sex="Female",
Age="Adult")

results <- predict(model,test)
results
results <- predict(model,test, "raw")
results

#ratio of probability scores
prob_survived/prob_not_survived
#ratio of actual probabilities
results[1,"Yes"]/results[1,"No"]
```