

Tutorial 10

DSA1101

Introduction to Data Science

November 9, 2018

Exercise 1. Logistic regression in R

In tutorial 7, we looked 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 trained a naïve Bayes classifier using this dataset, and predict survival. This week, we will use logistic regression to predict survival, and compare the performances of the two classifiers visually using an ROC curve.

- (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) Perform logistic regression of ‘Survived’ on all the feature variables.

```
1 Survival_logistic <- glm (Survived~.,
2 data=Titanic_dataset,
3 family=binomial(link="logit"))
```

- (c) Perform naïve Bayes classification of ‘Survived’ based on all the feature variables.

```
1 library(e1071)
2 Survival_Nbayes <- naiveBayes(Survived~.,
3 data=Titanic_dataset)
```

(d) Observe and compare the ROC curves for the two classifiers.

```
1 library(ROCR)
2 pred = predict(Survival_logistic, type="response")
3 predObj = prediction(pred, Titanic_dataset$Survived)
4 rocObj = performance(predObj, measure="tpr", x.measure="fpr")
5 plot(rocObj)
6
7
8 nb_prediction <- predict(Survival_Nbayes, Titanic_dataset, type='raw')
9 score <- nb_prediction[, 2]
10 pred_nb <- prediction(score, Titanic_dataset$Survived)
11 roc_nb = performance(pred_nb, measure="tpr", x.measure="fpr")
12 plot(roc_nb, add = TRUE, col = 2)
13
14
15 legend("bottomright", c("logistic regression", "naive Bayes"), col=c("
    black", "red"), lty=1)
```