# PREDICTION COMPETITION 5

## Q1

Training data : Error estimation using 10 fold cross validation : 0.1779548

Test data : Error estimation of 'svm' using 10-fold cross validation: 0.2038239

```r
1   library(readr)
2   dat <- read.csv("Desktop/aug_train (1).csv", header = TRUE)
3   dat2<- dat
4   library(e1071)
5   library(dplyr)
6
7   ### data cleaning :
8   dat$city = as.factor(dat$city)
9   dat$city_development_index = as.factor(dat$city_development_index)
10  dat$gender = as.factor(dat$gender)
11  dat$relevent_experience = as.factor(dat$relevent_experience)
12  dat$enrolled_university = as.factor(dat$enrolled_university)
13  dat$education_level = as.factor(dat$education_level)
14  dat$major_discipline = as.factor(dat$major_discipline)
15  dat$experience= as.factor(dat$experience)
16  dat$company_size = as.factor(dat$company_size)
17  dat$company_type = as.factor(dat$company_type)
18  dat$last_new_job = as.factor(dat$last_new_job)
19  dat$city<- factor(dat$city,levels = c("city_103","city_21","city_16","city_114","city_160"))
20  dat$enrollee_id = as.factor(dat$enrollee_id)
21  ### creating a data frame for easier use :::
22  frame<- tbl_df(dat)
23  ### dealing with missing values to reduce the imbalance in the data::
24  frame$gender[frame$gender == ""] <- NA
25  frame$enrolled_university[frame$enrolled_university == ""] <- NA
26  frame$education_level[frame$education_level == ""] <- NA
27  frame$major_discipline[frame$major_discipline == ""] <- NA
28  frame$company_size[frame$company_size == ""] <- NA
29  frame$company_type[frame$company_type == ""] <- NA
30  frame$last_new_job[frame$last_new_job == ""] <- NA
31
32
```

```
33  train <- sample(1:nrow(frame),1000)
34  training <- frame[train,]
35  testing <- frame[-train,]
36  #### creating a SVM function using target to see whether one is a potential
37  #### candidate or not.
38
39  svmf = svm(target~., data = frame[train,], kernel="radial",cost=1, gamma= 5.177323e-05)
40  #### tuning out :::
41  tune.out <- tune(svm , target~., data = frame[train,] , kernel = "linear",cost = 0.001)
42
43  xtest <- frame[-train,]
44  ytest<- frame$target
45  bestmod<- tune.out$best.model
46  ypred <- predict(bestmod,xtest)
47  ytest<- ytest[c(1:length(ypred))]
48  testdat<- data.frame(x = xtest,y= as.factor(ytest))
49  pred = predict(bestmod,newdata = frame[-train,])
50  prediction_model = predict(bestmod,newdata = frame[train,])
51
52
```

```
> summary(tune.out$best.model)

Call:
best.tune(method = svm, train.x = target ~ ., data = frame[train, ], kernel = "linear",
    cost = 0.001)


Parameters:
   SVM-Type:  eps-regression
 SVM-Kernel:  linear
       cost:  0.001
      gamma:  5.177323e-05
    epsilon:  0.1


Number of Support Vectors:  147
```
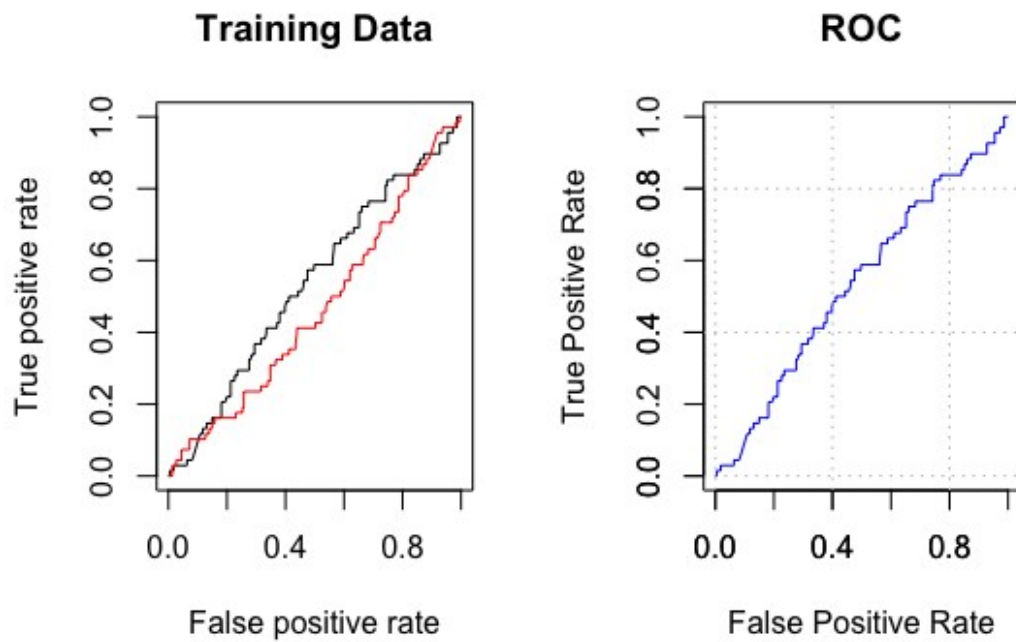
```
> summary(tune.out)

Error estimation of 'svm' using 10-fold cross validation: 0.1779548
```

**Q2**

**A)**

```
53
54   ##### QUESTION 2 :::
55 - rocplot <- function(pred , truth , ...) {
56     predob <- prediction(pred , truth)
57     perf <- performance(predob , "tpr", "fpr")
58     plot(perf , ...)
59 - }
60
61   library(ROCR)
62   svmfit.opt <- svm(target~.,data=frame[train,],kernel="radial",cost = 1, gamma=5, decision.values= T)
63   fitted<- attributes(predict(svmfit.opt,frame[train,],decision.values = TRUE))$decision.values
64   par(mfrow=c(1,2))
65   var <- frame[train,"target"]
66   len_f <- length(fitted)
67   rocplot(fitted,var[1:len_f,1], main = "Training Data")
68   rocplot(-fitted , var[1:len_f,1], add = T, col = "red")
69
70
71   ### ROC :::
72   rocplot(fitted, var[1:len_f,1], type = "l", col = "blue", xlab = "False Positive Rate"
73          , ylab = "True Positive Rate", main = "ROC")
74   axis(1, seq(0.0,1.0,0.4))
75   axis(2, seq(0.0,1.0,0.4))
76   abline(h=seq(0.0,1.0,0.4), v=seq(0.0,1.0,0.4), col="gray", lty=3)
77
```

**Training Data**

True positive rate / False positive rate

**ROC**

True Positive Rate / False Positive Rate

B)

Calculating the AUC using glm and the pROC package to get the auc.

```
78  ### Calculating the AUC :::
79  glm.fit = glm(target~., data = training, family="binomial")
80  predicted <- predict(glm.fit,data = frame[-train,],type= "response")
81  library(pROC)
82  testing <- frame[-train,]
83  AUC <- auc(testing$target[1:len_f],predicted) |
84
```

```
> AUC
Area under the curve: 0.5655
>
```

D)

Calculating the variable importance using the gbm method and boosting ::

```
86  #### part d)
87  boosting <- gbm(target~.-city_development_index-enrollee_id, data= testing,
88             distribution = "gaussian", n.trees= 1000, interaction.depth = 4, shrinkage = 0.01)
89  summary(boosting)
90
91
```

```
> summary(boosting)
                                          var    rel.inf
city                                     city 40.4247585
company_size                     company_size 23.3045664
experience                         experience 11.9304117
education_level               education_level  8.1528005
company_type                     company_type  6.4705238
last_new_job                     last_new_job  3.0812341
major_discipline             major_discipline  2.1827220
enrolled_university     enrolled_university  1.7676838
training_hours                   training_hours  1.4975057
relevent_experience relevent_experience  0.6497003
gender                                 gender  0.5380932
```
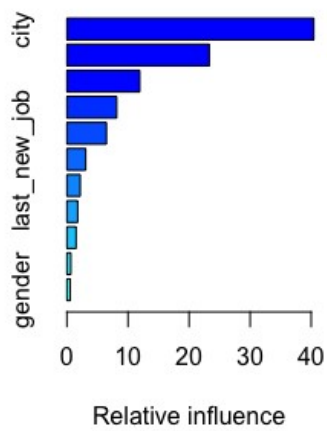
The most important variables are as follows from the training data^
However, the testing data gives a better importance :

```
> summary(boosting)
                                          var    rel.inf
experience                         experience 35.0245521
city                                     city 19.0175342
company_size                     company_size 16.6230017
last_new_job                     last_new_job  8.4786199
company_type                     company_type  5.7174332
education_level               education_level  5.5987691
training_hours                   training_hours  3.9712284
enrolled_university     enrolled_university  2.7384052
major_discipline             major_discipline  1.9324182
gender                                 gender  0.6377761
relevent_experience relevent_experience  0.2602621
```

(training data variable importance)

Confusion matrix using the best model prediction and the target which determines whether a potential candidate or not.

```
> table(predict= ypred, candidate= ytest)
                   candidate
predict              0 1
  0.0396450554191474 0 1
  0.0400642377264664 0 1
  0.0401501786651821 1 0
  0.0401992322490509 0 1
  0.0402438334891641 0 1
  0.0404269597208246 1 0
```