

Project 3: Facial Emotion Recognition

Group 9

3/30/2020

Abstract

	baseline model	improved model
test accuracy	40.4%	53.6%
training time	1816.61s	533.71s
testing time	13.199s	0.54s

The goal of this project is to predict facial emotion. After trying multiple machine learning models such as Gradient Boosting Machine (GBM), Linear discriminant analysis (LDA), Support Vector Machine (SVM), Random Forest (RF), XGBoost (XGB) and Convolutional Neural Network (CNN), we found that using XGB model has better effect. The comparison between baseline model (GBM) and improved model (XGB) is shown as below.

Step 1: set up work directory

Before reproducing the project, please set your work directory.

```
setwd("/Users/rachel/Documents/GitHub/Spring2020-Project3-ads-spring2020-project3-group9/doc")
```

```
set.seed(5)
```

```
seed <- .Random.seed
```

```
train_dir <- "../data/train_set/"
```

```
train_image_dir <- paste(train_dir, "images/", sep="")
```

```
train_pt_dir <- paste(train_dir, "points/", sep="")
```

```
train_label_path <- paste(train_dir, "label.csv", sep="")
```

```
run.feature.train=TRUE # process features for training set
```

```
run.test=TRUE # run evaluation on an independent test set
```

```
run.feature.test=TRUE # process features for test set
```

Step 2: Import data and train-test data split

We split 80% of data into training and 20% into testing. For baseline model - GBM, and other models such as SVM, XGBoost, RF, and LDA, we only used the fiducial points extracted from the images.

```
### set train-test split
```

```
info <- read.csv(train_label_path)
```

```
n <- nrow(info)
```

```
n_train <- round(n*(4/5), 0)
```

```
train_idx <- sample(info$Index, n_train, replace = F)
```

```
test_idx <- setdiff(info$Index, train_idx)
```

```

### read fiducial points
n_files <- length(list.files(train_image_dir))
readMat.matrix <- function(index){
  return(round(readMat(paste0(train_pt_dir, sprintf("%04d", index), ".mat"))[[1]],0))
}

### load fiducial points
fiducial_pt_list <- lapply(1:n_files, readMat.matrix)
save(fiducial_pt_list, file="../output/fiducial_pt_list.RData")
#load("../output/fiducial_pt_list.RData")

```

Step 3: Feature Extraction

Based on the fiducial points, we calculated the pairwise distances between fiducial points and utilized each distance as a feature. We didn't construct or remove any feature in this step as the fiducial points are extracted features from the original images. That said, the information we can obtain from the fiducial points is less than that from the images. If we continue removing some points or information, there is less data for models to train. But we conducted the principal component analysis(PCA) to reduce the dimensions in further steps.

```

### calculate the pairwise distances between fiducial points
source("../lib/feature.R")
tm_feature_train <- NA
if(run.feature.train){
  tm_feature_train <- system.time(dat_train <- feature(fiducial_pt_list, train_idx))
}

tm_feature_test <- NA
if(run.feature.test){
  tm_feature_test <- system.time(dat_test <- feature(fiducial_pt_list, test_idx))
}

save(dat_train, file="../output/feature_train.RData")
save(dat_test, file="../output/feature_test.RData")

#load(file="../output/feature_train.RData")
#load(file="../output/feature_test.RData")

```

Step 4: Baseline Model - Gradient Boosting Machine (GBM)

Step 4.1: Build model and fit on train data

To keep this main file simple, we removed the cross validation and parameters tuning step to reduce the reproduction time. After conduction cross validation, the hyperparameters we used for GBM model are listed as followed.

- n.trees = 300 – number of iterations/trees
- shrinkage = 0.15 – learning rate
- interaction.depth = 2 – number of splits on a tree
- n.minobsinnode = 10 – minimum number of observations in trees' terminal nodes

```

### load models built
source("../lib/train_gbm.R")      ### train model

### fit train data
tm_train=NA

```

```
tm_train_gbm <- system.time(fit_train_gbm <- train_gbm(feature_df = dat_train))
save(fit_train_gbm, file="../output/gbm_train.RData")
```

Please see the details of model as below.

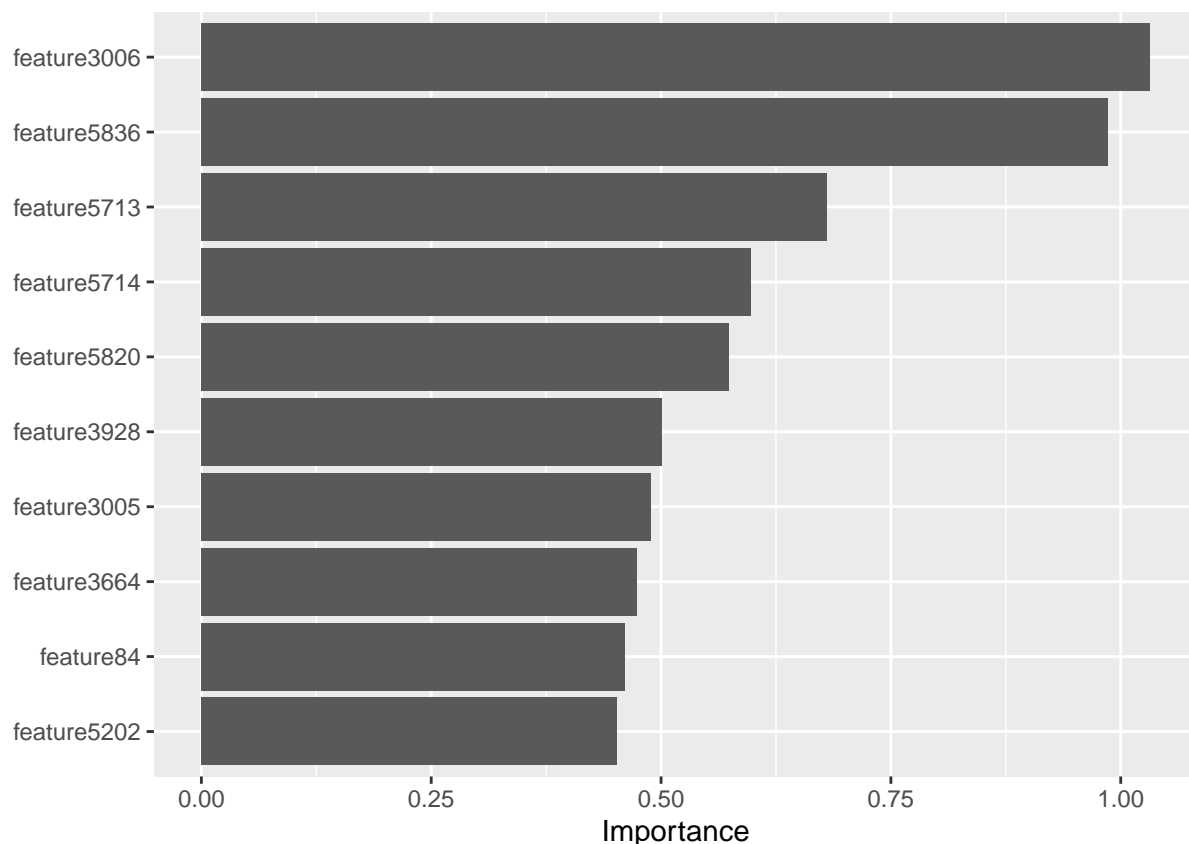
```
fit_train_gbm
```

```
## gbm(formula = emotion_idx ~ ., distribution = "multinomial",
##      data = feature_df, n.trees = 300, interaction.depth = 2,
##      n.minobsinnode = 10, shrinkage = 0.15)
## A gradient boosted model with multinomial loss function.
## 300 iterations were performed.
## There were 6006 predictors of which 3604 had non-zero influence.
```

Step 4.2: Visualize variables importances of train model

After running our baseline model, we wanted to understand the variables that influence the emotion prediction largely. So, we plotted the most influential variables. Here, we utilized the default method - relative influence to compute the variables importance, which shows the average improvement obtained by each variable across all trees that use the variable.

```
vip::vip(fit_train_gbm)
```



Step 4.3: Predict test data and evaluate

Then, we predicted the test data and evaluated the performance of model.

```
### load models built
source("../lib/test_gbm.R")      ### test model
```

```

### predict test data
tm_test_gbm=NA
if(run.test){
  load(file="../output/gbm_train.RData")
  tm_test_gbm <- system.time(pred <- test_gbm(model_best = fit_train_gbm, feature_test = dat_test))
}

```

Evaluation on GBM model

```

pred_model = predict(fit_train_gbm, dat_test, n.trees = 300, type = "response")
pred_gbm = apply(pred_model, 1, which.max)
accu_gbm <- mean(dat_test$emotion_idx == pred_gbm)
cat("The accuracy of model: gradient boosting machine", "is", accu_gbm*100, "%.\n")

```

The accuracy of model: gradient boosting machine is 40.4 %.

```

pred_gbm = as.factor(pred_gbm)
confusionMatrix(pred_gbm, dat_test$emotion_idx)

```

Confusion Matrix and Statistics

```

##
##           Reference
## Prediction  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22
##           1 16  0  3  0  0  0  0  0  0  1  0  0  0  0  0  1  0  1  0  0  0  0
##           2  0 15  0  0  0  0  0  1  3  0  0  0  0  0  0  0  1  0  0  1  1  0
##           3  3  0 11  4  0  0  1  0  0  2  2  1  1  0  0  0  0  1  0  0  1
##           4  0  0  2 11  0  0  0  0  0  0  2  0  2  0  0  0  0  0  0  0  0  1
##           5  0  1  0  0 11  0  0  2  0  0  0  0  0  0  0  0  5  2  1  0  0  0
##           6  1  0  1  4  0  8  1  0  0  0  0  0  1  0  0  0  0  0  0  0  0  1
##           7  0  1  0  1  1  1  6  1  0  1  0  0  0  0  1  0  1  1  0  0  1  0
##           8  0  6  0  0  0  0  0 17  0  0  0  0  0  0  0  0  1  0  0  0  0  0
##           9  0  4  0  0  0  0  0  0 11  0  0  0  0  0  0  0  0  0  0  1  1  1
##          10  0  0  1  3  0  1  0  0  0  7  3  2  7  0  0  1  0  0  0  0  0  2
##          11  0  0  1  2  0  4  0  0  1  4 10  1  1  0  0  0  0  0  0  0  0  0
##          12  0  0  1  4  0  2  0  0  0  1  5  5  4  0  1  1  0  0  1  0  0  5
##          13  0  0  2  4  0  1  0  0  0  4  2  5  6  0  0  0  0  0  0  0  2  0
##          14  0  0  0  0  0  0  0  0  0  0  0  1  1 10  4  0  0  0  1  1  2  0
##          15  0  0  0  0  0  0  0  0  0  0  0  0  0  3  3  0  0  0  2  0  0  1
##          16  0  0  0  1  0  0  0  0  0  0  0  0  1  0  0 15  0  1  0  0  0  0
##          17  0  0  0  0  3  0  1  2  0  0  0  0  0  1  1  1 10  4  1  2  0  0
##          18  0  0  0  0  2  0  1  0  0  0  0  0  0  1  1  0 10  4  2  0  0  1
##          19  0  0  1  0  0  0  9  0  0  1  0  0  1  0  6  0  3  0  6  4  4  3
##          20  0  2  0  0  1  0  0  0  1  0  1  0  0  2  0  1  1  1  0  8  2  2
##          21  0  0  1  0  0  0  2  0  0  0  0  1  0  2  2  0  1  1  4  3  6  0
##          22  0  1  5  0  0  1  1  0  0  1  0  4  1  2  0  2  0  0  5  0  0  6

```

Overall Statistics

```

##
##           Accuracy : 0.404
##           95% CI   : (0.3607, 0.4485)
##           No Information Rate : 0.068
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.3752
##

```

```
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##          Class: 1 Class: 2 Class: 3 Class: 4 Class: 5 Class: 6
## Sensitivity      0.8000    0.5000    0.3793    0.3235    0.6111    0.4444
## Specificity      0.9875    0.9851    0.9660    0.9850    0.9772    0.9813
## Pos Pred Value   0.7273    0.6818    0.4074    0.6111    0.5000    0.4706
## Neg Pred Value   0.9916    0.9686    0.9619    0.9523    0.9854    0.9793
## Prevalence       0.0400    0.0600    0.0580    0.0680    0.0360    0.0360
## Detection Rate   0.0320    0.0300    0.0220    0.0220    0.0220    0.0160
## Detection Prevalence 0.0440    0.0440    0.0540    0.0360    0.0440    0.0340
## Balanced Accuracy 0.8938    0.7426    0.6727    0.6543    0.7941    0.7129
##
##          Class: 7 Class: 8 Class: 9 Class: 10 Class: 11 Class: 12
## Sensitivity      0.2727    0.7391    0.6875    0.3182    0.4000    0.2500
## Specificity      0.9791    0.9853    0.9855    0.9582    0.9705    0.9479
## Pos Pred Value   0.3750    0.7083    0.6111    0.2593    0.4167    0.1667
## Neg Pred Value   0.9669    0.9874    0.9896    0.9683    0.9685    0.9681
## Prevalence       0.0440    0.0460    0.0320    0.0440    0.0500    0.0400
## Detection Rate   0.0120    0.0340    0.0220    0.0140    0.0200    0.0100
## Detection Prevalence 0.0320    0.0480    0.0360    0.0540    0.0480    0.0600
## Balanced Accuracy 0.6259    0.8622    0.8365    0.6382    0.6853    0.5990
##
##          Class: 13 Class: 14 Class: 15 Class: 16 Class: 17
## Sensitivity      0.2308    0.4762    0.1579    0.6818    0.3030
## Specificity      0.9578    0.9791    0.9875    0.9937    0.9657
## Pos Pred Value   0.2308    0.5000    0.3333    0.8333    0.3846
## Neg Pred Value   0.9578    0.9771    0.9674    0.9855    0.9515
## Prevalence       0.0520    0.0420    0.0380    0.0440    0.0660
## Detection Rate   0.0120    0.0200    0.0060    0.0300    0.0200
## Detection Prevalence 0.0520    0.0400    0.0180    0.0360    0.0520
## Balanced Accuracy 0.5943    0.7277    0.5727    0.8378    0.6344
##
##          Class: 18 Class: 19 Class: 20 Class: 21 Class: 22
## Sensitivity      0.2667    0.2500    0.4000    0.3158    0.2500
## Specificity      0.9629    0.9328    0.9708    0.9647    0.9517
## Pos Pred Value   0.1818    0.1579    0.3636    0.2609    0.2069
## Neg Pred Value   0.9770    0.9610    0.9749    0.9727    0.9618
## Prevalence       0.0300    0.0480    0.0400    0.0380    0.0480
## Detection Rate   0.0080    0.0120    0.0160    0.0120    0.0120
## Detection Prevalence 0.0440    0.0760    0.0440    0.0460    0.0580
## Balanced Accuracy 0.6148    0.5914    0.6854    0.6402    0.6008

cat("Time for training model GBM = ", tm_train_gbm[1], "s \n")

## Time for training model GBM = 1806.874 s

cat("Time for testing model GBM = ", tm_test_gbm[1], "s \n")

## Time for testing model GBM = 13.104 s
```

Step 4.4: Summary

Overall, the predictive accuracy of GBM is not bad. Also, it doesn't need data pre-processing and handles multiple categorical value well. However, it is computational expensive to build many trees and tune multiple parameters based on cross validation. The result of model is less interpretable and easily understood as we built models on fiducial points distances. To improve the accuracy and eliminate other disadvantages of GBM, we also tried other models such as LDA, SVM, RF, and XGBoost.

Step 5: Improved Models - Linear Discriminant Analysis (LDA), Support Vector Machine (SVM), Random Forest, Extreme Gradient Boosting (XGBoost), Convolutional Neural Network (CNN)

Step 5.1: Linear Discriminant Analysis (LDA) with Principal Component Analysis (PCA)

PCA

Extract first principal component from feature data

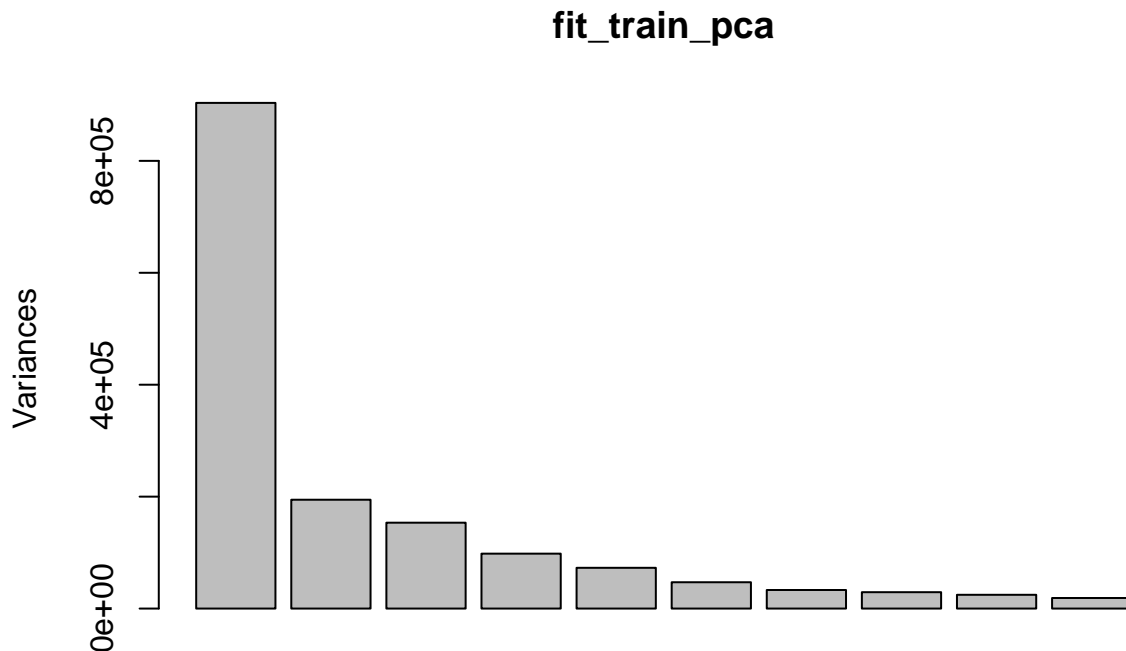
```
#load(file="../output/feature_train.RData")
#load(file="../output/feature_test.RData")

## Convert dat_train to numeric data.frame
dat_train.new <- matrix(0, ncol = ncol(dat_train) - 1, nrow = nrow(dat_train))
for (i in 1:(ncol(dat_train) - 1)) {
  dat_train.new[,i] <- as.numeric(dat_train[[i]])
}
dat_train.new <- as.data.frame(dat_train.new)

## PCA for training features
source("../lib/train_pca.R")
tm_train_pca <- NA
tm_train_pca <- system.time(fit_train_pca <- train_pca(dat_train.new))
save(fit_train_pca, file="../output/pca_train.RData")
```

Visualization of proportion of variance to choose number of principle components

```
screepplot(fit_train_pca)
```



```
## The proportion of variance for first 500 PCs
sum((fit_train_pca$sdev[1:500])^2) / sum((fit_train_pca$sdev)^2)
```

```
## [1] 0.9997662
```

```
## Combine pc_features with the emotion index
dat_train_pca <- data.frame(fit_train_pca$x[,1:500], emotion_idx = dat_train[,6007])
```

Apply PC model to test data

```
## Extract PC from test data
source("../lib/test_pca.R")
dat_test.new <- dat_test
colnames(dat_test.new) <- c(colnames(dat_train.new), "emotion_idx")
tm_test_pcs <- NA
tm_test_pca <- system.time(dat_test.new <- test_pca(fit_train_pca, dat_test.new))
```

```
## Combine pc_features with the emotion index
dat_test_pca <- data.frame(dat_test.new[,1:500], emotion_idx = dat_test[,6007])
```

```
save(dat_train_pca, file="../output/feature_train_pca.RData")
save(dat_test_pca, file="../output/feature_test_pca.RData")
```

```
#load("../output/feature_train_pca.RData")
#load("../output/feature_test_pca.RData")
```

Evaluation on PCA model

```
cat("Time for training model PCA = ", tm_train_pca[1], "s \n")
```

```
## Time for training model PCA = 86.123 s
```

```
cat("Time for testing model PCA = ", tm_test_pca[1], "s \n")
```

```
## Time for testing model PCA = 4.27 s
```

LDA

Train the model with the entire training set using the selected model (model parameter) via cross-validation.

```
source("../lib/train_lda.R")
tm_train_lda = NA
tm_train_lda <- system.time(fit_train_lda <- train_lda(dat_train_pca, par_best))
save(fit_train_lda, file="../output/lda_train.RData")
```

Predicted the test data using LDA model

```
source("../lib/test_lda.R")
tm_test_lda = NA
tm_test_lda <- system.time(pred_lda <- test_lda(fit_train_lda, dat_test_pca))
```

Evaluation on LDA model

```
accu_lda <- mean(dat_test_pca$emotion_idx == pred_lda)
cat("The accuracy of model: LDA", "is", accu_lda*100, "%.\n")
```

```
## The accuracy of model: LDA is 46.6 %.
```

```
confusionMatrix(pred_lda, dat_test_pca$emotion_idx)
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```

## Prediction  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22
##           1 17  1  1  2  2  0  0  0  0  0  0  2  1  0  0  0  1  0  0  1  0
##           2  0 22  0  0  0  0  0  0  3  0  0  0  0  0  0  0  0  0  1  0  0
##           3  0  1 19  1  0  1  0  0  1  5  0  3  2  0  1  1  0  2  0  1  0  2
##           4  0  0  0 13  0  0  0  0  0  0  0  4  0  0  0  0  0  1  0  1  1
##           5  0  0  0  0 10  0  1  0  0  0  0  0  2  0  0  3  1  0  0  0  0
##           6  1  0  0  2  0  7  1  0  0  0  2  0  0  0  1  0  0  0  1  0  1
##           7  0  0  0  1  0  0  8  0  0  1  0  0  0  0  1  0  1  1  0  2  0
##           8  0  1  0  0  1  0  1 20  0  0  0  0  0  0  0  0  0  0  0  0  0
##           9  0  3  0  0  0  0  0  0 11  0  0  0  1  0  0  0  0  0  0  1  1
##          10  2  0  1  2  0  1  0  0  0  8  2  1  0  0  0  1  0  0  1  0  0  1
##          11  0  0  0  1  0  1  0  0  0  0 10  1  1  1  0  0  0  0  0  0  0  0
##          12  0  0  0  2  0  4  0  0  0  0  4  6  5  0  0  1  0  0  1  0  0  2
##          13  0  0  3  5  0  3  0  0  0  5  5  4  7  0  0  0  0  0  1  0  1  2
##          14  0  0  0  0  0  0  1  1  0  0  0  0  1  9  5  0  0  0  1  0  1  0
##          15  0  0  0  0  1  0  0  0  0  0  0  1  0  2  5  0  0  0  0  0  0  0
##          16  0  0  0  1  0  1  1  0  0  0  1  1  0  0  0 12  2  0  0  0  0  2
##          17  0  1  0  0  0  0  2  2  0  0  0  0  0  1  0  0 17  2  1  1  1  0
##          18  0  0  0  0  4  0  2  0  0  0  0  0  0  3  1  7  7  3  0  0  2
##          19  0  0  1  1  0  0  2  0  0  1  0  0  1  2  4  0  1  1  7  2  3  3
##          20  0  0  1  0  0  0  1  0  0  0  2  0  1  1  1  0  3  0  2  5  1  0
##          21  0  1  2  2  0  0  2  0  0  0  0  1  0  2  0  0  0  0  2  5  6  0
##          22  0  0  1  1  0  0  0  0  1  2  1  0  1  0  0  4  0  0  3  4  1  7
##
## Overall Statistics
##
##           Accuracy : 0.466
##           95% CI : (0.4216, 0.5108)
##           No Information Rate : 0.068
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.4399
##
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: 1 Class: 2 Class: 3 Class: 4 Class: 5 Class: 6
## Sensitivity      0.8500    0.7333    0.6552    0.3824    0.5556    0.3889
## Specificity      0.9771    0.9915    0.9554    0.9850    0.9855    0.9813
## Pos Pred Value   0.6071    0.8462    0.4750    0.6500    0.5882    0.4375
## Neg Pred Value   0.9936    0.9831    0.9783    0.9562    0.9834    0.9773
## Prevalence       0.0400    0.0600    0.0580    0.0680    0.0360    0.0360
## Detection Rate   0.0340    0.0440    0.0380    0.0260    0.0200    0.0140
## Detection Prevalence 0.0560    0.0520    0.0800    0.0400    0.0340    0.0320
## Balanced Accuracy 0.9135    0.8624    0.8053    0.6837    0.7705    0.6851
##
##           Class: 7 Class: 8 Class: 9 Class: 10 Class: 11 Class: 12
## Sensitivity      0.3636    0.8696    0.6875    0.3636    0.4000    0.3000
## Specificity      0.9854    0.9937    0.9876    0.9749    0.9895    0.9604
## Pos Pred Value   0.5333    0.8696    0.6471    0.4000    0.6667    0.2400
## Neg Pred Value   0.9711    0.9937    0.9896    0.9708    0.9691    0.9705
## Prevalence       0.0440    0.0460    0.0320    0.0440    0.0500    0.0400
## Detection Rate   0.0160    0.0400    0.0220    0.0160    0.0200    0.0120
## Detection Prevalence 0.0300    0.0460    0.0340    0.0400    0.0300    0.0500

```


## Balanced Accuracy	0.6745	0.9316	0.8376	0.6693	0.6947	0.6302
##	Class: 13	Class: 14	Class: 15	Class: 16	Class: 17	
## Sensitivity	0.2692	0.4286	0.2632	0.5455	0.5152	
## Specificity	0.9388	0.9791	0.9917	0.9812	0.9764	
## Pos Pred Value	0.1944	0.4737	0.5556	0.5714	0.6071	
## Neg Pred Value	0.9591	0.9751	0.9715	0.9791	0.9661	
## Prevalence	0.0520	0.0420	0.0380	0.0440	0.0660	
## Detection Rate	0.0140	0.0180	0.0100	0.0240	0.0340	
## Detection Prevalence	0.0720	0.0380	0.0180	0.0420	0.0560	
## Balanced Accuracy	0.6040	0.7038	0.6274	0.7633	0.7458	
##	Class: 18	Class: 19	Class: 20	Class: 21	Class: 22	
## Sensitivity	0.4667	0.2917	0.2500	0.3158	0.2917	
## Specificity	0.9546	0.9538	0.9729	0.9647	0.9601	
## Pos Pred Value	0.2414	0.2414	0.2778	0.2609	0.2692	
## Neg Pred Value	0.9830	0.9639	0.9689	0.9727	0.9641	
## Prevalence	0.0300	0.0480	0.0400	0.0380	0.0480	
## Detection Rate	0.0140	0.0140	0.0100	0.0120	0.0140	
## Detection Prevalence	0.0580	0.0580	0.0360	0.0460	0.0520	
## Balanced Accuracy	0.7107	0.6227	0.6115	0.6402	0.6259	

```
cat("Time for training model LDA = ", tm_train_lda[1], "s \n")
```

```
## Time for training model LDA = 2.571 s
```

```
cat("Time for testing model LDA = ",tm_test_lda[1], "s \n")
```

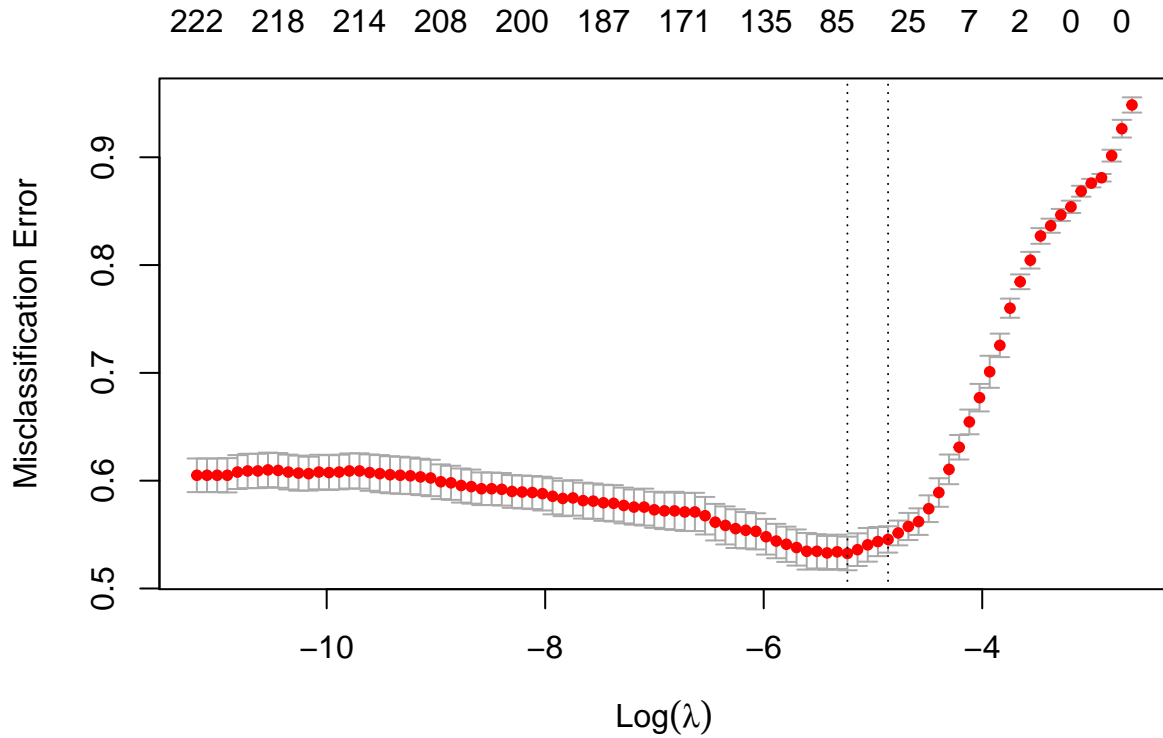
```
## Time for testing model LDA = 0.028 s
```

Prediction performance matters, so does the running times for constructing features and for training the model, especially when the computation resource is limited.

Step 5.5: Logistic Regression Model (LR) with Principal Component Analysis (PCA)

Train the model with the entire training set using the selected model (model parameter) via cross-validation.

```
source("../lib/train_lr.R")
tm_train_lr=NA
tm_train_lr <- system.time(fit_train_lr <- train_lr(dat_train_pca))
plot(fit_train_lr)
```



```
opt_lambda <- fit_train_lr$lambda.1se
fit_train_lr <- fit_train_lr$glmnet.fit
save(fit_train_lr, file="../output/lr_train.RData")
```

Predicted the test data using LR model

```
source("../lib/test_lr.R")
tm_test_lr = NA
tm_test_lr <- system.time(pred_lr <- test_lr(fit_train_lr, dat_test_pca))
```

Evaluation on LR model

```
accu_lr <- mean(dat_test_pca$emotion_idx == pred_lr)
cat("The accuracy of model: LR", "is", accu_lr*100, "%.\n")
```

```
## The accuracy of model: LR is 39.16129 %.
```

```
a <- matrix(as.vector(dat_test_pca$emotion_idx), 500, 93)
confusionMatrix(as.factor(pred_lr), as.factor(a))
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##
```

```
##          Reference
## Prediction   1   10   11   12   13   14   15   16   17   18   19   2   20
##          1 1240   47    0    0   66   59   13  122    3  145    3    0    0
##          10   73  428   50   82  195   79   11    9    0    0   36   26    5
##          11   91  129 1079  166  132   17    2   30   30   18   29    0   38
##          12    6  113  301  420  415    0    1   20    0    0   31    0   73
##          13   42  265  134  177  503  135   41    0    0    0    0    0   12
##          14   11   79  148  131  116  922  518    3   41   12  241   14  147
##          15    4   20    0   73  100  325  440    2   23   74   90    0    0
##          16   47   28    1   46   21    0    0  992  159   22    5    0   53
##          17   55    0    0    0    2   80   15   11 1201  127  159   28   83
```

```

##      18      6      0      0      0      1      2 148 102 511 312 76 0 0
##      19      0     88      0      4     31 127 82 7 47 68 364 10 28
##      2      0      0      7      0      0      0 0 0 95 0 0 1837 114
##      20      0     70    143     49     49 29 74 83 201 120 198 0 286
##      21      0     30     38     12     89 5 129 0 0 1 181 51 394
##      22     49    122    105     86     30 1 1 121 79 75 159 0 289
##      3    178    401    151    222    225 39 58 245 176 130 137 144 111
##      4     22    130    119     72    345 26 39 80 0 0 189 0 33
##      5     36     12      0      0     26 40 134 39 355 164 114 28 11
##      6      0      1     30    309 40 47 0 126 0 0 20 14 52
##      7      0     83      0      0      1 3 0 31 37 4 131 2 11
##      8      0      0      0      0      5 0 0 23 92 116 6 433 38
##      9      0      0     19     11     26 17 61 0 19 7 63 203 82
##
##      Reference
## Prediction 21 22 3 4 5 6 7 8 9
##      1      5     35    106    192    26 48 9 0 0
##     10     22     82    153    182     0 85 44 0 13
##     11     12     36    162    216    14 157 0 6 8
##     12      0    265     59    190     0 388 0 0 16
##     13      6    203    127    255     0 71 1 0 0
##     14    199    145     41    122    153 90 94 54 45
##     15      8      5     57     0 22 0 67 0 0
##     16     93     89     20     69     6 90 124 0 0
##     17     16     76     50     0 20 0 202 131 0
##     18      0     65     78     4 193 0 82 0 0
##     19    266     78     12     82     0 0 177 0 0
##      2      25      0      0      0      0 0 63 89 402
##     20    101     92      0      5      0 0 7 0 0
##     21    567     95     71      0     19 6 94 0 14
##     22     46    509    120      0      0 17 0 0 64
##      3      82    190 1404 312 45 165 125 71 36
##      4      31     35    103 1146 0 29 11 0 0
##      5      41     40     24     33 1152 0 208 73 0
##      6      0    179     74    207     0 478 96 0 21
##      7      99      4     21    147     0 50 522 0 34
##      8      0      0      4      0     24 0 65 1696 123
##      9    148      9     11      0      0 0 55 19 712
##
## Overall Statistics
##
##      Accuracy : 0.3916
##      95% CI : (0.3872, 0.3961)
##      No Information Rate : 0.068
##      P-Value [Acc > NIR] : < 2.2e-16
##
##      Kappa : 0.3616
##
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##      Class: 1 Class: 10 Class: 11 Class: 12 Class: 13 Class: 14
## Sensitivity      0.66667 0.209189 0.46409 0.225806 0.20802 0.47209
## Specificity      0.98031 0.974198 0.97073 0.957930 0.96668 0.94603

```

```

## Pos Pred Value      0.58518  0.271746  0.45489  0.182768  0.25507  0.27721
## Neg Pred Value      0.98603  0.963984  0.97176  0.967422  0.95699  0.97612
## Prevalence          0.04000  0.044000  0.05000  0.040000  0.05200  0.04200
## Detection Rate      0.02667  0.009204  0.02320  0.009032  0.01082  0.01983
## Detection Prevalence 0.04557  0.033871  0.05101  0.049419  0.04241  0.07153
## Balanced Accuracy    0.82349  0.591693  0.71741  0.591868  0.58735  0.70906
##
## Class: 15 Class: 16 Class: 17 Class: 18 Class: 19 Class: 2
## Sensitivity          0.249010  0.48485  0.39133  0.22366  0.163082  0.65842
## Specificity          0.980551  0.98036  0.97571  0.97189  0.974993  0.98181
## Pos Pred Value      0.335878  0.53190  0.53236  0.19747  0.247451  0.69795
## Neg Pred Value      0.970635  0.97639  0.95778  0.97589  0.958516  0.97828
## Prevalence          0.038000  0.04400  0.06600  0.03000  0.048000  0.06000
## Detection Rate      0.009462  0.02133  0.02583  0.00671  0.007828  0.03951
## Detection Prevalence 0.028172  0.04011  0.04852  0.03398  0.031634  0.05660
## Balanced Accuracy    0.614780  0.73261  0.68352  0.59777  0.569038  0.82012
##
## Class: 20 Class: 21 Class: 22 Class: 3 Class: 4 Class: 5
## Sensitivity          0.153763  0.32088  0.22805  0.52058  0.36243  0.68817
## Specificity          0.972648  0.97253  0.96919  0.92596  0.97083  0.96926
## Pos Pred Value      0.189781  0.31570  0.27176  0.30213  0.47552  0.45534
## Neg Pred Value      0.965017  0.97316  0.96139  0.96911  0.95428  0.98813
## Prevalence          0.040000  0.03800  0.04800  0.05800  0.06800  0.03600
## Detection Rate      0.006151  0.01219  0.01095  0.03019  0.02465  0.02477
## Detection Prevalence 0.032409  0.03862  0.04028  0.09994  0.05183  0.05441
## Balanced Accuracy    0.563206  0.64670  0.59862  0.72327  0.66663  0.82872
##
## Class: 6 Class: 7 Class: 8 Class: 9
## Sensitivity          0.28554  0.25513  0.79289  0.47849
## Specificity          0.97287  0.98520  0.97906  0.98334
## Pos Pred Value      0.28217  0.44237  0.64610  0.48700
## Neg Pred Value      0.97331  0.96637  0.98990  0.98277
## Prevalence          0.03600  0.04400  0.04600  0.03200
## Detection Rate      0.01028  0.01123  0.03647  0.01531
## Detection Prevalence 0.03643  0.02538  0.05645  0.03144
## Balanced Accuracy    0.62921  0.62017  0.88598  0.73092

cat("Time for training model LR = ", tm_train_lr[1], "s \n")

## Time for training model LR = 419.8 s

cat("Time for testing model LR = ", tm_test_lr[1], "s \n")

## Time for testing model LR = 0.349 s

```

Step 5.3: Support Vector Machine (SVM)

Then, we tried Support Vector Machine Model. An SVM training algorithm builds a model that assigns new examples to one category or the other, making it a non-probabilistic binary linear classifier.

The global argument we need to consider is stated below:

- kernel: the kernel used in training and predicting
- degree: parameter needed for kernel of type polynomial (default: 3)
- gamma: parameter needed for all kernels except linear (default: 1/(data dimension))

We tried linear, polynomial and radial basis kernels separately, explored different parameters (degree for polynomial kernel and gamma for radial basis kernel) for this model and found the linear kernel for SVM model fit the data best.

```
source("../lib/train_svm.R")
tm_train_svm <- NA
tm_train_svm <- system.time(fit_train_svm <- train_svm(dat_train))
save(fit_train_svm, file="../output/svm_train.RData")
```

Predicted the test data using SVM model

```
source("../lib/test_svm.R")
tm_test_svm <- NA
tm_test_svm <- system.time(pred_svm <- test_svm(fit_train_svm, dat_test))
```

Evaluation on SVM model

```
accu_svm <- mean(dat_test$emotion_idx == pred_svm)
cat("The accuracy of model: SVM", "is", accu_svm*100, "%.\n")
```

The accuracy of model: SVM is 48.6 %.

```
confusionMatrix(pred_svm, dat_test$emotion_idx)
```

Confusion Matrix and Statistics

```
##
##           Reference
## Prediction  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22
##           1 18  0  3  1  0  0  0  0  0  0  0  0  0  0  1  0  2  0  0  0  0
##           2  0 19  0  0  0  0  0  1  2  0  0  0  0  0  0  1  0  0  0  0  0
##           3  0  0 12  0  0  0  0  0  0  2  0  0  4  0  0  2  0  0  0  0  3
##           4  1  0  1 16  0  0  0  0  0  2  1  2  4  0  0  0  0  1  1  1  1
##           5  0  0  0  0 11  0  0  0  0  0  0  0  0  1  0  1  0  0  0  0  0
##           6  0  1  0  1  0 11  0  0  0  0  1  1  0  0  0  1  0  0  0  0  1
##           7  0  0  0  0  0  0  6  0  0  0  0  0  0  0  0  1  0  0  0  0  0
##           8  0  5  0  0  0  0  1 20  1  0  0  0  0  0  0  0  1  0  0  0  0
##           9  0  5  0  0  0  0  0  0 12  0  0  0  0  0  0  0  0  0  2  0  0
##          10  0  0  2  3  0  0  0  0  0  9  4  1  4  0  0  1  0  0  0  1  0  2
##          11  0  0  2  5  0  3  0  0  0  3 12  4  1  1  0  0  0  0  0  0  0  0
##          12  1  0  2  2  0  2  0  0  0  2  2  6  4  0  0  0  0  0  0  0  3
##          13  0  0  2  5  0  1  0  0  0  3  4  1  6  1  0  0  0  0  0  1  0
##          14  0  0  0  0  0  0  0  0  0  1  0  0  0 10  6  1  0  0  2  1  2  0
##          15  0  0  0  0  1  0  1  0  0  0  0  0  1  5  8  0  0  0  3  0  0  2
##          16  0  0  0  1  0  0  0  0  0  0  0  1  1  0  0 13  1  1  0  0  0  0
##          17  0  0  0  0  1  0  5  2  0  0  0  0  0  1  1  2 21  4  1  2  1  1
##          18  0  0  2  0  4  0  1  0  0  0  0  0  0  0  0  6  6  1  0  0  1
##          19  0  0  1  0  0  0  5  0  0  0  0  0  1  1  2  0  1  0  6  1  3  1
##          20  0  0  0  0  1  0  3  0  1  0  0  0  0  1  0  0  1  1  3  7  3  1
##          21  0  0  1  0  0  0  0  0  0  0  0  0  0  0  1  0  0  0  3  4  7  1
##          22  0  0  1  0  0  1  0  0  0  0  1  4  0  1  0  1  0  0  4  1  1  7
```

Overall Statistics

```
##
##           Accuracy : 0.486
##           95% CI   : (0.4414, 0.5308)
##           No Information Rate : 0.068
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.4605
##
```

```
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##          Class: 1 Class: 2 Class: 3 Class: 4 Class: 5 Class: 6
## Sensitivity      0.9000    0.6333    0.4138    0.4706    0.6111    0.6111
## Specificity      0.9854    0.9915    0.9766    0.9678    0.9959    0.9876
## Pos Pred Value   0.7200    0.8261    0.5217    0.5161    0.8462    0.6471
## Neg Pred Value   0.9958    0.9769    0.9644    0.9616    0.9856    0.9855
## Prevalence       0.0400    0.0600    0.0580    0.0680    0.0360    0.0360
## Detection Rate   0.0360    0.0380    0.0240    0.0320    0.0220    0.0220
## Detection Prevalence 0.0500    0.0460    0.0460    0.0620    0.0260    0.0340
## Balanced Accuracy 0.9427    0.8124    0.6952    0.7192    0.8035    0.7993
##
##          Class: 7 Class: 8 Class: 9 Class: 10 Class: 11 Class: 12
## Sensitivity      0.2727    0.8696    0.7500    0.4091    0.4800    0.3000
## Specificity      0.9979    0.9832    0.9855    0.9623    0.9600    0.9625
## Pos Pred Value   0.8571    0.7143    0.6316    0.3333    0.3871    0.2500
## Neg Pred Value   0.9675    0.9936    0.9917    0.9725    0.9723    0.9706
## Prevalence       0.0440    0.0460    0.0320    0.0440    0.0500    0.0400
## Detection Rate   0.0120    0.0400    0.0240    0.0180    0.0240    0.0120
## Detection Prevalence 0.0140    0.0560    0.0380    0.0540    0.0620    0.0480
## Balanced Accuracy 0.6353    0.9264    0.8678    0.6857    0.7200    0.6312
##
##          Class: 13 Class: 14 Class: 15 Class: 16 Class: 17
## Sensitivity      0.2308    0.4762    0.4211    0.5909    0.6364
## Specificity      0.9620    0.9729    0.9730    0.9895    0.9550
## Pos Pred Value   0.2500    0.4348    0.3810    0.7222    0.5000
## Neg Pred Value   0.9580    0.9769    0.9770    0.9813    0.9738
## Prevalence       0.0520    0.0420    0.0380    0.0440    0.0660
## Detection Rate   0.0120    0.0200    0.0160    0.0260    0.0420
## Detection Prevalence 0.0480    0.0460    0.0420    0.0360    0.0840
## Balanced Accuracy 0.5964    0.7245    0.6970    0.7902    0.7957
##
##          Class: 18 Class: 19 Class: 20 Class: 21 Class: 22
## Sensitivity      0.4000    0.2500    0.3500    0.3684    0.2917
## Specificity      0.9691    0.9664    0.9688    0.9792    0.9685
## Pos Pred Value   0.2857    0.2727    0.3182    0.4118    0.3182
## Neg Pred Value   0.9812    0.9623    0.9728    0.9752    0.9644
## Prevalence       0.0300    0.0480    0.0400    0.0380    0.0480
## Detection Rate   0.0120    0.0120    0.0140    0.0140    0.0140
## Detection Prevalence 0.0420    0.0440    0.0440    0.0340    0.0440
## Balanced Accuracy 0.6845    0.6082    0.6594    0.6738    0.6301
```

```
cat("Time for training model SVM = ", tm_train_svm[1], "s \n")
```

```
## Time for training model SVM = 93.083 s
```

```
cat("Time for testing model SVM = ", tm_test_svm[1], "s \n")
```

```
## Time for testing model SVM = 8.667 s
```

Step 5.4: Random Forest

Random Forest is another approach that we consider as a candidate of our advanced model. It is a non-parametric model that would give us more flexibility but along with the risk of overfitting. Random forest is a popular algorithm for classification in a long time, and it dominant a lot of other approaches in application.

However, tuning a random forest model is not easy, the global argument we need to consider is stated below:

- `ntrees`: the number of trees in the forest
- `mtry`: the number of features to consider when looking for the best split
- `max_depth`: the maximum depth of each tree

To tune these three parameters using grid search and cross validation is super time costing, since we are dealing with a multi-dimensional data. So we decide only tune the `mtry` parameter which we consider as the most important parameter that affect the model performance.

Thus, after tuning the `mtry` parameter, we get the optimal parameter as `mtry = 77`. We then plug in this optimal parameter into our training process.

```
## Fit the train model and record the running time
source("../lib/train_rf.R")
tm_train_rf <- NA
tm_train_rf <- system.time(fit_train_rf <- train_rf(dat_train))
save(fit_train_rf, file="../output/rf_train.RData")
```

The training process is already encapsulated in the `train_rf.R` file and the training time is recorded as `tm_train_rf`.

```
## Use the fitted model to make prediction
source("../lib/test_rf.R")
tm_test_rf <- NA
tm_test_rf <- system.time(pred_rf <- test_rf(fit_train_rf, dat_test = dat_test))
accu_rf <- mean(dat_test$emotion_idx == pred_rf)
```

Evaluation on RF model

```
accu_svm <- mean(dat_test$emotion_idx == pred_rf)
cat("The accuracy of model: RF", "is", accu_rf*100, "%.\n")
```

The accuracy of model: RF is 43.2 %.

```
confusionMatrix(pred_rf, dat_test$emotion_idx)
```

Confusion Matrix and Statistics

```
##
##              Reference
## Prediction  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22
##      1    16  1  3  2  0  1  0  0  0  0  0  1  0  0  0  2  0  1  1  0  0  0
##      2     0 19  0  0  0  0  2  0  4  0  0  0  0  0  0  0  0  0  0  0  0  1
##      3     2  0 15  2  0  1  1  0  0  4  1  1  1  0  0  2  0  2  2  0  1  6
##      4     0  0  1 13  0  1  0  0  0  1  2  0  7  0  0  0  0  0  0  0  0  1
##      5     0  0  0  0 12  0  1  1  0  0  0  0  0  0  1  0  4  1  0  0  0  0
##      6     0  0  3  1  0  8  0  0  0  1  3  1  1  0  0  2  0  0  0  0  0  0
##      7     0  0  0  0  0  0  5  0  0  0  0  0  0  0  1  0  1  0  2  0  0  0
##      8     0  7  0  0  0  0  2 21  0  0  0  0  0  0  0  0  3  1  0  1  0  0
##      9     0  2  0  0  0  0  0  0 12  0  1  0  0  0  0  0  0  0  1  3  2  0
##     10     1  0  2  2  0  0  0  0  0  7  2  0  2  0  0  0  0  0  0  0  0  2
##     11     0  0  0  4  0  2  0  0  0  2  8  2  0  1  0  0  0  0  0  0  0  0
##     12     0  0  1  4  0  3  0  0  0  4  6  8  8  0  0  0  0  0  0  0  0  4
##     13     1  0  1  5  0  1  0  0  0  3  1  4  2  0  0  1  0  0  1  0  1  0
##     14     0  0  0  0  2  0  1  0  0  0  1  2  1 14  4  0  0  0  4  1  4  1
##     15     0  0  0  0  1  0  2  0  0  0  0  0  0  3  7  0  1  0  1  0  0  1
##     16     0  0  0  1  0  0  0  0  0  0  0  0  0  0  0 14  0  0  0  0  0  1
##     17     0  0  0  0  0  0  4  0  0  0  0  0  0  0  0  1 16  4  2  3  0  0
##     18     0  0  0  0  3  0  1  0  0  0  0  0  0  1  0  0  7  5  2  0  0  1
##     19     0  0  1  0  0  1  2  0  0  0  0  0  1  0  5  0  0  0  3  1  2  1
```

```

##          20  0  1  0  0  0  0  0  1  0  0  0  0  1  1  1  0  1  1  1  2  2  1
##          21  0  0  1  0  0  0  0  0  0  0  0  0  1  1  0  0  0  0  4  8  6  1
##          22  0  0  1  0  0  0  1  0  0  0  0  1  1  0  0  0  0  0  0  1  1  3
##
## Overall Statistics
##
##          Accuracy : 0.432
##          95% CI : (0.3881, 0.4767)
##          No Information Rate : 0.068
##          P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.4043
##
## Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##          Class: 1 Class: 2 Class: 3 Class: 4 Class: 5 Class: 6
## Sensitivity      0.8000    0.6333    0.5172    0.3824    0.6667    0.4444
## Specificity      0.9750    0.9851    0.9448    0.9721    0.9834    0.9751
## Pos Pred Value   0.5714    0.7308    0.3659    0.5000    0.6000    0.4000
## Neg Pred Value    0.9915    0.9768    0.9695    0.9557    0.9875    0.9792
## Prevalence       0.0400    0.0600    0.0580    0.0680    0.0360    0.0360
## Detection Rate    0.0320    0.0380    0.0300    0.0260    0.0240    0.0160
## Detection Prevalence 0.0560    0.0520    0.0820    0.0520    0.0400    0.0400
## Balanced Accuracy 0.8875    0.8092    0.7310    0.6772    0.8250    0.7098
##
##          Class: 7 Class: 8 Class: 9 Class: 10 Class: 11 Class: 12
## Sensitivity      0.2273    0.9130    0.7500    0.3182    0.3200    0.4000
## Specificity      0.9916    0.9706    0.9814    0.9770    0.9768    0.9375
## Pos Pred Value   0.5556    0.6000    0.5714    0.3889    0.4211    0.2105
## Neg Pred Value    0.9654    0.9957    0.9916    0.9689    0.9647    0.9740
## Prevalence       0.0440    0.0460    0.0320    0.0440    0.0500    0.0400
## Detection Rate    0.0100    0.0420    0.0240    0.0140    0.0160    0.0160
## Detection Prevalence 0.0180    0.0700    0.0420    0.0360    0.0380    0.0760
## Balanced Accuracy 0.6095    0.9418    0.8657    0.6476    0.6484    0.6687
##
##          Class: 13 Class: 14 Class: 15 Class: 16 Class: 17
## Sensitivity      0.07692    0.6667    0.3684    0.6364    0.4848
## Specificity      0.95992    0.9562    0.9813    0.9958    0.9700
## Pos Pred Value   0.09524    0.4000    0.4375    0.8750    0.5333
## Neg Pred Value    0.94990    0.9849    0.9752    0.9835    0.9638
## Prevalence       0.05200    0.0420    0.0380    0.0440    0.0660
## Detection Rate    0.00400    0.0280    0.0140    0.0280    0.0320
## Detection Prevalence 0.04200    0.0700    0.0320    0.0320    0.0600
## Balanced Accuracy 0.51842    0.8114    0.6749    0.8161    0.7274
##
##          Class: 18 Class: 19 Class: 20 Class: 21 Class: 22
## Sensitivity      0.3333    0.1250    0.1000    0.3158    0.1250
## Specificity      0.9691    0.9706    0.9771    0.9667    0.9874
## Pos Pred Value   0.2500    0.1765    0.1538    0.2727    0.3333
## Neg Pred Value    0.9792    0.9565    0.9630    0.9728    0.9572
## Prevalence       0.0300    0.0480    0.0400    0.0380    0.0480
## Detection Rate    0.0100    0.0060    0.0040    0.0120    0.0060
## Detection Prevalence 0.0400    0.0340    0.0260    0.0440    0.0180
## Balanced Accuracy 0.6512    0.5478    0.5385    0.6413    0.5562

```



```
cat("Time for training model RF = ", tm_train_rf[1], "s \n")
```

```
## Time for training model RF = 434.312 s
```

```
cat("Time for testing model RF = ",tm_test_rf[1], "s \n")
```

```
## Time for testing model RF = 0.608 s
```

Step 5.5: Extreme Gradient Boosting (XGBoost)

Because the constructing XGBoost models take quite a long time, we used random search to tune the parameters. With the result of cross validation on the training set, the hyperparameters we used for XGBoost model are listed as followed.

- booster = “gblinear” – the booster to use, can be gbtrees or gblinear, gblinear has better performance and efficiency
- objective = “multi:softmax” – set xgboost to do multiclass classification using the softmax objective, the prediction outputs the class with maximum probability
- eval_metric = “mlogloss” – evaluation metrics for validation data
- lambda = 1.46 – L2 regularization term on weights
- lambda_bias = 0.234 – L2 regularization term on bias
- alpha = 0.0198 – L1 regularization term on weights

```
### load models built
source("../lib/train_xgb.R")      ### train model

### use the optimal parameters
xgb_param_list <- list(
  xgb_para = list(alpha = 0.0198, lambda = 1.46, lambda_bias = 0.234),
  nround = 100
)

### fit train data

tm_train_xgb = NA
tm_train_xgb <- system.time(fit_train_xgb <- train_xgb(feature_df = dat_train, par = xgb_param_list))
save(fit_train_xgb, file="../output/xgb_train.RData")
```

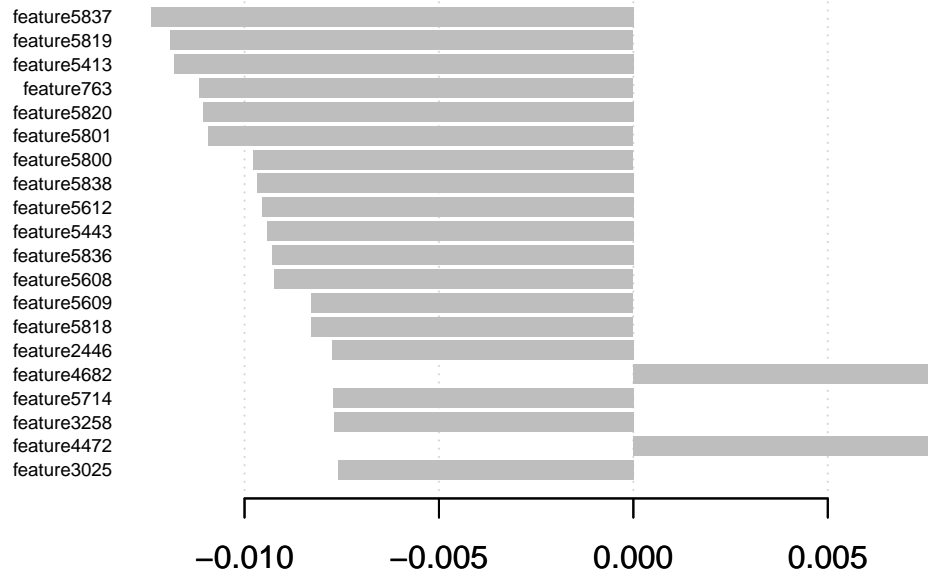
Please see the details of model as below.

```
fit_train_xgb
```

```
## ##### xgb.Booster
## raw: 516.7 Kb
## call:
##   xgb.train(params = par$xgb_para, data = dtrain, nrounds = par$nround,
##     nthread = 6)
## params (as set within xgb.train):
##   alpha = "0.0198", lambda = "1.46", lambda_bias = "0.234", booster = "gblinear", objective = "multi
## xgb.attributes:
##   niter
## callbacks:
##   cb.print.evaluation(period = print_every_n)
## # of features: 6006
## niter: 100
## nfeatures : 6006
```

After running our baseline model, we wanted to understand the variables that influence the emotion prediction largely. So, we plotted the most influential variables.

```
xgb_importance_mat <- xgb.importance(
  feature_names = colnames(dat_train[, -which(names(dat_train) == 'emotion_idx')]),
  model = fit_train_xgb)
xgb.plot.importance (importance_matrix = xgb_importance_mat[1:20])
```



Then, we predicted the test data and evaluated the performance of model.

```
### load models built
source("../lib/test_xgb.R")      ### test model

### predict test data
tm_test_xgb=NA
if(run.test){
  load(file="../output/xgb_train.RData")
  tm_test_xgb <- system.time(pred_xgb <- test_xgb(xgb_model = fit_train_xgb, dat_test = dat_test))
}
```

Evaluation on XGB model

```
confusionMatrix(factor(pred_xgb$max_prob), factor(dat_test$emotion_idx), mode = "everything")
```

Confusion Matrix and Statistics

```
##
##           Reference
## Prediction  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22
##           1 20  0  2  2  0  1  0  0  0  0  0  0  0  0  1  0  1  0  0  0  0
##           2  0 21  0  0  0  0  0  0  3  0  0  0  0  0  0  0  0  0  0  0  0
##           3  0  0 17  1  0  0  0  0  0  1  0  1  2  0  0  3  0  0  1  0  0  3
##           4  0  0  0 18  0  1  0  0  0  1  2  1  6  0  0  0  0  1  0  0  2
##           5  0  0  0  0 11  0  2  0  0  0  0  0  0  0  1  0  1  0  0  0  0
##           6  0  0  0  2  0 11  0  0  1  0  0  1  0  1  0  1  0  0  0  0  1
##           7  0  1  0  0  0  0  8  0  0  0  0  0  0  1  0  0  0  1  0  0  0
##           8  0  4  0  0  1  0  1 21  0  0  0  0  0  0  0  0  0  1  0  1  0
##           9  0  3  0  0  0  0  0  1 12  0  0  0  0  0  0  0  0  0  3  0  0
##          10  0  0  2  3  0  0  0  0  0 11  3  2  3  1  0  1  0  1  0  0  1  2
```

```

##      11  0  0  3  2  0  0  0  0  0  0  10  2  0  0  0  0  0  0  0  0  1
##      12  0  0  0  0  0  3  0  0  0  0  2  5  7  5  0  0  0  0  0  0  5
##      13  0  0  2  5  0  1  0  0  0  0  3  3  3  8  1  0  0  0  0  0  2  0
##      14  0  0  0  0  1  0  0  0  0  0  0  0  0  13  3  0  0  0  1  1  1  0
##      15  0  0  0  0  0  0  2  0  0  0  0  1  0  3  7  0  0  0  2  0  0  0
##      16  0  0  0  0  0  0  1  0  0  0  0  1  1  0  0  16  1  0  0  0  1  1
##      17  0  0  0  0  1  0  3  1  0  0  0  0  0  1  0  0  20  3  1  0  1  0
##      18  0  0  0  0  4  0  1  0  0  0  0  0  1  0  1  0  8  7  2  0  0  0
##      19  0  0  1  0  0  1  2  0  0  1  0  0  0  0  4  0  1  1  9  2  1  1
##      20  0  1  0  0  0  0  1  0  0  0  0  0  0  0  0  0  2  0  0  6  2  2
##      21  0  0  1  0  0  0  1  0  0  1  1  0  0  0  2  0  0  0  4  6  10  1
##      22  0  0  1  1  0  0  0  0  0  2  1  1  0  0  1  0  0  0  3  1  0  5
##
## Overall Statistics
##
##           Accuracy : 0.536
##           95% CI   : (0.4912, 0.5804)
##       No Information Rate : 0.068
##       P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa   : 0.5132
##
## Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: 1 Class: 2 Class: 3 Class: 4 Class: 5 Class: 6
## Sensitivity      1.0000    0.7000    0.5862    0.5294    0.6111    0.6111
## Specificity      0.9854    0.9936    0.9745    0.9700    0.9917    0.9855
## Pos Pred Value   0.7407    0.8750    0.5862    0.5625    0.7333    0.6111
## Neg Pred Value   1.0000    0.9811    0.9745    0.9658    0.9856    0.9855
## Precision        0.7407    0.8750    0.5862    0.5625    0.7333    0.6111
## Recall           1.0000    0.7000    0.5862    0.5294    0.6111    0.6111
## F1               0.8511    0.7778    0.5862    0.5455    0.6667    0.6111
## Prevalence       0.0400    0.0600    0.0580    0.0680    0.0360    0.0360
## Detection Rate   0.0400    0.0420    0.0340    0.0360    0.0220    0.0220
## Detection Prevalence 0.0540    0.0480    0.0580    0.0640    0.0300    0.0360
## Balanced Accuracy 0.9927    0.8468    0.7804    0.7497    0.8014    0.7983
##
##           Class: 7 Class: 8 Class: 9 Class: 10 Class: 11 Class: 12
## Sensitivity      0.3636    0.9130    0.7500    0.5000    0.4000    0.3500
## Specificity      0.9937    0.9832    0.9855    0.9603    0.9832    0.9583
## Pos Pred Value   0.7273    0.7241    0.6316    0.3667    0.5556    0.2593
## Neg Pred Value   0.9714    0.9958    0.9917    0.9766    0.9689    0.9725
## Precision        0.7273    0.7241    0.6316    0.3667    0.5556    0.2593
## Recall           0.3636    0.9130    0.7500    0.5000    0.4000    0.3500
## F1               0.4848    0.8077    0.6857    0.4231    0.4651    0.2979
## Prevalence       0.0440    0.0460    0.0320    0.0440    0.0500    0.0400
## Detection Rate   0.0160    0.0420    0.0240    0.0220    0.0200    0.0140
## Detection Prevalence 0.0220    0.0580    0.0380    0.0600    0.0360    0.0540
## Balanced Accuracy 0.6787    0.9481    0.8678    0.7301    0.6916    0.6542
##
##           Class: 13 Class: 14 Class: 15 Class: 16 Class: 17
## Sensitivity      0.3077    0.6190    0.3684    0.7273    0.6061
## Specificity      0.9578    0.9854    0.9834    0.9874    0.9764
## Pos Pred Value   0.2857    0.6500    0.4667    0.7273    0.6452

```

## Neg Pred Value	0.9619	0.9833	0.9753	0.9874	0.9723
## Precision	0.2857	0.6500	0.4667	0.7273	0.6452
## Recall	0.3077	0.6190	0.3684	0.7273	0.6061
## F1	0.2963	0.6341	0.4118	0.7273	0.6250
## Prevalence	0.0520	0.0420	0.0380	0.0440	0.0660
## Detection Rate	0.0160	0.0260	0.0140	0.0320	0.0400
## Detection Prevalence	0.0560	0.0400	0.0300	0.0440	0.0620
## Balanced Accuracy	0.6327	0.8022	0.6759	0.8574	0.7913
##	Class: 18	Class: 19	Class: 20	Class: 21	Class: 22
## Sensitivity	0.4667	0.3750	0.3000	0.5263	0.2083
## Specificity	0.9649	0.9685	0.9833	0.9647	0.9769
## Pos Pred Value	0.2917	0.3750	0.4286	0.3704	0.3125
## Neg Pred Value	0.9832	0.9685	0.9712	0.9810	0.9607
## Precision	0.2917	0.3750	0.4286	0.3704	0.3125
## Recall	0.4667	0.3750	0.3000	0.5263	0.2083
## F1	0.3590	0.3750	0.3529	0.4348	0.2500
## Prevalence	0.0300	0.0480	0.0400	0.0380	0.0480
## Detection Rate	0.0140	0.0180	0.0120	0.0200	0.0100
## Detection Prevalence	0.0480	0.0480	0.0280	0.0540	0.0320
## Balanced Accuracy	0.7158	0.6717	0.6417	0.7455	0.5926

```
cat("Time for training model XGB = ", tm_train_xgb[1], "s \n")
```

```
## Time for training model XGB = 568.264 s
```

```
cat("Time for testing model XGB = ",tm_test_xgb[1], "s \n")
```

```
## Time for testing model XGB = 0.527 s
```

Overall, the predictive accuracy of XGBoost is beyond 0.5. Also, the accuracy of test set is close to the accuracy of the cross validation on training set. And with the linear booster, the model performs better and more efficient than tree booster. However, it is computational expensive to tune multiple parameters based on cross validation. The result of model is less interpretable and easily understood as we built models on fiducial points distances.

Step 6: Comparison and Summary

From above data, we can choose XGB model as our advanced model and there are some advantages of this model:

- High efficiency in both training and testing process;
- Higher prediction accuracy compared to baseline model;
- Robustness.