

Music Trending Analysis Based on Spotify Data

1st Trent Zhang

dept. Statistics

University of Illinois at Urbana-Champaign

Urbana, IL

zz90@illinois.edu

Abstract—First, we tried few statistical model to know what are the largest influencers on a song's success based on the features extracted from spotify. For top 20 % popular musics, our regression model gives reasonable out-of-sample predictions: ±4.86 based on root-mean-squared-error. Then, we tried multinomial logistic regression, linear discriminative analysis, quadratic discriminative analysis, k-nearest neighbors and decision tree to classify whither a song is popular or not. The best classification model is KNN (K=2) gives out-of-sample predictive over all error rate: 0.391 based on 0-1-error, and error rate of 0.3953 on popular songs.

Index Terms—Trend Prediction, Music Analysis, Regression Analysis

I. INTRODUCTION

A. Background

In 2017, the music industry generated \$8.72 billion in the United States alone. Thanks to growing streaming services (Spotify, Apple Music, etc) the industry continues to flourish. Popular songs secure the lion's share of revenue. The top 10 artists in 2016 generated a combined \$362.5 million in revenue. The question of what makes a song popular has been studied before with varying degrees of success. Every song has key characteristics including lyrics, duration, artist information, temp, beat, loudness, chord, etc.

B. Problem Description

We decided to further investigate by asking three key questions:

- 1) Are there certain characteristics for hit songs?
- 2) What are the largest influencers on a song's success?
- 3) Can we predict the popularity of new songs?

II. DATA PREPARATION

A. Data Source

Predicting how popular one track will be a hard task. To answer these questions, we made use of Spotify Audio Features (Audio features for 115k tracks collected from the official Spotify Web API) provided by tomigelo from kaggle [1], which covered the tracks that released between January 2018 to April 2019.

B. Data Content

This dataset has altogether 130633 row and 17 column. Each track (row) has values for artist name, track name, track id and the audio features [2] itself. Additionally, there is also a popularity feature which is changing over time so it might not be up-to-date when we access the data. All features are listed as TABLE I

TABLE I
FEATURES OF ORIGINAL DATA

	Features		Features
1	artist_name	11	loudness
2	track_id	12	mode
3	track_name	13	speechiness
4	acousticness	14	tempo
5	danceability	15	time_signature
6	duration_ms	16	valence
7	energy	17	popularity
8	instrumentalness		
9	key		
10	liveness		

C. Data Split

The total data is randomly split into train data (70%) and validation data (30%).

III. DESCRIPTIVE ANALYSIS

We depict scatter, density and correlation plot of our numeric variables as Fig. 1, and have some insights from it as follows:

- The more energy, the more loudness;
- Most of popular tracks are loud;
- Most of tracks are non-live, loud, non-speech, 4/4 time signature;
- The popularity data is significantly right skewed.

IV. REGRESSIVE MODEL ANALYSIS OF TOP 20% POPULAR SONG

In this section, in order to know what are the largest influencers on a song's success. We discuss the results of the multiple regression model, ridge model, Lasso model and elastic-net model based on Top 20% Popular Song. Cross validation is used to evaluate the result of model.

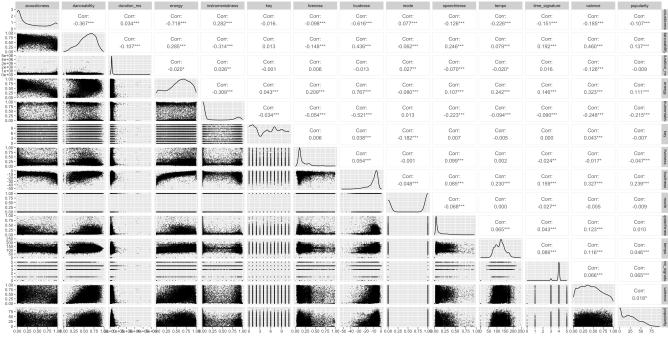


Fig. 1. Scatter, density and correlation plots

A. The Multiple Regression Model

We use linear regression with all variable and get a model with RMSE as TABLE III. Fig. 2 shows the prediction versus the observation when we use all train data to fit the model. Fig. 3 shows the prediction versus the observation when we use the model to test data.

From the figure, test error and we can know that Linear regression is not good at predicting popularity of one track.

TABLE II
COEFFICIENT OF LINEAR REGRESSION

	coefficients
(Intercept)	85.10
acousticness	-0.14
danceability	4.26
duration_ms	-0.00
energy	-2.05
instrumentalness	-1.26
key	-0.07
liveness	1.50
loudness	0.01
mode	0.10
speechiness	0.78
tempo	-0.02
time_signature	0.53
valence	-0.00

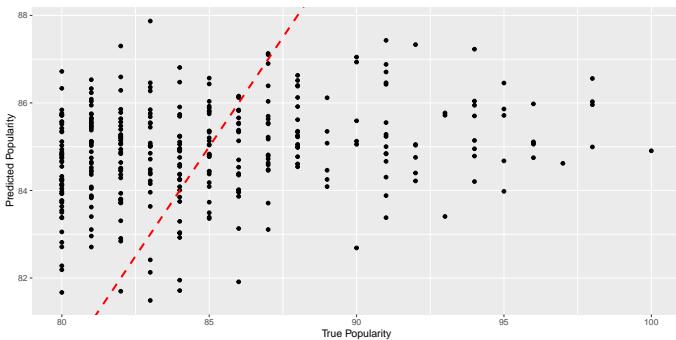


Fig. 2. Prediction Versus The Observation in train

B. Ridge with Cross Validation

But we noticed that there are many variables in this model and the some of the variable is statistically insignificant.

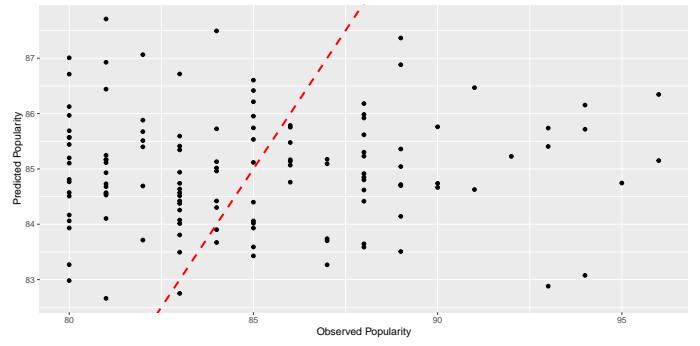


Fig. 3. Prediction Versus The Observation in test

TABLE III
RMSE TABLE OF LINEAR REGRESSION

rmse.lm.train	rmse.lm.test
4.04	4.86

Therefore, variable selection is needed.

We first use Ridge regression. The result is shown as Fig. 4. From the result, we can see that the error bar range is Large and the minimum RMSE is at $\log \lambda = -1.65$. Obviously, the ridge regression makes substantial improvements in improving test RMSE compared to the linear regression model which mainly because of the penalty.

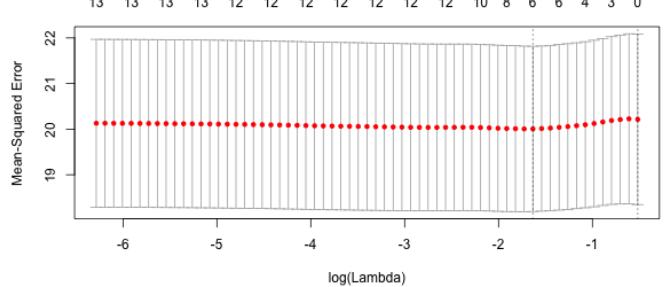


Fig. 4. Total Number of Coefficient and RMSE in Ridge¹

C. LASSO with Cross Validation

LASSO is a regression analysis method which performs both variable selection and regularization in order to improve prediction accuracy and interpretability of the statistical model it produces [3]. By LASSO we can greatly shrink the number of variable and balance the accuracy of the model. However, there is a parameter λ needed to tune, in order to make the objectivity of choosing λ , here we use 10th Cross Validation [4].

¹Red point is the mean RMSE of 10 Cross Validations, grey error bar is the uppermost and lowermost RMSE of 10 fold Cross Validation, the two dashed lines are value of lambda that gives minimum RMSE and largest value of lambda such that error is within 1 standard error of the minimum. The following similar figures' elements are the same as this figure.

TABLE IV
COEFFICIENT OF RIDGE REGRESSION

	s0
acousticness	-0.06
danceability	1.86
duration_ms	-0.00
energy	-0.73
instrumentalness	-0.87
key	-0.02
liveness	0.54
loudness	-0.03
mode	-0.01
speechiness	0.35
tempo	-0.01
time_signature	0.27
valence	0.10

The result is shown as Fig. 5. From the result, we can see that the error bar range is Large and the minimum RMSE is at $\log \lambda = -1.65$. Obviously, the LASSO makes substantial improvements in improving test RMSE compared to the linear regression model which mainly because of the penalty. Besides, LASSO can also make some coefficient to 0. From the coefficient TABLE V we can see the following imformation:

- 1) The higher danceability, the higher popularity;
- 2) The lower energy, the higher popularity.

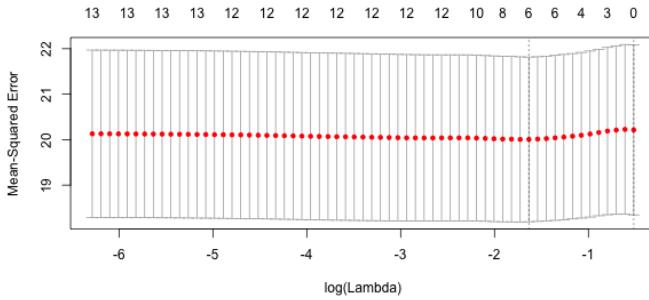


Fig. 5. Total Number of Coefficient and RMSE in Lasso

D. Elastic-net with Cross Validation

Elastic Net first emerged as a result of critique on lasso, whose variable selection can be too dependent on data and thus unstable. The solution is to combine the penalties of ridge regression and lasso to get the best of both worlds. Elastic Net aims at minimizing the following loss function:

$$\frac{\sum_{i=1}^n (y_i - x_i' \beta)^2}{2n} + \lambda \left(\frac{1-\alpha}{2} \sum_{j=1}^m \hat{\beta}_j^2 + \alpha \sum_{j=1}^m |\hat{\beta}_j| \right) \quad (1)$$

where α is the mixing parameter between ridge ($\alpha = 0$) and lasso ($\alpha = 1$).

The result is shown as Fig. 6. From the result, we can see that the best Elastic-net is similar to LASSO model.

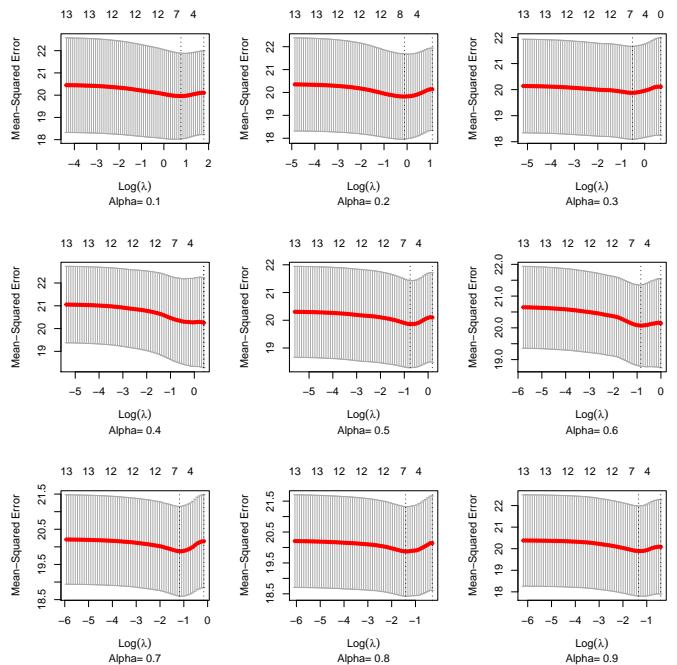


Fig. 6. Total Number of Coefficient and RMSE in Elastic-net

TABLE V
COEFFICIENT OF LASSO REGRESSION

	s0
acousticness	0.00
danceability	2.90
duration_ms	-0.00
energy	-0.58
instrumentalness	0.00
key	-0.00
liveness	0.00
loudness	0.00
mode	0.00
speechiness	0.00
tempo	-0.02
time_signature	0.00
valence	0.00

TABLE VI
COEFFICIENT OF ELASTIC-NET REGRESSION

	Coefficient
acousticness	0.00
danceability	2.34
duration_ms	-0.00
energy	-0.25
instrumentalness	0.00
key	0.00
liveness	0.00
loudness	0.00
mode	0.00
speechiness	0.00
tempo	-0.01
time_signature	0.00
valence	0.00

V. CLASSIFICATIVE MODEL ANALYSIS

In this chapter, we use multinomial logistic regression, linear discriminative analysis, quadratic discriminative analysis, tree method, random forest and k-nearest neighbors [5] to classify our data. All model result is trained by 10-fold cross validation within the train data. Before classification, we divide the popularity data to 3 classes as TABLE VII.

TABLE VII
DIVID THE DATA TO 3 CLASSES

	Unpopular	Normal	Popular
Popularity	0-40	40-70	70-100
Train Count	26621	63406	1437

In this table, we noticed that our data is highly imbalanced, if we use this data directly, our classification model might be highly biased towards the category which has the highest samples. Hence, we balanced the data, the way we use to balance the data subsample.

TABLE VIII
BALANCED DATA

	Unpopular	Normal	Popular
Popularity	0-40	40-70	70-100
Train Count	1437	1437	1437

A. Multinomial Logistic Regression

Logistic regression is commonly used in classic classification analysis, since we have multiple class in the data, multinomial logistic regression is introduced here.

The error analysis plot is depict as Fig.7. From error analysis plot we can see that the picked model's error rate is 0.54. The popular song is well predicted, however, normal song and unpopular song are bad predicted.

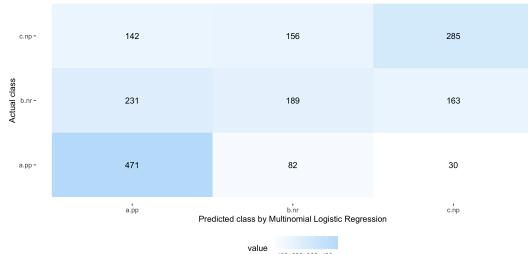


Fig. 7. Error Analysis² Of Multinomial Logistic Regression

B. Linear Discriminative Analysis (LDA)

After modeling LDA, we depicted and error analysis plot as Fig.8. From error analysis plot we can see that the picked model's error rate is 0.516 which is almost the same as multinomial logistic regression. The prediction is almost no

²The black line in the left figure is the average line of 10 accuracies in cross validation, and the overall accuracy in the right figure is from the best model picked from cross validation. All following similar figure are the same.

improvement, which means that both multinomial logistic regression and LDA can't greatly classify the data.

As we all know, these two way generates a linear boundary of classification, however, when the real boundary is non-linear these 2 ways can't produce a great result. So we next try some non-linear boundary models.

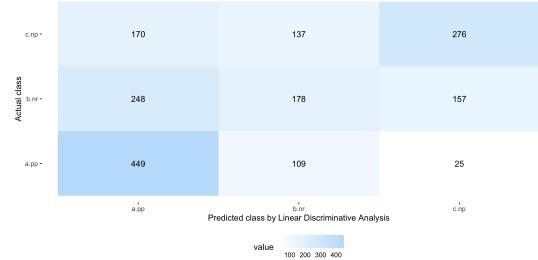


Fig. 8. Error Analysis Of Linear Discriminative Analysis

C. Quadratic Discriminative Analysis (QDA)

In order to model a non-linear boundary, naturally, QDA comes to our mind. Since QDA assumes a quadratic decision boundary, it can accurately model a wider range of problems than the linear methods.

We train the data with QDA model, and depicted error analysis plot as Fig.9. From error analysis plot we can see that the picked model's error rate is 0.516. The overall prediction is the same as LDA which means that QDA is no better than previous ways.

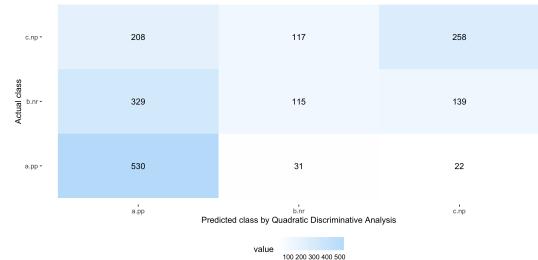


Fig. 9. Error Analysis Of Quadratic Discriminative Analysis

D. Decision Tree

Besides, we train the data with Decision Tree model. Then we depict error analysis plot as Fig.10. From error analysis plot we can see that the picked model's error rate is 0.509 and it only predicted popular class greatly. Obviously, Decision Tree have no significant improvement than previous ways.

E. K-Nearest Neighbors (KNN)

What if we use an even more flexible method than QDA? Will the result be better? Hence, KNN is introduced as our last classification method. We train the data with KNN model, get the best $k = 2$ by cross validation. Then we depict error analysis plot as Fig.11. From error analysis plot we can see that the picked model's error is 0.391, and error rate on popular



Fig. 10. Error Analysis Of Decision Tree Analysis

songs is 0.3953. The error rate has been decreased by 10% than prior methods, however, the accuracy is still very low.

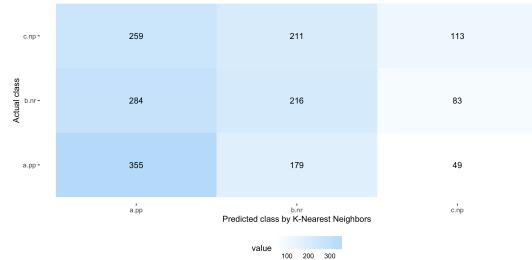


Fig. 11. Error Analysis Of 2-Nearest Neighbors Analysis

VI. CONCLUSION

Going into this endeavour, we were uncertain if it is even possible to predict, better than random, if a song will be popular or not. After testing our model on new songs pulling from Spotify, we observe that it is significantly simpler to correctly predict a bad song rather than a hit. It may have been easier to predict non hit songs because our data was skewed, with limited hit songs. Besides, We didn't consider lyrics, to predict a song's popularity.

From the result, we can see that all regression model are not so good. The error bar range in Fig. 4, Fig. 5 and Fig. 6 is large.

VII. FUTURE WORKS

Moving forward, we would like to explore how additional features such as lyrics, artist popularity, artist location or release date can influence a song's popularity. In addition, we may consider using the full dataset to see if we can improve our models. Another alternative is to use Spotify API to collect our own data.

REFERENCES

- [1] Tomigelo, "Spotify audio features," Dec 2018. [Online]. Available: <https://www.kaggle.com/tomigelo/spotify-audio-features/version/1>
- [2] "Get audio features for a track." [Online]. Available: <https://developer.spotify.com/documentation/web-api/reference/tracks/get-audio-features/>
- [3] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [4] R. Kohavi *et al.*, "A study of cross-validation and bootstrap for accuracy estimation and model selection," in *Ijcai*, vol. 14, no. 2. Montreal, Canada, 1995, pp. 1137–1145.
- [5] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An introduction to statistical learning*. Springer, 2013, vol. 112.

APPENDIX A R SOURCE CODE

```

1 library(glmnet)
2 library(ggplot2)
3 library(Metrics)
4 library(FNN)
5 library(caret)
6 library(gridExtra)
7 library(xtable)
8
9 rm(list = ls())
10 set.seed(5665)
11
12 error_analysis <- function(pred,actual, plot_name) {
13   # pred = predict(model, DT.test)
14   error = ifelse(pred == actual, 1, 0)
15   table.error = table(data.frame(pred, actual))
16   table.error.rate = table.error
17   errorrate = sum(error) / length(error)
18   t = reshape2::melt(table.error.rate)
19   ep = ggplot(t, aes(t[, 1], t[, 2], fill = value, label = round(value, 3))) +
20     # x and y axes => Var1 and Var2
21     geom_tile() + # background colours are mapped according to the
22     # value column
23     geom_text() +
24     scale_fill_continuous(high = "#0e1fa", low = "#fffff") +
25     theme(legend.position = "bottom",
26           panel.background = element_rect(fill = "white")) +
27     scale_x_discrete(label = abbreviate) + scale_y_discrete(label =
28       abbreviate) +
29     xlab(paste("Predicted class by", plot_name)) + ylab("Actual class")
30   # +ggtitle(paste("Overall Accuracy=", round(model$results$Accuracy,
31   # +4)))
32   # cvep = ggplot(model$resample, mapping = aes(Resample, Accuracy)
33   #   ) +
34   # geom_point() +
35   # geom_hline(yintercept = mean(model$resample$Accuracy))
36   summa = list()
37   summa[["error"]] = error
38   summa[["table.error"]] = table.error
39   summa[["table.error.rate"]] = table.error.rate
40   summa[["errorrate"]] = errorrate
41   ggsave(
42     file = paste("error", gsub(" ", "", plot_name), ".pdf", sep = ''),
43     path = "Thesis/figure",
44     # grid.arrange(cvep, ep, ncol = 2),
45     ep,
46     width = 10,
47     height = 5
48   )
49   return(summa)
50 }
51
52
53 # r.dt=data.table::fread('/Users/Trent/Downloads/track_features/tf_
54 #   ↪ 000000000000.csv', header = T, sep = ',')
55 # select_index=sample(nrow(r.dt),100000)
56 # dt = as.data.frame(r.dt[select_index,])
57 # rownames(dt.train) <- NULL
58
59 # dt=read.csv("data/track_features/tf_mini.csv")
60 # dt=dt[2:30]
61 # index.pop=which(dt[,3]>98)
62 # dt=dt[index.pop,]
63 # index.year=which(dt[,2]==2018)
64 # dt=dt[index.year,c(1,3:29)]
65
66 t = read.csv("data/SpotifyAudioFeaturesApril2019.csv")
67 data.nnum = t[, 1:3]
68 data = t[, 4:17]
69
70 train = sample(nrow(data), 0.7 * nrow(data))
71 data.train = data[train,]

```

```

data.test = data[-train,]

ggplot(data, aes(popularity)) + geom_histogram(bins=40)

library(GGally)
temp=ggpairs(data[sample(nrow(data), 0.1 * nrow(data)), ], lower = list(
  ↪ continuous = wrap(
  "points", alpha = 0.3, size = 0.1
)))
ggsave(file = 'ggpairs.png',
       path = "Thesis/figure",
       temp,
       width = 20,
       height = 10
)

#regression of popular type
dt.train = data.train[which(data.train$popularity>79,)]
dt.test = data.test[which(data.test$popularity>79,)]

lm1 = lm(popularity ~ ., data = dt.train)
summary(lm1)
x = dt.train$popularity
y = lm1[["fitted.values"]]
predict_and_observe_plot_data = data.frame(x, y)
p = ggplot(data = predict_and_observe_plot_data, aes(x, y)) +
  geom_point() +
  geom_abline(
    intercept = 0,
    slope = 1,
    color = "red",
    linetype = "dashed",
    size = 1
  ) +
  xlab("True Popularity") +
  ylab("Predicted Popularity")
ggsave(
  file = 'predict_and_observe_plot.pdf',
  path = "Thesis/figure",
  gridExtra::arrangeGrob(p),
  width = 10,
  height = 5
)
rmse.lm.train = rmse(dt.train$popularity, lm1[["fitted.values"]])
error.lm.train = dt.train$popularity - lm1[["fitted.values"]]

pred.lm = predict(lm1, dt.test[, -14])
x = dt.test$popularity
y = pred.lm
predict_and_observe_plot_data = data.frame(x, y)
p = ggplot(data = predict_and_observe_plot_data, aes(x, y)) +
  geom_point() +
  geom_abline(
    intercept = 0,
    slope = 1,
    color = "red",
    linetype = "dashed",
    size = 1
  ) +
  xlab("Observed Popularity") +
  ylab("Predicted Popularity")
ggsave(
  file = 'testerror.pdf',
  path = "Thesis/figure",
  gridExtra::arrangeGrob(p),
  width = 10,
  height = 5
)
rmse.lm.test = rmse(dt.test$popularity, pred.lm)
error.lm.test = dt.test$popularity - pred.lm

```

```

139 x = as.matrix(dt.train[,-14])
140 y = as.matrix(dt.train$popularity)
141 x.t = as.matrix(dt.test[,-14])
142 y.t = as.matrix(dt.test$popularity)
143
144 cvfit.ridge = cv.glmnet(x, y, alpha = 0)
145 pdf('Thesis/figure/cvfitridge.pdf',width = 600, height = 300)
146 plot(cvfit.ridge)
147 dev.off()
148 xtable(as.array(glmnet(x, y, alpha = 0, lambda = cvfit.ridge$lambda.min)
149   ↪ $beta))
150
151 cvfit.lasso = cv.glmnet(x, y, alpha = 1)
152 pdf('Thesis/figure/cvfitlasso.pdf',width = 600, height = 300)
153 plot(cvfit.lasso)
154 dev.off()
155 xtable(as.array(glmnet(x, y, alpha = 1, lambda = cvfit.lasso$lambda.min)
156   ↪ $beta))
157
158 pdf('Thesis/figure/cvfitpla_net.pdf')
159 par(mfrow=c(3,3))
160 pla_net=function(k){
161   cvfit.net = cv.glmnet(x, y, alpha=k)
162   plot(cvfit.net, sub=paste("Alpha=",k))
163 }
164 for (k in 0.1*1:9) {
165   pla_net(k)
166 }
167 dev.off()
168 xtable(as.array(glmnet(x, y, alpha = 0.5, lambda = 0.5)$beta))
169
170 f_knn_tst_rmse = function(k) {
171   rmse(y.t, knn.reg(x, x.t, y, k)$pred)
172 }
173 f_knn_trn_rmse = function(k) {
174   rmse(y, knn.reg(x, x, y, k)$pred)
175 }
176 # k = c(1, 5, 10, 25, 50, 250)
177 k = c(1, 3:10)
178 knn_tst_rmse = sapply(k, f_knn_tst_rmse)
179 knn_trn_rmse = sapply(k, f_knn_trn_rmse)
180 # determine "best" k
181 best_k = k[which.min(knn_tst_rmse)]
182 knn_results = data.frame(k,
183   round(knn_trn_rmse, 2),
184   round(knn_tst_rmse, 2))
185 colnames(knn_results) = c("K", "Train RMSE", "Test RMSE")
186
187 pred.best.knn = knn.reg(x, x.t, y, 25)$pred
188 error.best.knn = y.t - pred.best.knn
189 plot(error.best.knn)
190
191
192 ######classification
193 DT.train = data.train
194 DT.test = data.test
195 DT.train$popularity = as.factor(ifelse(
196   DT.train$popularity < 10,
197   "c.unpopular",
198   ifelse(DT.train$popularity < 70, "b.normal", "a.popular")
199 ))
200
201 DT.test$popularity = as.factor(ifelse(
202   DT.test$popularity < 10,
203   "c.unpopular",
204   ifelse(
205     DT.test$popularity < 40,
206     "b.normal",
207     ifelse(DT.test$popularity < 70, "b.normal", "a.popular")
208   )))
209
210
211   "b.normal",
212   ifelse(DT.test$popularity < 70, "b.normal", "a.popular")
213 )
214 })
215 #c.unpopular
216 table(DT.test$popularity)
217 table(DT.train$popularity)
218 # data is not balanced, we first balance our train data
219 DT.test=caret::downSample(x=DT.test[, -ncol(DT.test)],y=DT.test$popularity,yname="popularity")
220 DT.train=caret::downSample(x=DT.train[, -ncol(DT.train)],y=DT.train$popularity,yname="popularity")
221
222 table(DT.test$popularity)
223 table(DT.train$popularity)
224 #####Logistic Regression
225 ctrl = trainControl(method = "cv", 3)
226 model.multinom = train(
227   popularity ~.,
228   data = DT.train,
229   method = "multinom",
230   trControl = ctrl,
231   preProcess = c("center", "scale"))
232 )
233 pred.multinom = predict(model.multinom, DT.test, type = "prob")
234 error.multinom = error_analysis(predict(model.multinom, DT.test),DT.test[, 14], plot_name = "Multinomial Logistic Regression")
235 error.multinom
236 #####LDA
237 ctrl = trainControl(method = "cv", 10)
238 model lda = train(
239   popularity ~.,
240   data = DT.train,
241   method = "lda",
242   trControl = ctrl,
243   preProcess = c("center", "scale"))
244 )
245 pred.lda = predict(model.lda, DT.test, type = "prob")
246 error.lda = error_analysis(predict(model.lda, DT.test),DT.test[, 14], plot_name = "Linear Discriminative Analysis")
247 error.lda
248 #####QDA
249 ctrl = trainControl(method = "cv", 10)
250 model.qda = train(
251   popularity ~.,
252   data = DT.train,
253   method = "qda",
254   trControl = ctrl,
255   preProcess = c("center", "scale"))
256 )
257 pred.qda = predict(model.qda, DT.test, type = "prob")
258 error.qda = error_analysis(predict(model.qda, DT.test),DT.test[, 14], plot_name = "Quadratic Discriminative Analysis")
259 error.qda
260 #####KNN
261 pred.knn = knn(DT.train[, 1:13], DT.test[, 1:13], DT.train[, 14], k = 2,
262   prob = TRUE)
263 error.knn = error_analysis(as.vector(pred.knn),DT.test[, 14], plot_name = "K-Nearest Neighbors")
264 error.knn
265 error.knn$table.error
266
267 f_knn = function(k) {
268   error.knn = error_analysis(as.vector(knn(DT.train[, 1:13], DT.test[, 1:13], DT.train[, 14], k, prob = TRUE)),DT.test[, 14], plot_name = "K-Nearest Neighbors")
269   print(error.knn$table.error)
270 }
271 k = c(1:10)
272 sapply(k, f_knn)
273
274 View(t[row.names(DT.test)[which(as.vector(pred.knn) != DT.test$popularity & DT.test$popularity == "c.unpopular")],c(1,3,17)])

```

```
275 ↵ ]}  
276 #####Dicision tree  
277 model.tree = rpart::rpart(popularity ~ ., data = DT.train)  
278 pred.tree = predict(model.tree, DT.test[, 1:13], type = "class")  
279 error.tree = error_analysis(as.vector(pred.tree),DT.test[, 14], plot_name =  
    ↵ "Dicision Tree")  
280 error.tree  
281 #####Random forest  
282 model.forest= randomForest::randomForest(popularity ~ ., data=DT.train)  
283 pred.forest= predict(model.forest, DT.test[,1:13], type = "class")  
284 error.forest=error_analysis(as.vector(pred.tree),DT.test[, 14],plot_name="  
    ↵ Random Forest")  
285 error.forest
```

..../Code/Final project.R