

Final Report

Team Name: GPA 4.0

Member: Violet(Bingcui) Guo, Jonathan(Yiran) Liu

I. Introduction

A video's early growth pattern can indicate the eventual success of the video and also the overall health of the channel. The goal of the project was to predict the percentage change in views on a video between the second and sixth hour since its publishing with features of thumbnail image, video title and some other features like the duration of the video.

II. Methodology

A. Data Preprocessing

After we loaded the data into R, we first checked for missing values. Since there were no missing values found in the dataset, we didn't perform any missing values handling techniques. Carefully reviewing the Feature_Descriptions document, we found "PublishedDate" particularly valuable for the rich information this predictor provides because from common sense, the published weekday, month, and the time of a day that the video is published would affect the number of views on the video. With a format of Month/Day/Year Hour:Minute, we decided to dissect the "PublishedDate" into five distinct predictors with the help of function `as.Date` from the library(lubridate) (We've asked for permission to use the package via email) (Vitalie, 2020). The five predictors are, "years", "months", "dates", "weekdays" and "myduration". While the first four are self-explanatory, "myduration" is the amount of minutes past after 0am (So 1:00 corresponds to 60 and 23:30 corresponds to $60*23+30=1410$)

Removing highly correlated predictors is essential for later performing Random Forest to conduct feature selection. When there is a high correlation between two predictors, any of the two can be used as the predictor. But if one of them is used in the model, the importance of the other one is reduced since the impurity they can remove is removed by the selected predictor. This became a problem when we took the most important predictors to build the model.

To remove highly correlated predictors, we used the function `findCorrelation` from library(mlbench) (Friedrich and Evgenia, 2012) and function `vsetdiff` from library(vecsets)(Carl. W., 2018) (We've asked for permission to use the packages via email). Firstly, we extracted the numerical predictors in the training data sets, and then we used the `cor` function to get the correlation matrix of all the pairs. And then, I used `findCorrelation` function to find pairs with a correlation coefficient of over 0.7 and used the function `vsetdiff` to remove one predictor of each pair.

We decided to set the cutoff to 0.7 because cutoff was a parameter to decide what we deem as high enough correlation to toss out certain predictors. In other words, we would only keep predictors that have correlations below the cutoff value. We set the cutoff to 0.7 and that is because statisticians in general (Diana and Phoebe, 2014), consider a correlation of 0.7 as highly correlated. Furthermore, scaling the predictors was not a concern since we have decided to only use Bagging and Random Forest to build the model and tree-based models do not require feature scaling.

B. Feature Selection

Due to the large number of predictor variables in the given data set, the presence of non-informative variables can add uncertainty to the predictions and reduce the overall effectiveness of the model. Thus, feature selection became one of our most important tasks. We came up with two strategies for feature selection. Our first approach is Lasso. The reason that we preferred Lasso over Ridge is Lasso makes unimportant predictors to have zero coefficients, which resulting in an easier choice of identifying important predictors as compare to Ridge, which only shrinks the coefficients of unimportant predictors to a very small value. And our second approach is building a model with Random Forest in which all predictors from the dataset are included, and then we use the importance function to acquire the most important predictors. The reason why we used random forest to do feature selection is that the tree-based strategies used by random forests naturally rank by how well they improve the purity of the node. With the help of the histogram of the frequency of feature importance plot below(Figure 1), we know that we want the predictors that are of high importance and we don't want that many predictors. So from the plot, we could see that the predictors that are of value are in the tail of the distribution. So we dropped the predictors in the first 2 rectangles and used the predictors in the tail by setting the benchmark of 8 and we got 30 predictors. We ran a base default random forest with these 30 predictors by using 70% training data and got a RMSE result on the remaining 30% data of 1.47996. And if we increase the number of predictors to 31, 32, ...36, we get worse results of 1.48205, 1.48403, ... 1.50005. So with more predictors, the model gets worse in terms of interpretability. So we chose 30 predictors as the predictors for the final model.

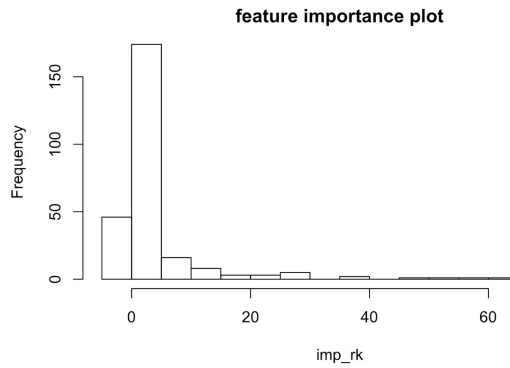


Figure 1

With the above two approaches, we had one group of features from Lasso and another group of features from Random Forest. We decided to use these two groups of features to build two models that depend on them respectively. And then, we have decided to compare oob for the . The details will be explained in the section of Model Selection Process.

C.Model Selection Process

We had separated the dataset from the training.csv into our own training data (70 percent) and testing data (30 percent) at the very beginning. In this way, we could build models from our training data solely to evaluate their performances and compare them with each other. This evaluation was based on a RMSE calculating function we wrote consulting the formula given on Kaggle. By doing so, not only we were able to get a general sense of our models' performance prior to submitting on Kaggle, but also we could choose the best model. Since the models we built were only based on 70 percent of the data from the original dataset in training.csv, after selecting our best model, we then used the whole dataset to build the same model, which in return, did give us a better score after submitting on Kaggle. It only took 7 submissions to get a score of 1.37665 on the public board and we credited this efficiency to our evaluation approach. We have two groups of predictors from Feature Selection; the predictors selected by Lasso and the predictors selected by Random Forest Importance function. We chose the best group by building two different Random Forest models and then compare the Out-of-bag (OOB) error. The following was the result (Figure 2 and Figure 3). Why we build random forests as a base model is because it utilizes an ensemble method which adds more diversity to the trees and works well in high-dimensional data. And we used OOB error allows us to estimate the prediction performance improvement by evaluating predictions on those observations which were not used in the building of the next base learner.

Model built with predictors selected by Lasso

	RMSE <dbl>	Rsquared <dbl>	mtry <dbl>
1	1.920184	0.4723282	2
2	1.633841	0.6179697	52
3	1.646280	0.6121305	103

Figure 2

Model built with predictors selected by Random Forest Importance

	RMSE <dbl>	Rsquared <dbl>	mtry <dbl>
1	1.559503	0.6480864	2
2	1.494743	0.6767067	16
3	1.486695	0.6801788	30

Figure 3

So as we can see, the models built with predictors selected by Random Forest Importance gave a better RMSE. So we had decided to use the predictors selected by Random Forest Importance. In order to achieve the best Random Forest performance possible, the tuning of hyperparameters, namely, "mtry" (Number of variables randomly sampled as candidates at each split) and "ntree" (Number of trees to grow), was necessary. With the plot of squared error vs ntree below(Figure4 and Figure5), it was obvious that the error (square RMSE) was the lowest when "ntree"=1100. Thus, we set the parameter "ntree" to be 1100 in our final model. To verify this decision, we had tried to build Random Forest models with "ntree"'s equal to "1100", "1500" and "2000". And evaluating the model on the test dataset, we get the results of 1.459922, 1.461391, and 1.461919 respectively. The reason that these test MSE's were lower than our public score was that we only used our own training data to train the model, which only contained 70 percent of the whole dataset in training.csv.

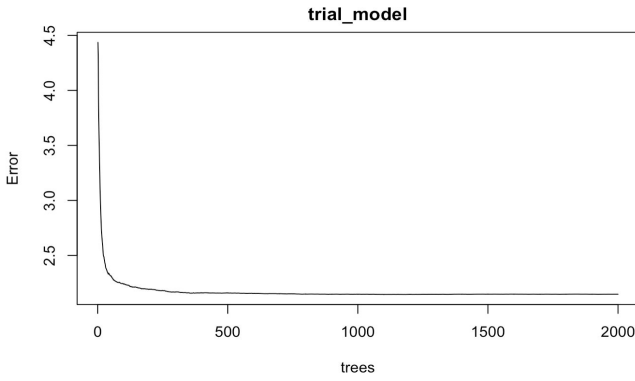


Figure4

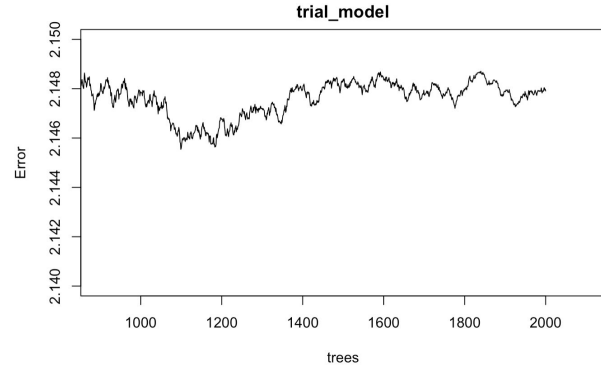


Figure5

From the following plot of Out of bag RMSE vs number of predictors (Figure6), we can see when mtry is 30, we get the lowest MSE.

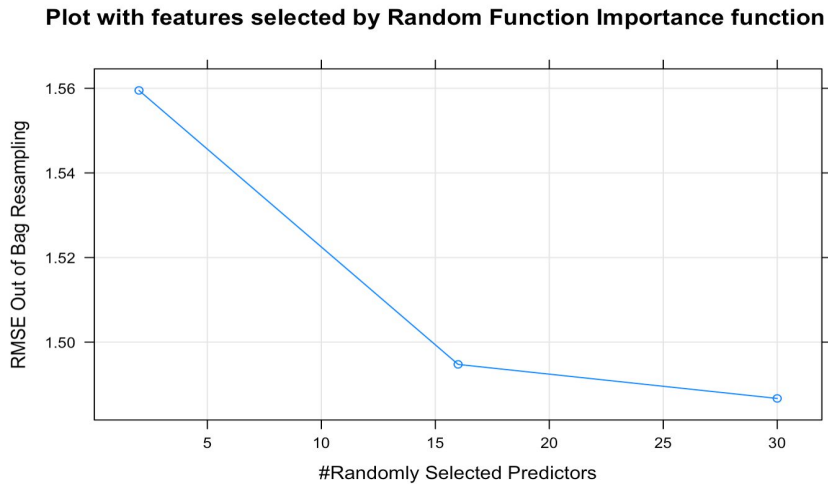


Figure6

Besides Random forest, we have also used Ridge Regression to build a model because ridge allows us to “shrink” unimportant coefficients, which makes it work well on new datasets and use complex models while avoiding overfitting at the same time. By using cross-validation, we found the best lamda. The model built by Ridge Regression with the best lamda we found did not give us a good performance at all. The RMSE is 1.70454, which is less than the RMSE we received from using Random Forest. After we had made the decision of using the Bagging model with “mtry” equals to 30 and “ntree” equals to 1100 as our final model. We trained the model with the whole dataset from training.csv and used this model to predict. We got 1.45467 by using RMSE on our laptop by using 70% data and got a RMSE of 1.37665 by using all the data on Kaggle public board and 1.39300 on private board.

III. Result

Our final model is built by Bagging with “mtry” equals to 30 and “ntree” equals to 1100 and we got a RMSE of 1.37665 on the Kaggle public board, which is a very decent score.

IV. Conclusion

We received a score of 1.37665 on the Kaggle public board and 1.39300 on the Kaggle private board, which is very similar. Thus, our model has a good interpretability. We credit this good performance to our tactic of using creative preprocessing method, creating our own testing dataset, so we can evaluate our models performance prior to submission and by using Out-of-bag error and Cross-Validation, we improved our model and avoided over-fitting. Ideas on further improving our model includes tuning more parameters such as maxnodes and nodesize and using some more advanced feature selection methods not seen in class, such as recursive feature elimination. We hope to learn more and improve our data science skills in the future.

V. Statement of Contribution

Jonathan: Data Preprocessing, Data Modeling, Report Writing

Violet: Data Preprocessing, Feature Selection, Data Modeling, Report Writing

VI. Reference

Carl, Wi., 2018. *Package 'vecsets'*. Retrieved 2018

(<https://cran.r-project.org/web/packages/vecsets/vecsets.pdf>)

Diana, M. and Phoebe B. 2014. *Scatterplots and Correlation*. Retrieved 2014

(https://www.westga.edu/academics/research/vrc/assets/docs/scatterplots_and_correlation_notes.pdf).

Evgenia, D and Friedrich, L. 2012. *Package 'mlbench'*. Retrieved 2012

(<https://cran.r-project.org/web/packages/mlbench/mlbench.pdf>)

Vitalie, S., 2020. *Package 'lubridate'*. Retrieved 2020

(<https://cran.r-project.org/web/packages/lubridate/lubridate.pdf>).

VII. R Appendix

```
``{r}
##read in data
data <- read.csv("training.csv")
```

```
``
```

```
``{r}
##load library
library(lubridate)
library(caret)
```

```
``
```

```
``{r}
dim(data)
```

```
``
```

```

25 ▾ ```{r}
26 ## data transformation
27 years1 <- c()
28 months1 <- c()
29 days1 <- c()
30 hours1 <- c()
31 minutes1 <- c()
32 weekday1<- c()
33 ▾ for (i in 1:7242){
34   mydate <- strptime(data[, "PublishedDate"][i] , format = "%m/%d/%Y%H:%M")
35   years1 <- c(years1, year(mydate))
36   months1 <- c(months1, month(mydate))
37   days1 <- c(days1, day(mydate))
38   hours1 <- c(hours1, hour(mydate))
39   minutes1 <- c(minutes1, minute(mydate))
40   weekday1 <- c(weekday1, weekdays(mydate))
41 ▴ }
42 ▴ ```
43
44
45 ▾ ```{r}
46 ##remove id and publshied date
47 data$myduration <- 60*hours1 + minutes1
48 data <- data[, -c(1,2)]
49 ▴ ```
50
51 ▾ ```{r}
52 ##add transformed columns to the data
53 data$years <- years1
54 data$months <- months1
55 data$dates <- days1
56 data$weekdays <- weekday1
57 data$weekdays <- as.factor(data$weekdays)
58 ▴ ```

```

```

61 `r}
62 ## train_test split
63 size <- dim(data)[1]
64 train_idx <- sample(1:size, floor(size*0.7))
65 train <- data[train_idx,]
66 test <- data[-train_idx,]
67
68
69 `r}
70
71 `r}
72 ##remove highly correlated predictors
73 # Splitting numeric and categorical variables
74
75 numeric_idx <-c(3:247, 261:262)
76 categorical_idx <- c(248:259)
77 total_colnames <- colnames(data)
78
79 set.seed(1)
80 library(mlbench)
81 library(caret)
82 library(vecsets)
83 # calculate correlation matrix
84 correlationMatrix <- cor(train[,numeric_idx],use="complete.obs")
85 is.na(train)
86
87 train<-train[,colSums(is.na(train))==0]
88
89
90 #remove columns with std=0
91 predictors <- names(which(is.na(correlationMatrix[1,])))
92 predictors # all numeric
93 idx <- match(predictors,total_colnames)
94 numeric_idx2 <- numeric_idx[!numeric_idx %in% idx]
95 numeric_idx2
96
97

```



```

100
101 # find attributes that are highly corrected
102 highlyCorrelated <- findCorrelation(correlationMatrix, cutoff=0.7)
103 numeric_highly_cor <- vsetdiff(numeric_idx2, highlyCorrelated)
104 length(numeric_highly_cor)
105 numeric_highly_cor
106
107 ##name of all the highly correlated data
108 total_colnames[numeric_highly_cor]
109 ```
110
111
112 ```{r}
113 ## lasso feature selection
114 lasso.mod <- glmnet(train.mat, train$growth_2_6, family = "gaussian", alpha = 1,
115                    lambda=grid,
116                    standardize = TRUE)
117 ```
118
119 ```{r}
120 ## cv lasso
121 cv.lasso <- cv.glmnet(train.mat, train$growth_2_6, alpha=1, lambda=grid,
122                      standardize=TRUE, nfolds=10)
123 ```
124 ```{r}
125 cv.lasso$lambda.min
126 ```
127
128 ```{r}
129 ## see important coefficients
130 d <- predict(lasso.mod, s=cv.lasso$lambda.min, type="coefficients")
131 ```
132

```

```

133 ▾ ```{r}
134 ##see important features
135 sig <- dimnames(d)[[1]][which(d!=0)]
136 sig1 <- sig[-1]
137 sig
138 ▲ ```
139
140 ▾ ```{r}
141 ## look at lasso prediction error by using linear regression
142 ## and we know that the number of predictors that lasso selects out are too much
143 ## and the lasso prediction is so abd
144 predict(lasso.mod, s = cv.lasso$lambda.min, type="coefficients")
145 testx <- model.matrix(growth_2_6~.,test)
146 testy <- test$growth_2_6
147 lasso.pred = predict(lasso.mod, newx= testx, s=cv.lasso$lambda.1se)
148 lasso.err = mean((testy - lasso.pred)^2)
149 ▲ ```
150
151 ▾ ```{r}
152 lasso.err
153 ▲ ```
154
155 ▾ ```{r}
156 ##select the predictors after removing all the highly correlated data
157 data <- data[, -c(total_colnames[numeric_idxe])]
158 ▲ ```
159
160 ▾ ```{r}
161 ## do train_test split again by using the new data
162 size <- dim(data)[1]
163 train_idx <- sample(1:size, floor(size*0.7))
164 train <- data[train_idx,]
165 test <- data[-train_idx,]
166 ▲ ```
167

```



```

168
169 ```{r}
170 ##feature selection rf
171 rf_feature_select <- randomForest(growth_2_6~., data=train, mtry=ncol(train)-1,
172 n.trees=1000, importance=T)
173
174 ```{r}
175 ## look at importance features
176 imp_rk <- sort(importance(rf_feature_select)[, 1],decreasing=TRUE)
177 ```
178
179 ```{r}
180 ##look at the most importance features
181 varImpPlot(rf_feature_select, pch=19, cex=0.6)
182 ```
183
184 ```{r}
185 ## generate variable importance histogram plot to see the distribution
186 ## of the importance score and the frequency the variables with that score appear
187 var_index <- importance(rf_all)[,1]
188 hist(var_index,density = FALSE,main = "feature importance plot",xlab = "feature
189 importance value", xlim = c(-50,150),breaks = 20)
190
191
192 ```{r}
193 ## find the optimal number of trees
194 library(randomForest)
195 trial_model <- randomForest(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid_
196 high+avg_growth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_Su
197 bscribers_Base_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growth
198 _mid_high+cnn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Durati
199 on+num_words+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mean
200 _green+hog_342, data=train, ntree=2000, mtry=30)
201 plot(trial_model)
202 plot(trial_model,ylim=c(2.14,2.15),xlim=c(900,2100))

```

```

192 ~~~{r}
193 ## test out the optimal number of predictors|
194 trial_model1 <- randomForest(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid
_high+avg_growth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_S
ubscribers_Base_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growt
h_mid_high+cnn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Durat
ion+num_words+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mea
n_green+hog_342+hog_641, data=train, ntree=2000, mtry=31)
195 trial_model1$results
196
197 trial_model2 <- randomForest(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid
_high+avg_growth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_S
ubscribers_Base_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growt
h_mid_high+cnn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Durat
ion+num_words+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mea
n_green+hog_342+hog_641+cnn_30, data=train, ntree=2000, mtry=32)
198 trial_model2$results
199
200
201 trial_model3 <- randomForest(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid
_high+avg_growth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_S
ubscribers_Base_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growt
h_mid_high+cnn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Durat
ion+num_words+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mea
n_green+hog_342+hog_641+cnn_30+week, data=train, ntree=2000, mtry=33)
202 trial_model3$results
203
204 trial_model4 <- randomForest(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid
_high+avg_growth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_S
ubscribers_Base_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growt
h_mid_high+cnn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Durat
ion+num_words+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mea
n_green+hog_342+hog_641+cnn_30+week+year, data=train, ntree=2000, mtry=34)
205 trial_model4$results
206
207 trial_model5 <- randomForest(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid
_high+avg_growth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_S
ubscribers_Base_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growt
h_mid_high+cnn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Durat
ion+num_words+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mea
n_green+hog_342+hog_641+cnn_30+week+year, data=train, ntree=2000, mtry=35)
208 trial_model5$results
209
210 trial_model6 <- randomForest(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid
_high+avg_growth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_S
ubscribers_Base_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growt
h_mid_high+cnn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Durat
ion+num_words+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mea
n_green+hog_342+hog_641+cnn_30+week+year, data=train, ntree=2000, mtry=36)
211 trial_model6$results
212 ~~~

```



```

199 train_control <- trainControl(
200     method = "oob",
201
202     savePredictions = TRUE
203 )
204
205 ## find OOB error to compare rf and bagging
206 treefit1<- train(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid_high+avg_gr
  owth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_Subscribers_B
  ase_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growth_mid_high+c
  nn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Duration+num_word
  s+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mean_green+hog_
  342,
207     data = train, method = 'rf',
208     trControl = train_control, ntree=2000)
209
210 plot(treefit1, main="Plot with features selected by Random Function Importance
  function")
211
212 treefit2<- train(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid_high+avg_gr
  owth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_Subscribers_B
  ase_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growth_mid_high+c
  nn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Duration+num_word
  s+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mean_green+hog_
  342,|
213     data = train, method = 'rf',
214     trControl = train_control, ntree=1100)
215
216 plot(treefit2, main="Plot with features selected by Random Function Importance
  function")
217
218 treefit3<- train(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid_high+avg_gr
  owth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_Subscribers_B
  ase_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growth_mid_high+c
  nn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Duration+num_word
  s+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mean_green+hog_
  342,
219     data = train, method = 'rf'.

```

```

219         data = train, method = 'rf',
220         trControl = train_control, ntree=1500)
221 ## plot the results using diffnet ntrees
222 par(mfrow=c(3,1))
223 plot(treefit1, main="Plot with features selected by Random Function Importance
224      function ntree=2000", xlab="mtry")
225 plot(treefit2, main="Plot with features selected by Random Function Importance
226      function ntree=1000", xlab="mtry")
227 plot(treefit3, main="Plot with features selected by Random Function Importance
228      function ntree=1500", xlab="mtry")
229 treefit1$results
230 treefit2$results
231 treefit3$results
232
233 treefit<- train(growth_2_6~cnn_10+cnn_17+cnn_86+cnn_89+cnn_25+cnn_12+views_2_hours+c
234 nn_68+punc_num_..21+cnn_88+punc_num_..28+Duration+num_uppercase_chars+num_words+hog_
235 454+num_digit_chars+num_chars+punc_num_..1+sd_blue+hog_453+mean_red+sd_green+num_upper
236 case_words+hog_859+hog_651+punc_num_..11+hog_359+doc2vec_4+hog_668+punc_num_..14+h
237 og_669+doc2vec_10+num_stopwords+hog_182+hog_295+hog_132+hog_177+hog_702+hog_40+hog_6
238 49+hog_738+hog_62+hog_716+doc2vec_2+doc2vec_8+hog_832+hog_657+hog_829+hog_791+hog_33
239 6+hog_94+punc_num_..15+hog_674+doc2vec_9+hog_350+doc2vec_11+hog_492+hog_452+pct_nonz
240 ero_pixels+hog_746+hog_60+hog_774+hog_844+hog_724+hog_640+doc2vec_15+hog_1+punc_num_
241 ..8+hog_341+hog_342+hog_78+hog_855+hog_125+doc2vec_3+hog_677+hog_782+hog_316+hog_849
242 +hog_705+hog_815+hog_378+hog_61+hog_195+hog_116+cnn_39+punc_num_..3+punc_num_..7+hog
243 _641+hog_676+hog_279+hog_523+hog_665+punc_num_..2+hog_133+hog_856+hog_655+max_red+pu
244 nc_num_..16+punc_num_..20+doc2vec_13+Num_Subscribers_Base_low+avg_growth_low_mid+cou
245 nt_vids_low,
246
247         data = train, method = 'rf',
248         trControl = train_control, ntree=2000)
249
250 plot(treefit, main="Plot with features selected by Lasso")
251 treefit$results
252
253 ```

```

```

241 `r}
242 ## function to get RMSE
243 RMSE <- function(y1,y2){
244   rmse <- sqrt(mean((y1-y2)^2))
245   return(rmse)
246 }
247 `r}
248
249 `r}
250 pred_test <- predict(trial_model, test)
251 pred_test1 <- predict(treefit1, test)
252 pred_test2 <- predict(treefit2, test)
253 pred_test3 <- predict(treefit3, test)
254 `r}
255
256 `r}
257 ##cross validation to see the result
258 RMSE(test$growth_2_6, pred_test2)
259 RMSE(test$growth_2_6, pred_test1)
260 RMSE(test$growth_2_6, pred_test2)
261 RMSE(test$growth_2_6, pred_test3)
262
263 `r}
264
265 `r}
266 ##read in test data
267 test_id <- read.csv("test.csv")
268 `r}
269
270 `r}
271 ## read in test data
272 test_data <- read.csv("test.csv")
273 `r}
274 `r}
275 ## transform test data
276 years3 <- c()
277 months3 <- c()
278

```



```

277 months3 <- c()
278 days3 <- c()
279 hours3 <- c()
280 minutes3 <- c()
281 weekday3<- c()
282 for (i in 1:3105){
283   mydate <- strptime(test_data[, "PublishedDate"][i] , format = "%m/%d/%Y%H:%M")
284   years3 <- c(years3, year(mydate))
285   months3 <- c(months3, month(mydate))
286   days3 <- c(days3, day(mydate))
287   hours3 <- c(hours3, hour(mydate))
288   minutes3 <- c(minutes3, minute(mydate))
289   weekday3 <- c(weekday3, weekdays(mydate))
290 }
291 ```
292
293 ```{r}
294 test_data$years <- years3
295 test_data$months <- months3
296 test_data$dates <- days3
297 test_data$weekdays <- weekday3
298 test_data$myduration <- hours3 * 60 + minutes3
299 ```
300
301 ```{r}
302 ## get our output
303 df1<-data.frame(test_id[,1],pred_labal2)
304 colnames(df1) <- c("id", "growth_2_6")
305
306 write.csv(df1,"submission_bagging_final2.csv", row.names = FALSE)
307
308
309 # ridge regression trial
310
311 library(glmnet)
312
313 i.exp <- seq(10, -2, length = 100)
314 grid <- 10^i.exp
315

```

69:1 (Top Level) ↕

R Markdown ↕


```

12
13 i.exp <- seq(10, -2, length = 100)
14 grid <- 10^i.exp
15
16 x <- model.matrix(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid_high+avg_g
rowth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_Subscribers_
_Base_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growth_mid_high+
cnn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Duration+num_wor
ds+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mean_green+hog
_342,train)[,-which(names(train) %in% "growth_2_6")]
17 y <- train$growth_2_6
18
19 x1 <- model.matrix(growth_2_6~avg_growth_low_mid+cnn_10+Num_Views_Base_mid_high+avg_
growth_low+cnn_86+cnn_89+cnn_12+cnn_17+Num_Subscribers_Base_mid_high+Num_Subscribers
_Base_low_mid+views_2_hours+cnn_25+myduration+count_vids_low_mid+avg_growth_mid_high
+cnn_88+cnn_68+count_vids_mid_high+punc_num_..28+num_uppercase_chars+Duration+num_wo
rds+cnn_19+num_digit_chars+punc_num_..21+num_chars+punc_num_..1+months+mean_green+ho
g_342,test)[,-which(names(test) %in% "growth_2_6")]
20
21
22 lasso.mod <- glmnet(x, y, family = "gaussian", alpha = 1,
23 lambda = grid, standardize = TRUE)
24
25
26 cv.output1 <- cv.glmnet(x, y, family = "gaussian", alpha = 1,
27 lambda = grid, standardize = TRUE,
28 nfolds = 10)
29
30 best.lambda.cv1 <- cv.output1$lambda.1se
31
32 Ridge_result<-predict(lasso.mod, s = best.lambda.cv1,
33 newx=x1)
34
35
36 RMSE(test$growth_2_6,Ridge_result)
37
38 #RMSE is 1.70454
39 ^
40

```