

STATS 101C Final Project

Limit DNE: Hana Yerin Lim, Yanhua Lin, Yingzhen Zhao

December 2020

1 Introduction

In this project, we aim to find the optimal factors to predict the percentage change in views of a video between the second and sixth hour since its publishing. The training data set includes 7242 videos, 258 predictors and 1 response variable called *growth_2.6*. There are 3105 videos in the test data set needed to predicted,

Our goal is to closely predict the true value of *growth_2.6* in the test data set. In order evaluate this problem, we need to minimize the root mean squared error (RMSE) metric measuring the distance between our predicted growth percentage and the true growth percentage, which is computed as follows:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (g_i - \hat{g}_i)^2}$$

2 Methodology

This section discusses data preprocessing and the results we obtained from establishing different models.

2.1 Preprocessing of the data

- Creating additional columns

From the feature description, there were several binary variables with different levels and one categorical variable called *PublishedDate*. We created several additional features regarding to these columns. Variable names and description are as follow:

1. *Num_Subscribers_Base_high*, *Num_Views_Base_high*, *avg_growth_high*, *count_vids_high*:

In addition to the existing levels of number of subscribers, number of total views, average growth from 2 to 6 hour, and the number of other videos that consisted of low, low mid, and mid high, we added additional high levels of corresponding data.

2. *Days*: We split the *PublishedDate* column, which contains the date and time of the video published on YouTube from April to September, into date and time columns respectively. We pulled the date and use 4/1/2020 to be the first day in order to observe how the growth rate changes day by day.
3. *Hours*: From the *PublishedDate* column, we pulled the time and convert it into hours in decimal form from midnight.
4. *midnight*, *morning*, *noon*, *afternoon*, *evening*: Divided the 24 hour a day into 5 chunks and made them into binary columns.

(For more details of the creation, please see the appendix on pg.1-2)

- Removing not important and highly correlated predictors

An important factor to consider is the existence of multicollinearity. This may impair the accuracy of the prediction and therefore needs to be avoided. First, we remove columns with zero variances, since all values within each of these columns are identical and they will not contribute to the prediction. Then, in order to remove highly correlated predictors, we created a while loop and we tested with different values of correlation thresholds on numeric variables and found out that 0.8 was the best threshold that kept enough amount of predictors. Although *views_2.hours* and *Duration* columns were removed after cleaning up the predictors, they are added back because they are considered to be important predictors from the description. (For more details, please see appendix on pg. 2)

After creating and cleaning the data set, we now have 7242 observations and 183 variables with 182 predictors and 1 response.

2.2 Description of the models

- Model and Predictors Selection

In order to compare the RMSE scores (calculated) from different models, we partitioned the data set into training and testing sets by 0.8/0.2 ratio. We fit LASSO regression, boosting, bagging and random forests to the training set using 182 predictors. Methods and parameters used for these models are described as follow:

LASSO: Our team started off by trying LASSO prediction because of its improvement on prediction accuracy and interpretability of the statistical model. We selected the best value for lambda using K-fold cross-validation using 10 folds. Despite obtaining the best predictor selections with the best value of λ of 0.010, the RMSE score only turned out to be 1.650821. (reference appendix on pg.4)

Boosting: We preform boosting on the training set with 1000 trees for a range of values of the shrinkage parameter λ . We then obtained a best λ 0.066 and a smallest RMSE 1.628784, which is slightly better than while using LASSO regression. (reference appendix on pg.4-5)

Bagging: Afterwards, we performed bagging method using the out-of-bag error estimate with mtry value of 182(numbers of predictors) and ntree value of 1000. We got the RMSE of 1.475989, which is a big improvement compared to LASSO and boosting. We then hold onto this value to compare with the random forest method, which utilizes even a smaller number of variables.

Random Forest: Performing a random forest is similar to bagging, as we fixed all the values of parameters and using out-of-bag error estimate except we used 60 (numbers of predictors/3 = 182/3) for mtry argument. The testing set RMSE associated with the random forest for regression tree is 1.478593 which is similar to what we obtained from bagging but much smaller than that of LASSO and boosting.

★ Table 1: RMSEs obtained from 4 models (All based on predictors p = 182)

Model	RMSE	Comment
LASSO	1.650821	Best $\lambda = 0.010$
Boosting	1.628784	Best $\lambda = 0.066$
Bagging	1.475989	mtry = 182
RandomForests	1.478593	mtry = 182 / 3 = 60

(For more details of the coding, please see appendix on pg. 4-5)

Different RMSEs associated with different models are shown in Table 1. Since the RMSEs obtained from bagging and random forest based on 182 predictors were very similar, we decided to do more comparison with these two models. We extracted all the predictors in the order of importance using the importance() function and the dplyr arrange function. We tried different range of thresholds from 3 to 11, and we noticed that 8, 9, and 10 were the best among this range. We then performed both bagging and random forest models once again to experiment with the best threshold values (8,9,10) of variable importance to get the most significant predictors. The results of RMSEs are summarized in Table 2.

★ Table 2: Bagging vs RandomForest with different number of predictors

Model	Threshold %IncMSE >	Number of Predictors p	mtry	RMSE
Bagging	10	28	28	1.434410
Bagging	9	31	31	1.433804
Bagging	8	32	32	1.438456
RandomForest	10	33	11	1.423562
RandomForest	9	35	11	1.423923
RandomForest	8	41	13	1.421445

(For more details of the coding, please see appendix on pg. 5-8)

After analyzing the results, we concluded that random forest yielded an improvement over bagging in this case. Thus, we decided to choose random forest as our final model and worked further on improving our final RMSE score using different subsets of predictors.

- Processing Random Forests

Having decided to use Random Forest as our final model, we tried to find the optimal subset of predictors. After running the initial randomForest() with mtry of p/3 and listing 182 predictors in the order of the importance, we not only tested different values of threshold, but also experimented different values of mtry. The results of RMSEs with corresponding threshold and mtry values are shown in Table 3.

★ Table 3: The results of processing Random Forest Model

Model	Threshold %IncMSE >	Number of Predictors p	mtry	RMSE
RandomForest	8	41	p/3 = 13	1.421445
RandomForest	8	41	14	1.422817
RandomForest	8	41	12	1.425424
RandomForest	9	35	p/3 = 11	1.423923
RandomForest	9	35	12	1.422035
RandomForest	9	35	10	1.421007
RandomForest	10	33	p/3 = 11	1.423562
RandomForest	10	33	12	1.412062
RandomForest	10	33	10	1.414207
RandomForest	11	29	p/3 = 10	1.424145

(For more details of the coding, please see appendix on pg. 6-8)

We can observe from Table 3 that, We can observe that the model with 33 predictors and mtry of 12 generates the best RMSE.

3 Results

Using our final model as the random forest with 33 predictors and mtry of 12 with ntree of 1000 yielded our best evaluation metric value of 1.35764 on Kaggle public leaderboard. The model that is obtained under many experiments performs well since it generated the lowest RMSE. When we applied our final model to the test data, we not only beat all of the benchmarks, but also outperformed other teams as we placed first in the competition.

4 Conclusions

We believe that there are several reasons to the successful outcome. Not only we had the ability to choose the right choice of the model, select important and necessary features, and remove highly correlated predictors, but also we were able to be creative and add additional features into our dataset.

We approached our analysis by crossing out the worse models. Despite the interpretability, LASSO wasn't the final selection because the assumption is based on linear regression and the response are not strongly correlated with other variables. The weakness of boosting method is its sensitivity to overfitting of the data contains too many noises and difficulty of tuning compared to random forest. We had to take the number of important features and the value of mtry into account in order to evaluate the better method between bagging and random forest, which turned out to be random forest model (Table 2).

After trying different types of models, we were confident that random forest worked the best after narrowing down our methods (Table 1 and Table 2). After selecting random forest model to dive into further exploration, our remaining task was to tune random forest with different values of predictors and mtry to find the best condition for the lowest RMSE value (Table 3).

On top of finding the right threshold value to reduce into a major set of predictors and removing highly correlated variables, another way to strengthen our model was to add additional variables of the binary columns for high levels of the existing binary columns and time intervals, hours and Days columns. In fact, 8 of the features out of 11 new features were part of the final 33 predictors.

Compared to the Kaggle public leaderboard score of 1.35764, we yielded 1.37819 from the private leaderboard. Adding the remaining test data caused a slight RMSE change, but we consider this two scores are similar. This is because random forest method with enough important predictors and sufficiently large value of trees limits overfitting as well as errors due to bias and therefore yield useful results.

On the other hand, there are still some possible ways to improve our model. For example, there's also a lot of inspiration for creating new variables based on the existing columns. We can also consider performing some of the advanced techniques such as oversampling and undersampling, and testing wider range of ML techniques that we haven't seen in class, such as XGBoosting method.

Overall, we are satisfied with our results, but we acknowledge that there could a possibility that our result could score lower RMSE metric by applying more advanced techniques.

5 Statement of Contributions

All the members collaborated on every parts of the project.

Predictor Selection: Yanhua Lin, Hana Yerin Lim, Yingzhen Zhao

Processing Methods: Hana Yerin Lim, Yanhua Lin, Yingzhen Zhao

Report: Yingzhen Zhao, Hana Yerin Lim, Yanhua Lin