

Twitter Sentiment Classification

Giuseppe Esposito, Syed Ali Abbas

Politecnico di Torino

Student id: s302179, s292423

s302179@studenti.polito.it, s292423@studenti.polito.it

Abstract—The goal of this report is to provide a possible approach to solve the classification problem of tweets sentiment. Our work takes in account the case normalization carefully to get a good data preprocessing. The results obtained are quite accurate also before the hyperparameter tuning.

I. PROBLEM OVERVIEW

The aim of this work is to solve a problem of supervised machine learning so we have 2 collection of tweets:

- Development set: which contains 225000 comments classified as positive sentiment or negative sentiment.
- Evaluation set: which contains 75000 tweets without label assignment.

Most of our observations will be based on the following table:

Variable	Cardinality
sentiments	2
ids	224716
Time zone	1
Flag	1
User	10647
Text	223106

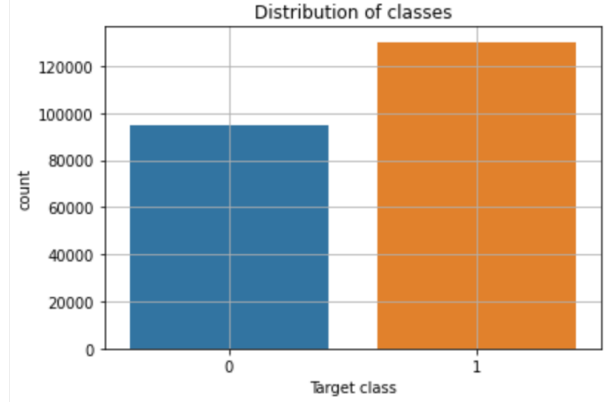
Table:1

Upon a first analysis it is noticeable saying that the analysis that we are carrying out is related to tweets which has no specific topic since the flag field has just one unique string value: "NOQUERY". Because of the aim of the analysis we chose **as target the sentiment field** where there are 2 unique binary values:

- 1: positive sentiment
- 0: negative sentiment

Moreover we can notice that it is a West-American community based analysis because in the Date field all the values has the same time zone that is PDT¹ so it is useful to analyze the general distribution of sentiment.

¹Pacific Daylight Time



Looking at the sentiment distribution we can say that it is quite unbalanced because there are way more positive sentiments than negative one.

Furthermore, it is interesting looking at the correlation matrix. The most noteworthy observation that we can infer is that there exist a dependence between the user names and the sentiment and this means that some people have a specific behaviour which leads them to write tweets with a specific sentiment.

II. PROPOSED APPROACH

You can use citations as follows: [1] (you can add BibTeX citations in the *bibliography.bib* file).

A. Preprocessing

Now let us have a look to the following table:

Variable	Cardinality
ids	224716
User	10647
Text	223106

Table 2

A noteworthy observation is the comparison between these 3 values, in fact there are more tweets than rows and we expected that the difference between these two cardinalities would have given us the number of duplicate data but then caught up with the idea that different people could have written the same comment so the way to obtain the cardinality of duplicate data was meaningless. Then we thought that, since the ids are uniquely associated to the tweets by definition, the actual number of duplicate data is

C. Hyperparameters tuning

At this point we had to split the DataSet in a training set and in a test set in order to tune our model. We decided to keep the 80% for the training set and the 20% for the test set. In order to find the best combination of training set and test set, we split the dataset by means of KFold method of Sklearn.model_selection (with the "shuffle" argument equal to True) and then iterated over the resulting indices of the dataset to check which one gave us the best score. However we also tried to deal with this in another way: we firstly used the train_test_split making the data order random and then, when we computed the Grid Search, we set the parameter cv = 5 (in order to have the same split in training set, 80%, and test set, 20%, as before) and we realised that this approach gave us the best score. To compute model tuning we took in account different combinations of parameters (which would have made sense to use in order to get a satisfactory baseline score) computing a grid search.

Model	Parameter configuration
Random Forest	"nEstimators": [100,200,300], "maxFeatures":["auto", 'sqrt', 'log2']
LinearSVC	"C": [0.1,1,10] "dual": [True,False] "fitIntercept" : [True,False] "maxIter" : 1000
Decision Tree	"splitter": ['best', 'random'], "maxFeatures": ['auto', 'sqrt', 'log2']

Table II-C

We repeated this process for each model and we tried the parameters shown in TableII-C.

III. RESULTS

Now we are going to discuss the main results of the analysis.

Model	Optimal Param. config.	F1-score
Random Forest	"nEstimators": 300 "maxFeatures": 'sqrt'	F1-score
LinearSVC	"nEstimators": "C":0.1, "dual":True "fitIntercept" :True "maxIter" : 1000	F1-score
Decision Tree	"splitter": ['best', 'random'], "maxFeatures": "sqrt"	F1-score

Table II-C

It is noteworthy noticing in Table III that the worst score is preformed by Decision Tree and a reason is embedded in the definition of decision tree algorithm because this model was optimized for classification tasks which provides more than 2 classes.

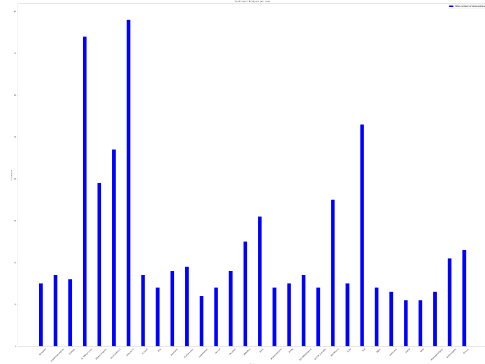
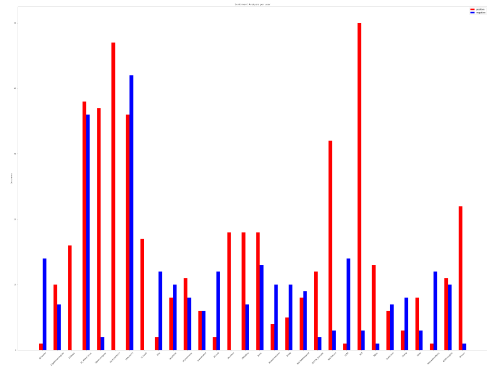
Another reason should be sought in the number of features

because when it is very high that is our case (120649) decision trees are not the best option because the Random Forest outperforms it, since it randomly selects observations and specific features to build multiple decision trees and then each decision tree 'votes' for the best class assignment.

Nevertheless, looking at the Random Forest score we can observe that it is not the best one and the reason of this is the same of the Decision Tree algorithm: it is optimized for classification tasks which provides more than 2 classes.

We can finally encounter that Linear Support Vector Machine reaches the best score and the reason lies in the way it manages sparse matrices, the high dimensionality of the problem and the scalability of the model.

IV. DISCUSSION



Overall it is interesting saying that the differences among the three scores are quite small and this means that the three models works pretty well with the approach proposed.

Another measure that is interesting to take in account is the time expenditure that is hugely higher in Random Forest model building because it is computing a number of trees that is equal to the number of estimators and it takes much time.

We have also thought to a possible future implementation that could make the analysis more accurate so that the score would be improved. The changing on our approach will be the following:

- Analyse the correlation between the sentiment and the user. Since we are comparing a categorical field and a binary one we have shown this correlation in Figure IV
- We can notice that for some users there exist this correlation because some of them have definitely a negative or a positive behaviour. In order to make this assumption more meaningful we could set a threshold on the maximum number of observations as shown in figure IV.
- This means that the accuracy of our model is badly affected by the sentiments, related to the user, that are different from his usual behaviour.
- Then the idea could be: for each user extract the corresponding tweet text which has a different behaviour from the usual one.
- then compute a tf-idf only for these "strange" tweets (Output = Vector-2)
- check the words with the highest tf-idf value that will be the words that will most negatively affect the accuracy of the model
- In the end compute the tf-idf vectorizer on the whole column of text (Output = Vector-1)
- then remove the words that in Vector-1 have the highest value of tf-idf from the Vector-2 if they are still in there (maybe a possible strategy could be: substituting them to the list of stop words since we have not removed them).
- so that we compute a sort of biased tf-idf vectorizer in order to reduce the dimensionality of the problem and the analysis will be also more meaningful.

REFERENCES

- [1] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep learning*, vol. 1. MIT press Cambridge, 2016.

can we put also references to the slides of the course?