

Project 2 on Machine Learning

Text classification

Team Yoor

Sergei Volodin¹, Baran Nama¹, and Omar Mehio¹

¹EPFL

{sergei.volodin,baran.nama,omar.mehio}@epfl.ch

Abstract—A classification dataset consisting of Tweets is being studied. First, the data is thoroughly explored using visual aids. Several basic Natural Language Processing methods are applied. Results are evaluated using cross-validation. Model overview is given and the best model is chosen.

I. INTRODUCTION

This paper investigates into improving the quality of sentiment analysis on Tweets dataset [7]. It consists of $N_1 = N_2 = 1250000$ positive/negative tweets, each of them representing a message in English and numerical alphabet Σ with no longer than 140 characters. This way, each of $N = N_1 + N_2 = 2500000$ tweets is assigned to one of the classes $\mathcal{C} = \{:(, :)\}$. The task is to minimize the classification error. In other words, if $\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N$ is the dataset with tweets $x_i \in \Sigma^*$ being messages and $y_n \in \mathcal{C}$ being class labels, the goal is to train a classifier $f: \Sigma^* \rightarrow \mathcal{C}$ which minimizes the loss function $l(y, \hat{y}) = [y \neq \hat{y}]$.

The task of sentiment analysis of tweets was thoroughly studied [2], [3], [4], [5]. These solutions usually give accuracy around 80%. Several techniques were applied, mostly consisting of two steps. First, the words are converted to dense vectors using Glove, word2vec, cbow or skip-gram models. After that, the resulting word vectors are used to construct features for the whole tweet. At the end, the vector is feeded into a classifier, such as SVM or Logistic Regression. Two latter steps might be replaced with a neural network accepting variable-length input such as RNN or CNN. Moreover, the embeddings themselves might be trained using backpropagation while training the classifier.

Claim: it is possible to find a model which fits the data better than the current state-of-the-art using expert knowledge on the Tweets dataset.

Next sections describe in details our approaches and compare them to various baselines.

II. MODELS AND METHODS

This paragraph describes the data being studied. Tweets are short messages no longer than 140 characters long [1]. The table I represents a few examples from the small dataset. Being considered an informal way of communication, tweets

often contain misspelled words and letter repetitions, grammatical and other writing mistakes. Besides plain text with punctuation, tweets also contain hashtags, two types of tags, *user* and *url*, which are tokens for replaced user mentions and URL links, respectively. In addition, tweets sometimes contain emoticons. Despite the fact that objects in the training dataset mostly comply with its classes, some tweets cannot be determined as positive/negative even by a human (appear neutral). Moreover, some of the tweets are clearly mislabeled in the training data. Dataset contains repeating tweets.

Representation of tweets was through the occurency matrix, where each row corresponds to a tweet, each column corresponds to a term in the vocabulary and each cell corresponds to the number of times this word occurs in this tweet. Through preprocessing the aim is to maximise the knowledge that can be extracted from the datasets in order to build this matrix. Hence the focus was on increasing the occurrence of words through correction of misspelled words and stemming techniques.

Usually tweeters repeat a letter in a word to stress on a certain emotion, as an example: "i am soooooo angryyy". This problem was tackled by two approaches. The first approach is called stemming. A formal definition of stemming: is the process by which we find the origin of the word. Linguistically any word can be changed to different forms by either adding affixes or prefixes, thus one word can have multiple forms, all corresponding to the same origin. Through stemming, the aim is to minimise the presence of these different forms and replacing them with the original stem. Algorithm We do this by replacing every repeated letter by two instances of that letter. We check if the word is valid (contained in a dictionary) if it is the case then we found the correct form of the word, otherwise we reoccur on the new word. Our end result is a set of tweets with most of the words having repeated characters replaced by their stems.

The second approach adopted is misspelling correction. The process goes as follows: a token in the tweet is proposed to an english dictionary, the dictionary responds by either indicating that the word is valid or not. This dictionary

Class	Row	Message	Comment
:(171	<user> 5k i could cry i'm so unfit !	User tag
:(99804	if i feel like this tomorrow i'm going to the er	Contraction, missed comma
:)	524	its the weekend ! ! <user> s coming home andddd <user> s baby shower , excitedddd	Letter repetitions
:)	99615	<user> aren't you just an adorable granny <url>	URL tag
:)	443	<user> :d :d :d :d :d wish jocelyn had a twitter . #kudos for her too .	Hashtag, emoticons
:)	468	retweet if you was born in the 90 ' s ! #90's babies	Grammar
:(14	<user> i'm white . #aw	Appear neutral
:(99594	<user> <user> <user> they're there tonight ! ! !	Appear positive

Table I
EXAMPLES OF TWEETS IN THE SMALL DATASET

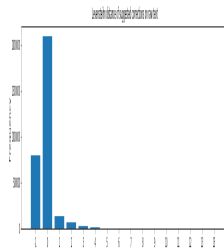


Figure 1. Exploratory data analysis using charts

provides the most relevant suggestions to correct this word. As a metric to modify a word , the levenshtein distance was use. Levenshtein distance is defined as : "the distance between two strings is the minimal number of insertions, deletions, and substitutions of one character for another that will transform one string into the other". Thus to explore how many misspelled-wrong words are in the dataset , we plotted the distribution in Fig []. All most of the data correction is covered by setting our levensteihn threshold to 4.

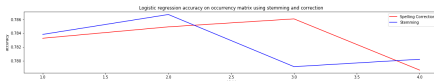


Figure 2. Exploratory data analysis using charts

The metric used to compare both methods was prediction accuracy of logistic regression. As a first step the combined dataset of both the positive and negative tweets was exposed to stemming and correction. Then for each method a vocabulary of words is constructed. Finally, the occurrence matrix is built from both the tweets and the dictionary. K-fold is used to split the matrix and then logistic regression is run using the occurrence matrix as the data matrix, labels are split into 0 denoting positive sentiment and 1 denoting, a negative sentiment. As observed in figure x , the spelling correction method is superficial to the stemming method

with reaching a maximum of 78.6

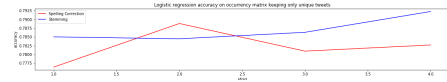


Figure 3. Exploratory data analysis using charts

Another property observed in the dataset is the presence of redundant tweets. To view their impact ,repeated instances of a tweet are eliminated and the newly formed unique dataset is exposed to the previous methods. And as observed once again the correction spelling method overpowers the stemming technique. Although with very small improvement in accuracy prediction (78.88

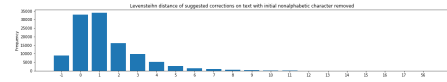


Figure 4. Exploratory data analysis using charts

According to fig ,we can observe that many tokens in our datasets start with non alphabetic characters. Those tokens were eliminated from the previous experiments as all dictionaries can't recognize non alphanumeric words. Thus a new set of tweets was constructed having all tokens that contain characters, but start with a non-alphabetic character to be stripped of that character or characters. AS before, the distribution of levensteihn distances for suggested words was plotted,which results in a more uniform distribution. It is suspected that this is due to the informal language used in hashtags. The newly constructed dataset is applied

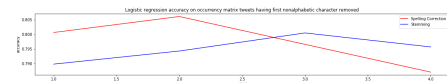


Figure 5. Exploratory data analysis using charts

to previous methods .Eliminating redundant tweets , since as observed this gave us a slight improvement. Then the

dataset is exposed it to stemming and misspelling correction. Accuracy prediction then improves to 0.802 , the best score ever reached.

Now with the optimal model having initial nonalphabetic characters removed , redundant tweets eliminated and misspelling correction applied with levensteihn threshold set to 5. The dataset is tackled from a different angle, the effect of the frequency of which one token appears in the document is explored. Thus two approaches are considered : removing most frequent tokens that are common in both datasets and removing nonfrequent words from the vocabulary. The results of the first experiment came as a surprise. As suspected that common words would make it harder for the classifier to distinguish between a positive sentiment tweet and a negative one. A reasoning for that is due to the limited number of characters provided to users to express themselves and thus most of the tweet gets distributed among tagging other users, using punctuation and stopwords to express feelings. Also eliminating non frequent tokens from our vocabulary has helped improving accuracy , as observed in fig.

A several models were considered. First, the baseline was the GloVe [6] model, which considers training embeddings on the co-occurrence matrix. Unfortunately, the model did train quite slow on the dataset gave meaningless embeddings (closest words for a given word did not correspond to it in meaning). Therefore, it was not considered.

Secondly, a word count model was considered as a replacement. Here, a tweet is represented by a high-dimension sparse vector, each item corresponds to a number of occurrences of this word in the tweet. This vector was further classified using several techniques, such as SVM and Neural Network. This method gave higher accuracy and allowed for tweaking its parameters with relatively small training time (a Google Cloud instance with 8 CPU was used)

III. RESULTS

Model	Train accuracy	Validation accuracy	Kaggle accuracy
GloVe RNN	6	87837	787
GloVe CNN	7	78	5415
GloVe FC	545	778	7507
WC	545	18744	7560
W	88	788	6344

Table II

COMPARISON OF DIFFERENT METHODS (WC=WORD COUNT, GLOVE – PRETRAINED GLOVE VECTORS)

The following models were considered: CNN, RNN, Logistic Regression, Fully Connected Neural Network. The best model demonstrated This (does not) correspond to already conducted experiments [99, 98, 97]. Our contribution consists of running *method* with *xxx* modified with *yyy* and this does (not) give an improvement of 0.01231%

IV. DISCUSSION

Our experiments lack *zzz*, which can be improved by doing also *ttt*

V. SUMMARY

We have shown that it is possible to predict tweets using *aba* better than state-of-the-art.

REFERENCES

- [1] Twitter
- [2] Go, Alec, Lei Huang, and Richa Bhayani. "Twitter sentiment analysis." *Entropy* 17 (2009): 252.
- [3] Kouloumpis, Efthymios, Theresa Wilson, and Johanna D. Moore. "Twitter sentiment analysis: The good the bad and the omg!." *Icwsn* 11.538-541 (2011): 164.
- [4] Tapan Sahni, Chinmay Chandak, Naveen Reddy, Manish Singh. "Efficient Twitter Sentiment Classification using Subjective Distant Supervision"
- [5] Alec Go, Richa Bhayani, Lei Huang. "Twitter Sentiment Classification using Distant Supervision"
- [6] Jeffrey Pennington, Richard Socher, Christopher D. Manning. "GloVe: Global Vectors for Word Representation"
- [7] Competition and data downloads on kaggle.com