

TWEET SENTIMENT EXTRACTION

JIA HUANG

ABSTRACT. With all of the tweets circulating every second it is hard to tell whether the sentiment behind a specific tweet will impact a company, or a person's, brand for being viral (positive), or devastate profit because it strikes a negative tone. Capturing sentiment in language is important in these times where decisions and reactions are created and updated in seconds. But, which words actually lead to the sentiment description?

In this competition we've extracted support phrases from Figure Eight's Data for Everyone platform. The dataset is titled Sentiment Analysis: Emotion in Text tweets with existing sentiment labels, used here under creative commons attribution 4.0. international licence. Your objective in this competition is to construct a model that can do the same - look at the labeled sentiment for a given tweet and figure out what word or phrase best supports it.

CONTENTS

1. Introduction	2
2. Data Exploratory Analysis	2
3. Model Construction and Training	4
4. Conclusions	5
List of Todos	6

Date: 2020-10-15.

1991 *Mathematics Subject Classification.* Artificial Intelligence.

Key words and phrases. Emotional extract.

1. INTRODUCTION

With all of the tweets circulating every second it is hard to tell whether the sentiment behind a specific tweet will impact a company, or a person's, brand for being viral (positive), or devastate profit because it strikes a negative tone. Capturing sentiment in language is important in these times where decisions and reactions are created and updated in seconds. But, which words actually lead to the sentiment description? In this competition you will need to pick out the part of the tweet (word or phrase) that reflects the sentiment.

Help build your skills in this important area with this broad dataset of tweets. Work on your technique to grab a top spot in this competition. What words in tweets support a positive, negative, or neutral sentiment? How can you help make that determination using machine learning tools?

In this competition we've extracted support phrases from Figure Eight's Data for Everyone platform. The dataset is titled Sentiment Analysis: Emotion in Text tweets with existing sentiment labels, used here under creative commons attribution 4.0. international licence. Your objective in this competition is to construct a model that can do the same - look at the labeled sentiment for a given tweet and figure out what word or phrase best supports it.

2. DATA EXPLORATORY ANALYSIS

Now let's take a look at the training data provided by this competition and do some exploratory data analysis to get a detailed understanding of the training set.

	textID	text	selected_text	sentiment
0	cb774db0d1	I'd have responded, if I were going	I'd have responded, if I were going	neutral
1	549e992a42	Sooo SAD I will miss you here in San Diego!!!	Sooo SAD	negative
2	088c60f138	my boss is bullying me...	bullying me	negative
3	9642c003ef	what interview! leave me alone	leave me alone	negative
4	358bd9e861	Sons of ****, why couldn't they put them on t...	Sons of ****,	negative

Flag.1

(1) Text and selected_text of one record in training set are null, so delete this data.

(2) Check the distribution of the sentiment category in the training set data.

🔥 (None)-develop (2020-10-15)

2

Committed by: jiaHuang

Out[21]:

	sentiment	text
1	neutral	11117
2	positive	8582
0	negative	7781

Flag.2

We can see that across the entire dataset, there are 111,117 tweets for neutral emotions, 8,582 tweets for positive emotions, and 7,781 tweets for negative emotions.

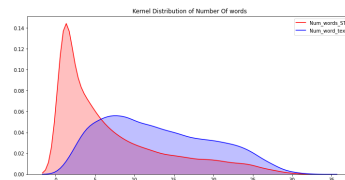
(3) Look at the Jaccard similarity between Text and selected_text.

	textID	text	selected_text	sentiment	jaccard_score
0	cb774db0d1	I'd have responded, if I were going	I'd have responded, if I were going	neutral	1.000000
1	549e992a42	Sooo SAD I will miss you here in San Diego!!!	Sooo SAD	negative	0.200000
2	088c60f138	my boss is bullying me...	bullying me	negative	0.166667
3	9642c003ef	what interview! leave me alone	leave me alone	negative	0.600000
4	358bd9e861	Sons of ****, why couldn't they put them on t...	Sons of ****,	negative	0.214286

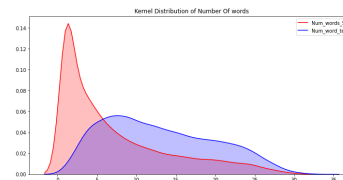
Flag.3

As can be seen from the above output results, such as I'd have responded, if I were going. This kind of words has the highest similarity with its emotional polarity, with a similarity value of 1, followed by words like leave me alone, with a negative emotional polarity, with a similarity value of 0.6.

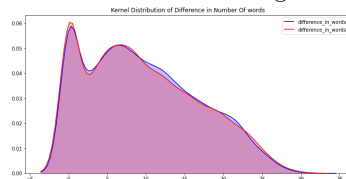
(4) The kernel distribution graph of word length.



Flag.4



Flag.5



Flag.6

(5) The word cloud shows how often words appear in different categories of Twitter.



3. MODEL CONSTRUCTION AND TRAINING

🔥 (None)-develop (2020-10-15)

the start, so we can use has trained model parameters, used after fine-tuning in our NLP tasks. Here I used Huggingface's open source Transformers library, which provides a lot of pre-trained models with a very good and easy-to-use interface that can be directly called and built using PyTorch or TensorFlow.

- Construct tokenizer, and transform the text of training set and test set into token.
By making a full connection layer of 768x1 to the output vector to get the Head, the output becomes the vector of BatchxMAX_LENx1, then rehaspe is the vector of BatchxMAX_LEN, and then Softmax.
- Divide the training set into 5 copies, and take 4 copies of training and 1 copy for verification.
- Three epochs are trained each time, and parameters of the epoch with the lowest Loss in verification set are taken. We get a new model at the end of each training session, so we end up with 5 models.
- When making predictions, the prediction results of these five models will be averaged.
- The predicted values of the five trained models were averaged.

	textID	text	sentiment	selected_text
0	f87dea47db	Last session of the day http://twitpic.com/67ezh	neutral	last session of the day http://twitpic.com/67ezh
1	96d74cb729	Shanghai is also really exciting (precisely -...	positive	exciting
2	eee518ae67	Recession hit Veronique Branquinho, she has to...	negative	such a shame!
3	01082688c6	happy bday!	positive	happy bday!
4	33987a8ee5	http://twitpic.com/4w75p - I like it!!	positive	i like it!!
5	726e501993	that's great!! weee!! visitors!	positive	that's great!!
6	261932614e	I THINK EVERYONE HATES ME ON HERE lol	negative	hates
7	afa11da83f	soooooo wish i could, but im in school and my...	negative	blocked
8	e64208b4ef	and within a short time of the last clue all ...	neutral	and within a short time of the last clue all o...

Flag.3

4. CONCLUSIONS

Through this project, we further understand the related contents of big data analysis, and also have a preliminary understanding of natural language processing. I believe this project experience will be of great help in the following study.

LIST OF TODOS

(A. 1) SCHOOL OF COMPUTER SCIENCE,, XI'AN SHIYOU UNIVERSITY, SHAANXI 710065, CHINA
Email address, A. 1: xxx@tulip.academy