

Report for Final Project : Analysis of Rumor

Chen Ying, Zhu Simo, Liu Guoding, Mao Haining

January, 2019

Contents

1	Introduction	2
2	Preliminary Process (Data Reading)	2
3	EDA & Feature Extraction	4
3.1	Exploratory Data Analysis (EDA)	4
3.2	Feature Extraction	6
4	Model Selection and Model Training	8
4.1	Comparison between several models	8
5	Evaluation and Analysis	9
5.1	Feature Importances & Confusion Matrix	9
5.2	Analysis of the Result	9
6	Further Improvemets	9

1 Introduction

Text data is around us everywhere, and how to deal with those words, sentences or passages has long been a hot topic for natural language processing fields. With the rapid development and popularization of social media services, rumors are spreading with unprecedented rapidity and have a tremendous impact on human society. Meanwhile, the development of artificial intelligence technologies provides a promising approach for social media platform to automatically detect rumors.

In daily life, people usually distinguish the authenticity of microblog events based on their common sense or through news websites and public communities, however, the reports of this kind of website media are incomplete and have certain time delay. Therefore, the automatic identification of rumor events can help us better prevent rumors and assist management agencies in rumor intervention and governance.

Under that situation, we hope to set up a model to distinguish those rumors from normal ones with only contents rather than making predictions based on all comments when users forward the message. In the following part, we will explain how the semantic and syntactic information influence the authenticity of a given context, how different representing methods and classification models performs on this classification problem, and, finally, how our primary features can be explanatory for deeper analysis of rumor.

2 Preliminary Process (Data Reading)

The documents in the dataset are scores of JSON (JavaScript Object Notation) documents, the information contained in one json document has the construct below:

JSON	原始数据	头
保存 复制 全部折叠 全部展开	▼ 过滤 JSON	
multi:	multi:	null
▼ text:	text:	"朝鲜日报最新披露中日军力对比称：日本空军自卫队飞行员每周平均飞行训练为168小时，这一训练时间是中国空军的4倍。日本虽然没有核武器，但是拥有先进的核技术，一旦发生战争组装一枚核武器的时间只需要4小时，而中国导弹从内地发射打到日本本土需要4周时间。差距明显，这就是中国只敢抗议的根本原因。"
▼ user:	user:	verified: false description: true gender: "m" messages: 5659 followers: 3730 location: "海外 美国" time: 1316463415 friends: 636 verified_type: -1 has_url: false comments: 140 pics: 0 source: "微博 weibo.com" likes: 1 time: 1351433795 reposts: 217

JSON	原始数据	头
保存 复制 全部折叠 全部展开	▼ 过滤 JSON	
multi:	multi:	null
▼ text:	text:	"人间惨剧：今天下午约14点，宁波妇儿医院，一妇女携带一婴儿在住院楼跳楼，后抢救无效死亡。具体情况有关部门正在调查。据现场网友称妇女因小孩病重，加上负担不起昂贵的医疗费，带着刚满月的宝宝从12楼跳楼身亡。【蜡烛】底层民众的医疗费用猛于虎，国人的生命其何等脆弱！【泪】"
▼ user:	user:	verified: false description: true gender: "f" messages: 5653 followers: 227833 location: "上海" time: 1312112304 friends: 907 verified_type: -1 has_url: false comments: 55 pics: 1 source: "微博 weibo.com" likes: 0 time: 1347334462 reposts: 225

The useful information of it contains:

- Text
- Weibo Features
 - (a) Whether the weibo has URL
 - (b) Number of comments
 - (c) Pics
 - (d) Sources of this weibo
 - (e) Likes
 - (f) The time when this weibo is sent
- User Features
 - (a) Whether the user has description
 - (b) Whether the user is verified and the verified type
 - (c) The gender of user
 - (d) Number of followers
 - (e) Number of friends
 - (f) The location of the user
 - (g) The time when the user joined weibo
 - (h) Number of messages user had sent

To get this information, we preprocess these json documents. In this process, we use the 'os' & 'json' module to read all the json data and extract information we need. Since the dataset has been processed, there is no missing data (in other words, all the json documents have complete information) or duplicate records, the presentation of features are also consistent, which saves us much effort.

Then we encode each feature by the method **ordinary encoding**. First, we use list to collect data, and then transform it into ndarray to adapt (except text data, which will be processed in the feature extraction part). In the meantime, we transform the data type into int (some examples are below).

```
1 has_url = np.array(has_url).astype(int)
2 verified = np.array(verified).astype(int)
3 description = np.array(description).astype(int)
4 gender = np.array(list(map(trans_gender, gender)))
5 followers = np.array(followers)
6 friends = np.array(friends)
7 category = np.array(category)
```

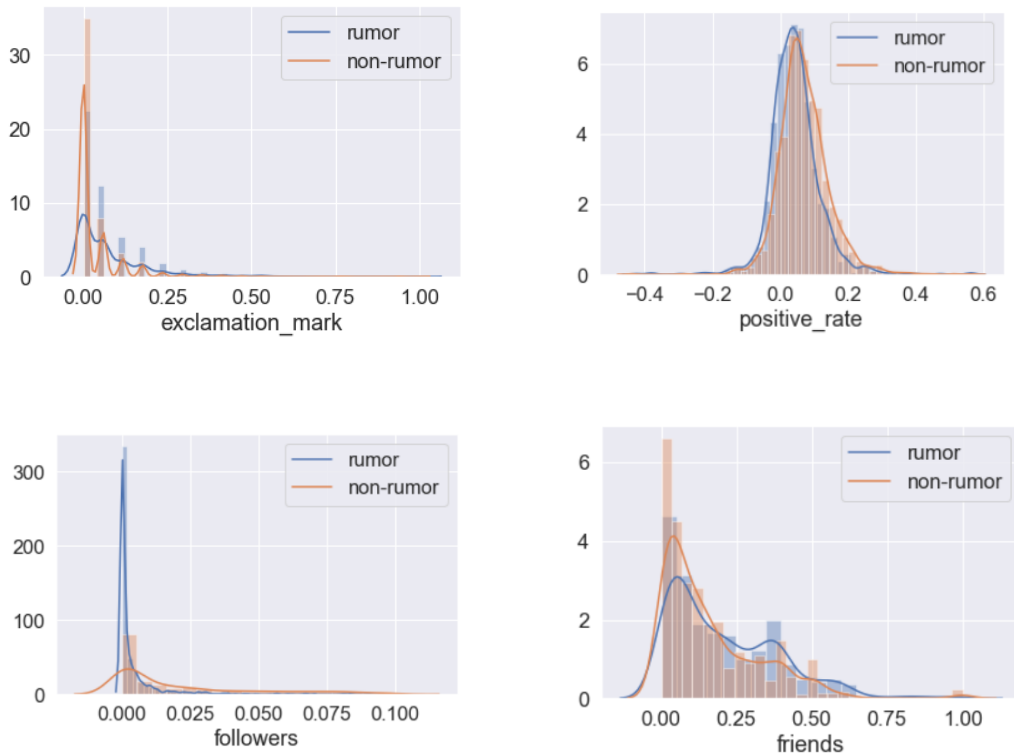
3 EDA & Feature Extraction

In this project, we try to extract features from the information in each json document. Before this project, we read some papers, among which the paper written by Prof. Liu Zhiyuan impressed us most. He asserts that there are several methods to detect rumors according to the features of the time series, but in that method, we can't detect rumors as soon as they are sent. Additionally, stopping the rumor early will help stop the potential social loss. This idea determines the features we finally choose. It should be attached great importance here that in this project we want to extract **the early characteristics of rumors**. The information of rumor after the time when rumor is sent, such as comments of this rumor isn't considered into this classification model.

First we use Exploratory Data Analysis methods to see the general features of rumors in this dataset and then we select the features which will be used in the model training process.¹

3.1 Exploratory Data Analysis (EDA)

In EDA process, we first explore the single feature of rumors as well as non-rumors and then compare them in one picture. For example we count the number of exclamation marks(**the coordinate has been normalized**) in each text and then describe it through frequency histogram to figure out the characteristic of distribution.



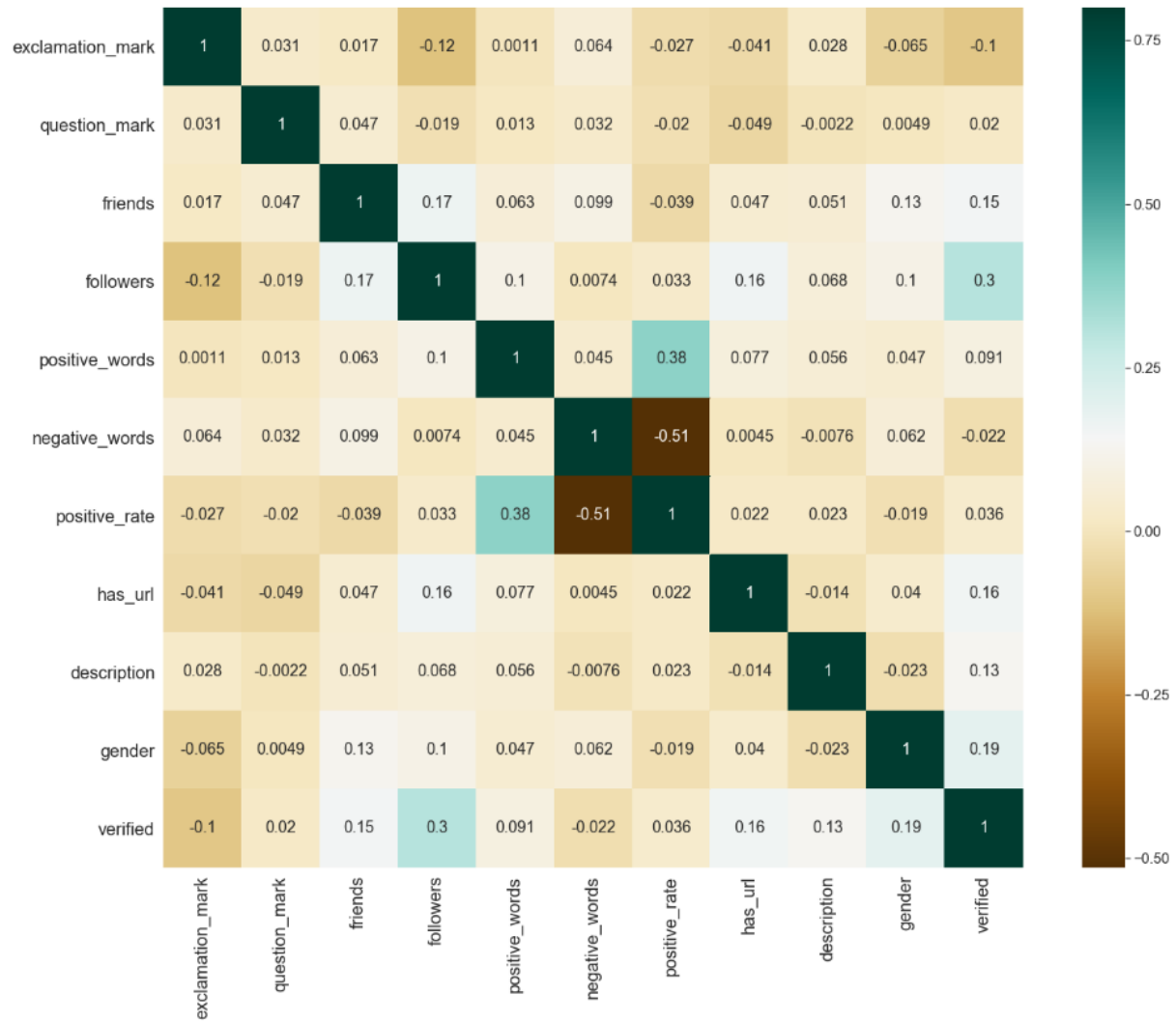
¹In the real process, we must construct feature first and then use EDA to see the characteristics. After this we determine whether we will use these features. The order this report describe is little different from the real order.

In this process, we find that rumors are more frequently expressed with more exclamation marks as well as question marks. Besides, rumors tend to present negative attitude. Meanwhile, we find the number of friends of rumor-mongers' weibo account is statistically larger.

Then we explore **the correlation between different features** (variables). At the same time, we can see the features more clearly. Actually, none of these features are predominant feature which can almost determine the classification with a relatively high coefficient of association, so the combination of different features may be a good method.



We also find an interesting thing that the positive rate of the rumors is not always low, instead, it is very low or very high. That reminds us the rumors may be too positive or too negative, so decision tree may be a good choice. However, generally speaking, the rumors are always radical.



3.2 Feature Extraction

After the exploration process, considering the construct of the dataset and the characteristics of the rumor texts, we select the features below as the features we aim to extract and use as analytical basis in this classification model.

- Text Features
 - (a) Number of exclamation marks
 - (b) Number of question marks
 - (c) Number of positive words
 - (d) Number of negative words
 - (e) Positive rate of the text
- Weibo Features
 - (a) Whether the weibo has URL

- User Features
 - (a) Whether the user has description
 - (b) Whether the user is verified
 - (c) The gender of user
 - (d) Number of followers
 - (e) Number of friends

Since the features (except text features) have been processed in the preliminary process, this time we only need to extract the features from weibo text message mainly through NLP (Natural Language Processing) methods.

This process is hard to handle. We refer to papers in the field of NLP and conclude that **rumors are always radical**. So we thought about those features which can reflect whether the text is “radical”. According to the general knowledge of linguistics (~~though quite poor~~), finally we select these features to construct our rumor classification model:

- Number of exclamation marks
- Number of question marks
- Number of positive words
- Number of negative words
- Positive rate of the text

We use these features for the reason that if one text is radical, it will tend to use more exclamation marks and more question marks². Besides, rumors tend to use more positive or negative words (radical, in summary). In general the tonal of the text will be radical so we try to construct a feature which will reflect this point.

We know the character of one sentence will be represented by **the “key word” of this sentence**. The “key word” in this project is defined:

The words appeared in this sentence but don’t appear in most of sentences we use in daily life.

So we want to get the “key word” and evaluate the text by the “key word” of the text. Then we use module jieba to split the raw text message (before this we use regular expression to clean the text) and get the most importance “key words” of the text (jieba has a corpus and can return the importance words according to the TF). The number of importance words can be set, which forms a hyper-parameter.

Once we get the “key words”, we use SnowNLP module to evaluate the emotion of each word. The value of neutral word is 0.5. The value will be larger if the word is more positive and smaller if the word is more negative. So we count the positive words and negative words of these “key

²Just as the explanation we’ve proposed before in subsection 3.1

words” according to the rule: if the value is larger than 0.6, the word is positive word. If the value is smaller than 0.4, the word is negative word.

```

1 for x, w in jieba.analyse.textrank(sent, topK=top_word, withWeight=True,
   allowPOS=allowpos):
2     rate += (SnowNLP(x).sentiments - 0.5) * w
3     if (SnowNLP(x).sentiments < 0.4):
4         negative_num += 1
5     if (SnowNLP(x).sentiments > 0.6):
6         postive_num += 1
7     ##加上最后一个0.000000001是因为有的微博没有关键词，所以为了防止分母为0加上一个微小数
8 rate = rate / (len(jieba.analyse.textrank(sent, topK=top_word, withWeight=True
   , allowPOS=allowpos)) + 0.000000001)

```

At the same time, we calculate the positive rate of the whole sentence according to the “key words”. The algorithm is:

- First get the evaluate value of each word and then subtract it by 0.5. Then get the average of this subtracted value. This value is the whole positive rate of the whole sentence.
- Meanwhile, we count the number of exclamation mark Number of question mark and then transform the data into array.
- Then we get the whole features we want to extract. All the features have been transformed into array and all the data types are int or double.

4 Model Selection and Model Training

4.1 Comparison between several models

In this model, we want to select a model that can classify whether a document is a rumor or not. Because we cannot easily tell the answer, we need a model that is highly interpretable. Here are some alternatives:

- KNeighborsClassifier
- Logistic regression
- Decision tree
- Random Forest
- Adaboost

The model of deep learning is so complex that we can hardly interpret it, so we didn’t use neural networks.

Then we use cross-validation, split the dataset into training set (70%) and testing set (30%). Using the training-set to train the model and test it on the testing set, we can calculate the accuracy of the model. Initially we do not adjust the parameters of the model.

Random Forest is the best model for this problem.

In the process of adjusting the parameters of the LR model, we find that decreasing regular terms significantly increases model accuracy. This means we are dealing with an under-fitting problem. So Random Forest perform better than LR. Then we use GridSearchCV to get better parameters. At last we find that when we do not limit the max-depth and max-leaf-nodes, set estimated quantity to 118 and criterion to gini, we can get the highest accuracy.

Limited by the quality of the original dataset, the accuracy of the model is 0.789.

Then we can calculate the importance of each feature to tell which feature we should care more.

We can find some interesting results from the importance of feature. As predicted previously, all features are weak. The most important feature is the number of followers. This is a counter-intuitive result that neither validation nor positive word rate matter. This may be because the dataset contains rumors of many low-level users who have few followers. Another interesting finding is that although both positive and negative words are weak features, the positive rate of a sentence is a relatively strong feature. This could offer us some insights when identifying rumors.

5 Evaluation and Analysis

All contexts contain two kinds of information: semantic information and syntactic information. Study of the former can lead to judgement for each single word and phrase, and study of the latter focus more on the organization of texts. Above all, the key problem is the way to extract features from the text. In our project, we have found some interesting points.

5.1 Feature Importances & Confusion Matrix

5.2 Analysis of the Result

6 Further Improvemets