

# Rumor detecting

Masterarbeit  
zur Erlangung des akademischen Grades  
M.Sc. Internet Technologies and Information Systems

CHENG LI

Leibniz Universität Hannover  
Fakultät für Elektrotechnik und Informatik  
Institut für verteilte Systeme  
Fachgebiet Wissensbasierte Systeme  
Forschungszentrum L3S

Hannover  
13.11.2016

Examiner I	Prof. Dr. Wolfgang Nejdl
Examiner II	Dr. Eirini Ntoutsis

---

## Abstract

In this thesis, we propose a approach of identifying rumors in Twitter.

Titter is a mircoblogiging service that which are used by millions users. Users can publish and exchange information with short tweets whatever when and where. This makes it a ideal media for spreading breaking news and false rumors.

So automatic detecting rumors on social media has become a trending topic. But early researches mostly focused on rumors during one or several events like earthquake or terrorist attack [26] [33][31]. But in our work, we study on the general rumors.

And most of previous work for rumor detection modeled on static features like the content of tweets or propagation features, but they ignored that those features can change during the information's propagation over time.

We build up a system with random forest and the Dynamic Series-Time Structure (DSTS) [23] with the temporal features and their varieties over time. And we test the Spike Model [18], SIS Model and SEIZ Model [13] as temporal features in our model. To improve the time series model's performance at early stage of the event we develop a single tweet's credibility scoring model with zhou's CNN+LSTM model [39] which can predict a single tweet whether is rumor related or not with 81.20%accuracy. Finally we got the accuracy of the time series rumor detecting model get 91% accuracy within 48 hours.

Our experiments using the events from Twitter and our model demonstrates better performance on detecting rumors.



## **Zusammenfassung**



## **Acknowledgment**

I am dedicating this thesis to my parents. They have been strong and supportive of my choices during my masters and have pushed me to work harder everyday. I am very grateful to my mentor Dr. Avishek Anand for helping me across the finish line and helping me make defining career choices. My sincere gratitude to Prof. Dr. Wolfgang Nejdl for having faith in me and guiding me through my bachelors and masters. A big thank you to the LearnWeb team including Dr. Ivana Marenzi for giving me a project to contribute towards. I am especially grateful to M. Sc. Sergej Zerr for taking me under his wing when I first arrived and I owe a lot of what I have learned to him. A special mention to the ITIS program for providing international students like me a world-class education. Last but not the least, to my friends - Pragyan, Zeon, Saniya, Rohit, Irene and the others who have helped me through good times and bad.

Without you all my masters would not have been a success.





## Contents

<b>Abstract</b>	<b>iii</b>
<b>Acknowledgment</b>	<b>vii</b>
<b>1 Introduction</b>	<b>1</b>
1.0.1 Contributions . . . . .	1
1.0.2 Thesis Outline . . . . .	2
<b>2 Background and Related Work</b>	<b>3</b>
2.1 Twitter . . . . .	3
2.1.1 Retweet . . . . .	3
2.1.2 Mentions . . . . .	3
2.1.3 Hashtags . . . . .	3
2.1.4 Favorite . . . . .	3
2.1.5 Verified User . . . . .	4
2.1.6 Followers . . . . .	4
2.1.7 Following . . . . .	4
2.1.8 Twitter API . . . . .	4
2.2 Credibility of Tweets . . . . .	4
2.3 Definition of Rumor . . . . .	4
2.3.1 Rumor Event . . . . .	5
2.3.2 Time Period of an Event . . . . .	5
2.4 Machine learning . . . . .	6
2.4.1 Machine learning Overview . . . . .	6
2.4.2 Decision tree . . . . .	7
2.4.3 Random Forest (RF) . . . . .	7
2.5 Neural Network . . . . .	8

2.5.1	Convolutional Neural Network . . . . .	10
2.5.2	Recurrent Neural Network . . . . .	10
2.5.3	Long Short-Term Memory . . . . .	13
2.5.4	GRU Cell . . . . .	13
2.5.5	Dropout . . . . .	13
2.6	K Fold Cross Validation . . . . .	14
2.7	Levenberg-Marquardt Method . . . . .	15
2.8	Spark . . . . .	15
2.9	Tensorflow . . . . .	15
2.10	Related Work . . . . .	15
<b>3</b>	<b>Data Collection</b>	<b>19</b>
<b>4</b>	<b>Single Tweet's Creditability Scoring Model</b>	<b>21</b>
4.1	Single Tweet's Creditability Model with handcrafted features . . . . .	21
4.1.1	Features . . . . .	21
4.1.2	Classification Methods . . . . .	23
4.2	Single Tweet's Creditability Model without handcrafted features . . . . .	26
4.2.1	Data Preparing . . . . .	26
4.2.2	Embedding Layer . . . . .	26
4.2.3	Tested Results . . . . .	26
<b>5</b>	<b>Time Series Rumor Detection Model</b>	<b>29</b>
5.1	Dynamic Series-Time Structure (DSTS) . . . . .	29
5.1.1	Time Stamps Generation . . . . .	29
5.1.2	Dynamic Series-Time Structure (DSTS) . . . . .	30
5.2	Features . . . . .	30
5.2.1	Text Features . . . . .	30
5.2.2	Twitter Features . . . . .	30
5.2.3	User Features . . . . .	30
5.2.4	Epidemiological Modeling Features . . . . .	31
5.2.5	SpikeM model Features . . . . .	37
5.2.6	Crowd Wisdom Features . . . . .	43
5.2.7	CreditScore Features . . . . .	43
5.3	Classification Models . . . . .	45
5.4	Feature Evolution . . . . .	46
5.4.1	Time Series VS Static Feature . . . . .	46
5.4.2	Feature Analyzing Over Time . . . . .	47
<b>6</b>	<b>Outlook</b>	<b>49</b>
6.1	Conclusion . . . . .	49
6.2	Future Work . . . . .	51

<b>Bibliography</b>	<b>53</b>
<b>List of Figures</b>	<b>61</b>
<b>List of Tables</b>	<b>63</b>



Twitter is a microblogging service which are used by millions users. Users can publish and exchange information with short tweets within 140 characters. It is cheap and can be accessed through several like website, email or mobile phone. That makes it a ideal media for spreading breaking news and false rumors. A study by the Pew Research Center showed that the people in USA under age of 30 consider Internet as the major resource of news and Internet became the second important media overall(wenxian).

But these advantages make Twitter which became one of the most importance resources of breaking news, at the same time becomes into a ideal media for spreading unverified information. On Twitter everyone can be a journalist and publish news or rumors without any substantiation which must be done by traditional journalists before news' publishing.

Rumor could be defined as a statement whose truth value is unverified or deliberately false(wenxian). And they could be harmful to the government, market and society. One case is some hacked accounts spread a rumor about Obama had been injured in white house. The S&P crashed and wiped off 130 Billion dollars of stock value (wenxian).

So a method of detecting rumors on Twitter that could detect rumors in the early stage of propagation can be very useful.

### **1.0.1 Contributions**

In this thesis, we make the following contributions:

- We develop a model of classification of single tweet with high credibility or low credibility. We call it single tweet's creditability scoring model. Considering it could be set up online for the early rumor events detection in the further and it should response as quick as possible. So we use only the features which can

be extracted from one tweet in the Twitter's interface. We test 2 models. One is random decision forest with handcrafted features. Because the features are limited, it gets only 64% accuracy. Second model what we tested is based on zhou's[39] work, it is a hybrid CNN and LSTM model for text classification. The result of this model we called it credit score.

- We develop a time series model for detecting rumor events. We used Dynamic Series-Time Structure (DSTS)(wenxian) to capture the changes of features over time. And we tested 3 time series model as feature: modified Spike Model(wenxian), SIS model and SEIZ model(wenxian). We add the results of credit scoring model as a feature into time series model to improve the performance in the early stage. And we approved that credit score is one of the best feature in the rumor detection task.
- We study how features changes over time during the spreading of rumors.

### 1.0.2 Thesis Outline

(xiugai) The rest of this these is organized as follows: In Chapter ?? we try to understand the research practices of humanities scholars when working with web archives. We also seek to find problems faced by the scholars when accessing web archives using keyword search. Based on our findings we introduce Historical Query Intents(HQIs) in Chapter ?. In the same chapter we also motivate and develop a novel retrieval model suited for HQIs. The experiments to test the performance of this model against its competitors is covered in the same chapter. The penultimate Chapter ?? details the architecture of the system that was built to demonstrate the efficacy of the retrieval models to users. Finally, in Chapter 6 we add some concluding remarks and describe future work.

## Background and Related Work

### 2.1 Twitter

Twitter is a microblogging service. Now there are more than 140 million active users<sup>1</sup>. User can publish short messages within 140 characters aka tweets.

#### 2.1.1 Retweet

A retweet is re-posting of a tweet by other users. One way of retweet is using RT at the beginning of the tweet which is retweeted. The other way is using "retweet button" which is officially launched by Twitter after 2015. The difference between these two retweet is tweets are retweeted by "retweet button" can't be searched by Twitter's searching interface, but manually retweets with RT keyword can. So in our work retweet means the tweets are retweeted manually. The number of how many times of this tweet has been retweeted is showed behind it.

#### 2.1.2 Mentions

Mentions are in form like "@username" which are added in the text of tweet. The users are mentioned will receive the notification of this tweet on their homepage.

#### 2.1.3 Hashtags

Mentions are in form like "#topic". It means this tweet belongs to some topics.

#### 2.1.4 Favorite

Favorites means how many users like this tweet. It is showed on the interface of Twitter

---

<sup>1</sup><https://blog.Twitter.com/2012/Twitter-turns-six>

### 2.1.5 Verified User

Verified User means this account of public interest is authentic by Twitter. It is showed by a blue icon behind the name of the poster.

### 2.1.6 Followers

The followers of a user are accounts who receive this user's posting. The total number of followers can be seen in the profile of poster.

### 2.1.7 Following

The following are other accounts who follow this user. The total number of following can be seen in the profile of poster.

### 2.1.8 Twitter API

Twitter API is provided by Twitter<sup>2</sup> for developer. But the search API only return a sampling of recent Tweets published in the past 7 days<sup>3</sup>. We need the full stories of the events, so we crawled the data directly from the searching interface<sup>4</sup>.

## 2.2 Credibility of Tweets

Titter has been used for reporting breaking news when emergency events happen like disaster (wexian). But the people are likely to trust the news which post on traditional news website more than the news with same headline but posted on twitter. And Thomson's work shows that different crowds of people trust tweets basing on different kinds of sources[34].

## 2.3 Definition of Rumor

The definition of rumor in our work is unverified information spreading on Twitter over time. It is a set of tweets including the the sources, spreaders, retweets, doubting tweets and the denying tweet.

---

<sup>2</sup><https://dev.twitter.com/overview/api>

<sup>3</sup><https://dev.twitter.com/rest/public/search>

<sup>4</sup><https://twitter.com/search-home>



### 2.3.1 Rumor Event

If a rumor didn't widely spread and it could be harmless. So we more focused on the rumors which are widely spread and contain one or more bursty pikes during propagation. We call it "rumor event".

### 2.3.2 Time Period of an Event

The beginning of a rumor event is hard to define. As far as I know there is related work which give us a approach to define the beginning of one rumor event. One reason is a rumor may be created several years ago and kept existing in Twitter, but it didn't attract the people's attention. However it can be triggered by other events and suddenly quickly spread as a bursty event.

For example, a rumor<sup>5</sup> claimed that Robert Byrd was member of KKK. This rumor has been circulating in Internet for a while, as shown in figure 2.1 that almost every day there are several tweets talking about this rumor. But this rumor was triggered by a picture about Robert Byrd kissing Hillary Clinton in 2016 and twitter users suddenly notice this rumor and this rumor was bursty spread. So what we are really interested in is the hours around the bursty pike.

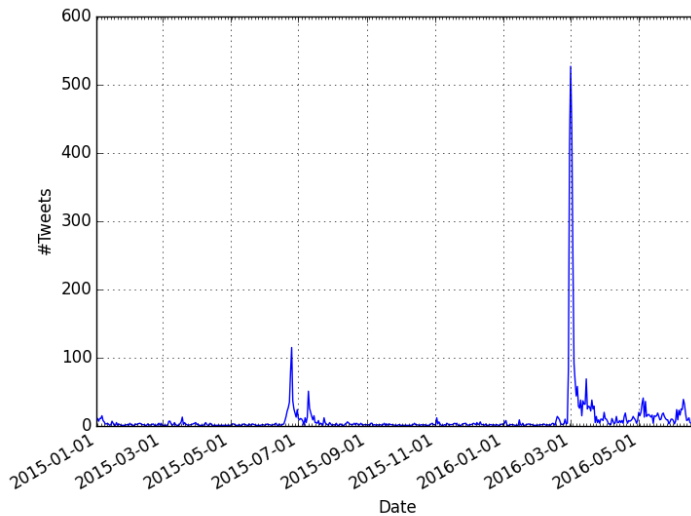


Figure 2.1: large scale tweet volume of event Robert Byrd

We defined the hour which has the most tweet's volume as  $t_{\max}$  and we want to detect the rumor event as soon as possible before its burst, so we define the time of the first tweet before  $t_{\max}$  within 48 hours as the beginning of this rumor event,

<sup>5</sup><http://www.snopes.com/robert-byrd-kkk-photo/>

marked as  $t_0$ . And the end time of events we defined as  $t_{\text{end}} = t_0 + 48$ . We show the tweet volume in figure 2.2 of the above rumor events after defined 48 hours time period.

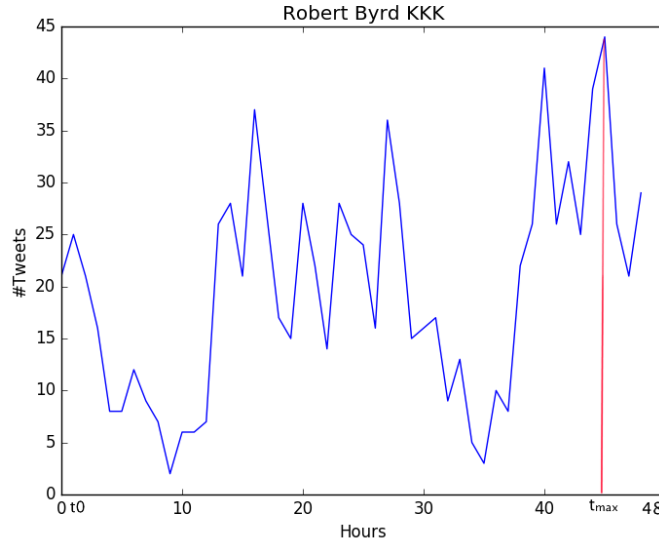


Figure 2.2: tweet volume of the rumor event of Robert Byrd after selected time period

## 2.4 Machine learning

### 2.4.1 Machine learning Overview

Machine learning covers vast numbers of algorithms and has been successful applied in different field. The challenges of ML are finding the best model which is suitable for this task, fitting the parameters and selecting the features.

Normally we split the ML methods into Supervised learning, Unsupervised learning and Reinforcement learning[30].

The supervised learning is the most popular method of ML. It needs a set of inputs and a set of desired outputs. And the algorithm will learn to produce the correct output based on the new input.

The unsupervised learning needs a set of inputs but no outputs. The algorithm will generate the outputs like clusters or patterns. The unsupervised learning task can be used for example when people can't label the outputs.

### 2.4.2 Decision tree

A decision tree uses a tree model of decisions and their possible consequences as shown in figure ?? . DT is an non-parametric supervised learning method used for classification and regression. Because the trees structure is very simple to visualize, the result of DT is easy to understand. But Decision trees are unstable [6]. But DT is the basic unit of random forest.

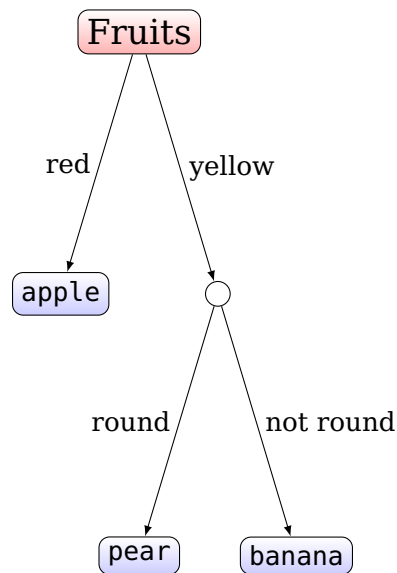


Figure 2.3: Decision tree

### 2.4.3 Random Forest (RF)

Classification is a supervised data mining technique. Our work can be considered as a classification task. And the classification model is random forest.

RF is an algorithm of supervised learning which developed by Leo Breiman [7]. It's built by a set of classification trees [8]. Each tree is trained by a small bootstrap sample of training set and while prediction each tree votes one single candidate. RF generates the result of prediction By taking the majority vote.

For example we got task to classify pears and apples. We have the features are whether the fruit round, whether the fruit has seeds, whether the fruit red and whether the fruit is juicy. We build up a 3 trees random forest as the graph 2.4. Tree 1 randomly takes the subset of the features red and seeds and votes to apple. Tree 2 randomly takes the subset of the features red and seeds and votes to pear. Tree 3 randomly takes the subset of the features round and red and votes to apple. The most vote is apple, so the output of RF is apple.

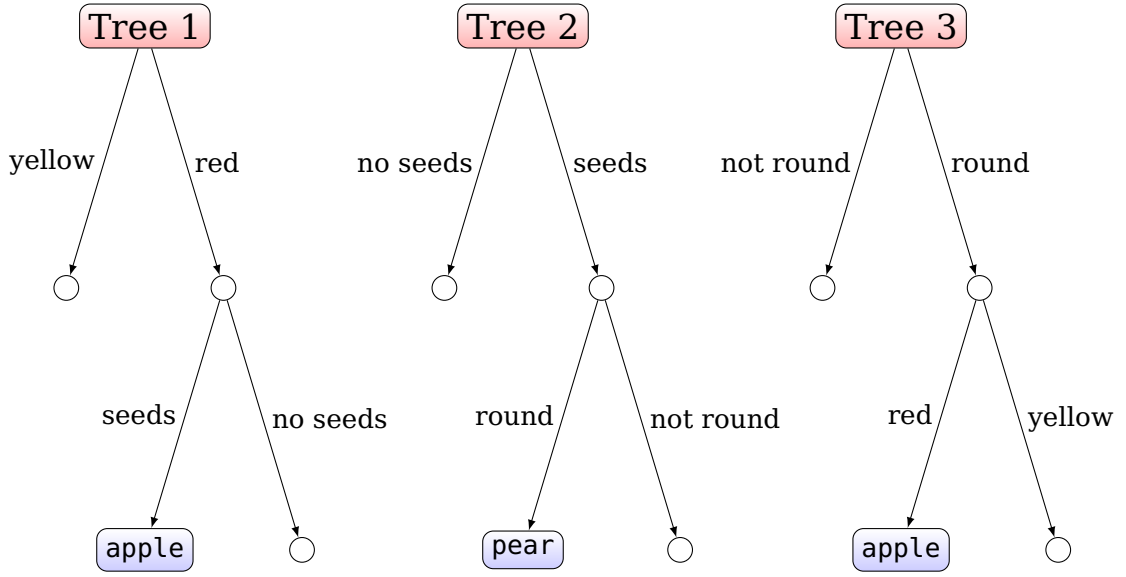


Figure 2.4: An example of random forest with 3 trees

Because RF uses a random subset of features instead of the best features in every node, so it can avoid the overfitting[7].

Another benefit of RF is that it can return the **features' importance**. RF is built up by a subset of training data. If there are  $N$  trees, RF will take  $N$  times bootstrap sampling. So some features we didn't select to use for training the model. The data we didn't selected we call it out-of-bag (OOB) data. In the above example the feature whether the fruit is juicy didn't join the construction the Trees, so it is a OOB data. These data can be use for validating the model and the output is OOB error  $E_{\text{oob}}(G) = \frac{1}{N} \sum_{n=1}^N \text{err}(y_n, G_n^-(X_n))$  with  $X_n$  are features only in OOB. At last we get the importance of feature by permuting one feature to random numbers.  $\text{importance}(i) = E_{\text{oob}}(G) - E_{\text{oob}}^p(G)$ ; where  $E_{\text{oob}}^p(G)$  is the OOB error after permuting a feature. The more performance drop down, that means the more important this feature is. We use this method to measure the performance of features and evolute the models later.

## 2.5 Neural Network

Figure 2.5 is a simple model of forwarding neural network which is a computing systems that processes information by their dynamic state response to external inputs [39]. In the figure 2.5 the round nodes are called neurons. Each neuron has a active function like sigmod or tanh 2.6 Neurons are connected each others with weight. The process of training the model is actually the process learning the weights of the connections of the network. With the back-propagation algorithm we 1) compute

the loss from the network output and 2) update the weights by passing the error backwards through the network [29].

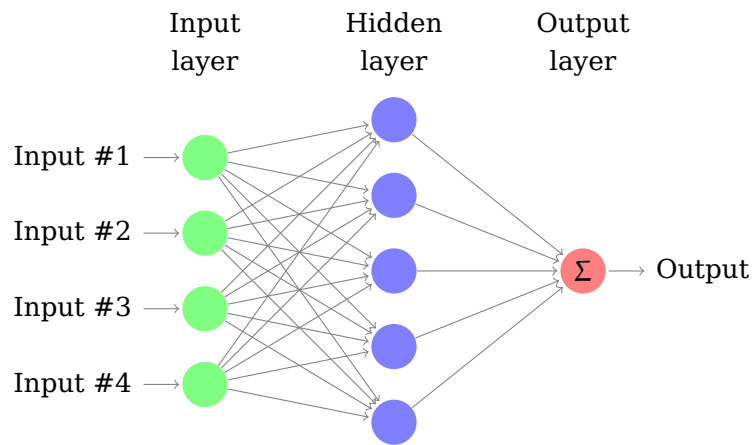


Figure 2.5: An example of 1 layer neural network

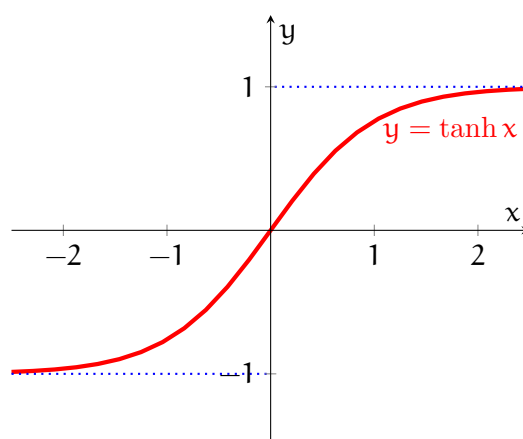


Figure 2.6: tanh function

### 2.5.1 Convolutional Neural Network

CNN is a kind of neural network which contains a convolutional layers. First CNN is developed by LeCun et al. developed CNNs [20] for computer vision to resolve the problem of images' shift or rotation. With the convolutional layer CNN make that is possible to learn from the local features. Figure 2.7 is an example p of CNN in application of text classification. The word "wait" and its following word "for" are mapped together into different classes though different convolutional filters. That makes convolutional neural network does not learn from the single word form sentence but from the word and its context. That makes sense because one word can in different context have different meaning. For example "Apple Tree" and "Apple Company", the same apple but contains different meaning. Pooling layer often follows after the convolutional layers in CNN models whose mainly task is to down-sample from the output of convolutional layers.

### 2.5.2 Recurrent Neural Network

A Recurrent neural network (RNN) is one of of class of neural network. The different between RNN and other neural network is there are feedback connections between the hidden layers showing as figure 2.8, that makes the input from the past also can influence the network. The activation status of the hidden layer neurons presents as an internal state which can be considered as a kind of memory. That makes recurrent neural network has the ability to process the sequence input. Humen understand the words of a sentence based on the previous words, the feed forwarding neural network will ignore the previous inputs.

There are several types of applications of RNN. First one is multi-input and single-output 2.9, application in for example text classification. Second is single-input

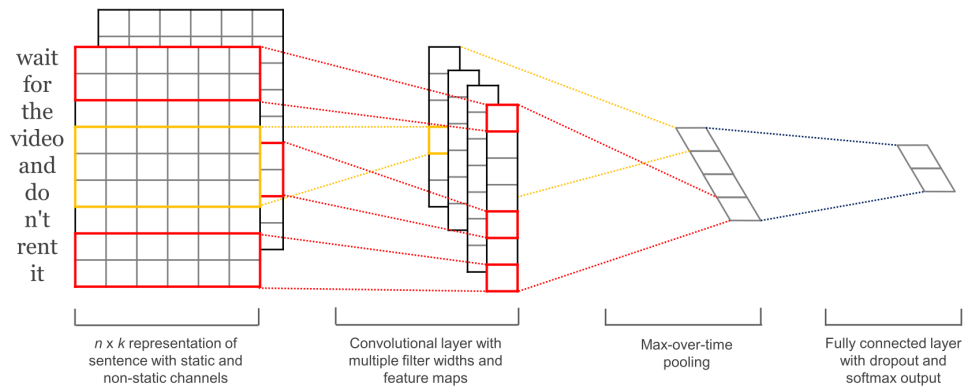


Figure 2.7: CNN for Text Classification (source from[15])

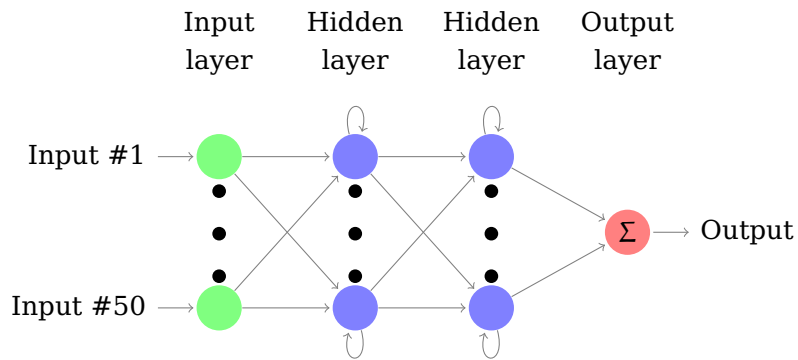


Figure 2.8: 2 layer Recurrent neural network

and multi-output [2.10](#), application in for example input an image and generate a sentence for description. Third is multi-input and multi-output [2.11](#), application in for example the translation system.

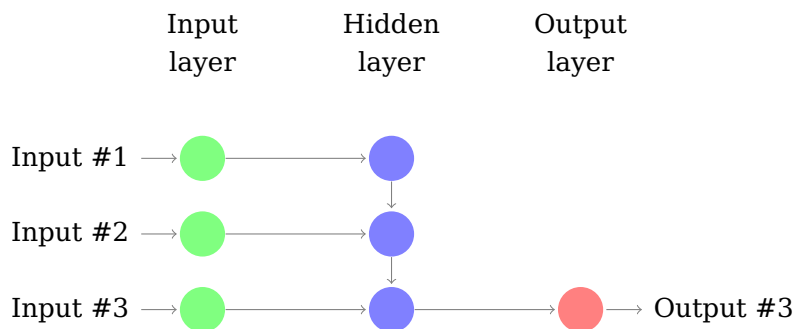


Figure 2.9: multi-input single-output Recurrent neural network

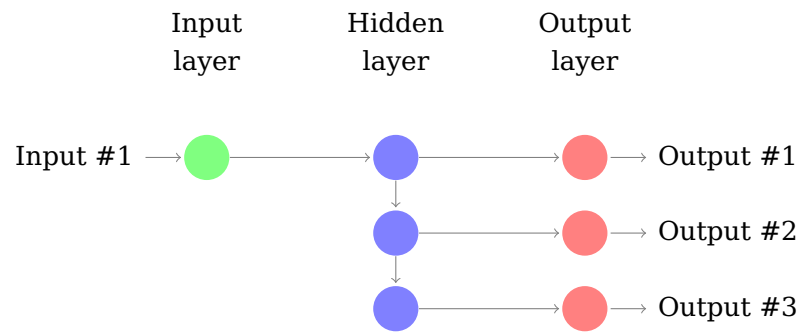


Figure 2.10: single-input multi-output Recurrent neural network

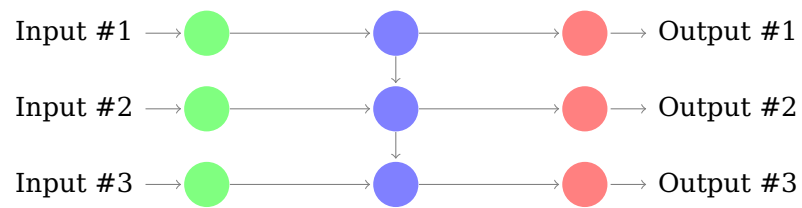


Figure 2.11: multi-input multi-output Recurrent neural network



### 2.5.3 Long Short-Term Memory

But recurrent neural networks are difficult to train for a long sequence input, because there are two problems "**gradient explode**" and "**gradient vanish**" [28]. Gradient vanish means that the gradient of the earlier inputs goes exponentially fast towards zero. That makes recurrent neural networks can't learn the long-term dependencies in sequences. Gradient explode means the gradients become exponentially large while training.

People can change the active function like ReLU to avoid the vanishing gradient problem, but more popular solution is using Long Short-Term Memory which is first invented by Hochreiter and Schmidhuber [12]. An LSTM memory cell contains self-connected memory cell with its a memory cell, an input gate, an output gate and a forget gate units. LSTM uses a series gate units to control the cell receiving data flow, extracting data flow or forgetting the current state at each time step.

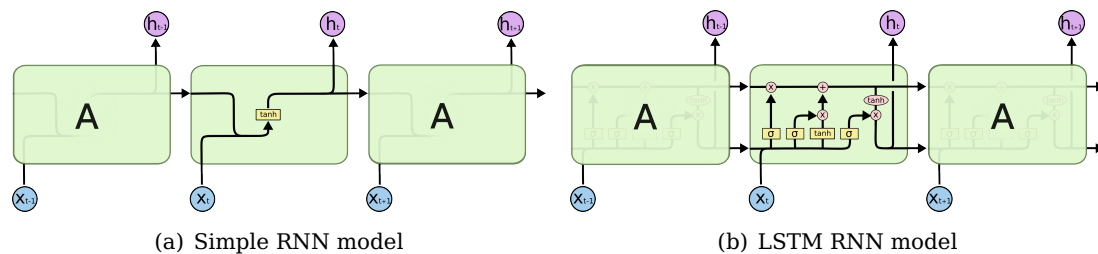


Figure 2.12: Two models of RNNs (a) Normal model of RNN with tanh Unit (b) RNN with LSTM Cells (source [27])

### 2.5.4 GRU Cell

GRU Cell is a special case of LSTM which is invited by Cho, et al [10]. The forget and input gates are merged into a update gate. The cell state and hidden state are combined as one state. These changes make the structure of GRU simpler than LSTM.

### 2.5.5 Dropout

To avoid the overfitting we used dropout method[?]. The main idea of dropout is to dropout some units in hidden layer in neural network while training. When training the model some of the units are removed temporarily with probability  $p$  as figure 2.14. But while testing the dropout units still join the test. The aim of dropout is to avoid one best state keeping exist. By disabling some units of the network dropout makes

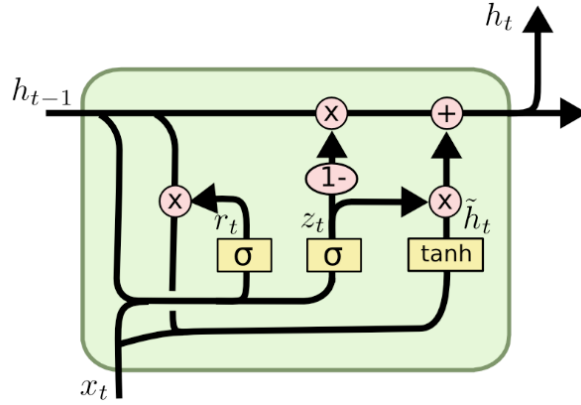


Figure 2.13: GRU Cell (source [27])

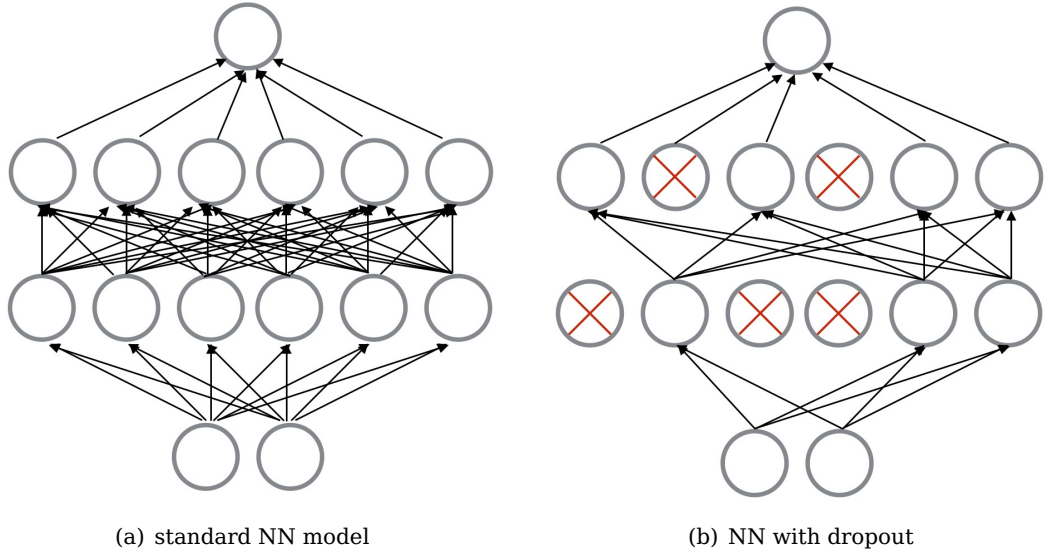


Figure 2.14: Dropout Neural Net Model. (a) a standard neural network (b) after deploying dropout

## 2.6 K Fold Cross Validation

To limit the overfitting and improve the robustness of the model we often use cross validation technology. K fold cross validation is that the dataset was shuffled and split into k subsets. We train the model with some of subsets and others are used for testing.

In our work we split the data not equally, we consider an event as a unit. We use 10 times cross validation, so we split 260 event into 10 subsets. And we test all models with the same shuffled sequence, aka same training sets and test sets.

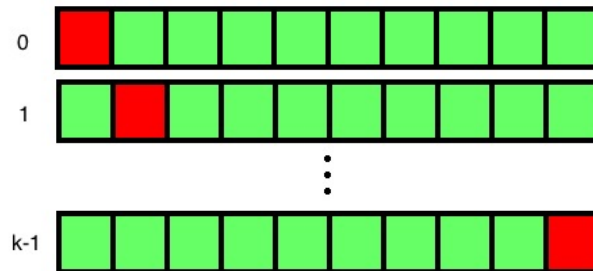


Figure 2.15: K Fold Cross Validation (green rectangle is the training set, red rectangle is the testing set)

## 2.7 Levenberg-Marquardt Method

## 2.8 Spark

## 2.9 Tensorflow

## 2.10 Related Work

People research on Twitter for a long time and there are a lots directions to study this complex social network like event detection[16], spam detection [1] [35] or sentiment detection[4]. These work are similar with our task rumor detection but still contain many differences.

People first studied on rumors in psychology area for years [2] [32]. But Twitter becomes an important platform for people sharing and exchanging information including rumors, the detection of misinformation on twitter turns to be a trending researching topic.

Researcher first began from studying rumor spreading during several special events like natural disasters[26] [33] or terrorist attacks [31]. But these results are not general enough to other types of rumors. Our data includes sports, politic, entertainment, also natural disasters and other kinds of rumors.

Carlos Castillo researched the information credibility on Twitter[25] [9][11]. But his work is more about people opinion (trustful or not) to a tweet not the credibility of tweet itself. But it still a good start of resolving problem. Lots of other works are based on Castillo's work [37] [21] and used different set of features, for example yang [37] add the location of poster and the client type of user (web client or mobile phone) as two important features. And

There are rencent researches based on the propagation structure of rumor on Sina Weibo[36] [3], but Twitter doesn't give us such detail of propagation like Weibo, so these works are useless for us.

Liu's work [36] focused on one type of rumor "the false photos" on social network not the general type of rumors.

Xiaomo claimed their system is the first real time rumor debunking system on Twitter [21]. But their work are similar with above other previous works, they use the static features of rumors and ignored the feature's changing over time during the event spreading on the social network which can be seen in figure 2.16 and 2.17. The fraction of tweets containing url with top 5000 domain and the fraction of poster living in large city which are both constantly changing. But one feature called "crowd wisdom" is an new idea and we use it in our work.

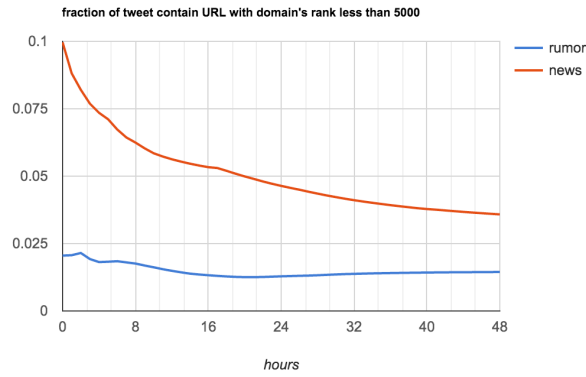


Figure 2.16: The fraction of tweets containing url with top 5000 domain

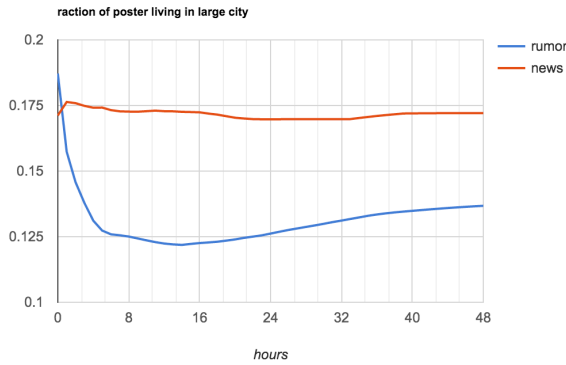


Figure 2.17: The fraction of poster living in large city

Other researchers used propagation models like SpikeM [18] and SEIZ models [13] to capture the tweet's volume changing over time, but they didn't use other features into time series model. And these features can't show the effect in the early stage of rumor spreading. We can provide it in the last section(wenxian).

J Ma used Recurrent Neural Networks for rumor detection [22], as the same disadvantage of the other Neural Networks models, the process of model's training is black box to us, so we can't know way we got this result. In our case we can't get the evidences to convince people why the detected event is more likely to be a rumor or rather than a breaking news.

Another work from J Ma used a time series model called Dynamic Series-Time Structure [23] to capture the variation of features. But they only used the features of social content and they didn't further explain how the performances of features change over time or how do they improve the performance of the model in first few hours after the events beginning. We use their Series-Time Structure with our extended data and we will show how do these features change during time.



We collected rumor stories from a rumor tracking website **snopes.com**. We crawled 4300 stories from the website and we manually constructed 130 queries. The approach of constructing queries is mainly following the work(wenxian). The regular expression of a query is (Object&Subject&Description(Description1||Description2|...)) for example a story is about Obama removing a flag in pearl harbor. Object is Obama, subject is flag and its synonym like flags, flagpole. Description is remove and its synonym removes, removed, removal, removing, "token down" or a url about this rumor "Departed.co". In this case there is a proper noun "pearl harbor" is also useful. Finally we transfer the regex to Twitter's query: obama (flag OR flags OR flagpole) (remove OR removes OR removed OR removal OR removing OR "token down" OR "Departed.co") pearl harbor.

We collect news stories for the data of (wenxian). We crawled average 50% of them and we selected events with the most coverage and minimal 100 tweets.

After we crawled and parsed the whole timeline of an events. We detect the 48 hours time period of the burst in the way we mentioned in section 2.3.2. We crawled the homepage of poster of the tweets in the event time period total 133,396 users. We also extracted 11,038 domains which are contained in the tweets in the 48 hours time period and we crawled these domains' catalogs in bluecoat.com<sup>1</sup>, ranks in alexa.com<sup>2</sup> and WOT score in wot.com<sup>3</sup>.

We use Beautiful Soup as the html parsing library to parse the Twitter timeline pages and the users' homepage<sup>4</sup>. Beautiful Soup is a Python library for pulling data out of HTML and XML files. For increasing the speed of parsing html and extracting features from raw data, we use Spark<sup>5</sup> technology, because it can simply

<sup>1</sup><http://sitereview.bluecoat.com/sitereview.jsp#/?search=bbc.com>

<sup>2</sup><http://www.alexa.com/siteinfo/bbc.com>

<sup>3</sup><https://www.mywot.com/en/api>

<sup>4</sup><https://www.crummy.com/software/BeautifulSoup/bs4/doc/>

<sup>5</sup><http://spark.apache.org/>

manage multithread and its mapReduce in memory technology makes the process much faster.



## Single Tweet's Creditability Scoring Model

Most of the pervious researches are focusing on event level rumors and claims taht the task of classification for individual tweet is not reliable [21] [23] [38], because a single tweet is to short to contains engough features. But Carlos Castillo designed a single tweet's creditability rank system Tweetcred [11]. So we think it is still enable to build up a single tweet creditability model. And we were inspired by the (wenxian) work that we may use nn() without handcrafted features to build up the single tweet's creditability model which in later experiments outputs a better performance.

### 4.1 Single Tweet's Creditability Model with handcrafted features

First we follow Castillo's [11] idea to implement a single tweet's creditability model with handcrafted features and we test it with decision trees, decision forest and SVM. But the performance of each model are not good enough.

#### 4.1.1 Features

We use a collection of features major from Castillo's Tweetcred system[11] totally 27 features in table 4.1. These features can be extracted directly from Twitter interface without third part website.

##### Text Features

The Text features capture the content of the text of the tweets. There are 16 Text features. **Sentiment Features** are included in text features. We used the

python natural language Toolkit (nltk) <sup>1</sup> to analyze the tweets' sentiment and extract the features: the NumPositiveWords, NumNegativeWords and Polarityscores. Polarity scores is a float for sentiment strength of one tweet<sup>2</sup>  $\text{Polarity\_scores} = \frac{1}{N} \sum_0^N \text{Polarity}(\text{token}_n)$ .

### User Features

We selected total 5 features of the poster. These features are extracted directly from the twitter interface as in figure 4.1. ReputationScore is defined as the ratio between #Friends over # Followers.  $\text{ReputationScore} = \frac{\# \text{Friends}}{\# \text{Friends} + \# \text{Followers}}$ .



Figure 4.1: Sample of Users' Information on Twitter Interface

### Twitter Features

Twitter Features are the features of twitter's special functions. It includes Hash-tag, Mention, Number of URLs, Number of retweeted and whether this tweet is retweet (contains RT as keywords or "QuoteTweet-innerContainer u-cf js-permalink js-media-container" as the CSS class of the tweet in html).

<sup>1</sup><http://www.nltk.org/>

<sup>2</sup><http://www.nltk.org/api/nltk.sentiment.html>

#### 4.1 Single Tweet's Creditability Model with handcrafted features

Category	Feature	Description
Twitter Features	Hashtag	Whether the tweet contains #hashtag
	Mention	Whether the tweet mentions others @user
	NumUrls	# url in the tweet
	Retweets	How many times this tweet has been retweeted
	Isretweet	Whether this tweet is retweeted from others
Text Features	LengthofTweet	The length of tweet
	NumofChar	# of individual characters
	Capital	Fraction of characters in Uppercase
	Smile	Whether this tweet contains :->, :-), ;->, ;-)
	Sad	Whether this tweet contains :-<, :-(, ;->, ;-(
	NumPositiveWords	# of positive words
	NumNegativeWords	# of negative words
	Polarityscores	polarity scores of the Tweet
	Via	Whether this tweet contains via
	Stock	Whether this tweet contains \$
	Question	Whether this tweet contains ?
	Exclamation	Whether this tweet contains !
	QuestionExclamation	Whether this tweet contains multi Question or Exclamation mark
	I	Whether this tweet contains first pronoun like I, my, mine, we, our
	You	Whether this tweet contains second pronoun like U, you, your, yours
	HeShe	Whether this tweet contains third pronoun like he, she, they, his, etc.
User Features	Followers	# of followers
	Friends	# of friends
	Tweets	# of tweets which are posted by this user
	Description	Whether this user has description
	Verified	Whether this user is a verified user
	ReputationScore	Ratio between #Friends over (# Followers + #Friends)

Table 4.1: Features for Single Tweet's Creditability Scoring Model

#### 4.1.2 Classification Methods

Most of existing researches use SVM, Random Forest and DT. We test our data also with these 3 models. We implement these 3 model with scikit-learn library<sup>3</sup>. We shuffled the 260 events and split them into 10 subset, we uses them for 10 times cross-validation. We show the parameters after optimization for each model them in table 4.2.

Model	Parameters	Value
Random Forest	Number of Trees	200
SVM	kernel	radial basis function
	penalty parameter of the error term	2.0
	gamma	$\frac{1}{27}$
Decision Trees	criterion	gini

Table 4.2: Parameters of Classification models

<sup>3</sup>scikit-learn.org/

We show the results in the table 4.3. The best model with the highest accuracy is RF, but it reaches only 64.87% and the other two models are even worse, so it is clear to see with manually handcrafted features, one single tweet is difficulty to be classified.

<b>Model</b>	<b>Accuracy</b>
Random Forest	<b>0.6487</b>
SVM	0.5802
Decision Trees	0.5774

Table 4.3: Prediction Accuracy of Different Single Tweet's Creditability Scoring Models

And we rank the features using the features importance which we mentioned in section 2.4.3, showing in table 4.4. The best feature is polarity scores of sentiment. It means that there is a big bias between the rumors tweets and the tweets real events. It was mentioned by previous work [2] where he gathered a large rumors collection during WW2 which are printed in the Boston Herald's Rumor Clinic. He summarized rumors as several types: pipe-dream, fear and aggression. The most researches believe that rumors mostly contain negative sentiment and polarity [32][17]. In our study average polarity score of news event is -0.066 and average polarity score of rumors is -0.1393 ,it means that rumors contain more negative sentiment.

And we usually think the verified users may have less possibility to be involved in the rumors' spreading, but the result shows that the verified users may be not really trustful like we thought. And "Isretweeted" feature is neither a good feature which means the possibility of people retweeting the rumors or true events are similar.

Feature	Feature Importance
PolarityScores	0.1460686474
Capital	0.09638447209
LengthofTweet	0.09283739724
Tweets	0.08750049577
Friends	0.08065591431
ReputationScore	0.08002109553
Followers	0.07938657292
NumOfChar	0.07659755102
Stock	0.04920394972
NumNegativeWords	0.03068379335
Exclamation	0.02304551015
numUrls	0.02124370609
NumPositiveWords	0.01976939973
Hashtag	0.01851408745
Mention	0.01596532677
Question	0.01486070376
Retweets	0.01349486577
I	0.0109471116
You	0.00998103276
HeShe	0.00774915859
UserDescription	0.007402174886
Via	0.005545157727
QuestionExclamation	0.005422123705
Isretweet	0.003240079497
Verified	0.003081752983
Smile	0.0003979192278
Sad	0

Table 4.4: Features Importance

## 4.2 Single Tweet's Creditability Model without handcrafted features

As we showed in last section. The result of Single Tweet's Creditability Model with handcrafted features no matter with SVM, DT or RF is not so convincing. We may miss some important features we didn't mine out of the tweet.

Inspired by the Lai and J Ma's works [19] [22] we test neural networks as the classifier which does not need to extract features from the data. Based on the previous work we tested it with (wenxian) models: Basic tanh-RNN [4.2\(a\)](#), 1-layer GRU-RNN [4.2\(b\)](#), 1-layer LSTM [4.2\(c\)](#), 2-layer GRU-RNN [4.2\(d\)](#), Fasttext [4.2\(e\)](#) and CNN+LSTM [4.2\(f\)](#) model as figure [4.2\(d\)](#) model. Basic tanh-RNN, 1-layer GRU-RNN, 1-layer LSTM-RNN and 2-layer GRU-RNN models are based on the work of Lai and J Ma's works [22] [19]. Fasttext comes from Joulin's work [14] which is a fast text classification model no need of GPU acceleration. The hybrid model of CNN and LSTM is Zhou's idea [39].

### 4.2.1 Data Preparing

We test the above models with same dataset and same shuffled sequence as the RF, SVM model. The difference is that we feed only the text of each tweet to the neural network models. And we randomly take 20% data from the training set to be the valid set for optimizing the parameters.

### 4.2.2 Embedding Layer

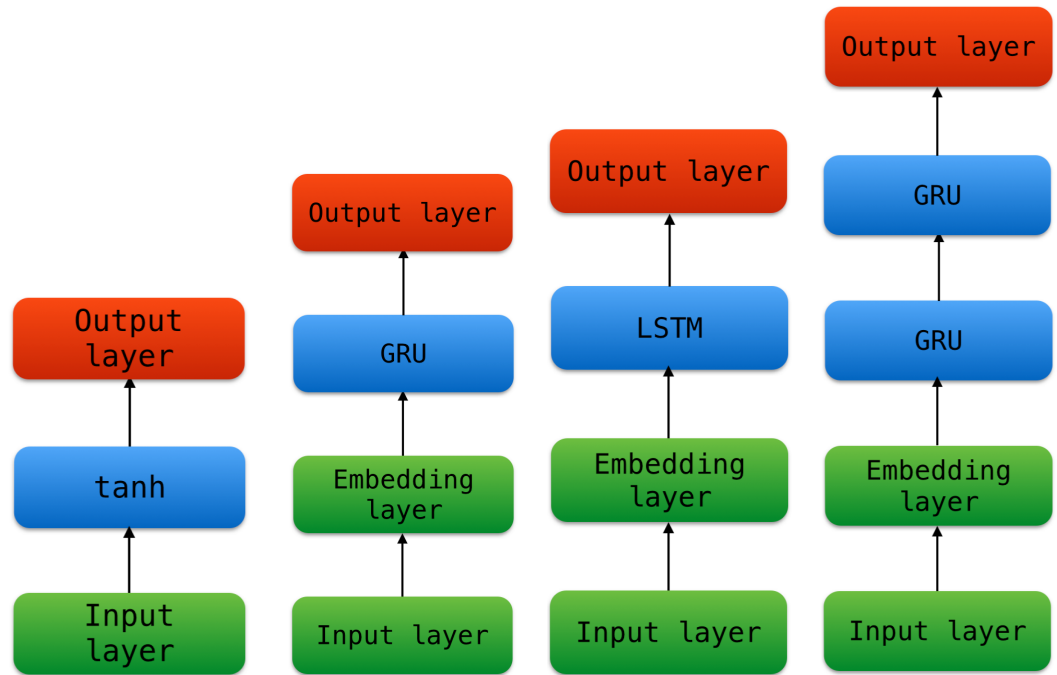
The first layer is embedding layer which is set up the same to all models. The embedding size is 50. The output of the embedding layer is the vectors representing the words.

### 4.2.3 Tested Results

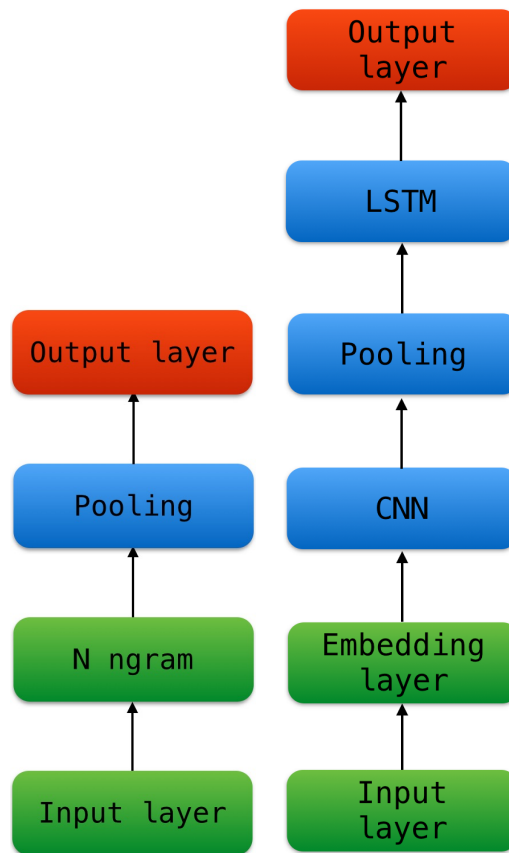
The best accuracy result is CNN+LSTM model as shown in table [4.5](#). To avoid overfitting we use 10-fold cross validation and dropout technology.

Model	Accuracy
CNN+RNN	<b>0.8119</b>
2-layer GRU	0.7891
GRU	0.7644
LSTM	0.7493
Basic RNN with tanh	0.7291
FastText	0.6602

Table 4.5: Prediction Accuracy of Different Single Tweet's Creditability Scoring Models



(a) Sample RNN Model (b) 1-layer GRU-RNN (c) 1-layer LSTM-RNN (d) 2-layer GRU-RNN



(e) FastText (f) Hybrid CNN + LSTM

Figure 4.2: neural network model for single tweet classification 1



## Time Series Rumor Detection Model

As we showed in figures 2.16 and 2.17, the features change during the the events' spreading over time. In order to capture these changes of the each features we use Dynamic Series-Time Structure (DSTS) which was presented in Ma's work [23]. In the an Event  $E_i$  there is a set of tweets  $tw_{ij}$  and we split them into different time intervals according to the creation time so that we can analyze the features in time series. We test the different classifiers with this model, we compare it with static features and in the end we rank all features and show their performance over time.

### 5.1 Dynamic Series-Time Structure (DSTS)

#### 5.1.1 Time Stamps Generation

For an event  $E_i$  we define  $timeFirst_i$  as the start time of the event,  $timeLast_i$  as the time of last tweet of the event. We split the each tweet  $tw_{ij}$  into  $N$  time intervals according to the creation time. The length of each time interval we define as follow:

$$Interval(E_i) = \frac{\lceil (timeLast_i - timeFirst_i) \rceil}{N} \quad (5.1)$$

And the index of time interval  $TS(t_{ij})$  where a tweet  $tw_{ij}$  which is created in time  $t_{ij}$  should fall into, we define as follow :

$$TS(t_{ij}) = \frac{\lfloor (t_{ij} - timeFirst_i) \rfloor}{Interval(E_i)} \quad (5.2)$$

In our work  $Interval(E_i)$  as we defined in section 2.3.2 is one hour and  $N$  is constant 48 hours for each event.

### 5.1.2 Dynamic Series-Time Structure (DSTS)

Now we have all the time intervals of an event  $E_i$  and we can generate a vector  $V(E_i)$  of features in each time interval. And in order to capture the changes of feature over time we should not only model the features in individual time intervals but also we should model their difference between two time intervals. So the model of DSTS is represented as:

$$V(E_i) = (\mathbf{F}_{i,0}^D, \mathbf{F}_{i,1}^D, \dots, \mathbf{F}_{i,N}^D, \mathbf{S}_{i,1}^D, \dots, \mathbf{S}_{i,N}^D) \quad (5.3)$$

where the  $\mathbf{F}_{i,t}^D$  is the feature vector in time interval  $t$  of event  $E_i$ .  $\mathbf{S}_{i,t}^D$  is the difference between  $\mathbf{F}_{i,t}^D$  and  $\mathbf{F}_{i,t+1}^D$ .  $V(E_i)$  is the time series feature vector of the event  $E_i$ .

$$\mathbf{F}_{i,t}^D = (\tilde{f}_{i,t,1}, \tilde{f}_{i,t,2}, \dots, \tilde{f}_{i,t,D}) \quad (5.4)$$

$$\mathbf{S}_{i,t}^D = \frac{\mathbf{F}_{i,t+1}^D - \mathbf{F}_{i,t}^D}{\text{Interval}(E_i)} \quad (5.5)$$

We use Z-score to normalize feature values which is implemented by sklearn.

$$\tilde{f}_{i,t,k} = \frac{f_{i,t+1,k} - \bar{f}_{i,k}}{\sigma(f_{i,k})} \quad (5.6)$$

where  $f_{i,t,k}$  is the  $k$ -th feature in time interval  $t$  of the event  $E_i$  in time interval  $t$ .  $\bar{f}_{i,k}$  is the mean of the feature  $k$  of the event  $E_i$  and  $\sigma(f_{i,k})$  is the standard deviation of the feature  $k$  over all time intervals. We skip this step when we use random forest and Decision Trees because they do not need feature normalization.

## 5.2 Features

We use a collection of features based on previous works [9][11][37][21][22][25][23][36][13]. We extracted totally 50 features in table 5.4. These features are not only extracted from Twitter interface but also other website like bluecoat.com we mentioned them in section 3.

### 5.2.1 Text Features

### 5.2.2 Twitter Features

### 5.2.3 User Features

The list of large city <sup>1</sup>

---

<sup>1</sup><http://www.demographia.com/db-worldua.pdf>

### 5.2.4 Epidemiological Modeling Features

Jin's work is as far as we know the first people using epidemiological model to analyze rumors' prorogation on twitter [13]. They fits the volume of the rumors and news events into two models SIS (Susceptible, Infected, Susceptible) and SEIZ (susceptible, exposed, infected, skeptic).

**SIS** is one of the most popular epidemiological model. To adapt to the scenario of Twitter, we define a user who posts a tweet of relevant event as **(I)** infected, a user who didn't we define as **(S)** susceptible. But unlikely as a normal epidemiological scenario infected nodes can be cured and return to be susceptible, the user once posts a tweet of the certain events, he will be classified into the infected component forever. He can't be return susceptible class. At time  $t$  the total number of population is  $\Delta N(t) = I(t) + S(t)$  where  $I(t)$  is the size of infected population and  $S(t)$  is the size of susceptible population. As shown in Figure 5.1, SIS model works as follow:

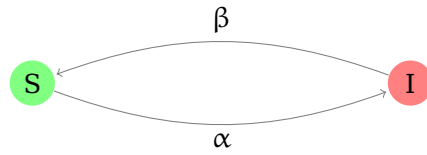


Figure 5.1: SIS Model

- A user who posts tweets about the certain event is regarded as infected.
- A susceptible user has not tweeted about the certain event
- A susceptible user may see the a tweet about the certain event from a infected users and he immediately retweets or posts a tweet about this events, and in that he turns himself to infected.
- Susceptible user will remain susceptible until he contacts (via tweet) with infected person.

we show SIS model mathematical as follow:

$$\frac{d[S]}{dt} = -\beta SI + \alpha I \quad (5.7)$$

$$\frac{d[I]}{dt} = \beta SI - \alpha I \quad (5.8)$$

SIS model assumes that a susceptible user once exposed to a infected user turns to infected immediately. That is one reason of this model why it didn't fit to Twitter. If fact when twitter users see a tweet they have their normal senesce to judgment

the truth of the information and they can decide whether further spreading the tweet or ignoring them.

Another popular model is SIR which contains one more term than the SIS. The definitions of **(S)** and **(I)** are the same of SIS but the term **(R)** stands for recover. Once a susceptible user is recovered, he will be removed from the susceptible component and he can't be infected again. But we can't get a reasonable explanation of the term R if we model the an event spreading on Twitter.

Because of the Shortcomings of above two model, they test another model called SEIZ which reference from [5] .

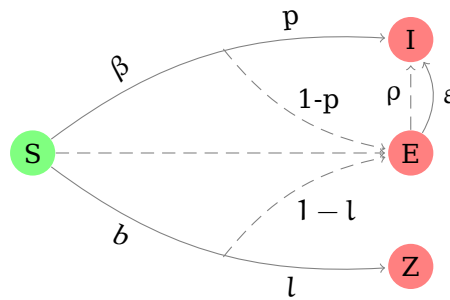


Figure 5.2: SEIZ Model

To adapt to Twitter context, the compartments of the SEIZ model can be mapped like this: **(S)** Susceptible is a user who has not been exposed to the event aka he didn't see any tweets about the certain event yet, **(I)** infected means a user has posted tweets about the certain events, **(Z)** skeptic is a user who has been exposed to the certain event but he decides to ignore it and **(E)** exposed is a user who has been exposed to the certain event but he will post the tweets after some delay.

We show the model in figure 5.2. And the SEIZ works as follow:

- People recruit from **(S)** Susceptible compartment to Skeptics with rate  $b$ . But with probability  $l$  some of them directly deny the events and turn to **(Z)** skeptic compartments. Others with probability  $1-l$  probability turn to **(E)** exposed compartment.
- People recruit from **(S)** Susceptible compartment to Infected with rate  $\beta$ . But with probability  $p$  some of them directly believe the events and repost it and turn them to be **(Z)** skeptic compartments. Others with probability  $1-p$  probability turn to **(E)** exposed compartment.
- People from **(E)** exposed compartment have  $\rho$  probability contacting again with the Infected and turn them to **(I)** infected compartment. And others have  $\varepsilon$  probability turn into **(I)** infected compartment by themselves for example external shock.

And we show the model mathematical like: we show SEIZ model mathematical as follow:

$$\frac{d[S]}{dt} = -\beta S \frac{I}{N} - bS \frac{Z}{N} \quad (5.9)$$

$$\frac{d[E]}{dt} = (1-p)\beta S \frac{I}{N} + (1-l)bS \frac{I}{N} - \rho S \frac{Z}{N} - \varepsilon E \quad (5.10)$$

$$\frac{d[I]}{dt} = p\beta S \frac{I}{N} + \rho S \frac{Z}{N} + \varepsilon E \quad (5.11)$$

$$\frac{d[Z]}{dt} = lbS \frac{Z}{N} \quad (5.12)$$

Symbol	Definition
$\beta$	S-I contact rate
$b$	S-Z contact rate
$\rho$	E-I contact rate
$\varepsilon$	Incubation rate
$1/\varepsilon$	Average Incubation Time
$bl$	Effective rate of S $\rightarrow$ Z
$\beta\rho$	Effective rate of S $\rightarrow$ I
$b(1-l)$	Effective rate of S $\rightarrow$ E via contact with Z
$\beta(1-p)$	Effective rate of S $\rightarrow$ E via contact with I
$l$	S $\rightarrow$ Z Probability given contact with skeptics
$1-l$	S $\rightarrow$ E Probability given contact with skeptics
$p$	S $\rightarrow$ I Probability given contact with adopters
$1-p$	S $\rightarrow$ E Probability given contact with adopters

Table 5.1: Parameters of SEIZ

The author presents an index of SEIZ called  $R_{SI}$  as equation 5.13. It contains all rate values of SEIZ and related to the flux ratio of the **(E)** exposed compartment, the ratio of entering **(E)** to leaving **(E)**. If  $R_{SI}$  is bigger than 1 means the influx of exposed compartment is bigger than the efflux. This index may be a good candidate of feature to analyze rumor spreading on Twitter.

$$R_{SI} = \frac{(1-p)\beta + (1-l)b}{\rho + \varepsilon} \quad (5.13)$$

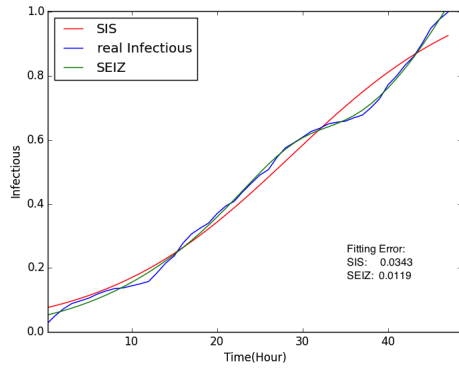
We use Levenberg-Marquard algorithm which we present in section 2.7 to learn the parameters of the SIS and SPEI. The fitting data is the tweet volume of the 260 events (130 rumors and 130 news). In time each interval from  $t_0$  to  $t_n$ , we fit the sequenced tweet volume from the beginning time the  $t_0$  to the current time interval  $t_n$  of an event to SIS and SEIZ model and learn. From SIS we get two feature  $\beta_n, \alpha_n$  and from SEIZ we get 7 features  $\beta_n, b_n, l_n, p_n, \varepsilon_n, \rho_n, RSI_n$ . We add them into our DSTS.

FittingFunction<sub>SIS</sub>(TweetVolume<sub>0</sub>, ..., TweetVolume<sub>n</sub>) $\rightarrow \beta_n, \alpha_n$

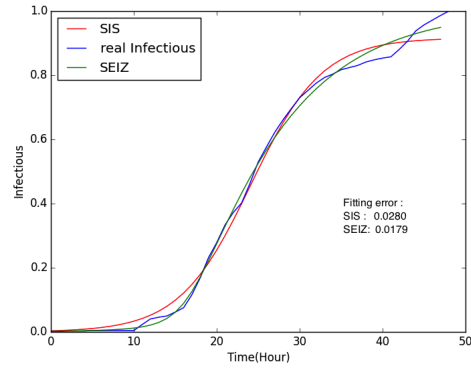
FittingFunction<sub>SEIZ</sub>(TweetVolume<sub>0</sub>, ..., TweetVolume<sub>n</sub>) $\rightarrow \beta_n, b_n, l_n, p_n, \varepsilon_n, \rho_n, RSI_n$

We show 4 examples as following two rumors in figure 5.4(a) 5.4(b) and two news in figure 5.4(c) 5.4(d). It is obvious that SEIZ is more appropriate than SIS to model in our Twitter application, because the fitting error of SPEI is less than SIS.

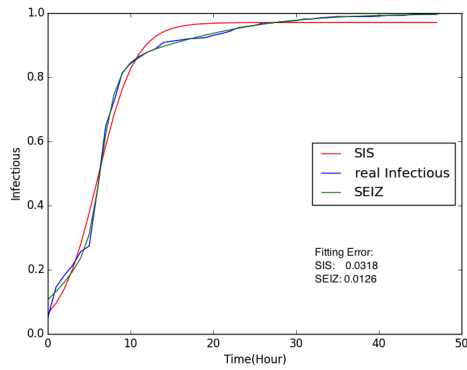
But fitting the models needs enough input data. We don't have enough data to learn the parameters at the first few hours. We show the performance of fitting these two model with only the first 10 hours tweet volume in figure 5.4. As we can see excepting the first one, the fitting result of other three is not good enough.



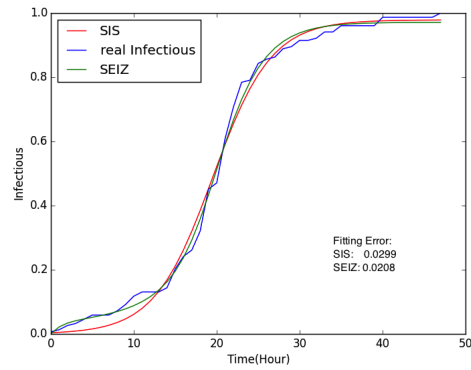
(a) SIS and SEIZ model for rumor 1



(b) SIS and SEIZ model for rumor 2

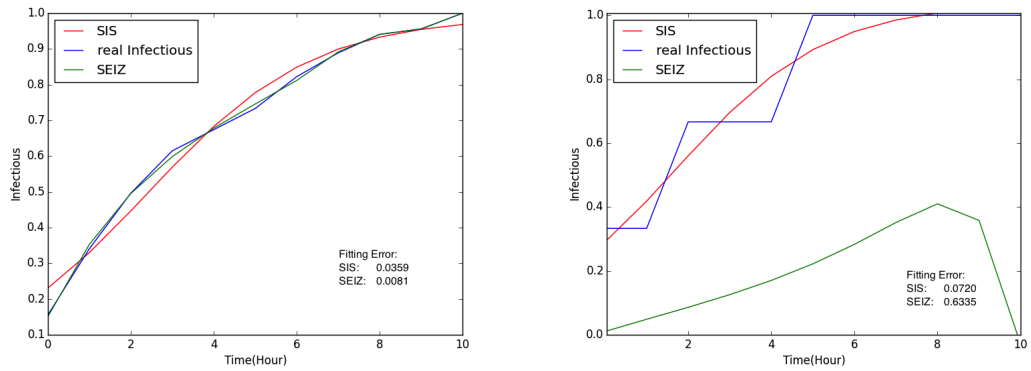


(c) SIS and SEIZ model for news 1

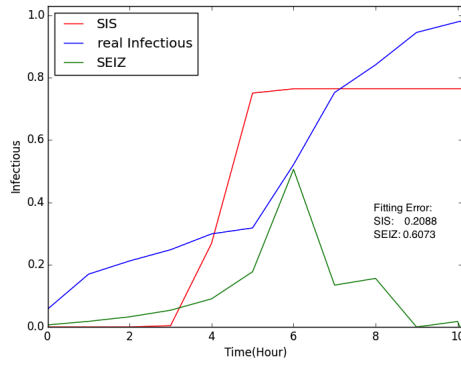


(d) SIS and SEIZ model for news 2

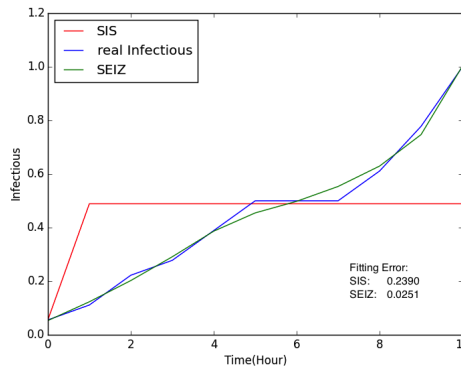
Figure 5.3: Fitting results of SIS and SEIZ model of (a) Rumor: Robert Byrd was a member of KKK (b) Rumor: CNN altered a photograph of a shooter making him look white (c) News: Doctor announces Michael Schumacher is making process (d) News: Two U.S. sailors are arrested over an alleged rape of a Japanese woman on Okinawa



(a) SIS and SEIZ model for rumor 1 with 10 hours data (b) SIS and SEIZ model for rumor 2 with 10 hours data



(c) SIS and SEIZ model for news 1 with 10 hours data



(d) SIS and SEIZ model for news 2 with 10 hours data

Figure 5.4: Fitting results of SIS and SEIZ model with only first 10 hours tweet volume data (same 4 stories as above)



### 5.2.5 SpikeM model Features

Kwon is another approach [18] discussing differences between the rumors' propagation pattern and the news events' propagation pattern on twitter. He adjusted the SpikeM Model and also used the parameters as features.

SpikeM first was introduced by Yasuko Matsubara [24] which can describe the pattern of information diffusion. We present it as follows:

$$\Delta B(n+1) = p(n+1) \cdot (U(n) \cdot \sum_{t=n_b}^n (\Delta B(t) + S(t)) \cdot f(n+1-t) + \varepsilon) \quad (5.14)$$

$$p(n) = 1 - \frac{1}{2} P_a (\sin(\frac{2\pi}{P_p}(n + P_s))) + 1 \quad (5.15)$$

$$U(n+1) = U(n) - \Delta B(n+1) \quad (5.16)$$

where

$$f(\tau) = \beta \cdot \tau^{-1.5} \quad (5.17)$$

and initial conditions:

$$\Delta B(0) = 0, U(0) = N \quad (5.18)$$

In addition, adding an external shock  $S(n)$ , a spike generated at beginning time  $n_b$ . Mathematically, it is defined as follows:

$$S(n) = \begin{cases} 0 & (n \neq n_b) \\ S_b & (n = n_b) \end{cases} \quad (5.19)$$

As the definition:

$$B(n) + U(n) = N \quad (5.20)$$

The term of  $\sum_{t=n_b}^n (\Delta B(t) + S(t))$  is the total number of informed users at time  $n$ , so  $\Delta B(n+1) = p(n+1) \cdot (U(n) \cdot \sum_{t=n_b}^n (\Delta B(t) + S(t)) \cdot f(n+1-t) + \varepsilon)$  means that at time  $n+1$  an infected node  $n$  randomly select a node  $m$  of all nodes and if the node  $m$  is susceptible the probability of  $m$  turning to infected is  $\beta$ , so it is a standard SI model. SpikeM extends the SI model from

- a power-law decay term  $f(\tau) = \beta \cdot \tau^{-1.5}$  in equation 5.24. So the earlier infected nodes has less strength of infection in a power-law decay pattern.
- a periodic interaction function in equation 5.22. It stands for that people have a periodic interaction patterns, like people go to sleep at night so they post much more tweets in the day. Parameters  $P_p$ ,  $P_a$ , and  $P_s$  are the period, strength, and shift of the periodic interaction function.
- $\varepsilon$  is the background noise term.

Symbol	Definition
N	total population of available bloggers
$n_d$	duration of sequence
n	time-tick ( $n=0, \dots, n_d$ )
U(n)	count of <u>un</u> -informed bloggers
B(n)	count of informed <b>b</b> loggers
$\Delta B(n)$	count of informed <b>b</b> loggers at time n
f(n)	<u>in</u> fectiveness of a blog-post, at age n
$\beta$	strength of infection
$\beta \cdot N$	"first-burst" size of infection
S(n)	volume of external <u>s</u> hock at time n
$n_b$	starting time of <b>b</b> reaking news
$S_b$	strength of external shock at birth (time $n_b$ )
$\epsilon$	background noise
$P_a$	strength of periodicity
$P_p$	period
$P_s$	phase shift of periodicity

Table 5.2: Parameters of SpikeM

But the SpikeM can't fit to the events with multi-pike like the figure 2.2. So the author think the term external shock  $S(n)$  in equation 5.25 should not occur once but more. So they extend the SpikeM model by adding a periodic interaction function the term external shock  $S(n)$ .

$$\Delta B(n+1) = p(n+1) \cdot (U(n) \cdot \sum_{t=n_b}^n (\Delta B(t) + \bar{S}(t)) \cdot f(n+1-t) + \epsilon) \quad (5.21)$$

$$p(n) = 1 - \frac{1}{2} P_a (\sin(\frac{2\pi}{P_p} (n + P_s))) + 1) \quad (5.22)$$

$$U(n+1) = U(n) - \Delta B(n+1) \quad (5.23)$$

$$f(\tau) = \beta \cdot \tau^{-1.5} \quad (5.24)$$

The external shock  $S(n)$  is added a a periodic interaction function

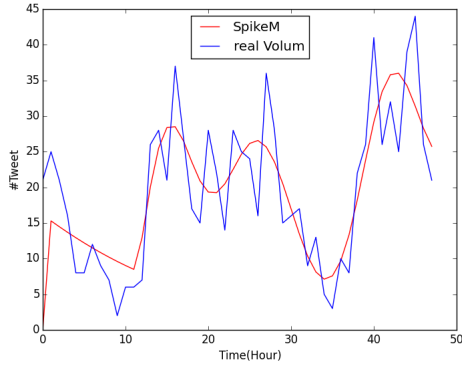
$$\bar{S}(t) = S(t) + q(t) \quad (5.25)$$

$$q(t) = q_a(\sin(\frac{2\pi}{q_p}(t + q_s))) + 1 \quad (5.26)$$

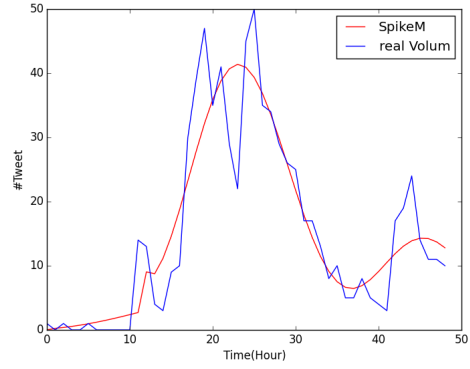
Symbol	Definition
$q_a$	strength of periodicity of the external shock
$q_p$	period of periodicity of the external shock
$q_s$	phase shift of periodicity of the external shock

Table 5.3: New Parameters of extended SpikeM

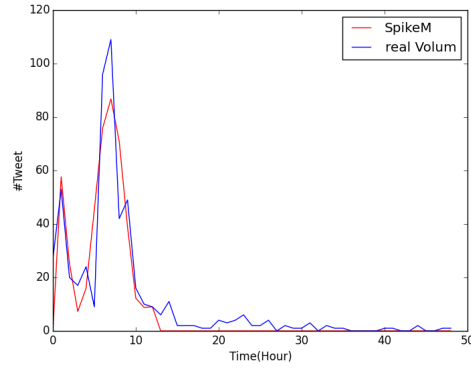
As the same approach of fitting SIS model, we learn the parameters of SpikeM model with Levenberg-Marquard algorithm. We fit the sequenced tweet volume from the beginning time the  $t_0$  to the current time interval  $t_n$  of an event to the model and use the output parameters as the features adding into DSTS. We use the  $p_a, p_p, p_s$  and  $q_a, q_p, q_s$  as feature. We show 4 examples of the SpikeM fitting result in figure 5.5. But same problem as fitting SIS or SEIZ, if we test only within 10 hours data the result seems much worse than the result with full 48 hours showing in figure 5.6.



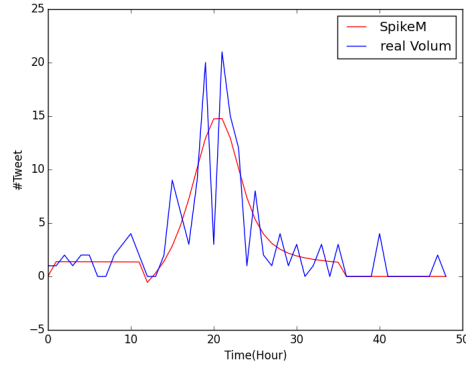
(a) SIS and SEIZ model for rumor 1



(b) SIS and SEIZ model for rumor 2

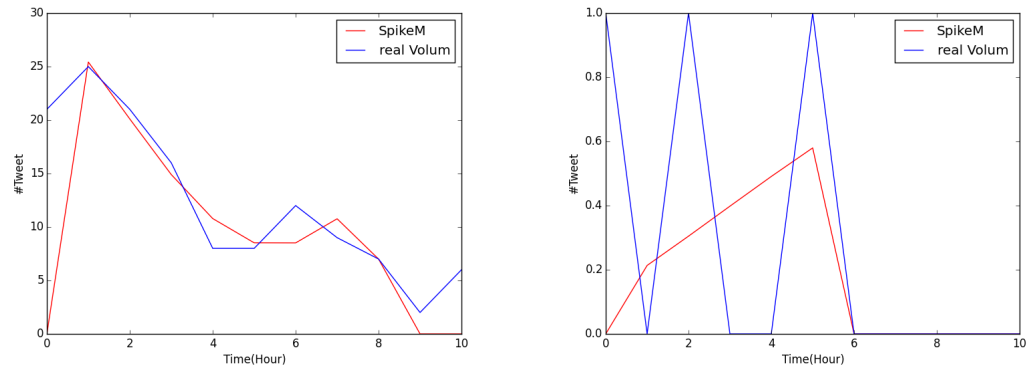


(c) SIS and SEIZ model for news 1

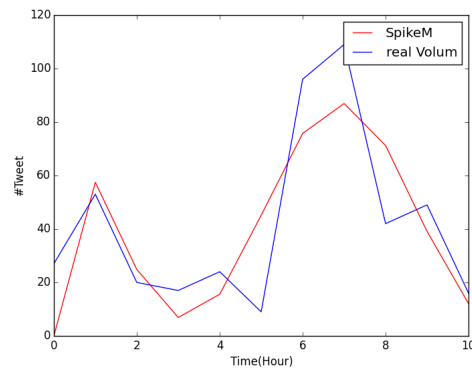


(d) SIS and SEIZ model for news 2

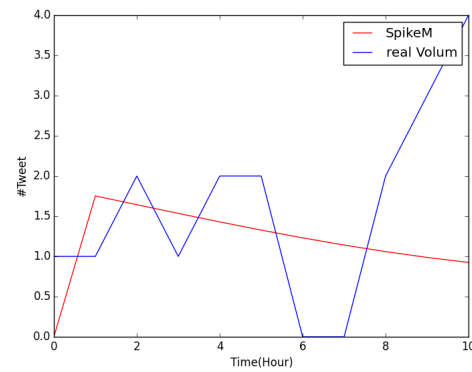
Figure 5.5: Fitting results of SpikeM model of (a) Rumor: Robert Byrd was a member of KKK (b) Rumor: CNN altered a photograph of a shooter making him look white (c) News: Doctor announces Michael Schumacher is making process (d) News: Two U.S. sailors are arrested over an alleged rape of a Japanese woman on Okinawa



(a) SIS and SEIZ model for rumor 1 with 10 hours data (b) SIS and SEIZ model for rumor 2 with 10 hours data



(c) SIS and SEIZ model for news 1 with 10 hours data



(d) SIS and SEIZ model for news 2 with 10 hours data

Figure 5.6: Fitting results of SpikeM model with first 10 hours data (same stories as above)

### 5.2.6 Crowd Wisdom Features

The idea come from Liu's work [21] but not same. The core idea is using the public's common sense to detecting the rumors. If there are more people denying or doubting the truth of an event, this event are more likely to be a rumor. In the Liu's work he uses a extensive list of positive, negative and negation keywords and a set of rules like "negative words without negation words means the poster denies the event". And he uses the ratio number of positive poster (supporter) to the negative poster (deny the events).

Our work is simpler than than his work. We have only a set of negative words, we call it "debunking words" like hoax, rumor, not true, etc. In our test, this is a good feature, but it needs 17 hours to "warm up". It is logical because crowds can debunk rumors but crowds use time to react the event.

### 5.2.7 CreditScore Features

This feature is new feature. We our trained single tweet's creditability model the predict the tweets of the events. If the output is rumor we label it 1 otherwise 0 and we calculate the average score of the events in certain hour. We call this feature creditsScore. We will show it late, this is the best of our dataset, it improves the performance of time series model especially in the first 12 hours. When the event begin bursting in the early stage, people can only rely on the information from the single tweet itself. Because there is no clear propagation structure (pattern) or the wise men or journalists deny the event yet. Our neural network model "sees and check" the text of a single tweet and "give us the advise" when we have no other features in the first few hours since the beginning of the events.

The result shows this feature is the best feature.

Category	Feature	Description
Twitter Features	Hashtag	% of the tweets containing #hashtag
	Mention	% of the tweets mentioning others @user
	NumUrls	# of url in the tweet
	Retweets	average times of tweets have been retweeted
	IsRetweet	% of tweets are retweeted from others
	ContainNEWS	% of tweets containing URL and its domain's catalogue is News
	WotScore	average WOT score of domain in URL
Text Features	Isretweet	% of tweets being retweeted from others
	LengthofTweet	average length of tweets
	NumOfChar	average number of individual characters of tweets
	Capital	average fraction of characters in Uppercase of tweets
	Smile	% of tweets containing :->, :-), :->, :-)
	Sad	% of tweets containing :-<, :-(. :->, :-<
	NumPositiveWords	average number of positive words
	NumNegativeWords	average number of negative words
	PolarityScores	average polarity scores of the Tweets
	Via	% of tweets containing via
	Stock	% of tweets containing \$
	Question	% of tweets containing ?
	Exclamation	% of tweets containing !
	QuestionExclamation	% of tweets containing multi Question or Exclamation mark
	I	% of tweets containing first pronoun like I, my, mine, we, our
User Features	You	% of tweets containing second pronoun like U, you, your, yours
	HeShe	% of tweets containing third pronoun like he, she, they, his, etc.
	UserNumFollowers	average number of followers
	UserNumFriends	average number of friends
	UserNumTweets	average number of users posted tweets
	UserNumPhotos	average number of users posted photos
	UserIsInLargeCity	% of users living in large city
	UserJoinDate	average days since users joining Twitter
Epidemiological Features	UserDescription	% of user having description
	UserVerified	% of user being a verified user
	UserReputationScore	average ratio of #Friends over (#Followers + #Friends)
	$\beta_{SIS}$	Parameter $\beta$ of Model SIS
	$\alpha_{SIS}$	Parameter $\alpha$ of Model SIS
	$\beta_{SEIZ}$	Parameter $\beta$ of Model SEIZ
	$b_{SEIZ}$	Parameter $b$ of Model SEIZ
	$l_{SEIZ}$	Parameter $l$ of Model SEIZ
SpikeM Model Features	$p_{SEIZ}$	Parameter $p$ of Model SEIZ
	$\epsilon_{SEIZ}$	Parameter $\epsilon$ of Model SEIZ
	$\rho_{SEIZ}$	Parameter $\rho$ of Model SEIZ
	$R_{SI}$	Parameter $R_{SI}$ of Model SEIZ
	$P_s$	Parameter $P_s$ of Model Spike
	$P_a$	Parameter $P_a$ of Model SpikeM
Crowd Wisdom Features	$P_p$	Parameter $P_p$ of Model SpikeM
	$Q_s$	Parameter $Q_s$ of Model SpikeM
	$Q_a$	Parameter $Q_a$ of Model SpikeM
	$Q_p$	Parameter $Q_p$ of Model SpikeM
Credit Score Features	CreditScore	average CreditScore

Table 5.4: Features of Time Series Rumor Detection Model



### 5.3 Classification Models

Same reason as the single tweet's Creditability we test the time series model also with 3 popular models Decision Trees, SVM, Random Forest and one more model the multilayer perceptron (MLP). We show the optimized parameters in the table 5.5. After 10-fold cross validation with same shuffled sequence testing above models, the result is shown the figure 5.6. The random forest is still the best model for our task, so we use RF as the test model for the further features' Evolution.

Model	Parameters	Value
Random Forest	Number of Trees	350
SVM	kernel	radial basis function
	penalty parameter of the error term	3.0
	gamma	$\frac{1}{50}$
Decision Trees	criterion	gini
MLP	alpha	0.0001
	activation function	ReLU
	hidden layer sizes	2 layer(50 nodes each layer)
	weight optimization	adam

Table 5.5: Parameters of Classification models

Model	Accuracy in hours			
	6	12	24	48
Random Forest	<b>0.8615</b>	<b>0.8615</b>	<b>0.8692</b>	<b>0.9076</b>
MLP	0.7423	0.7423	0.7692	0.8192
SVM	0.7423	0.7884	0.7769	0.7538
Decision Trees	0.7807	0.8115	0.75	0.7385

Table 5.6: Prediction Accuracy of Different Single Tweet's Creditability Scoring Models

## 5.4 Feature Evolution

### 5.4.1 Time Series VS Static Feature

First we compare our time series model and the normal static feature model. We show the result is in figure 5.7 the full 48 hours details in Appendix table .2 and in figure . As we can see from the result that the accuracy of time Series is better than the static. But after 24 hours the advantage of the time series is very limited.

Hour	Time series model	Static model	Difference
1	0.82	0.78	0.03
6	0.86	0.8	0.06
12	0.87	0.83	0.04
18	0.88	0.83	0.05
24	0.87	0.85	0.01
30	0.88	0.85	0.03
36	0.87	0.86	0.01
42	0.88	0.88	0
48	0.89	0.87	0.02

Table 5.7: Accuracy: Time series VS static Features

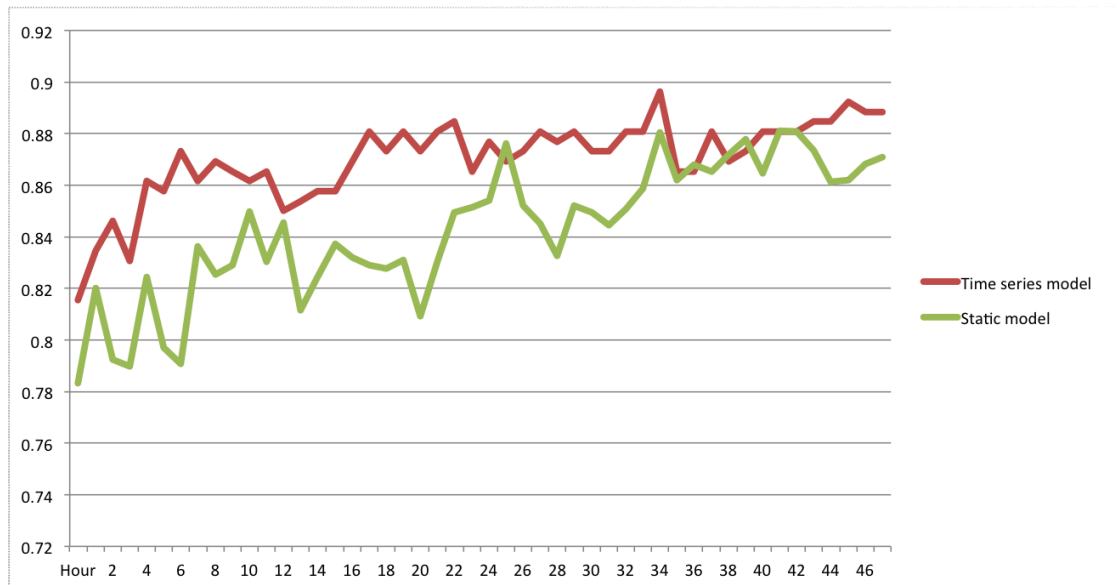


Figure 5.7: Accuracy: Time series VS static Features

### 5.4.2 Feature Analyzing Over Time

We rank the features' importance using the method we introduced in section 2.4.3, the full result is shown in Appendix table .1. First we split the features in 7 catalogues as table 5.4 Tweet\_Feature, User\_Feature, Text\_Feature, CreditScore, SpikeM, Epidemiological Features, CrowdWisdom and the BestSets. The BestSets is top 23 features on average in table 5.8.

Best Feature set
CreditScore
ContainNEWS
NumChar
UserTweetsPerDays
QuestionExclamation
UserReputationScore
WotScore
Question
UserJoin_date
Mention
LengthOfTweet
DebunkingWords
UserFollowers
Exclamation
UserVerified
Hashtag
Capital
You
UserNumPhoto
numUrls
UserFriends
NumPositiveWords
Via
$P_a$
UserIsInLargeCity
PolarityScores
UserDescription
$R_{SI}$
$\beta_{SIS}$
$\epsilon_{SEIZ}$

Table 5.8: Best Features



## 6.1 Conclusion

With more nations working actively to preserve the web due to the importance it plays in our lives, scholars from various fields have started looking to web archives as a data source waiting to be analyzed. We find that scholarly research with web archives carried out until today can be broadly classified as link based and content based studies. The link structure of web archives has been studied on a quantitative scale by Internet research institutes and computer scientists. They either study the link structure over time to analyze the evolution of the web or try to correlate the linking structure to real world phenomena. In this thesis we first looked at the usage of web archives by scholars. According to Brügger there are 4 major steps when conducting any type of research on web archives: corpus creation, analysis, dissemination and storage. An ideal system for web archive research should support all 4 steps however we find that adequate steps are only now being made to help researchers with the first step: corpus creation. Humanities scholars are usually interested in doing small scale qualitative and quantitative studies for which they need to explore web archives for relevant material. But the access to web archives with current web archive access systems like the wayback machine and rudimentary keyword search is not satisfactory. There are many other factors as well that have led to the lack of humanities studies on web archives. We attempted to illicit these through literature surveys and group discussions with humanities scholars conducting web archive research. Some of the interesting outcomes of this study were:

- The absence of abilities to explore the web archive discourage scholars from pursuing their ideas. Search engines are the predominant way of exploring the web today whereas archives seem to be left behind.
- The wayback machine is not an effective tool for building a corpus unless you

know the exact set of URLs you are interested in. Full text search should be provided for researchers who want to explore the archive.

- More details of the crawl should be exposed so that scholars can motivate their corpus of study better. The details of the crawl strategy should be added to the meta data of all documents.
- Scholars are used to working with well curated archives. Web archives pose greater burden for scholars in corpus creation since they need to curate the materials themselves.
- Due to the inconsistent nature of web archives, an error margin should be defined and an approximation of confidence as well even for small scale qualitative studies.
- Algorithms and features used to help scholars find new documents should be made transparent to users so that they can effectively document the corpus creation process.
- A web archive search system should allow for the corpus and its history of creation to be exported in standard formats researchers are used to working with.
- Scholars are interested in leveraging digital methods to analyze their data but find it hard to find tools flexible enough for their needs.

Scholars also saw potential in the usage of web archives as a playground for combining qualitative and quantitative studies. Based on our discussions with the scholars who attended the summer school in Aarhus University, we found that many would be interested in the ability to scale results up. They also bemoaned the lack of analysis tools currently available to work with web archives. The lack of technical expertise was clearly affecting a scholar's ability to conduct his research even though they knew what kind of analysis they wanted to do. We proposed that computer scientist should bridge this gap by either providing tools to support scholars or work with them on scaling up their hypothesis using pattern recognition and data mining techniques.

Initiatives are being made kick-start humanities research in web archives but to do the necessary infrastructure to first access web archives has to be in place. These initiatives are currently being carried out by national libraries and institutions like the Internet Archive and the IIPC. The BUDDAH project, a joint effort between the British Library and leading Universities in the UK, is working towards building a web archive search system by working in tandem with humanities scholars working with the UK web archive. We held discussions with members of the BUDDAH project and scholars to identify the pitfalls with the current web archive search system

they were working with. The British Library web archive search prototype already provides its users with keyword search and filtering based on domain and crawl date. As noted by Costa et. al. current IR techniques like keyword search do not fare well in web archive search. We observed this first hand with the British Library’s choice to rank search results by crawl date. This makes it very difficult for users to easily comprehend a large result set. We decided to better understand a scholar’s search intent by analyzing written descriptions of proposed studies on web archives by humanities scholars involved in the BUDDAH project. We found that scholars are interested in studying results over time and covering a variety of aspect for their topic.

To this end we introduced the notion of a historical query intent over longitudinal collections like web archives. Historical query intent was modeled as a 2 dimension diversification problem where one dimension is time and the other aspects. In this thesis we limit ourselves to a subset of web archives: the news archive. We proposed a new evaluation metric Tia-SBR to evaluate the performance of retrieval models for this task. We also adapted existing diversity metrics to take time into account. To evaluate the retrieval models we built a new temporal test collection based on 20 years of the *New York Times* collection. We strengthened existing 1 dimensional diversity methods by linearizing the aspect-time space into a single set. We also consider temporal diversity retrieval models in our experiment. We introduced HistDiv which shows improvements over temporal and non-temporal methods for most of the time-aware diversification methods. We also outperform all competitors in Tia-SBR showing the suitability of our approach for historical query intents. We observe that HistDiv works well for topics which have aspects that span across multiple time intervals and have fluctuating importance at different times. It trades-off nicely between important aspects and important times which we perceive as important in historical search. HistDiv does not perform quite as well for queries which only one dominant aspect at a certain time window. We also described the architecture of the History Search system which consists of: a REST API to access the retrieval models used in our experiments, a user interface designed specifically to show results as newspaper articles and a time line to provide an overview and ability to filter.

## 6.2 Future Work

The performance of the HistDiv algorithm encourages us to look beyond the traditional approaches and devise methods more suited to archives. HistDiv currently considers time as a set of disjointed windows. It will be interesting to see if bursts rather than time windows are more effective for historical query intents. In our current implementation, subtopics are mined using wikiminer which is a relatively naive approach. In the future, we want to experiment with different types of aspect

mining like LDA for instance. We also want to add more dimensions to the model like geographic locations. Just like aspect utility decays based on time, we can also decay it based on location. The absolute distance between lat long coordinates can be used or a taxonomy based approach.

The results also indicate where HistDiv performs badly and other algorithms perform well. We believe that if the user is presented with a choice of algorithms, she can effectively select the algorithm best suited to her need. A step beyond this would be to train a classifier that selects the appropriate retrieval model based on features from the query's aspect and temporal profile.

The test collection we built is still relatively small when compared to collections by TREC. We must strive to add to this test collection, by increasing the workload, or explore different avenues for evaluating retrieval models meant for archive search. For our test collection it will also be interesting to get actual historians to provide topics. More research is also needed to construct test collections for other query intents for archives.

Having performed admirably in a news archive test collection, it will be interesting to test HistDiv and the other methods on a web archive or at least a subset which is not just news. Existing web crawl datasets only span a small period of time but it might be interesting nonetheless to experiment with smaller time window sizes for these collections. Another important temporal feature that we want to consider is the temporal references within an article. It is reasonable to assume that major bursts are usually covered by summary articles after the end of the burst. These articles are valuable but may get neglected because they do not occur within the burst based on our current algorithm. HistDiv's ability to find documents from important time windows and aspects can also be applied in Temporalia's query classification task to interesting effect.

There is also much to be done to help scholars better engage with web archives. In this thesis we have only scratched the surface of the issues surrounding intelligent access to web archives. We must work together with scholars in order to develop effective algorithms and user friendly systems so that the researchers of tomorrow can conduct their studies with far greater ease than today.



## Bibliography

- [1] F. Ahmed and M. Abulaish. An mcl-based approach for spam profile detection in online social networks. In *2012 IEEE 11th International Conference on Trust, Security and Privacy in Computing and Communications*, pages 602–608. IEEE, 2012.
- [2] G. W. Allport and L. Postman. The psychology of rumor. 1947.
- [3] Y. Bao, C. Yi, Y. Xue, and Y. Dong. A new rumor propagation model and control strategy on social networks. In *Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, pages 1472–1473. ACM, 2013.
- [4] L. Barbosa and J. Feng. Robust sentiment detection on twitter from biased and noisy data. In *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*, pages 36–44. Association for Computational Linguistics, 2010.
- [5] L. M. Bettencourt, A. Cintrón-Arias, D. I. Kaiser, and C. Castillo-Chávez. The power of a good idea: Quantitative modeling of the spread of ideas from epidemiological models. *Physica A: Statistical Mechanics and its Applications*, 364:513–536, 2006.
- [6] L. Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.
- [7] L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [8] L. Breiman, J. Friedman, C. J. Stone, and R. A. Olshen. *Classification and regression trees*. CRC press, 1984.

- [9] C. Castillo, M. Mendoza, and B. Poblete. Information credibility on twitter. In *Proceedings of the 20th international conference on World wide web*, pages 675–684. ACM, 2011.
- [10] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [11] A. Gupta, P. Kumaraguru, C. Castillo, and P. Meier. Tweetcred: Real-time credibility assessment of content on twitter. In *International Conference on Social Informatics*, pages 228–243. Springer, 2014.
- [12] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [13] F. Jin, E. Dougherty, P. Saraf, Y. Cao, and N. Ramakrishnan. Epidemiological modeling of news and rumors on twitter. In *Proceedings of the 7th Workshop on Social Network Mining and Analysis*, page 8. ACM, 2013.
- [14] A. Joulin, E. Grave, P. Bojanowski, and T. Mikolov. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*, 2016.
- [15] Y. Kim. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*, 2014.
- [16] D. Kimmey. Twitter event detection. 2015.
- [17] S. Kwon, M. Cha, K. Jung, W. Chen, and Y. Wang. Aspects of rumor spreading on a microblog network. In *International Conference on Social Informatics*, pages 299–308. Springer, 2013.
- [18] S. Kwon, M. Cha, K. Jung, W. Chen, and Y. Wang. Prominent features of rumor propagation in online social media. In *2013 IEEE 13th International Conference on Data Mining*, pages 1103–1108. IEEE, 2013.
- [19] S. Lai, L. Xu, K. Liu, and J. Zhao. Recurrent convolutional neural networks for text classification. In *AAAI*, pages 2267–2273, 2015.
- [20] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- [21] X. Liu, A. Nourbakhsh, Q. Li, R. Fang, and S. Shah. Real-time rumor debunking on twitter. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, pages 1867–1870. ACM, 2015.

- 
- [22] J. Ma, W. Gao, P. Mitra, S. Kwon, B. J. Jansen, K.-F. Wong, and M. Cha. Detecting rumors from microblogs with recurrent neural networks.
- [23] J. Ma, W. Gao, Z. Wei, Y. Lu, and K.-F. Wong. Detect rumors using time series of social context information on microblogging websites. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, pages 1751–1754. ACM, 2015.
- [24] Y. Matsubara, Y. Sakurai, B. A. Prakash, L. Li, and C. Faloutsos. Rise and fall patterns of information diffusion: model and implications. In Q. Yang, D. Agarwal, and J. Pei, editors, *KDD*, pages 6–14. ACM, 2012.
- [25] M. Mendoza, B. Poblete, and C. Castillo. Twitter under crisis: can we trust what we rt? In *Proceedings of the first workshop on social media analytics*, pages 71–79. ACM, 2010.
- [26] O. Oh, K. H. Kwon, and H. R. Rao. An exploration of social media in extreme events: Rumor theory and twitter during the haiti earthquake 2010. In *ICIS*, page 231, 2010.
- [27] C. Olah. Understanding lstm networks. *Net*: <http://colah.github.io/posts/2015-08-Understanding-LSTMs>, 2015.
- [28] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. *ICML (3)*, 28:1310–1318, 2013.
- [29] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning representations by back-propagating errors. *Cognitive modeling*, 5(3):1, 1988.
- [30] S. J. Russell, P. Norvig, J. F. Canny, J. M. Malik, and D. D. Edwards. *Artificial intelligence: a modern approach*, volume 2. Prentice hall Upper Saddle River, 2003.
- [31] K. Starbird, J. Maddock, M. Orand, P. Achterman, and R. M. Mason. Rumors, false flags, and digital vigilantes: Misinformation on twitter after the 2013 boston marathon bombing. *iConference 2014 Proceedings*, 2014.
- [32] C. R. Sunstein. *On rumors: How falsehoods spread, why we believe them, and what can be done*. Princeton University Press, 2014.
- [33] Y. Tanaka, Y. Sakamoto, and T. Matsuka. Transmission of rumor and criticism in twitter after the great japan earthquake. In *Annual Meeting of the Cognitive Science Society*, page 2387, 2012.
- [34] R. Thomson, N. Ito, H. Suda, F. Lin, Y. Liu, R. Hayasaka, R. Isochi, and Z. Wang. Trusting tweets: The fukushima disaster and information source credibility on twitter. *Proc. of ISCRAM*, 10, 2012.

- [35] A. H. Wang. Don't follow me: Spam detection in twitter. In *Security and Cryptography (SECRYPT), Proceedings of the 2010 International Conference on*, pages 1–10. IEEE, 2010.
- [36] K. Wu, S. Yang, and K. Q. Zhu. False rumors detection on sina weibo by propagation structures. In *2015 IEEE 31st International Conference on Data Engineering*, pages 651–662. IEEE, 2015.
- [37] F. Yang, Y. Liu, X. Yu, and M. Yang. Automatic detection of rumor on sina weibo. In *Proceedings of the ACM SIGKDD Workshop on Mining Data Semantics*, page 13. ACM, 2012.
- [38] Z. Zhao, P. Resnick, and Q. Mei. Enquiring minds: Early detection of rumors in social media from enquiry posts. In *Proceedings of the 24th International Conference on World Wide Web*, pages 1395–1405. ACM, 2015.
- [39] C. Zhou, C. Sun, Z. Liu, and F. Lau. A c-lstm neural network for text classification. *arXiv preprint arXiv:1511.08630*, 2015.

## **Appendix**

time	CreditScore	ContainNEWS	NumChar	UserTweetsPerDays	QuestionExclamation	UserReputationScore	UrlRankIn5000	WotScore	Question	UserJoin_date	Mention	LengthOfTweet	DebunkingWords	UserFollowers	Exclamation	UserVerified	Hashtag	Capital	You	UserNumPhoto	numUrls	UserFriends	NumPositiveWords	Via	P <sub>α</sub>	UserIsInLargeCity	PolarityScores	UserDescription	R <sub>SI</sub>	β <sub>SI</sub>	ε <sub>SEIZ</sub>	NumNegativeWords	P <sub>s</sub>	b <sub>SEIZ</sub>	β <sub>SEIZ</sub>	Retweets	p <sub>SEIZ</sub>	α <sub>SI</sub>	Q <sub>s</sub>	P <sub>p</sub>	p <sub>SEIZ</sub>	I	Q <sub>p</sub>	l <sub>SEIZ</sub>	Q <sub>α</sub>	HeShe	IsRetweet	Stock	Smile	Sad	
1	2	8	3	0	25	1	14	4	15	5	13	6	33	7	27	24	20	10	28	9	22	11	18	30	29	19	12	23	45	42	39	17	31	43	46	21	40	41	47	44	49	16	38	37	48	26	32	34	35	36	
2	1	5	2	0	18	3	11	4	12	7	15	6	43	8	36	22	21	14	40	9	30	10	23	34	25	24	13	38	31	28	41	16	29	27	20	19	39	32	48	26	35	17	47	33	49	37	42	44	45	46	
3	1	4	2	0	17	5	12	3	10	8	11	6	44	9	38	24	23	15	41	7	30	13	20	32	25	22	14	39	21	26	27	18	36	37	35	19	29	33	48	28	40	16	47	31	49	34	42	43	45	46	
4	1	5	3	0	16	2	13	4	14	7	11	9	43	6	39	34	21	12	30	8	37	10	23	32	26	25	15	41	28	33	31	17	19	22	27	24	38	40	48	18	29	20	47	35	49	36	42	44	45	46	
5	0	4	3	1	16	2	13	10	11	8	5	6	42	7	38	17	20	14	33	9	36	12	22	28	35	30	15	41	31	24	23	18	32	25	34	21	29	19	48	37	40	27	47	26	49	39	43	44	45	46	
6	0	4	2	1	14	3	11	9	8	10	7	6	41	5	35	29	22	15	34	12	40	16	17	21	28	25	18	42	24	30	13	19	36	20	27	26	23	37	48	33	31	32	47	39	49	38	43	44	45	46	
7	0	6	2	1	7	3	11	12	10	8	4	5	23	9	22	28	21	15	39	13	37	14	20	19	27	26	18	42	24	38	17	16	34	35	25	31	29	41	48	30	33	32	47	40	49	36	43	44	45	46	
8	0	4	2	1	5	3	7	9	8	10	13	6	24	11	18	15	21	14	41	12	34	19	20	31	28	17	22	40	23	29	27	16	25	35	33	26	32	37	48	30	36	38	47	39	49	42	43	44	45	46	
9	0	8	2	1	4	3	7	5	10	6	13	9	25	11	15	14	23	12	32	16	28	19	22	24	26	33	20	37	31	30	18	17	41	35	29	21	27	36	48	39	34	38	47	40	49	42	43	44	45	46	
10	0	8	3	1	2	4	5	10	11	6	12	7	18	9	17	19	20	13	32	14	26	24	23	27	28	36	16	39	41	30	21	15	35	28	34	31	22	25	48	38	37	33	47	40	49	43	42	44	45	46	
11	0	5	4	1	2	3	7	6	13	8	10	9	19	11	17	12	15	14	24	16	18	20	25	31	33	32	23	40	26	29	36	21	37	39	22	27	38	30	48	35	42	28	47	34	49	41	43	44	45	46	
12	0	6	3	1	2	4	7	5	13	8	9	11	18	12	16	15	17	10	21	20	23	14	28	27	34	25	19	41	30	33	36	31	40	26	37	24	32	22	48	35	29	39	47	38	49	42	43	44	45	46	
13	0	7	3	1	2	4	9	5	12	6	8	11	15	13	22	14	17	10	35	18	20	16	31	27	36	24	23	43	40	30	26	34	37	28	32	21	29	19	42	33	25	38	39	41	44	45	46	47	48	49	
14	0	6	3	1	2	4	5	9	12	7	8	11	15	13	17	14	22	10	26	19	18	16	27	36	34	29	21	39	24	23	25	33	41	28	32	20	40	42	31	35	30	44	38	37	43	46	45	47	48	49	
15	0	5	4	1	2	3	6	13	8	7	11	10	12	16	9	15	27	14	20	26	19	18	22	25	32	34	30	40	24	37	17	21	38	33	41	28	35	44	29	39	23	43	36	42	31	46	45	47	48	49	
16	0	4	3	2	1	5	7	10	11	9	12	6	13	17	8	15	21	14	20	22	23	24	43	36	18	30	29	39	33	37	16	28	41	26	38	32	31	34	35	44	19	42	40	27	25	45	46	48	47	49	
17	0	4	3	2	1	5	11	10	7	8	14	6	12	17	9	16	18	13	15	21	20	19	44	26	35	43	28	31	22	32	36	34	37	39	40	38	30	23	25	41	27	33	24	29	42	45	46	47	48	49	
18	0	5	3	2	1	4	9	8	7	10	12	11	14	16	6	15	17	13	18	24	20	19	40	26	30	32	29	35	23	28	25	31	41	27	38	37	33	36	34	43	39	44	22	42	21	45	46	47	48	49	
19	0	4	3	2	1	7	8	9	6	11	12	10	13	15	5	18	17	14	16	24	19	20	36	26	22	38	29	43	28	33	35	30	23	32	25	41	27	39	34	40	37	42	21	31	44	45	46	47	48	49	
20	0	3	4	2	1	7	5	6	8	13	12	11	10	16	9	17	18	15	14	21	20	23	35	33	28	30	41	29	22	27	39	36	37	43	32	31	25	24	38	40	42	46	19	26	34	45	44	48	47	49	
21	0	4	3	2	1	7	6	8	10	14	12	13	5	17	9	18	15	16	11	27	20	19	30	35	24	33	29	28	25	22	44	37	26	31	36	38	40	21	39	32	42	43	23	41	34	46	45	47	48	49	
22	0	4	3	2	1	7	6	8	11	14	12	13	5	18	9	17	15	16	10	20	21	22	26	33	24	36	30	28	35	25	38	27	23	37	41	43	31	19	40	42	29	39	34	44	32	46	45	47	48	49	
23	0	4	3	2	1	6	7	10	5	12	13	14	9	15	8	17	16	18	11	19	24	23	34	35	25	43	33	32	27	38	29	39	29	21	41	28	42	37	40	22	30	36	46	20	26	31	44	45	47	48	49
24	0	3	5	2	1	6	4	14	7	13	10	11	9	17	8	16	15	19	12	20	18	26	25	24	33	32	30	22	36	31	35	37	23	39	41	40	42	38	21	27	34	45	29	43	28	44	46	47	48	49	
25	0	3	4	2	1	10	7	9	5	13	11	14	6	17	8	19	18	15	12	22	16	30	23	20	37	26	29	24	21	32	31	40	41	28	39	42	36	35	27	38	43	44	25	34	33	45	46	47	48	49	
26	0	4	6	2	1	8	3	9	7	12	14	13	5	18	11	16	15	17	10	22	20	24	19	32	23	28	36	21	26	37	31	34	27	33	43	44	29	35	25	39	41	38	40	42	30	46	45	47	48	49	
27	0	3	8	4	1	7	2	9	6	15	12	13	5	17	14	16	11	18	10	22	21	29	19	25	23	24	33	20	26	30	36	31	34	42	27	44	38	40	28	32	41	37	45	39	35	43	46	47	48	49	
28	0	1	7	9	2	8	3	6	4	14	10	15	5	17	13	16	11	21	12	23	19	27	20	26	35	24	29	18	32	30	22	38	25	28	41	42	44	31	36	37	43	40	34	39	33	46	45	47	48	49	
29	0	2	8	9	1	5	3	6	4	14	10	16	7	17	15	12	13	22	11	23	21	24	19	25	18	32	33	20	31	34	27	36	28	44	39	35	41	38	26	29	42	40	37	43	30	46	45	47	48	49	
30	0	3	8	7	1	9	2	5	4	14	10	16	6	17	13	15	12	20	11	22	21	31	18	23	24	25	29	19	39	26	27	32	36	38	33	34	43	37	35	28	44	30	42	41	40	45	46	47	48	49	
31	0	3	9	4	2	7	1	5	6	16	10	15	8	20	13	14	12	21	11	22	18	24	23	29	17	34	28	19	25	26	30	37	31	36	35	32	44	39	27	46	43	38	41	40	33	45	42	47	48	49	
32	0	2	6	4	3	7	1	9	5	16	10	15	8	17	12	14	13	21	11	22	18	24	23	30	19	26	33	20	28	25	32	42	31	38	34	27	36	44	40	41	43	39	35	29	37	45	46	47	48	49	
33	0	2	4	7	3	8	1	6	11	13	10	15	5	17	12	9	14	19	16	22	20	30	21	24	25	26	42	18	29	23	36	37	33	35	32	27	43	31	28	40	34	39	41	38	46	44	45	47	48	49	
34	0	1	5	3	4	6	2	8	12	16	9	15	7	13	11	10	17	23	14	24	20	27	21	25	19	28	36	18	37	26	32	33	22	39	34	31	45	29	38	40	35	42	43	41	30	44	46	47	48		

Hour	Time series model	Static model	Difference
1	0.82	0.78	0.03
2	0.83	0.82	0.01
3	0.85	0.79	0.05
4	0.83	0.79	0.04
5	0.86	0.82	0.04
6	0.86	0.8	0.06
7	0.87	0.79	0.08
8	0.86	0.84	0.03
9	0.87	0.83	0.04
10	0.87	0.83	0.04
11	0.86	0.85	0.01
12	0.87	0.83	0.04
13	0.85	0.85	0
14	0.85	0.81	0.04
15	0.86	0.82	0.03
16	0.86	0.84	0.02
17	0.87	0.83	0.04
18	0.88	0.83	0.05
19	0.87	0.83	0.05
20	0.88	0.83	0.05
21	0.87	0.81	0.06
22	0.88	0.83	0.05
23	0.88	0.85	0.04
24	0.87	0.85	0.01
25	0.88	0.85	0.02
26	0.87	0.88	-0.01
27	0.87	0.85	0.02
28	0.88	0.85	0.04
29	0.88	0.83	0.04
30	0.88	0.85	0.03
31	0.87	0.85	0.02
32	0.87	0.84	0.03
33	0.88	0.85	0.03
34	0.88	0.86	0.02
35	0.9	0.88	0.02
36	0.87	0.86	0
37	0.87	0.87	0
38	0.88	0.87	0.02
39	0.87	0.87	0
40	0.87	0.88	0
41	0.88	0.86	0.02
42	0.88	0.88	0
43	0.88	0.88	0
44	0.88	0.87	0.01
45	0.88	0.86	0.02
46	0.89	0.86	0.03
47	0.89	0.87	0.02
48	0.89	0.87	0.02

Table .2: Time series VS static Features in detail





## List of Figures

2.1	large scale tweet volume of event Robert Byrd . . . . .	5
2.2	tweet volume of the rumor event of Robert Byrd after selected time period . . . . .	6
2.3	Decision tree . . . . .	7
2.4	An example of random forest with 3 trees . . . . .	8
2.5	An example of 1 layer neural network . . . . .	9
2.6	tanh function . . . . .	10
2.7	CNN for Text Classification (source from[15]) . . . . .	11
2.8	2 layer Recurrent neural network . . . . .	11
2.9	multi-input single-output Recurrent neural network . . . . .	11
2.10	single-input multi-output Recurrent neural network . . . . .	12
2.11	multi-input multi-output Recurrent neural network . . . . .	12
2.12	Two models of RNNs (a) Normal model of RNN with tanh Unit (b) RNN with LSTM Cells (source [27]) . . . . .	13
2.13	GRU Cell (source [27]) . . . . .	14
2.14	Dropout Neural Net Model. (a) a standard neural network (b) after deploying dropout . . . . .	14
2.15	K Fold Cross Validation (green rectangle is the training set, red rectangle is the testing set) . . . . .	15
2.16	The fraction of tweets containing url with top 5000 domain . . . . .	16
2.17	The fraction of poster living in large city . . . . .	16
4.1	Sample of Users' Information on Twitter Interface . . . . .	22
4.2	neural network model for single tweet classification 1 . . . . .	28
5.1	SIS Model . . . . .	31
5.2	SEIZ Model . . . . .	32

5.3	Fitting results of SIS and SEIZ model of (a) Rumor: Robert Byrd was a member of KKK (b) Rumor: CNN altered a photograph of a shooter making him look white (c) News: Doctor announces Michael Schumacher is making process (d) News: Two U.S. sailors are arrested over an alleged rape of a Japanese woman on Okinawa . . . . .	35
5.4	Fitting results of SIS and SEIZ model with only first 10 hours tweet volume data (same 4 stories as above) . . . . .	36
5.5	Fitting results of SpikeM model of (a) Rumor: Robert Byrd was a member of KKK (b) Rumor: CNN altered a photograph of a shooter making him look white (c) News: Doctor announces Michael Schumacher is making process (d) News: Two U.S. sailors are arrested over an alleged rape of a Japanese woman on Okinawa . . . . .	41
5.6	Fitting results of SpikeM model with first 10 hours data (same stories as above) . . . . .	42
5.7	Accuracy: Time series VS static Features . . . . .	46

## List of Tables

4.1	Features for Single Tweet's Creditability Scoring Model . . . . .	23
4.2	Parameters of Classification models . . . . .	23
4.3	Prediction Accuracy of Different Single Tweet's Creditability Scoring Models . . . . .	24
4.4	Features Importance . . . . .	25
4.5	Prediction Accuracy of Different Single Tweet's Creditability Scoring Models . . . . .	27
5.1	Parameters of SEIZ . . . . .	34
5.2	Parameters of SpikeM . . . . .	38
5.3	New Parameters of extended SpikeM . . . . .	39
5.4	Features of Time Series Rumor Detection Model . . . . .	44
5.5	Parameters of Classification models . . . . .	45
5.6	Prediction Accuracy of Different Single Tweet's Creditability Scoring Models . . . . .	45
5.7	Accuracy: Time series VS static Features . . . . .	46
5.8	Best Features . . . . .	47
.1	All feature's importances over 48 hours . . . . .	58
.2	Time series VS static Features in detail . . . . .	59