

# Twitter information extraction system

Yao Chen                      Torsten Suel  
Department of Computer Science  
NYU Polytechnic School of Engineering  
6 MetroTech Center  
Brooklyn, NY 11201  
[sunnychen@nyu.edu](mailto:sunnychen@nyu.edu)

**Abstract:** This paper talks about how to extract users' specific goal from their new year resolutions which are posted on twitter. Finally, the new year resolutions will be classified and ranked according to the total number of each class. The snowball algorithm, LSC algorithm, stemming, lingpipe API to extract users' specific goals from twitter message. The tf-idf algorithm and vector-space model algorithm is used to classify those goals.

**Keywords:** Information extraction, twitter, snowball, stemming, lingpipe, tf-idf, vector-space model

## 1. Introduction:

In twitter message analysis area, there are many work about the sentimental analysis. But there is little software which is used to extract information from twitter message. The reason is twitter messages are not complete English sentence. The existing technique is hard to analysis those information. However, the twitter messages are considerably valuable in many area, such as economic, political field. To solve this problem, this paper filter the twitter message and study on special topic.

In this paper, the “new year resolution” related twitter message is chosen as data source. The system extracts the users' resolution from twitter message and get the top 3 resolution. Techniques in information retrieval and machine learning area are evolved in this system.

Related work “Snowball: Extracting Relations from Large Plain-Text Collections” studies on the texture material and extract the “organization” and “location” information. This study is based on Brin's DIPRE method which is used to extract the “book” and “authority” information. Those two system are partly similar with this paper. However, those study are about information extracted on nouns. Under this background, the work is much more easier for the named-entity tagger algorithm. But, the named-entity tagger algorithm can not be used in this paper, because the users' goal are unusually complex components include verbs, nouns and Adjective.

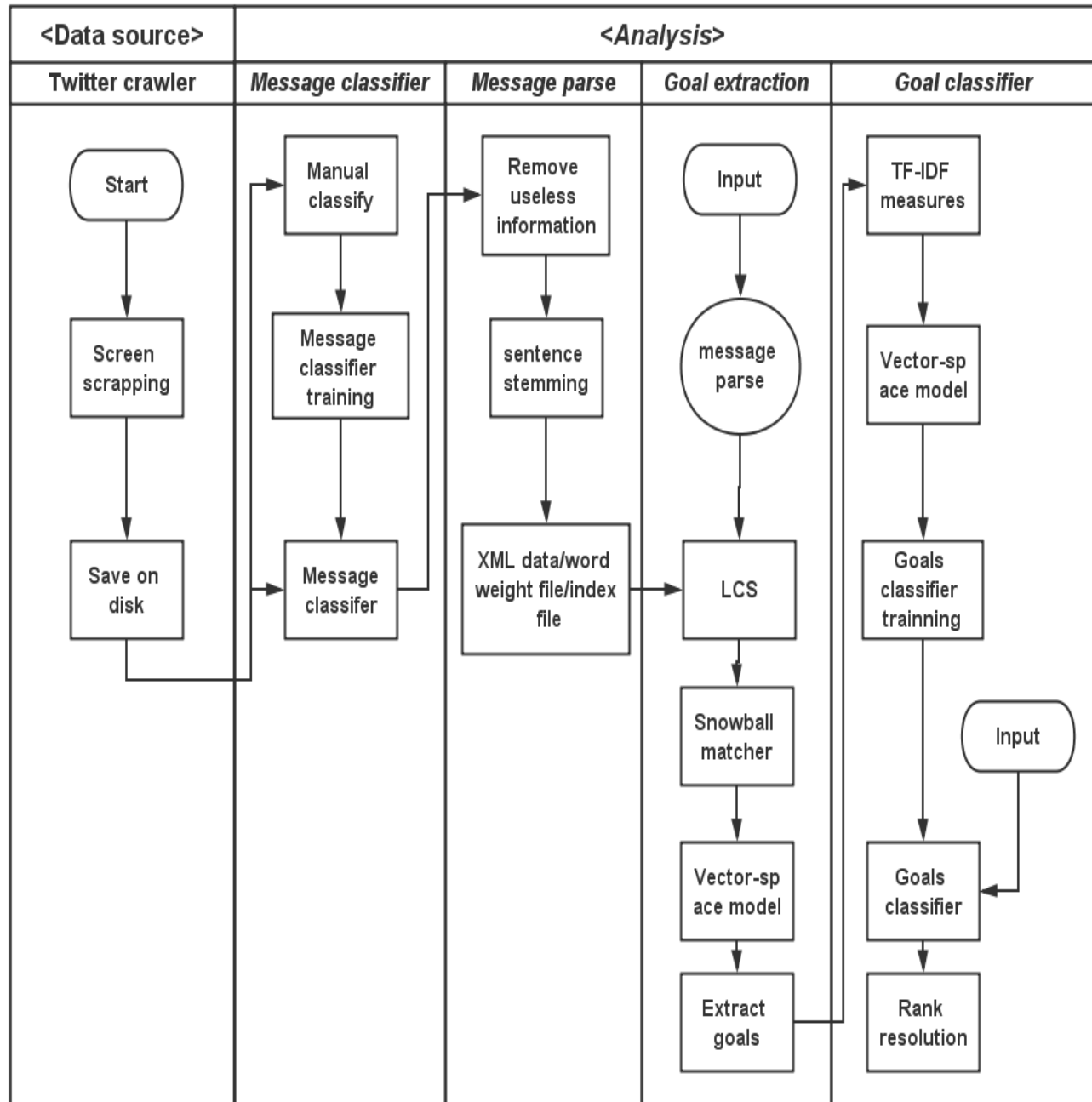
The sentimental classifier algorithm is good tool to classify information, such as Naive Bayes. The drawback of those algorithm is only constant classification of the data

source. In this paper, the “new year resolution” message has strong diversity which should be classified into variable number of classifications. “The Evaluation of Sentence Similarity Measures” study is used to solve this problem.

This system can be used in many area, such as advertisement, market analysis. For advertisement, the company will sell their product to their real customers who's new year resolution is related to the products. The users also will can use the analysis result to find friends who may already achiever their goal and get help.

## **2. Twitter message analysis system**

The system is can be divide into two parts, data source and Analysis. The data source part will get the data source from web site. The Analysis part will extract information and get the final rank information. It also divided to five phases: Twitter crawler, Message classifier, Message parse, Goal extraction and Goal classifier. The details of each part will explain as following.



### Twitter crawler

In this phase, the system will crawler twitter message from twitter site. The search key word will be set to “new year resolution” and time will from 12-25-2013 to 01-05-2014. The screen scrapping technique will be used to get information, since the twitter API has limitation on the messages. By using twitter API, only recently one week message can be crawled. The “new year resolution” messages are more dense around new year. So twitter API will can not be used for the message time limitation. I use the “Scraper” to crawl message from twitter website, which is an Google chrome extension.

## ***Message classifier***

Message classifier is used to classify message into two sets: resolution and not resolution. The machine learning and sentence sentimental algorithm is evolved in this part. For the training set, the data source has been divided into two sets manually. Then lingpipe lib will used to train the classifier, which has implement sentence sentimental algorithm and training process algorithm.

Implement class: ResolutionClassifier.java

## ***Message parse***

This parse is used to filter sentence and get the useful part of each sentence. According to the study on the data source, the useful information for the analysis includes web link and punctuation, such as period, comma. But other punctuations will be retained since it may be a mark for the resolution such as colon. Those punctuations will be add space to split them from words. For example, "My new year resolution: quit smoke". The colon here is sign for user's goal. The message will be transfer to "My new year resolution : quit smoke". The regular expression will be used here to remove or modify information.

Then the message will be stemmed. The stemming process will increase the information matching accuracy. "The Porter Stemming algorithm" and implement is used to deal with this problem.

In this parse, the each message will be divide into five part for the latter phase "Goal extraction" manually. The divided information is used to training data set. The related algorithm will be used in "Goal extraction" phase. The word frequency will record, which will be used in "snowball matcher" algorithm of "Goal extraction" phase. The account of sentence also should be record, which will be used in "Goal classifier" phase.

The output file includes source XML data, word frequency and index file. The source XML date will include source sentence, dealed sentence and their five parts in XML format. The formation should be as following.

```

- <twitter>
  - <message orientation="true">
    - <source>
      Maybe Dr. Ralph's new year resolution will be to call more snow days
    </source>
    - <res>
      mayb dr ralph " s new year resolut will be to call more snow dai
    </res>
    <left score="24">mayb dr ralph " s </left>
    <resolu score="208">new year resolut</resolu>
    <mid score="124"> will be to </mid>
    <key score="3">call more snow dai</key>
    <right score="0"/>
  </message>
+ <message orientation="true"></message>
+ <message orientation="true"></message>
+ <message orientation="true"></message>
+ <message orientation="true"></message>
+ <message orientation="true"></message>

```

The root node is "twitter" node, which contains many children node "message". Each message node represents one twitter message. The message node contains seven children node, source, res, left, resolu, mid, key, right. The source node is contain the raw twitter message. The res node is the result after information filter. The sentence should be divide into five part by "resolu" and "key" part. The "resolu" node contain the flag word, which has same meaning with "new year resolution". The "key" node is user specific goals. If the "resolu" part appears earlier than the "key" part in the source sentence, the orientation attribution of message node will be set as "true"; otherwise, it will be set as "false".

The index file records each word and the account of documents which contain the word. The word frequency record each word and the account of each word in the whole data set.

Implement class: XMLWrite.java, wordDictionary.java, stem.java, filterInfor.java

## Goal extraction

In this part, the user specific goals are extract from input twitter message. The system will compare the test twitter message with the data source. The data source is used as pattern to test the test message. The Longest common subsequence (LCS) algorithm will be used here to solve this problem. Every message in data source has been divided into five part as shown in last phase, left, resolu, mid, key, right part. The "resolu" parts have highest similarity in each sentence because the search key is "new year resolution". Firstly, we will compare the test twitter message with the "resolu" part of one sentence in the data source using LCS algorithm. We will use the result subsequence to divide the test twitter message into thire part, left, resolu, right part. Then according to the orientation of the data source sentence, we will pick the left or right part to do

the next step. If the orientation is true, the mid and key part is latter than the “resolu” part. So we will compare the right part of tested message with the mid part of the data source sentence using LCS algorithm; otherwise, the left part of tested message should be used to compare. We chose the “mid” part to test because the mid part has stronger similarity than the rest part. It is designed according to the “Snowball algorithm”.

Snowball algorithm this system is used to test the organization and location of sentence. For each tested sentence, it will be divided into five part, left, organization, mid, location and right part. The named-entity tagger algorithm is used here to find the organization and location part. Then the sentence will be divided into five part. Then the single pass clustering technique will be used to compare the similarity of two sentence. It is defined as following:

**Definition 2** *The degree of match  $Match(t_P, t_S)$  between two 5-tuples  $t_P = \langle l_P, t_1, m_P, t_2, r_P \rangle$  (with tags  $t_1$  and  $t_2$ ) and  $t_S = \langle l_S, t'_1, m_S, t'_2, r_S \rangle$  (with tags  $t'_1$  and  $t'_2$ ) is defined as:*

$$Match(t_P, t_S) = \begin{cases} l_P \cdot l_S + m_P \cdot m_S + r_P \cdot r_S & \text{if the tags match} \\ 0 & \text{otherwise} \end{cases}$$

According to the similarity between tested sentence and data source, snowball system will generate new pattern and add it into the data source. The pattern included in left, mid and right part of tested sentence.

In my system, users' specific goals can not be tested by the named-entity tagger algorithm. So I use the “resolu” and mid part to divide sentence. The result of this phase is four part sentence, the left, resolut, mid, rest part. Compare to five part structure, the rest part include the users' specific goals and right part. Because it is hard to distinguish them, which has little influence on the final rank result.

To compare the similarity between two sentence, I use the single pass clustering technique. I use the frequency of each word on the data set as the weight of each word. If “smoke” appears 60 times on the source data set, the weight of “smoke” will be 60. If the word appear on the middle part it will be two times of word weight. If the word appear on the left or middle part, the value will be  $0.2 * (\text{word weight})$ . The mid part always more important than other part for the users' goal identify. Secondly, the word weight should be add together to get the weight of each part. The “mid” and “left” part will be add together. Then there are four parts with weight on both tested sentence and sentence from data source set. The single pass clustering technique will be used to computer the final similarity score of two sentence.

The tested sentence will be compared with every sentence of the data set. The highest score will be show the highest similarity. Then the most similarity sentence structure will be used to divide the tested sentence. The tested sentence should use the “resolu” and “mid” part to divided into four part. The rest part contain the users' specific goals.

Implement class: GoalsExtract.java

## **Goal classifier**

After we get the users' specific goals, we should classifier those goal and get the rank of each class. The sentence sentimental algorithm is good tools to classify, but the number of classes should be set before training and testing. The user's specific goals are different from different data set. So the algorithm is not suitable. We will test the goal similarity to classify by using sentence similarity algorithm. There are three classes of measures for sentence similarity, word overlap measures, TF-IDF measures and Linguistic Measures. The word overlap measure is mainly compare the common subsequence word between two sentence. In our system, we should compare the rest parts which contain many other words besides the users' specific goals. So the word overlap is not suitable. The Linguistic measures should compare the syntactic structure between two sentence. But the users' goals may have great different syntactic structure with the same meaning. Therefore, the Linguistic measures is not suitable.

TF-IDF this algorithm compare the overlap word between two sentence, considering the word importance for specific topic. For example, *"Suppose we have a set of English text documents and wish to determine which document is most relevant to the query "the brown cow". A simple way to start out is by eliminating documents that do not contain all three words "the", "brown", and "cow", but this still leaves many documents. To further distinguish them, we might count the number of times each term occurs in each document and sum them all together; the number of times a term occurs in a document is called its' term frequency.*

*However, because the term "the" is so common, this will tend to incorrectly emphasize documents which happen to use the word "the" more frequently, without giving enough weight to the more meaningful terms "brown" and "cow". The term "the" is not a good keyword to distinguish relevant and non-relevant documents and terms, unlike the less common words "brown" and "cow". Hence an inverse document frequency factor is incorporated which diminishes the weight of terms that occur very frequently in the document set and increases the weight of terms that occur rarely."*

TF-IDF algorithm definition:

The **inverse document frequency** is a measure of how much information the word provides, that is, whether the term is common or rare across all documents. It is the logarithmically scaled fraction of the documents that contain the word, obtained by dividing the total number of **documents** by the number of documents containing the term, and then taking the logarithm of that **quotient**.

$$\text{idf}(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

with

- $N$ : total number of documents in the corpus
- $|\{d \in D : t \in d\}|$ : number of documents where the term  $t$  appears (i.e.,  $\text{tf}(t, d) \neq 0$ ). If the term is not in the corpus, this will lead to a division-by-zero. It is therefore common to adjust the denominator to  $1 + |\{d \in D : t \in d\}|$ .

Mathematically the base of the log function does not matter and constitutes a constant multiplicative factor towards the overall result.

Then **tf-idf** is calculated as

$$\text{tfidf}(t, d, D) = \text{tf}(t, d) \times \text{idf}(t, D)$$

A high weight in **tf-idf** is reached by a high term **frequency** (in the given document) and a low document frequency of the term in the whole collection of documents; the weights hence tend to filter out common terms. Since the ratio inside the **idf**'s log function is always greater than or equal to 1, the value of **idf** (and **tf-idf**) is greater than or equal to 0. As a term appears in more documents, the ratio inside the logarithm approaches 1, bringing the **idf** and **tf-idf** closer to 0.

Before computing, we should get the common words of two rest part by using the **hashMap**, because the sequence of word is not importance. To compute the **TFIDF** value of the common words, we will use the document account of each word, which we get from the "message parse" phase. After we get the **tfidf** value of each word in each users' specific goal, we will use the vector-space model to compute the similarity of two part. If the result of vector-spece is bigger than 0.05, the two goals belong to the same class.

The classify process starts from the first sentence in the date set. The first sentence will be used to compare with the rest sentence and get the similar goals. Then remove the first sentence and other sentence with the similar goals with the first sentence. After that, the second sentence will used to test the rest sentence and pick the second class out. This process will be iterate until the source data set is empty.

After this classify, we can use the users' specific goals part sets to training the sentence sedimentary classifier and get the classifier. If there are other parsed twitter message, it can be directly classify by the trained sentence classifier. Finally, we will get the total number of goals in each class and get the top 3 new year resolution.



Implement class:goalsCluster

### 3. Experimental setting

- Install the Scraper on chrome which link is in reference.  
Build ./data folder under the project folder.  
Put the download twitter message file in the data folder, which should be named as "resolution".
- Build folder under project folder.  
Built ./data/trainDirectory/res and ./data/trainDirectory/not folder. Put the resolution set and "not resolution" set into those direction. Those two set should be divided manually. It is used to training the resolution classifier.  
Build ./data/goalClassifierTrain folder. It will used to training the goal classifier.
- Get the uses' specific goals manually.  
Read each message and get the "new year resolution" phase and users' specific goals. Put those information into "resolutionKey" file with format as "new year resolution " part + tab + "users' specific goals" .
- import the project into Eclipse and run.

### 4. Experimental results

- Experience 1: Test the goals classifier:  
Search key word on twitter is "new year resolution **quit smoke** since:2013-12-25 until:2014-01-05 ", "new year resolution **lost weight** since:2013-12-25 until:2014-01-05" and "new year resolution **read book** since:2013-12-25 until:2014-01-05", Then I get 44 twitter messages for this search.  
Resolution classifier: 12 twitter messages are new year resolution  
Goals classifier: 10 twitter messages about lost weight and 2 message about read book.  
Conclusion: the goals are successful extract from the data source and all of them are belongs to right class.
- Experience 2: Test the resolution classifier:  
Search key word on twitter is "new year resolution since:2013-12-25 until:2014-01-05"  
Result set: get 918 twitter messages  
Wrong classifier: 16 new year resolutions are wrongly classified to the not-resolution set.
- Experience 3: Big data source:  
Search key word on twitter is "new year resolution since:2013-12-18 until:2013-12-25".  
Get 1468 twitter messages  
Search key word on twitter is "new year resolution since:2013-12-25 until:2014-01-05".  
Get 918 twitter messages

Search key word on twitter is "new year resolution since:2014-01-05 until:2014-01-11"  
Get 512 twitter messages  
Searchs key word on twitter is "new year resolution since:2012-12-25 until:2013-01-05"  
Get 888 results  
Searchs key word on twitter is "new year resolution since:2011-12-25 until:2012-01-05"  
Get 1815 results  
The total message number is 5601  
the resolution classify phase: 736ms  
the user's goal extraction phase: 4267ms  
the user's goal classify phase:1257ms

The resolution message number is 3518  
The users' goal extraction number is 1954  
All of thoes file are divided into 243 classes  
Top1 class: contain 350 messages  
the cluster resolution around topic:"Still looking for a New Year's resolution that's worth doing? Why not consider a commitment to assisting us with...  
<http://fb.me/1cgKHNR1X>"  
Top2 class: contain 283 messages, not sepicial resolution  
the cluster resolution around topic:"Palin New Year's Resolution: 'Be Even More Aggressive in Calling Out Media for Practicing Lapdog Laziness' <http://ift.tt/JSpCug> :  
" be even more aggress in call out media for practic lapdog lazi "  
Top3 class: contain 99 messages  
the cluster resolution around topic:"TN's New Year's Resolution: Hurt Workers, Undermine Prevailing Wage Law  
<http://thecontributor.com/prevailing-wage-no-more-tennessee> ... via  
@ContributorNews"  
Top4 class: contain 90 messages  
the cluster resolution around topic:"According to an inside source, Miley's new year's resolution is to stop smoking. <http://tumblr.co/ZkSLHu13ITifv>"  
Conclusion: The new year resolutions has great different between people. So each class do not have a central topic. The distribute of new year resolution is really sparse. The classifier of users' resolution are more suitable for the dense information such as experience 1.

## 5. Conclusion

This paper present twitter message extraction system for special topic "new year resolution". We introduced novel method to analysis the users' new year resolution and get the rank and distribution of those goals. In the future, the system can track how the

users archive their goals and analysis those information to get statistic data. The user suggestion is also good direction. The system will analysis the similarity between the achievement process of two users to get suggestion when the users run into trouble of their new year resolution. It will help people to become better themselves.

## 6. References

- Snowball: Extracting Relations from Large Plain-Text Collections. Eugene Agichtein, Luis Gravano, 2000.
- Sergey Brin. Extracting patterns and relations from the World-Wide Web. In Proceedings of the 1998 International Work-shop on the Web and Databases(WebDB'98), March 1998.
- The Evaluation of Sentence Similarity Measures . Palakorn Achananuparp, Xiaohua Hu<sup>1</sup> , and Shen Xiajiong
- Scraper:  
<https://chrome.google.com/webstore/detail/scraper/mbigbapnjcgafohmbkdlecaccpngjd>
- lingpipe lib: <http://alias-i.com/lingpipe/>
- The Porter Stemming algorithm: <http://tartarus.org/martin/PorterStemmer/>
- Single Pass Clustering Technique:  
<http://facweb.cs.depaul.edu/mobasher/classes/csc575/assignments/single-pass.html>
- TF-IDF: <http://en.wikipedia.org/wiki/Tf%E2%80%93idf>