

# Comment Extraction from Blog Posts and Its Applications to Opinion Mining

Huan-An Kao, Hsin-Hsi Chen

Department of Computer Science and Information Engineering  
National Taiwan University, Taipei, Taiwan  
E-mail: hhchen@ntu.edu.tw

## Abstract

Blog posts containing many personal experiences or perspectives toward specific subjects are useful. Blogs allow readers to interact with bloggers by placing comments on specific blog posts. The comments carry viewpoints of readers toward the targets described in the post, or supportive/non-supportive attitude toward the post. Comment extraction is challenging due to that there does not exist a unique template among all blog service providers. This paper proposes methods to deal with this problem. Firstly, the repetitive patterns and their corresponding blocks are extracted from input posts by pattern identification algorithm. Secondly, three filtering strategies, i.e., tag pattern loop filtering, rule overlap filtering, and longest rule first, are used to remove non-comment blocks. Finally, a comment/non-comment classifier is learned to distinguish comment blocks from non-comment blocks with 14 block-level features and 5 rule-level features. In the experiments, we randomly select 600 blog posts from 12 blog service providers. F-measure, recall, and precision are 0.801, 0.855, and 0.780, respectively, by using all of the three filtering strategies together with some selected features. The application of comment extraction to blog mining is also illustrated. We show how to identify the relevant opinionated objects – say, opinion holders, opinions, and targets, from posts.

## 1. Introduction

In recent years, blogs have become increasingly popular and have changed the style of communications on the Internet. Blogs allow readers to interact with bloggers by placing comments on specific blog posts. The commenting behavior not only implies the increasing popularity of a blog post, but also represents the interactions between an author and readers.

Due to the growing amount of blogs, many works such as blog search, summarization, opinion mining, *etc.*, have been investigated. Cao et al. (2008) showed that consideration of both post content and comment region achieves better retrieval performance in blog search. Hu et al. (2007) extracted sentences from post content and regarded them as summary of the blog post. Liu et al. (2007) mentioned that bloggers express their opinions on a particular subject through writing blog posts.

Identifying the boundary between post content and comment region, and extracting the comments in a region are fundamental for blog applications. Moreover, mining opinions in a blog post, author's opinions are not enough. It is necessary to consider both author's and readers' opinions toward the same topic.

To extract comments from blog posts is challenging. Each blog service provider has its own templates to present the information in comments. These templates do not have a general specification about what components must be provided in a comment or how many complete sub-blocks a comment is composed of.

This paper studies how to extract comments in blog posts and illustrates how to identify both author's and readers' opinions. Section 2 describes the system flow including the repetitive pattern identification, filtering strategies, and binary classification. Section 3 shows the experimental setup and evaluation. Section 4 applies the results of comment extraction to opinion mining.

## 2. Comment Extraction

### 2.1 System Flow

Given a blog post  $P$ , the task of comment extraction is to extract a set of comments  $C = \{c_1, c_2, \dots, c_n\}$  associated with  $P$ . A "site-level" approach gathers information from a designated blog service provider, parses the HTML contents, and identifies comment extraction rules manually. This approach suffers from human cost to formulate the rules and fail when a new blog site is first encountered. A "page-level" approach reads blog pages from different blog sites, and identifies the repetitive patterns embedded in the pages.

Figure 1 shows the architecture of our page-level approach. It includes an encoder which accepts an input post page, a repetitive pattern identifier which recognizes the repetitive patterns and the set of blocks, three filtering strategies which remove blocks with loop or overlap, and a comment/non-comment classifier which distinguishes comment blocks and non-comment blocks.

### 2.2 Repetitive Pattern Identification

HTML documents are composed of various kinds of tags carrying structure and presentation information, and text contents enwrapped by tags. Because our goal is to mine general comment structures, the information irrelevant to the document structures is not considered.

The input to pattern identification is an encoded string from an encoder. Each token in the string represents an HTML tag or a non-tag text. The algorithm scans the tokens. When encountering a token that is likely to be the head of a repetitive pattern (called a "rule" hereafter too), the subsequent tokens are examined if any rules can be formed.

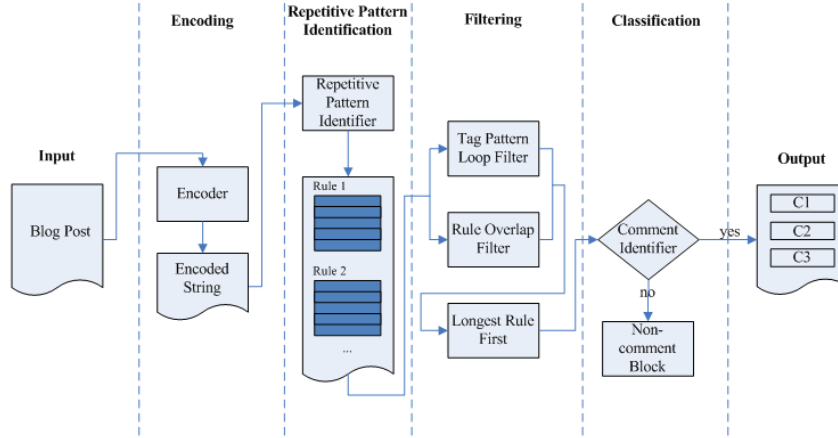


Figure 1: System Architecture

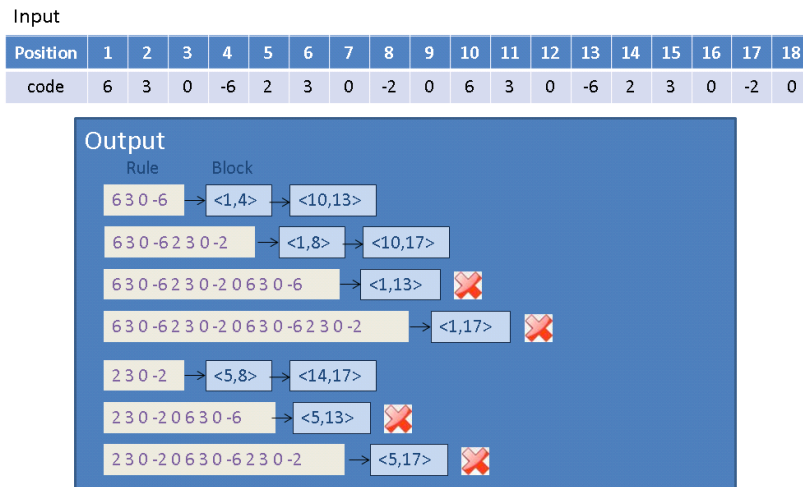


Figure 2: An Example of Repetitive Pattern Identification

HTML Tag	Code String	Class
Non-tag Text	0	string
STRUCTURE	1	begin
DIV CLASS	2 3	begin with attribute
DIV ID	2 4	begin with attribute
DT CLASS	5 3	begin with attribute
LI CLASS	6 3	begin with attribute
DD CLASS	7 3	begin with attribute
TR CLASS	8 3	begin with attribute
/STRUCTURE	-1	end
/DIV	-2	end
/DT	-5	end
/LI	-6	end
/DD	-7	end
/TR	-8	end
Other Tags	0	string

Table 1: Coding Scheme for Repetitive Pattern Identification

Figure 2 shows an encoded string '6 3 0 -6 2 3 0 -2 0 6 3 0 -6 2 3 0 -2 0' corresponding to an HTML document denoting a blog post. Table 1 lists all the tags and the corresponding codes used by the encoder. We mine seven rules and corresponding blocks from this string. We remove those rules (3<sup>rd</sup>, 4<sup>th</sup>, 6<sup>th</sup>, and 7<sup>th</sup>) with only one block, and keep the remaining repetitive patterns (1<sup>st</sup>, 2<sup>nd</sup> and 5<sup>th</sup>). Finally, 6 candidates are proposed and sent to the next stage.

### 2.3 Filtering Strategies

Since not all mined repetitive patterns are correct, non-comment blocks may be proposed wrongly. We present three filtering strategies, i.e., tag pattern loop filtering (M1), rule overlap filtering (M2) and the longest rule first (M3), to eliminate non-comment blocks. M1 and M2 are independent of each other, and M3 must be performed after M1 and M2. Figures 3-5 list an example

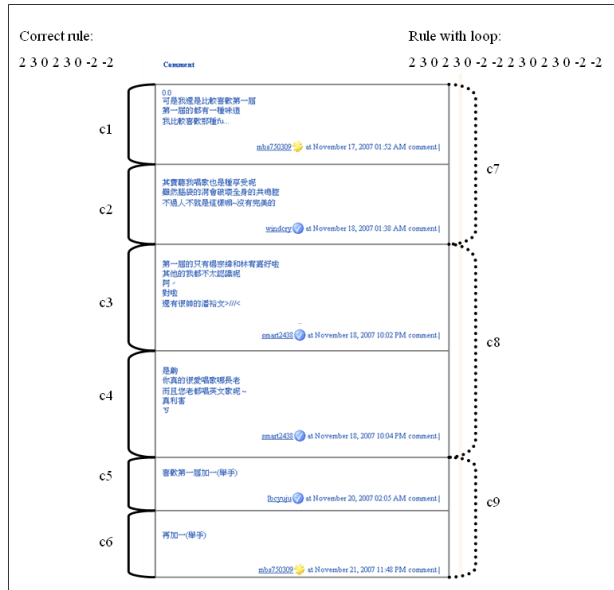


Figure 3: Correct Rule vs. Rule with Loop

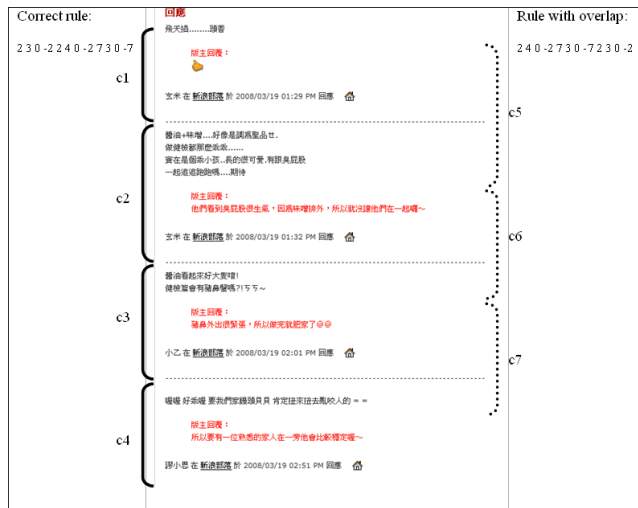


Figure 4: Correct Rule vs. Rule with Overlap

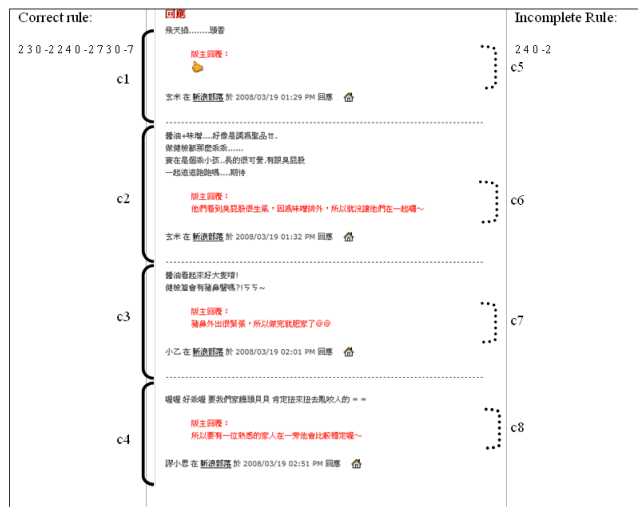


Figure 5: The Longest Rule First

for each case. The left braces list the correct boundaries of comments in the post shown in the middle. Some

incorrect blocks shown on the right braces are too large (Figure 3), overlap (Figure 4) or too small (Figure 5).

## 2.4 Binary Classification

A comment/non-comment classifier is based on features selected from both block-level and rule-level listed below.

### Block-Level Features

- (1) **Block length without tags.** Comments are shorter than post contents in blog posts on average.
- (2) **Block length with tags.** HTML tags are considered in the determination of block length.
- (3) **Number of words.** We consider block length in words instead of characters.
- (4) **Frequency of “comment” word.** The word “comment” often appears in comment blocks.
- (5) **Ratio of anchor tags.** Anchor tag contains a hyperlink to a page. This feature measures the anchor tag ratio in a block as (number of anchor tags) / (number of tags).
- (6) **Number of anchor tags.** This feature measures the number of anchor tags instead of ratio in a block.
- (7) **Ratio of stop words.** We postulate that stop word ratio may be lower in blogroll, categories, advertisements and other possible templates. In contrast, stop word ratio tends to be higher in comment blocks.
- (8) **Number of stop words.** This feature measures the number of stop words instead of ratio in a block.
- (9) **Ratio of punctuation marks.** We postulate that the punctuation ratio is higher in comment blocks than that in other blocks.
- (10) **Number of punctuation marks.** This feature measures the number of punctuation marks instead of ratio in a block.
- (11) **Block start position.** Comment region always appears after the post content. It seldom occurs in the top of the page.
- (12) **Block end position.** This feature is defined as the end position of a block divided by the length of whole post.
- (13) **Number of date and time expressions.** Reader responses always accompany with time and date expressions. This feature counts the occurrences of date and time expressions in a block.
- (14) **Occurrence of date and time expressions.** This feature equals to 1 if date and time expressions occur in a block, 0 otherwise.

### Rule-Level Features

- (15) **Rule start position.** This feature captures the starting position of a rule used to divide a blog post into a post region and a comment region.
- (16) **Rule end position.** Rule end position models the end position of a comment region.
- (17) **Density.** Density measures the ratio of total length of comment blocks divided by the length of a rule.
- (18) **Coverage.** The length of a region compared to the whole blog post may provide information about whether it is a comment region.
- (19) **Regularity.** The space between each adjacent comment block is almost the same.

Support Vector Machine (SVM) is adopted to learn a comment/non-comment classifier with the selected features.

## 3. Experiments

### 3.1 Experimental Setup

Total 12 blog service providers listed in Table 2 are used to collect a corpus *CommentExtract* 1.0. Total 50 blog posts are selected from each service provider. We manually labeled each comment in a given blog post by considering two comment styles: comment blocks without or with author reply. For the former, besides comment content, some related information such as commenter name, date and time are included. For the latter, we regard a comment block as a composite of both a reader comment and an author reply. A labeled comment block must include comment contents of readers and author. Table 3 lists the statistics of the *CommentExtract* 1.0 corpus.

### 3.2 Evaluation

A 12-fold cross validation is conducted. Each fold comes from a blog site. The data from 11 blog sites are used for training, and the remaining site is for testing. Each fold contains 50 posts from the same site. Table 4 compares different combinations of the three filtering strategies. Employing all strategies together achieves the best.

Provider	URL
Wretch	<a href="http://www.wretch.cc/blog">http://www.wretch.cc/blog</a>
Yam	<a href="http://blog.yam.com">http://blog.yam.com</a>
Pixnet	<a href="http://www.pixnet.net/blg">http://www.pixnet.net/blg</a>
Roodo	<a href="http://blog.roodo.com">http://blog.roodo.com</a>
Blogspot	<a href="http://www.blogger.com/home">http://www.blogger.com/home</a>
Xuite	<a href="http://blog.xuite.net">http://blog.xuite.net</a>
Sina	<a href="http://blog.sina.com.tw">http://blog.sina.com.tw</a>
Yahoo	<a href="http://tw.blog.yahoo.com">http://tw.blog.yahoo.com</a>
China Times	<a href="http://blog.chinatimes.com">http://blog.chinatimes.com</a>
Udn	<a href="http://blog.udn.com">http://blog.udn.com</a>
MSN	<a href="http://home.services.spaces.live.com">http://home.services.spaces.live.com</a>
Oui	<a href="http://www.oui-blog.com">http://www.oui-blog.com</a>

Table 2: Blog Service Providers

Number of blog posts	600
Number of blog posts with comments	482
Number of comments	3,505
Mean # comments per blog post	5.8
Comment Length (with tag)	
Mean	906.1
Maximum	7,593
Minimum	140
Median	756
Comment Length (without tag)	
Mean	290.4
Maximum	6,905
Minimum	41
Median	206

Table 3: Statistics of *CommentExtract* 1.0 Corpus

Strategy	Recall	Precision	F-measure
No Filter	0.607	0.166	0.221
M1	0.652	0.453	0.493
M2	0.663	0.206	0.256
M3	0.695	0.347	0.387
M1+M2	0.646	0.520	0.526
M2+M3	0.694	0.416	0.452
M1+M3	0.660	0.687	0.640
M1+M2+M3	<b>0.717</b>	<b>0.793</b>	<b>0.715</b>

Table 4: Comparisons of Different Filtering Strategies

Recall that we propose 14 block-level features and 5 rule-level features to discriminate comment blocks from non-comment blocks. To examine which features are critical, we remove a feature from the feature set one at a time, repeat the same training and testing procedure, and tell out the performance differences. In total, there are 19 experiments on 12-fold cross validation. Table 5 shows F-measure of classifiers after a feature being removed. Except that features 5 and 15 do not result in clear performance difference, removing features 4, 13, 14, 16 or 17 lower the average performance, and removing features 1, 2, 3, 6, 7, 8, 9, 10, 11, 12, 18, or 19 increases the average performance. The former features may be important for improving the performance because the performance decreases when these features are removed. When only they are used, the F-measure is improved from 0.715 to 0.781 compared to using all features.

We also employ feature scores for feature selection (Chen and Lin, 2005). Feature score measures the discrimination of two sets of real numbers. For each feature, its values of positive and negative instances in

training data can be used to assess if this feature is discriminative. Given a training vector  $x_i$ , its elements are all values extracted by the  $i$ -th feature. Now we have 19 features, so that  $i$  can be 1 to 19. Feature score for  $i$ -th feature  $F(i)$  is defined as follows.

$$F(i) = \frac{(\bar{x}_i^p - \bar{x}_i) + (\bar{x}_i^n - \bar{x}_i)}{\frac{1}{n_p - 1} \sum_{k=1}^{n_p} (x_{k,i}^p - \bar{x}_i^p)^2 + \frac{1}{n_n - 1} \sum_{k=1}^{n_n} (x_{k,i}^n - \bar{x}_i^n)^2}$$

where  $n_p$  and  $n_n$  denote the number of positive and negative examples, respectively;  $\bar{x}_i$ ,  $\bar{x}_i^p$ ,  $\bar{x}_i^n$  is the average of the  $i$ -th feature of the whole, positive, and negative data sets, respectively;  $x_{k,i}^p$  is the  $i$ -th feature of the  $k$ -th positive example, and  $x_{k,i}^n$  is the  $i$ -th feature of the  $k$ -th negative example.

Table 6 lists the feature score of each feature and the corresponding rank. The top-3 discriminative features are *occurrence of date and time expressions*, *number of date and time expressions*, and *density*. They also belong to the positive feature set selected by the approach of removing one feature at a time. The next top two features are *ratio of stop words* and *frequency of "comment" word*. Table 7 presents the performance of overall system and comment/non-comment classifier in the same 12-fold cross validation with different feature sets. When features 14, 13, 17, 7 and 4 (i.e., *occurrence of date and time expressions*, *number of date and time expressions*, *density*, *ratio of stop words*, and *frequency of "comment" word*) are adopted, F-measure is improved further to 0.801.

feature	Roodo	Wretch	Yahoo	Xuite	China	Yam	Pixnet	Sina	Oui	Blogspot	Udn	MSN	Average
All	0.790	0.639	0.751	0.459	0.965	0.783	0.581	0.900	0.662	0.841	0.338	0.871	0.715
<b>Removed Feature</b>													
1	0.866	0.956	0.802	0.533	<b>0.626</b>	0.811	0.591	<b>0.865</b>	0.742	<b>0.797</b>	0.363	<b>0.800</b>	0.729
2	0.806	0.724	0.774	<b>0.448</b>	<b>0.956</b>	0.795	<b>0.573</b>	<b>0.885</b>	0.717	<b>0.804</b>	0.401	<b>0.848</b>	0.728
3	0.810	0.746	0.763	<b>0.444</b>	<b>0.956</b>	0.796	<b>0.572</b>	<b>0.893</b>	0.726	<b>0.819</b>	0.342	<b>0.846</b>	0.726
4	<b>0.779</b>	0.798	0.792	0.604	<b>0.579</b>	0.804	<b>0.553</b>	<b>0.825</b>	0.740	<b>0.598</b>	<b>0.317</b>	<b>0.741</b>	<b>0.677</b>
5	0.856	0.909	0.789	0.492	<b>0.621</b>	0.811	0.585	<b>0.871</b>	0.746	<b>0.765</b>	0.400	<b>0.762</b>	0.717
6	0.808	0.787	0.763	<b>0.446</b>	<b>0.959</b>	0.798	<b>0.573</b>	0.899	0.726	<b>0.817</b>	0.338	<b>0.841</b>	0.729
7	<b>0.639</b>	0.972	0.757	0.458	<b>0.924</b>	0.801	0.587	<b>0.879</b>	0.733	<b>0.736</b>	0.471	0.901	0.738
8	0.807	0.746	0.763	<b>0.444</b>	<b>0.959</b>	0.798	<b>0.572</b>	0.899	0.724	<b>0.820</b>	<b>0.331</b>	<b>0.846</b>	0.726
9	<b>0.785</b>	0.746	0.762	<b>0.446</b>	<b>0.953</b>	0.796	<b>0.572</b>	0.899	0.721	<b>0.818</b>	0.357	<b>0.839</b>	0.724
10	0.810	0.787	0.763	<b>0.446</b>	<b>0.959</b>	0.798	<b>0.572</b>	<b>0.893</b>	0.724	<b>0.821</b>	0.339	<b>0.841</b>	0.729
11	0.804	0.853	0.769	<b>0.447</b>	<b>0.959</b>	0.796	<b>0.573</b>	0.896	0.725	<b>0.817</b>	0.340	<b>0.848</b>	0.735
12	0.813	0.825	0.769	<b>0.448</b>	<b>0.959</b>	0.798	<b>0.572</b>	0.898	0.724	<b>0.809</b>	0.345	<b>0.848</b>	0.734
13	<b>0.413</b>	<b>0.356</b>	<b>0.700</b>	<b>0.208</b>	<b>0.458</b>	<b>0.440</b>	<b>0.574</b>	<b>0.435</b>	<b>0.596</b>	<b>0.576</b>	<b>0.104</b>	<b>0.316</b>	<b>0.431</b>
14	0	0	0	0	0	0	0	0	0	0	0	0	0
15	<b>0.762</b>	0.816	0.770	<b>0.430</b>	<b>0.953</b>	<b>0.774</b>	<b>0.570</b>	<b>0.886</b>	0.716	<b>0.788</b>	<b>0.316</b>	<b>0.808</b>	0.716
16	<b>0.722</b>	<b>0.337</b>	<b>0.716</b>	<b>0.436</b>	0.969	0.792	<b>0.558</b>	<b>0.891</b>	0.726	<b>0.646</b>	<b>0.331</b>	<b>0.813</b>	<b>0.661</b>
17	<b>0.618</b>	<b>0.318</b>	<b>0.712</b>	<b>0.413</b>	<b>0.890</b>	<b>0.706</b>	0.628	<b>0.723</b>	0.663	<b>0.769</b>	<b>0.320</b>	<b>0.685</b>	<b>0.620</b>
18	<b>0.712</b>	0.862	0.758	<b>0.431</b>	<b>0.960</b>	0.783	<b>0.563</b>	0.920	0.736	<b>0.776</b>	<b>0.328</b>	<b>0.846</b>	0.723
19	<b>0.557</b>	0.904	0.770	0.547	<b>0.946</b>	0.790	<b>0.561</b>	0.953	0.741	<b>0.813</b>	<b>0.305</b>	<b>0.825</b>	0.726

Table 5: Comparisons of Different Features Using Removing One Feature at a Time

Rank	Id	Feature	Feature Score
1	14	Occurrence of date and time expressions	2.908
2	13	Number of date and time expressions	0.630
3	17	Density	0.287
4	7	Ratio of stop words	0.250
5	4	Frequency of “comment” word	0.190
6	15	Rule start position	0.113
7	8	Number of stop words	0.102
8	3	Number of words	0.092
9	10	Number of punctuation marks	0.089
10	19	Regularity	0.087
11	1	Block length without tags	0.083
12	18	Coverage	0.076
13	5	Ratio of anchor tags	0.052
14	9	Ratio of punctuation marks	0.027
15	12	Block end position	0.022
16	11	Block start position	0.019
17	2	Block length with tags	0.014
18	6	Number of anchor tags	0.011
19	16	Rule end position	0.003

Table 6: Feature Scores and Rank of Each Feature


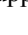
Number of Features	Overall			Classifier		
	Recall	Precision	F-measure	Recall	Precision	F-measure
19	0.668	0.793	0.682	0.717	0.793	0.715
18	0.597	0.793	0.627	0.636	0.793	0.656
17	0.604	0.793	0.632	0.643	0.793	0.661
16	0.622	0.795	0.644	0.662	0.795	0.674
15	0.640	0.797	0.661	0.683	0.797	0.692
14	0.644	0.797	0.663	0.687	0.797	0.693
13	0.661	0.799	0.677	0.706	0.799	0.709
12	0.650	0.795	0.669	0.693	0.795	0.699
11	0.686	0.798	0.703	0.736	0.798	0.737
10	0.696	<b>0.800</b>	0.712	0.748	<b>0.800</b>	0.746
9	0.628	0.784	0.629	0.671	0.784	0.657
8	0.639	0.785	0.640	0.682	0.785	0.668
7	0.648	0.786	0.651	0.693	0.786	0.681
6	0.656	0.786	0.656	0.701	0.786	0.686
5	<b>0.793</b>	0.780	<b>0.766</b>	<b>0.855</b>	0.780	<b>0.801</b>
4	0.774	0.737	0.730	0.836	0.737	0.763
3	0.782	0.758	0.743	0.845	0.758	0.776
2	0	0	0	0	0	0
1	0	0	0	0	0	0

Table 7: Comparison of Different Number of Features Using Feature Scores

#### 4. Application on Opinion Mining

After comment extraction, the opinions in each comment and the amount of comments which indicates the polarity tendency in each post can be presented to users. Typical opinionated information contains three basic components: an opinion holder, an opinion, and a target. The author of a blog post is the opinion holder of the post content and the reader who writes a comment is the opinion holder of this comment. The opinions are actually viewpoints or attitudes expressed in post content and each comment. A target can be a product, an event, a person, an organization, etc. It is usually specified in blog post.

We adopted opinion mining algorithms proposed by Ku and Chen (2007) to determine the opinion tendency of post content and the accompanying comments. They are

categorized into positive, negative, or neutral for further applications. Figure 6 shows a user interface of blog search. The search results are categorized into positive, negative, and neutral, and the numbers of positive, negative, and neutral blog posts are also presented. For a blog post, the result shows its link, title, and snippets. Besides, the numbers of positive, negative, and neutral comments in a blog post are also summarized. Figure 7 shows a blog post with opinion information. The left side of this figure lists the original blog post and its right side the opinions of the post content and each comment. We can easily tell out the opinions of both the author and the readers by the up and down symbols, i.e.,  supportive and  not supportive.





Calendar for May 2008 showing days of the week and dates.

May 26, 2008

### 「兩岸交流」是「我國科技發展策略」？

作者：blackbox

看到這樣的新聞，心理覺得很沉重，原來，不只是經貿、文化、政治，連「我國未來科技發展策略」，都不是優先加強與先進的歐美日交流，而是透過「加強兩岸科技交流」！

中國時報 2008.05.24  
兩岸科技交流 首長下月會面

。。。對於我國未來科技發展策略，李羅權特別提到頂尖領域重點突破和加強兩岸科技交流。

。。。在兩岸科技交流方面，推動兩岸政府科技高層互訪，被國科會列為近期工作重點。據透露，中國國家自然科學基金委員會（NSFC）主任陳宜瑜（相當於我方國科會主委）將在六月以科學家的身分率團來台訪問，屆時兩岸科技首長很可能透過學術研討會交流溝通，首度進行「非官方場合」的高峰會。。。

和中國、越南、甚至非洲國家做科技交流是有一定的價值，只不過，有必要成為「我國未來科技發展策略」嗎？這到底是切確交流，還是把高科技成果拱手讓人？這樣一來，台灣還有多少優勢？

列已開發國家之林的台灣（或是中華民國）的「未來科技發展策略」，為什麼不是加強和科技先進國家科技交流？卻選擇與仍是開發中國家的中國？

這無涉政治上的統獨或本土論述，只是單單要求國科會講求科技專業，不要泛政治化地把與開發中國家的科技交流，當作我國未來科技發展策略。

千萬，不要讓台灣科技上的驕點，在泛政治化的表相中，就這就那地消失殆盡了！  
妙子的故鄉剪影：「兩岸交流」是「我國科技發展策略」？ - 樂多日誌

請按下列按鈕推薦這篇文章

推薦 44

妙子的Twitter with Friends

Twitter widget showing recent tweets from users like 'singjerv' and 'olinp'.

Options A Spring widgets

最新的記事

- 用心的TAIWAN信箋  
「兩岸交流」是「我國科技發展策略」？  
手忙腳亂  
「台灣人」如何能打中國  
「馬融」  
「有情的心，美好的仗」領書定點公佈  
「非廣告」上綠色和平接受訪問  
來自長仔的告白  
去不生英九，萬古如長夜？  
用微距離賞花  
「有情的心，美好的仗」助印

美食精華區

- 再訪好吃的土東市場 阿吉師  
榮白葡萄酒  
幸福美饗 - 哈露奇山旅行咖啡  
鐘哈咖啡：咖啡學：秘史、精品豆與烘焙入門 即將出版 - Yam天空部落  
東京「女性專屬酒吧」大受女性歡迎  
這是個Premium 部落格

相機精華區

- 超級望遠鏡，倫敦看紐約！  
教你認識115個副檔名  
影像管理真麻煩  
Canon 200mm f/2L IS 香港開賣 (HK\$50,290)  
masaru's vision: 經典迷你 RF 相機 - OLYMPUS XA

記事分類

- 攝影 (7)
- 其它 (70)
- 站務公告 (67)
- 媒體 (45)

### 回應文章

可能政府發現了中國些微的起伏  
認為應該降低買進了

6 positive comment  
2 negative comment  
0 neutral comment

馬也說講要把糧苗賣到中國  
幫幫忙，話可以說錯，要事

positive opinion  
score:0.27222

322之後，西方媒體大量報  
IP合作案的西方朋友私下問  
國的保護，在某些敏感尖端

positive opinion  
score:0.02781

另外，一位西方科技官員的  
術帶到中國去進行合作，他的

comment#3  
score:2.11502

這些國家若要和中國交流，  
中國教授留學生很多，有些

comment#5  
positive opinion

台灣天天被統媒用「經營立  
把自己的視野鎖在中國了！

comment#6  
score:0.01523

真不知道  
這兩岸科技交流是在交流這  
只要我們的技術夠好

comment#7  
negative opinion  
score:-0.65875

對於馬政府把所有雞蛋都  
實在是浪費等Y----

comment#8  
negative opinion  
score:-0.75011

截止台灣的科技即將淪陷，  
連文化也即將不保了，若不見，

comment#9  
negative opinion  
score:-0.75011

comment#10  
negative opinion  
score:-0.75011

comment#11  
negative opinion  
score:-0.75011

comment#12  
negative opinion  
score:-0.75011

comment#13  
negative opinion  
score:-0.75011

comment#14  
negative opinion  
score:-0.75011

comment#15  
negative opinion  
score:-0.75011

comment#16  
negative opinion  
score:-0.75011

comment#17  
negative opinion  
score:-0.75011

comment#18  
negative opinion  
score:-0.75011

comment#19  
negative opinion  
score:-0.75011

Figure 7: A Blog Post with Opinion Information