

## KEYWORD EXTRACTION FROM A SINGLE DOCUMENT USING WORD CO-OCCURRENCE STATISTICAL INFORMATION

Y. MATSUO

*National Institute of Advanced Industrial Science and Technology*  
*y.matsuo@aist.go.jp*

M. ISHIZUKA

*University of Tokyo*  
*ishizuka@miv.t.u-tokyo.ac.jp*

Received 18 July 2003

Revised 19 October 2003

Accepted 19 October 2003

We present a new keyword extraction algorithm that applies to a single document without using a corpus. Frequent terms are extracted first, then a set of co-occurrences between each term and the frequent terms, i.e., occurrences in the same sentences, is generated. Co-occurrence distribution shows importance of a term in the document as follows. If the probability distribution of co-occurrence between term  $a$  and the frequent terms is biased to a particular subset of frequent terms, then term  $a$  is likely to be a keyword. The degree of bias of a distribution is measured by the  $\chi^2$ -measure. Our algorithm shows comparable performance to *tfidf* without using a corpus.

*Keywords:* Keyword extraction; co-occurrence;  $\chi^2$ -measure.

### 1. Introduction

Keyword extraction is an important technique for document retrieval, Web page retrieval, document clustering, summarization, text mining, and so on. By extracting appropriate keywords, we can easily choose which document to read to learn the relationship among documents. A popular algorithm for indexing is the *tfidf* measure, which extracts keywords that appear frequently in a document, but that don't appear frequently in the remainder of the corpus. The term "keyword extraction" is used in the context of text mining, for example<sup>15</sup>. A comparable research topic is called "automatic term recognition" in the context of computational linguistics and "automatic indexing" or "automatic keyword extraction" in information retrieval research.

Recently, numerous documents have been made available electronically. Domain-independent keyword extraction, which does not require a large corpus, has many applications. For example, if one encounters a new Web page, one might like to know

the contents quickly by some means, e.g., by having the keywords highlighted. If one wants to know the main assertion of a paper, one would want to have some keywords. In these cases, keyword extraction without a corpus of the same kind of documents is very useful. Word count <sup>8</sup> is sometimes sufficient for document overview; however, a more powerful tool is desirable.

This paper explains a keyword extraction algorithm based solely on a single document. First, frequent terms are extracted. Co-occurrences of a term and frequent terms are counted. If a term appears frequently with a particular subset of terms, the term is likely to have important meaning. The degree of bias of the co-occurrence distribution is measured by the  $\chi^2$ -measure. We show that our keyword extraction performs well without the need for a corpus. In this paper, a term is defined as a word or a word sequence. We do not intend to limit the meaning in a terminological sense. A word sequence is written as a phrase.

This paper is organized as follows. The next section describes our idea of keyword extraction. We describe the algorithm in detail followed by evaluation and discussion. Finally, we summarize our contributions.

2. Term Co-occurrence and Importance

A document consists of sentences. In this paper, a sentence is considered to be a set of words separated by a stop mark (“.”, “?” or “!”). We also include document titles, section titles, and captions as sentences. Two terms in a sentence are considered to co-occur once. That is, we see each sentence as a “basket,” ignoring term order and grammatical information except when extracting word sequences.

We can obtain frequent terms by counting term frequencies. Let us take a very famous paper by Alan Turing <sup>20</sup> as an example. Table 1 shows the top ten frequent terms and the probability of occurrence, normalized so that the sum is to be 1 (i.e., normalized relative frequency). Next, a co-occurrence matrix is obtained by counting frequencies of pairwise term co-occurrences, as shown in Table 2. For example, term *a* and term *b* co-occur in 30 sentences in the document. Let *N* denote the number of different terms in the document. While the term co-occurrence matrix is an  $N \times N$  symmetric matrix, Table 2 shows only a part of the whole – an  $N \times 10$  matrix. We do not define diagonal components.

Assuming that term *w* appears independently from frequent terms (denoted as *G*), the distribution of co-occurrence of term *w* and the frequent terms is similar to the unconditional distribution of occurrence of the frequent terms shown in Table 1.

Table 1. Frequency and probability distribution.

Frequent term	a	b	c	d	e	f	g	h	i	j	Total
Frequency	203	63	44	44	39	36	35	33	30	28	555
Probability	0.366	0.114	0.079	0.079	0.070	0.065	0.063	0.059	0.054	0.050	1.0

<sup>a</sup>machine, <sup>b</sup>computer, <sup>c</sup>question, <sup>d</sup>digital, <sup>e</sup>answer, <sup>f</sup>game, <sup>g</sup>argument, <sup>h</sup>make, <sup>i</sup>state, <sup>j</sup>number

Table 2. A co-occurrence matrix.

	a	b	c	d	e	f	g	h	i	j	Total
a	—	30	26	19	18	12	12	17	22	9	165
b	30	—	5	50	6	11	1	3	2	3	111
c	26	5	—	4	23	7	0	2	0	0	67
d	19	50	4	—	3	7	1	1	0	4	89
e	18	6	23	3	—	7	1	2	1	0	61
f	12	11	7	7	7	—	2	4	0	0	50
g	12	1	0	1	1	2	—	5	1	0	23
h	17	3	2	1	2	4	5	—	0	0	34
i	22	2	0	0	1	0	1	0	—	7	33
j	9	3	0	4	0	0	0	0	7	—	23
...	...	...	...	...	...	...	...	...	...	...	...
u	6	5	5	3	3	18	2	2	1	0	45
v	13	40	4	35	3	6	1	0	0	2	104
w	11	2	2	1	1	0	1	4	0	0	22
x	17	3	2	1	2	4	5	0	0	0	34

u: imitation, v: digital computer, w: kind, x: make

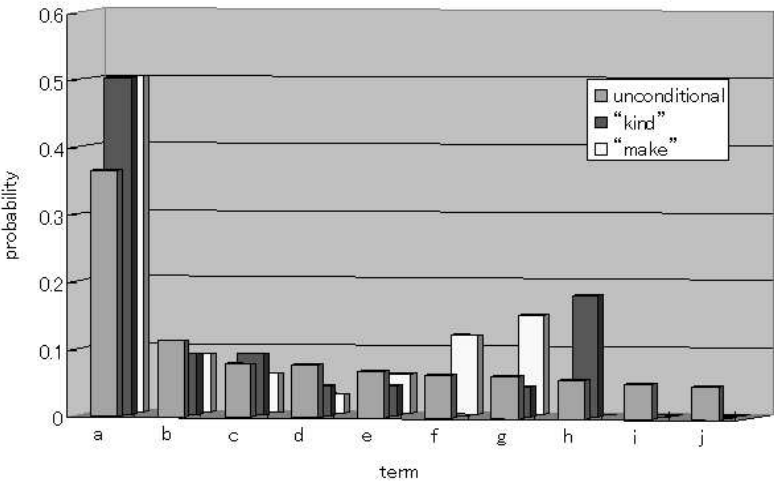


Fig. 1. Co-occurrence probability distribution of the terms “kind”, “make”, and frequent terms.

Conversely, if term  $w$  has a semantic relation with a particular set of terms  $g \in G$ , co-occurrence of term  $w$  and  $g$  is greater than expected, the distribution is said to be biased.

Figures 1 and 2 show the co-occurrence probability distribution of some terms and frequent terms. In the figures, unconditional distribution of frequent terms is shown as “unconditional”. A general term such as ‘kind’ or “make” is used relatively impartially with each frequent term, while a term such as “imitation” or “digital computer” shows co-occurrence with particular terms. These biases are derived from either semantic, lexical, or other relationships between two terms. Thus, a

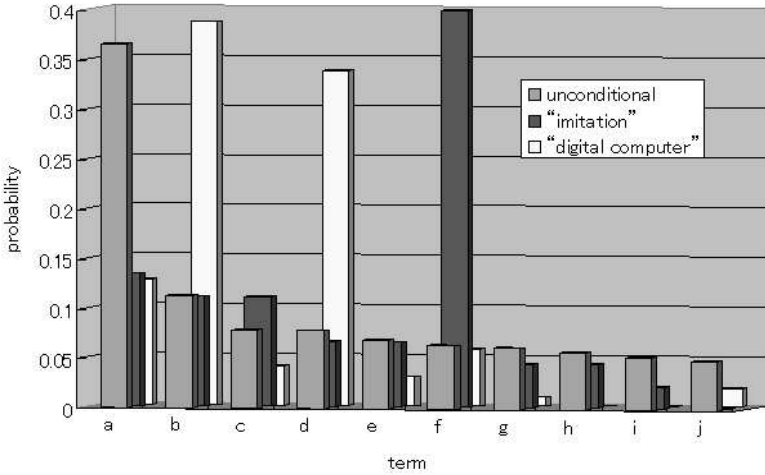


Fig. 2. Co-occurrence probability distribution of the terms “imitation”, “digital computer”, and frequent terms.

term with co-occurrence biases may have an important meaning in a document. In this example, “imitation” and “digital computer” are important terms, as we all know: In this paper, Turing proposed an “imitation game” to replace the question “Can machines think?”

Therefore, the degree of bias of co-occurrence can be used as an indicator of term importance. However, if term frequency is small, the degree of biases is not reliable. For example, assume term  $w_1$  appears only once and co-occurs only with term  $a$  once (probability 1.0). At the other extreme, assume term  $w_2$  appears 100 times and co-occurs only with term  $a$  100 times (with probability 1.0). Intuitively,  $w_2$  seems more reliably biased. In order to evaluate statistical significance of biases, we use the  $\chi^2$  test, which is very common for evaluating biases between expected frequencies and observed frequencies. For each term, frequency of co-occurrence with the frequent terms is regarded as a sample value; the null hypothesis is that “occurrence of frequent terms  $G$  is independent from occurrence of term  $w$ ,” which we expect to reject.

We denote the unconditional probability of a frequent term  $g \in G$  as the expected probability  $p_g$  and the total number of co-occurrences of term  $w$  and frequent terms  $G$  as  $n_w$ . Frequency of co-occurrence of term  $w$  and term  $g$  is written as  $freq(w, g)$ . The statistical value of  $\chi^2$  is defined as

$$\chi^2(w) = \sum_{g \in G} \frac{(freq(w, g) - n_w p_g)^2}{n_w p_g}. \quad (1)$$

If  $\chi^2(w) > \chi^2_\alpha$ , the null hypothesis is rejected with significance level  $\alpha$ . The term  $n_w p_g$  represents the expected frequency of co-occurrence; and  $(freq(w, g) - n_w p_g)$

Table 3. Terms with high  $\chi^2$  value.

Rank	$\chi^2$	Term	Frequency
1	593.7	digital computer	31
2	179.3	imitation game	16
3	163.1	future	4
4	161.3	question	44
5	152.8	internal	3
6	143.5	answer	39
7	142.8	input signal	3
8	137.7	moment	2
9	130.7	play	8
10	123.0	output	15
$\vdots$	$\vdots$	$\vdots$	$\vdots$
553	0.8	Mr.	2
554	0.8	sympathetic	2
555	0.7	leg	2
556	0.7	chess	2
557	0.6	Pickwick	2
558	0.6	scan	2
559	0.3	worse	2
560	0.1	eye	2

(We set the top ten frequent terms as  $G$ .)

represents the difference between observed and expected frequencies. Therefore, large  $\chi^2(w)$  indicates that co-occurrence of term  $w$  shows strong bias. In this paper, we use the  $\chi^2$ -measure as an index of biases, not for tests of hypotheses.

Table 3 shows terms with high  $\chi^2$  values and ones with low  $\chi^2$  values in Turing’s paper. Generally, terms with large  $\chi^2$  are relatively important in the document; terms with small  $\chi^2$  are relatively trivial. The table excludes terms whose frequency is less than two. However, we don’t have to define such a threshold, because low frequency usually indicates low  $\chi^2$  value (unless  $n_wp_g$  is very large, which is quite unusual.)

In summary, our algorithm first extracts frequent terms as a “standard”; then it extracts terms with high deviation from the standard as keywords.

3. Algorithm Description and Improvement

In the previous section, the basic idea of our algorithm is described. This section gives the precise algorithm description and two algorithm improvements: calculation of  $\chi^2$  value and clustering of terms. These improvements lead to better performance.

3.1. Calculation of  $\chi^2$  values

To improve the calculation of the  $\chi^2$  value, we focus on two aspects: variety of sentence length and robustness of the  $\chi^2$  value.

Int. J. Artif. Intell. Tools 2004.13:157-169. Downloaded from www.worldscientific.com by SWISS FEDERAL INSTITUTE OF on 02/05/14. For personal use only.

First, we consider the length of sentences. A document consists of sentences of various lengths. If a term appears in a long sentence, it is likely to co-occur with many terms; if a term appears in a short sentence, it is less likely to co-occur with other terms. We consider the length of each sentence and revise our definitions. We denote

- $p_g$  as (the sum of the total number of terms in sentences where  $g$  appears) divided by (the total number of terms in the document),
- $n_w$  as the total number of terms in sentences where  $w$  appears.

Again  $n_w p_g$  represents the expected frequency of co-occurrence. However, its value becomes more precise.<sup>a</sup>

Second, we consider the robustness of the  $\chi^2$  value. A term co-occurring with a particular term  $g \in G$  has a high  $\chi^2$  value. However, these terms are sometimes adjuncts of term  $g$  and not important terms. For example, in Table 3, a term “future” or “internal” co-occurs selectively with the frequent term “state,” because these terms are used in the form of “future state” and “internal state.” Though  $\chi^2$  values for these terms are high, “future” and “internal” themselves are not important. Assuming that “state” is not a frequent term,  $\chi^2$  values of these terms diminish rapidly.

We use the following function to measure robustness of bias values; it subtracts the maximal term from the  $\chi^2$  value,

$$\chi'^2(w) = \chi^2(w) - \max_{g \in G} \left\{ \frac{(\text{freq}(w, g) - n_w p_g)^2}{n_w p_g} \right\}. \quad (2)$$

Using this function, we can estimate  $\chi'^2(w)$  as low if  $w$  co-occurs selectively with only one term. It will have a high value if  $w$  co-occurs selectively with more than one term.

### 3.2. Clustering of terms

Some terms co-occur with each other and clusters of terms are obtained by combining co-occurring terms. Below we show how to calculate the  $\chi^2$  value more reliably by clustering terms.

A co-occurrence matrix is originally an  $N \times N$  matrix, where columns corresponding to frequent terms are extracted for calculation. We ignore the remaining columns, i.e., co-occurrence with low frequency terms, because it is difficult to estimate precise probability of occurrence for low frequency terms.

To improve extracted keyword quality, it is very important to select the proper set of columns from a co-occurrence matrix. The set of columns is preferably orthogonal; assuming that terms  $g_1$  and  $g_2$  appear together very often, co-occurrence

<sup>a</sup> $p_g$  is the probability of a term in a document to co-occur with  $g$ . Each term can co-occur with multiple terms. Therefore, the sum of  $p_g$  for all terms is not 1.0 but the average number of frequent terms in a sentence.

Table 4. Two transposed columns.

	a	b	c	d	e	f	g	h	i	j	...
c	<b>26</b>	<b>5</b>	—	<b>4</b>	<i>23</i>	<b>7</b>	<b>0</b>	<b>2</b>	<b>0</b>	<b>0</b>	...
e	<b>18</b>	<b>6</b>	<i>23</i>	<b>3</b>	—	<b>7</b>	<b>1</b>	<b>2</b>	<b>1</b>	<b>0</b>	...

Table 5. Clustering of the top 49 frequent terms.

C1:	game, imitation, imitation game, play, programme
C2:	system, rules, result, important
C3:	computer, digital, digital computer
C4:	behaviour, random, law
C5:	capacity, storage
C6:	question, answer
...	...
C26:	human
C27:	state
C28:	learn

of terms  $w$  and  $g_1$  might imply the co-occurrence of  $w$  and  $g_2$ . Thus, term  $w$  will have a high  $\chi^2$  value; this is very problematic. It is straightforward to extract an orthogonal set of columns, however, to prevent the matrix from becoming too sparse, we will cluster terms (i.e., columns).

Many studies address term clustering. Two major approaches <sup>6</sup> are:

**Similarity-based clustering** If terms  $w_1$  and  $w_2$  have similar distribution of co-occurrence with other terms,  $w_1$  and  $w_2$  are considered to be in the same cluster.

**Pairwise clustering** If terms  $w_1$  and  $w_2$  co-occur frequently,  $w_1$  and  $w_2$  are considered to be in the same cluster.

Table 4 shows an example of two (transposed) columns extracted from a co-occurrence matrix. Similarity-based clustering centers upon boldface figures and pairwise clustering focuses on italic figures.

By similarity-based clustering, terms with the same role, e.g., “Monday,” “Tuesday,” ..., or “build,” “establish,” and “found” are clustered <sup>13</sup>. In our preliminary experiment, when applied to a single document similarity-based clustering groups paraphrases and a phrase and its component (e.g., “digital computer” and “computer”). Similarity of two distributions is measured statistically by Kullback-Leibler divergence or Jensen-Shannon divergence <sup>2</sup>.

On the other hand, pairwise clustering yields relevant terms in the same cluster: “doctor,” “nurse,” and “hospital” <sup>19</sup>. A frequency of co-occurrence or mutual information can be used to measure the degree of relevance <sup>1,3</sup>.

Our algorithm uses both types of clustering. First we cluster terms by a similarity measure (using Jensen-Shannon divergence); subsequently, we apply pairwise clustering (using mutual information). Table 5 shows an example of term cluster-

ing. Proper clustering of frequent terms results in an appropriate  $\chi^2$  value for each term. We don't take the size of the cluster into account for simplicity. Balancing the clusters may improve the algorithm performance.

Below, co-occurrence of a term and a cluster implies co-occurrence of the term and any term in the cluster.

### 3.3. Algorithm

The algorithm follows. Thresholds are determined by preliminary experiments.

1. Preprocessing: Stem words by Porter algorithm<sup>14</sup> and extract phrases based on the APRIORI algorithm<sup>5</sup>. We extract phrases of up to 4 words with frequency more than 3 times. Discard stop words included in stop list used in the SMART system<sup>16</sup>.
2. Selection of frequent terms: Select the top frequent terms up to 30% of the number of running terms,  $N_{total}$ .
3. Clustering frequent terms: Cluster a pair of terms whose Jensen-Shannon divergence is above the threshold ( $0.95 \times \log 2$ ). Jensen-Shannon divergence is defined as

$$J(w_1, w_2) = \log 2 + \frac{1}{2} \sum_{w' \in C} \{h(P(w'|w_1) + P(w'|w_2)) - h(P(w'|w_1)) - h(P(w'|w_2))\}$$

where  $h(x) = -x \log x$ ,  $P(w'|w_1) = \text{freq}(w', w_1) / \text{freq}(w_1)$ . Cluster a pair of terms whose mutual information is above the threshold ( $\log(2.0)$ ). Mutual information between  $w_1$  and  $w_2$  is defined as

$$\begin{aligned} M(w_1, w_2) &= \log \frac{P(w_1, w_2)}{P(w_1)P(w_2)} \\ &= \log \frac{N_{total} \text{freq}(w_1, w_2)}{\text{freq}(w_1) \text{freq}(w_2)}. \end{aligned}$$

Two terms are in the same cluster if they are clustered by either of the two clustering algorithms. The obtained clusters are denoted as  $C$ .

4. Calculation of expected probability: Count the number of terms co-occurring with  $c \in C$ , denoted as  $n_c$ , to yield the expected probability  $p_c = n_c / N_{total}$ .
5. Calculation of  $\chi'^2$  value: For each term  $w$ , count co-occurrence frequency with  $c \in C$ , denoted as  $\text{freq}(w, c)$ . Count the total number of terms in the sentences including  $w$ , denoted as  $n_w$ . Calculate the  $\chi'^2$  value following

$$\chi'^2(w) = \sum_{c \in G} \left\{ \frac{(\text{freq}(w, c) - n_w p_c)^2}{n_w p_c} \right\} - \max_{c \in G} \left\{ \frac{(\text{freq}(w, c) - n_w p_c)^2}{n_w p_c} \right\}.$$

6. Output keywords: Show a given number of terms having the largest  $\chi'^2$  value.

In this paper, we use both nouns and verbs because verbs or verb+noun are sometimes important for illustrating the content of the document. Of course, we can apply our algorithm only to nouns.



Table 6. Improved results of terms with high  $\chi^2$  value.

Rank	$\chi^2$	Term	Frequency
1	380.4	digital computer	63
2	259.7	storage capacity	11
3	202.5	imitation game	16
4	174.4	machine	203
5	132.2	human mind	2
6	94.1	universality	6
7	93.7	logic	10
8	82.0	property	11
9	77.1	mimic	7
10	77.0	discrete-state machine	17

Table 6 shows the results for Turing's paper. Important terms are extracted regardless of their frequencies.

#### 4. Evaluation

For information retrieval, index terms are evaluated by their retrieval performance, namely recall and precision. However, we claim that our algorithm is useful when a corpus is not available due to cost or time to collect documents, or in a situation where document collection is infeasible.

Keywords are sometimes attached to a paper; however, they are not defined in a consistent way. Therefore, we employ author-based evaluation. Twenty authors of technical papers in artificial intelligence research have participated in the experiment. For each author, we showed keywords extracted from his/her paper by  $tf$ (term frequency),  $tfidf^b$ , *KeyGraph*, and our algorithm. *KeyGraph*<sup>11</sup> is a keyword extraction algorithm which requires only a single document as does our algorithm. It calculates term weight based on term co-occurrence information and was recently used to analyze a variety of data in the context of *Chance Discovery*<sup>12</sup>.

All these methods use word stem, elimination of stop words, and extraction of phrases. Using each method we extracted, gathered, and shuffled the top 15 terms. Then, the authors were asked to check terms which they thought were important in the paper. Precision can be calculated by the ratio of the checked terms to 15 terms derived by each method. Furthermore, the authors were asked to select five (or more) terms which they thought were indispensable for the paper. Coverage of each method was calculated by taking the ratio of the indispensable terms included in the 15 terms to all the indispensable terms. It is desirable to have the indispensable term list beforehand. However, it is very demanding for authors to provide a keyword list without seeing a term list. In our experiment, we allowed authors to add any

<sup>b</sup>The corpus is 166 papers in JAIR (Journal of Artificial Intelligence Research) from Vol. 1 in 1993 to Vol. 14 in 2001. The  $idf$  is defined by  $\log(D/df(w)) + 1$ , where  $D$  is the number of all documents and  $df(w)$  is the number of documents including  $w$ .

Table 7. Precision and coverage for 20 technical papers.

	<i>tf</i>	<i>KeyGraph</i>	<b>ours</b>	<i>tfidf</i>
Precision	0.53	0.42	<b>0.51</b>	0.55
Coverage	0.48	0.44	<b>0.62</b>	0.61
Frequency index	28.6	17.3	<b>11.5</b>	18.1

Table 8. Results with respect to phrases.

	<i>tf</i>	<i>KeyGraph</i>	<b>ours</b>	<i>tfidf</i>
Ratio of phrases	0.11	0.14	<b>0.33</b>	0.33
Precision w/o phrases	0.42	0.36	<b>0.42</b>	0.45
Recall w/o phrases	0.39	0.36	<b>0.46</b>	0.54

terms in the paper to the indispensable term list (even if they were not derived by any of the methods.).

Results are shown in Table 7. For each method, precision was around 0.5. However, coverage using our method exceeds that of *tf* and *KeyGraph* and is comparable to that of *tfidf*; both *tf* and *tfidf* selected terms which appeared frequently in the document (although *tfidf* considers frequencies in other documents). On the other hand, our method can extract keywords even if they do not appear frequently. The frequency index in the table shows average frequency of the top 15 terms. Terms extracted by *tf* appear about 28.6 times, on average, while terms by our method appear only 11.5 times. Therefore, our method can detect “hidden” keywords. We can use the  $\chi^2$  value as a priority criterion for keywords because precision of the top 10 terms by our method is 0.52, that of the top 5 is 0.60, while that of the top 2 is as high as 0.72. Though our method detects keywords consisting of two or more words well, it is still nearly comparable to *tfidf* if we discard such phrases, as shown in Table 8.

Computational time of our method is shown in Figure 3. The system is implemented in C++ on a Linux OS, Celeron 333MHz CPU machine. Computational time increases approximately linearly with respect to the number of terms; the process completes itself in a few seconds if the given number of terms is less than 20,000.

5. Discussion and Related Work

Co-occurrence has attracted interest for a long time in computational linguistics. For example, co-occurrence in particular syntactic contexts is used for term clustering<sup>13</sup>. Co-occurrence information is also useful for machine translation: for example, Tanaka et al. uses co-occurrence matrices of two languages to translate an ambiguous term<sup>19</sup>. Co-occurrence is also used for query expansion in information retrieval<sup>17</sup>.

Weighting a term by occurrence dates back to the 1950s in the study by Luhn<sup>8</sup>.

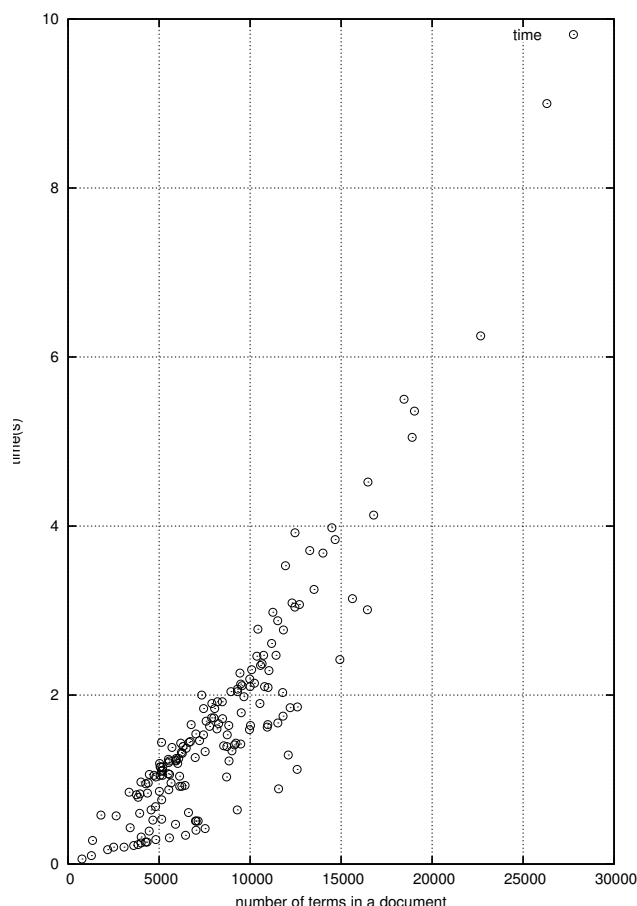


Fig. 3. Number of total terms and computational time.

More elaborate measures of term occurrence have been developed<sup>18,10</sup> by essentially counting term frequencies. Kageura and Umino summarized five groups of weighting measure<sup>7</sup>:

- (i) a word which appears in a document is likely to be an index term;
- (ii) a word which appears frequently in a document is likely to be an index term;
- (iii) a word which appears only in a limited number of documents is likely to be an index term for these documents;
- (iv) a word which appears relatively more frequently in a document than in the whole database is likely to be an index term for that document;
- (v) a word which shows a specific distributional characteristic in the database is likely to be an index term for the database.

Our algorithm corresponds to approach (v). Nagao used the  $\chi^2$  value to calculate

the weight of words<sup>9</sup>, which also corresponds to approach (v). But our method uses a co-occurrence matrix instead of a corpus, enabling keyword extraction using only the document itself.

From a probabilistic point of view, a method for estimating probability of previously unseen word combinations is important<sup>2</sup>. Several papers have addressed this issue, but our algorithm uses co-occurrence with frequent terms, which alleviates the estimation problem.

In the context of text mining, to discover keywords or keyword relationships is an important topic<sup>4,15</sup>. The general purpose of knowledge discovery is to extract implicit, previously unknown, and potentially useful information from data. Our algorithm can be considered a text mining tool in that it extracts important terms even if they are rare.

## 6. Conclusion

In this paper, we developed an algorithm to extract keywords from a single document. Main advantages of our method are its simplicity without requiring use of a corpus and its high performance comparable to *tfidf*. As more electronic documents become available, we believe our method will be useful in many applications, especially for domain-independent keyword extraction.

## References

- [1] K. W. Church and P. Hanks. Word association norms, mutual information, and lexicography. *Computational Linguistics*, 16(1):22, 1990.
- [2] I. Dagan, L. Lee, and F. Pereira. Similarity-based models of word cooccurrence probabilities. *Machine Learning*, 34(1):43, 1999.
- [3] T. Dunning. Accurate methods for the statistics of surprise and coincidence. *Computational Linguistics*, 19(1):61, 1993.
- [4] R. Feldman, M. Fresko, Y. Kinar, Y. Lindell, O. Liphstat, M. Rajman, Y. Schler, and O. Zamir. Text mining at the term level. In *Proceedings of the Second European Symposium on Principles of Data Mining and Knowledge Discovery*, page 65, 1998.
- [5] Johannes Fürnkranz. A study using n-grams features for text categorization. Technical Report OEFAI-TR-98-30, Austrian Research Institute for Artificial Intelligence, 1998.
- [6] Thomas Hofmann and Jan Puzicha. Statistical models for co-occurrence data. Technical Report AIM-1625, Massachusetts Institute of Technology, 1998.
- [7] K. Kageura and B. Umino. Methods of automatic term recognition. *Terminology*, 3(2):259, 1996.
- [8] H. P. Luhn. A statistical approach to mechanized encoding and searching of literary information. *IBM Journal of Research and Development*, 1(4):390, 1957.
- [9] M. Nagao, M. Mizutani, and H. Ikeda. An automated method of the extraction of important words from Japanese scientific documents. *Transactions of Information Processing Society of Japan*, 17(2):110, 1976.
- [10] T. Noreault, M. McGill, and M. B. Koll. *A Performance Evaluation of Similarity Measure, Document Term Weighting Schemes and Representations in a Boolean Environment*. Butterworths, London, 1977.

- [11] Y. Ohsawa, N. E. Benson, and M. Yachida. KeyGraph: Automatic indexing by co-occurrence graph based on building construction metaphor. In *Proceedings of the Advanced Digital Library Conference*, 1998.
- [12] Yukio Ohsawa. Chance discoveries for making decisions in complex real world. *New Generation Computing*, to appear.
- [13] F. Pereira, N. Tishby, and L. Lee. Distributional clustering of English words. In *Proceedings of the 31th Meeting of the Association for Computational Linguistics*, pages 183–190, 1993.
- [14] M. F. Porter. An algorithm for suffix stripping. *Program*, 14(3):130, 1980.
- [15] M. Rajman and R. Besancon. Text mining – knowledge extraction from unstructured textual data. In *Proceedings of the 6th Conference of International Federation of Classification Societies*, 1998.
- [16] G. Salton. *Automatic Text Processing*. Addison-Wesley, 1988.
- [17] H. Schutze and J. O. Pederson. A cooccurrence-based thesaurus and two applications to information retrieval. *Information Processing and Management*, 33(3):307–318, 1997.
- [18] K. Sparck-Jones. A statistical interpretation of term specificity and its application in retrieval. *Journal of Documentation*, 28(5):111, 1972.
- [19] K. Tanaka and H. Iwasaki. Extraction of lexical translations from non-aligned corpora. In *Proceedings of the 16th International Conference on Computational Linguistics*, page 580, 1996.
- [20] A. M. Turing. Computing machinery and intelligence. *Mind*, 59:433, 1950.