

WISDOM: A Web Information Credibility Analysis System

Susumu Akamine[†] Daisuke Kawahara[†] Yoshikiyo Kato[†]
Tetsuji Nakagawa[†] Kentaro Inui[†] Sadao Kurohashi^{†‡} Yutaka Kidawara[†]

[†]National Institute of Information and Communications Technology

[‡]Graduate School of Informatics, Kyoto University

{akamine, dk, ykato, tnaka, inui, kidawara}@nict.go.jp, kuro@i.kyoto-u.ac.jp

Abstract

We demonstrate an information credibility analysis system called WISDOM. The purpose of WISDOM is to evaluate the credibility of information available on the Web from multiple viewpoints. WISDOM considers the following to be the source of information credibility: information contents, information senders, and information appearances. We aim at analyzing and organizing these measures on the basis of semantics-oriented natural language processing (NLP) techniques.

1. Introduction

As computers and computer networks become increasingly sophisticated, a vast amount of information and knowledge has been accumulated and circulated on the Web. They provide people with options regarding their daily lives and are starting to have a strong influence on governmental policies and business management. However, a crucial problem is that the information available on the Web is not necessarily credible. It is actually very difficult for human beings to judge the credibility of the information and even more difficult for computers. However, computers can be used to develop a system that collects, organizes, and relativizes information and helps human beings view information from several viewpoints and judge the credibility of the information.

Information organization is a promising endeavor in the area of next-generation Web search. The search engine Clusty provides a search result clustering¹, and Cuil classifies a search result on the basis of query-related terms². The persuasive technology research project at Stanford University discussed how websites can be designed to influence people's perceptions (B. J. Fogg, 2003). However, as per our knowledge, no research has been carried out for supporting the human judgment on information credibility and information organization systems for this purpose.

In order to support the judgment of information credibility, it is necessary to extract the background, facts, and various opinions and their

distribution for a given topic. For this purpose, syntactic and discourse structures must be analyzed, their types and relations must be extracted, and synonymous and ambiguous expressions should be handled properly.

Furthermore, it is important to determine the identity of the information sender and his/her specialty as criteria for credibility, which require named entity recognition and total analysis of documents.

In this paper, we describe an information credibility analysis system called WISDOM, which automatically analyzes and organizes the above aspects on the basis of semantically oriented NLP techniques. WISDOM currently operates over 100 million Japanese Web pages.

2. Overview of WISDOM

We consider the following three criteria for the judgment of information credibility.

- (1) Credibility of information contents,
- (2) Credibility of the information sender, and
- (3) Credibility estimated from the document style and superficial characteristics.

In order to help people judge the credibility of information from these viewpoints, we have been developing an information analysis system called WISDOM. Figure 1 shows the analysis result of WISDOM on the analysis topic "Is bio-ethanol good for the environment?" Figure 2 shows the system architecture of WISDOM.

Given an analysis topic (query), WISDOM sends the query to the search engine TSUBAKI (Shinzato et al., 2008), and TSUBAKI returns a list of the top N relevant Web pages (N is usually set to 1000).

Then, those pages are automatically analyzed, and major and contradictory expressions and evaluative expressions are extracted. Furthermore, the information senders of the Web pages, which were analyzed beforehand, are collected and the distribution is calculated.

The WISDOM analysis results can be viewed from several viewpoints by changing the tabs using a Web browser. The leftmost tab, "Summary," shows the summary of the analysis, with major phrases and major/contradictory statements first.

¹ <http://clusty.com/>, <http://clusty.jp/>

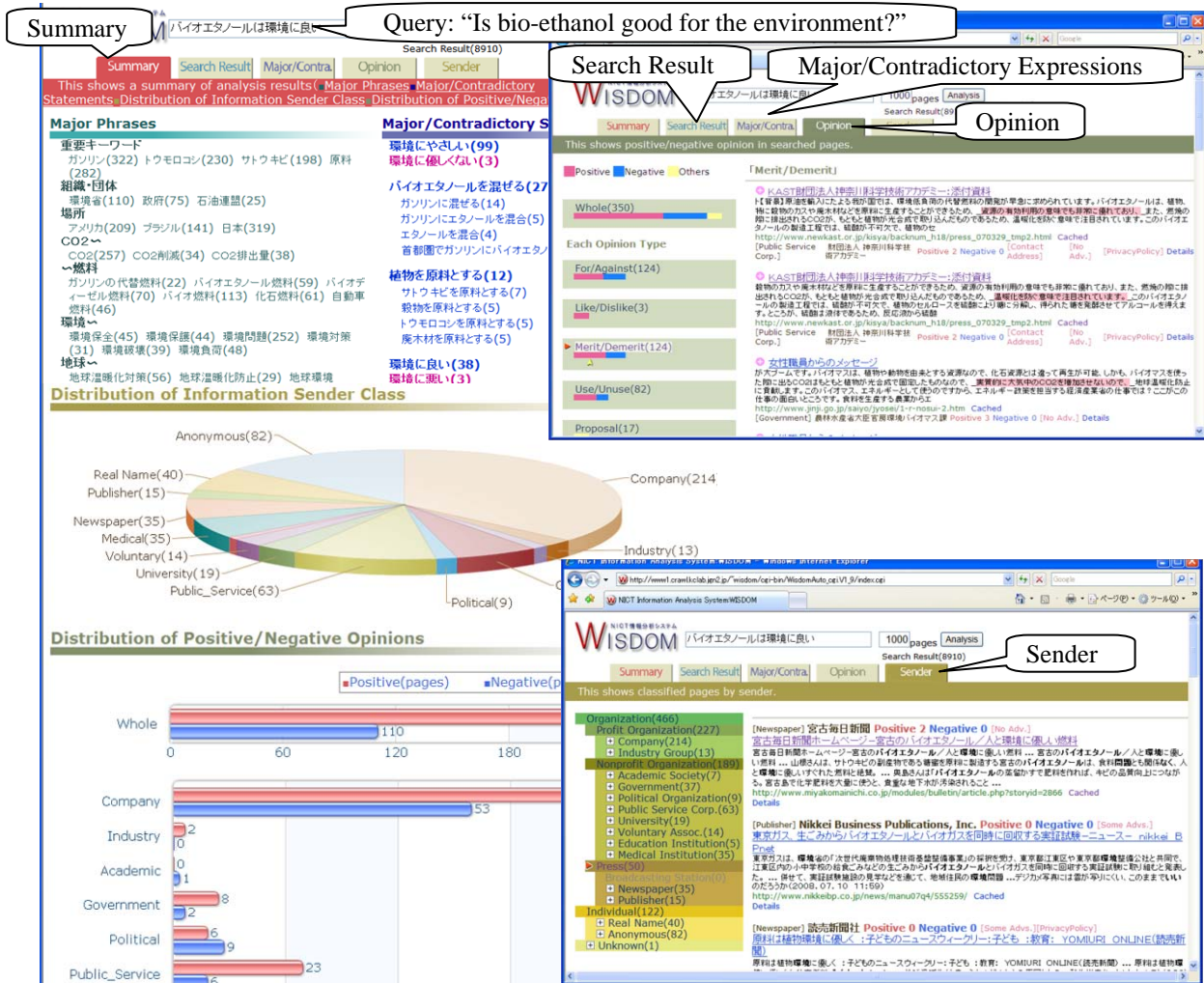


Figure 1. An analysis example of the information credibility analysis system WISDOM.

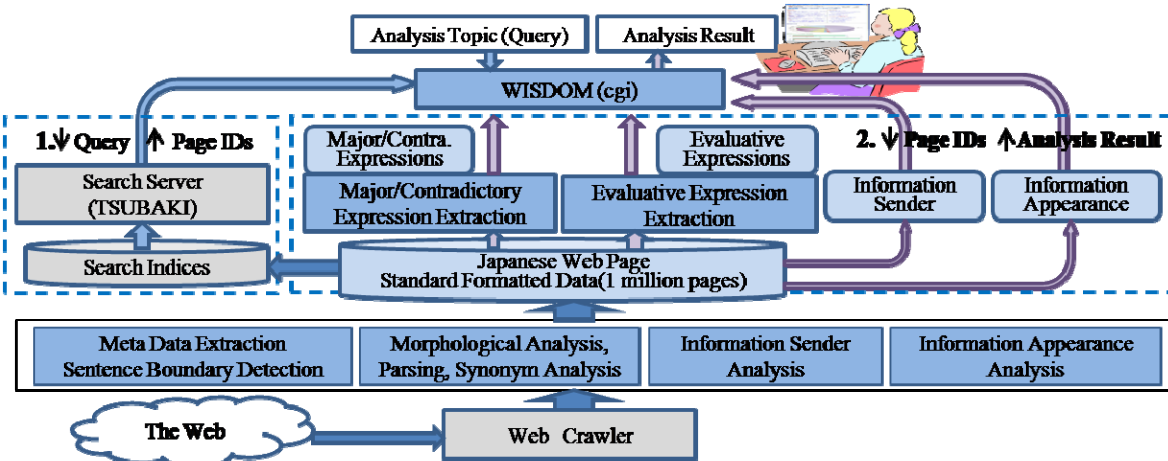


Figure 2. System architecture of WISDOM.

By referring to these phrases and statements, a user can grasp the important issues related to the topic at a glance. The pie diagram indicates the distribution of the information sender class spread over 1000 pages, such as company, industry group, and government. The names of the information senders of the class can be viewed by placing the cursor over a class region. The last bar chart shows the distribution of positive and

negative opinions related to the topic spread over 1000 pages, for all and for each sender class. For example, with regard to “Bio-ethanol,” we can see that the number of positive opinions is more than that of negative opinions, but it is the opposite in the case of some sender classes. Several display units in the Summary tab are cursor sensitive, providing links to more detailed information (e.g., the page list including a major state-

ment, the page list of a sender class, and the page list containing negative opinions).

The “Search Result” tab shows the search result by TSUBAKI, i.e., ranking the relevant pages according to the TSUBAKI criteria. The “Major/Contradictory Expressions” tab shows the list of major phrases and major/contradictory statements about the given topic and the list of pages containing the specified phrase or statement. The “Opinion” tab shows the analysis result of the evaluative expressions, classified according to for/against, like/dislike, merit/demerit, and others, and it also shows the list of pages containing the specified type of evaluative expressions. The “Sender” tab classifies the pages according to the class of the information sender, for example, a user can view the pages created only by the government.

Furthermore, the superficial characteristics of pages called as information appearance are analyzed beforehand and can be viewed in WISDOM, such as whether or not the contact address is shown in the page and the privacy policy is on the page, the volume of advertisements on the page, the number of images, and the number of in/out links.

As shown thus far, given an analysis topic, WISDOM collects and organizes the relevant information available on the Web and provides users with multi-faceted views. We believe that such a system can considerably support the human judgment of information credibility.

3. Data Infrastructure

We usually utilize 100 million Japanese Web pages as the analysis target. The Web pages have been converted into the standard formatted Web data, an XML format. The format includes several metadata such as URLs, crawl dates, titles, and in/out links. A text in a page is automatically segmented into sentences (note that the sentence boundary is not clear in the original HTML file), and the analysis results obtained by a morphological analyzer, parser, and synonym analyzer are also stored in the standard format. Furthermore, the site operator, the page author, and information appearance (e.g., contact address, privacy policy, volume of advertisements, and images) are automatically analyzed and stored in the standard format.

4. Extraction of Major Expressions and Their Contradictions

For the organization of information contents, WISDOM extracts and presents the major expressions and their contradictions on a given analysis topic (Kawahara et al., 2008). Major expressions are defined as expressions occurring at a high frequency in the set of Web pages on the analysis topic. They are classified into two: noun phrases and predicate-argument structures (statements). Contradictions are the predicate-argument structures that contradict the major expressions. For the Japanese phrase *yutori kyōiku*

(cram-free education), for example, *tsumekomi kyōiku* (cramming education) and *ikiru chikara* (life skills) are extracted as the major noun phrases; *yutori kyōiku-wo minaosu* (reexamine cram-free education) and *gakuryokuga teika-suru* (scholastic ability deteriorates), as the major predicate-argument structures; and *gakuryoku-ga koujousuru* (scholastic ability ameliorates), as its contradiction. This kind of summarized information enables a user to grasp the facts and arguments on the analysis topic available on the Web.

We use 1000 Web pages for a topic retrieved from the search engine TSUBAKI. Our method of extracting major expressions and their contradictions consists of the following steps:

1. Extracting candidates of major expressions:

The candidates of major expressions are extracted from each Web page in the search result. From the relevant sentences to the analysis topic that consist of approximately 15 sentences selected from each Web page, compound nouns, parenthetical expressions, and predicate-argument structures are extracted as the candidates of the major expressions.

2. Distilling major expressions:

Simply presenting expressions at a high frequency is not always information of high quality. This is because scattering synonymous expressions such as *karikyuramu* (curriculum) and *kyōiku katei* (course of study) and entailing expressions such as IWC and IWC *soukai* (IWC plenary session), all of which occur frequently, hamper the understanding process of users. Further, synonymous predicate-argument structures such as *gakuryoku-ga teika-suru* (scholastic ability deteriorates) and *gakuryoku-ga sagaru* (scholastic ability lowers) have the same problem. To overcome this problem, we distill major expressions by merging spelling variations with morphological analysis, merging synonymous expressions automatically acquired from an ordinary dictionary and the Web, and merging expressions that can be entailed by another expression.

3. Extracting contradictory expressions:

Predicate-argument structures that negate the predicate of major ones and that replace the predicate of major ones with its antonym are extracted as contradictions. For example, *gakuryoku-ga teika-shi-nai* (scholastic ability does not deteriorate) and *gakuryokuga koujou-suru* (scholastic ability ameliorates) are extracted as the contradictions to *gakuryoku-ga teikasuru* (scholastic ability deteriorates). This process is performed using an antonym lexicon, which consists of approximately 2000 pairs; these pairs are extracted from an ordinary dictionary.

5. Extraction of Evaluative Information

The extraction and classification of evaluative information from texts are important tasks with

many applications and they have been actively studied recently (Pang and Lee, 2008). Most previous studies on opinion extraction or sentiment analysis deal with only subjective and explicit expressions. For example, Japanese sentences such as *watashi-wa apple-ga sukida* (I like apples) and *kono seido-ni hantaida* (I oppose the system) contain evaluative expressions that are directly expressed with subjective expressions. However, sentences such as *kono shokuhin-wa kou-gan-kouka-ga aru* (this food has an anti-cancer effect) and *kono camera-wa katte 3-ka-de kowareta* (this camera was broken 3 days after I bought it) do not contain subjective expressions but contain negative evaluative expressions. From the viewpoint of information credibility, it appears important to deal with a wide variety of evaluative information including such implicit evaluative expressions (Nakagawa et al., 2008).

A corpus annotated with evaluative information was developed for evaluative information analysis studies. Fifty topics such as “Bio-ethanol” and “Pension plan” were chosen. For each topic, 200 sentences containing the topic word were collected from the Web to construct the corpus totaling 10,000 sentences. For each sentence, annotators judged whether or not the sentence contained evaluative expressions. When evaluative expressions were identified, the evaluative expressions, their holders, their sentiment polarities (positive or negative), and their relevance to the topic were annotated.

We developed an automatic analyzer of evaluative information using the corpus. We performed experiments of sentiment polarity classification using Support Vector Machines. Word forms, POS tags, and sentiment polarities from an evaluative word dictionary of all the words in evaluative expressions were used as features, and an accuracy of 83% was obtained. From the error analysis, we found that it was difficult to classify domain-specific evaluative expressions; we are now planning the automatic acquisition of evaluative word dictionaries.

6. Information Sender Analysis

The source of information (or information sender) is one of the important elements when judging the credibility of information. It is rather easy for human beings to identify the information sender of a Web page. When reading a Web page, whether it is deliberate or not, we attribute some characteristics to the information sender and accordingly form our attitudes toward the information. However, the state-of-the-art search engines do not provide facilities to organize a vast amount of information on the basis of the information sender. If we can organize the information on a topic on the basis of who or what type the information sender is, it would enable the user to grasp an overview of the topic or to judge the credibility of relevant information.

WISDOM automatically identifies the *site operators* of Web pages and classifies them into predefined categories of information sender

called *information sender class*. A site operator of a Web page is the governing body of a website on which the page is published. The information sender class categorizes the information sender on the basis of axes such as individuals vs. organizations and profit vs. nonprofit organizations. The list below shows the categories of information sender class.

- | | |
|---|-------------------------------|
| 1. Organization | 1. Organization (cont'd) |
| (a) Profit Organization | (c) Press |
| i. Company | i. Broadcasting Station |
| ii. Industry Group | ii. Newspaper |
| (b) Nonprofit Organization | iii. Publisher |
| i. Academic Society | 2. Individual |
| ii. Government | (a) Real Name |
| iii. Political Organization | (b) Anonymous,
Screen Name |
| iv. Public Service Corp.,
Nonprofit Organization | |
| v. University | |
| vi. Voluntary Association | |
| vii. Education Institution | |

WISDOM allows the user to organize the information on the basis of the information sender class assigned to each Web page. Technical details of the information sender analysis employed in WISDOM can be found in (Kato et al., 2008).

7. Conclusions

This paper has described an information analysis system called WISDOM. As shown in this paper, WISDOM already provides a reasonably nice organized view for a given topic and can serve as a useful tool for handling informational queries and for supporting human judgment of information credibility. WISDOM is freely available at <http://wisdom-nict.jp/>.

References

- B. J. Fogg. 2003. *Persuasive Technology: Using Computers to Change What We Think and Do (The Morgan Kaufmann Series in Interactive Technologies)*. Morgan Kaufmann.
- K. Shinzato, T. Shibata, D. Kawahara, C. Hashimoto, and S. Kurohashi 2008. TSUBAKI: An open search engine infrastructure for developing new information access methodology. In *Proceedings of IJCNLP2008*.
- D. Kawahara, S. Kurohashi, and K. Inui 2008. Grasping major statements and their contradictions toward information credibility analysis of web contents. In *Proceedings of WI'08*.
- B. Pang and L. Lee 2008. Opinion mining and sentiment analysis, *Foundations and Trends in Information Retrieval*, Volume 2, Issue 1-2, 2008.
- T. Nakagawa, T. Kawada, K. Inui, and S. Kurohashi 2008. Extracting subjective and objective evaluative expressions from the web. In *Proceedings of ISUC2008*.
- Y. Kato, D. Kawahara, K. Inui, S. Kurohashi, and T. Shibata 2008. Extracting the author of web pages. In *Proceedings of WICOW2008*.