

Annotating Semantic Relations Combining Facts and Opinions

Koji Murakami[†] Shouko Masuda^{†‡} Suguru Matsuyoshi[†]

Eric Nichols[†] Kentaro Inui[†] Yuji Matsumoto[†]

[†]Nara Institute of Science and Technology

8916-5, Takayama, Ikoma, Nara 630-0192 JAPAN

[‡]Osaka Prefecture University

1-1, Gakuen, Naka-ku, Sakai, Osaka 599-8531 JAPAN

{kmurakami, shouko, matuyosi, eric-n, inui, matsu}@is.naist.jp

Abstract

As part of the STATEMENT MAP project, we are constructing a Japanese corpus annotated with the semantic relations bridging facts and opinions that are necessary for online information credibility evaluation. In this paper, we identify the semantic relations essential to this task and discuss how to efficiently collect valid examples from Web documents by splitting complex sentences into fundamental units of meaning called “statements” and annotating relations at the statement level. We present a statement annotation scheme and examine its reliability by annotating around 1,500 pairs of statements. We are preparing the corpus for release this winter.

1 Introduction

The goal of the STATEMENT MAP project (Murakami et al., 2009) is to assist internet users with evaluating the credibility of online information by presenting them with a comprehensive survey of opinions on a topic and showing how they relate to each other. However, because real text on the Web is often complex in nature, we target a simpler and more fundamental unit of meaning which we call the “statement.” To summarize opinions for the statement map users, we first convert all sentences into statements and then, organize them into groups of agreeing and conflicting opinions that show the logical support for each group.

For example, a user who is concerned about potential connections between vaccines and autism would be presented with a visualization of the opinions for and against such a connection together with the evidence supporting each view as

shown in Figure 1.

When the concerned user in our example looks at this STATEMENT MAP, he or she will see that some opinions support the query “Do vaccines cause autism?” while other opinions do not, but it will also show what support there is for each of these viewpoints. So, STATEMENT MAP can help user come to an informed conclusion.

2 Semantic Relations between Statements

2.1 Recognizing Semantic Relations

To generate STATEMENT MAPs, we need to analyze a lot of online information retrieved on a given topic, and STATEMENT MAP shows users a summary with three major semantic relations.

AGREEMENT to group similar opinions

CONFLICT to capture differences of opinions

EVIDENCE to show support for opinions

Identifying logical relations between texts is the focus of Recognizing Textual Entailment (RTE). A major task of the RTE Challenge (Dagan et al., 2005) is the identification of [ENTAILMENT] or [CONTRADICTION] between Text (T) and Hypothesis (H). For this task, several corpora have been constructed over the past few years, and annotated with thousands of (T,H) pairs.

While our research objective is to recognize semantic relations as well, our target domain is text from Web documents. The definition of contradiction in RTE is that T contradicts H if it is very unlikely that both T and H can be true at the same time. However, in real documents on the Web, there are many examples which are partially contradictory, or where one statement restricts the applicability of another like in the example below.

(1) a. Mercury-based vaccines actually cause autism in children.

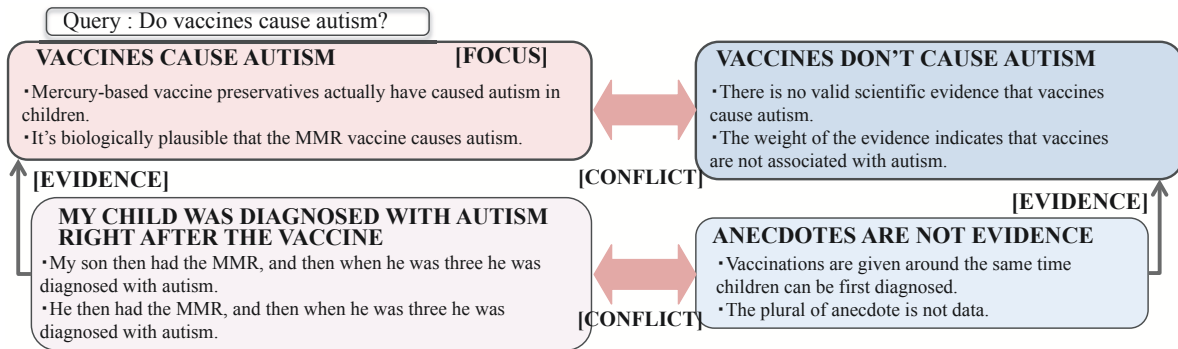


Figure 1: An example STATEMENT MAP for the query “Do vaccines cause autism?”

- b. Vaccines can trigger autism in a vulnerable subset of children.

While it is difficult to assign any relation to this pair in an RTE framework, in order to construct statement maps we need to recognize a contradiction between (1a) and (1b).

There is another task of recognizing relations between sentences, CST (Cross-Document Structure Theory) which was developed by Radev (2000). CST is an expanded rhetorical structure analysis based on RST (Mann and Thompson, 1988), and attempts to describe relations between two or more sentences from both single and multiple document sets. The CSTBank corpus (Radev et al., 2003) was constructed to annotate cross-document relations. CSTBank is divided into clusters in which topically-related articles are gathered. There are 18 kinds of relations in this corpus, including [EQUIVALENCE], [ELABORATION], and [REFINEMENT].

2.2 Facts and Opinions

RTE is used to recognize logical and factual relations between sentences in a pair, and CST is used for objective expressions because newspaper articles related to the same topic are used as data. However, the task specifications of both RTE and CST do not cover semantic relations between opinions and facts as illustrated in the following example.

- (2) a. There must not be a connection between vaccines and autism.
 b. I do believe that there is a link between vaccinations and autism.

Subjective statements, such as opinions, are recently the focus of many NLP research topics, such as review analysis, opinion extraction, opinion QA, or sentiment analysis. In the corpus constructed by the MPQA Project (Multi-Perspective Question Answering) (Wiebe et al., 2005), individual expressions are marked that correspond to

explicit mentions of private states, speech events, and expressive subjective elements.

Our goal is to annotate instances of the three major relation classes: [AGREEMENT], [CONFLICT] and [EVIDENCE], between pairs of statements in example texts. However, each relation has a wide range, and it is very difficult to define a comprehensive annotation scheme. For example, different kinds of information can act as clues to recognize the [AGREEMENT] relations. So, we have prepared a wide spectrum of semantic relations depending on different types of information regarded as clues to identify a relation class, such as [AGREEMENT] or [CONFLICT]. Table 1 shows the semantic relations needed for carrying out the annotation. Although detecting [EVIDENCE] relations is also essential to the STATEMENT MAP project, we do not include them in our current corpus construction.

3 Constructing a Japanese Corpus

3.1 Targeting Semantic Relations Between Statements

Real data on the Web generally has complex sentence structures. That makes it difficult to recognize semantic relations between full sentences, but it is possible to annotate semantic relation between parts extracted from each sentence in many cases. For example, the two sentences A and B in Figure 2 cannot be annotated with any of the semantic relations in Table 1, because each sentence include different types of information. However, if two parts extracted from these sentences C and D are compared, the parts can be identified as [EQUIVALENCE] because they are semantically close and each extracted part does not contain a different type of information. So, we attempt to break sentences from the Web down into reasonable text segments, which we call “statements.” When a real sentence includes several pieces of se-

Table 1: Definition of semantic relations and example in the corpus

Relation Class	Relation Label	Example
AGREEMENT	Equivalence	A: The overwhelming evidence is that vaccines are unrelated to autism. B: There is no link between the MMR vaccine and autism.
	Equivalent Opinion	A: We think vaccines cause autism. B: I am the mother of a 6 year old that regressed into autism because of his 18 month vaccinations.
	Specific	A: Mercury-based vaccine preservatives actually have caused autism in children. B: Vaccines cause autism.
CONFLICT	Contradiction	A: Mercury-based vaccine preservatives actually have caused autism in children. B: Vaccines don't cause autism.
	Confinement	A: Vaccines can trigger autism in a vulnerable subset of children. B: Mercury-based vaccine actually have caused autism in children.
	Conflicting Opinion	A: I don't think vaccines cause autism. B: I believe vaccines are the cause of my son's autism.

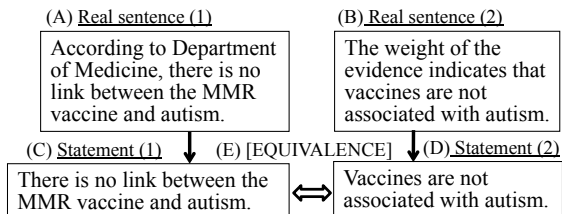


Figure 2: Extracting statements from sentences and annotating a semantic relation between them

semantic segments, more than one statement can be extracted. So, a statement can reflect the writer’s affirmation in the original sentence. If the extracted statements lack semantic information, such as pronouns or other arguments, human annotators manually add the missing information. Finally we label pairs of statements with either one of the semantic relations from Table 1 or with “NO RELATION,” which means that two sentences (1) are not semantically related, or (2) have a relation other than relations defined in Table 1.

3.2 Corpus Construction Procedure

We automatically gather sentences on related topics by following the procedure below:

1. Retrieve documents related to a set number of topics using a search engine
2. Extract real sentences that include major subtopic words which are detected based on TF or DF in the document set
3. Reduce noise in data by using heuristics to eliminate advertisements and comment spam
4. Reduce the search space for identifying sentence pairs and prepare pairs, which look feasible to annotate.

Dolan and Brockett (2005) proposed a method to narrow the range of sentence pair candidates and collect candidates of sentence-level paraphrases which correspond [EQUIVALENCE] in [AGREEMENT] class in our task. It worked well

for collecting valid sentence pairs from a large cluster which was constituted by topic-related sentences. The method also seem to work well for [CONFLICT] relations, because lexical similarity based on bag-of-words (BOW) can narrow the range of candidates with this relation as well.

We calculate the lexical similarity between the two sentences based on BOW. We also used hyponym and synonym dictionaries (Sumida et al., 2008) and a database of relations between predicate argument structures (Matsuyoshi et al., 2008) as resources. According to our preliminary experiments, unigrams of KANJI and KATAKANA expressions, single and compound nouns, verbs and adjectives worked well as features, and we calculate the similarity using cosine distance. We did not use HIRAGANA expressions because they are also used in function words.

4 Analyzing the Corpus

Five annotators annotated semantic relations according to our specifications in 22 document sets as targets. We have annotated target statement pairs with either [AGREEMENT], [CONFLICT] or [NO RELATION]. We provided 2,303 real sentence pairs to human annotators, and they identified 1,375 pairs as being invalid and 928 pairs as being valid. The number of annotated statement pairs are 1,505 ([AGREEMENT]:862, [CONFLICT]:126, [NO RELATION]:517).

Next, to evaluate inter annotator agreement, 207 randomly selected statement pairs were annotated by two human annotators. The annotators agreed in their judgment for 81.6% of the examples, which corresponds to a kappa level of 0.49. The annotation results are evaluated by calculating recall and precision in which one annotation result is treated as a gold standard and the other’s as the output of the system, as shown in Talbe 2.

Table 2: Inter-annotator agreement for 2 annotators

		Annotator A			TOTAL
		AGR.	CON.	NONE	
Anno- tator B	AGR.	146	7	9	162
	CON.	0	13	1	14
	NONE	17	4	10	31
	TOTAL	163	24	20	207

5 Discussion

The number of sentence pairs that annotators identified as invalid examples shows that around 60% of all pairs were invalid, showing that there is still room to improve our method of collecting sentence pairs for the annotators. Developing more effective methods of eliminating sentences pairs that are unlikely to contain statements with plausible relations is important to improve annotator efficiency. We reviewed 50 such invalid sentence pairs, and the results indicate two major considerations: (1) negation, or antonyms have not been regarded as key information, and (2) verbs in KANJI have to be handled more carefully. The polarities of sentences in all pairs were the same although there are sentences which can be paired up with opposite polarities. So, we will consider the polarity of words and sentences as well as similarity when considering candidate sentence pairs.

In Japanese, the words which consist of KATAKANA expressions are generally nouns, but those which contain KANJI can be nouns, verbs, or adjectives. Sharing KATAKANA words was the most common way of increasing the similarity between sentences. We need to assign a higher weight to verbs and adjectives that contain KANJI, to more accurately calculate the similarity between sentences.

Another approach to reducing the search space for statement pairs is taken by Nichols et al. (2009), who use category tags and in-article hyperlinks to organize scientific blog posts into discussions on the same topic, making it easier to identify relevant statements. We are investigating the applicability of these methods to the construction of our Japanese corpus but suffer from the lack of a richly-interlinked data source comparable to English scientific blogs.

6 Conclusion

In this paper, we described the ongoing construction of a Japanese corpus consisting of statement pairs annotated with semantic relations for handling web arguments. We designed an annotation

scheme complete with the necessary semantic relations to support the development of statement maps that show [AGREEMENT], [CONFLICT], and [EVIDENCE] between statements for assisting users in analyzing credibility of information in Web. We discussed the revelations made from annotating our corpus, and discussed future directions for refining our specifications of the corpus. We are planning to annotate relations for more than 6,000 sentence pairs in this summer, and the finished corpus will consist of around 10,000 sentence pairs. The first release of our annotation specifications and the corpus will be made available on the Web¹ this winter.

Acknowledgments

This work is supported by the National Institute of Information and Communications Technology Japan.

References

- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In *Proc. of the PASCAL Challenges Workshop on RTE*.
- Bill Dolan and Chris Brockett. 2005. Automaticaly constructing a corpus of sentential paraphrases. In *Proc. of the IWP 2005*, pages 9–16.
- William Mann and Sandra Thompson. 1988. Rhetorical structure theory: towards a functional theory of text organization. *Text*, 8(3):243–281.
- Suguru Matsuyoshi, Koji Murakami, Yuji Matsumoto, , and Kentaro Inui. 2008. A database of relations between predicate argument structures for recognizing textual entailment and contradiction. In *Proc. of the ISUC 2008*.
- Koji Murakami, Eric Nichols, Suguru Matsuyoshi, Asuka Sumida, Shouko Masuda, Kentaro Inui, and Yuji Matsumoto. 2009. Statement map: Assisting information credibility analysis by visualizing arguments. In *Proc. of the WICOW 2009*, pages 43–50.
- Eric Nichols, Koji Murakami, Kentaro Inui, and Yuji Matsumoto. 2009. Constructing a scientific blog corpus for information credibility analysis. In *Proc. of the Annual Meeting of ANLP*.
- Dragomir Radev, Jahna Otterbacher, and Zhu Zhang. 2003. CSTBank: Cross-document Structure Theory Bank. <http://tangra.si.umich.edu/clair/CSTBank>.
- Dragomir R. Radev. 2000. Common theory of information fusion from multiple text sources step one: Cross-document structure. In *Proc. of the 1st SIGdial workshop on Discourse and dialogue*, pages 74–83.
- Asuka Sumida, Naoki Yoshinaga, and Kentaro Torisawa. 2008. Boosting precision and recall of hyponymy relation acquisition from hierarchical layouts in wikipedia. In *Proc. of the LREC 2008*.
- Janyce Wiebe, Theresa Wilson, and Claire Cardie. 2005. Annotating expressions of opinions and emotions in language. *Language Resources and Evaluation*, 39(2-3):165–210.

¹<http://cl.naist.jp/stmap/corpus/ja>