

情報科学論文のための意味関係検索システム

建石 由佳 宮尾 祐介 相澤 彰子

国立情報学研究所

{yucca,yusuke,aizawa}@nii.ac.jp

1 はじめに

学術論文のオンライン出版化が広まるに伴って出版量が増え、読む側では自動検索が必須のツールとなり、知的な検索が望まれるようになってきている。例えば「CRF (conditional random field)」を検索するときも、CRF とは何か、CRF は何に使われるか、あるいはCRF を改善するにはどうしたらよいか、など、さまざまな検索要求があり、それに沿った結果に容易に辿りつけることが望ましい。しかし、現在の検索システムの多くはキーワード(フレーズを含む)の集合を検索キーとして、実際の意図とは関係なくそのキーワードあるいはその異形を含む文献をすべて結果とするため、結果が要求に沿うものであるかどうかは検索者が判断しなければならない。

ここで、検索結果を表示する際に何らかの意味的なグループ分けができれば判断の助けになるであろう。そのためには、キーワードと文中の他の語の意味関係を知ることが重要になると考えられる。例えば、「CRF を用いた形態素解析」ではCRF を利用したこと、「CRF における素性削減」ではCRF そのものを改良することに関する記述であることが読み取れる。ただし、「用いる」「おける」などの語との共起のみを手掛かりにするのでは、先の例と「Confusion Network を用いたCRF」のようにCRF そのものの仕組みに関する記述を区別できないので不十分である。

そこで、我々は文中の語と語を関係づけることによって検索対象のキーワードが文中でどのような意味的役割を果たしているかを検出し、それにより検索結果を文単位で分類することを試みる。この目的のために我々は情報科学論文に意味に基づく関係をアノテートしたコーパスを作成し、そのコーパスに基づく教師あり学習により論文中の文を分類するシステムを試作したので報告する。

2 関連研究

従来、関係をとらえるために、構文解析が用いられてきた。例えば、ACL Searchbench[4]ではHPSGパーサーが論文のトピックの検出と、Predicate-Argument関係に基づく結果の絞り込みに利用される。Guptaら

[2]は、論文のFocus(トピック)、Domain(Focusが適用される対象)、Technique(Focusを達成するための方法)を抽出するためにDependency Parsingにもとづく特定の動詞との依存関係を利用している。しかし、意味的な依存関係が直接構文上の依存関係ですべてとらえられるとは限らないため、我々は直接意味関係を抽出する方法をとる。日本では、[1]など、論文や特許文書の中のTechnology(技術)、Effect(技術を使用した結果を表すAttributeとValueのペア)に対応するエンティティを同定する研究があるが、エンティティの同定にとどまり、相互の関係付けは行われていない。

意味的關係を直接抽出する研究は、生命科学分野で盛んに行われており、タンパク質相互作用等のイベント、細胞内局在関係などの関係など様々な情報をアノテートしたコーパスが作られ、情報検索システムが作成されている(例えば[3]など)。生命科学分野の大きな特徴は、オントロジー(Gene Ontology¹など)、標準語彙集(UMLS²など)などの言語資源の整備が進んでいることである。そのため、多くの研究では特定の現象について、それにかかわる物質や作用などをオントロジー等へのマッピングを経由して抽出することが行われている。

我々の対象である情報科学分野はそのようなオントロジー等が存在せず、また、日々新しい技術が生み出されていくという情報科学(あるいは工学一般)の特質上、標準オントロジーの作成は困難と考えられる。したがって、文中に出てくるものごととものごとの関係を構造化するのに必要なアノテーションスキーマを実際のアノテーションを行いながら定義していくという方針を取った。

3 関係アノテーション付きコーパス

我々は情報処理学会論文誌のアブストラクトをベースとして、論文中の語と語の関係をほぼもれなく構造化したコーパスを作成した[6, 8]。検索システムの作成を最終目的とするため、抽出対象としての特定のトピックを仮定せず、論文に書かれた情報を網羅的に構

¹<http://www.geneontology.org/>

²<http://www.nlm.nih.gov/research/umls/>



図 1: 揺れの例

ノテータにはあらかじめ、同じスキーマでアノテートした別の 100 アブストラクトとマニュアルを渡したが、他の教示は行わなかった。また、自動前処理もおこなわなかった。F-Score による一致率はエンティティで 96.4%、関係で 79.1%であった。小規模な実験結果ではあるが、この結果は、現在のスキーマが、第三者にも安定したアノテーションを可能にすることを示唆する。また、アノテータからは、判定に迷うケースで実際のアノテート例を見ることが非常に役に立ったとのフィードバックを得ており、高い一致率と合わせて、アノテーションガイドラインに実例が重要であるということを示している。

30 件中最初の 10 件についてアノテータの結果とゴールドスタンダードを比較し、エラーの傾向を調べた。エンティティの不一致は 20 か所あったがそのうちの 16 か所は「問題がある」、「課題である」などにおける「問題」「課題」などを MEASURE 扱いするかどうかに関するものであった。これらの表現はスキーマの設計時にも揺れが多くみられ、ガイドラインにいくつかの例を追加したが、まだ基準が十分クリアでないようである。関係の不一致 (64 か所) についてはさまざまな物があったが、目立ったパターンとして、「文の場合は、単語数が多いとこのパターンに従う」の「文」と「単語数」の関係 (この場合は ATTRIBUTE 関係) のように、関係づけられるもの同士の関係が読点を挟んだり遠い場合の不一致が見られた (図 1)。この例は、英訳すると“A sentence follows this pattern when the number of words in it is large”のように代名詞が現れるため、その共参照関係 (我々のスキーマでは EQUIVALENCE 関係) を介して“sentence”と“number of words”を関係づけることができるが、和文では代名詞が陽に表れないために関係があることが意識されないのだという、日本語に特有の問題であると考えられる。

実験に使った 30 件と、例示に使った 100 件の計 130 件はリクエストベースで公開している。

4 論文検索システム

作成したコーパスを使って、SVM による関係抽出器³を制作し、その結果に基づいて論文中の文を分類するシステムを試作した。システムはウェブブラウザ上で動作⁴し、情報処理学会論文誌アブストラクト約

³詳細に関する報告は準備中

⁴現在 URL は公開していない

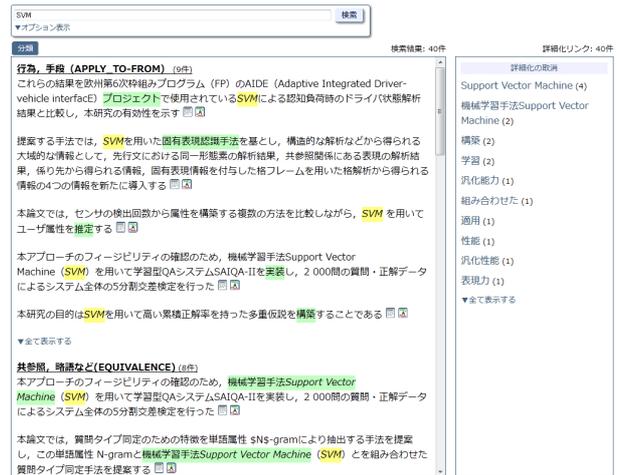


図 2: 検索システム：キーワードの役割による分類

3500 件からの検索が行える。図 2 と図 3 は現在のシステムのインターフェースを示している。現在、システムを使用しながら、機能、インターフェースの両面での問題点を洗い出している。

図 2 では、左上に示されているウィンドウに検索キーワードを入力すると、右側のウィンドウに検索キーワードと関連付けられた語が表示され、同時に左側のウィンドウにキーワードを含む文が表示される。キーワード (SVM) は図上では黄色であらわされている。結果はまとめて表示することも可能であるが、キーワード入力ウィンドウの下のボタンにより、関係により分類することが可能である。図 2 では左側のウィンドウに分類された結果が表示されており、関係づけられた相手の語が緑で表示されている。また、図 2 の画面で右側のウィンドウの特定の語をクリックすることにより、その語と関係づけられている場合のみを表示するような絞り込みができる (図 3)。検索キーワードには、正規表現を含めることができ、あいまい検索も可能にしている。また、検索結果の文ごとに、原文 (情報処理学会論文誌アブストラクト) と自動アノテーション結果 (図 4) へのリンクを持つ。

図 3 では、「SVM」 (左ウィンドウ黄色) と「構築」 (左ウィンドウ緑) が関係づけられている文を、関係ごとに表示している。最初の文 (「本研究の目的は SVM を用いて高い累積正解率を持った多重仮説を構築することである」) では SVM を利用して他のものを構築することが「SVM」と「構築」の間の APPLY_TO 関係により示されている。2 番目の文 (「... パラメトリックな確率モデルおよび SVM を構築する」) では SVM 自体を構築することが「SVM」が「構築」 (という行為) の OUTPUT であるという関係で示されている。現在の実装は、コーパス上の関係ラベルをそのまま使って関係とその方向性をもとに分類するものであるが、これが検索者にとって本当に使いやすい分類であるのか、また、使いにくい場合、どのような提示をするのがよいのか、ということが現在検討中の課題で



図 3: 検索システム：他の語との関係による分類

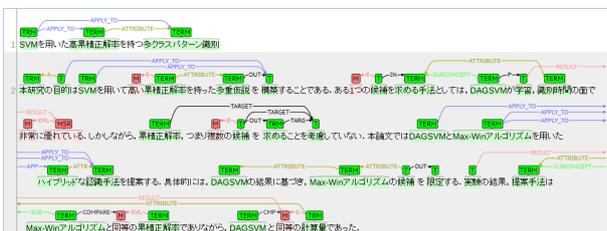


図 4: 図 3 の最初の文を含むアブストラクトのアノテーション (brat[5] による表示)

ある。

5 おわりに

情報科学論文中に意味関係をアノテートしたコーパスを作成し、それに基づく機械学習により、キーワードで示される語句を含む文を、キーワードと他の語との関係に基づいて分類し表示するシステムを試作した。

我々の枠組みは文中に示されるものと他のものや行為との関係を示すことにより、そのものが論文に示される研究の中で果たしている役割にもとづいた文の分類・検索を可能にしているといえる。ただし、我々の枠組みでは、文に記述されたことが新たに分かったことなのか、問題として指摘されていることなのか、すでに分かっていたことなのか、などの区別をしていない。これらは Zoning と呼ばれる枠組みで研究がおこなわれており ([7] など)、今後これらと組み合わせるなどして、より柔軟な検索を可能にすることを目指す。

参考文献

[1] Satoshi Fukuda, Hidetsugu Nanba, and Toshiyuki Takezawa. Extraction and visualization of technical trend information from

research papers and patents. In *Proceedings of the 1st International Workshop on Mining Scientific Publications*, 2012.

- [2] Sonal Gupta and Christopher D Manning. Analyzing the dynamics of research by extracting key aspects of scientific papers. In *Proceedings of 5th IJCNLP*, 2011.
- [3] Claire Nédellec, Robert Bossy, Jin-Dong Kim, Jung-Jae Kim, Tomoko Ohta, Sampo Pyysalo, and Pierre Zweigenbaum. Overview of bionlp shared task 2013. In *Proceedings of the BioNLP Shared Task 2013 Workshop*, pp. 1–7, Sofia, Bulgaria, August 2013. Association for Computational Linguistics.
- [4] Ulrich Schäfer, Bernd Kiefer, Christian Spurk, Jörg Steffen, and Rui Wang. The ACL anthology searchbench. In *Proceedings of the ACL-HLT 2011 System Demonstrations*, pp. 7–13, 2011.
- [5] Pontus Stenetorp, Sampo Pyysalo, Goran Topić, Tomoko Ohta, Sophia Ananiadou, and Jun’ichi Tsujii. brat: a web-based tool for NLP-assisted text annotation. In *Proceedings of the Demonstrations Session at EACL*, 2012.
- [6] Yuka Tateisi, Yo Shidahara, Yusuke Miyao, and Akiko Aizawa. Relation annotation for understanding research papers. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pp. 140–148, Sofia, Bulgaria, August 2013. Association for Computational Linguistics.
- [7] Simone Teufel, Advait Siddharthan, and Colin Batchelor. Towards discipline-independent argumentative zoning: evidence from chemistry and computational linguistics. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pp. 1493–1502, 2009.
- [8] 建石由佳, 仕田原容, 宮尾祐介, 相澤彰子. 情報科学論文からの意味関係抽出に向けたタグ付けスキーマ. 言語処理学会第 19 回年次大会講演論文集, pp. 702–705, 2013.