

## Wikipedia マイニングと Web オントロジ構築の新しい方向性

中山浩太郎† 原隆浩‡ 西尾章治郎‡

東京大学知の構造化センター† 大阪大学大学院情報科学研究科‡

[nakayama@cks.u-tokyo.ac.jp](mailto:nakayama@cks.u-tokyo.ac.jp) [hara@ist.osaka-u.ac.jp](mailto:hara@ist.osaka-u.ac.jp) [nishio@ist.osaka-u.ac.jp](mailto:nishio@ist.osaka-u.ac.jp)

### Wikipedia Mining and a New Paradigm on Web Ontology Construction

Kotaro NAKAYAMA† Takahiro HARA‡ Shojiro NISHIO‡

The Center for Knowledge Structuring, the University of Tokyo†

Dept. of Multimedia Eng., Graduate School of Information Science and Technology, Osaka University‡

#### 1. はじめに

Wikipedia は、Wiki をベースにした大規模 Web 百科事典であり、誰でも Web ブラウザを通じて記事内容を変更できることから、一般的な概念だけでなく、文化、歴史、数学、科学、社会、テクノロジーなどの幅広い分野をカバーし、普遍的な概念から新しい概念に至るまで、非常に膨大なコンテンツが網羅されている。その記事数は既に 200 万 (2007 年 12 月英語のみカウント) を超えており、世界最大の百科事典である Britannica の記事数が、全 60 巻で約 65,000 記事であることと比較した場合、実に 30 倍近い数の記事が網羅されていることになる。

Wikipedia は、この幅広いトピックの網羅性以外にも興味深い特徴をいくつか持つ。密なリンク構造や、URL により語彙の意味が一意に特定されている点、質の高いリンクテキスト、精査されたカテゴリ構造や情報ボックス (Infobox) なども Wikipedia の知識抽出のコーパスとして見たときの特徴である。筆者らは Wikipedia の特徴を利用して、概念同士の関係度を数値化した連想シソーラスの構築に関する研究を進め、その実装である「Wikipedia シソーラス」(<http://wikipedia-lab.org> 参照) を公開した。連想シソーラスは、概念間の関係度を数値化した辞書であり、情報検索や自然言語処理、情報フィルタリングなど幅広い分野で必要とされている。しかし、連想シソーラス辞書は、単に概念間の関係度を数値化しているだけであり、より明確な意味関係を抽出することで大規模 Web オントロジを構築することが本研究の目標である。Wikipedia からの Web オントロジ構築はセマンティック Web 研究や概念辞書構築の新しいパラダイムとなる可能性が高く、国内外で注目を集めている研究分野である。

#### 2. 技術的課題

前述のとおり、Wikipedia は知識抽出のコーパスとして有用であるが、①情報の信頼性、②スケーラビリティ、③情報統合などの技術的課題も残されている。特に、情報の信頼性に関する問題は、Wikipedia のように多数のユーザがコンテンツを編集するソーシャルメディアでは重要な問題であり、質の高い情報だけを抽出できるかがポイントになる。また、Web オントロジ構築は、単体だけでは意味がなく、アプリケーションやタスクを考慮して手法を開発・最適化する必要がある。たとえば、情報検索の分野で

は比較的精度よりも網羅性が必要とされる場合が多いが、機械翻訳のための要素技術として利用する場合は、精度が極めて重要になる。

#### 3. 研究の方向性

本研究では Wikipedia マイニングによって、概念間の明確な意味関係 (is-a 関係など) を定義した Web オントロジを自動構築する手法の構築を目指す。また、その際に上述の問題点を解決する基盤技術を開発していくことが研究のポイントとなる。

①情報の信頼性: Wikipedia の記事内容の信頼性を数値化するためのパラメータとして利用できる指標はいくつかある。バックワードリンク数や記事の更新状況、生存期間などがその最たる例である。他にどのような指標が Wikipedia の記事の信頼性に影響を与えるのかを調査・利用することで質の高い情報抽出を目指す。

②スケーラビリティ: Wikipedia のテキストを解析するには、構文解析が必要となるが、構文解析は多量の計算機リソースを消費するため、日に日に増加する Wikipedia 記事をすべて解析するためには計算量を抑える仕組みが必要となる。そのため、1) インクリメンタルな解析手法、2) 重要文だけの解析によるパフォーマンス向上、3) 非構文解析手法の検討などの研究を行う必要がある。

③情報統合: Wikipedia は専門的な概念や新しい単語に対する網羅性が高い。その一方で、概念辞書である WordNet や OpenCYC などは一般的な概念がよく網羅されているため、融合することで両者の利点を活かした概念辞書が構築できると考えられる。また、Web マイニング手法との融合手法により、さらに網羅性を高めるつもりである。

④アプリケーションへの応用: 本研究では、開発した大規模シソーラスおよび大規模オントロジを公開している。これらのリソースをほかの研究者に公開することで、アプリケーションへの応用の可能性を促進し、より有用な Web オントロジの構築を目指す。

#### 参考文献

- 1) 中山浩太郎, 原隆浩, 西尾章治郎: 人工知能研究の新しいフロンティア: Wikipedia, 人工知能学会誌, Vol. 22, No. 5 (2007 年 9 月).