

氏名 中渡瀬秀一

学位（専攻分野） 博士（情報学）

学位記番号 総研大甲第 1157 号

学位授与の日付 平成 20 年 3 月 19 日

学位授与の要件 複合科学研究科 情報学専攻  
学位規則第 6 条第 1 項該当

学位論文題目 テキストコーパスを用いた語の相互関係の発見に関する研究

論文審査委員 主査教授 相澤 彰子  
教授 高須 淳宏  
教授 武田 英明  
助教 井上 雅史  
名誉教授 上野 晴樹（国立情報学研究所）  
准教授 影浦 峠（東京大学）

## 論文内容の要旨

本論文は、日本語コーパスの文を解析して語の相互関係を発見する手法の提案とその有効性の検証について述べる。近年、ネットワークを通じて多くの電子化文書が蓄積され、それらを検索、分類、加工する利用が進んでいる。その際には、これらの文書は自然言語処理によって解析が行われる。そしてこの解析を行う時には語彙知識の電子化辞書が必要とされる。語の相互関係はこの語彙知識の一種であり、電子化辞書の重要な構成要素である。本論文では、そのような語の相互関係として類義関係、多義語、上位下位関係に焦点をあて、それらを日本語コーパスの解析に基づいて発見する手法の提案とその検証を行った。

コーパス中から語の相互関係を発見するときに、従来以下の2つの課題が存在していた。

第一は、ある2語の間に類義関係が存在するか否かは観点に依存しているため、どの観点による類義関係なのかをいかに明確にするかというものである。従来の類義語獲得の研究では、主に語の間の類似度が一意であることを前提にしたクラスタリングによる手法が用いられてきた。しかしこれらの研究は、2語の間には常に一意な類似度が与えられるため類似性の観点による違いを扱うことが困難であった。そして類義語集合の作成はこの類似度に基づくため、ある語A、B、Cに関してAB間の類似度とAC間の類似度の観点が同じでないときにも同一尺度として、それを基準としてクラスタリングするために不適切な類義語集合が得られるという問題があった。

第二は、語の間の上位下位関係を発見する手法における計算手順の導出をどの上位下位関係の定義から合理的に行うかというものである。この点で従来の上位下位関係獲得の研究では、コーパス中に存在する「AなどのB」といった定型表現からBとAの上位下位関係を抽出する手法のような表層的な手法が提案されていた。この場合、抽出対象となる定型表現にはその著者の意識した上位下位関係しか表現されないので、著者が意識しないかもしくは、自明であるが電子化辞書に必要な上位下位関係が得られないという問題があった。それに対して、コーパス中の語の共起関係に基づく抽出手法も提案されている。例えば語A、Bが類義関係にあり、Aの出現頻度がBより高ければAはBの上位語とする手法などである。しかしそれらは上位下位関係の定義と無関係に構成されていたために、理論的合理性を欠いていた。

第一の課題に対しては、同じ観点で類義関係にある語だけを直接グルーピングして類義語集合を得るというアプローチを検討した。そして新たな手法として語の共起関係から得られるグラフ構造から類義語集合なすグラフ構造を探索し列挙する方法を提案し、検証実験によって、提案手法が同じ観点の類義関係をグルーピングすることにおいて有効であることを示した。

第二の課題に対しては、まず名詞が指示する対象の範囲の包含関係によって語のIS-A関係を定めることを示した。次にある動詞に対して格関係をもつ対象

も名詞として表現されることに着目し、動詞と名詞の依存関係を用いてIS-A関係を発見する方法を検討した。そして動詞と名詞の依存性解析を行って得られる名詞ごとの動詞集合を比較する手法の提案を行った。さらに、この手法の有効性を検証するための実験を行い、本手法によって獲得されたIS-A関係が従来手法によるカバレッジを大幅に拡大することが確認した。

本論文は6章から構成される。

第1章では、本研究における背景と目的について説明し、課題とそれに対するアプローチを述べる。まず自然言語処理で用いられる電子化辞書や従来の伝統的な類義語辞書などの適用分野やその辞書構造について概観する。その上で本研究の取り組む課題として、コーパスの解析に基づく語の相互関係の発見を取り上げる。そして、最後に本論文の全体構成について述べる。

第2章では、本研究に関連する研究として、従来の類義語集合や階層関係の抽出手法の概観を行い、本研究の位置付けを明確にする。

第3章では、同じ観点で類義関係にある語をグルーピングする方法について議論する。

ここでは、コーパス中に含まれる複合名詞から修飾関係を抽出し、それらによる語と修飾関係を頂点と辺とするグラフ構造中から極大完全2部グラフ部分を類義語集合として探索し抽出する手法を提案した。評価実験では1ヶ月分の新聞記事コーパスから得られる修飾関係をグラフ化し、その中に含まれる極大完全2部グラフ（約4900個）を抽出した。得られた類義語集合の正解判定は人手で行い、その結果、2頂点同士からなる2部グラフにおける正解率は約30%であった。また観点の違いによる類義語集合の獲得に関しては、同じ語に対して違う観点で複数の類義語集合が得られることを確認した。これにより本研究の第一の課題に対し提案手法が有効であることを確認した。

第4章では、多義性を持つ語の発見方法について議論する。語の相互関係を考える場合、同じ語でもその語義ごとに分けて扱わなければならない。そこで本研究ではこの課題を解決するために、多義性を持つ語の発見方法を提案する。ここでは、多義語となる語はその語義によって異なる類義語集合に含まれることに注目する。そこで、本手法ではまず3章の手法で類義語集合を抽出し、次に多義性を調べたい語について、その語を含む類義語集合の数を調べて複数の類義語集合に属するなら多義性を持つと判断する。ところがコーパス中の修飾関係が少ない場合、同じ観点による類義語集合が複数得られるという問題がある。そこで、本手法は構成要素が類似した類義語集合は併合する手順を加えてこれに対処している。評価実験では1年分の新聞記事コーパスから得られる修飾関係をグラフ化し、極大完全2部グラフ（約1,300,000個）を抽出した。この実験では抽出グラフ数が多いため、高速なグラフ探索列挙アルゴリズムとして逆探索法を用いた。次に得られたグラフに併合操作を行い、形態素解析誤りなどを含むグラフを除去した後、複数の類義語集合に含まれる語を集計した（182語）。そしてこれらの語の多義性を人手で判定した。その結果、そのうち31語については多義語であることが確認された。これにより、多義語の発見手法の有

効性を確認した。

第5章では、IS-A関係の発見方法について議論する。ここでは動詞と名詞の依存関係を用いた発見方法を提案する。まずIS-A関係の定義を明確にし、これよりIS-A関係の発見手法を導いた。具体的には動詞と名詞の依存関係を解析し、名詞ごとの動詞集合を作成し、それらの包含を比較して、与えられた名詞の下位語候補となる順序リストを作成する手法を提案した。評価実験ではまず11年分の新聞記事コーパスを解析し、名詞ごとの動詞集合を作成した。次に評価に用いる上位語のサンプルリストを分類語彙表から作成した。そしてこのリストに含まれる名詞に対して提案手法を用いて下位語候補を作成した。この候補語リストの正解判定は人手で行なった。その結果、平均正解率は約34%であった。また従来手法とのカバレッジの違いを確かめるために、同じコーパスから定型表現を用いてIS-A関係を抽出し、本手法によって得られたIS-A関係と比較した。その結果、本手法で得られるIS-A関係の約6%が従来手法によって獲得されることを確認した。これにより、本手法の有効性が示された。

第6章では、本研究での結論をまとめる。さらに将来の展望として、類義語や上位下位関係の獲得手法の拡張可能性やこれらの手法によって得られる語彙知識の応用システムへの適用分野について述べる。

## 論文の審査結果の要旨

中渡瀬秀一さんの博士論文は、「テキストコーパスを用いた語の相互関係の発見に関する研究」と題され、日本語テキストコーパスを解析して語の相互関係を抽出するための手法および有効性の検証について述べたものである。

近年、電子化された文書情報が大量に流通するようになり、これらの検索や分類のために、計算機による自然言語処理が不可欠になっている。自然言語処理の適用では語彙的な知識を電子化した辞書が必要とされるが、処理対象となる電子文書が急激に大規模化・多様化する中で、辞書の整備は必ずしも十分に行われていない。このような背景のもと申請者の博士論文は、電子化辞書の構築支援を目的として、その重要な構成要素である「語の相互関係」を自動的に獲得する手法の確立を目指している。具体的には、類義関係、多義関係、上位下位関係という3種類の語の関係に注目して、あらかじめ蓄積した日本語の文書集合から、これらを自動的に獲得する手法の提案と検証を行っている。

博士論文における研究のポイントをまとめると以下のようになる。第一に、本研究では、「観点」を考慮した類義関係の抽出法を新たに提案した。ある二語の間に類義関係が存在するか否かは観点に依存するが、従来の類義語獲得の研究では、二語間の類似度はつねに一意の数値で与えられるため、類義関係の観点による違いを扱うことが困難であった。これに対して本研究では、まず語の共起関係から二部グラフを構成し、次に極大完全2部グラフを列挙する高速アルゴリズムを適用して類語集合候補を抽出する方法を提案し、検証実験によってその有効性を示した。第二に、本研究では、従来取り扱われることの少なかった「多義性」の抽出に取り組み、同一の語に対する複数の語義の獲得が可能であることを示した。従来、語の多義性については、あらかじめ辞書に登録された語義を文脈情報に基づき選択する語義あいまい性解消を中心に研究が進められてきた。これに対して本研究では、多義語自体を抽出する問題に取り組み、類義関係抽出の場合と同様に列挙した類語集合候補を要素の重なりに基づき併合する抽出法を提案し、検証実験によって実際に多義語が抽出できることを示した。第三に、本研究では、語の上位下位関係をテキストから獲得する手法を新たに提案した。具体的には、語のIS-A関係を名詞が指示する対象を包含関係によって定義し、動詞と名詞の依存関係を用いてこれを発見する方法を提案した。検証実験によって、従来手法では獲得できなかつたIS-A関係が獲得できることを確認し、特に暗黙のうちに書き手が想定する上位下位関係の抽出において、提案手法が有効であることを示した。

このように申請者の博士論文では、基本的な語の相互関係である類義・多義関係と階層関係について、計算機による獲得手法を提案し、実験により有効性を確認している。本研究で扱う語彙知識は、形態素解析や構文解析などの基本的な自然言語処理ツール、情報抽出や要約などの高度な意味処理、分類やクラスタリングを含む情報検索など、多くの言語処理の基本になる資源であり、言語コーパスの解析による自動獲得の試みは、様々な方面への貢献が期待できる。

以上のような博士論文について、平成 20 年 1 月 7 日に審査委員全員の出席のもと公開の博士論文発表会を開催した。まず出願者が提出博士論文の内容に関する 40 分間のプレゼンテーションを行い、その後、博士論文の内容について 20 分の質疑応答を行った。発表会において申請者は、研究の背景や着想にいたった経緯、提案手法と有効性について明確に発表するとともに、研究の内容に関する質問に的確に返答をして、問題点や関連分野の中での研究の位置づけについて深い洞察があることを示した。

6 名の審査委員による博士論文審査の結果、本論文は、博士（情報学）の学位論文として十分な価値があると、全員一致で認められた。