



新しい語・崩れた表記の処理

基
般

笹野 遼平¹ 鍛冶 伸裕²

¹ 東京工業大学精密工学研究所 ² 東京大学生産技術研究所

情報発信者の多様化

ブログ、ミニブログ、SNS (Social Networking Service) などのCGM (Consumer Generated Media) の一般化により、さまざまな情報がこれらのメディアから発信されるようになった。それに伴い、多様な発信者が書いたテキストを目にする機会が増えてきており、これらのテキストを分析することにより、消費者の生の評価やニーズを把握したり、通時的変化を分析することで流行を定量的に観察することが可能となった。

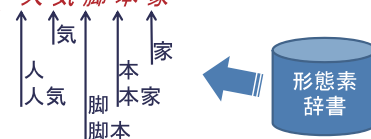
しかし、このような多様な発信者が書いたテキストには、「ググる」「婚活」といった辞書に載っていない新しい語や、「知らない」「もしも〜し」のような崩れた表記など、新聞記事を主な対象としてきた従来のテキスト処理技術では対応できない表現が多く出現する。本稿ではこれらの表現への対処の試みを紹介する。

現実世界のテキストと形態素解析

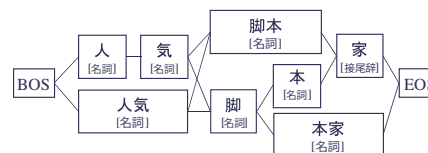
多くの自然言語処理アプリケーションでは最初に形態素解析と呼ばれる処理が行われる。形態素解析とは入力されたテキストを、意味を持つ最小の単位である形態素に分割する処理であり、一般に図-1に示すような手順で行われる。この際、手順①において作成される形態素の候補は、基本的に事前に準備した辞書に含まれる形態素から作成される。このため、辞書に含まれていない表現を含む文は正しく解析できない場合が多い。

入力テキスト: “人気脚本家”

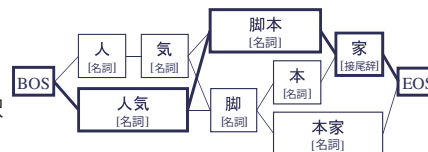
手順①: 各位置から始まる形態素の候補を形態素辞書から検索



手順②: 可能な形態素列を列挙した形態素ラティスを作成



手順③: 文として最も適切であると考えられる形態素列を選択



解析結果: 人気(名詞)-脚本(名詞)-家(接尾辞)

図-1 形態素解析の流れ

現在、一般的に使用されている形態素解析システムは新聞記事のような整ったテキストに対しては高い精度を実現している。しかし、「ググる」という新しい語や「知らない」「ありがとー」などといった崩れた表記を含む図-2に示すようなWeb上のテキストに対しては十分な精度を実現しているとはいえないのが現状であり、Webテキストに代表されるような現実世界のテキストに対応した言語処理技術が求められている。

新しい語の処理

我々が生活する社会においては、絶えず新しい語（以下では単に新語と呼ぶ）が作り出されている。たとえば「ググる」「婚活」などは、今から10年前には日本語に存在していなかった単語である。

かわらず、辞書に登録されていない片仮名列は新語であると考えることができるため、これによって新語抽出を行うことが可能となる。もし、新語の品詞推定を行いたいのであれば、2値分類ではなく多値分類問題を扱えばよい。

こうした分類問題に対しては、機械学習のコミュニティにおいて研究されているさまざまな教師あり学習手法を適用することができる。実際に研究の中で使用したのは SVM (Support Vector Machine) である。ご存知の読者も多いだろうが、SVM は最大マージン原理に基づく分類器であり、高い分類性能を発揮することから、自然言語処理を含めた多くの分野において広く用いられている。

SVM に限らず、一般的に分類器の学習を行うためには、分類対象（ここでは片仮名列）を特徴量ベクトルに変換する必要がある。新語抽出を行う場合であれば、片仮名列の直後に出現する文字 n -gram^{☆3} を特徴量として使うことができる。すなわち、各片仮名列はテキストに出現した文字 n -gram の異なり数と同じだけの次元数を持った特徴ベクトルへと変換される。特徴ベクトルの次元は個別の文字 n -gram に対応しており、今注目している片仮名列の直後にその文字 n -gram が出現すれば、対応するベクトルの要素は 1、そうでなければ 0 となる。このとき n の値は自由に定義できるが、実験では文字 1-gram から 5-gram までを用いた。

我々の実験において抽出された新語の一部を以下に示す（文献 1）からの抜粋）。

- コラボる, トイツる, ジコる, テソパる, デモる, タクる, ラチる, ヘチる
- イナタい, スンバラシい, ウツザい, ナヨい, ヘヴィい, ズブい

「コラボる」「トイツる」「イナタい」など、新語が抽出されていることが分かる。こうした新語は、従来の形態素解析技術では正しく解析することが困難であるとされていたが、このようにして構築された新語辞書を活用すれば、正しく解析することが可能

☆3 文字 n -gram とは連続する n 文字の文字列のこと。

タイプ	例 (括弧内は元の語)
長音記号の挿入*	でーす、もしも〜〜し
母音の挿入	やったあ、行けえええ
小書き文字の挿入*	見たああい、ねむうい(眠い)
促音・発音の挿入	すっばらしい、すんばらしい
長音記号による置換*	ありがとー、ねーさん(姉さん)
小書き文字による置換*	おいしい(おいしい)、かわいい(かわいい)
類似記号による置換	あやしい(怪しい)、こわばわわ(こんばんは)

表-1 Web テキストに出現する崩れた表記 (* は JUMAN7.0 で対応済みであることを表す)

になる。

崩れた表記の処理

Web テキストなど現実世界のテキストの処理を行う上で問題となる表現には、新語以外にも、「おいしい」や「もしも〜〜し」などといった崩れた表記がある。ここからは、このような崩れた表記を含んだテキストの処理を実現するための取り組みの 1 つとして、形態素解析システム JUMAN7.0²⁾ における取り組みを紹介する。

■ Web テキストに出現する崩れた表記

まず表-1 に、Web 上のテキストに出現する代表的な崩れた表記を示す。これらはいずれも辞書に登録されている語（以下では既知語と呼ぶ）に、長音記号や母音字、小書き文字が挿入されたり、一部の文字が小書き文字等に置換された表現であり、一般的な形態素解析システムでは正しく解析できない。

- おいしかったでーす

たとえば上記のようなテキストが入力された場合を考えると、「おいしい」(形容詞) や「です」(助動詞) という形態素が辞書に登録されていたとしても、「おいしかった」や「でーす」という表記に対応していないため正しく解析することができない。

形態素解析における辞書に含まれていない表現への対応策としては、前章までで紹介したように大量のコーパスから事前に新語や業界用語などといった未知語を獲得しておく手法と、統計情報や機械学習を用いて未知語モデルを学習する手法の 2 つがよく

タイプ	削除対象	削除する条件
長音記号の削除	ー、～	1. 直前が平仮名、または漢字 2. 一部の品詞（接頭辞、格助詞等）でない
小書き文字の削除	あ、い、う、え、お	直前の文字が平仮名で、かつ、削除対象がその平仮名を長音化させる文字である場合（e.g. 削除対象が「あ」で、かつ、直前の文字が「か」）
タイプ	置換対象	置換する文字
長音記号の置換	ー、～	直前の文字が 1. 「が」「ば」「ま」「ゃ」なら「あ」 2. い段、「え」「ね」以外のえ段なら「い」 3. う段またはお段なら「う」 4. 「え」「ね」なら「え」
小書き文字の置換	あ、い、う、え、お、わか、	1. 「あ」⇒「ア」 2. 「い」⇒「イ」 3. 「う」⇒「ウ」 4. 「え」⇒「エ」 5. 「お」⇒「オ」 6. 「わ」⇒「ワ」 7. 「か」⇒「カ」

表-2 崩れた表現の正規化ルール

用いられる。しかし、これらの手法はいずれも、既知語からの派生ではない“完全な未知語”を扱う場合に適した手法であり、本稿で対象としている崩れた表記を扱うのに適した手法であるとは言えない。

まず、未知語を事前に獲得する手法には多様な表記バリエーションに対応できないという問題がある。たとえば、「もしもし」という語に対して「もしも～し」「もしも～～し」などといった表記を仮に獲得することができたとしても、「もしも～～し」や「も～しも～し」という表記を獲得していなければ、これらの表記に対応することはできない。未知語モデルを学習する手法は、このような表記バリエーションに対応できる可能性はあるが、大量の学習データが必要となり、また、解析速度も一般的に大きく低下する。

さらに、いずれの手法にも元となった語に関する知識を活用できないという問題がある。人間が「おいしい」や「もしも～～し」などといった文字列をどのように理解するかを考えると、仮にそれらの文字列を見るのが初めてであったとしても、それらを完全に未知の文字列として語の区切りや意味を推定するのではなく、それぞれ「おいしい」、「もしもし」から変形した語として語を区切り、意味を理解していると考えられ、計算機による解析を行う場合も既知語に関連付けて解析するのが自然であると考えられる。

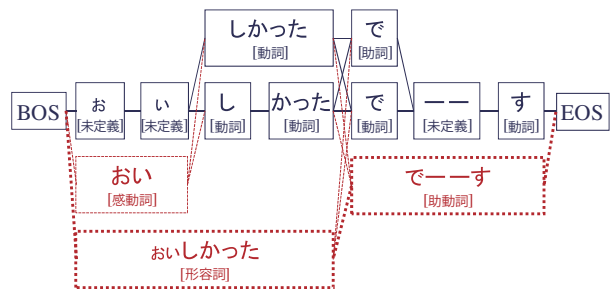


図-3 形態素解析における崩れた表記の認識

■形態素解析における崩れた表記の認識

そこでJUMAN7.0では、一般の未知語処理^{☆4}とは別に、形態素解析時に崩れた表記を既知語と関連付けることで、崩れた表記に対処している。具体的には、形態素解析において各位置から始まる形態素の辞書引きを行う際（図-1における手順①）、入力文字列に対する辞書引きに加えて、表-2に挙げたような正規化ルールに基づいて入力文字列を変形し、変形した文字列に対しても辞書引きを行うことで、崩れた表記に対応している。

たとえば、先述の「おいしかったでーす」というテキストの解析を行う場合、長音記号を削除した文字列「おいしかったです」や、小書き文字を置換した文字列「おいしかったでーす」に対しても辞書引きを行い、図-3において破線で示された「おい」（感動詞）、「おいしかった」（形容詞）、「でーす」（助動詞）などの形態素候補を形態素ラティスに追加し、最終的に「おいしかった」（形容詞）と「でーす」（助動詞）の2形態素に正しく分割できるようになっている。JUMAN7.0で採用されている手法の主な特長として次の4点がある。

1. 任意の文字数の挿入に対処できる
2. 学習データを必要としない
3. 解析速度にほとんど影響がない
4. 既知語に関する知識を活用できる

まず、入力文字列を規則的に変形した上で辞書引きを行っているため、多くの表記バリエーションに

☆4 JUMAN7.0では「ググる」のような一般の新語・専門用語等への対応は人手で行うのではなく自動獲得によって行うべきであるとの考えに基づき、コーパスから自動的に構築した辞書を付属している。詳細は文献3)を参照のこと。

タイプ	1万文あたり の改善数	1万文あたり の悪化数	解析速度 の低下率	解析が改善 する入力例
長音記号 の挿入	108	0	2.0%	・もしも〜し ・ぜーんぶ
小書き文字 の挿入	36	0	0.4%	・行くぞお ・コチラでえす
長音記号 による置換	51	1	0.8%	・うらやまし〜 ・いー感じ
小書き文字 による置換	137	2	0.7%	・ばあちゃん ・書いていい

表-3 崩れた表記の解析精度と速度低下率, および, 解析が改善する入力例

対応できる。たとえば、「もしも〜し」や「もしも〜〜し」、「も〜しも〜し」などといった表記があったとしても、いずれも「もしもし」という文字列に直した上で辞書引きを行うため、いずれの表記に対しても同じように対処することが可能である。

また、基本的に辞書引きの方法を改良しているだけであるため、学習データを必要としない。JUMAN7.0では、通常の辞書引きにより生成された形態素候補を、正規化ルールを適用することで新たに追加された形態素候補より優先するため、後者に一定のペナルティコストを与えているが、このコストの設定に必要となる事例はごく少数である。

さらに、正規化ルールが適用された場合のみ新たな辞書引きを行うようにすることにより、解析速度の低下を最小限に抑えることができる。実際にWebテキストを解析した場合の解析速度の低下率を、1万文あたりの改善数、悪化数、改善例とともに表-3に示す。

JUMAN7.0で採用されている手法には、崩れた表記を辞書に載っている一般的な表記と関連付けることから、既知語に関する既存の知識を活用できるという特長もある。すなわち、崩れた表記で出現した場合であっても、JUMANが使用している辞書に付与されているカテゴリやドメイン、反義語などの情報にとどまらず、大規模なシソーラスなど元となった語に関する種々の言語リソースを活用することが可能となる。

本章で紹介した手法は、文字列の削除・置換を行

うためのコストを設定する必要はあるものの、基本的に辞書引きを改善することにより崩れた表記への対処を行っている。このため、各形態素の生起コストや接続コストの設定・推定法とは独立しており、ChaSenやMeCabなどのような機械学習により各種コストを推定するシステムにも応用することが可能であることを最後に付け加えておく。

応用事例：言語学研究の支援

ここまで、辞書に含まれていない表現の解析に関する研究紹介を行ってきた。それでは、これらの表現を正しく解析できるようになったとして、その先にはどのような応用可能性が考えられるのであろうか。ここでは、新語処理技術の応用事例として、言語学研究支援に関する試みを紹介する⁴⁾。

言語学とは、言うまでもなく、種々の言語現象の理論化を目的とした学問である。言語現象の分析を行う上で、大きな問題となるのが、言語データの包括的な収集が困難なことである。たとえば、言語学者が新語に関する研究を行おうとした場合、世の中で使われている新語に関するデータ収集を行う必要が生じる。しかしながら、これが非常に困難な作業であることは論を俟たない。

こうした問題解決のため、時系列Webテキストと新語処理の技術を利用することにより、新語の言語学的分析の支援を行う研究を進めている。Webテキストには新語の用例が豊富に存在する。そうしたテキストを言語処理技術を用いて解析し、新語の用例を大規模に収集することによって、分析作業の網羅性および効率性を向上させることが狙いである。

言語学研究支援の一環として、新動詞の通時変遷を分析するための基盤構築に取り組んできた。

図-4は、新動詞「ファブる」と「バルビる」の使用頻度の時間変化を比較したものである^{☆5}。この図から、新動詞「ファブる」は世の中に広まりつつ

☆5 それぞれ「ファブリーズ（消臭剤の名称）を使う」「バルビレッジ（ゲームの名称）をプレイする」という意味の新語である。

あるが、逆に「バルビる」の流行は一過性であったことを見とることができる。

このように大規模な Web テキストと新語処理技術を組み合わせることによって、新語の盛衰という、これまで観測困難であった言語データを簡単に取得できるようになることがお分かりいただけたかと思う。これは新語処理技術の応用の一例であるが、CGM を中心とする新世代のテキストメディアと言語処理技術の融合によって、新たな価値や知識が創出される可能性を感じ取っていただければと思う。

実世界の言語処理に向けて

本稿では、新しい語や崩れた表記の扱いに焦点を当てて、従来技術では対処しきれなかった言語表現の形態素解析処理に関する取り組み、およびその応用事例を紹介した。

こうした表現は、従来の自然言語処理研究においては例外的な言語現象とみなされ、議論の中心となることは少なかった。しかし、これは強い言い方をすれば、新語や崩れた表記が存在しない、現実世界から乖離した言語が、自然言語処理の対象として暗黙のうちに想定されてきたということでもある。ブログを始めとする CGM の出現は、こうした問題について我々研究者が再考する良い機会であると言える。

CGM テキストは、少なくとも当面の間は自然言語処理において重要な研究対象であり続ける可能性

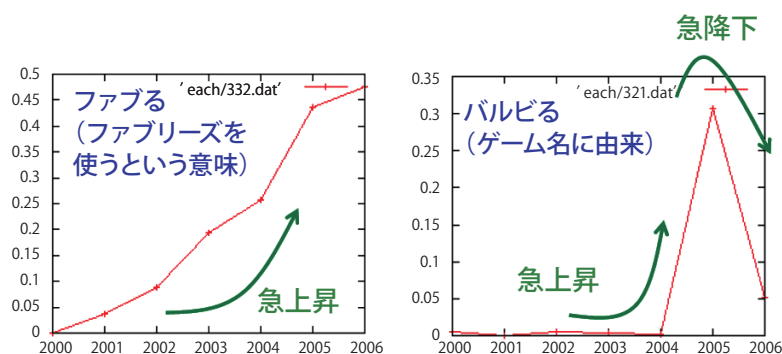


図-4 新動詞「ファブる」と「バルビる」の使用頻度の通時変化の比較。縦軸は使用頻度、横軸は時間を表す。

が高い。そのため、現在の自然言語処理技術が処理対象として想定している言語と、現実世界において使用されている言語の差異を埋めていくことが、今後ますます重要になってくると考えられる。新語や崩れた表記に関する一連の研究が、そのような潮流を形成する一助となればと思う。

参考文献

- 1) 鍛冶伸裕, 福島健一, 喜連川優: 大規模ウェブテキストからの片仮名用言の自動獲得, 電子情報通信学会論文誌, Vol. J92-D, No.3, pp.293-300 (2009).
- 2) 黒橋禎夫, 河原大輔: 日本語形態素解析システム JUMAN version 7.0 使用説明書, 京都大学大学院情報学研究所 (2012).
- 3) 村脇有吾, 黒橋禎夫: 形態論的制約を用いたオンライン未知語獲得, 自然言語処理, Vol.17, No.1, pp.55-75 (2010).
- 4) 宇野良子, 鍛冶伸裕, 喜連川優: 新動詞の成立にみる意味と形の変化の相関—「ファブる」と「モブる」の分析から—, 日本認知言語学会論文集第 10 巻 (2010).

(2011 年 11 月 18 日受付)

笹野 遼平 (正会員) | sasano@pi.titech.ac.jp

2009 年東京大学大学院情報理工学系研究科博士課程修了。博士 (情報理工学)。2010 年より東京工業大学精密工学研究所助教。自然言語処理、特に述語項構造解析、照応解析の研究に従事。

鍛冶 伸裕 (正会員) | kaji@tkl.iis.u-tokyo.ac.jp

2005 年東京大学大学院情報理工学系研究科博士課程修了。情報理工学博士。現在、東京大学生産技術研究所特任助教。自然言語処理の研究に従事。