

近代の歴史的資料を対象とした機械学習による文境界推定

白井 良介^{1,a)} 松村 雪桜^{1,b)} 小木曾 智信^{2,c)} 小町 守^{1,d)}

受付日 2019年5月7日, 採録日 2019年11月7日

概要: 本稿では、機械学習を用いて近代の歴史的資料に対して文境界を検出する手法を提案する。近代の歴史的資料は明確な文境界が必ずしも存在しないため、これまで人手作業による文境界の付与が行われてきたが、膨大な資料に対してなかなか作業が進んでいない現状がある。そこで我々は機械学習を用いて文境界を検出する手法を提案する。この手法により膨大な量の資料に対して文境界の一次的なアノテーションを施すことができることに加えて、形態素解析の精度を向上させたことが本研究の貢献である。また、モデルの訓練に日本語の近代語のデータを使用して、複数の機械学習手法を比較して近代の歴史的資料を対象とした文境界推定を行うのは本研究が初めてである。

キーワード: 近代文語, 文境界推定, 近代語コーパス, UniDic

Machine Learning-based Sentence Boundary Detection for Modern Japanese Texts

RYOSUKE SHIRAI^{1,a)} YUKIO MATSUMURA^{1,b)} TOSHINOBU OGISO^{2,c)} MAMORU KOMACHI^{1,d)}

Received: May 7, 2019, Accepted: November 7, 2019

Abstract: In this study, we propose a method to detect sentence boundaries for modern Japanese texts using machine learning. For modern Japanese texts, sentence boundaries are not explicitly marked so that human annotation is inevitable, but the annotation process is far from complete due to enormous number of materials. Therefore, we propose a method to detect sentence boundaries using machine learning. The main contribution of this study is that this method can support the annotation task as a primary annotation. We also show that the accuracy of morphological analysis can be improved by performing sentence boundary detection. Moreover, this is the first work to detect sentence boundaries targeting modern Japanese texts by using modern Japanese data for model training and comparing multiple machine learning methods.

Keywords: modern historical language, sentence boundary detection, modern historical corpus of Japanese, UniDic

1. はじめに

文境界推定は多くの自然言語処理の分野において必要不可欠となる要素技術である。形態素解析や固有表現抽出、係り受け解析などのタスクでは、文書ではなくそれぞれの

文に対して解析を行うため、正しい文境界が定まっていることが前提になっている。

現代の日本語の書き言葉の文境界は、文章中に挿入された住所の記述や括弧類などの一部の例外は存在するものの、句点^{*1}や感嘆符、疑問符を手がかりにすることで正しい文境界を付与することが比較的容易である。その一方で、Twitter や Facebook などのソーシャルメディアの投稿やマイクロブログなどのウェブテキストのように自由記述形式のものや、話し言葉の書き起こし、また歴史的資料においては手がかりが曖昧であり、ルールベースで文境界を付

¹ 首都大学東京
Tokyo Metropolitan University, Hino, Tokyo 191-0065, Japan

² 国立国語研究所
The National Institute for Japanese Language and Linguistics, Tachikawa, Tokyo 190-8561, Japan

a) ryosuke@komachi.live

b) matsumura-yukio@ed.tmu.ac.jp

c) togiso@ninjal.ac.jp

d) komachi@tmu.ac.jp

*1 本来、“句点”は“記号の形”と“記号の機能”の2つの意味を持つ術語であるが、本稿では文境界を表す“記号の機能”の意味で使用する。



図 1 本研究の位置付け

Fig. 1 Role of this study.

与することが困難な場合が存在する。特に、近代の歴史的資料に対して文境界を付与することは近代語の知識のある専門家の手によらなければ難しく、膨大な量の資料の前に作業がなかなか進まないでいるのが現状である。現在、近代の歴史的資料に対しては専門家らによる人手のアノテーションが行われている [23] が、まだアノテーションのなされていない膨大な量の資料が存在する。

国立国語研究所での近代の歴史的資料に対するアノテーションでは、生のテキストからはじまり、最終的に人手による修正を経た高精度な形態素解析済みのコーパスの作成までを行い、日本語史研究用の資料として広く研究者に提供している。同研究所では、原本からデジタルデータへの書き起こしがされたものに対して句読点などを仮の文境界として形態素解析をしたノンコアデータ、それに対して専門家らによる人手の修正がなされたコアデータの 2 種類のデータを用意している。本研究の対象としている文境界推定は、書き起こしがされたデータに対する前処理の段階で使用され、ノンコアデータの代わりとして使用されることを目指している。このノンコアデータは国立国語研究所が公開しているコーパス検索アプリケーション『中納言』でも用いられるため、一般のコーパス利用者にとっても役立つことが期待される。本研究の位置付けを図 1 に示した。専門家が翻刻作業をすることで原本である紙から電子の生テキストへの翻刻の際に同時に文境界を付与するということが可能であるが、紙から電子への翻刻は専門家でない作業者が担当しており、同時に行うことは現実的には難しい。そのため文境界推定のタスクが果たす役割は大きいといえる。

そこで、本研究では、近代の歴史的資料を対象に機械学習による文境界推定を行う。ルールベースに対して複雑な素性を扱うことができる機械学習を用いた文境界推定を行うことで、膨大な量の資料に対して人手の修正が行われる前段階の一次的なアノテーションを改善することができるということが本研究の貢献である。本研究で主たる対象としているのは明治・大正期の文語論説文であり、近代に成立した新聞・雑誌などのマスメディアを通して人々に広く読まれた文章がこの文体で書かれている。また、昭和前期までの法律文や行政文書などもこれに近い文体で書かれて

いるため、将来的には近代語の研究だけでなく近代日本を研究するうえで重要な多数の資料での利用も見込まれる。また、モデルの訓練に日本語の近代語のデータを使用して、複数の機械学習手法を比較して近代の歴史的資料を対象とした文境界推定を行うのは本研究が初めてである。

本研究で文境界推定を行う資料は、1895 年 (明治 28 年) から 1928 年 (昭和 3 年) に博文館より発行された総合雑誌『太陽』を対象とし、データは『太陽コーパス』[14] の文語コアデータを用いた。句点で文境界を付与するルールベースのものと同様に句読点で文境界を付与するルールベースの 2 つをベースラインとし、『太陽コーパス』のみを用いて学習したモデル、『太陽コーパス』に『太陽』と同時代の資料のコーパスである 3 種類の近代文語コーパスを加えて学習したモデルを用いて、文境界推定との異なり具合と、近代語への文境界推定の精度を確認した。機械学習の手法としては、提案手法として短単位の素性テンプレートを用いた条件付き確率場 (Conditional Random Fields: CRF) [4] と、文字単位の GRU (Gated Recurrent Unit) [1] を双方向に用いた Bi-GRU (Bi-directional GRU) を使用した。ベースライン (句点ルールベース) の適合率 94.34%・再現率 34.81%・F 値 50.85 ポイントの精度と比較して、『太陽』に 3 種類の近代文語コーパスを加えて学習した CRF を用いた手法では適合率 83.75%・再現率 73.68%・F 値 78.40 ポイント、Bi-GRU を用いた手法では適合率 75.07%・再現率 62.01%・F 値 67.08 ポイントと F 値を大きく向上させることができた。

上記の実験に加えて、本研究の提案手法で文境界を付与することが具体的にどう役立つかということを確認するために、文境界推定によって得られた文境界を与えて、形態素解析の精度を比較した。『太陽コーパス』に 3 種のコーパスを追加した文境界推定実験で付与した文境界が、ルールベースと比べて 0.02 ポイント高い F 値を得ることができた。ブートストラップ検定を行ったところベースラインに対して統計的に有意 ($p < 0.001$) であり、形態素解析の前処理としても役立つことを示した。

加えて、実際の文境界修正作業を模したアノテーション支援実験も行った。結果として、文書に対してルールベースの文境界が付与されているものに比べて、『太陽』に 3 種類の近代文語コーパスを加えて学習した提案手法による文境界が付与されているほうが文境界修正作業の時間を大きく短縮することができた。

2. 関連研究

現代の書き言葉を対象にした文境界推定は、いくつかの研究が行われている。たとえば、英語では文境界を表すピリオドと “Mr.” などのように文境界を表さないピリオドが存在するため、書き言葉に対する文境界推定を行う必要がある [6]。

日本語の文境界推定の研究として行われているのは、推定の対象として主に Twitter や Facebook などのソーシャルメディアの投稿やマイクロブログなどのウェブテキストのように自由記述形式の書き言葉を対象としたものや、話し言葉の書き起こしを対象としたものである。これらは日本語の書き言葉の文境界を表す句点などの目印が必ずしも付与されていないため、文境界の推定が必要である。

文境界推定の方法には主に機械学習が用いられている。福岡ら [12] は Web およびニュースグループから集めたテキストに対して SVM (Support Vector Machine) を用いて文境界推定を行った。難波ら [13] は Twitter に投稿された Tweet を対象として、文境界推定を系列ラベリング問題として扱い、CRF を用いた文境界推定を行った。CRF の素性には単語と品詞と文字種を使用して実験を行い、同時に文節境界推定と係り受け推定を行った。また日本語以外でも、Rudrapal ら [7] は英語やヒンドゥー語で書かれた Twitter の Tweet や Facebook のメッセージを対象として、CRF、ナイーブベイズ、SVM を用いて文境界推定を行った。話し言葉の書き起こしでは、下岡ら [26] は日本語話し言葉コーパス (CSJ) [21] を対象にして SVM を用いた文境界推定を行った。文境界推定をテキストチャンキングの問題として扱い、テキストチャンカとして SVM に基づく YamCha [2] を用いた。これらの研究のように、ウェブテキストやスピーチの書き起こしではルールに基づく処理が困難なため、機械学習による文境界推定が行われている。本研究の対象である近代の歴史的資料についても、文境界の手かがりとなるものが必ずしも存在しないため、同様に文境界推定を行うことが必要である。

近年の機械学習のスタンダードであるニューラルネットワークを用いた研究では、Straka ら [8] が開発している UDPipe^{*2} というソフトウェアがある。これは文字単位の GRU (Gated Recurrent Unit) [1] を双方向に用いた Bi-GRU (Bi-directional GRU) を系列ラベリング手法として採用し、生テキストに対して単語分割、タグ付け、係り受けまでの解析をサポートしており、単語分割が行われるのと同時に文境界推定も行われている。本研究では上記の関連研究でも使用されている CRF と Bi-GRU の両手法を用いた文境界推定を行った。

生テキストを対象としている CoNLL 2018 Shared Task^{*3} [9] では、生テキストから係り受けまでの解析を共通タスクとしており、古代ギリシア語、ゴート語、ラテン語、古代教会スラヴ語、古フランス語など、それぞれの現代語とは異なる時代の言語も対象言語として含まれている。その中には日本語の近代語も含まれており、文境界推定も行われている。しかし、モデルの訓練には日本語の現代語のデータである Japanese-GSD を使用して F 値

0.23% ときわめて低い精度であり、現代語で機械学習されたモデルを単純に適用するだけでは解析が難しいことが報告されている [5]。本研究では、先行研究のように現代語のモデルを単純に適用するのではなく、モデルの訓練に現代語のデータを適応するのではなく、近代語のデータを直接使用してモデルを訓練することで、実用的な精度で解析することを目指す。また、複数の機械学習手法を比較して近代の歴史的資料に対して文境界推定を行うことは本研究が初めてである。

日本語の近代の歴史的資料に関する先行研究では、形態素解析については小木曾ら [20] によってすでに実現されているが、現代語とは異なる近代語の表記の不完全さにより、一部にしか付与されていない濁点の付与や、必ずしも明示されない文境界 (句読点) の整備という前処理が必要となっていた。このうち、濁点の付与については岡ら [10] によって自動化が可能になっていたが、文境界の整備についてはもっぱら人手に頼っており、その自動化が課題となっていた。本研究は、この課題を解決し、近代語のテキストデータを自動処理で形態素解析まで行って日本語学や人文情報学における研究利用に供する一連の流れを可能にするものである。

3. 機械学習を用いた文境界推定

本研究では、文境界推定を文頭の形態素に対応する B ラベルと文頭でない形態素に対応する I ラベルを予測する BI ラベルの系列ラベリング問題としてとらえ、人手でアノテーションされたデータを用いて CRF と Bi-GRU による機械学習を行うことで、文境界を自動で付与する手法を提案する。

3.1 4つの文パターン

近代の歴史的資料において文境界を推定することが困難な理由としては、4つの文パターンが混在していることがあげられる。表 1 に『太陽コーパス』の文語コアデータにおける各パターンの例文を示し、表 2 に今回実験に用いた各コーパスにおける文パターンの割合と統計情報を示した。句読点混合パターン以外にも句点パターンと読点パターン、そして句読点なしパターンが存在し、後者の3パターンはルールベースで解析することができない。『太陽コーパス』以外のコーパスについては 4.1 節で詳しく述べる。

3.2 CRF

機械学習の手法として CRF を用いる。実装には CRF++^{*4} を使用した。素性には近代文語 UniDic [20]^{*5} で

^{*4} <https://taku910.github.io/crfpp/>

^{*5} 近代文語 UniDic には『太陽コーパス』、『明六雑誌コーパス』、『国民之友コーパス』、『女性雑誌コーパス』に加えて、1つのコーパスとしては扱われていない多くの近代論説文の語彙が収録されている。

^{*2} <http://ufal.mff.cuni.cz/udpipe>

^{*3} <http://universaldependencies.org/conll18/>

表 1 近代の歴史的資料における 4 つの文パターン

Table 1 Four sentence patterns in modern historical materials.

パターンと例文 (文境界を “—” で示す)	
句読点混合パターン：句点と読点で現代語の書き言葉と同じように付与されている	例：— 一は歐羅巴の海岸線が甚だ複雑なる事にして、一は其上に位する國民の種類甚だ夥多なる事なり。—
句点パターン：句点を読点の役割としても付与している全句点パターン	例：— おや。二個貰ったのか。—
読点パターン：読点を句点の役割としても付与している全読点パターン	段落終わりのみ句点を付与している例外パターンも存在する
例 1：— 記者曰、君は徳太郎と稱し、慶應三年十二月を以て江戸芝神明町に生る、—	例 2：— 豈に戒めざる可けんや、— 豈に懼れざる可けんや。— (段落終)
句読点なしパターン：そもそも句点と読点が付与されていない	例：— 請ふ其の昨年度の形勢を觀察せん — 今昨年五月末日に於ける船舶の統計は左の如し —

表 2 各近代語コーパスにおける文パターンの割合と統計情報

Table 2 Percentage and statistical information of sentence patterns in each modern Japanese corpus.

コーパス	句読点混合	句点	読点	句読点なし	短単位数	文書数	文数
『太陽コーパス』文語	30.8%	3 文のみ	35.7%	33.4%	71,850	33	3,686
『明六雑誌コーパス』	0.0%	0.0%	3.3%	96.7%	179,522	198	9,563
『国民之友コーパス』	11.0%	1 文のみ	21.8%	67.1%	32,154	24	1,479
『女性雑誌コーパス』文語	30.2%	2 文のみ	31.7%	38.7%	39,779	64	2,148

表 3 素性テンプレート

Table 3 Feature template.

N-gram	観測するトークン
uni-gram	$x_{t-2}, x_{t-1}, x_t, x_{t+1}, x_{t+2}$
bi-gram	$x_{t-2}x_{t-1}, x_{t-1}x_t, x_t x_{t+1}, x_{t+1}x_{t+2}$
tri-gram	$x_{t-2}x_{t-1}x_t, x_{t-1}x_t x_{t+1}, x_t x_{t+1}x_{t+2}$

定義される素性のうち、1. 書字形出現形 (orth), 2. 品詞 (pos), 3. 活用形 (cForm), 4. 語彙素表記 (lemma) の 4 種類を用いた。それぞれ、現在のトークンを x_t としたとき、現在のトークンと前後 2 トークンずつの uni-gram, bi-gram, tri-gram の素性を利用する。詳しくは表 3 に示した。

3.3 Bi-GRU

機械学習のもう 1 つの手法として Bi-GRU を用いる。文字単位の GRU を双方向に用いた Bi-GRU を実装している UDPipe というソフトウェアを使用した。UDPipe は CoNLL-U フォーマットのコーパスからモデルを構築し、解析結果を出力するソフトウェアである。構築したモデルを用いて、生テキストに対して単語分割、タグ付け、係り受けまでの解析をサポートしており、この単語分割が行われるのと同時に文境界推定も行う。

4. 近代語に対する文境界推定実験

近代語の文境界推定において、コーパスの形態素に対して系列ラベリングを適用し、様々な学習データのパターンから『太陽コーパス』のコアデータのうち文語データにお

ける文境界推定の性能を比較した。形態素として近代文語 UniDic の短単位 [16] を用いた。評価には B ラベル推定の再現率、適合率、F 値を用いた。CRF の実装には CRF++ を使用した。実験時のパラメータにはツールのデフォルト値を用いた。文字単位の Bi-GRU による文境界推定の実装には UDPipe の単語分割機能により出力される文境界を使用した。実験時のパラメータは予備実験の結果より dimension を 64 に、segment size を 200 に変更し、そのほかはデフォルト値を用いた。

4.1 データ

実験対象の近代語資料として『太陽』の人手で修正が行われているコアデータを用いた。『太陽』は当時最もよく読まれた総合雑誌であり、政治・経済・世界情勢から科学・思想、文学作品までの様々な記事ジャンルが揃っている。『太陽コーパス』には文語・口語の両データが存在するが、近代の資料には文語体で記述されたものが多いことを考慮して、より多くの資料に対して文境界を推定できるモデルを構築するために文語データのみを用いて 5 分割交差検証を行った。5 分割は全 33 文書からなる文語データをランダムに 7 文書または 6 文書ずつ抽出することにより行った。

また、学習データの不足を考慮して、追加の学習データとして『太陽コーパス』と同じく近代語コーパスである、『明六雑誌コーパス』[23]、『国民之友コーパス』[23]、『女性雑誌コーパス』[22] の 3 種を用いることにした。『女性雑誌コーパス』については、『太陽コーパス』と同様に文語・口語の両データが存在するため、文語データのみを用

表 4 文境界推定 実験結果

Table 4 Results of sentence boundary detection.

文境界推定の手法名	単語分割手法	文分割手法	学習データ	適合率	再現率	F 値
句点ルール	MeCab	句点	—	94.34%	34.81%	50.85
句読点ルール	MeCab	句読点	—	42.92%	61.89%	50.67
Bi-GRU	Bi-GRU	Bi-GRU	『太陽』のみ	72.90%	63.94%	68.13
Bi-GRU	Bi-GRU	Bi-GRU	『太陽』+ 3 コーパス	75.41%	66.68%	70.82
句読点ルール + CRF	MeCab	CRF	『太陽』のみ	95.00%	34.51%	50.63
句読点ルール + CRF	MeCab	CRF	『太陽』+ 3 コーパス	82.87%	73.76%	78.05
Bi-GRU+CRF	Bi-GRU	CRF	『太陽』のみ	95.00%	34.51%	50.63
Bi-GRU+CRF	Bi-GRU	CRF	『太陽』+ 3 コーパス	82.26%	74.65%	78.38

いた。それぞれのコーパスの総短単位数と総文数を表 2 に示した。

本研究の文境界推定は生テキストに対して用いることを想定しているため、CRF を用いた実験では人手で修正された形態論情報が付与されているコアデータではなく、自動解析結果であるノンコアデータに相当するものをテストデータとして使用する必要がある。そこで、2つの疑似的なノンコアデータを作成した。1つは句読点を文境界として文分割を行って MeCab を使用して形態素情報を付与したもので、もう1つは UDPipe の単語分割によって得られた文境界で文分割を行って同じく MeCab の部分的解析機能を使って単語境界以外を推定したものである。どちらも辞書には近代文語 UniDic を使用した。

4.2 手法

ベースラインとして、句点を文境界とする1つめのルールベース手法(句点ルール)と、句読点を文境界とする2つめのルールベース手法(句読点ルール)、UDPipeによる文字単位の Bi-GRU を使用した手法(Bi-GRU)を用意した。提案手法として句読点ルールの形態素解析結果を使った CRF と、UDPipe の形態素解析結果を使った CRF で実験を行った。機械学習を用いた手法では文書を入力とした。

評価には B ラベル推定の適合率、再現率、F 値を用いた。また、各文書の最初のトークンが B ラベルであることは自明なので、該当するトークンは評価対象から外した。

4.3 実験結果

表 4 に適合率・再現率・F 値を示した。実験の結果、ベースラインのうち精度がより高かった句点ルール手法と比較して『太陽コーパス』に3種のコーパスを追加した句点ルール + CRF 手法の実験では 27.20 ポイント高い F 値を得ることができ、『太陽コーパス』に3種のコーパスを追加した Bi-GRU+CRF 手法の実験では 27.53 ポイント高い F 値を得ることができた。

4.4 考察

句読点ルール手法を除き、いずれの実験結果でも適合率

表 5 BCCWJ ルールベース適用結果

Table 5 Results of sentence boundary detection for BCCWJ by rule-based method.

文書	適合率	再現率	F 値
OC (知恵袋)	81.06%	86.39%	83.64
OW (白書)	97.69%	63.18%	76.74
OY (プロゲ)	80.98%	60.63%	69.35
PB (書籍)	97.64%	87.29%	92.18
PM (雑誌)	97.62%	73.91%	84.12
PN (新聞)	99.56%	71.99%	83.56

と比較した際の再現率が低く、文境界があるべき場所を正しく検出するのは難しいという傾向があることが分かった。

4.4.1 現代語との比較

句点ルール手法の結果における再現率は全体で2番目に低く、近代の資料に対しては現代語と同じように句点を文境界として扱う手法では文境界推定のカバー率を上げるのが難しいことが分かる。ここで比較対象としている現代語の例として、現代日本語書き言葉均衡コーパス(BCCWJ) [15] のコアデータに対して句点・感嘆符・疑問符を文境界として扱った際の B ラベル推定の再現率・適合率・F 値を表 5 に示した。現代語では BCCWJ のいずれのジャンルにおいても再現率6割以上になっているが、近代の資料に対しては、機械学習を用いなければ適合率を高く保ったまま再現率を上げることができないことが分かる。

そのほかの特徴として、近代文語では動詞“あり”などのラ行変格活用の語が頻出するが、現代語であればともに“ある”となる終止形と連体形がそれぞれ“あり”と“ある”で異なる一方、近代文語では終止形と連用形がともに“あり”で同形となる。そのため、文境界認定のうえでは現代語とは異なって、終止と連用中止の区別が付けにくいという特徴がある。

4.4.2 データ量の比較

学習用データとして『太陽』のみを用いた手法と、『明六雑誌』・『国民之友』・『女性雑誌』それぞれのコーパスを追加して用いた手法では、後者の方が再現率・F 値が高くなっている。このことから同じ近代の文語体で書かれているデータを追加することで再現率を上げることができるこ

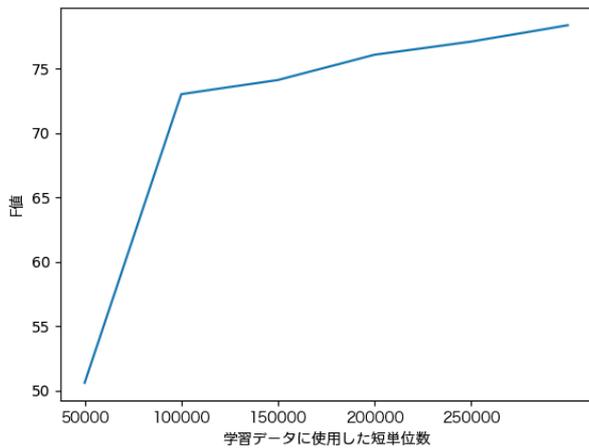


図 2 Bi-GRU+CRF モデルにまず『太陽』のデータを加え、順次 3 コーパスを加えていった場合の文境界推定の学習曲線

Fig. 2 Learning curve for sentence boundary detection using Bi-GRU+CRF model in the case where the data of “THE SUN” is added first, and then 3 corpus are added sequentially.

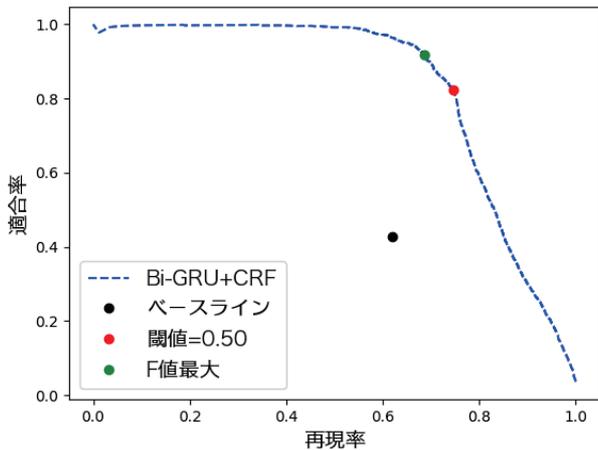


図 3 Bi-GRU+CRF モデルの PR 曲線

Fig. 3 Precision-Recall curve by Bi-GRU+CRF model.

とが確認された。

『太陽』のみを学習に用いた場合では、句点ルール手法とほとんど変わらない精度であることが分かる。また、句読点ルール手法と『太陽』のみで学習した Bi-GRU では後者のほうが F 値が 20.15 ポイント高いにもかかわらず、『太陽』のみを学習に用いた CRF のモデルで文境界推定を行うと両者の精度に違いは現れなかった。このことから、前処理も精度に影響を与えるが、CRF のモデルによる影響のほうがより大きいことが確認できる。加えて、『太陽』にそれぞれのコーパスの短単位を追加していったときの学習曲線を図 2 に示す。この曲線が示すように、同じく近代の文語体で書かれているデータを追加することで F 値を上げることができることが分かる。

図 3 に一番精度が高かった Bi-GRU+CRF 手法の PR 曲線を示す。適合率を句読点ルール手法と同じ 42.92% になるように調整した (閾値 = 0.02) とし、Bi-GRU+CRF

表 6 False Negative (FN) の頻出のエラー 上位 5 件

Table 6 Top 5 most frequent error in False Negative.

間違えたトークン	全 FN に占める割合
と	8.41%
全角空白スペース	3.88%
◎	2.04%
其	1.51%
今	1.05%

手法の再現率は 85.11% となり、句読点ルール手法の再現率を 23.22% 上回った。また、F 値が最大 (閾値 = 0.89) となるように閾値を調整しなくても、デフォルトの閾値でもルール手法より大きく適合率と再現率が向上していることから、Bi-GRU+CRF 手法が句読点ルール手法に対して精度的に優れていることが確認できる。なお、閾値には CRF の解析の周辺確率を用いた。

4.5 エラー分析

最も精度が高い『太陽コーパス』に 3 種のコーパスを追加した Bi-GRU+CRF 手法の実験結果の中で生じた全エラー 1,522 個のうち、False Negative が 927 個 (60.91%)、False Positive が 595 個 (39.09%) であった。再現率を高くするために改善が必要な False Negative (FN) の中から割合の高いものを表 6 に示す。エラーについて考察を述べる際に、“—” で文境界を表す。

個別のトークンとして最も割合の高い、“と” は、“夫れは君の意見に任せる — と言ひます” のように直前が文境界となり“と” が B ラベルになる場合と、“波蘭統監に任ずと —”，“狩野氏と志筑氏” のように直前が文境界とならず“と” が I ラベルになる場合があり、“任せる”，“任ず” のように終止形の後ろに“と” が出現する場合でも推定するラベルに異なりがある。“波蘭統監に任ずと —” のような終止形で終わる文では、終わったあとの文末に“と” を終端記号のように付与していることが特徴である。“と” には格助詞、接続助詞、係助詞など様々な用例があり、品詞推定によってこれらの用例は区別できるが、品詞の同定にも曖昧性があるため、品詞を素性に使用していたにもかかわらずエラーが多くなってしまったと考えられる。

2 番目に割合の高い“全角空白スペース”は、段落始めに頻出する形態素であるが、1 つで出現することもあれば数個続いて出現することもあり、小見出しでは文末に付与されることもあるため、B ラベルと I ラベルが混在しやすい特徴がある。また、歴史的資料によく見られる決まりごとである、皇室関係者の名前を記す際には敬意を表して該当する用語の前に空白を付する関字の影響もあり、識別が困難であったと考えられる。

3 番目に割合の高い“◎”は、主に文書中の小見出しのような文に付与されている記号であるが、今回は最も頻出の

表 7 形態素解析 実験結果

Table 7 Results of morphological analysis.

文境界推定の手法名	単語分割手法	文分割手法	学習データ	適合率	再現率	F 値
句読点ルール	MeCab	句読点	—	94.85%	94.88%	94.86
句読点ルール + CRF	MeCab	CRF	『太陽』 + 3 コーパス	*94.88%	*94.89%	*94.88
Bi-GRU+CRF	Bi-GRU	CRF	『太陽』 + 3 コーパス	*94.87%	*94.88%	*94.88

エラーの“と”の直後や2番目に頻出の“全角空白スペース”の前後に出現する回数が多く、互いに影響しあって揺れが生じたと考えられる。

残りの頻出エラーである“其”，“今”については、はっきりとしたエラーの傾向が発見できず、出現回数の多い短単位であるため必然的にエラーの発生数が多くなってしまい、合計した結果、エラー頻出率の上位に入ってしまった可能性があると考えられる。

また、エラーに出現する回数が3回以下ときわめて少ないエラーが34.30%にもなる。図2から分かるように、微量ではあるが学習データの増加が精度向上に効果的であることが示されているため、同時代の学習データを増やして改善していくことが期待される。

5. 推定した文境界を用いた検証実験

本研究の提案手法で文境界を付与することが具体的にどのように役立つかということを確認するために、前項で推定した文境界を与えて、形態素解析の精度の比較を行った。加えて、アノテーション支援実験により作業効率の向上を確かめた。

5.1 形態素解析の精度比較実験

4章で推定した文境界を与えて、形態素解析の精度の比較実験を行った。

5.1.1 データ

正解データには『太陽コーパス』文語コアデータの形態素情報を用い、前章で推定した3種類の文境界推定手法による形態素解析結果の比較実験を行った。

3種類とは、句読点ルールによるルールベース、句読点ルール+CRF手法で付与した文境界、Bi-GRU+CRF手法で付与した文境界、である。現在、国立国語研究所のノンコアデータには句読点ルール手法の文境界が付与されているため、ルールベースには句点ルールベースの句点ルールではなく句読点ルールを使用した。

5.1.2 手法

形態素解析にはMeCab[3]を用いた。辞書には近代文語UniDicの学習に使用しているコーパス群から文境界推定実験時の5分割交差検証で用いた『太陽コーパス』のデータについて、近代文語UniDicの収録語彙はデフォルトのまま、学習に使用しているコーパス群から、文境界推定実験時の5分割交差検証で用いた5分割した『太陽コーパ

ス』のデータに対して、それぞれ解析対象とするデータのみを除きパラメータ推定をやり直したものを使用した。つまり、5分割されたそれぞれのデータに対して解析用の辞書を5つ作成した。

また、比較実験のツールにはMevAL^{*6}を使用した。

5.1.3 実験結果

表7に形態素解析実験の適合率・再現率・F値を示した。ブートストラップ検定を行い統計的に有意であるかということも確認を行った。表中の“*”はベースラインに対して統計的に有意($p < 0.001$)であることを示している。実験の結果、ベースラインの句読点ルール手法と比較してBi-GRU+CRFによる文境界推定実験で付与した文境界と句読点ルール+CRFによる文境界推定実験で付与した文境界がそれぞれ0.02ポイント高いF値を得ることができた。また、前述のようにどちらの手法もベースラインに対して統計的に有意であるため、形態素解析の前処理としても役立つことを示した。

5.1.4 考察・エラー分析

エラー分析には形態素解析結果の中から、品詞(pos)、活用形(cForm)、語彙素表記(lemma)を使用した。表8に句読点ルール手法、句読点ルール+CRF手法、Bi-GRU+CRF手法で付与した文境界で形態素解析を行った実験のエラー上位10件と、それらの出現数を示した。活用形が存在しないものについては“*”で表記した。終止形になるべき形態素が連用形と誤検出されているエラーについて、Bi-GRU+CRF手法では出現数が減っていることが確認できる。句読点ルール+CRF手法では句読点ルールと比べて精度は向上しているものの、頻出のエラーについては同程度しか対応できていないことが分かった。

近代語の中で特に活用を間違いやすいものとして、表8の上から1番目と5番目に見られる助動詞“ず”と動詞“あり”がある。表9にそれらの活用表を示す。括弧内は助動詞に接続する際の特殊な活用である。どちらも連用形と終止形が同じ形をしているため、前後の形態素情報が連用形か終止形かの判別の重要な手がかりとなる。提案手法により文境界が決まるとEOSが分かり、後ろの形態素情報がEOSであると分かることは終止形であるという手がかりになるため、連体形と終止形の区別ができるようになる。そのためBi-GRU+CRF手法において形態素解析精度を向上させたと考えられる。4章にて句読点ルール手法の精度

^{*6} <https://teru-oka-1933.github.io/meval/>

表 8 形態素解析エラー上位 10 件と手法ごと出現数 (品詞+活用形+語彙素)

Table 8 Top 10 morphological analysis errors and occurrences by each method. (pos+iForm+lemma)

正解	解析誤り	句読点 ルール	Bi-GRU +CRF	句読点 ルール+CRF
助動詞+終止形-一般+ず	助動詞+連用形-一般+ず	80	70	80
助動詞+連用形-ニ+なり-断定	助詞-格助詞+*+に	59	60	60
名詞-普通名詞-一般+*+者	名詞-普通名詞-サ変可能+*+物	55	55	55
助詞-格助詞+*+に	助動詞+連用形-ニ+なり-断定	51	52	52
動詞-非自立可能+終止形-一般+有る	動詞-非自立可能+連用形-一般+有る	35	29	35
助詞-接続助詞+*+も	助詞-係助詞+*+も	35	35	35
副詞+*+又	接続詞+*+又	28	28	28
動詞-一般+連用形-一般+つく	動詞-非自立可能+連用形-一般+付く	26	26	26
助動詞+連用形-ニ+だ	助動詞+連用形-ニ+なり-断定	22	22	22
接頭辞+*+低	形容詞-一般+語幹-一般+低い	21	21	21

表 9 助動詞“ず”・動詞“あり”活用表

Table 9 Usage table of the auxiliary verb “Zu” and the verb “Ari”.

	未然形	連用形	終止形	連体形	已然形	命令形
助動詞“ず”	ず(ざら)	ず(ざり)	ず	ぬ(ざる)	ね(ざれ)	ざれ
動詞“あり”	あら	あり	あり	ある	あれ	あれ

表 10 アノテーション支援実験に使用するデータ

Table 10 Data used for annotation support experiment.

データ名	句読点の有無	短単位数
A	なし	514
B	なし	582
C	あり	521
D	あり	529

表 11 各実験協力者が文境界を修正したデータと順序

Table 11 Data and order for each experiment cooperator to correct sentence boundaries.

実験協力者	データと順序
W	Ar → Bm → Cr → Dm
X	Ar → Bm → Cr → Dm
Y	Am → Br → Cm → Dr
Z	Am → Br → Cm → Dr

よりも Bi-GRU 手法の精度が高いという実験結果が示されているが、その結果が形態素解析の精度改善の度合いにも影響を与えていると考えられる。

5.2 アノテーション支援実験

次に、句読点ルール手法で文境界が付与されたノンコアデータと提案手法で文境界を付与したデータとで、どの程度作業効率に影響を及ぼすかということを確認するためにアノテーション支援実験を行った。このノンコアデータが実際にアノテーションの前処理として使われている形式である。

5.2.1 データ

表 10 にそれぞれのデータの詳細を示す。句読点の有無については、出現しないデータと出現するデータを 2 つずつ用いることにした。本実験では、『太陽コーパス』文語コアデータの異なる文書から抽出した 4 つのデータ (A, B, C, D とする) を用いる。

5.2.2 実験計画

4 人の協力者に句読点ルール手法もしくは Bi-GRU+CRF 手法で文境界が付与された 4 つの文書を与えて、それぞれ

の文書に対して文境界の修正を行わせた。表 11 に各人がそれぞれどのデータを対象に実験を行ったか示す。データ名の後ろについている“r”は句読点ルール手法で文境界を付与していることを表し、“m”は Bi-GRU+CRF 手法で文境界を付与していることを表す。また、同時に文境界の修正にかかる時間も計測した。計測時間から、どの程度アノテーションを支援することができているかを確認する。

5.2.3 実験結果

各データごとの文境界修正にかかった時間と、各データごとの提案手法-ルールベース間で短縮できた時間を表 12 に示す。“+”はルールベースと比べて多くかかった時間、“-”はルールベースと比べて短縮できた時間を表す。

5.2.4 考察

データ全体として、提案手法による文境界が付与されているデータのほうがルールベースのものに比べて合計して 7 分 37 秒の時間を短縮することができた。また、提案手法を用いた結果をデータ別に見ていくと、句読点がない A, B の両方で短縮に成功し、句読点がある C, D では、C においては短縮できなかったが、D では短縮に成功している。このことから、本研究による文境界の付与は特に句読点が

表 12 各データごとの文境界修正にかかった時間

Table 12 Time spent correcting sentence boundaries with each data.

実験協力者	Ar	Bm	Cr	Dm	合計
W	4分41秒	1分49秒	2分26秒	1分11秒	10分07秒
X	6分41秒	4分50秒	4分04秒	3分20秒	18分55秒
合計	11分22秒	6分39秒	6分30秒	4分31秒	29分02秒
実験協力者	Am	Br	Cm	Dr	合計
Y	4分41秒	7分08秒	4分02秒	3分22秒	19分13秒
Z	5分36秒	5分35秒	3分51秒	3分00秒	18分02秒
合計	10分17秒	12分43秒	7分53秒	6分22秒	37分15秒
時間差 (m - r)	-1分05秒	-6分04秒	+1分23秒	-1分51秒	-7分37秒

ない文書に対して大きく役立つことが確認できた。

6. おわりに

本研究では、近代語の歴史的資料に対する CRF と Bi-GRU を用いた機械学習による文境界推定を行った。専門家によるアノテーションを待っている膨大な量のデータに付与することができる一次的な文境界としての活用が期待される。エラー分析の結果、出現する回数が3回以下ときわめて少ないエラーが33.9%にもなるため、同時代の学習データを増やして改善していくことが期待される。現在、国立国語研究所では近代語のコーパスとして『教科書コーパス』[11]、『東洋学芸雑誌コーパス』[19]、『読売新聞コーパス』[25]の構築が行われている。それらのコーパスが完成し、学習に使用できるデータが増えることでエラーの改善につなげることができると考えられる。データが増えることで提案手法の形態素情報を作る段階に使用している Bi-GRU の精度がさらに向上することが期待され、そこから高精度な形態素情報を用いた CRF による文境界推定を行うことができると考えられるためである。

また、本研究の提案手法で文境界を付与することが具体的にどう役立つかということを確認するために、文境界推定によって得られた文境界を与えて、形態素解析の精度を比較した。提案手法である、『太陽』+3コーパスの学習データから Bi-GRU で単語分割をした文の形態素情報を用いた CRF による文境界推定実験が、ルールベース手法で付与した文境界と比べて0.02ポイント高いF値を得ることができた。また、形態素解析の前処理としても役立つことを示した。

形態素解析の精度比較に加えて、アノテーション支援実験を行い、本研究による文境界付与が実際に人手で文境界を修正する際にも時間の短縮につながることを示した。

また、近代の文献は「国立国会図書館デジタルコレクション」*7で多数の電子化画像が開かれている一方、近代文献の OCR の研究が進められており*8 [17], [18], [24], 本研究は将来的に幅広い文献データの処理に適用できる可能性

がある。

謝辞 本研究は国立国語研究所の共同研究プロジェクト「通時コーパスの構築と日本語史研究の新展開」の研究成果の一部を報告したものである。

参考文献

- [1] Cho, K., van Merriënboer, B., Bahdanau, D. and Bengio, Y.: On the Properties of Neural Machine Translation: Encoder-Decoder Approaches, *Proc. SSST*, pp.103–111 (2014).
- [2] Kudo, T. and Matsumoto, Y.: Chunking with Support Vector Machines, *Proc. NAACL*, pp.192–199 (2001).
- [3] Kudo, T., Yamamoto, K. and Matsumoto, Y.: Applying Conditional Random Fields to Japanese Morphological Analysis, *Proc. EMNLP*, pp.230–237 (2004).
- [4] Lafferty, J.D., McCallum, A. and Pereira, F.C.N.: Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data, *Proc. ICML*, pp.282–289 (2001).
- [5] Qi, P., Dozat, T., Zhang, Y. and Manning, C.D.: Universal Dependency Parsing from Scratch, *Proc. CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pp.160–170 (2018).
- [6] Read, J., Drīdan, R., Oepen, S. and Solberg, L.J.S.: Sentence Boundary Detection: A Long Solved Problem?, *Proc. COLING*, pp.985–994 (2012).
- [7] Rudrapal, D., Jamatia, A., Chakma, K., Das, A. and Gamback, B.: Sentence Boundary Detection for Social Media Text, *Proc. ICON*, pp.254–260 (online), DOI: 10.13140/RG.2.1.4481.8002 (2015).
- [8] Straka, M., Hajič, J. and Straková, J.: UDPipe: Trainable Pipeline for Processing CoNLL-U Files Performing Tokenization, Morphological Analysis, POS Tagging and Parsing, *Proc. LREC*, pp.4290–4297 (2016).
- [9] Zeman, D., Hajič, J., Popel, M., Potthast, M., Straka, M., Ginter, F., Nivre, J. and Petrov, S.: CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies, *Proc. CoNLL*, pp.1–21 (2018).
- [10] 岡 照晃, 小町 守, 小木曾智信, 松本裕治: 統計的機械学習を用いた歴史的資料への濁点付与の自動化, 情報処理学会論文誌, Vol.54, No.4, pp.1641–1654 (2013).
- [11] 服部紀子, 間淵洋子, 近藤明日子, 小木曾智信: 国定教科書のコーパス構築と公開, 日本語学会 2018 年度秋季大会予稿集, pp.582–585 (2018).
- [12] 福岡健太, 松本裕治: Support Vector Machines を用いた日本語書き言葉の文境界推定, 言語処理学会年次大会発表論文集, pp.1221–1224 (2005).

*7 <http://dl.ndl.go.jp/>

*8 <http://codh.rois.ac.jp/collaboration/>

- [13] 難波悟史, 門内健太, 但馬康宏, 菊井玄一郎: マイクロブログに対する文境界推定および係り受け解析, 言語処理学会年次大会発表論文集, pp.107-111 (2015).
- [14] 国立国語研究所: 『太陽コーパス—雑誌「太陽」日本語データベース—』, 博文館新社 (2005).
- [15] 国立国語研究所コーパス開発センター: 『現代日本語書き言葉均衡コーパス』利用の手引第 1.1 版, 国立国語研究所コーパス開発センター (2011).
- [16] 国立国語研究所コーパス開発センター (近藤明日子): 近代文語 UniDic 短単位規定集 Ver.1.1 (2016).
- [17] 美馬秀樹, 丹治 信, 増田勝也, 太田 晋: 近代文献のデジタルアーカイブ化とテキストマイニング—岩波書店「思想」を題材に, 情報処理学会研究報告人文科学とコンピュータ (CH), Vol.2012, No.4, pp.1-8 (2012).
- [18] 増田勝也: 言語情報と字形情報を用いた近代書籍に対する OCR 誤り訂正, じんもんこん 2016 論文集, Vol.2016, pp.57-62 (2016).
- [19] 南雲千香子, 近藤明日子: 『東洋学芸雑誌』コーパスの構築, 通時コーパス活用班合同研究集会 (2017).
- [20] 小木曾智信, 小町 守, 松本裕治: 歴史的日本語資料を対象とした形態素解析, 自然言語処理, Vol.20, No.5, pp.727-748 (2013).
- [21] 古井貞熙, 前川喜久雄, 井佐原均: 科学技術振興調整開放的融合研究推進精度—大規模コーパスに基づく『話し言葉工学』の構築, 日本音響学会誌, Vol.56, No.1, pp.752-755 (2000).
- [22] 田中牧郎: 『近代女性雑誌コーパス』の概要, 『日本学術振興会科学研究費補助金研究成果報告書基盤研究 (B) 「20 世紀初期総合雑誌コーパス」の構築による確立期現代語の高精度な記述』, pp.55-62 (2006).
- [23] 近藤明日子: 『明六雑誌コーパス』『国民之友コーパス』の構築, 日本語の研究, Vol.12, No.4, pp.167-174 (2016).
- [24] 永野雄大, 幡谷龍一郎, 増田勝也, 持橋大地: CNN を用いた近代文献画像からのテキスト領域抽出, 情報処理学会研究報告コンピュータビジョンとイメージメディア (CVIM), No.12, pp.1-6 (2018).
- [25] 間淵洋子: 明治・大正期『読売新聞』コーパスの構築と課題, 言語処理学会第 24 回年次大会発表論文集, pp.500-503 (2018).
- [26] 下岡和也, 内元清貴, 河原達也, 井佐原均: 日本語話し言葉の係り受け解析と文境界推定の相互作用による高精度化, 自然言語処理, Vol.12, No.3, pp.3-17 (2005).



白井 良介

2017 年明治大学文学部文学科卒業.
2019 年首都大学東京大学院システムデザイン研究科情報通信システム学域博士前期課程修了.



松村 雪桜

2017 年首都大学東京システムデザイン学部システムデザイン学科情報通信システムコース卒業. 2019 年同大学院システムデザイン研究科情報通信システム学域博士前期課程修了.



小木曾 智信 (正会員)

1995 年東京大学文学部日本語日本文学 (国語学) 専修課程卒業. 1997 年東京大学大学院人文社会系研究科日本文学文化研究専攻修士課程修了. 2001 年同博士課程中途退学. 2014 年奈良先端科学技術大学院大学情報科学研究科博士課程修了. 博士 (工学). 2001 年明海大学講師, 2006 年独立行政法人国立国語研究所研究員を経て, 2009 年人間文化研究機構国立国語研究所准教授, 2016 年より教授. 専門は日本語学, 自然言語処理. 言語処理学会, 日本語学会各会員.



小町 守 (正会員)

2005 年東京大学教養学部基礎科学科学史・科学哲学分科卒業. 2007 年奈良先端科学技術大学院大学情報科学研究科博士前期課程修了. 2008 年日本学術振興会特別研究員 (DC2) を経て, 2010 年博士後期課程修了. 博士 (工学). 同年同研究科助教を経て, 2013 年より首都大学東京システムデザイン学部准教授. 大規模なコーパスを用いた意味解析および統計的自然言語処理に関心がある. 人工知能学会, 言語処理学会, ACL 各会員.