

語彙的結束性のない文章と新出単語の主語に着目した 文章改善点指摘システムの検討

正村 真悟¹ 児玉 哲也² 大西 諒太¹ 澤野 弘明¹

概要：論文執筆のためのソフトウェアは、誤字脱字や句読点の統一など、単文もしくは文章全体のパターンを評価している。これにより文章作成の基本的なルールが守られるが、文章の本質的な可読性つまり、理解のしやすさに影響する文章を指摘することができない。そこで本研究では、複数文のつながりにおける可読性を向上させて、理解しやすい文章を作成するための改善点を指摘するシステムを提案する。具体的には、前後2文で比較した語彙的結束性がない文の指摘と、新出単語の主語の指摘による文章改善を支援する。提案システムのユーザビリティ評価では平均点より高く使いやすいことが示された。

1. はじめに

学術論文の執筆において、正確で、かつわかりやすい文章の作成が必須技術として求められる。学術論文の執筆経験がない学生は、誤字・脱字、送りがなの間違いをなくすといった、基本的な文章作成の技術から習得する。文章作成に関する教材 [1], [2], [3] も 30 年以上も前からいくつか発刊されているが、根本的な指摘箇所は変わっていない。すなわち、学生を指導する立場である教官は、毎年入れかわる学生に対して同様の内容を教授しているといえる。

学生が教材を利用して、基本的な文章作成技術を身につけることが望ましいが、執筆した文章に対して、学生自身が客観的な視点で修正箇所を発見する技術を獲得するには時間が掛かる。また、指導教官による個別の添削指導は、学生に「気付き」を与えるためには効果的であるが、大人数に対して対応する場合は負荷が大きくなる。さらに複数の指導教官が添削する場合、執筆及び添削経験の違いにより、指導教官ごとに添削内容が異なる場合もある。そこで本研究では、自然言語処理による文章指摘手法に着目する。

学術論文の執筆に利用される、Microsoft Word のような文書作成ソフトウェアには、文法誤りや誤字脱字、英単語のスペルを指摘する校正機能が搭載されている。この校正機能では単文の指摘に留まるが、即座に利用できるようにユーザビリティ性は高い。また、複数の文章における表記揺れ、例えば「サーバ」、「サーバー」のような違いを指摘

できるシステム、RedPen [4] が伊藤によって提案されている。この RedPen では、その他に文の長さや二重否定に関する指摘も可能である。また坂本らは、意味を有する最小の言語単位である形態素に着目して、学術論文で指摘される形態素が現れた場合に、文の誤りを指摘する手法 [5] を提案している。坂本らの手法では、辞書のように指摘項目が登録されているために、登録していない形態素に関しては対応できないという課題がある。ただし、これらの手法による指摘箇所を学生自身で追加・修正すれば、文書作成の基本的なルールが守られた文章となるが、文章の本質的な読みやすさ、すなわち可読性自体は向上できない。

文の修飾関係における係り受けに着目した、可読性を評価する手法が祖ら [6] や松本ら [7] によって提案されている。どちらの手法も係り受け関係を木構造で表現して、数値的もしくは画像で文章の可読性を提示している。これら手法では執筆者に指摘箇所を提示できないが、可読性の可視化により「気付き」を与えるきっかけを提供できる。

可読性を向上させるために、上記のようなルールベースとは異なるアプローチとして、機械学習やビッグデータによって改善する手法も提案されている。日本語文章の大量の言語資料をニューラルネットワークに学習させて、不自然な箇所を検知する手法 [8] を鈴井らが提案している。また、英語の学術論文において、使用される頻度が高い何語かのかたまり（語連鎖）を推薦する手法 [9] が Mizumoto らによって提案されている。これらの手法では入力されるデータに依存するために、高度な説明が要求される学術論文には使用できない状況が少なからず存在する。また、これまでに述べた関連研究では、一文の指摘に留まっており、連続する文（文章）の可読性に関する指摘に関しては筆者

¹ 愛知工業大学
Aichi Institute of Technology

² 奈良先端科学技術大学院大学
Nara Institute of Science and Technology

らが調査した限り存在していない。そこで本研究では、文章の可読性を評価して指摘する手法を検討する。特に、連続する二文の可読性を評価する「結束性」[10]に着目する。また、新出単語の名詞が主語である場合に事前知識のない読者にとっては可読性が低いため、新出単語の主語についても指摘する。第2節に結束性と新出単語の主語に着目した、文章改善点指摘項目を判定する手法を述べ、第3節では提案手法を実装したシステムについて示す。第4節で実験と考察を述べた後、最後に本論文をまとめる。

2. 文章改善点指摘項目の判定手法

2.1 結束性

文章の可読性の指標の一つに、前後二文のつながりを示す結束性と呼ばれる評価指標が用いられる。結束性のある二文が多い文章は可読性が高く、結束性のある二文が少ない文章は可読性が低いといわれている。結束性には文法的結束性と語彙的結束性の2種類があり、まず文法的結束性について説明する。文法的結束性とは、前後2文が人称代名詞や指示語などによって、意味的につながる性質である。文法的結束性のある2文の例を図1に示す。図1の2文目の代名詞「これ」は、1文目の名詞「結果」を指しており、文法的に繋がりがある結束性を示している。

つぎに語彙的結束性について述べる。語彙的結束性とは、文と文が共通単語や同義語を持ち、意味的につながる性質である。語彙的結束性のある2文の例を図2に示す。図2の2文目にある名詞「手法A」は1文目にも登場しており、図2の1文目と2文目には、名詞「手法A」という共通単語を含んでいる。この共通単語による繋がりを持つ2文には語彙的結束性が存在する。このような語彙的結束性の判定であれば、前後2文が持つ単語の比較により機械的に判定できる。以降、本稿における結束性は、共通単語について着目した語彙的結束性を指す。

2.2 語彙的結束性判定のアルゴリズム

本節では語彙的結束性を判定するアルゴリズムを提案する。提案するアルゴリズムを以下に示す。

- (1) 複数の文章のテキストを入力する。

1 文目: 手法 A では B のような結果になった。
2 文目: **これ**は、C だからと考えられる。

図 1 文法的結束性のある 2 文の例

1 文目: 本論文では**手法 A** について提案する。
2 文目: **手法 A** では手法 B を使用する。

図 2 語彙的結束性のある 2 文の例

- (2) 句点に注目してテキストを 1 文ずつに分割する。
- (3) 分割した 1 文を形態素に分割する。
- (4) 得られた形態素のうち、名詞と動詞の単語を抽出する。
このとき、動詞を終止形に変形する。
- (5) 前後 2 文で名詞及び動詞の単語を比較する。比較した 2 文に共通する単語が存在している場合、語彙的結束性ありと判定する。

2.3 新出単語の主語の抽出手法

本節では新出単語の主語に着目した文章改善点指摘手法について示す。新出単語の主語が登場する場合、事前知識のない読み手にはその主語の名詞を理解することが難しく、可読性が下がる場合がある。新出単語の主語が登場する複数文章の例を図3に示す。ここで、図3の冒頭2文は、図2の例の2文を示している。図3の3文目の主語には名詞「手法C」が使用されている。名詞「手法C」は1文目及び2文目に出現していないため、「手法C」について知る術がなく、可読性が低い文章といえる。このような新出単語の主語を抽出するアルゴリズムを以下に示す。

- (1) 複数の文章のテキストが入力される。
- (2) 句点に注目してテキストを 1 文ずつに分割する。
- (3) 分割した 1 文を形態素に分割する。
- (4) 名詞とその名詞の直後の助詞に「は」「が」「も」が使用されている場合、主語と定義する。
- (5) 主語がこれまでに文章中に登場していない名詞の場合、新出単語の主語と抽出される。

3. 文章改善点指摘システム

3.1 Dropbox を利用した文章改善点指摘システム

前節までに提案した文章改善点の指摘項目を、自動的に判定するシステムを実装する。さて、文章校正機能には、即時性や簡潔性が求められる。そこで提案システムにおいても、利用者にとって簡潔に動作する仕組みを目指す。

提案システムには、複数のファイル間で同期可能なオンラインクラウドストレージ Dropbox^{*1}を利用する。Dropbox は、利用者がクラウド上にファイルをアップロードすると、Webhook という機能で外部 Web サービスに対して HTTP の POST パラメータを送信することが可能である。提案システムでは、この Webhook を利用して、利用者が添削用のファイルをアップロードするだけで、ファイルに対し

1 文目: 本論文では手法 A について提案する。
2 文目: 手法 A では手法 B を使用する。
3 文目: **手法 C** は正村らによって開発された。

図 3 新出単語の主語のある文章の例

^{*1} <http://dropbox.com>

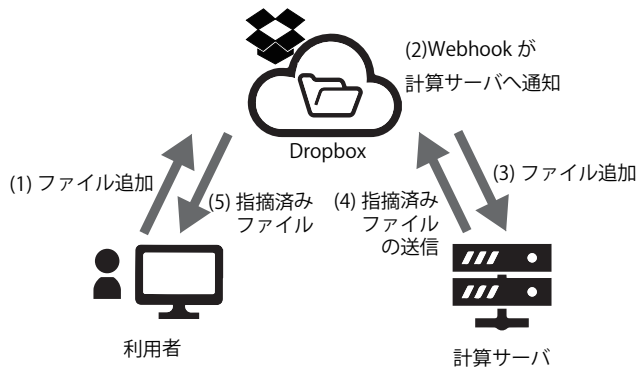


図 4 文章改善点指摘システム

て自動的に文章改善点の指摘項目を追記するという仕組みを持つ。提案システム（図 4）の流れを以下に示す。

- (1) 利用者が Dropbox のフォルダにテキストファイルを追加する。
- (2) Webhook が計算サーバへ通知する。
- (3) 利用者が追加したファイルに対して計算サーバで指摘項目を判定して、追記する。
- (4) 計算サーバが指摘済みファイルを Dropbox へ送信する。
- (5) 指摘済みファイルを利用者に提示する。

4. 文章改善点指摘システムの実験・考察

4.1 実験環境

提案システムの開発言語には、PHP 7.1.23, Python 3.7.4 を使用した。入力形式は学術論文の執筆に利用される \LaTeX ファイルである。形態素解析には、mecab [11] を用いた。 \LaTeX 形式の入力テキストの例を図 5 に、提案システムによる出力結果を図 6 に示す。図 6 に示すように指摘システムは利用者が執筆した \LaTeX ファイルに影響を与えないようするため、文末にコメント化した指摘項目を入力ファイルに追記する。また、指摘項目の最後には、検査した文章数・検証数・結束率といった、可読性の指標を記載する。

4.2 ユーザビリティ評価

論文執筆中の A 大学の学生 6 人の実験協力者に対して、提案システムのユーザビリティ性を評価する実験を行った。提案システムには第 2.2 節に示した語彙的結束性判定のアルゴリズムのみを実装した。評価手法には定量的にユーザビリティ性を評価できる SUS (System Usability Scale) [12] を用いた。SUS は、被験者が 10 項目の質問に対して、「とてもそう思う」(4 点) から「まったくそう思わない」(0 点) の 5 段階で回答する。10 項目の質問には、肯定的な項目 (奇数番号) と否定的な項目 (偶数番号) の質問がそれぞれ 5 問ずつ用意され、否定的な項目 (偶数番号) の場合は、4 から回答の得点を引き、0 から 4 点で評価される。SUS における質問 10 項目を表 1 に示す。

```
\section{ はじめに }
```

工場生産における外観検査工程を自動化するために画像処理技術が用いられている。人間の学習過程を模倣したディープラーニングが製造現場で導入されているが、学習時に大量の教師データが必要であり、教師データが少ない場合に十分な効果は発揮できない \cite{bib:sakakibara}。また生産ラインでは高速な検査が求められるため、画像処理アルゴリズムのハードチップ化が難しいディープラーニングをそのまま導入することは難しい。また、画像処理アルゴリズムを進化的アルゴリズムにより自動最適化する手法 \cite{bib:muroi} も室井らによって提案されているが、進化的アルゴリズムの結果は、準最適解にとどまるという課題がある。しかし、最適解を求めるためには、自動最適化による計算コストがかかる。そこで、分散処理によって解決する手法を提案する。本項では、その前段階である、画像処理アルゴリズムの自動生成とその実験と結果について述べる。

図 5 \LaTeX 形式の入力テキストの例 (7 文)

```
%-----
%「はじめに」節の 5 分目と 6 文目に結束性がありません。
%5 文目:しかし、最適解を求めるためには、自動最適化による
計算コストがかかる。
%6 文目:そこで、分散処理によって解決する手法を提案する。
%-----
%文章数:7 検証数:6
%
%結束率は 83%です。
```

図 6 提案システムの出力結果

10 項目ある質問の合計得点 (40 点満点) を 2.5 倍して 100 点満点の結果と、SUS の標準平均である 68 点を比較する。標準平均を上回るシステムは基準よりもユーザビリティ性の高いシステムとされる。具体的には、80.3 点以上で A ランク、68 点より高く 80.3 点未満であれば B ランク、68 点以上で C ランク、51 点より高く 68 点未満であれば D ランク、51 点未満であれば E ランクと判定される。

SUS によるユーザビリティ評価の結果を表 2 に示す。SUS 評価の結果、平均点が 75.4 点となり標準平均である 68 点を上回り B ランクとなった。この結果から提案システムは使いやすいシステムであることが示された。

4.3 今後の展望

本論文で提案した、新出単語の主語を抽出する仕組みには、名詞及びその名詞直後の助詞「は」「が」「も」を利用している。日本語の科学技術論文の場合、主語が省略される場合も存在するため、全文に対して有効な手法とはいえない。また、関連研究を示す文章の新出単語の主語に人名が利用される場合もあり、慣習的に問題がない場合もある。そのため、今後の展望としては、章・節のタイトルや段落

表 1 SUS 評価のための質問 10 項目

| 番号 | アンケート内容 |
|-----|----------------------------|
| 1 | しばしば利用したいと思いましたか？ |
| 2* | 説明が必要なほど複雑だと感じましたか？ |
| 3 | 容易に使いこなすことができると感じましたか？ |
| 4* | 利用するのに専門的な人の説明が必要だと感じましたか？ |
| 5 | 内容に統一性が十分にあると感じましたか？ |
| 6* | 内容に一貫性のないところが多くあったと感じました。 |
| 7 | 大体の人は、利用方法をすぐに理解すると思いました。 |
| 8* | とても操作しづらいと感じましたか？ |
| 9 | 利用できる自信はありますか？ |
| 10* | 利用前に知っておくべきことが多くあると思いました。 |

*: 否定的な項目

ごとの文章群で評価する手法を検討する予定である。

5. おわりに

本稿では、可読性向上のために語彙的結束性のない 2 文および、新出単語の主語を指摘する手法を提案した。提案手法を自動的に判定するシステムを、オンラインストレージサービス Dropbox を利用して実装した。Dropbox を利用することで即時的に文書改善点を指摘することが可能になった。また、語彙的結束性判定手法のみを実装した提案システムを構築して、SUS によるユーザビリティ評価を実施した。評価実験の結果、標準平均を上回り、提案システムの基準が B ランクとなり、利用者にとって使いやすいシステムであることが示された。今後の課題として、主語判定方法の改善や、章・節タイトル及び段落ごとの文章群で文を評価する手法の検討が挙げられる。

参考文献

- [1] 中島利勝, 塚本真也: “知的な科学・技術文章の書き方 - 実験レポート作成から学術論文構築まで-”, 株式会社コロナ社 (1996)
- [2] 阿部圭一: “明文術 - 伝わる日本語の書き方”, NTT 出版株式会社 (2006)
- [3] 塚本真也: “知的な科学・技術文章の徹底演習”, コロナ社 (2007)
- [4] 伊藤敬彦: “自動文書検査ツール RedPen”, 信学技報, Vol. 114, No. 211, pp. 69-74 (2014)
- [5] 坂本俊介, 須藤崇志, 丸山広, 中村太: “形態素解析を利用し

- た文章校正手法の提案”, 情処研報, Vol. 2009-DD-72, No. 17, pp. 1-6 (2009)
- [6] 祖国威, 加納敏行: “構文的な分かりやすさを評価する可読性評価技術”, 言語処理学会第 16 回年次大会発表論文集, pp. 1082-1085 (2010)
- [7] 松本章代: “科学的文章の推敲・校正を支援する教育システムの構築, 東北学院大学教養学部論集, Vol. 167, pp. 53-62 (2014)
- [8] 鈴木克徳, 若林啓: “ニューラルネットワークを用いた日本語学習者の文章における不自然箇所検知”, 第 10 回データ工学と情報マネジメントに関するフォーラム, 5 pages (2018)
- [9] A. Mizumoto, S. Hamatani, Y. Imao: “Applying the Bundle-Move Connection Approach to the Development of an Online Writing Support Tool for Research Articles, *Language Learning*, (2017)
- [10] 池上嘉彦: “テキストとテキストの構造”, 国立国語研究所, 談話の研究と教育 1, pp. 7-42 (1983)
- [11] T. Kudo, K. Yamamoto & Y. Matsumoto: “Applying Conditional Random Fields to Japanese Morphological Analysis”, *Proc. of the 2004 Conf. on Empirical Methods in Natural Language Processing (EMNLP-2004)*, pp. 230-237 (2004)
- [12] J. Break: “SUS: A Retrospective”, *J. Usability Studies*, Vol. 8, Issue 2, pp. 29-40 (2013)

表 2 A 大学の評価結果

| 被験者 | アンケート番号 (得点) | | | | | | | | | | 合計点 |
|-----|--------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|------|
| | 1 | 2* | 3 | 4* | 5 | 6* | 7 | 8* | 9 | 10* | |
| A | 4 | 1 | 4 | 3 | 3 | 2 | 3 | 2 | 4 | 1 | 67.5 |
| B | 3 | 2 | 3 | 3 | 3 | 2 | 3 | 3 | 3 | 2 | 67.5 |
| C | 3 | 4 | 4 | 3 | 2 | 2 | 4 | 4 | 4 | 3 | 82.5 |
| D | 4 | 3 | 4 | 4 | 3 | 3 | 4 | 4 | 3 | 3 | 87.5 |
| E | 4 | 2 | 3 | 2 | 4 | 4 | 2 | 2 | 2 | 2 | 67.5 |
| F | 4 | 3 | 2 | 3 | 2 | 3 | 4 | 4 | 4 | 3 | 80.0 |
| 平均 | 3.7 | 2.5 | 3.3 | 3.0 | 2.8 | 2.7 | 3.3 | 3.2 | 3.3 | 2.3 | 75.4 |

*: 否定的な項目