

# 語彙の出現位置と頻度による文体類似度を用いた 文章の執筆者数推定

渡邊 充博<sup>†</sup> Eint Sandi Aung<sup>‡</sup> 山名 早人<sup>§</sup>

<sup>†</sup> 早稲田大学大学院基幹理工学研究科 〒169-8555 東京都新宿区大久保3-4-1

<sup>‡</sup> 早稲田大学理工学術院 〒169-8555 東京都新宿区大久保3-4-1

<sup>§</sup> 国立情報学研究所 〒101-8430 東京都千代田区一ツ橋2-1-2

E-mail: <sup>†</sup><sup>‡</sup> {mwatanabe, esdaung, yamana}@yama.info.waseda.ac.jp

**あらまし** ブログやSNSの普及によって個人の情報発信の機会は増加した一方で、フェイクニュースと呼ばれる虚偽の情報を含んだニュースが発信、拡散されることが問題となっている。情報内容の真偽の自動的な判定を行うにあたって、著者の先行研究では、「信頼性の高い文章ほど編集に関わった人数が多い傾向にある」という点に着目した。文章の執筆者数を信頼性の測定指標とするため、文体の類似度を用いて文章中の文体変化を検出し執筆者数を推定する手法を提案した。本研究では、執筆者数の推定精度向上のために、語彙出現位置に応じた重みを付与して、類似度算出手法の改良を行い、複数人による記述を含む文章に対して文体の変化点と執筆者数を推定した。2人によって記述された文章における執筆者数の正解率は81.8%、MAEは0.183人となり、執筆者数が正解し、かつ文体変化点を正しく推定できた確率は65.5%となった。

**キーワード** 情報信頼性, 文体分析

## 1. はじめに

ブログやSNSの普及に伴って、多様な情報をリアルタイムかつ手軽に収集できるようになり、情報の発信や共有も容易になった。SNSにおける共有機能は、有用な情報の伝達に役立つ一方、情報の正しさが確認されないまま、誤った情報を拡散させる場合がある。フェイクニュースと呼ばれる虚偽の情報を含んだニュースは、事実であるニュースより短時間で広範囲に拡散される [1]。誤情報の拡散を防ぐためには、ユーザ自身が情報内容の真偽を見極める必要がある。しかし、内容の真偽判定には関連する知識や手間が必要となり、情報が誤りであると周知されるまでにも時間が掛かる。例えば、Twitterではフェイクニュースが拡散された後に、内容が事実と反すると共有・拡散されるまでには10～20時間程度を要する [2]。したがって、短時間で自動的に実行できる信頼性判定システムが求められている。信頼性の自動判定には、ページ内のコンテンツに基づく特徴量やページ同士のネットワークに基づく特徴量が用いられる。本稿では特徴量として文章の執筆者数に着目した。

新聞社や出版社では、原稿上に不適切な表現がないか、記述内容に整合性があるかどうかを校閲者が確認する。Web上の文章でも、Wikipedia<sup>1</sup>のように共同でページを編集できる場合、記述の誤りや不足を修正して文章の品質を向上させることができる。Wikipediaでは「秀逸な記事」と呼ばれる高品質な記事がユーザ

による推薦と投票から選出されている。秀逸な記事は、一般の記事と比較して編集回数と編集者数が多い傾向にある [3]。つまり、文章の執筆や編集に関わった人数が内容の信頼性や品質を反映しているといえる。しかしながら、編集の履歴を保持するWikipediaと異なり一般の文章では何人によって執筆されたかは明らかでない。よって執筆者数を信頼性推定の評価指標として用いるためには、執筆者数の推定が必要となる。

塩浦ら [4]は、執筆者数と文章の信頼性との関係に着目し、執筆者数の推定を行った。[4]では、文章を分割し文体を基準にクラスタリングを行って人数を推定した。クラスタリングを用いた [4]の手法に対して筆者は、文章内の文体の変化を定量的に扱うために、文体の類似度を用いて文体の変化を検出し執筆者数を推定する手法を提案した [5]。本稿では推定精度向上のため、語彙の出現位置に応じた重みを付与するよう類似度算出手法を改良した。

本稿は以下の構成をとる。2節で情報の信頼性および執筆者数推定を扱う関連研究を述べ、3節で語彙出現位置に応じた重みを付与した文体の類似度から文体の変化を検出し執筆者数を推定する手法を提案する。4節で提案手法の評価実験の結果および考察を述べる。最後に5節でまとめを述べる。

<sup>1</sup> <https://en.wikipedia.org/>

## 2. 関連研究

### 2.1. 情報の信頼性と品質

情報の信頼性は人間が知覚する情報の品質のことである [6]. 信頼性は複数の要素から複合的に決定付けられる. Foggら [6]は情報の信頼性に関する従来の研究を統合し, 信頼性を仮定的信頼性, 表面的信頼性, 評判の信頼性, 経験的信頼性の4種類に整理した.

Foggら [7]は, ユーザがWebサイトの信頼性を評価する際に注目する要素を調査するために, 100件のWebサイトに対するユーザの合計2,440件のコメントを収集した. これらのうち46.1%のコメントではデザインに着目して評価が行われていたことから, ユーザは表面的信頼性を重視する傾向にあることが明らかになった.

情報の内容は信頼性があるかという観点と, 事実であるかという観点から分類することができる [6] [8]. 信頼できる内容でも, セールストークや巧妙な詐欺のように事実でない場合もある. つまり「情報の正しさや有用性を示す品質」と「知覚される信頼性の高さ」は必ずしも一致しない. しかしながら「高品質の情報」は信頼性も高いことが示されている.

### 2.2. 文章の編集と品質

情報の品質を測定する方法には, 人手に基づく方法と評価指標を用いて自動的に行う方法がある.

Giles [9]は, Wikipediaの記事の正確さを調査した. 科学の幅広い分野に関連する42のテーマについて, Wikipediaの記事とブリタニカ百科事典の項目を専門家が検証した. 1記事あたりWikipediaでは平均3.9件, ブリタニカ百科事典では平均2.9件の誤りや不適当な記述が指摘され, 両者の正確さには差が見られないと結論付けた.

Chesney [10]は, Wikipediaの記事に対して55人の専門家・研究者による品質評価を行った. 与えられた30件の記事について, 自身の専門分野に関連する記事を読んだ専門家のユーザは, 専門分野でないユーザより, 記事が信頼できると評価した. また, 記事の内容に誤りがあったのは13%であり, 完全に信頼できるわけではなくとも, Wikipediaの記事の精度が高いと示した.

Wikipediaでは, ユーザによる推薦と投票から「秀逸な記事」と呼ばれる高品質な記事が選出されている. 品質評価のための指標評価として秀逸な記事と一般の記事の区分を用い, 指標評価が行われている. 品質評価において, 編集の特徴を利用したものとしては, 記事の差し戻し回数 [11]や安定性 [12]といった編集頻

度のほか, 質の高い記述は削除されにくいことから文章の残存率 [13]が評価指標として提案されている.

Wilkinsonら [3]は, 秀逸な記事1,211件を含むWikipediaの150万件の記事の編集回数を分析して, 記事の扱うトピックの人気や注目されやすさと記事の品質との関係を調査した. 記事の人気度として, ページの閲覧回数と相関のあるPageRankを用いた. 人気度と閲覧回数の関係を分析すると, 人気度の高い記事ほど編集回数と編集者数が多い傾向にあることが明らかになった. なお, 編集回数と編集者数は記事の存在期間で正規化した値を用いた. また, 同じ値のPageRankを持つ記事群内で比較すると, 秀逸な記事は一般の記事より有意に多い編集回数と編集者数となっていることが明らかになった.

### 2.3. 品質評価のための執筆者数推定

文章の品質評価を目的として執筆や編集に関わった人数を推定する取り組みには, 塩浦ら [4]による手法がある. [4]では, まず文章をスライディングウィンドウで分割し, 各ウィンドウ内の単語に対して品詞n-gramの頻度ベクトルを求める. 次に頻度ベクトルをx-meansを用いてクラスタリングし, 得られたクラスタの数が推定された執筆者数となる. また, 頻度ベクトルの重み付けには $\log(3n)$ を用いると最も高いクラスタリング精度を示すと報告されている.

クラスタリングを用いた [4]の手法では文章内におけるウィンドウの連続性が反映されないこと, 文章内の文体の変化を定量的に扱えないことから, 筆者は品詞n-gramの出現頻度に基づく類似度を用いる手法を提案した [5]. 文章をスライディングウィンドウで分割し, ウィンドウの半分より前と後の部分における品詞n-gramの出現頻度を求め特徴量とした. 前半部と後半部の特徴量ベクトル間の類似度が変化する位置を文体変化が発生した位置として検出し, 執筆者数を推定した. [5]では, 2人の異なる執筆者による記述を含む文章に対して, ウィンドウ長を400単語とした場合, 執筆者数を正しく推定できた確率として67.1%, 推定執筆者数の平均絶対誤差0.328人, 文体変化点と執筆者数を同時に正しく推定できた確率として46.7%という結果が得られた.

## 3. 語彙の出現位置と頻度による文体類似度を用いた文章の執筆者数推定手法の提案

本稿では, 語彙出現位置に応じた重みを付与して類似度算出を行うように筆者の既存手法 [5]を改良して, 執筆者数を推定する手法を提案する. 提案手法では, 文章中で執筆者が変わったとする位置を仮定し, その

位置で文章を分割する．当該分割点の前後近傍間の文体変化を求め，大きな文体変化が認められる点を執筆者が変わった点として出力する．具体的には，分割点の前後近傍それぞれにおいて，出現する単語の品詞n-gramを求め，分割点からの距離に応じて重み付けを行う．

提案手法は次の6つの処理から構成される．

1. 文章中の単語を品詞名に置き換え品詞列を取得する．
2. 文章全体の品詞列を前後2つの部分品詞列に分割する．
3. 分割した前半部と後半部の文体の類似度を算出する．
4. 分割位置を変化させて3を実行し，類似度が極小となる点を抽出し，文体変化点の候補とする．
5. 設定した閾値に基づいて文体変化点を絞り込む．
6. 文体変化点の数から執筆者数を決定する．

### 3.1. 品詞n-gramの取得

手法の適用結果が文章中の単語やトピックに依存しないよう，単語の品詞情報を利用する．執筆者数の推定を行う文章に対して形態素解析を適用し，単語の品詞情報を取得する．形態素解析器にはMeCab，辞書にはNEologdを用いる．MeCabが出力する品詞情報の2階層目までを取得する．取得される品詞情報の例を表3.1に示す．助詞と助動詞については，品詞名に置き換えずに原型を用いる．

表3.1 取得される品詞情報の例

原型	1階層目	2階層目
雨	名詞	一般
が	助詞	格助詞
降る	動詞	自立
.	記号	句点
,	記号	読点

得られた品詞列を図3.1のように2つに分割して前半部と後半部の部分品詞列に分解する．前半部と後半部でそれぞれ品詞n-gramを取って，含まれる語彙の出現頻度ベクトルを求める．nには1から3を用いる．なお，本モデルは，前半部と後半部で2名の執筆者を前提としている点ではない点に注意する必要がある．すなわち，分割位置に近い部分品詞列に高い重みを与えることによって，分割位置前後での執筆者の違いを判定することを目指しており，分割位置をずらしていくことにより，複数名によって記述された文章での執筆者数推定が可能である．

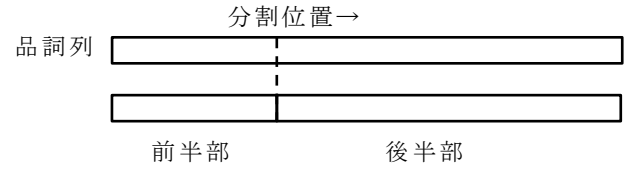


図3.1 品詞列の分割

### 3.2. 語彙の出現頻度算出

文章中の品詞のn-gramの全体の単語数がm個である文章の品詞列 $p = \{p_1, p_2, \dots, p_m\}$ が与えられたとき，品詞 $p_s$ の直後の位置sで分割すると，前半部 $p_a = \{p_1, p_2, \dots, p_s\}$ と後半部 $p_b = \{p_{s+1}, p_{s+2}, \dots, p_m\}$ に分割される．前（後）半部においてk回出現する語彙 $v = \{v_1, v_2, \dots, v_k | v_i \in \{\text{前(後)半部の品詞n-gram}, 1 \leq n \leq 3\}\}$ について， $v_i$ の出現位置 $x(v_i)$ と分割位置sとの距離を重みとして，頻度を式3.1のように求める．重みには $\lambda$ を定数として距離 $|s - x(v_i)|$ が大きくなるにつれ減衰する $\exp(-\lambda|s - x(v_i)|)$ を用いる．なお， $|s - x(v_i)|$ は，前半部においては， $|x(p_s) - x(v_i)|$ ，後半部においては $|x(p_{s+1}) - x(v_i)|$ として計算するものとする．

$$\begin{aligned} \text{position\_weighted\_freq}(v) \\ = \sum_{i=1}^k \exp(-\lambda|s - x(v_i)|) \end{aligned} \quad (3.1)$$

次に，個々の語彙に対して，対象文書中での出現頻度に応じて重要度(inverse in a document frequency)を与える．通常のidfとは異なり，与えられた一文章中に出現する稀な語彙に大きな重みを与える(式3.2)．sumを対象文章中(前半部+後半部)の延べ語彙数( $n$ -gram,  $1 \leq n \leq 3$ の延べ語彙数)， $\text{freq}(v)$ を対象文章中(前半部+後半部)の語彙vの出現回数として，式3.2のように対象文章における当該語彙vの出現割合が小さいほど大きな重みを与える．

$$\begin{aligned} \text{idf}(v) &= \log\left(\frac{\text{sum}}{\text{freq}(v)}\right) \text{ if } \text{freq}(v) \geq 1 \\ \text{idf}(v) &= 0 \text{ if } \text{freq}(v) = 0 \end{aligned} \quad (3.2)$$

最後に，文章に含まれる語彙に対してそれぞれ式3.1と式3.2を求め乗算して特徴量とする(式3.3)．1-gram, 2-gram, 3-gramの全ての特徴量を順に並べて前半部における語彙頻度ベクトルaを得る．後半部も同様に頻度ベクトルを求め，bとする．

$$\text{position\_weighted\_freq}(v) \times \text{idf}(v) \quad (3.3)$$

### 3.3. 文体変化点の推定

続いて文体変化点を推定する．前半部 $\mathbf{a}$ と後半部 $\mathbf{b}$ の品詞 $n$ -gramの頻度ベクトル間のコサイン類似度を求める(式3.4)．

$$\cosSim(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{a}| |\mathbf{b}|} = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2} \cdot \sqrt{\sum_{i=1}^n b_i^2}} \quad (3.4)$$

文章の分割位置を動かしながら前半部と後半部の頻度ベクトル同士の類似度を求めることで，文章中の類似度の変化を得る．

前半部 $\mathbf{a}$ と後半部 $\mathbf{b}$ で文体が異なるほど類似度は小さくなる．分割位置を一定間隔で動かした場合に，分割位置 $s_j$ における類似度 $sim_{s_j}$ が，1つ前の分割位置での類似度 $sim_{s_{j-1}}$ より小さく，かつ1つ後の分割位置での類似度 $sim_{s_{j+1}}$ より小さい値となった場合， $s_j$ を文体変化点の候補とする．

類似度の閾値 $th$ ，および文体変化点同士の最小間隔を $minInt$ [単語]を設定し，文体変化候補点における類似度が $th$ より小さく，かつ類似度が前後 $minInt$ [単語]の範囲内で最も小さい場合，文体変化点を残し，当てはまらない場合は削除する．図3.2では， $s_j$ と $s_{j+2}$ の位置に文体変化点の候補があると考えたとき， $s_{j+2}$ における類似度 $sim_{s_{j+2}}$ は $th$ より小さいが， $minInt$ 単語より短い距離にある $s_j$ における類似度 $sim_{s_j}$ の方が小さいため $s_{j+2}$ の点は候補から削除する．

候補の絞り込みの後，残った文体変化点の数を推定執筆者数とする．

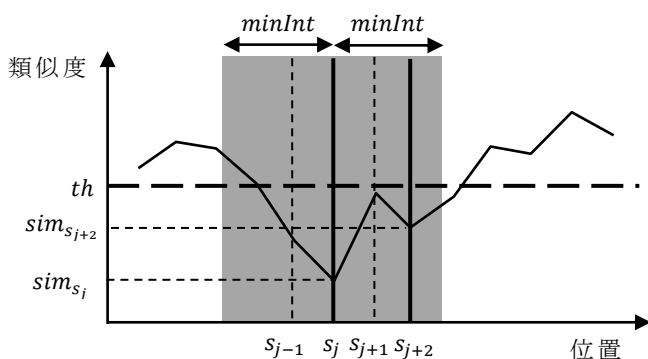


図 3.2 文体変化候補点の絞り込み

## 4. 実験と評価

### 4.1. データセット

1人のみによって記述された文章を複数連結するこ

とで，複数人によって記述された文章を生成し，データセットとした．小説の文章は1人のみによって書かれているとみなして，青空文庫<sup>2</sup>の収録作品の執筆者の中から執筆者を無作為に選び，選んだ執筆者の作品の中から無作為に作品を選出した．各作品の文章の長さを揃えるために，各作品の先頭から500単語以上を含み，最初に終わる文までを文群として抽出した．

文群から同一の執筆者による文群を連結することで1人によって記述された文章を生成した．次に異なる執筆者の文群からそれぞれ1文群を選択して連結し一つの文章とすることで2人によって記述された文章を生成した．3人以上によって記述された文章も連結する文群の数を増やすことで生成した．

### 4.2. パラメータの決定

まず，提案手法で用いるパラメータの値を決定する．対象とするパラメータは，出現位置の重みの減衰係数 $\lambda$ ，類似度の閾値 $th$ ，文体変化点同士の最小間隔 $minInt$ の3つである．1人または2人によって記述された文章を用いて文体変化点の推定を行い，最良の推定結果となるパラメータを求める．

一執筆者に対して5作品，10人の執筆者から合計50作品を用いて，1人または2人で記述された文章を生成した．50作品は平均で520.28単語，19.34文を含む．同一執筆者の異なる2作品を連結した文章(=執筆者数1人)と，異なる執筆者の2作品を連結した文章(=執筆者数2人)を各100組，合計200組用いた．

各文章に対してパラメータを変えながら提案手法を適用した．前半部と後半部の単語数に極端な差が出ないように，最初の分割位置は先頭から100単語目として，30単語ずつ分割位置を動かした．最後の分割位置は文章の終端からの距離が100単語以上となる位置とした．推定された文体変化点が実際の文体変化点から50単語以内にある場合を正解として，執筆者数が正解するという条件と文体変化点が正解するという条件を同時に満たす確率を求めた．この確率を全文章数(=200組)で平均した値が最も高くなる時のパラメータを求めた．

5分割交差検証を用いて，探索対象とする文章の範囲を変えながら，執筆者数と文体変化点が同時に正解となる確率を最大化するパラメータの組を5通り求め，各パラメータ $\lambda$ ， $th$ ， $minInt$ の平均を取って， $\lambda$ : 0.006， $th$ : 0.772， $minInt$ : 168[単語]とした．

<sup>2</sup> <https://www.aozora.gr.jp/>

4.3. 文体変化点の推定

(1) 複数人によって記述された文章

次に複数人によって記述された文章に対して、4.2項で決定したパラメータを用いて、文体変化点および執筆者数を推定した。4.2項で対象とした執筆者とは異なる執筆者10人による各10作品、合計100作品を用いた。100作品は平均で517.66単語、18.69文を持つ。これらの作品の文章を連結して生成した複数の執筆者(執筆者数2人～4人)の記述を含む文章、各400件について提案手法を適用した。

分割位置のスライド幅を30単語として、推定執筆者数と実際の執筆者数が一致した確率および、推定執筆者数と実際の執筆者数との平均絶対誤差MAEを求めた。また、推定執筆者数と実際の執筆者数が一致し、かつ文体変化の推定点が実際の変化点の前後50単語以内であった確率を求めた。結果を表4.1に示した。また、実際の執筆者数別の推定執筆者数の内訳を図4.1に示した。2人によって記述された文章では、執筆者数の正解率は81.8%、MAEは0.183人となり、執筆者数が正解し、かつ文体変化点を正しく推定できた確率は65.5%となった。

表4.1 文体変化点(誤差が前後50単語以内)と執筆者数の推定結果 (執筆者数2～4人)

実際の執筆者数 [人]	執筆者数を正しく推定できた確率	左記に加え、文体変化点を正しく認識できた確率	MAE [人]
2	0.818	0.655	0.183
3	0.748	0.415	0.253
4	0.650	0.253	0.373

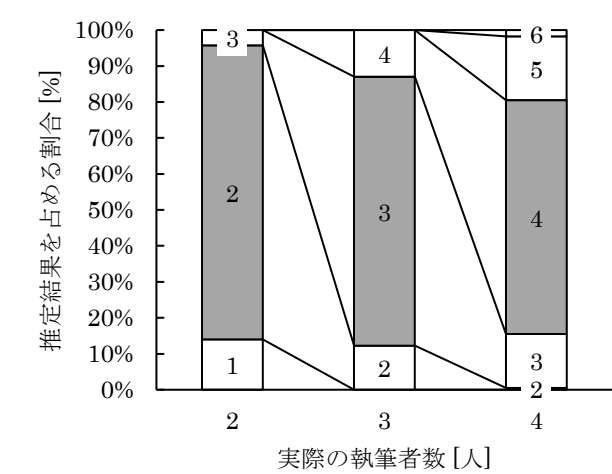


図4.1 実際の執筆者数別の推定執筆者数の割合 (実際の執筆者数2～4人、棒グラフ中の数は推定執筆者数)

ウィンドウのクラスタリングによる既存手法 [4]は、文章の執筆者数が既知としてk-meansを用いると、執筆者を提案手法より高い精度で識別できる(表4.2)。ただし、[4]で執筆者数が未知の場合にx-meansを用いて執筆者の判別を行うと、3人以上の記述が文章に含まれる場合、提案手法の方がより高い精度で執筆者数を推定できる(表4.3)。

表4.2 執筆者数が既知の場合の文体変化点の正解率

実際の執筆者数 [人]	提案手法	既存手法 [4] (k-means)
2	0.655	1.000
3	0.415	0.986
4	0.253	0.962

表4.3 執筆者数が未知の場合の執筆者数の正解率

実際の執筆者数 [人]	提案手法	既存手法 [4] (x-means)
2	0.818	0.956
3	0.748	0.183
4	0.650	0.005

また、ウィンドウと類似度を用いた筆者の既存手法 [5]と比較すると、実際の執筆者数が2人の場合、執筆者数を正しく推定できた確率は14.7%、執筆者数と文体変化点を正しく推定できた確率は18.8%向上した(表4.4)。

表4.4 提案手法と既存手法 [5]との比較 (実際の執筆者数2人)

	執筆者数を正しく推定できた確率	左記に加え、文体変化点を正しく認識できた確率	MAE [人]
提案手法	0.818	0.655	0.183
既存手法 [5]	0.671	0.467	0.328

(2) 1人によって記述された文章

作品同士を連結する前の文章100作品分と、同一の執筆者の作品を2～4作品連結した文章各200組について4.2項で決定したパラメータを用いて文体変化点と執筆者数の推定を行った。結果を表4.5と表4.6に示した。また、文章の連結数別の推定執筆者数の内訳を図4.2に示した。

表4.5 文体変化点と執筆者数の推定結果  
(執筆者数1人)

同一執筆者の 文章の連結数	執筆者数を 正しく推定 できた確率	MAE [人]
1	1.000	0.000
2	0.200	0.835
3	0.015	1.890
4	0.000	3.085

表4.6 文章の連結位置と推定文体変化点の位置

同一執筆者の 文章の連結数	推定執筆者数と連結数が一致し、 かつ推定文体変化点が連結位置から 50単語以内であった確率
1	-
2	0.475
3	0.170
4	0.115

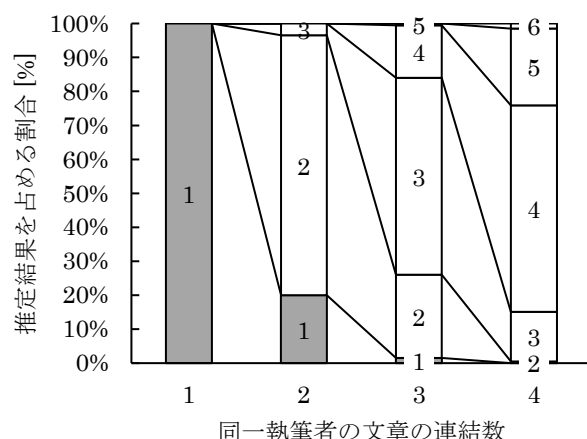


図4.2 同一執筆者の連結数別の推定執筆者数の割合  
(実際の執筆者数1人、  
棒グラフ中の数は推定執筆者数)

図4.2より、連結後の文章が長くなるほど、検出される文体変化点が増える．表4.5と表4.1の結果を比較すると、同一執筆者の文章の連結位置を文体変化点と認識する確率は、複数人の執筆者による文章の連結位置を文体変化点と認識する確率よりも総じて低い．連結位置の前後で執筆者が同じ場合と異なる場合で類似度の値に差があるためと考えられる．執筆者間の文体差が類似度に反映されているかを4.4項で検証した．

#### 4.4. 執筆者間の文体差の検証

執筆者間で文章の類似度に差が見られるかどうかを検証するために、4.3項で用いた作品から同一執筆者の2作品を連結した文章と、異なる執筆者の2作品を連結した文章を各400組用いて、文章中の類似度変化を求めた．4.3項の実験では文単位で文章を連結したが、

この実験では各作品の先頭から数えてちょうど500単語までを用いて、連結後の各文章は500単語目に連結位置が来るように揃えた．100単語間隔で分割位置を動かしながら前半部と後半部との間の類似度を求め、図4.3のように箱ひげ図を作成した．

同一執筆者の作品を連結した場合、500単語目の位置における類似度の平均は0.686、異なる執筆者の作品を連結した場合、0.651となり、分割位置の前後で執筆者が異なる場合の方が低い類似度が算出される傾向にある．ただし、執筆者が同じ場合と異なる場合の類似度の分布は重なっている部分が多い．

文章の端近くで分割して類似度を求めた場合、分割した前半部と後半部で含まれる語彙数や出現回数に偏りが生まれるため、類似度が小さな値を示すと考えられる．文章の端から離れるほど、類似度は大きくなるが、分割位置が文章同士を連結した500単語目に近づくにつれて、類似度の値は小さくなる．よって、作品間の文体の差は検出できていないと言える．提案手法において用いたのは、品詞と助詞、助動詞の出現頻度であるから、作品ごとに品詞と助詞、助動詞の使い方に違いが見られると考えられる．しかしながら、執筆者間の文体の差は十分に検出できていないと言える．

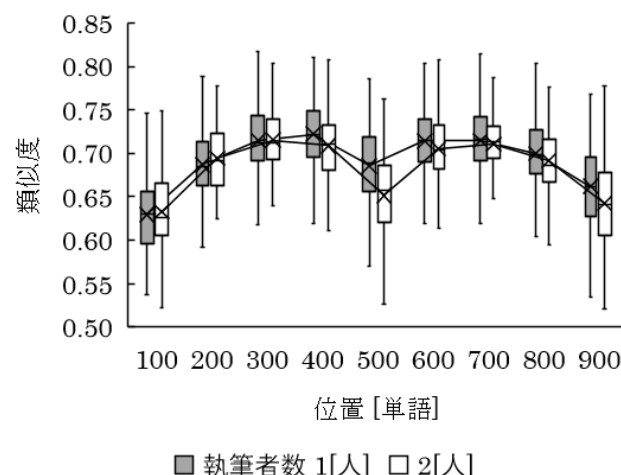


図4.3 同一執筆者(1人)と異なる執筆者(2人)の作品を  
連結した文章の類似度の箱ひげ図

#### 5. まとめ

本稿では、文章の信頼性や品質を反映する文章の執筆者数を推定する手法として、品詞n-gramの出現位置と出現頻度に基づく類似度を用いる手法を提案した．提案手法では、文章内で執筆者が変わる分割点を仮定し、分割点の前後近傍それぞれにおいて、出現する単語の品詞n-gramを取得した．分割点からの距離に応じ

て重み付けを行って語彙の出現頻度ベクトルを求め、分割点前後の語彙頻度ベクトル同士の類似度を計算した。類似度が極小となり、大きな文体変化が認められる点を執筆者が変わった点として、執筆者数を推定した。

評価実験では、1人によって記述された文章を連結することで複数人によって記述された文章を生成し、文体の変化点と執筆者数を推定した。2人によって記述された文章では、執筆者数の正解率は81.8%、平均絶対誤差は0.183人となり、執筆者数が正解し、かつ文体変化点を正しく推定できた確率は65.5%となった。対して、1人によって記述された文章については推定執筆者数に0.838人の誤差が見られ、執筆者ごとの差異を説明する特徴量抽出や類似度の算出手法を検討する必要がある。

### 謝 辞

本研究の一部は JSPS科研費(17KT0085)の助成を受けている。

### 参 考 文 献

- [1] S. Vosoughi, D. Roy and S. Aral, "The spread of true and false news online," *Science*, vol. 359, issue 6380, pp. 1146-1151, 2018.
- [2] C. Shao, G. L. Ciampaglia, A. Flammini and F. Menczer, "Hoaxy: A Platform for Tracking Online Misinformation," in *Proc. of WWW '16*, pp. 745-750, 2016.
- [3] D. Wilkinson and B. Huberman, "Cooperation and Quality in Wikipedia," in *Proc. of WikiSym '07*, pp. 157-164, 2007.
- [4] 塩浦尚久, 山名早人, "日本語の文章を対象にした執筆者人数推定", DEIM Forum 2019 論文集, B5-1, 2019.
- [5] 渡邊充博, E. S. Aung, 山名早人, "文体変化と文体類似度を用いた文章の執筆者数推定", DEIM Forum 2020 論文集, G1-3, 2020.
- [6] B. J. Fogg and H. Tseng, "The Elements of Computer Credibility," in *Proc. of CHI '99*, pp. 80-87, 1999.
- [7] B. J. Fogg, C. Soohoo, D. R. Danielson, L. Marable, J. Stanford and E. R. Tauber., "How do users evaluate the credibility of Web sites? a study with over 2,500 participants," in *Proc. of DUX '03*, pp. 1-15, 2003.
- [8] A. Wierzbicki, "Web Content Credibility", Springer, 2018.
- [9] J. Giles, "Internet encyclopaedias go head to head," in *Nature*, vol. 438, no.15, pp. 900-901, 2005.
- [10] T. Chesney, "An empirical examination of Wikipedia's credibility," *First Monday*, vol. 11, no. 11, 2006.
- [11] F. B. Viegas and M. Wattenberg, "Talk Before You Type: Coordination in Wikipedia," in *Proc. of HICSS '07*, pp. 78-87, 2007.
- [12] P. Dondio, S. Barrett, S. Weber and J. M. Seigneur, "Extracting Trust from Domain Analysis: A Case Study on the Wikipedia Project," in *Proc. of ATC '06*, vol. 4158, pp. 362-373, 2006.
- [13] B. T. Adler, L. d. Alfaro, I. Pye and V. Raman, "Measuring Author Contributions to the Wikipedia," in *Proc. of WikiSym '08*, no. 15, pp. 1-10, 2008.