

Wikipedia の編集履歴に基づく記事の信頼性導出

A Reliability Measure of Articles in Wikipedia based on Edit History

鈴木 優^{1*}

Yu Suzuki¹

金本 径卓¹

Keitaku Kanemoto¹

川越 恭二¹

Kyoji Kawagoe

¹ 立命館大学大学院 理工学研究科

¹ Graduate School of Science and Technology, Ritsumeikan University

Abstract:

In this paper, we propose an editor-based credibility calculation method for articles in Wikipedia. Wikipedia is the encyclopedia which can edit by everyone who access the Wikipedia. Therefore, when an author edit uncertain or uncredible description to an article, the edited article is uncredible. We assume that the credibilities of articles are based on the credibility of the authors. This means that uncredible authors frequently write uncredible descriptions, whereas credible authors frequently write credible descriptions. We assume that uncredible articles are suddenly edit by another credible authors for correcting descriptions. Then, the remain ratio of articles by authors should depends on the credibility of authors. In our experimental evaluation, we confirmed that our proporsed method performs better accuracy.

1 はじめに

現在、様々な事項に関する最新の情報を調べる手法として、主に Wikipedia^{*1} が用いられつつある。Wikipedia とは、誰もが自由に編集することが可能であることが特徴となっている百科事典である。個人が所有する様々な知識を Wikipedia へ記述することによって、多くの利用者同士で知識を共有することが可能となっている。

一方、誰もが自由に編集することが可能であるという特徴から、必ずしも正しい情報、信頼できる情報が記載されているわけではないという点が問題点として挙げられている。例えば、悪意を持った利用者によって意図的に誤った記事を記載する場合や、間違った情報を持った利用者によって記事を記載する場合も考えられる。ある記事に関する正しい知識を持たない利用者がそのような記事を閲覧した場合には、誤った情報を利用者が得てしまう点は問題である。そのため、記事に対して信頼性の度合いを算出することは重要である。

Wikipedia における特徴として、ある著者は信頼できる記事を数多く記述する傾向にあり、別の著者は信頼で

きない記事ばかりを記述することが挙げられる。つまり、一人の著者が信頼できる記事と信頼できない記事をどちらも記述する場合は他の場合と比較して小さい。そのため、もし著者に対して信頼度を算出することが可能であれば、複数の著者によって記述された記事であっても、その記事を編集した著者の信頼度を組み合わせることによって、記事の信頼性を算出することができると考えた。

ここで、著者に対する信頼性を算出する方法について考える。信頼できる記事を数多く記述する著者は、他の著者によってその記事を修正する回数が少ない著者であると言える。これは、信頼できる記事は修正を行う必要が無いが、信頼できない記事は修正を行う必要があるためである。つまり、その著者が編集した記事における著者の記述部分の残存率は、その著者の信頼性と関連があると考えた。

そこで本研究では、著者の記事の残存率を基に記事の信頼性を算出する手法についての提案を行う。提案手法は二つの部分に分割されており、一つは Wikipedia を記述した各著者の信頼度を算出する部分、もう一つは著者の信頼度を基に記事の信頼度を算出する部分である。

* 連絡先：立命館大学大学院 理工学研究科
〒525-8577 滋賀県草津市野路東 1-1-1
E-mail: suzuki@ics.ritsumeik.ac.jp

^{*1} <http://ja.wikipedia.org>

2 関連研究

2.1 記事の編集履歴に基づく信頼度算出手法

Wikipedia の記事の信頼度を求めるために、信頼度を独自に定義することによってその定義された信頼度を算出する研究が行われている。

Adler ら [1] は、Wikipedia の記事内容の変化によって Wikipedia の著者の評価を行うシステムの提案を行っている。これは、著者の編集した記事が後の編集過程で維持されることにより、記事を書いた著者の評価は増加する一方で、編集した記事内容が縮小されることやすぐに取り消された場合に評価が減少する。この点においては、本研究と同様の考え方をを行っている。しかし、Adler らは記事内容の追記と置換を個別に考えることにより、信頼度の算出を行っており、本研究では置換を追記と削除の組み合わせととらえることにより、著者の信頼度を算出している点で異なる。

文献 [2] では、ある記事から次に編集される記事へと遷移する過程における信頼度の変化を Dynamic Bayesian Networks を適用することによって算出している。これは、記事の編集履歴を用いることにより、Wikipedia における記事に対する信頼度を求めることが本研究に類似しているが、本研究では著者に着目した信頼度を求める点において異なるといえる。

さらに、Wikipedia の記事の信頼度に応じて色分けを行い記事の信頼度を利用者に対して視覚的に表示する方法も提案されている。Cross[3] は、編集回数と編集されてからの経過時間に応じた記事の信頼度を算出することによって、同一記事内の文書の色分けを行っている。これにより、記事内の記述のどの部分の記述が信頼できるかを利用者は把握することができる。

Priedhorsky ら [4] は、Wikipedia における記事の閲覧回数と閲覧者が記事に与える影響との関係について研究を行っている。これは、記事の閲覧回数と記事に与える影響から著者の記事に与える貢献度合を求めることを行っている。

一方 Stein ら [5] は、Wikipedia における著者の貢献度合を算出することを行っている。これは、著者の評価を行う点で本研究と類似するが、この研究で用いられている Wikipedia には 'excellent' 等といった記事の評価ラベルが存在し、著者の編集履歴に存在する評価ラベルから著者の貢献度を行っている。そのため、評価ラベルの存在しない本研究で対象となる日本語版の Wikipedia においては適用できない。

Lih[6] は、Wikipedia における記事の編集回数と記事

を編集した異なる著者数との関係から Wikipedia の記事の質の評価を行った。提案されている手法では、記事内容の解釈を必要とすることなく記事の編集回数と記事を編集した異なる著者数とによって Wikipedia の記事と実際に起こる出来事との関連性を導出している。本研究とは、Wikipedia における記事内容が信頼できるかどうかという指標についてでなく Wikipedia と現実世界の出来事との関連性を評価するという提案において異なっている。

Kittur ら [7] は、機械学習により Wikipedia の編集履歴から記事内容の誤りとその誤りを修正する場合における傾向を、編集履歴全体の傾向から知る手法について提案している。本研究では、著者から算出される信頼度について提案しているが、Wikipedia 全体における傾向から誤りを修正しようと考えている点で本研究とは着目点異なる。

2.2 リンク構造に基づく Wikipedia の信頼度算出手法

記事の編集履歴ではなく Wikipedia における記事間のリンク構造から信頼度を算出するような研究もされている。

McGuinness ら [8] は、Wikipedia の引用構造に着目した信頼度の算出方法について研究を行っている。Wikipedia の記事内に出現する単語のうち、単語を表す Wikipedia の記事へのリンクが貼られている単語とリンクが貼られていない単語の違いに着目している。

ある単語からその単語の説明がされている記事へのリンクが貼られている場合、リンクは被リンク側の記事に対して信頼があるものとみなしていると仮定している。そのため、Wikipedia の記事の中において、出現する単語とその出現する単語のうちどれだけの割合でその単語の説明記事へのリンクが貼られているかによって信頼度の算出を行っている。

しかし、Wikipedia におけるリンクは通常の Web ページのリンクと異なり、他の記事に対して積極的にリンクを貼るような構造である。そのため、必ずしもリンクが他の記事に対する信頼関係を明示するものとは限らない。そこで、本研究では記事の内容を直接表す編集履歴から Wikipedia の信頼度の算出を行っている。

2.3 文書間の影響度合に関する研究

文献 [9], [10] では、掲示板の書込み間における影響度合を考慮することにより、掲示板における会話分割手法が提案されている。提案手法では、掲示板における書込

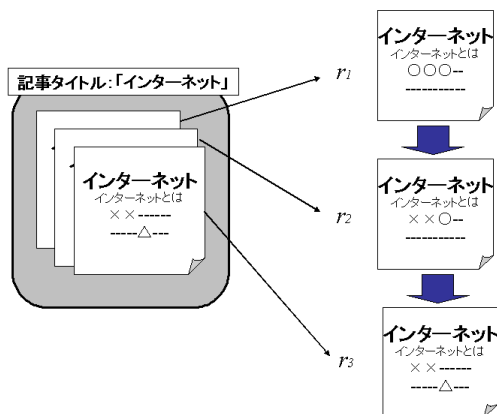


図1 記事における編集履歴の構造

み中の話題が次の書込みに引き継がれているかどうかという点から会話の分割を試みている。本研究では、ある編集が次の編集に対して維持されている割合を考えている点が、掲示板における書込み間における影響度合を考える点と類似している。

しかし、掲示板は前の書込みによる影響を受けることにより、次々と話題が転換していくという構造をとる。一方、Wikipediaは著者が編集する前の記事内容を基準として、記事内容に変更が必要である部分に対して編集を行うといった構造をとる。つまり、掲示板では書込み自体は独立しているため、単純に書込み自体を比較することにより書込み同士の影響度を算出することができる。これに対して、Wikipediaでは編集前の記事に対して追記するだけでなく、置換や削除の行われた編集部分を特定し、編集された部分とその後の編集過程において維持されているかを考慮している点において異なる。

3 提案手法

本章では、Wikipediaにおける信頼度の定義とその算出方法について述べる。

3.1 概要

Wikipediaにおける記事は、図1に示されるように一つのトピックに関する記事の作成とその記事を読んだ閲覧者による編集の積み重ねによって構成されている。記事の閲覧者が記事内容の正誤を判断し、閲覧者自身が編集作業を行うことにより著者の役割を果たすこととなる構造をとっている。

そこで、本稿では記事の内容がどの程度信頼できるかを判断するには、記事の編集作業に携わった著者に着目

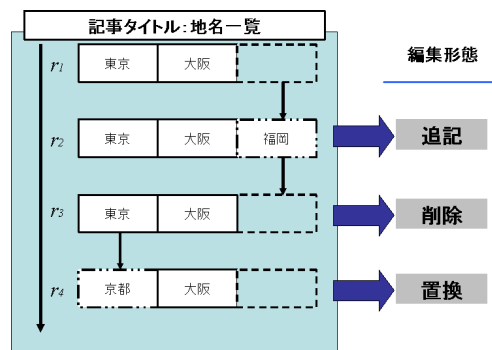


図2 記事の編集形態

すればよいのではないかと考えた。つまり、的確な編集を行うことのできる著者によって編集が行われた記事は、その内容についても信頼がおけるのではないかと考えたためである。また、著者が行った編集は、その後の記事内容に対して影響を及ぼすため編集後の記事内容を見た閲覧者は不適切な編集がなされれば、その部分について適切であるように編集を行うと考えられる。

そのため、本稿では著者の信頼度を算出するために、各著者のもつ記事の編集履歴から著者の行った編集部分を特定し、編集部分を用いることによって著者の信頼度を算出できると考えた。提案手法ではまず、Wikipediaの記事が N 回の編集から構成された場合、その記事を構成する i 回目版 $r_i (i = 1, 2, \dots, N)$ から r_i において行われた著者の編集部分を特定する。

ここで、 r_i において編集された部分とその後の r_{i+1} 以降の編集過程において維持されていれば適切な編集がなされているため、信頼度の高い編集であると考えた。よって、提案手法における信頼度の算出は、 r_i において編集された部分が維持されている割合から求めることとした。

以上の処理を著者の保持する編集履歴に存在する全ての編集に対して行い、Wikipediaの著者に対する信頼度を算出する。さらに、算出された著者の信頼度を用いることによって、記事の編集履歴から記事の編集に携わった著者を特定し、記事の信頼度を得ることとした。

3.2 Wikipediaの編集形態への対応

本節では、 r_{i-1} から r_i へ記事が編集された場合における編集の形態に対する考え方について述べる。本研究では、信頼における編集とは、編集した後の過程において編集した記事が維持されることであると考えている。

このように考えたのは、Wikipediaの誰もが編集可能であるという性質上不適切な編集がなされた場合には、

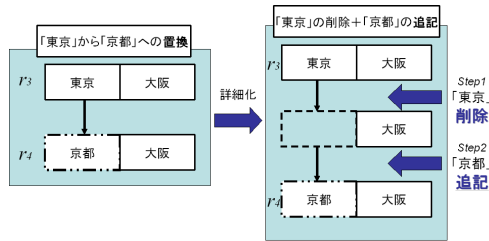


図3 記事内容の置換

その後記事内容が不適切であると判断した利用者によって、記事内容が適切であるように編集し直されるためである。したがって、記事内容が適切でない場合には r_i によってなされた編集部分が変化する割合が大きくなると考えられる。その一方で、記事内容が正しい場合には記事内容に変化を加える必要がないと考えられるため、記事内容が変化する割合が小さくなると考えられる。よって、本研究では r_i において編集された部分が維持されている割合を記事の信頼度として算出することとした。そこで、記事の編集を行った部分を特定すること、及びその編集部分が維持されている割合をどのように算出するのか重要となってくる。

そこで我々は記事の編集された部分を特定するために、Wikipedia における記事の編集形態には以下の三つがあると考えた。

- 記事内容の追記
- 記事内容の削除
- 記事内容の置換

記事内容の追記とは、 r_i の著者が r_{i-1} を閲覧することによって r_{i-1} の内容が不足していたり、さらに追記すべき内容があると判断した場合にその内容を r_{i-1} に対して追加することによって、 r_i を作成する編集をいう。図2における r_1 から r_2 への編集のように、 r_1 に対して「福岡」を加えるような編集のことを本研究では記事内容の追記と考える。

これに対して、記事内容の削除とは r_i の著者が r_{i-1} を閲覧することによって、 r_{i-1} の内容が冗長であったり不適切であると判断した場合にその内容を r_{i-1} から削除する編集のことをいう。図2における r_2 から r_3 への編集において、「福岡」を取り除くような編集を本研究における記事内容の削除ととらえることとする。さらに、記事内容の置換とは記事内容の追記だけの場合でも削除だけの場合でもなく、 r_{i-1} の内容が適切であっても表現として正しくない場合等、正しくない部分を r_i において別の語で置き換える編集形態をいう。

しかし、我々は記事内容の置換は本質的には記事内容

の削除を行い、その後に記事内容を追記することであると考えた。例えば、 r_3 を r_4 への編集を詳細化すると「東京」を削除した後に「京都」を削除した部分に対して追記することと考えられる(図3)。よって提案手法では、記事内容の追記と削除に対して編集の信頼度を求めることにより、記事内容の置換にも対応することとした。

3.3 記事内容の追記における維持割合の算出

本節では、記事内容が追記された場合における追記部分を特定する手法とその追記した部分が後の編集過程において維持されている割合を算出する方法について述べる。

Wikipedia における編集では、記事中における単語を変化させるだけでなく、句読点やリンク先の追加または削除といった細かな編集もなされる。そのため、このような編集は記事内容に影響を与えないと考えられるため、記事内容に対する信頼度に影響を与えないと考えられる。そのため、本研究では記事内容の信頼度を考える際に、記事内容の変化の記事に含まれる単語の変化として考えることとした。

そこで、 r_i から r_{i-1} の間で記事内容が追記される場合において、 r_{i-1} に含まれる単語を $F(r_{i-1})$ 、 r_i に含まれる単語を $F(r_i)$ とする。編集において追記された単語群 $A(i-1, i)$ は、 r_i から r_{i-1} と r_i において共通に存在する単語を除いた単語であるから、

$$A(i-1, i) = \overline{F(r_i)} \cup \overline{F(r_{i-1} \cap r_i)} \quad (1)$$

と表すことができる。

次に、追記部分 $A(i-1, i)$ を用いることにより r_i から k 番目の編集 $r_k (i < k)$ との間における追記部分の維持割合についての算出手法について述べる。

式(1)により記事内容の追記部分が特定できるため、 r_i と r_k における追記部分の変化割合は、追記した単語群 $A(i-1, i)$ が記事 r_k に存在する割合と考えることができる。

例えば、図4において r_1 から r_2 において著者は「福岡」及び「名古屋」を追記したことから、 r_3 において「福岡」が削除され、 r_4 においては「福岡」が追記されると共に「名古屋」が削除されている。そのため、 r_1 から r_2 の間でなされた追記部分の維持割合は、 r_3 と r_4 においてそれぞれ $1/2$ であることがわかる。

したがって、追記部分の維持割合 $I_{add}(r_i, r_k)$ は以下の式により求めることができる。

$$I_{add}(r_i, r_k) = \frac{|A(i-1, i) \cap F(r_k)|}{|A(i-1, i)|} \quad (2)$$

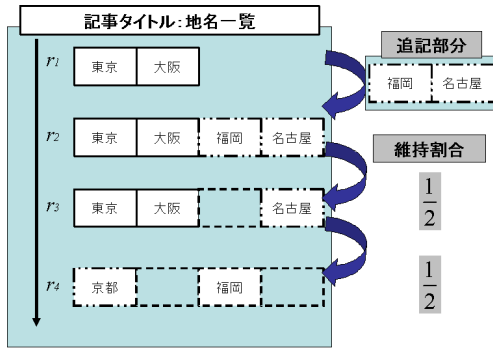


図4 追記部分が維持される様子

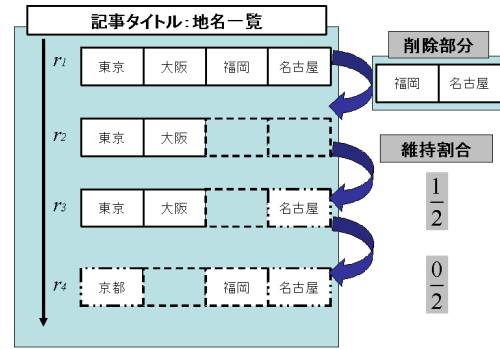


図5 削除部分が維持される様子

3.4 記事内容の削除における維持割合の算出

本節では、まず記事内容が削除された場合における削除部分を特定する手法について述べる。次に、その削除した部分が後の編集過程において残留している割合を算出する方法について説明する。

r_{i-1} から r_i において著者が削除した部分の単語は、 r_{i-1} から r_{i-1} と r_i において共通して出現する単語を除いた単語群 $D(i-1, i)$ は以下の式で表すことができる。

$$D(i-1, i) = \overline{F(r_{i-1})} \cup \overline{F(r_{i-1} \cap r_i)} \quad (3)$$

また、 r_i から r_k の間で削除した部分が元に戻されている割合は、 $D(i-1, i)$ に含まれる単語が r_k において出現している割合であると考えられる。

つまり、図5において著者は r_1 から r_2 において、「福岡」及び「名古屋」の二つの単語を削除していることがわかる。そして、 r_3 において「名古屋」の部分が追記により元に戻され、さらに r_4 において「福岡」の部分も元に戻されている。したがって、 r_1 から r_2 の間で削除された部分が維持されている割合は r_3 では $1/2$ 、 r_4 では $0/2$ であることがわかる。

以上により、削除部分の維持割合 $I_{del}(r_i, r_k)$ は、削除部分 $D(i-1, i)$ が r_k において共に出現していない割合と考えることができるため、以下の式により算出される。

$$I_{del}(r_i, r_k) = 1 - \frac{|D(i-1, i) \cap F(r_k)|}{|D(i-1, i)|} \quad (4)$$

3.5 編集に対する信頼度の定義

本節では、以上で述べた記事内容の維持割合から記事内容の編集に対しての信頼度を定義する。

我々は、記事内容の編集部分とその後の編集過程において維持されれば、その編集は適切であると判断できると考えた。なぜなら、記事を編集する必要があると判断した者によっても、著者が編集された部分については修正を加える必要がないと判断されるためである。また、利用者は記事の編集履歴を閲覧することができるため、不適切な編集部分を特定することによって、記事の修正を行うこともできる。したがって、編集部分が維持されている割合が高いということはその編集部分について閲覧者の信頼を得たものととらえることができる。

編集部分が維持されている割合とは、3.3節と3.4節において述べたように $I_{add}(r_i, r_k)$ 及び $I_{del}(r_i, r_k)$ により算出される。また、記事内容の置換は記事内容の追記と削除の組合せとして考えた。したがって、記事の編集が行われた場合の信頼度とは、記事内容の追記の場合と削除の場合における編集部分が維持されている割合の総和であると考えることが可能である。

よって、上述した三つの編集形態を考慮した r_i に対する編集の信頼度 $Rel(r_i, r_k)$ は、 r_i から r_k に至るまでの $I_{add}(r_i, r_j)$ と $I_{del}(r_i, r_j)$ ($i+1 \leq j \leq k$) のそれぞれの総和として、以下の式により表される。

$$Edit_rel(r_i, r_k) = \sum_{j=i+1}^k I_{add}(r_i, r_j) + \sum_{j=i+1}^k I_{del}(r_i, r_j) \quad (5)$$

3.6 著者の信頼度の算出

本節では、Wikipediaにおける著者の信頼度の算出方法について述べる。

本研究では、著者の信頼度を著者が行った記事の編集に対する信頼度を用いて算出する。我々は、著者のもつ編集履歴から著者の編集した記事を抽出し、編集内容が

それぞれ信頼がおけるものかどうかを式 (5) によって評価する。そして、著者の行った編集内容は著者の信頼度を表すと考えたため、各編集に対する信頼度により著者の信頼度を表すこととした。

著者 a が過去に編集した編集履歴において、 M 回の編集を行った場合における著者の信頼度 $Rel_author(a)$ は以下の式により、求めることができる。

$$Author_rel(a) = \frac{\sum_{j=1}^M Edit_rel_j(r_i, r_k)}{M} \quad (6)$$

式 (6) は、 a の行った編集に対する信頼度の平均値である。よって、著者 a の行った編集が維持されるような場合、すなわち信頼のおける編集を a が行うことによって、著者に対する信頼度も向上することがわかる。

3.7 著者の信頼度を用いた記事の信頼度の算出

本節では、Wikipedia における記事内容の信頼度を算出する方法について述べる。

我々は信頼度の高い著者は、信頼度の高い記事を作成する可能性が高いと考えた。そのため、記事の編集に関した著者の評価を用いることによって、記事内容の信頼度を算出することが可能であると考えた。そこで記事の信頼度を算出する際に、記事の編集履歴からその記事の編集を行った著者を抽出し、抽出された著者を式 (6) によって評価する。そして、記事 d の信頼度 $Article_rel(d)$ は、記事の編集履歴から記事 d を書いてきた著者 a_i の信頼度の平均によって表す。したがって、記事 d が K 人の著者によって編集されている場合における記事の信頼度は、

$$Article_rel(d) = \frac{\sum_{i=1}^K Author_rel(a_i)}{K} \quad (7)$$

となる。

4 評価実験

本提案手法によって算出される記事の信頼度は、著者の信頼度を用いて表されるため著者の信頼度を評価する評価実験を行った。評価を行うにあたり、実際に Wikipedia の編集履歴データを用いることとした*2。

4.1 予備実験

評価実験を行う前に、式 (5) における k の値を決定を行う必要がある。 k は、 r_i によってなされた編集部分が

表 1 予備実験の結果

k	101 点平均適合率
1	0.727
2	0.750
3	0.745

r_i から何回先までの編集における編集部分の維持割合を考慮すればよいのかを決定する値である。また、 k の値が大きくなると r_{i+1} 以後の編集による影響も大きくなってくると考えられるため、 $1 \leq k \leq 3$ において信頼度を算出する際に最も適した値を採用することとする。そこで、Wikipedia の編集履歴データから記事の編集履歴を 100 件抽出し、 k を 1 から 3 値まで変化させることによって、 $Edit_rel(r_i, r_k)$ を算出した。 k の値の決定にあたっては、情報検索における評価手法である 101 点平均適合率を用いることとした。

101 点平均適合率を求めるにあたり、再現率 α 及び適合率 β を以下のように定義する。 R を提案手法により検出された正解編集数、 N を評価対象となる全編集数、 C を評価対象となる編集集中における正解とする全編集数とすると、

$$\alpha = \frac{R}{N} \quad (8)$$

$$\beta = \frac{R}{C} \quad (9)$$

となる。

また、正解集合の作成にあたり記事の編集に対して信頼の高い編集を決定することは困難である。しかし、明らかに誤った内容である編集や故意に記事内容と無関係な記述を行っている編集は、編集内容を実際に見ることによって判断が可能であるため、これらの信頼が低いと判断される編集を正解集合として作成した。そして、各編集に対する $Edit_rel(r_i, r_k)$ の値が低い方からランキングを行うことによって、101 点平均適合率を算出した。なお、記事の変化割合を算出する際の単語として名詞及び未知語を用いた。

予備実験の結果を表 1 に示す。結果より、 $k = 2$ の場合に平均適合率が 0.75 と最も高い値を示したために、本稿では $k = 2$ として評価実験を行うこととした。

4.2 実験の目的

実験では、記事の信頼度を算出するための前提である著者の信頼度 $Author_rel(a)$ の有効性を検証する。また、 $Author_rel(a)$ は著者の編集履歴から算出される

*2 <http://download.wikimedia.org/jawiki/>

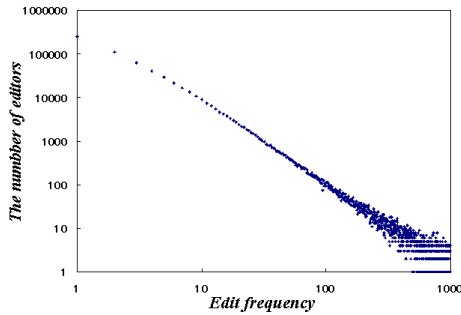


図6 編集回数と著者数の関係

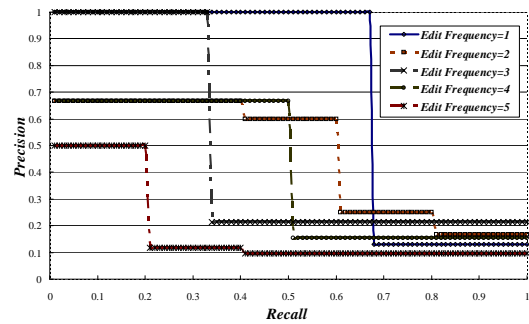


図7 101点再現率-適合率曲線

$Edit_rel(r_i, r_k)$ の平均値であるため、著者 a の行った編集回数に影響されると考えられる。そのため、 a の行った編集回数ごとに精度の変化を検証することを行う。

4.3 実験方法

実験方法として予備実験同様に、101点平均適合率により評価を行った。また、あらかじめ実験対象となる編集履歴データ 657555 人の著者を抽出し、編集回数と著者数との関係の統計を作成した。図6に見られるように、7回までに約80%の著者が集中していることがわかる。このように、編集回数が1回や2回といった著者によって半数を超える著者数を占めることから、Wikipediaは編集回数の比較的多い特定の著者によって記事の編集がなされているわけではなく、様々な著者によって編集が行われていることがわかる。

したがって、実験では編集回数が1回から5回の著者をそれぞれ100人ずつランダムに抽出し、編集回数ごとに平均適合率を算出することとした。正解集合となる著者として、3人の被験者により各著者の編集内容を閲覧し、2人以上の被験者が不適切な編集であると判断した編集を編集履歴の一つでも含む著者を正解集合とした。その結果、正解である著者として選ばれた著者の数は1回から5回の編集回数の順にそれぞれ100人のうち3人、5人、3人、2人、5人となった。

4.4 結果と考察

実験結果として、再現率適合率曲線と平均適合率のグラフを図7及び図8に示す。

編集回数が1回から5回の間において、1回の場合に平均適合率が0.69と最も高い値を示した。これは、著者の行った1回の編集だけによって、著者の信頼度が決定されるためである。つまり、1度不適切な編集を行えばその編集がその後すぐに修正されることによって、信

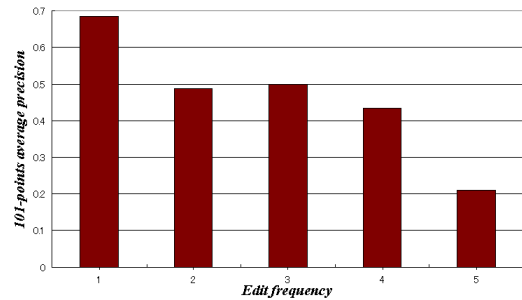


図8 101点平均適合率

頼度は低下するために信頼のない著者を判断することが可能となる。

これに対して、編集回数が2回以上の場合において平均適合率が減少したのは、著者のもつ編集履歴が複数となったことに起因すると考えられる。というのも、同じ著者の保持する複数の編集履歴中においても、不適切な編集と適切と考えられる編集が存在する。そのため、1度不適切な編集を行ったとしても、他の編集履歴中に適切な編集が存在すれば、不適切な編集による信頼度の低下が抑えられたためであると考えられる。

したがって、正解とされる著者の編集のうち不適切な編集が全ての編集において行われる場合とそうでない場合とで著者の信頼度が異なってくるのがわかる。しかし、図6にも示されるように、編集回数が少ない方に著者の人数が分布している傾向にあるため、1回に近い値において精度が高い点において本提案手法は有効であると考えられる。

また、著者が編集を行った部分について間違った内容を書いてはいないが、記事内容として議論の必要な編集がなされると、その部分について繰り返し編集がなされる場合がある。このような場合、最終的に編集部分が維持されたとしても、繰り返し編集が行われるため、信頼度が低下しランキングの上位に出現する結果となった。

したがって、信頼がある編集でも繰り返し編集される場合には、本提案手法により信頼度を適切に判断することは困難であると考えられる。

5 おわりに

本稿では、Wikipediaにおける信頼度の定義を行い、信頼度の算出方法についての提案を行った。Wikipediaにおける著者の信頼度を算出するために、著者が行った編集に対する記事の編集部分を評価することによって、信頼度を算出した。信頼度の算出を行うため、著者がWikipedia上において行った編集行為の結果が直接表現されているため、編集履歴を用いた。

本研究では、Wikipediaの編集を記事の追記、削除、置換の三つの形態として分類できると考えた。この方法により、編集履歴から著者の行った編集部分を記事の編集履歴の前後を比較することによって特定した。そして、特定された編集部分からその後の記事の編集履歴をたどることによって、著者の編集履歴の維持割合を特定した。その割合から、著者に対する信頼度を算出した。

評価実験において、提案手法が有効であるかどうかを確かめるために、実際のWikipediaの編集履歴を用いて評価実験を行った。著者の編集回数ごとに分類して評価実験を行った結果、提案手法が大規模なデータにおいて有効であることを確かめた。

今後の課題として、適切な編集回数を設定する方法について考えている。実験結果において、著者自身が行った編集操作の結果、その著者の信頼度が低下する場合があった。また、繰り返し編集部分が削除されたり追記される場合においても適切な信頼度を算出することができなかった。そのため、編集内容についても考慮した信頼度の算出を行う必要があると考える。

また、信頼度の提示方法についても課題であると考えている。本研究では、評価実験のように過去の編集時点における信頼度を算出することができるが、新たな編集を行った場合には即座に信頼度を算出することができない。また、大規模な計算が必要であるため、計算量の削減が必要である。さらに信頼度に応じて記事の色分けを行うなど、利用者に対して直感的に理解可能な信頼度提示手法を検討する必要があると考えている。

謝辞

本研究の一部は、文部科学省科学研究費補助金(若手研究(B), 40388111)によります。ここに記して謝意を表します。

参考文献

- [1] B. Thomas Adler and Luca de Alfaro. A content-driven reputation system for the wikipedia. *Proceedings of the 16th international conference on World Wide Web*, pp. 261 – 270, 2007.
- [2] Honglei Zeng, Maher A Alhossaini, Li Ding, Richard Fikes, and Deborah L McGuinness. Computing trust from revision history. *The 2006 International Conference on Privacy, Security and Trust (PST 2006)*, 2006.
- [3] T.Cross. Puppy smoothies:improving the reliability of open, collaborative wikis, 2006.
- [4] Reid Priedhorsky, Jilin Chen, Shyong (Tony) K. Lam, Katherine Panciera, Loren Terveen, and John Riedl. Creating, destroying, and restoring value in wikipedia. *Association for Computing Machinery GROUP '07 conference proceedings*, 2007.
- [5] Klaus Stein and Claudia Hess. Does it matter who contributes? a study on featured articles in the german wikipedia. *Proceedings of the 18th conference on Hypertext and hypermedia*, pp. 171 – 174, 2007.
- [6] Andrew Lih. Wikipedia as participatory journalism: Reliable sources? metrics for evaluating collaborative media as a news resource. *In Proceedings of the 5th International Symposium on Online Journalism*, 2004.
- [7] Aniket Kittur, Bongwon Suh, Bryan A. Pendleton, and Ed H. Chi. He says, she says: conflict and coordination in wikipedia. *Proceedings of the SIGCHI conference on Human factors in computing systems*, pp. 453 – 462, 2007.
- [8] Deborah McGuinness, Honglei Zeng, Paulo Pinheiro da Silva, Li Ding, Dhyanes Narayanan, and Mayukh Bhaowal. Investigations into trust for collaborative information repositories: A wikipedia case study. *The Workshop on the Models of Trust for the Web (MTW'06)*, 2006.
- [9] 荒川智之, 鈴木優, 川越恭二. 電子掲示板における会話分割手法. データベースと Web 情報システムに関するシンポジウム (DBWeb2006) 論文集, pp. 27–34, 2006.
- [10] 松尾豊, 大澤幸生, 石塚満. 電子掲示板における会話からのハイライト部分の抽出. 第 46 回人工知能基礎論研究会, 2002.