ソーシャルストリームの読み飛ばしを考慮した コンテクスト適合型インタフェース

†九州大学芸術工学部 〒815-8540 福岡県福岡市南区塩原 4-9-1 ‡九州大学大学院芸術工学研究院 〒815-8540 福岡県福岡市南区塩原 4-9-1

E-mail: † mari7toki@gmail.com, ‡ ushiama@desighn.kyushu-u.ac.jp

あらまし

Twitter 等のソーシャルストリームを閲覧する際、ユーザは全てのコンテンツを読む訳ではなく、価値が低いと推定されるものを読み飛ばすことが多い。しかし、コンテンツの価値は常に同一ではなく、ユーザのコンテクストによって変化する場合がある。例えば、忙しいときに読み飛ばされたコンテンツの中にも、時間的に余裕があるときにはユーザが価値を感じるものが含まれる可能性がある。本論文では、時間的に余裕がないコンテクスト「busy」とそれ以外の「free」を想定し、ユーザが「busy」のときに読み飛ばされたコンテンツから、ユーザにとって価値の高いコンテンツを抽出し、「free」のときに再提示することにより、ユーザがソーシャルストリームから効果的に情報を取得可能とするインタフェースを提案する。そして、コンテクストが free であるときの、ブラウザ上でのユーザの振る舞いによってユーザの興味プロファイルを生成し、価値の高いコンテンツを抽出する手法を提案する。そして、被験者実験に基づいて、提案手法の有効性を定量的に検証する。

キーワード ソーシャルコンテンツ, 読み飛ばし, コンテクストアウェアネス, コンテンツの価値, 個人化, ソーシャルストリーム, SNS, ユーザの振る舞い, tf-idf, インタフェース

A Context-Adaptive Interface for Social Stream Suitable for Skip Reading

Marina TOKI[†] and Taketoshi USHIAMA[‡]

† Graduate School of Design, Kyushu University 4-9-1 Minami-ku, Fukuoka-shi, Fukuoka, 815-8540 Japan ‡ Faculty of Design, Kyushu University 4-9-1 Minami-ku, Fukuoka-shi, Fukuoka, 815-8540 Japan E-mail: † mari7toki@gmail.com, ‡ ushiama@desighn.kyushu-u.ac.jp

Keyword Social Contents, skipping, Context Awareness, Value of Contents, Personalization, Social Stream, SNS, user's behavior, tf-idf, interface

1. はじめに

Twitter や Facebook に代表される SNS (ソーシャルネットワーキングサービス) が爆発的に普及し、友人・知人とのコミュニケーション,情報収集等,様々な目的のために日常的に利用されるようになった. ユーザが SNS に投稿した「ソーシャルコンテンツ」は、そのユーザをフォローしているユーザに自動的に配送される. ユーザは複数のユーザをフォローするために、情報が断続的に配信される. このように、多くのSNS では、断続的に、様々なユーザから様々な種類のコンテンツが配信され、システムは配信されたコンテンツを一列に並べた形式でユーザに提示する. SNS 上でユーザに配信されたコンテンツ列は、一種のストリームデータと考えることができ、「ソーシャルストリーム」と呼ばれる.

一方,スマートフォンやタブレット型コンピュータの普及で,

場所と時間を問わず SNS を利用可能となった. ソーシャルストリームには次々に新しいコンテンツが追加されるため,多くのユーザは,日常生活の隙間時間を利用してソーシャルストリームを閲覧するようになった.

一般に、ソーシャルストリーム上には多種多様なコンテンツが含まれているため、その中には、ユーザにとって価値の高いコンテンツもあれば、ユーザにとって価値の低いコンテンツも存在する。また、ソーシャルストリームは大量のコンテンツから構成され、ユーザはコンテンツーつ一つの内容を正確に確認することが困難である。そこで多くのユーザは、忙しい朝や隙間時間等、時間的な余裕がないときには、ソーシャルストリームの中から目についたものだけを拾い読みすることがある。しかし、内容を正確に確認せず読み飛ばされたコンテンツの中に、ユーザにとって価値の高い情報が含ま

れている可能性がある。このようなコンテンツは、ユーザに時間的な余裕があれば、有効に活用できるかもしれない. 我々は、ユーザにとって価値が高いにも関わらず、ユーザが 閲覧する際のコンテクストが適切でなかったために、活かされなかったコンテンツを、「見落としコンテンツ(Slipped Contents)」と呼ぶ.

本研究では、ユーザのコンテクストを、時間的に余裕がない状況である「busy」と、時間的に余裕がある状況の「free」の2種類に分類する.このとき、見落としコンテンツは、「ユーザが内容を確認しなかったコンテンツの内、ユーザのコンテクストが busy のときに提示され、ユーザにとって価値の高いコンテンツ」と定義できる. 見落としコンテンツの定義を図示したものを図 1 に示す.

本論文で提案するインタフェースでは、ソーシャルストリームを読むユーザの振る舞いから、ユーザのコンテクストがbusyであるか、freeであるかを自動的に区別する.そしてbusyのときに流れたコンテンツの中から、ユーザの興味プロファイルに適合する価値の高いコンテンツ(見落としコンテンツ)を抽出する.

本研究では、ユーザの興味プロファイルを、ソーシャルストリームの閲覧履歴から自動的に作成する。そのために、一般的な TF-IDF 法にコンテクストの要素を取り入れることで、ソーシャルコンテンツにより適した手法である TF-IDF-RT 法を開発した。

そして、コンテクストが free のときに、ソーシャルストリームの再構成により、抽出した価値の高いコンテンツをユーザに再提示する.これにより、ユーザがコンテクストに応じて適切なタイミングで適切な情報を受け取ることで、コンテンツの価値を最大限に活かすことを支援する.

本論文の構成は以下の通りである。第2章で関連研究を紹介する。第3章では、提案システム「Context Catcher」の概要を述べる。第4章ではシステムを実現するための以下の手法について述べる。ユーザの振る舞いデータの取得方法と、コンテクスト判定の手法、さらにコンテクストと滞留時間に着目した興味プロファイル作成手法と、コンテンツ価値判定手法について述べる。第5章で被験者実験について述べ、第6章で考察を行い、第7章でまとめを述べる。

15	読まれなかったコンテンツ							
		コンテクスト						
		busy	free					
コンテンツ	高い	見落としコンテンツ						
ツ 価 値	低い							

図 1 見落としコンテンツの位置づけ

2. 関連研究

これまでにも、膨大かつ多様な情報の中から、効率的に価値の高い情報を選別する研究がなされている[1]. 例えば Facebook で利用されている、EdgeRank アルゴリズムでは、ユーザと投稿者との「親密度」という指標をもとに、表示するコンテンツをランク付けしている[2]. しかし EdgeRank では主にユーザの交友関係に重きを置いており、ユーザ自身の嗜好とコンテクストを考慮していない. そのため、ユーザの興味のある投稿が読み飛ばされ、充分に活用されない可能性がある.

松尾ら[3]は、web 上のユーザの振る舞いから興味を把 握し,個人化した情報提示を行う手法を提案している.具 体的には、ユーザが閲覧した文書の履歴から、ユーザにとっ て重要度の高い語を抽出し,ブラウジングでの読み飛ばしを 防ぐ.この手法では、ユーザが閲覧する度にユーザにとって 「身近な語」を収集し、文書を閲覧する際に「身近な語」との 共起語を重要語とする. 興味をユーザの過去の閲覧文書 から抽出し,個人化に利用する点は,本研究と共通してい る. しかし, この手法では, 一般的な web 文書を対象として いるため、ユーザの閲覧時のコンテクストの違いは考慮して いない. また, ソーシャルストリーム上のコンテンツは, ユーザ の閲覧時のコンテクストに強く依存して, その価値が変化す ると考えられるため、「読み飛ばし」が起こる要因が異なると 考えられる. さらに、松尾らの手法では、コンテンツを「読み 飛ばされないように」重要語をハイライトするのに対し、本研 究では「読み飛ばされた」重要なコンテンツを、ユーザに時 間的な余裕があるときに補助的に再提示することで、ユーザ のコンテクストに応じてコンテンツの有効的な活用の支援を 目的としている.

ユーザのコンテクストに着目した研究として、中野ら[4]は、発想支援において、オフタイム(休息のための時間)を有効に活用することに着目し、オンタイム(業務や勉学に充てる時間)で取得したユーザの暗黙的な興味情報に基づき抽出したテレビ番組をオフタイムに埋め込むことにより、ユーザの発想の広がりを促進するシステムを提案している。このシステムでは、オンタイムに於ける検索クエリ等からユーザの興味を推定し、ユーザの興味に関するテレビ番組を、オフタイム(発想しようと意識していないコンテクスト)に提示する。これにより、オンタイムで得られなかった発想が生まれ、同じコンテンツでもコンテクストに応じてユーザに与える影響が異なることを目指している。この研究では、具体的に有効な結果は出なかったものの、意外性のある興味深い番組抽出に成功している。

林ら[5]は、ソーシャルコンテンツを含む CGDC(Customer Generated Digital Contents)特有の特徴を考慮し、ドキュメントの発生順序を考慮したTF-IDF法の提案を行っている. 時間とともに発生するソーシャルコンテンツには、通常のTF-IDF 法では最新の状況を的確に反映することができな

いとし、文書数の増加に応じて IDF 項の値が変化するよう、イベントの影響により変化した特徴語を TF-IDF 法に取り入れている。本手法では「変化」のあったコンテンツを重要だと捉えているが、我々の手法では、ユーザが興味のあるコンテンツが対象であり、時間経過によっても変化のない情報も、重要なものだと考える。

3. システムの概要

我々は「同一のコンテンツでも、ユーザのコンテクストによって情報の価値は異なる」という仮説をたてた。本論文で提案する SNS クライアント「Context Catcher」は、この仮定に基づきユーザが忙しい時間に閲覧した見落としコンテンツを抽出し、時間的な余裕がある時に再提示することにより、コンテンツが本来有する価値を生かすというアプローチをとる。なお、本研究では SNS の中でも、投稿数が多く、読み飛ばされやすいと考えられる Twitter を対象とする.

本システムの処理の流れを図2に示す.

- ① ユーザがシステムを利用しているとき、システムはユーザの振る舞いからコンテクスト(busy/free)を判定する.
- ② ユーザが free のときは、ユーザの振る舞いからユーザの 興味を表す2つの興味プロファイルを生成する. 興味プロファイルは、興味のある事項を表す「興味単語プロファイル」と、興味のあるユーザを表す「興味ユーザプロファイル」から構成される.
- ③ ユーザが busy のとき,ユーザが読み飛ばしたコンテンツを抽出し,興味プロファイルに合致する価値の高いコンテンツを見落としコンテンツと判定し,見落としコンテンツデータベースに蓄積する.
- ④ ユーザが free のとき, 見落としコンテンツデータベースに 蓄積された見落としコンテンツを新しいソーシャルストリームと融合して再提示する.

以上の流れによって、ユーザは能動的な操作なしに、見落としていた価値の高い情報を活用可能となる.

本システムは、PHP と JavaScript を利用して記述され、 JavaScript が動作するWeb ブラウザ上で動作する. また、データベース管理システムとしては MySQL を利用している.

図 3 に Context Catcher のインタフェースのスナップショットを示す. 本研究では、ユーザの自然な SNS の利用法のまま、コンテンツ価値を向上させることを理想としている. その

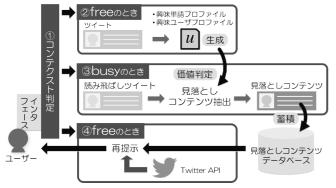


図 2 システムの処理

ため、特殊な設備を使用せず、インタフェースデザインも最も一般的であると考えられる公式の Twitter クライアントのデザインを踏襲している.



図 3 Context Catcher のインタフェース

4. 提案手法

4.1. コンテクスト判定

4.1.1. ユーザの振る舞いデータ

本研究では、ソーシャルストリーム閲覧時のユーザの振る 舞いを抽出し、以下の目的のために利用する.

- 1. ユーザのコンテクスト(busy/free)の判別(図 2①)
- 2. ユーザの興味の抽出(図 2②)
- 3. 見落としコンテンツの抽出(図 2③)

ユーザの振る舞いを抽出目的として,これまでにも多くの 手法が提案されている. 例えば、Web ページにおいて、ユー ザが興味を持った箇所を推定するために,アイトラッキング が用いられることが多い[6]. これは、ユーザが今注視してい るスクリーンの位置と時間を, 眼球の動きをアイトラッカー等 の特別な装置を利用して取得し、分析する手法である. アイ トラッキングを利用することにより、ユーザの興味を抽出可能 である. しかし、全てのユーザがアイトラッカーを利用すること を前提とすることは現実的ではない. そこで, 本研究では特 別な装置を利用せずにユーザの振る舞いを抽出することを 目指す.このため、ブラウザ操作を元にユーザの振る舞いを 取得する. 今回対象とするソーシャルストリームでは,一次 元的にストリーム化されたコンテンツをスクロールによって閲 覧することが多い. そこで、特にスクロール操作に注目する. 他に、ソーシャルストリームを構成する個々のコンテンツ(ツイ ート)に対して、「表示領域の外に出た時刻」、「外部リンクの アクセス」,「リツイート」,「お気に入り登録」等をユーザの振 る舞いに関するデータとして取得する. 振る舞いデータはそ れぞれイベントとして管理され,発生時刻とその識別番号 (id)をもとに RDB 上に記録される.

4.1.2. 滞留時間の推定

一般的に、ユーザは画面中央のツイートを読むと考えられるため、個々のツイートが表示領域の外に出た時刻に基づいて、中央に位置するツイートの滞留時間を推定する(図4). 具体的には、振る舞いデータのうち、ソーシャルストリームを構成するツイートがスクロールによりブラウザ画面の外部に移動したことを表す displayEnd イベントを用いる. いま、ツイート Tw_j が、スクロールによってブラウザの表示領域の外部に出た時刻を displayEnd_jと表すことにする. また、表示領域上に n 個のツイートが表示されているとする. このとき、ユーザは表示領域内に表示されているツイートのうち、上からn/2 番目に表示されているツイートを注目していると考える. そこで、 Tw_j を読むのにかかった時間 rt_j を、ツイート rt_j で、 rt_j rt_j rt

なお、ツイートを閲覧する際には、ツイートの本文を読むだけではなく、ツイートに含まれる URL をクリックして Web ページを閲覧したり、リツイートや、お気に入り登録等の振る舞いを行うことがある。したがって、上記で推定される滞留時間には、ツイート本文の閲覧だけでなく、上記のような振る舞いに要する時間が含まれる可能性がある。また、ツイート長によっても、閲覧に要する時間は異なると考えられる。そこで、上記の要因を考慮して、滞留時間を正規化する。また、ユーザの注目点が中心からずれることがあるため、窓関数を適応し平滑化を行う。本手法によって計測された滞留時間の例を図5に示す。

4.1.3. コンテクストの判別

システムを利用した際のユーザのコンテクスト(busy または free)は、システムを起動してソーシャルストリームを閲覧する セッション毎に判別する. 具体的には、そのセッションにおけるユーザの平均的な滞留時間および平均的なスクロール速度を利用して閾値に基づいて判別する.

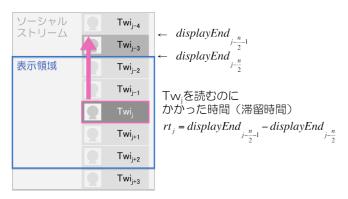


図 4 ツイート単位の滞留時間の計算方法

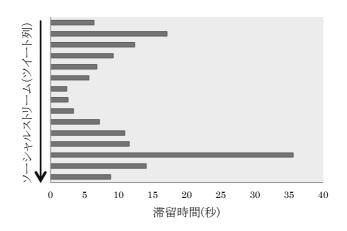


図 5 推定された滞留時間の例

4.2. コンテンツ価値判定

見落としコンテンツの抽出のために、コンテンツの価値判定が必要である。コンテンツの価値判定のために、本システムでは、ユーザのコンテクストが free のときのソーシャルストリームに対するユーザの振る舞いを利用して、ユーザの興味を表す興味プロファイルを生成する。次に、ユーザのコンテクストが busy のときに提示されたツイートのうち、見落としコンテンツとして再提示するべきものを、興味プロファイルを利用して抽出する。

本研究ではユーザの興味は、ツイートの内容自体に対する興味と、投稿者に関する興味の 2 種類に分類できると考え、前者を興味単語プロファイルで表し、後者を興味ユーザプロファイルで表現する

4.2.1. 興味単語プロファイルの生成

文書中の単語の重要度を推定するために多くの手法が提案されている。代表的な推定手法としてTF-IDF法がある。TF-IDF法では、ある文書では出現頻度が高いが、他の文書に出現する頻度が低い単語は、その文書の特徴を強く表していると考え、文書中に出現する単語の重要度を、対象とする文書jに於ける単語iの出現頻度 $tf_{i,j}$ と、単語iを含む文書の出現頻度 df_{i} の逆数の積として求める。TF-IDF法による単語の重み付けは以下の式(1)で形式的に定義される。ある文書jにおける単語iが、その文書中で出現頻度が高く、かつ文書集合中で出現頻度の低いものならば、文書jに於ける単語iの重みtfidfi,jの値が高くなる。なお、iNは文書集合に含まれる文書数を表す。

$$tfidf_{i,j} = tf_{i,j} \times log\left(\frac{N}{df_i}\right)$$
 ... (1)

TF-IDF 法による重み付けは文書検索等、様々な分野で広く利用されている.しかし、Twitter において、個々のツイートを一つの文書としてとらえ、個々の文書に出現する単語の重要度をTF-IDF 法をそのまま適用して求めることは適切ではない.ツイートは 140 字以内であるため、一つの単一のツイート中に同じ単語が複数回出現することは稀である.し

たがって、ほとんどのツイートに於いて、そこに含まれる単語数は1となり、idf値が小さい単語が重要度が多くなる. つまり、出現頻度が少ない単語の重要度が高くなり、利用者の興味を正確に反映することが困難になると考えられる.

この問題を解決するために、ユーザが興味を示しているだろうお気に入りツイートやリツイートしたツイートに含まれる単語を、「ユーザが興味のある単語」とする方法が考えられる。しかし、全てのユーザが興味のあるコンテンツ全てにお気に入りやリツイートといった操作をする訳ではない。また、ソーシャルコンテンツには「重要だが一度読めば十分である」ものも多く含まれている。そのため、お気に入りツイートに於ける重みは、ユーザにより大きく異なる要素となる。そこで、我々はユーザのツイートの閲覧時間を考慮する。

一般的に、ユーザがツイート群を閲覧するとき、興味のあるツイートをより時間をかけて読むと考えられる。また、free の時に読むツイート群の方が、busy の時に読むツイート群よりも、ユーザはツイートの好き嫌いを選別していると考えられる。このことから我々は、free のときに読んだツイート群の中から、より時間をかけて読んだツイートを「興味のあるツイート」と考える。そして、TF-IDF 法を拡張した TF-IDF-RT 法を提案する。単語 i のツイートj に於ける重み $w_{i,j}$ の定義を式 2 に示す・

$$w_{i,j} = tfidf_{i,j} \times rt_j \qquad \cdots (2)$$

この式では、ツイートj中の単語iに対しTF-IDF法で抽出された重みに、ツイートjを読むのに要したと推定される時間 (滞留時間) rt_j を掛け合わせている。これにより、お気に入りやリツイートといった明示的な操作がなくても、ツイートとそこに含まれる単語への着目度を、時間という観点から考慮可能になる。

なお、今回対象とした Twitter では、一つの投稿は 140 字以内であり、1文書の文書長が短いため、滞留時間が閾値を超えた場合にツイート中の単語 $w_{i,j}$ に関しては $tf_{i,j}$ =1 とし、閾値未満未満のツイート中の単語に関しては $tf_{i,j}$ =0 とした。今回は、滞留時間の閾値は経験的に 2 秒と設定した.

いま、注目するユーザのタイムラインを構成するツイート集合を TL とすると、ユーザの興味単語プロファイル \mathbf{u}_{word} は、それぞれの単語 i について、 TL に含まれる全てのツイート j の重み $\mathbf{w}_{i,j}$ の和として定義する(式 3). \mathbf{u}_{word} は、ユーザがどの単語にどの程度価値を感じるかというプロファイルである。この際、数字・漢数字のみやひらがな一文字等の興味単語として意味を持たない語、記号、「www」や「RT」等の Twitterでよく見られる意味のない語はストップワードとした.

$$\boldsymbol{u_{word}} = \begin{pmatrix} \sum_{j \in \text{TL}} w_{1j} \\ \sum_{j \in \text{TL}} w_{2j} \\ \vdots \end{pmatrix} \cdots (3)$$

例として,著者が 10 日間システムを利用し作成された興 味単語プロファイルの中で,重みが大きい上位 15件の単語 を図 6 に示す. なお, ユーザが意識的に興味があると思っている単語と, 無意識的に興味を持っている単語との相違がある. 図 6 に示されているのは, ユーザが「意識的に選択している単語」(趣味やトレンドの単語)ではなく, ユーザが「どういった単語を用いるツイートを好むのか」という傾向である.

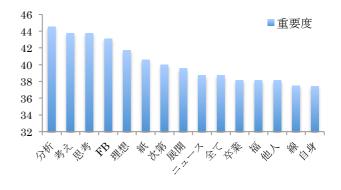


図 6 興味単語プロファイルの例

4.2.2. 興味ユーザプロファイルの作成

ユーザが興味を持つツイートの重要な要素として、「どのユーザが発信したツイートか」という、発信者の情報がある。多くの Twitter クライアントのインターフェースデザインでは、ツイートの左側に発言者のプロフィールアイコンが表示されるため、ユーザはアイコンでツイートを読むかどうかを判別することがある。ユーザが発信者の重みを抽出するために、発信者 k のツイートを、ユーザがどの程度時間をかけて読んでいるかという度合い $uval_k$ を用いる(式 4). ここで、ユーザがソーシャルストリームに表示した k のツイートの集合をTW(k)、ツイート J の滞留時間を rt_j 、ツイート J の総単語数をword(j)とする.

$$uval_k = \sum_{j \in TW(k)} \frac{rt_j}{word(j)}$$
 ... (4)

興味ユーザプロファイル u_{user} は、式(4)で求められた $uval_k$ を全てのユーザに関してベクトルとして表現したものである.

$$\mathbf{u}_{user} = \begin{pmatrix} uval_1 \\ uval_2 \\ \vdots \end{pmatrix} \cdots (5)$$

4.2.3. コンテンツ価値判定

興味単語プロファイル u_{word} と、興味ユーザプロファイル u_{user} を利用して、ユーザのコンテクストが busy のときに読み飛ばされたツイートの中から、見落としコンテンツを抽出する

興味単語プロファイル \mathbf{u}_{word} と、興味ユーザプロファイル \mathbf{u}_{user} が与えられたとき、ツイートtwに対するコンテンツ価値tw valuetw を以下の式で計算する. なお、tw はツイート tw と興味プロファイルとの類似度を表す関数であり、tw なである.

value(
$$tw$$
, u_{word} , u_{user}) = $sim(tw$, u_{word}) + $asim(tw$, u_{user})

読み飛ばされたツイートそれぞれに対してコンテンツ価値を計算し、その値が高いものを、見落としコンテンツと判定する.

例として、著者が 10 日間システムを利用し2つの興味プロファイルを作成し、後日機械的に収集した対象ツイート 773 件に、次の 3 パターンでプロファイルを適用し推薦された上位ツイートを表 $1\sim3$ に示す.なお、下線を引いた語はプロファイルに適合した単語である.

- ① 興味単語プロファイルと興味ユーザプロファイルを併用 して適用(表 1)
- ② 興味単語プロファイルのみ適用(表 2)
- ③ 興味ユーザプロファイルのみ適用(表 3)

表1の1位のユーザは興味ユーザプロファイルの値は高くはないが、興味単語プロファイルに該当する単語を多数含むツイートであった。2・3位は著者の興味ユーザプロファイルの値が比較的高い友人である。ツイート自体にも著者の最近の関心事が反映されており、「理想」、「部屋」、「卒」、「研」といった比較的上位の興味単語が含まれている。また、表1の3位のツイートは表2では上位にランクインしなかったものだが、興味ユーザプロファイルを併用することで上位に上がっている。実際に、著者は以前表1の3位のツイートにリプライを返しており、十分に興味ツイートであったと考えることができる。また、表3の興味ユーザプロファイルのみの推薦では、興味ユーザプロファイルで上位にあたるユーザの発言がランダムに入るため、必ずしも興味が引かれる内容ではなかった。

以上より、表 3 のように興味ユーザプロファイルのみを利用すると、「特定ユーザの発言は全て」が「良い」ツイートとして扱われてしまう.一方、表 2 のように興味単語プロファイルのみでは、文脈によらず興味単語が含まれてさえいえば上位にランクインしてしまうため、これも適切とは言えない.また、単語数の多いツイートが比較的上位に入っているという特徴が見られた.

表 1 ①興味単語・興味ユーザプロファイルを併用して適用し推薦されたツイート

順位	ユーザ	コンテンツ 価 値	ツイート
1	F	0.492	ー握りのエリートの後ろには <u>何</u> 千何万という <u>俺</u> 達のような <u>人間</u> がいるんだ!そして社会の歯車,使い捨て,操り人形などと言われながらも,自分の夢だったり家族や恋人の幸せだったり,…
2	G	0.483	【軽量鉄骨造】には防音性求めちゃ <u>駄</u> 目なのか.【RC造】がいいのか.【RC造 防音】でググっても、安心できん. 監禁 犯が、分厚い壁の角 <u>部屋</u> &隣空き <u>部</u> 屋に住んでて、誰も気づかなかった. て のがあった気が…
3	D	0.449	この <u>時期</u> にやることではございませんが、 <u>明日</u> こっそり <u>わたし</u> の卒研作品見にきてやってもいいぞっていう方いらっしゃいますかね。

表 2 ②興味単語プロファイルのみを適用し推薦されたツイート

			, , , , , , , , , , , , , , , , , , , ,			
順位	ユーザ	コンテンツ 価 値	ツイート			
1	F	0.157	一握りのエリートの後ろには <u>何</u> 千何万という <u>俺</u> 達のような <u>人間</u> がいるんだ!そして社会の歯車,使い捨て,操り人形などと言われながらも, <u>自分の夢</u> だったり <u>家族</u> や恋人の幸せだったり,…			
2	G	0.155	【軽量鉄骨造】には防音性求めちや <u>駄目</u> なのか.【RC 造】がいいのか.【RC 造 防音】でググっても、安心できん. 監禁犯が、分厚い壁の角部屋&隣空き部屋に住んでて、誰も気づかなかった. てのがあった気が…			
3	Н	0.142	RT @KometsubuKomeo: 『髑髏城の七 <u>人</u> 』テレビではまず中々出会えない熱量 だったと思います. 魂,って <u>やつ</u> ですか ね?生の <u>人間</u> (スクリーンやけど)が <u>今</u> そこ でまさに演じて造ってるって <u>感じ</u> がしまし た			

表 3 ③興味ユーザプロファイルを適用し推薦されたツイート

順位	ユーザ	コンテンツ 価 値	ツイート
1	A	1	"マンチェスターユナイテッド対リヴァプー ルのレビュー 並びに香川のプレースタイ ルの <u>話</u> - pal-9999の日記" http://t.co/Mp1sg7Rm
2	В	0.777	汚れた壁のフリーテクスチャ素材20 http://t.co/tFTn7AUv
3	В	0.777	レスポンシブ・ウェブデザインの流行パタ ーンまとめ http://t.co/HWAKD3Zb

5. 実験

5.1. 実験内容

提案手法のなかで、ツイートの滞留時間に基づいたユーザの興味プロファイルの生成手法及び、生成した興味プロファイルを用いたツイートの価値判定手法の有効性を評価するための実験を行った.被験者は、19~25歳の9人である.被験者には5日間、パソコンでTwitterを閲覧する際に、実験用システムを用いソーシャルストリーム(タイムライン)を閲覧してもらった.なお、被験者のコンテクストをfreeのときに限定するために、システムの利用は時間のある時のみに限定し、その他(数秒しか開かないとき、携帯端末での利用時等)では、普段利用している Twitter クライアントを用いるよう指示した.その他、ツイート、お気に入り、リツイート、リプライなどの操作もシステムから行ってもらった.

次に、実験で得られたデータを用い各被験者の興味単語プロファイルと興味ユーザプロファイルの作成を行った。実験最終日より7日経過した後に、新たにアンケート用に被験者それぞれのソーシャルストリームから700~800件のツイート(対象ツイート)を取得した。提案手法と既存手法を含めた5つの手法によって対象ツイート全てのコンテンツ価値を計算し、推薦ツ

イートを選出した. 今回対象とした5つの手法は以下の通りである.

- ① 興味単語プロファイルと興味ユーザプロファイルを併用した推薦手法
- ② 興味単語プロファイルのみを用いた推薦手法
- ③ 興味ユーザプロファイルのみを用いた推薦手法
- ④ TF-IDF 法を用いる推薦手法
- ⑤ TF法を用いる推薦手法

被験者には、推薦ツイートの中から同一ツイートを 除いてランダムな順番で提示し、「面白い」と感じるか どうかを、5段階で評価してもらった.

5.2. 実験結果

被験者 9 名の内, 2 名はデータ不足のため検証不可能とし除外した. 以下, 実験結果の検証は有効なデータ量が得られた 7 名のデータを元に進める.

7 名の被験者の評価を元に、各手法の有効性を比較するため、推薦ツイートの面白さについて有意な差があるかを t 検定を用いて調べた結果を図 7 および表 4 に示す、提案手法①と②の推薦結果への評価が最も値が高くなっている、提案手法①と既存手法④、提案手法①と既存手法④、提案手法①と既存手法④、提案手法②と既存手法④、提案手法②と既存手法④、提案手法②と既存手法⑥に於いて、有意水準 5%で有意差が見られた.

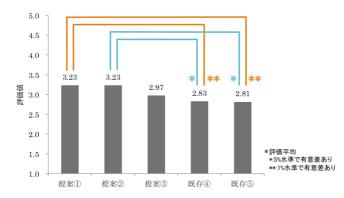


図 7 全被験者の各手法の評価

表 4 t 検定における p 値

	提案①	提案①	提案②	提案②	提案③	提案③
	と	と	と	と	と	と
	既存④	既存⑤	既存④	既存⑤	既存④	既存⑤
p 値	0.0088	0.0089	0.0101	0.011	0.2372	0.2047

6. 考察

6.1. 提案手法の有効性

上記の実験結果に於いて,「興味単語プロファイル と興味ユーザプロファイルを併用」した提案手法①, 「興味単語プロファイルのみを用いた」提案手法②と もに、既存手法と有意な差が見られた.しかし「興味 ユーザプロファイルのみを用いた」提案手法③ととえユーザが自己を見られなかった.これは、たとえコーザが普段注目する発信者であっても、発信者ではコーザにとって興味のあるツイートをする訳で自せる。一方、提案手法①の場合、普及に大きな示している。一方、提案手法②と既存手法の内推薦ができたと言えると既存手法の有意差、提案手法②と既存手法のが出なかったが、これはコンテインを表したが出なかったが、これはコンテータの値が小さないが出なかったが、これはコンデーをである。第5章式6)際の、興味ユーザがある。ユーザ毎にパラメータの値が変化すること検討することが必要である。

以上より、検討すべき点はあるものの、情報収集を 目的とするユーザに対し、滞留時間を用いた興味単語 の抽出とコンテンツ推薦は既存手法に対し有効である と考えられる.

6.2. 滞留時間の推定に於ける課題と改善案

本手法では, ユーザのスクロール操作に注目し, ど のツイートにどのくらいの時間をかけて読むか(滞留 時間)を元に興味単語の抽出を行った.しかしこの滞 留時間は, ユーザが「画面中心のツイートを読む」と いう前提であり、その誤差による影響は無視できない. 実験後,被験者数名に聞き取り調査を行ったところ, 被験者が注目していたツイートと滞留時間の長かった ツイートとに差が見られた者もいた. 今回は単純な補 正と窓関数での平滑化を行ったが、明らかに特定のツ イートへのアクションとして認識可能な「お気に入り」, 「リツイート」,「返信」などの操作時に、その対象ツ イートが画面のどの辺りに位置するかを抽出すること で,ユーザごとにソーシャルストリームを読む際の「目 線の位置」を最適化することが可能だと考える. 振る 舞いデータを学習させることで、ユーザごとの目線の 位置を最適化し、滞留時間の精度向上を目指したい.

また、滞留時間が長かったものの、「偶然見ただけで興味はなかった」との回答を得たツイートもあった。これは我々の「滞留時間が長いツイートは、ユーザが興味をもったツイート」という前提と異なるものであり、滞留時間の推定による興味の抽出の限界でもある。この点に関しては、興味プロファイル作成時に滞留時間の閾値を設けることで、「偶然読んだ」ツイートへの誤差を減らすことで改善が期待できる。また、システムの継続的な利用により、長期的に出現する興味単語が、「偶然読んだ」ツイートに含まれる単語の影響を低減させることが可能と考えられる。

6.3. 抽出された重要単語

既存手法の TF-IDF 法の問題点として, 4章で「出現頻度の高い語ほど重要度が低く評価される恐れがある」と述べた. そこで, 提案手法である興味単語プロファイルによる抽出単語と, TF-IDF 法による抽出単語の上位 10 件を表 5 に示す. 例として被験者 B と F の結果を用いた.

被験者によらず、既存手法では「年」「日」「さん」「今日」などの、ツイート内のデータを表すだけでツイート自体の重要単語とは言えない単語が抽出されているのに対し、提案手法ではツイートの内容をふまえたものが抽出された。被験者 B はアップル社の新製品や、IT 系の新技術に着目しているユーザであり、抽出単語は適当である。なお、被験者 B の閲覧ツイートと照らし合わせたところ、「時代」「型」「便利」などの単語は、新製品や事件のニュースツイートに頻繁に含まれる単語であり、被験者 B の好むツイートの「傾向」を表す単語と言える。

表!	5	提案手	法と既存	∈法 による抽	出単語の比較
----	---	-----	------	---------	--------

	ユーザ B				ユーザ F			
手法	興味単語 プロファイル		tfidf 法		興 味 単 語 プロファイル		tfidf 法	
	単語	重要 度	単語	重要 度	単語	重要度	単語	重要度
1	more	37.08	年	10.49	ひとり	40.77	笑	12.15
2	技術	37.08	Brand	7.71	it	40.77	人	11.99
3	時代	35.65	円	7.63	便利	40.77	今日	10.40
4	型	34.03	デイリー	7.26	とこ	39.34	日	10.07
5	便利	33.24	日	7.11	精度	34.33	中	10.02
6	目	31.59	さん	7.04	種類	34.08	年	8.33
7	CES	31.35	新	7.03	games	34.00	私	8.28
8	警戒	31.07	月	6.71	金	32.41	分	8.04
9	ダウ	31.07	株	6.65	プレゼン	32.34	何	7.98
10	アップル	30.26	*	6.61	Harden	31.13	ノミスギタ	7.75

7. まとめ

従来、蓄積された膨大な情報の中から、ユーザに対して価値が高い情報を検索・推薦する手法は活発に研究されてきた。しかし、ソーシャルメディアの爆発的な普及のために、飽和状態の供給され続ける情報の中から、ユーザに適切なコンテンツをいかに提示するかが重要な課題となっている。本研究では、「同じコンテンツでも、ユーザのコンテクストによって情報の価値は異なる」という仮説をもとに、読み飛ばしが日常的に行われているソーシャルストリームにおける、有用なコンテンツを有効的に活用する新しいアプローチを提案した。時間に余裕のない busy のときにソーシャルストリー

ム上で読み飛ばされたコンテンツの中から、ユーザにとって価値の高い「見落としコンテンツ」を抽出し、ユーザの適切なコンテクスト(free)において再提示することで、コンテンツ価値の向上を目指した。

本研究では、アイトラッカー等の特別な装置を用いずに、スクロールなどのユーザのブラウザ操作のみを元に、ソーシャルストリーム上のコンテンツへのユーザの着目度を「滞留時間」という指標で捉える手法を開発した。また、滞留時間を用いユーザの興味を表す興味単語プロファイル、興味ユーザプロファイルの作成を行った。それに基づいて、読み飛ばされたソーシャルコンテンツ中から価値の高いコンテンツを推薦する手法を開発した。被験者実験により、この推薦手法は既存手法に対し有効であると分かった。滞留時間による興味単語抽出手法は、ソーシャルストリームだけでなく、スマートフォン端末におけるニュース配信サイトやeコマースサイトに於いても活用できると考えられる。

コンテンツが流れて行くことが前提であるソーシャルストリームでは、重要度の高いコンテンツの推薦精度の向上よりもむしろ、ユーザにとって無意識的に(自然な形で)受け取ることが可能なコンテンツ提示のデザインが必要であると考えている。今後は、ユーザにより受け入れやすい形でコンテンツを提供するきっかけをつくるべく、開発手法に基づきインタフェースの完成を目指す。

参考文献

¹ Michael S. Bernstein, Bongwon Suh, Lichan Hong, Jilin Chen, Sanjay Kairam, Ed H. Chi, "Eddi: Interactive Topic-based Browsing of Social Status Streams", UIST '10 Proceedings of the 23nd annual ACM symposium on User interface software and technology, Pages 303-312, 2010

 3 松尾豊,福田隼人,石塚満,"ユーザ個人の閲覧履歴からのキーワード抽出によるブラウジング支援",電子情報通信学会技術研究報告. PRMU,パターン認識・メディア理解 101(712),85-92,2002-03-07

⁴中野利彦, 西本一志, "オンタイムの興味情報をオフタイムに埋め込むアイディア・インキュベーション支援システム", 情報処理学会研究報告. GN, [グループウェアとネットワークサービス] 2007(32), 1-6, 2007-03-22

 5 林春男, 佐藤翔輔, "膨大な情報から必要とされる情報を報せるビジネスツールとしての TRENDREADER", 情報管理 54(1), 2-12, 2011

⁶梅本和俊,山本岳洋,中村聡史,田中克己,"ユーザの視線を利用した検索意図推定とそれに基づく情報探索支援",日本データベース学会論文誌 Vol.10,一般論文(Research Papers),61-66,2011-06

² http://whatisedgerank.com/