

# 修 士 論 文

マイクロブログにおける  
インタラクションと投稿内容に着目した  
ユーザ推薦に関する研究

Research on User Recommendation in  
Microblogs Based on Users' Interactions  
and Contents

指導教員 豊田 正史 准教授



東京大学 情報理工学系研究科  
電子情報学専攻

氏 名 48-126409 岡本 大輝

提 出 日 平成26年2月6日

## 概要

ソーシャルネットワークサービス等では、ユーザ推薦によって、他ユーザと効率的に、そして円滑に、コミュニケーションを取ることが可能になる。ユーザ推薦には、ソーシャルネットワークのグラフ構造に着目した手法や、投稿内容に基づく手法等が提案されているが、これらはソーシャルメディアの環境によって様々な併用方法が考えられる。本論文では、マイクロブログサービスの1つである Twitter 上において、ユーザ推薦の実験を行う。実験にあたり、Twitter が提供するリプライやリツイートといったインタラクション機能を元に、ソーシャルグラフを構築する。さらに、インタラクションの種類や回数、ユーザ同士の投稿内容の類似度といった特徴を活用して、Random Walk 手法を提案し、その精度を検証する。

# 目次

<b>第 1 章</b>	<b>はじめに</b>	<b>1</b>
1.1	ソーシャルネットワークサービスの普及と課題	1
1.2	本研究の目的	1
1.3	本論文の構成	2
<b>第 2 章</b>	<b>関連研究</b>	<b>3</b>
2.1	特徴量に着目したユーザ推薦	3
2.1.1	ユーザ特徴に着目した手法	3
2.1.2	ユーザ推薦における問題点	4
2.2	ネットワーク構造を用いたユーザ推薦	4
2.2.1	グラフ上の近縁性を表す指標	5
2.2.2	Random Walk	5
2.2.3	ネットワークに属性を付与したユーザ推薦	6
<b>第 3 章</b>	<b>マイクロブログにおけるユーザ推薦</b>	<b>7</b>
3.1	マイクロブログの機能	7
3.2	ユーザ推薦における投稿内容の活用	8
<b>第 4 章</b>	<b>提案手法</b>	<b>10</b>
4.1	Weighted PageRank	10
4.1.1	インタラクション回数を用いたウエイト	11
4.1.2	投稿内容の類似度を用いたウエイト	11
4.1.3	複合ウエイト	12
4.2	データセット	12

---

4.2.1	グラフの構築	12
4.2.2	ユーザの分類	13
4.2.3	$tf \cdot idf$ ベクトルの作成	13
4.2.4	グラフとユーザの分析	14
<b>第 5 章</b>	<b>実験と結果の考察</b>	<b>21</b>
5.1	推薦の精度の評価方法	21
5.2	実験結果	22
5.3	実験結果の考察と検証	27
<b>第 6 章</b>	<b>おわりに</b>	<b>34</b>
	<b>謝辞</b>	<b>36</b>
	<b>参考文献</b>	<b>37</b>
	<b>発表文献</b>	<b>40</b>

# 目 次

4.1	対象ユーザと候補ユーザの経路距離とユーザ分布 . . . . .	18
4.2	候補ユーザへの経路上ユーザのインタラクション回数 . . . . .	19
4.3	候補ユーザへの経路上の投稿内容類似度 . . . . .	20
5.1	リプライによるエッジを正例とした際の精度 . . . . .	25
5.2	リツイートによるエッジを正例とした際の精度 . . . . .	26
5.3	Mean Average Precision の差によるユーザ数 . . . . .	31
5.4	各ユーザ群の発した名詞の種類と回数の特徴 . . . . .	33

# 第1章 はじめに

## 1.1 ソーシャルネットワークサービスの普及と課題

近年，社会生活における web の重要性が増し，人々は web 上でも社会活動や経済活動を行うようになった．そうした中で，web ユーザに対する情報推薦は，web サービスにおいて非常に重要な技術となっており，その精度を高める研究は数多くなされている．例えばユーザに商品を推薦することで，購買行動が活発になる．また，ユーザにユーザを推薦することで，ユーザ間にコネクションが生じ，サービスの利用率を高めることが可能となる．

一方で，web 自体の肥大化と多様化に伴い，web ユーザの数が増加し，web の使い方も様々に変化・多様化している．そうした中で，推薦の手法もまた，様々なサービスやユーザの目的に適応する必要がある．このような環境の中で，これまで盛んに研究されてきた，ソーシャルネットワークのグラフ構造に着目した手法に加え，グラフに依拠しない様々な指標を併用し，状況や目的に適した推薦手法の提案が必要とされる．

## 1.2 本研究の目的

本稿では，ソーシャルネットワークにおけるユーザ推薦を目的とし，マイクロブログサービスである Twitter の API を利用し，そのデータを用いた実験を行う．実験においては，Twitter のグラフ構造と，グラフに因らない特徴として，ユーザが投稿するテキストの内容の両方に着目し，それらを併用する．また，実験の結果や使用したデータを分析し，ユーザがサービスを利用する際の特徴や，グラフ構造と

### 1.3. 本論文の構成

---

の相関性を議論する。そして、ユーザの用途や特徴によって、最適な推薦手法にどのような違いが現れるのか、その考察を行う。

## 1.3 本論文の構成

本論文の構成は次のとおりである。

**第2章** ユーザ同士の関係から構築されるグラフに着目したユーザ推薦手法や、ユーザの特徴に基づいたユーザ推薦手法について、先行研究の紹介を行う。

**第3章** Twitterを始めとしたマイクロブログサービスについて説明し、今回の研究の意義と目的について述べる。

**第4章** インタラクションに応じた最適な推薦手法の分析を目的とし、今回実施した実験の詳細と提案手法について述べる。

**第5章** 実験の評価方法と結果を述べ、投稿内容類似度の作用について考察する。

**第6章** 全体のまとめと今後の課題について述べる。

## 第2章 関連研究

近年、Web 上においても、様々なコンテンツやユーザを、効率的に、かつ的確に推薦する手法について、多種多様な研究が行われている。本章では、まず、広く推薦に用いられる手法として、コンテンツに着目した手法について論じる。次に、SNS におけるユーザ推薦をリンク予測問題へと帰着させることを示し、ソーシャルグラフに着目した手法について論じる。

### 2.1 特徴量に着目したユーザ推薦

広く一般的な推薦手法の一つとして、ユーザの特徴量に着目した手法が存在する。推薦の候補 (ユーザ推薦の場合はユーザ) の特徴を元に、推薦の対象たるユーザとの親和性の高い候補を推薦するものである。

#### 2.1.1 ユーザ特徴に着目した手法

ユーザの特徴に着目し、機械学習等によってユーザ同士の接触を予測するアプローチがある。

Hannon らは、Twitter におけるユーザ推薦において、ユーザのツイート内容と、フォロワー/フォロー者のリストを利用した content based の手法を提案した。彼らは、ユーザのツイートに含まれる名詞のベクトルと、ユーザのフォロワー/フォロー者のリストを、それぞれユーザの特徴量と見做し、機械学習によって特徴量に重み付けを行い、類似するユーザを推薦した [10]。

## 2.2. ネットワーク構造を用いたユーザ推薦

---

Wangらはさらに、マイクロブログの機能とユーザのアクティビティに着目し、他ユーザの投稿の拡散回数や拡散力を特徴量として、ユーザへのツイート推薦を行った [17].

いずれの実験も結果的に、2-hop(隣接ユーザの隣接ユーザ)までは特徴量による作用が及ぶが、より遠方のユーザに対する特徴の作用に関しては言及していなかった.

### 2.1.2 ユーザ推薦における問題点

ユーザ推薦において、特徴量に着目する手法は一般的ではない。その理由を以下に挙げる.

ユーザ推薦に対し、アイテム推薦では、推薦する候補の特徴量が定義されている場合も多く、特徴量に従って様々な手法を適用することができる [6]. 一方ユーザ推薦、特に SNS においては、ユーザの特徴が定義されない、または少ない場合が多い [18]. 前項で示したいずれの手法も、各ユーザに対し適宜特徴量を定義していた [10] [17].

また、ユーザ推薦では推薦の対象と候補がともにユーザであり、莫大な数の組み合わせが予測される。2013 年末時点で、Twitter のアクティブユーザは 230,000,000 を超えており [8], 1 ユーザに対してその全てを一様な候補と捉えることは現実的ではない.

## 2.2 ネットワーク構造を用いたユーザ推薦

SNS におけるユーザ推薦を始めとして、特定の環境下のユーザに対して、同一のサービスを利用するユーザを推薦する際には、一般的にそのネットワーク構造に着目したアプローチが用いられる [9]. ネットワーク構造に着目する場合、ユーザ推薦はネットワーク内のリンク予測問題に帰着する.

以下では、データ構造の表現に則り、推薦候補および被推薦の対象となるオブジェクトをノード、オブジェクト同士のリンク関係をエッジと表現する。また、ノード  $A$  とエッジを形成しているノードの集合を  $\Gamma_{(A)}$  とする.

### 2.2.1 グラフ上の近縁性を表す指標

ユーザ推薦においては、グラフ構造を利用した様々な手法が考案されている。

最もシンプルな推薦手法として、2つのノードに共通する隣接ノードの数を数えて、数が多い順に推薦するというものがある [16] [14].

$$Score(A, B) = |\Gamma(A) \cap \Gamma(B)| \quad (2.1)$$

また、共通する隣接ノードの数に、分母として両者の隣接ノードの集合の論理和をとった Jaccard's coefficient という指標がある [16]. 推薦候補に隣接ノードが多い場合、無用にスコアが高くなることを防ぐ。

$$Score(A, B) = \frac{|\Gamma(A) \cap \Gamma(B)|}{|\Gamma(A) \cup \Gamma(B)|} \quad (2.2)$$

さらに、ノード A とノード B に共通する隣接ノードであるノード C の隣接ノードが少ないほど、C による A と B の近縁性が高まるとした Adamic Adar の評価指標がある [1]. Jaccard's coefficient と同じく、多数の隣接ノードを保有するノードによって、近縁性が不当に高まることを抑制することができる。

$$Score(A, B) = \sum_{\Gamma(A) \cap \Gamma(B)} \frac{1}{\log(\Gamma(C))} \quad (2.3)$$

推薦の対象となるノードから離れたノードも評価できる指標として、ノード A からノード B へ、距離  $l$  で移動する経路の集合  $paths_{A,B}^{(l)}$  を用いた Katz の指標がある [12].

$$Score(A, B) = \sum_{l=1}^{\infty} \beta^l \cdot |paths_{A,B}^{(l)}| \quad (2.4)$$

一般的に、 $\beta$  は 0 以上 1 未満の値をとり、 $\beta$  が小さいほど、A、B 間の短距離経路が重視される。

### 2.2.2 Random Walk

推薦される対象 (target) となるノードを  $s$  として、 $s$  の隣接ノードへとランダムで移動する動点の挙動を考える。動点はエッジに従ってさらに隣接ノードへの移動を

## 2.2. ネットワーク構造を用いたユーザ推薦

---

繰り返し、最終的に動点が存在する確率が高いノードを  $s$  に推薦する [4] [5].  $s$  と同じ距離に存在する複数の候補ノードの中でも、より  $s$  に近い対象を推薦する手法として、よく用いられている.

さらに Backstorm らは、facebook において友人関係をエッジとして、エッジの形成時間、インタラクションの回数、友人申請の方向などをエッジの特徴量とした. それらの特徴量を元に機械学習を行い、エッジにウエイトを付与してランダムウォークを行う Supervised Random Walk を提案した [2].

### 2.2.3 ネットワークに属性を付与したユーザ推薦

予測されるリンクに属性を付与し、ユーザの特徴や目的と属性を関連付けて推薦する手法は、Random Walk に限らず研究されている. Leskovec らは、ユーザ間のインタラクションに positive , negative の属性を付与し、属性を参照することで精度の向上に寄与することを示した [13]. また、神畠は推薦技術において、精度の向上とともに利用者の目的に合わせた推薦の必要性について言及している [20] [19] [21] [22].

Random Walk を用いた推薦手法は、ネットワーク構造を用いて近接性の高いユーザをランキングすることにより推薦を行うが、これまでに提案された手法では、ユーザの投稿内容までは考慮されていなかった. また、投稿内容を用いた推薦手法では、ネットワーク構造が十分に考慮されていなかった.

## 第3章 マイクロブログにおけるユーザ 推薦

本章では，マイクロブログの機能と特徴について説明し，本論文の実験の目的について述べる．

### 3.1 マイクロブログの機能

本論文で扱う Twitter は，マイクロブログと呼ばれるサービスの一種である．マイクロブログは，他のウェブサービスに無い独特な機能をいくつか備えている．本稿ではそれらの機能について説明する．

- 短文投稿

アメリカの Twitter，中国の weibo 等のマイクロブログサービスは，ユーザによる短文投稿を主な機能とする．投稿の内容，または投稿すること自体をツイートと表現し，ユーザはあらかじめ決められた文字数の中で，自由かつ手軽にテキストを投稿することができる．投稿内容は多岐にわたり，現在の状況や位置，情報の伝達，意見や気分の表明などにまで及ぶ．短文を投稿するという特徴から，投稿内容は簡潔で即時性の高いものが多く，ユーザの位置の推定や [7]，現実世界の出来事を認識する試みもある [15]．

- ユーザのフォロー

マイクロブログのユーザが，特定の他ユーザの投稿を容易に閲覧できるようになる機能である．Twitter を利用した研究には，他ユーザのフォローを推薦するものもある [10]．

### 3.2. ユーザ推薦における投稿内容の活用

---

- 投稿によるインタラクション

マイクロブログには、特定のユーザや投稿に対して特殊な投稿を行うことで、他ユーザとインタラクションを行う機能が存在する。本論文では、このインタラクションを用いたユーザ推薦を行う。

1. リプライ

投稿文内に、@マークの後にユーザ ID を伴う (@USERID) ことで、そのユーザに対して言及を行う投稿のことを、リプライ (reply) という。リプライを受け取ったユーザは、常時見ているユーザの投稿とは別に、自分へのリプライとしてその投稿に注目することが可能となる。リプライは、ある投稿に対して (その投稿者に対して) 行われる場合や、相手のユーザの投稿と無関係に行われる場合などがあるが、いずれのリプライ投稿にも明確な相手が存在する。

マイクロブログの短文投稿という特性から、リプライは同一ユーザ間で短時間に繰り返されることがある。

2. リツイート

リツイート (retweet) とは、他のユーザの投稿を自らの投稿と同様の扱いで再投稿・拡散する機能である。あるユーザが別のユーザの投稿をリツイートした際には、投稿主のユーザに対してリツイートしたユーザを通知する仕組みがある。リツイートは、ユーザが興味を持った話題や意見を、自身をフォローするユーザに拡散する場合に行われるほか、投稿主に対する対話や意思表示の手段として行われる場合もある [3] [11]。

## 3.2 ユーザ推薦における投稿内容の活用

マイクロブログは短文投稿を主な機能とするサービスである。それと同時に、投稿を用いて他のユーザとコミュニケーションを行うこともできる。他のユーザとのコミュニケーションが可能であるため、ユーザを推薦する技術も存在し得る。現在、Twitter では他ユーザを推薦する機能が実装・提供されている。

### 3.2. ユーザ推薦における投稿内容の活用

---

しかしながら、マイクロブログの主たる機能である短文の投稿を活用したユーザ推薦の研究は多くない。投稿内容の活用がマイクロブログニオケルユーザ推薦において有効であることは、Hannon らの研究 [10] でも示されているが、フォローユーザの推薦に留まり、限定的なものであった。

そこで本論文では、Twitter におけるユーザ推薦を目的とし、ソーシャルネットワークとユーザの投稿内容を併用したユーザ推薦の実験を行う。その後、実験の結果を検証し、投稿内容に基づいた推薦が有効に作用する状況や用途について考察を深める。

## 第4章 提案手法

本章では、マイクロブログにおいて投稿内容を活用したユーザ推薦の有用性を検証を目的として、今回行った実験の詳細を説明する。

### 4.1 Weighted PageRank

本論文では、Random Walk による推薦手法を基にして、投稿内容も考慮するように拡張した推薦手法を提案し、実験と検証を行う。

本手法では、推薦対象となるユーザ  $s$  を起点とし、他ユーザとのリンクをたどる Random Walk を行い、各ステップにおいて一定の確率で  $s$  に戻る、Random Walk with Restart 手法 [4] を基にする。また、遷移確率の決定にあたり、ユーザ間のインタラクション回数や投稿内容の類似度などを考慮して、遷移確率に重み付けを行う。遷移確率行列を  $M$  とすると、 $k$  回遷移後の存在確率のベクトル  $\vec{r}_k$  の更新式は、以下で表わされる。

$$\vec{r}_{k+1} = (1 - \alpha)M\vec{r}_k + \alpha\vec{r}_0 \quad (4.1)$$

ただし、 $\vec{r}_0$  は、ユーザ  $s$  のみ 1 で残りは 0 となるようなベクトルである。 $\alpha$  は、各ステップにおいて  $s$  に戻る確率を表し、今回の実験では Supervised Random Walk 手法 [2] において用いられていた 0.3 を値として用いた。

遷移確率の決定にあたり、インタラクションに係る 2 種類と、投稿内容の類似度に係る 3 種類の、計 5 種類の特徴量を用いる。それぞれから算出したウエイトと、適切に選んだ 2 種類の複合ウエイトを、各エッジのウエイトとして実験を行う。

### 4.1.1 インタラクション回数を用いたウエイト

ユーザ間のインタラクション回数は親密性を表すため、回数が多いほど高い確率で遷移するような重み付けを行う。ユーザ間で行われたインタラクションの回数を元に、全てのエッジにリプライ、リツイートに基づいた2種類の重みを設定する。

ユーザ  $u$ , ユーザ  $v$  間のリプライ回数を  $n_{reply(u,v)}$ , リツイート回数を  $n_{retweet(u,v)}$  として、ウエイト  $w_{i-reply(u,v)}$  と  $w_{i-retweet(u,v)}$  を以下のように定義する。

$$w_{i-reply(u,v)} = 1 - \frac{1}{1 + n_{reply(u,v)}} \quad (4.2)$$

$$w_{i-retweet(u,v)} = 1 - \frac{1}{1 + n_{retweet(u,v)}} \quad (4.3)$$

### 4.1.2 投稿内容の類似度を用いたウエイト

訓練期間でリプライまたはリツイートによってエッジが形成された2ユーザで、両ユーザの  $tf \cdot idf$  ベクトルのコサイン類似度を、投稿内容の類似度に基づくウエイト  $W_{t(u,v)}$  として定義する。

$$w_{t(u,v)} = \cos(\vec{t}_{(u)}, \vec{t}_{(v)}) \quad (4.4)$$

なお、インタラクションの有無によって作成された3種類の  $tf \cdot idf$  ベクトルを、それぞれ式4.4の  $\vec{t}_u, \vec{t}_v$  とすることで、以下の3種類のウエイトを生成する。

- 全ての投稿の類似度に基づいたウエイト  $w_{t-all(u,v)}$
- リプライ投稿の類似度に基づいたウエイト  $w_{t-reply(u,v)}$
- リツイート投稿の類似度に基づいたウエイト  $w_{t-retweet(u,v)}$

### 4.1.3 複合ウエイト

上記のインタラクションを用いたウエイト2種類、投稿内容の類似度を用いたウエイト3種類から、それぞれ1種類ずつを用いた複合ウエイトを作成し、推薦に利用する。インタラクションを用いたウエイトを  $w_{i(u,v)}$ 、投稿内容の類似度を用いたウエイトを  $w_{t(u,v)}$  とすると、遷移確率行列の要素  $w_{(u,v)}$  は、以下の式で定義される。

$$w_{(u,v)} = \begin{cases} (1 - \beta)w_{i(u,v)} + \beta w_{t(u,v)} \\ 0 \quad (u, v \text{ 間にエッジが無い場合}) \end{cases} \quad (4.5)$$

こうして定義された遷移確率行列において、あるユーザがエッジを形成するユーザに対して、これらの重みに比例した確率で遷移を行うよう正規化を行う。

なお、今回の実験では、複合でない5種類のウエイトを用いた実験を行い、その結果を元に、 $(W_i, W_t)$  の組み合わせを決定する。

## 4.2 データセット

Twitter 社が提供する API を用いて、ツイートを公開しているアカウントからツイートを収集した。著者の使用している Twitter アカウントを起点として、ソーシャルグラフ上で近傍に存在するユーザがフォローしているユーザを順次集めて、最終的に 250,000 ユーザの tweet データを取得した。その中で、特に 2013 年 8 月から同年 9 月まで (以下、訓練期間) のデータと、2013 年 10 月 (以下、テスト期間) のデータを抽出した。

### 4.2.1 グラフの構築

今回の実験は、マイクロブログサービスでのユーザ推薦の精度の向上と、着目するインタラクションの違いによる、投稿内容を活用した推薦の差異の分析を主目的とする。Twitter におけるユーザ同士のインタラクションは、いくつかの形態が存在するが、今回の実験では、リプライまたはリツイートのうち、いずれかが 2 ユーザ間に存在する場合に、エッジを形成するものとした。

## 4.2. データセット

---

### 4.2.2 ユーザの分類

構築したグラフをもとに、以下の条件を満たすユーザを推薦の候補ユーザとする。

- 自動投稿アカウントではない
- train 期間で 100 件以上のツイートを取得できた
- train 期間でデータセット内のユーザとインタラクションをしている

さらに、候補ユーザのうち、テスト期間において

- 新たに 10 以上のユーザとリプライを交わす
- 新たに 10 以上のユーザとリツイートを交わす

以上の 2 条件を満たすユーザを、積極的にインタラクションを行うユーザとして、ランダムに 1000 ユーザを抽出し、推薦の対象ユーザとした。

### 4.2.3 $tf \cdot idf$ ベクトルの作成

訓練期間において、各ユーザが投稿したツイートに形態素解析を行い、ユーザ毎に名詞セットを抽出した。リプライとリツイートは投稿によって行われることから、名詞セットにその特性を反映させるために、ツイートをインタラクションの有無で分類し、以下の 3 種類のそれぞれ名詞セットを抽出した。

- 全てのツイート内容から抽出した名詞セット
- リプライツイート内容から抽出した名詞セット
- リツイートの内容から抽出した名詞セット

## 4.2. データセット

---

また、各名詞セットから  $tf \cdot idf$  ベクトルを生成し、各ユーザーの特徴量とする。ユーザー  $v$  における名詞  $noun$  の  $tf \cdot idf_{(v,noun)}$  は、以下の式で表される。

$$tf \cdot idf_{(v,noun)} = \frac{count_{(v,noun)}}{\sum_N count_{(v,N)}} \cdot \log \frac{|U|}{|u : u \ni noun|} \quad (4.6)$$

なお、 $count_{(v,noun)}$  は、ユーザー  $v$  の一つの名詞セットに含まれる  $noun$  の数、 $|U|$  は候補ユーザー数、 $|u : u \ni noun|$  は名詞  $noun$  を発していた候補ユーザーの数である。これにより、各ユーザーごとに以下の3種類の  $tf \cdot idf$  ベクトルを生成した。

- 全てのツイートから作成した  $tf \cdot idf$  ベクトル  $\vec{t}_{all}$
- リプライツイートから作成した  $tf \cdot idf$  ベクトル  $\vec{t}_{reply}$
- リツイートから作成した  $tf \cdot idf$  ベクトル  $\vec{t}_{retweet}$

### 4.2.4 グラフとユーザーの分析

以上の条件で構築したグラフと、ユーザーの分析を行った。

#### 4.2.4.1 エッジに関する分析

図 4.1 に、訓練期間、テスト期間でそれぞれ構築したグラフのエッジ数を示す。実験においては、テスト期間に新たに生成されるエッジを正例エッジ、正例エッジで隣接することになるユーザーを正例ユーザーとして、推薦の精度の検証に用いる。

また、候補ユーザーと対象ユーザーのそれぞれで、インタラクションの種類別に隣接ユーザー数を算出した。その結果を表 4.2 に示す。

集計の結果から、

1. テスト期間のリプライの6割以上は訓練期間でインタラクションを行ったユーザー同士のものである

## 4.2. データセット

---

2. リツイートは交流の無かったユーザ同士でも多数行われている

という2点を述べることができる。また、選択した対象ユーザが平均より活発にインタラクションを行っていることもうかがえる。

## 4.2. データセット

---

表 4.1: 各期間のエッジ数

	訓練期間	テスト期間	テスト期間に新たに生成されるエッジ
ユーザ数	209,743		
リプライによるエッジ	2,569,214	1,450,987	454,361
リツイートによるエッジ	8,521,830	3,866,848	2,339,750
全てのエッジ(重複を除外)	10,134,638	4,889,299	2,552,946

表 4.2: ユーザ分類別度数

	候補ユーザ	対象ユーザ
ユーザ数	161927	1000
リプライによる平均隣接ユーザ	27.45	86.04
リツイートによる平均隣接ユーザ	94.16	156.16
平均隣接ユーザ	111.10	210.45

### 4.2.4.2 対象ユーザを起点とした経路探索

訓練期間のグラフ上で、1000人の対象ユーザを基点として経路探索を行い、候補ユーザまでの距離と、正例ユーザまでの距離を分析した(図4.1, 表4.3).

併せて、候補ユーザまでに至るまでの経路に存在するユーザの、リプライ回数とリツイート回数の平均値を計算した(図4.2).

さらに、対象ユーザから各候補ユーザに至る最短経路において、経路上の各エッジの両端に存在するユーザの  $tf \cdot idf$  ベクトルのコサイン類似度も計算した(4.3). ユーザ  $A$ , ユーザ  $B$  の  $tf \cdot idf$  ベクトルのコサイン類似度は、以下の式で表される.

$$\cos(\vec{t}_{(A)}, \vec{t}_{(B)}) = \frac{\vec{t}_{(A)} \cdot \vec{t}_{(B)}}{|\vec{t}_{(A)}| |\vec{t}_{(B)}|} \quad (4.7)$$

なお、分析に用いる  $tf \cdot idf$  ベクトルは、全てのツイートから作成したもの ( $\vec{t}_{all}$ ) とした. 2ユーザの  $tf \cdot idf$  ベクトルのコサイン類似度は、ユーザ同士の投稿内容の類似度を示す.

図4.1と表4.3から、今回構成したグラフでは3hopに候補ユーザのピークが存在するが、リプライ、リツイートのいずれの正例ユーザも、訓練期間においては約90また、正例ユーザへ至る経路上では、平均より多くのインタラクションが行われている傾向があり、特にリプライで顕著であることがわかる.

以上から、ユーザ推薦の実験をするにあたり、今回構築したグラフの構造と、ユーザ同士のインタラクションの回数を参照することの有効性が期待される.

図4.3から、正例ユーザへの経路では、平均より投稿内容類似度が高くなることがわかる. また、インタラクションの回数とは違い、正例の条件とするインタラクションの種類によって大きな差がなく、リツイートによる正例ユーザへの経路の平均類似度が、リプライによる正例ユーザへの経路より僅かに高くなった.

これにより、ユーザ推薦において投稿内容の類似度が高いユーザを推薦することで、精度が向上する可能性があるかと述べることができる.

## 4.2. データセット

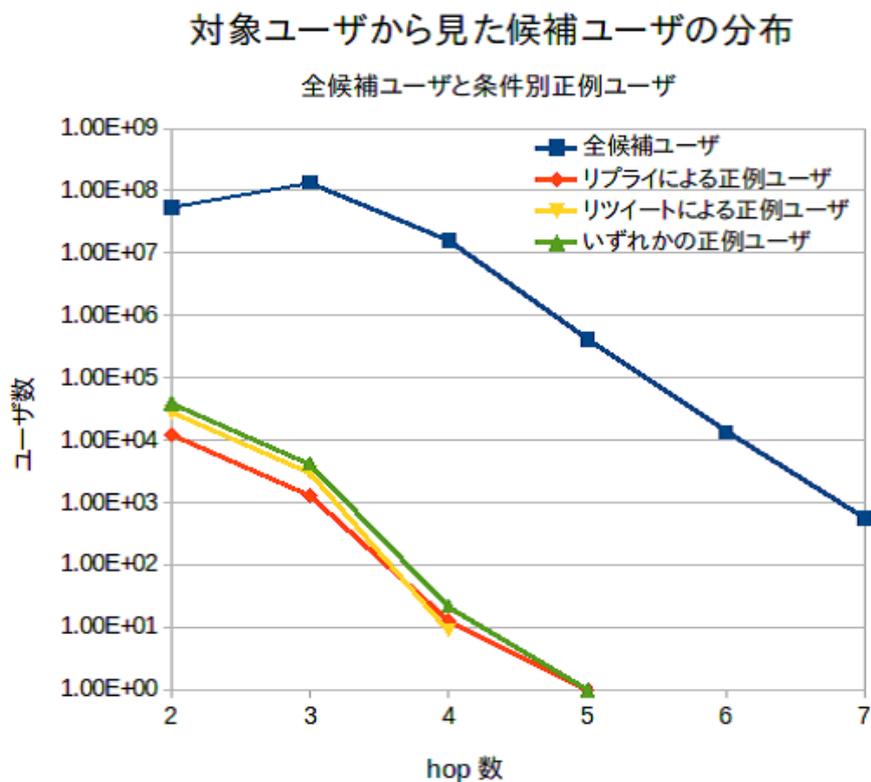


図 4.1: 対象ユーザと候補ユーザの経路距離とユーザ分布

表 4.3: 対象ユーザを起点とした候補ユーザの分布

	hop数による累積率 (%)				平均総数
	2	3	4	5	
全候補ユーザ	24.343	92.206	99.795	99.993	209,742
リプライによる 正例ユーザ	90.239	99.895	99.993	100	13.391
リツイートによる 正例ユーザ	90.161	99.971	100	100	31.426
全ての 正例ユーザ	89.999	99.946	99.998	100	42.428

## 4.2. データセット

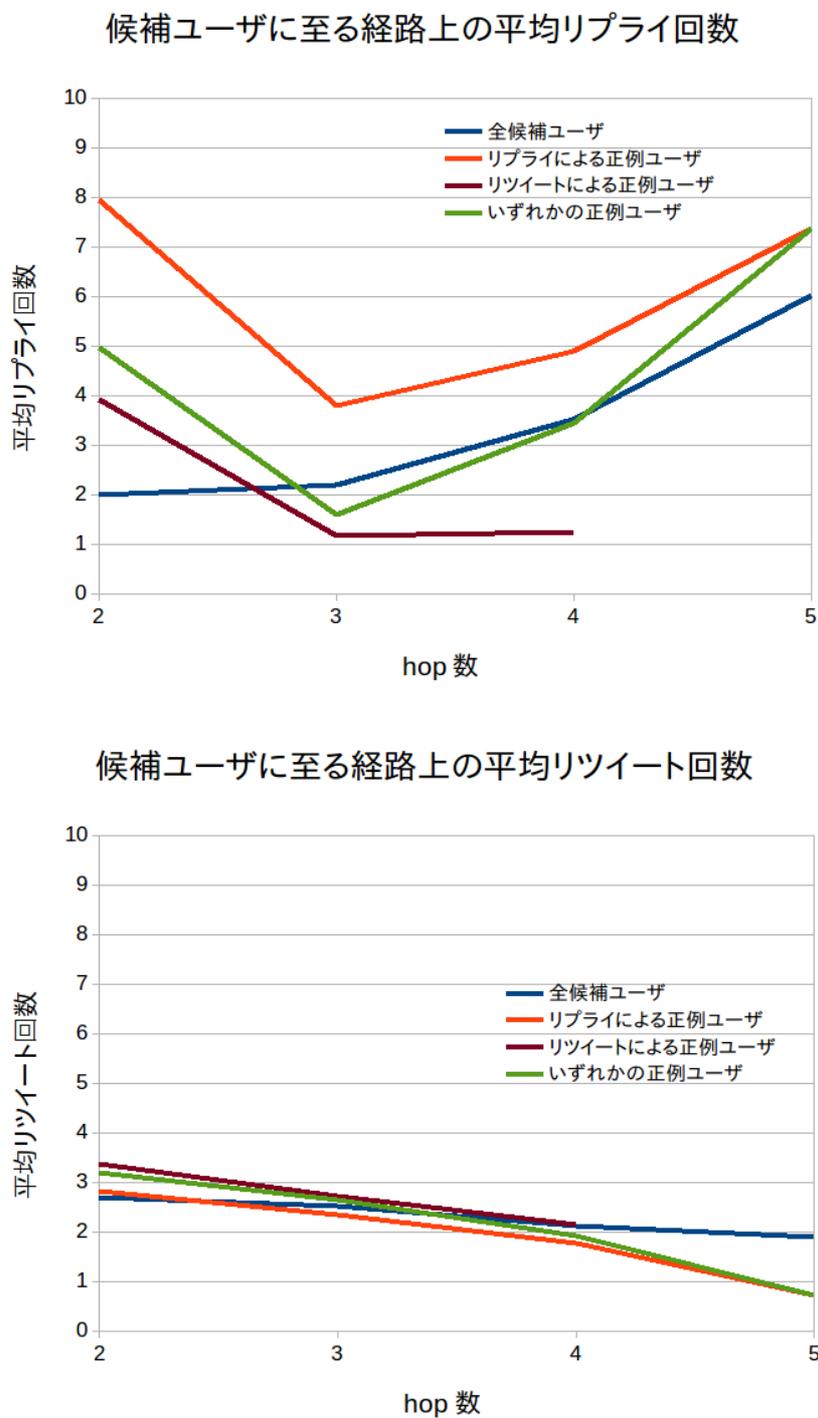


図 4.2: 候補ユーザへの経路上ユーザのインタラクション回数

候補ユーザに至る経路上の投稿内容類似度平均

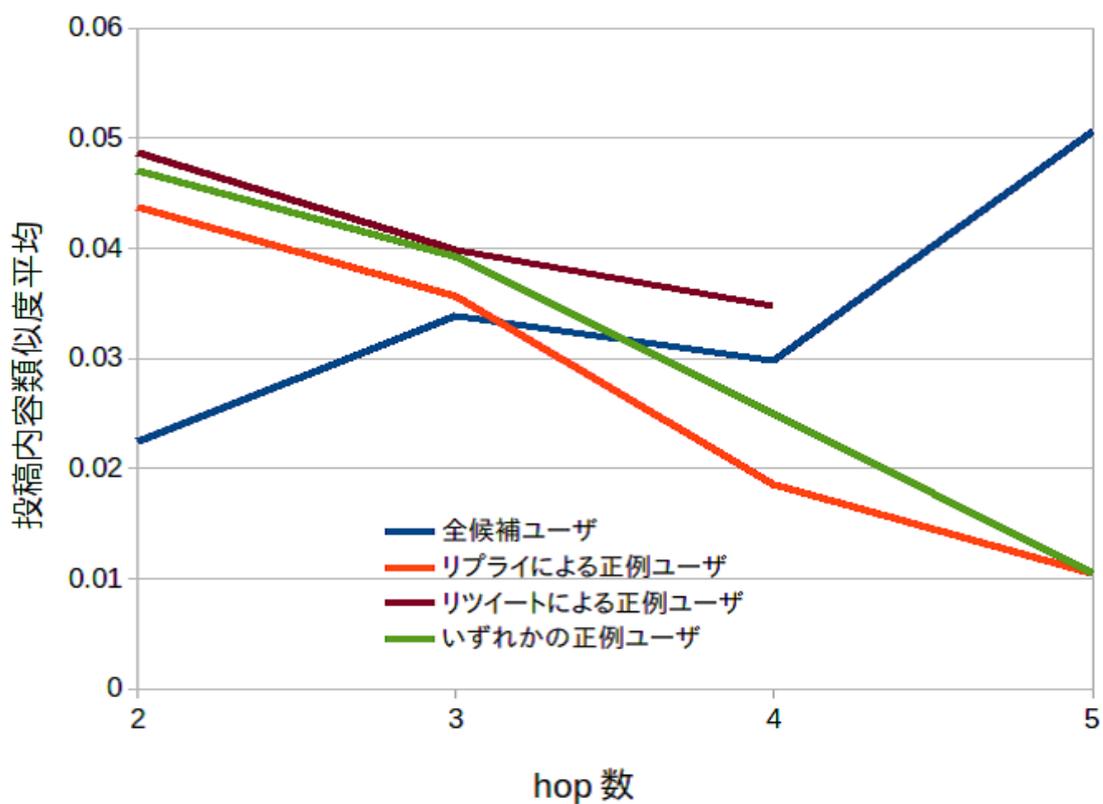


図 4.3: 候補ユーザへの経路上の投稿内容類似度

## 第5章 実験と結果の考察

### 5.1 推薦の精度の評価方法

本稿では、Twitter のユーザに対し、その目的に応じて最適なユーザ推薦を行うことを目的とする。グラフを活用したユーザ推薦では、テスト期間のデータから作成されたグラフで新たに得られるエッジを正例エッジとして、手法に基づく予測と比較して精度を検証することが一般的であるが、今回の実験では、正例の条件を

1. リプライによるもの
2. リツイートによるもの
3. それらのいずれか

の3パターンで変え、各ウエイトによる精度の変化と最適な推薦手法を検証する。

エッジに加える様々なウエイトを、それぞれ変化させながら、Prec@10, Mean Average Precision の2つの指標を用いて、推薦の精度を評価する。また、Mean Average Precision の説明の過程で、Average Precision にも言及する。

- Precision@10 (Prec@10)  
1 ユーザに対する推薦候補として、スコアが高くなった上位 10 ユーザの中の精度。今回のデータでは、推薦のターゲットユーザがテスト期間に新たに 10 以上のエッジを形成するため、推薦の精度の評価としてこの値を使用する。
- Average Precision (AP)  
1 ユーザに対する推薦候補のスコアランキングを上位から検証し、正例となる候補が出現した時点での精度の平均値。

## 5.2. 実験結果

---

$$AP_{(N_{TRUE})} = \frac{1}{N_{TRUE}} \sum_{K=1}^{N_{TRUE}} \frac{K}{\text{rank of } K\text{th } TRUE} \quad (5.1)$$

但し,  $N_{TRUE}$  は1ユーザが有する全ての正例の数である.

- Mean Average Precision (MAP)

推薦候補のスコアランキングを上位から検証し, 正例となる候補が出現した時点での Average Precision の平均値. ランキングの上位に正例が集中していると高い値となる.

$$MAP = \frac{1}{N_{TRUE}} \sum_{K=1}^{N_{TRUE}} AP_{(K)} \quad (5.2)$$

## 5.2 実験結果

本項ではまず, 前章で定めた5つのウエイトについて, それぞれ単独で用いて推薦を行った場合の精度を評価する. また, 比較の為に全てのエッジのウエイトを均一としたウエイト ( $w_{baseline}$ ) を用いた推薦の結果も記載する.

表5.1, 表5.2から, 正例となるインタラクションの回数を利用したウエイトを用いることで, ベースラインより良い精度となることがわかる. また, 投稿内容の類似度を用いた結果に着目すると, いずれの正例条件であっても全ての投稿内容を用いたウエイトで精度が高くなった.

## 5.2. 実験結果

---

表 5.1: 各ウエイト単独による推薦 (Precision@10)

	正例の条件		
	リプライ	リツイート	いずれかの インタラクション
$W_{i-reply}$	<b>0.0819</b>	0.0490	<b>0.1108</b>
$W_{i-retweet}$	0.0233	<b>0.0785</b>	0.0940
$W_{t-all}$	0.0484	0.0739	0.1083
$W_{t-reply}$	0.0472	0.0539	0.0878
$W_{t-retweet}$	0.0375	0.0719	0.0979
$W_{baseline}$	0.0537	0.0624	0.1003

表 5.2: 各ウエイト単独による推薦 (Mean Average Precision)

	正例の条件		
	リプライ	リツイート	いずれかの インタラクション
$W_{i-reply}$	<b>0.1093</b>	0.0444	0.0887
$W_{i-retweet}$	0.0347	<b>0.0694</b>	0.0740
$W_{t-all}$	0.0663	0.0659	0.0873
$W_{t-reply}$	0.0624	0.0502	0.0726
$W_{t-retweet}$	0.0511	0.0635	0.0779
$W_{baseline}$	0.0766	0.0616	<b>0.0892</b>

## 5.2. 実験結果

---

実験の結果から、正例となるエッジの条件別に2種類の複合ウエイトを提案し、式4.5の $\beta$ の値を変えながら精度を検証した。正例とするエッジの条件、および組み合わせたウエイトの種類は以下のとおりである。

- リプライによるエッジを正例とする場合  
リプライの回数に基づくウエイト ( $W_{i-reply}$ ) と、全ての投稿内容に基づくウエイト ( $W_{t-all}$ ) を用いる。
- リツイートによるエッジを正例とする場合  
リツイートの回数に基づくウエイト ( $W_{i-retweet}$ ) と、全ての投稿内容に基づくウエイト ( $W_{t-all}$ ) を用いる。

以下、図5.1にリプライによるエッジを正例とした際の実験結果を、図5.2にリツイートによるエッジを正例とした際の実験結果を、それぞれ示す。いずれも、 $\beta$ が大きいほど投稿内容の類似度の重要性が高まる。

## 5.2. 実験結果

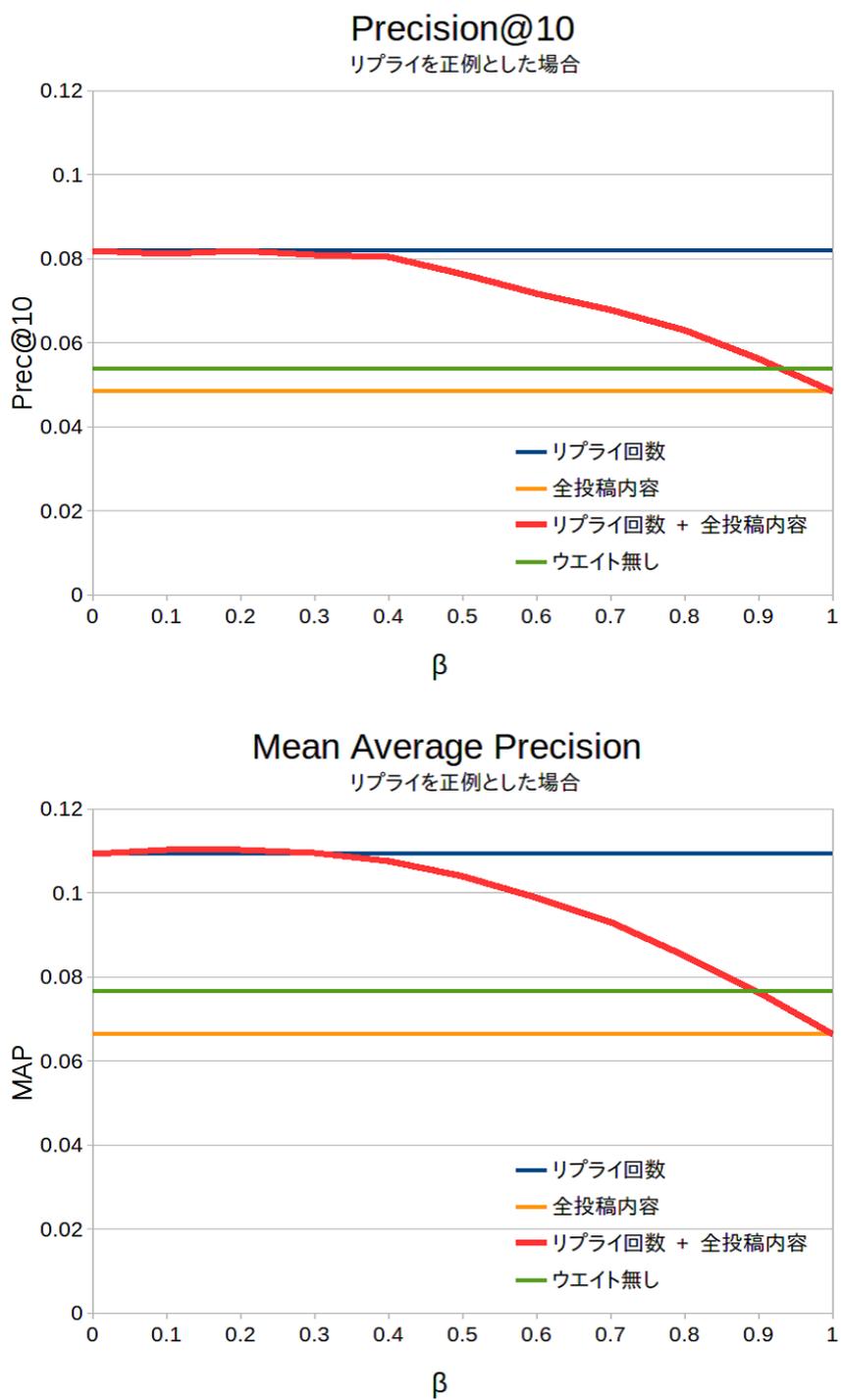


図 5.1: リプライによるエッジを正例とした際の精度

## 5.2. 実験結果

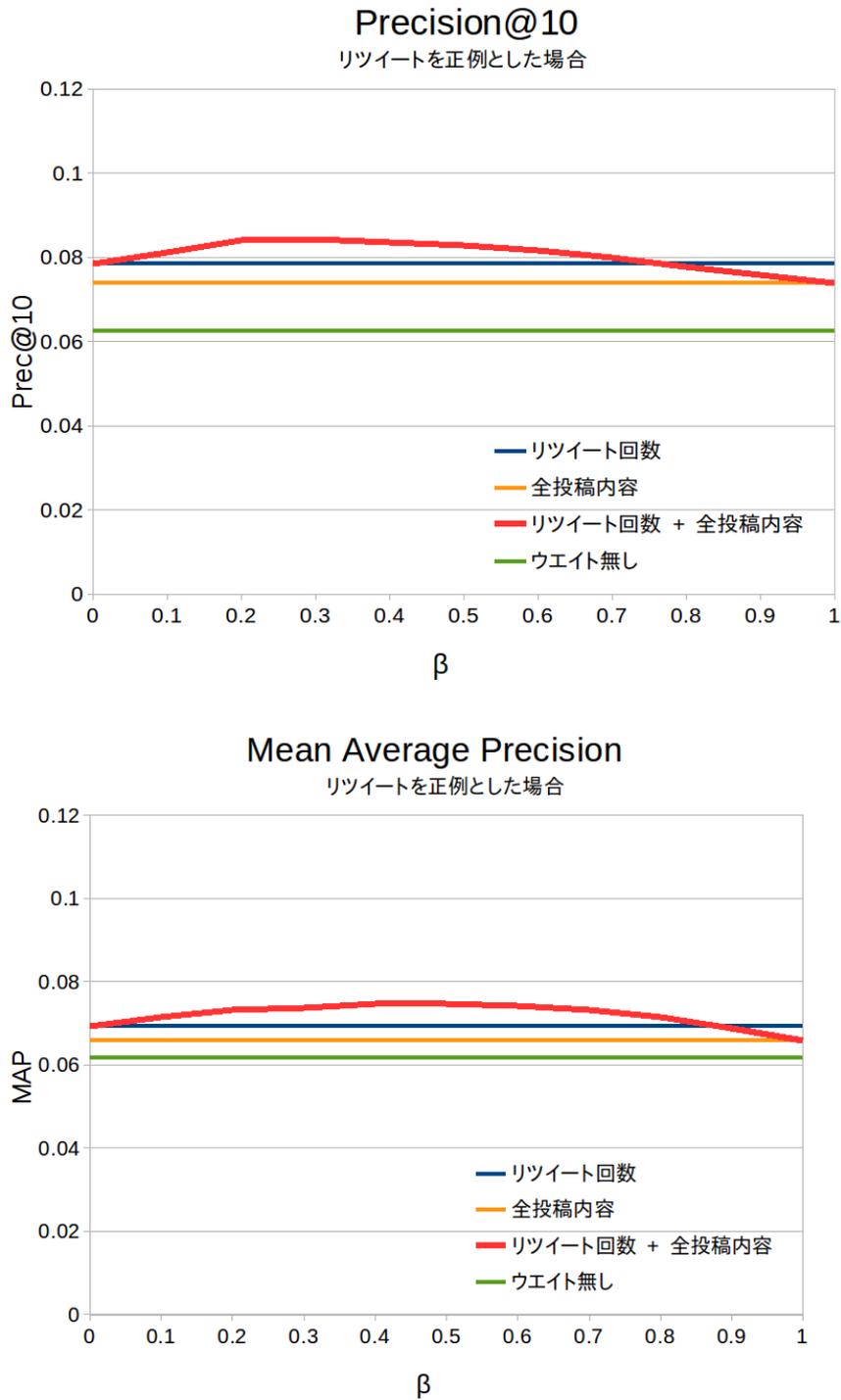


図 5.2: リツイートによるエッジを正例とした際の精度

## 5.3 実験結果の考察と検証

本項では、前項までの実験結果と 4.2.4 で行った分析結果を元に、各ウエイトの特性や作用、および Twitter ユーザの特徴についての考察を述べる。

- インタラクション回数によるウエイト

表 5.1, 表 5.2 から、正例とするインタラクションの回数をウエイトとして用いた場合に、精度が向上することがわかる。一方で、正例ではないインタラクションのウエイトを用いても、それぞれ精度が低下してしまっている。

これにより、Twitter のユーザは、リプライとリツイートを明確に違う機能として活用し、頻繁にリプライを交わす相手のツイートを頻繁にリツイートするとは限らないと述べることができる。

また、表 5.1 で、リプライ、リツイートをそれぞれ正例エッジの条件とした場合の比較を行う。Precision@10 は一般的に、候補中に占める正例の数が多いほど高くなる傾向がある。表 4.3 によると、リツイートによる正例ユーザの数はリプライによる正例ユーザの 2 倍以上存在する。しかしながら、均一なウエイト ( $W_{baseline}$ ) を用いた場合には、リツイートを正例とした場合の精度の比は 1.2 未満であり、リプライ回数によるウエイトでリプライを予測する場合と、リツイート回数によるウエイトでリツイートを予測する場合を比べると、リプライを用いた場合が高精度となっている。

このことから、リプライ関係によるユーザの親近性は、リプライ回数に顕著に反映され、将来のリプライを予測する際には、既知のリプライ回数を参照することが有効であると述べることができる。

- 投稿内容類似度によるウエイト

ツイート内容の類似度による推薦の結果にも着目する。正例とするインタラクションによらず、ツイートを分類せずに全てのツイート内容の類似度を用いた場合に、精度が高くなることがわかる。ここで、ベースラインとの比較を試みる。リプライによるエッジを正例とする場合には、どの類似度を用いても

### 5.3. 実験結果の考察と検証

---

ベースラインより精度が低下した。一方で、リツイートによるエッジを正例とする場合には、全てのツイート内容、リツイートの内容のいずれの類似度を用いても、ベースラインより精度が向上している。図 4.3 によると、リツイートによるエッジで正例となるユーザへの経路では、リプライによる正例となるユーザへの経路より、隣接ユーザ同士の類似度が高い傾向があった。

以上のことから、Twitter のユーザは、他のユーザからリツイートを行う際には投稿内容を加味しており、その傾向はリプライよりも顕著であると推察できる。

- 複合ウエイト

まず、図 5.1 に着目する。いずれもリプライ回数によるウエイトと投稿内容類似度によるウエイトを結合した場合の結果を示しているが、リプライを正例として予測する場合には、投稿内容による改善が見られないことがわかる。一方、図 5.2 に着目する。リツイート回数と投稿内容類似度によるウエイトを結合し、リツイートを正例として予測した際の精度であるが、両者を単体で用いた場合より改善していることがわかる。

ここで、複合ウエイトを用いてリツイートを予測する際に見られた改善の有意性を検証するために、Wilcoxon の符号付順位和検定と、4 分割交差検定を実施する。

1. Wilcoxon の符号付順位和検定

各候補ユーザに対し、

- リツイート回数によるウエイトで推薦を行った場合
- 複合ウエイト ( $\beta = 0.4$ ) で推薦を行った場合

以上の 2 回の推薦で得られた Mean Average Precision の、差の絶対値の順位と分布をもとに、有意性を評価する。

検定にあたり、Mean Average Precision の差による対象ユーザのヒストグラムを示す (図 5.3)。

### 5.3. 実験結果の考察と検証

---

検定の結果、検定統計量  $Z_0$  は 8.065 , 正規分布に従った有意確率は  $7.323e - 16$  となる。これは、「複合ウエイトを用いた推薦の精度はリツイートの回数によるウエイトを用いた推薦の精度と差はない」とする帰無仮説を棄却するに値する低確率である。

#### 2. 4分割交差検定

1000 人の対象ユーザをランダムに 4 分割し、3 つを訓練事例として、ウエイトの結合割合  $\beta$  を学習し、残る 1 つに適用する交差検定 4 回を行った。得られた結果を表 5.3 に示す。

Precision@10, Mean Average Precision のいずれでも、交差検定によって得られた平均値は、結合前のウエイトから得られた値より大きくなった。この検定においても、ウエイトの結合による精度の向上の有意性が示されている。

以上 2 つの検定から、リツイートを予測する場合に、リツイートの回数だけでなくユーザの投稿内容の類似度を用いると、精度が向上すると述べることができる。一方で、リプライを予測する場合の結果についても検定を行ったが、有意な精度改善は見られなかった。これにより、リプライを予測する際には、リプライの回数だけでなくユーザの投稿内容の類似度にも着目することで、精度が改善することが示された。

#### ● 精度の変化によるユーザ分類と分析

複合ウエイトを用いて、リプライ、リツイートのいずれを予測する実験でも、ユーザによって精度の向上・低下に違いが見受けられた。そこで、

1. インタラクション回数による推薦で高精度のユーザ
2. 投稿内容類似度による推薦で高精度のユーザ
3. 複合ウエイトによる推薦で高精度のユーザ

### 5.3. 実験結果の考察と検証

---

以上のように対象ユーザを3分類し、各ユーザ群の特徴を分析した。その結果を、表 5.4, 表 5.5 に示す。

まず、リプライを予測する場合(表 5.4)に着目し、各ユーザ群を比較する。すると、リプライによるエッジを有する数が多く、リプライの総数が少ないユーザほど、リプライ回数によるウェイトを用いることで精度が高くなるという傾向が見受けられる。これらのユーザは、多数のユーザを相手にリプライでインタラクションをするが、一人あたりのインタラクション回数は少ない、という特徴を示している。

次に、リツイートを予測する場合(表 5.5)に着目する。こちらは、ユーザ群同士でリツイートの回数に顕著な特徴は見受けられないが、発した名詞の種類と回数に大きな違いが見受けられる。それによると、特定の単語を何度も発するユーザは、リツイート回数に基づく推薦で高精度となり、広く多くの単語を発するユーザは、投稿内容類似度に基づく推薦で高精度となる。その様子を図 5.4 に示す。各ユーザの発した名詞の種類と回数に着目し、一人のユーザが発する名詞を出現回数順に並べて横軸とし、累積回数を縦軸として正規化を施した。曲線が左上に曲がり、曲線下部の面積が大きくなるほど、特定の単語を頻繁に発していることになる。これによると、リツイート回数を利用した推薦で高精度となったユーザは、特定の単語を集中して発しており、投稿内容類似度による推薦で高精度となったユーザは、単語の偏りが少ないことがわかる。

推薦の精度の結果を元に、分類したユーザの特徴を分析したが、ユーザ群によって特徴に違いがあることが分かった。これは即ち、推薦の前にユーザの特徴を分析し、インタラクション回数、投稿内容類似度の特徴量を、ユーザ毎に柔軟に用いることで、推薦の精度の向上が期待できるということを示唆している。

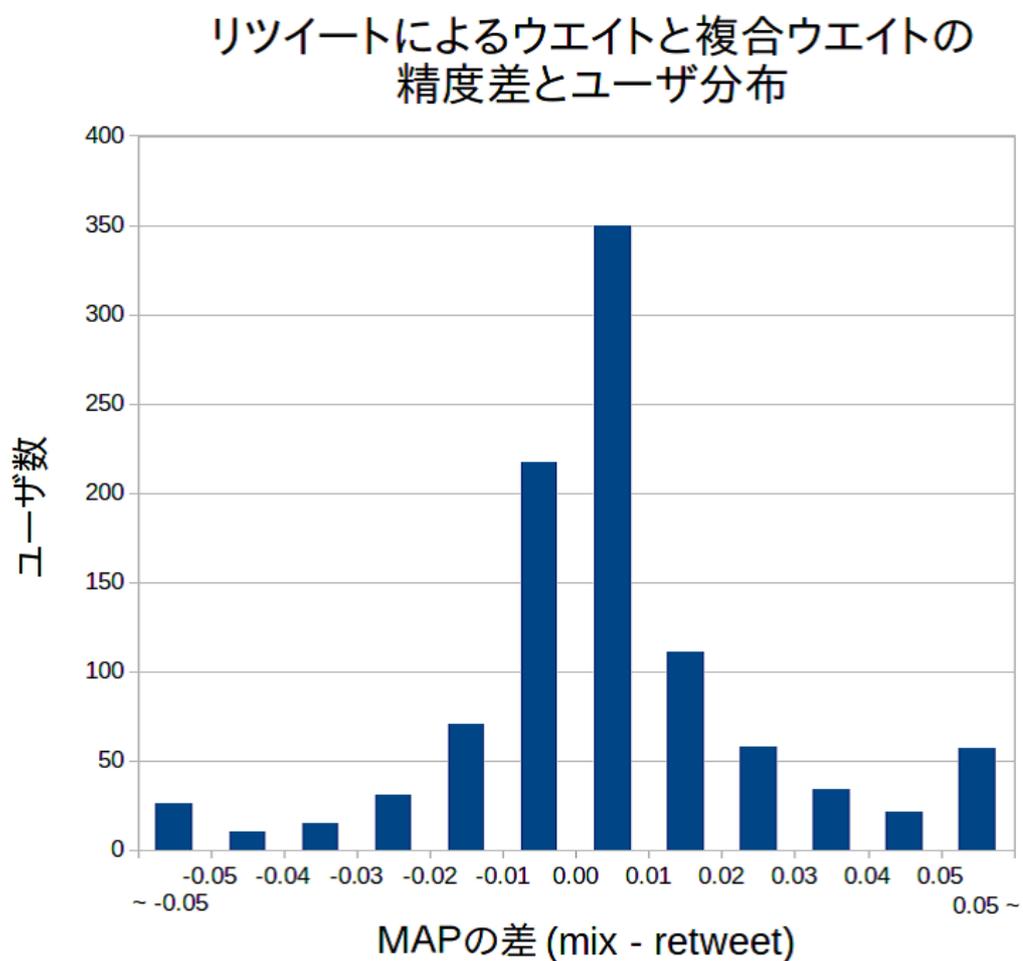


図 5.3: Mean Average Precision の差によるユーザ数

表 5.3: 4 分割交差検定による結果と比較

	precision@10	Mean Average Precision
$W_{i-retweet}$	0.0784	0.0693
$W_{t-all}$	0.0738	0.0659
$\beta W_{i-retweet} + (1 - \beta) W_{t-all}$ (4 分割交差検定)	<b>0.0839</b>	<b>0.0745</b>
$W_{baseline}$	0.0624	0.0616

### 5.3. 実験結果の考察と検証

---

表 5.4: 最適手法別ユーザ群の訓練期間の特徴(リプライを予測する場合)

高精度となった推薦手法で 用いたウエイト	リプライ回数	投稿内容類似度	複合ウエイト
平均リプライ回数	744.75	769.80	734.46
リプライによる 隣接ユーザ数の平均	78.51	72.45	75.77
名詞を発した回数の平均	6974.27	7987.09	7854.42
発した名詞の種類 の平均	736.33	788.59	732.12

表 5.5: 最適手法別ユーザ群の訓練期間の特徴(リツイートを予測する場合)

高精度となった推薦手法で 用いたウエイト	リツイート回数	投稿内容類似度	複合ウエイト
平均リツイート回数	253.56	268.22	272.90
リツイートによる 隣接ユーザ数の平均	131.15	140.52	137.02
名詞を発した回数の平均	8226.09	7480.52	8300.25
発した名詞の種類 の平均	731.37	731.03	741.50

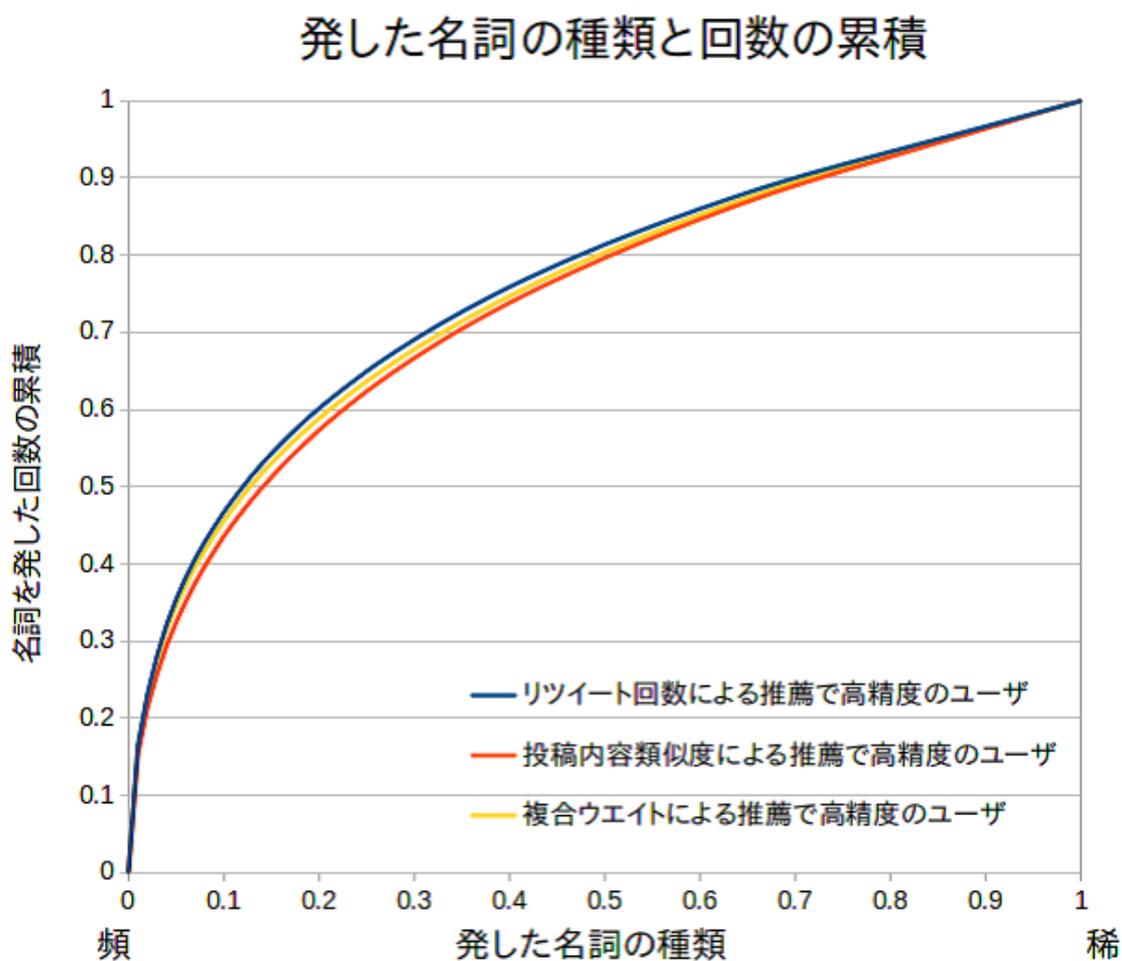


図 5.4: 各ユーザ群の発した名詞の種類と回数の特徴

## 第6章 おわりに

本論文では、ソーシャルネットワークにおけるインタラクションに基づくユーザ推薦の精度の向上を目的として、インタラクションの回数とユーザの投稿内容の類似度に着目した Random Walk に基づく推薦手法を提案した。この手法では、ユーザのインタラクションを元にソーシャルグラフを構築し、インタラクション回数と投稿内容の類似度をエッジのウエイトとして、ウエイトを遷移確率として Random Walk によるユーザ推薦を行った。また、推薦するユーザとの関係がリプライによるものか、リツイートによるものかを区別し、ウエイトに用いる特徴量が各条件下で推薦の精度に与える影響を検証した。

実験に際し、まず、構築したソーシャルグラフの特徴と、推薦の対象とするユーザを起点とした経路の分析を行った。その結果、将来インタラクションによる関係を結ぶユーザへの経路では、他の候補ユーザへの経路に比べ、インタラクション回数、投稿内容類似度のいずれの特徴量も高い値を示すことが示された。それを踏まえて、各特徴量から生成したウエイトを単独で用いた実験を行い、正例とするインタラクションの回数と、ユーザの全ての投稿内容の類似度を用いることが有用であることを示した。さらにそれを踏まえ、高精度となった正例条件と特徴量を元に、インタラクション回数と投稿内容類似度のそれぞれに基づくウエイトを結合し、結合の際の割合による精度の変化を検証した。その結果、リプライを正例として予測する際には、リプライの回数によるウエイトのみを利用し、投稿内容の類似度によるウエイトを加味することによる精度の向上が見られないことが判明した。一方で、リツイートを予測する際には、リツイートの回数によるウエイトのみではなく、投稿内容の類似度によるウエイトを併せることで、精度が向上した。また、Wilcoxon の符号付順位和検定を行うことで、精度向上が有意であることを示した。

---

以上のことから、Twitterにおいてユーザが他者とインタラクションを行う際には、インタラクションに応じた特徴があることが示された。まず、ユーザ同士のリプライ関係の強さは、その回数に顕著に表れ、リプライを予測する際には回数を活用することで、精度が向上する。一方、リツイートを予測する際には、その回数のみならず、ユーザの投稿内容の類似度を活用することで、精度が向上する。これにより、ユーザはTwitterにおいて、他者とリプライを交わす際には、相手のユーザとの既存の関係性が重視され、他者のツイートをリツイートする際には、既存の関係性に加えて投稿の内容が重視されている、と述べることができる。

さらに、ユーザの特徴量と、推薦に用いる最適な特徴量の相関が示唆されている。推薦に用いる特徴量を決定する際に、インタラクションの種類に加えて、ユーザの特徴量を参照することの有用性にも研究の余地がある。

今回の研究により、インタラクションに応じた特徴量の選択と活用が、よりユーザの需要に即した推薦が可能であることが示された。

# 謝辞

はじめに、指導教官の豊田正史准教授に厚く感謝いたします。豊田准教授には幾度となく貴重なご助言を賜り、ご指導していただきました。学徒として非常に未熟な私に、実験のアプローチに始まり、論文の執筆や、プレゼン資料の作成など、研究や発表に係る様々なノウハウを教えていただきました。しばしば私の不手際により、実験用器材に大きな負荷をかけてしまったり、准教授ご自身に多大な時間を割かせたりしてしまうことがありました。この場を借りてお詫びするとともに、改めて多大なる感謝の意を表します。

次に、国立情報学研究所長も兼務されてお忙しい中、折に触れ我々学生に貴重なご助言や叱咤激励を下された喜連川優教授に、深く感謝いたします。昨年紫綬褒章を受章され、学内に留まらず様々な機会に、喜連川教授のお名前を拝聴する機会がありました。私はその度に畏敬するとともに、喜連川教授の下で研究に励む自覚と自負を新たにし、研究に邁進することができました。ここに深く感謝するとともに、今後も変わらぬご活躍を祈念致します。

そして、転任された中野美由紀特任准教授、鍛冶伸裕特任准教授、吉永直樹特任准教授、伊藤正彦助教、横山大作助教を始めとした研究室スタッフの皆様には、毎週のミーティングで貴重なご助言を賜りましたことを感謝いたします。また、学生全般をサポートしていただいた秘書の方々にも、併せて感謝いたします。

最後に、なかなか連絡ができない私を心配し、遠く福岡から支えて下さった両親に感謝いたします。

2014年2月6日

## 参考文献

- [1] L. A. Adamic and E. Adar. Friends and neighbors on the web. *Social Networks*, 25(3):211–230, 2003.
- [2] L. Backstrom and J. Leskovec. Supervised random walks: predicting and recommending links in social networks. *Web search and data mining (WSDM)*, pages 635–644, 2011.
- [3] D. Boyd, S. Golder, and G. Lotan. Tweet, tweet, retweet: Conversational aspects of retweeting on twitter. *System Sciences (HICSS)*, pages 1–10, 2010.
- [4] S. Brin and L. Page. The anatomy of a large-scale hypertextual web search engine. *Computer Networks and ISDN Systems*, 30:107–117, 1998.
- [5] S. Brin and L. Page. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford University, January 1998.
- [6] P. Brusilovsky, A. Kobsa, and W. Nejdl, editors. *Content-Based Recommendation Systems*. Springer Berlin Heidelberg, 2007.
- [7] Z. Cheng, J. Caverlee, and K. Lee. You are where you tweet: a content-based approach to geo-locating twitter users. *CIKM '10 Proceedings of the 19th ACM international conference on Information and knowledge management*, pages 759–768, 2010.
- [8] R. C. et al. United states securities and exchange commission, October 2013.

- 
- [9] L. Getoor and C. P. Diehl. Link mining: a survey. *SIGKDD Explorations Newsletter*, 7(2):3–12, December 2005.
- [10] J. Hannon, M. Bennett, and B. Smyth. Recommending twitter users to follow using content and collaborative filtering approaches. *Recommender Systems (RecSys)*, pages 199–206, 2010.
- [11] H. Kwak, C. Lee, H. Park, and S. Moon. What is twitter, a social network or a news media? *World Wide Web (WWW)*, pages 591–600, 2010.
- [12] K. L. A new status index derived from sociometric analysis. *Psychometrika*, 18(1):39–43, 1953.
- [13] J. Leskovec, D. Huttenlocher, and J. Kleinberg. Predicting positive and negative links in online social networks. *World wide web (WWW)*, pages 641–650, 2010.
- [14] M.E.J.Newman. Clustering and preferential attachment in growing networks. *Physical Review E*, 64, 2001.
- [15] T. Sakaki, M. Okazaki, and Y. Matsuo. Earthquake shakes twitter users: real-time event detection by social sensors. *World Wide Web (WWW)*, pages 851–860, 2010.
- [16] G. Salton and M. J. McGill. Introduction to modern information retrieval. McGraw-Hill, 1986.
- [17] B. Wang, C. Wang, J. Bu, C. Chen, W. V. Zhang, D. Cai, and X. He. Whom to mention: Expand the diffusion of tweets by @ recommendation on micro-blogging systems. *World Wide Web (WWW)*, pages 1331–1340, 2013.
- [18] 鹿島久嗣. ネットワーク構造解析 - 機械学習によるアプローチ -. slide, 2006.
- [19] 神寫敏弘. 推薦システムのアルゴリズム (1). *人工知能学会誌*, 22(6):826–837, 2007.

- 
- [20] 神寫敏弘. 推薦システム recommender system. Slide, 2008.
- [21] 神寫敏弘. 推薦システムのアルゴリズム (2). 人工知能学会誌, 23(1):89–103, 2008.
- [22] 神寫敏弘. 推薦システムのアルゴリズム (3). 人工知能学会誌, 23(2):248–263, 2008.

## 発表文献

1. 岡本大輝, 豊田正史, 喜連川優. マイクロブログにおける対話ネットワークと投稿内容を併用したユーザ推薦に関する一考察. 電子情報通信学会データ工学研究会, 電子情報通信学会技術報告 Vol. 112, DE2013-29, pp.169-173 (2013.07).
2. 岡本大輝, 豊田正史, 喜連川優. 5. マイクロブログにおける対話手段と投稿内容に着目したユーザ推薦に関する研究と分析. 日本データベース学会情報処理学会データベースシステム研究会, B1-5 (2014.03). (to appear)