

フォロー別フィルタによるツイートフィルタリングの提案

山村 悟[†] 佐藤 哲司^{††}

[†] 筑波大学情報学群知識情報・図書館学類 〒305-8550 茨城県つくば市春日 1-2

^{††} 筑波大学大学院図書館情報メディア研究科 〒305-8550 茨城県つくば市春日 1-2

E-mail: [†]s0911666@u.tsukuba.ac.jp, ^{††}satoh@ce.slis.tsukuba.ac.jp

あらまし Twitter は、リアルタイム性の高い情報やユーザの嗜好を反映した情報など、多くの有益な情報が投稿されており、フォローをすることでフォロー先ユーザ (followee) の投稿を簡単に閲覧できる。このため、生活者の情報を遅滞なく共有できる新たな情報メディアとして注目されているが、followee のツイートが全てユーザの選好ツイートであるとは限らない。本研究では、ユーザの選好に合致するツイートのみをタイムラインに表示することを目的に、フォロー先毎にフィルタを設けて followee に対するユーザの選好を考慮するフィルタリング手法を提案する。followee のツイートだけを用いてフィルタを学習したのでは、十分な学習データが得られていないフォロー直後に、十分なフィルタリング性能が得られないことから、提案法では、あらかじめ分類済みの文書集合を用いた初期学習を行い、ユーザの判定した followee のツイートをを用いて逐次的に追加学習を行う。追加学習について 3 つの手法を考案し評価実験を行った結果、追加学習の際に重みをつける手法が最も提案法に適していることを明らかにした。

キーワード Twitter, 情報フィルタリング, ペイジアンフィルタ

Suggestion of Tweet Filtering using Filters by Follow

Satoru YAMAMURA[†] and Tetsuji SATOH^{††}

[†] College of Knowledge and Library Sciences, School of Informatics University of Tsukuba
1-2, Kasuga, Tsukuba, Ibaraki, 305-8550 Japan

^{††} Graduate School of Library and Information Sciences, University of Tsukuba
1-2, Kasuga, Tsukuba, Ibaraki, 305-8550 Japan

E-mail: [†]s0911666@u.tsukuba.ac.jp, ^{††}satoh@ce.slis.tsukuba.ac.jp

Abstract In Twitter, a lot of fresh information which contain user's preference is posted. Users can read tweets by following other user. Twitter has attracted as new information media, but followee's tweets are not necessarily the tweets of all the user's preference. In this paper, we propose filtering method that make and train filter for each followee. to display only tweets of user preferences in the timeline The proposed method train the filter sequentially after initial training using a set of documents classified because if it train the filter using only followee's tweet filter don't work initially. We devised methods of additional training. Result of evaluation, we demonstrated that method which give a weight to additional training data is the most suitable for the proposed system.

Key words Twitter, Information Filtering, Bayesian Filter

1. はじめに

代表的なマイクロブログである Twitter^(注1) は、ツイートと呼ばれる短い記事を投稿し、ユーザ同士で閲覧することができるサービスである。一般的に長文の記事を投稿する従来のブログサービスと比べて、ツイッターは 140 文字以内の短文を投稿するため、記事の作成に要する時間が非常に短く、手軽に投稿で

きるという特徴がある。そのため、イベントに対する反応のようなりリアルタイム性の高い情報や、新商品の感想のようなユーザの嗜好を反映した情報が投稿されやすく、新たな情報メディアとして注目されている。

Twitter では、ユーザは他のユーザを自由にフォローでき、フォローされたユーザのことを followee と呼ぶ。ユーザのタイムラインには、全ての followee のツイートが時系列順で表示される。田中ら [1] は、フォローする際のユーザの意図を、ユーザ指向、内容指向、相互性という 3 つの軸で分類している。ユーザ

(注1) : <http://twitter.com/>.

指向のフォローは、followee 自体に興味がある場合のフォローであり、内容指向のフォローは、followee の投稿するツイートの内容に興味がある場合のフォローである。相互性のフォローは、コミュニケーションを目的としたフォローである。例えば、友人や有名人をフォローする場合は、followee 自体に興味があるためユーザ指向のフォローである。この場合、ユーザにとって followee のツイートは、内容に関わらず閲覧したいツイートである可能性が高い。本研究では、ユーザが閲覧したいツイートのことを選好のツイートと呼ぶ。一方で、特定の話題に関する情報を得るためにフォローした場合は、followee の発信する内容に興味があるため内容指向のフォローである。この場合、followee のツイートでも、特定の話題以外のツイートであれば選好のツイートでない可能性が高い。

内容指向のフォローにおいて、ユーザの選好は followee によって異なることが一般的である。例えば、「雨が降ってきた」というような天気に関するツイートの場合、同じ地域の followee による投稿であれば選好であっても、別の地域の followee による投稿であれば選好でないことが多い。followee のツイートの中にも選好でないツイートは存在するにもかかわらず、フォローを行うと全てのツイートがユーザのタイムラインに表示される。Twitter では、投稿の容易さからツイートが頻繁に投稿される傾向があるため、選好でないツイートによって選好のツイートが埋没し、ユーザがタイムラインを閲覧する際の負担が大きくなるという問題がある。そのため、選好のツイートだけをタイムラインに表示することが望まれている。

本研究では、followee のツイートをフィルタリングし、選好のツイートのみをタイムラインに表示することで、内容指向のフォローを支援することを目的とする。ユーザの選好は followee によって異なるという点に着目し、各フォローごとに学習型のフィルタを作成する。フィルタは、あらかじめ分類済みの文書集合を用いて初期学習することで、フォロー直後からでも十分に使用可能な状態にする。さらに、followee のツイートをを用いて再学習を行うことで、逐次的にフィルタを更新し、followee に対するユーザの選好を反映したフィルタリングを行う。

2. 先行研究

Twitter を対象とした研究は盛んに行われている。岩木ら [2] はユーザの行動とユーザ同士のリンク構造によって有用な記事を発見する手法を提案している。ユーザ同士のつながりの強さを返信回数によって推定し、返信内容の特徴語を抽出することでユーザの興味に近いツイートの判別を行なっている。濱田ら [3] は、ユーザの選好が、傍観者、先駆者、追従者の 3 種類あると仮定し、ツイートをそれぞれの選好に分類している。分類の際に助詞、助動詞のような品詞に着目して重みをつけることで、精度が向上することを示している。竹中ら [4] は、ツイートをラベル付けを行うことができるハッシュタグという機能に着目し、ハッシュタグのついたツイートで学習した分類器を用いて、ハッシュタグのついていないツイートにハッシュタグを付けることで、ツイートの分類を行なっている。小坂ら [5] は、カテゴリ分類済みの文章によって学習した分類器を用いて、

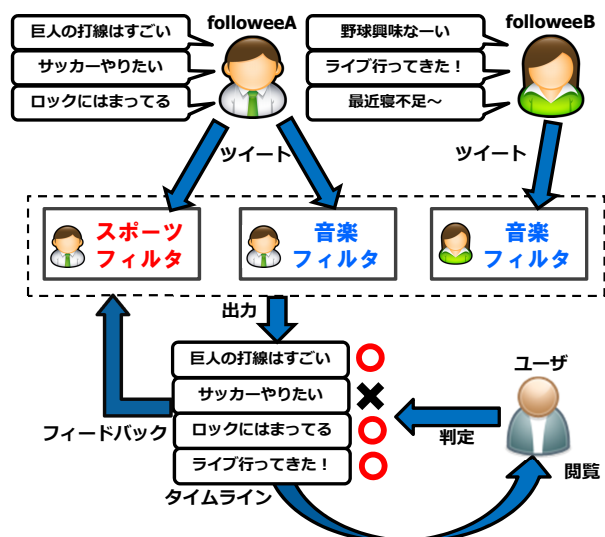


図 1 提案法の概要

ツイートをいくつかのカテゴリに分類している。分類器を構築する際に、ツイート以外の文章を用いて学習を行なっても、ツイートの分類ができることを示している。西田ら [6] は、特定の話題に関する圧縮済みのツイート集合に、新たなツイートを追加した際の圧縮されやすさに基づいて、ツイートの話題分類を行なっている。データ圧縮による分類手法は、形態素解析に依存しないため、新語や口語の多く含まれるツイートを精度よく分類できるという利点があるが、分類に時間がかかる。Sriram ら [7] は、ツイートをニュース、イベント、意見、お得情報、メッセージの 5 つの目的別に分類するシステムを提案している。BOW(Bag Of Words) に加えて、スラングや記号の有無、ユーザ情報などを素性を使用して分類精度を向上させている。

ツイートを分類する研究は数多く存在し、ツイート分類手法の応用としてツイートのフィルタリングを行うことができる。本研究では各フォロー関係ごとにフィルタを作成して、学習と分類をフィルタごとに行うことで、followee に対するユーザの選好の違いを考慮する点において、他の研究と大きく異なる。

3. フォロー別フィルタリングの提案

3.1 概要

本研究では、各フォロー別に異なる学習型フィルタを作成することで、followee に対するユーザの選好の違いに基づいてフィルタリングする手法を提案する。提案法の概要を 1 に示す。フィルタは followee に対するユーザの選好ごとに作成する。同じ選好であっても、followee が異なれば別のフィルタを作成し、同じ followee であっても、選好が異なれば別のフィルタを作成する。

学習型のフィルタは、事前に入力されたデータによって学習を行い、新たに入力されたデータを出力するか否かを確率によって判断する。本研究で作成するフィルタの学習には、ユーザが選好かどうか判定した followee のツイートをを用いるのが有効であると考えられる。なぜなら、followee に対するユーザの選好に基づいたフィルタリングをするには、followee のどのよ

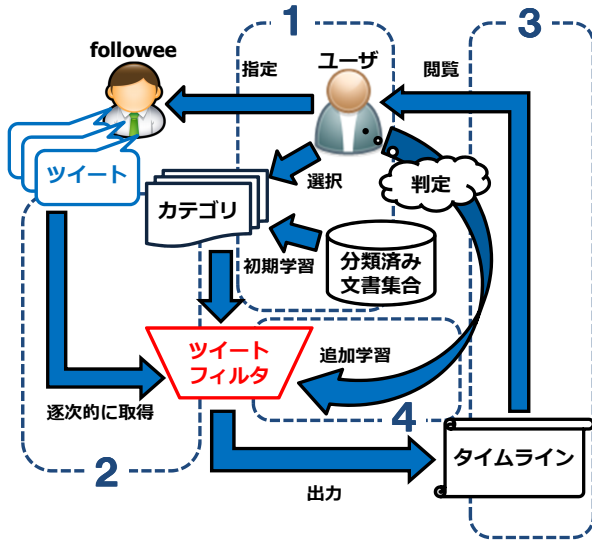


図2 ツイートフィルタリングの流れ

うなツイートがユーザの選好であるかをフィルタに学習させる必要があるからである。しかし、ユーザが判定した followee のツイートだけを学習データとして用いる場合、十分な量の学習データが蓄積するまではフィルタが上手く機能しないという問題がある。そこで、事前に分類済みの文書集合を用いて、フィルタ作成時に初期学習を行う。しかし、ユーザの選好が無数に存在し、それぞれの選好について事前に文書を分類しておくことは不可能である。そのため、本研究では内容の近い選好をカテゴリとして統合し、カテゴリごとに分類した文書を用いて初期学習を行う。カテゴリごとに分類された文書集合を用いて初期学習することで、フォロー直後であっても、カテゴリに関するツイートであるかどうかという基準でフィルタリングを行うことができる。さらに、followee のツイートによって追加学習することで、followee に対するユーザの選好を学習してフィルタの性能を向上する。

提案手法における処理の流れを図2に示す。

(1) ユーザは、followee を指定するとともに選好カテゴリを選択してフィルタを作成する。フィルタは分類済みの文書集合を用いて事前に初期学習しておく。

(2) followee のツイートを逐次的に取得し、フィルタによって選好確率を計算する。選好確率とは、取得したツイートがユーザの選好である確率である。

(3) 選好確率が閾値を越えたツイートをタイムラインに出力する。ユーザは出力されたツイートを閲覧し、選好のツイートであるかどうかの判定を行う。

(4) ユーザに判定されたツイートによって追加学習を行い、フィルタにユーザの選好を学習させる。

3.2 フィルタリング手法

文書フィルタリングは、文書分類に関する様々なアルゴリズムを応用することで実現できる。提案法ではユーザの判定のたびに何度も追加学習を行うため、学習が高速な手法によってフィルタリングする必要がある。そこで、NaiveBayes を文書フィルタリングに応用した手法であり、学習を高速に行うことのでき

るベジアンフィルタを提案法に採用する。ベジアンフィルタは一般的に spam メールフィルタによく用いられている。spam メールフィルタリングにおいては誤検出^(注2)は充分に少なくすることが重要であり、ベジアンフィルタは閾値を調整することで誤検出を減らすことができる。Twitter でも、選好ツイートをタイムラインに表示しないことは、非選好ツイートをタイムラインに表示することよりも、ユーザに与える影響が大きいと考えられるため、閾値の調整ができることは提案法にとって都合が良い。以上の理由より、本研究ではフィルタリング手法としてベジアンフィルタを採用する。

ベジアンフィルタでは、ツイート中の単語の出現に独立性を仮定し、各単語の選好確率の積によってツイートの選好確率を表す。ツイート T が k 個の単語を含むとき、以下のように定義する。

$$T = \{w_1, w_2, \dots, w_k\} \quad (1)$$

ツイート T の選好確率 $P(T)$ は以下の式で計算する。 $P(T)$ が閾値 t を上回ったツイート T をタイムラインに出力する。

$$P(T) = \frac{G}{G+B}$$

$$G = n_{good}/n_{all} \times \prod_{i=1}^k \frac{g(w_i)}{\sum_{w' \in V} g(w')}$$

$$B = n_{bad}/n_{all} \times \prod_{i=1}^k \frac{b(w_i)}{\sum_{w' \in V} b(w')}$$

n_{all} : 学習データの文書総数

n_{good} : 学習データにおける選好文書の総数

n_{bad} : 学習データにおける非選好文書の総数

$g(w_i)$: 選好文書における w_i の出現頻度

$b(w_i)$: 非選好文書における w_i の出現頻度

3.3 フィルタの初期学習

本研究では、フィルタ作成時にカテゴリを選択することで初期学習を行う。初期学習のためには、カテゴリごとに選好と非選好に分類された文書集合が必要となる。分類の基準は文書とカテゴリの関連性である。例えば、スポーツカテゴリであれば、スポーツに関連する文書を選好、関連しない文書を非選好として分類する。また、ベジアンフィルタにおいては、学習データは多いほどフィルタリングの結果が向上する。そこで、コミュニティQA(CQA)のデータを用いて初期学習を行う。CQAは、ユーザが投稿した質問に対して他のユーザが回答を投稿することで知識共有を行うサービスである。CQAでは一般的にカテゴリ毎に大量の記事が存在しているため、分類済みの文書として初期学習に用いるデータに適している。また、小坂ら[5]はツイートをカテゴリ分類する際のSVMの学習データとしてCQAデータを利用している。本研究では、大学共同利用機関法人国立情報研究所が提供をしている、Yahoo!知恵袋のデータ^(注3)を用いて初期学習を行う。フィルタ作成時に選択するカテゴリに

(注2) : 正常なメールを誤って spam として検出してしまうこと。メールのフィルタでは spam メールを通過させるよりも、正常なメールを遮断することの方が致命的であるため、誤検出は充分に少なくする必要がある。

(注3) : <http://research.nii.ac.jp/tdc/chiebukuro.html>

表 1 初期学習で選択するカテゴリ

初期学習カテゴリ	CQA カテゴリ
エンターテインメント	テレビゲーム全般/トレーディングカード/携帯型ゲーム全般/ 芸能人/テレビ、ラジオ/映画/アニメ/コミック/読書/音楽 おもちゃ、ホビー/楽器全般
スポーツ	スポーツ/オリンピック/ゴルフ/サッカー/スキー/プロ野球/ マラソン、陸上競技/モータースポーツ/格闘技、武術全般 卓球/釣り
IT	インターネット接続/ホームページ作成/動画共有/メール/LAN/ プログラミング/ウイルス対策、セキュリティ対策 Windows 系/Office 系 (Word, Excel)/画像処理、制作/インターネット
家電	家電、AV 機器/テレビ、DVD、ホームシアター/Windows 全般/ 携帯電話、モバイル/au/ソフトバンク/ドコモ デジタルカメラ/パソコン/掃除機、洗濯機/携帯オーディオプレーヤー
食事	料理、グルメ、レシピ/お酒、ドリンク/レシピ/ 菓子、スイーツ/ 料理、食材/飲食店
暮らし	住宅/役所、手続き/保険/税金/園芸、ガーデニング/ペット/ 家事/法律相談/消費者問題/家計、節約
美容	コスメ、美容/ヘアケア、ヘアスタイル/メイクアップコスメ/ ダイエット/ファッション/メンズ全般/レディース全般
健康	メンタルヘルス/カウンセリング、治療/健康、病気、病院/ 病気、症状、ヘルスケア/病院、検査
社会	国際情勢/政治、社会問題/ニュース、事件/事件、事故/株式/ 外国為替、FX/経済、景気
仕事	会計、経理、財務/労働問題、働き方/労働条件、給与、残業/ 就職、転職/就職活動/転職/アルバイト、フリーター
学校	大学/受験、進学/小・中学校、高校/中学校/小学校/高校/宿題/ 幼児教育、幼稚園、保育園/習い事
科学	天文、宇宙/化学/数学/物理学/工学/ 生物、動物、植物数学、サイエンス
文化	文学、古典/美術、芸術/歴史/日本史/世界史/宗教/ 絵画、手芸、工芸/伝統文化、伝統芸能
交通	車、高速道路/鉄道、列車、駅/飛行機、空港/テーマパーク/国内/ 海外/ホテル、旅館/観光地、行楽地
車	自動車/バイク/自転車、サイクリング/運転免許
恋愛	恋愛相談、人間関係の悩み/恋愛相談/性の悩み、相談

対応付けた CQA カテゴリの記事を取得する。例えば、スポーツカテゴリを選択した場合、CQA カテゴリの中で、スポーツカテゴリに対応する「プロ野球」、「オリンピック」といったカテゴリから取得した記事を選好の文書とし、スポーツカテゴリに対応しない「政治、社会問題」、「インターネット」といったカテゴリから取得した記事を非選好の文書として初期学習を行う。フィルタ作成時に選択するカテゴリと CQA カテゴリの対応付けを行った結果を表 1 に示す。初期学習の具体的な方法を以下に示す。

(1) 初期学習数が n の時、学習データの文書総数を $n_{all} = n$ とする。

(2) 学習データにおける選好文書の総数を $n_{good} = n/2$ とする。選択したカテゴリに対応する CQA カテゴリ数を m とすると、各 CQA カテゴリから n_{good}/m 件の記事を選好文書として取得する。

(3) 選好文書に対して、形態素解析器 MeCab^(注4)で形態素解析し、{ 名詞、動詞、形容詞 } のいずれかの品詞に該当する単語を取り出して得られた単語集合を $W = \{w_1, w_2, \dots, w_k\}$ とする。選好文書における単語 w_i の出現頻度 $g(w_i)$ に 1 を加算する。

(4) 学習データにおける非選好文書の総数を $n_{good} = n/2$ とする。選択したカテゴリ以外のカテゴリ数を l 、各カテゴリに対応する CQA カテゴリ数を m とすると、各 CQA カテゴリから $(n_{bad}/l)/m$ 件の記事を非選好文書として取得する。

(5) 選好文書の場合と同様に形態素解析によって、非選好文書から単語集合 W を取り出し、非選好文書における単語 w_i

の出現頻度 $b(w_i)$ に 1 を加算する。

3.4 フィルタの追加学習

初期学習したフィルタは、カテゴリに関連するツイートをタイムラインに表示することができる。しかし、実際のユーザの選好がカテゴリより具体的であった場合、初期学習だけでユーザの選好のツイートのみを出力することはできない。例えばスポーツカテゴリで初期学習したフィルタは、スポーツに関連するツイートを出力するが、ユーザの選好が野球であった場合には、サッカーやテニスについてのツイートは非選好である。そこで、ユーザは出力されたツイートが選好であるか判定を行い、判定結果をもとにフィルタを追加学習することで、フィルタにユーザの選好を学習させる。

フィルタの追加学習手法としてまず考えられるのは、ユーザが判定したツイートを新たな学習データとして加えるという方法である。つまり、ユーザが判定したツイートを $T = \{w_1, w_2, \dots, w_k\}$ とすると、判定が選好であった場合には、選好文書における単語 w_i の出現頻度 $g(w_i)$ と、選好文書の総数 n_{good} に 1 を加算し、判定が非選好であった場合には、非選好文書における単語 w_i の出現頻度 $b(w_i)$ と、非選好文書の総数 n_{bad} に 1 を加算する。いずれの場合も、学習データの文書総数 n_{all} に 1 を加算する。しかし、フィルタは初期学習によって多数の学習データを保持しているため、単純にその中にツイートを加える手法では追加学習の効果は薄いと考えられる。選好を学習するために大量の判定が必要になると、ユーザの負担は大きくなる。初期学習時のフィルタの性能は重要なため、初期学習に用いる学習データ数を減らすことは避けたい。本研究では、この問題に対処するために、以下に示す 3 種類の実現方法を検討する。

(注4) : <http://mecab.sourceforge.net/>

重み付き追加学習

ツイートを単純に学習データとして追加する手法の問題点は、初期学習データが多いため追加学習の効果が薄いことである。例えば、フィルタが初期学習で 10,000 件の文書を学習する場合、ユーザの判定が初期学習の内容以上の影響力を持つまでに、単純には 10,000 回の判定が必要となる。そこで、追加学習時に加算する値を σ として、最適な値を求めることで追加学習の効果を高めることができると考えられる。この手法を重み付き追加学習と呼ぶ。 σ の値は実験によって求める。

階層フィルタリング

本研究では、ユーザの選好を統合する目的でカテゴリを用意しているため、選好 \in カテゴリという関係になっている。そのため、初期学習したフィルタは選好のツイートに対して、再現率は高く適合率が低いフィルタリングをされると考えられる。例えば、選好が野球でスポーツカテゴリを選択した場合、サッカーやゴルフのツイートまで出力されるため適合率は低い、野球のツイートも出力するため再現率は高い。そこで、初期学習したフィルタとは別に、ユーザの判定したツイートのみで学習したフィルタを作成し、それぞれを初期学習フィルタ、追加学習フィルタと呼ぶ。まず、初期学習フィルタでフィルタリングを行い、選好とされたツイートを追加学習フィルタでさらにフィルタリングすることとする。この手法を階層フィルタリングと呼ぶ。なお、追加学習数が 0 の時は、初期学習フィルタによってフィルタリングを行う。階層フィルタリングは、フィルタリングを 2 回行うため再現率が低下することが考えられる。そこで、閾値 θ の値を低く設定し直す必要がある。 θ の値は実験によって求める。

並列フィルタリング

階層フィルタ同様に、初期学習フィルタと追加学習フィルタを別に作成する。追加学習フィルタは、followee に対するユーザの選好を学習しているため、十分な学習データ数があれば初期学習フィルタより性能が優れていると考えられる。しかし、初期状態の学習データ数は 0 であり、学習データが少ないうちは上手く機能しないという問題がある。そこで、初期学習フィルタと追加学習フィルタの両方でフィルタリングを行い、追加学習フィルタによってフィルタリング可能なツイートであれば追加学習の結果を優先し、それ以外の場合には初期学習フィルタの結果を優先することとする。この手法を並列フィルタリングと呼ぶ。追加学習フィルタリングによってフィルタリング可能な状態の判断には、選好確率 $P(T)$ の値を利用する。ベイジアンフィルタにおいて、 $P(T)$ が 0 や 1 に近い極端な値の時はツイートに対してフィルタが機能していると考えられる。逆に、ツイートの中に未知語や低頻度語が多く含まれてフィルタが上手く機能しない場合には $P(T)$ は 0.5 に近い値になる。そこで、追加学習フィルタによる $P(T)$ の値によって優先するフィルタを切り替える。具体的には、 $|0.5 - P(T)| > \delta$ の場合のみ追加学習フィルタの結果を優先し、それ以外は初期学習フィルタの結果を用いる。なお、追加学習数が 0 の時は、初期学習フィルタによってフィルタリングを行う。 δ の値は実験によって求める。

4. フォロー別フィルタリングの評価

本章では、CQA データを用いた初期学習手法の評価と追加学習手法の比較を行う。評価尺度として、それぞれの実験で用意したテストセットのツイートをフィルタリングした際の、適合率 (*Precision*)、再現率 (*Recall*)、F 値 (*F-measure*) を以下のように計算する。

$$\begin{aligned} Precision &= \frac{\text{出力した選好ツイート数}}{\text{出力したツイート数}} \\ Recall &= \frac{\text{出力した選好ツイート数}}{\text{テストセット中の全ての選好ツイート数}} \\ F\text{-measure} &= \frac{2 \times Precision \times Recall}{Precision + Recall} \end{aligned} \quad (3)$$

4.1 CQA データによる初期学習手法の評価

CQA データを用いて初期学習したフィルタによって、ツイートをフィルタリングした際の性能を F 値によって評価し最適な閾値 θ を決定する

4.1.1 実験データ

社会カテゴリと非社会カテゴリで CQA データを 10,000 件ずつ初期学習データとして使用する。テストセットには、社会カテゴリと非社会カテゴリに 300 件ずつ分類した 600 件のツイートを使用する。public_timeline^(注5)には膨大な量のツイートが存在するため、その中から特定のカテゴリのツイートを分類することは難しい。そこで、ハッシュタグを利用してカテゴリに近い内容のツイート集合を取得し、その中から人手でツイートを分類した。具体的なツイートの分類の手順を以下に示す。

- (1) ハッシュタグを利用してツイートを取得する。使用したハッシュタグは「#seiji」である。
- (2) 取得したツイートの中から人手で 300 件のツイートを社会カテゴリに分類する。
- (3) public_timeline からツイートを取得する。
- (4) 取得したツイートの中から人手で 300 件のツイートを非社会カテゴリに分類する。

本実験ではテストセットのツイートからハッシュタグは全て除去した。そのためフィルタリングの際にハッシュタグの情報は一切使用しない。

4.1.2 実験結果

閾値 θ を 0.05 から 0.95 まで変化させてテストセットのツイートをフィルタリングし、F 値を比較した結果を表 2 に示す。

実験の結果、閾値 0.85 の時に F 値が最大で 0.916 と高いため、本実験のように選好がカテゴリとほとんど一致する場合には、CQA データにより初期学習したベイジアンフィルタで十分にツイートをフィルタリングできることを明らかにした。

4.2 追加学習手法の比較

3.4 節で説明したフィルタの追加学習手法の比較実験を行う。

(注5) : Twitter 上で公開されているツイートの全てを一覧表示するタイムライン

閾値 θ	再現率	適合率	F 値
0.05	0.987	0.630	0.769
0.10	0.977	0.669	0.794
0.15	0.970	0.694	0.809
0.20	0.967	0.718	0.824
0.25	0.963	0.743	0.839
0.30	0.953	0.759	0.845
0.35	0.943	0.775	0.851
0.40	0.943	0.802	0.867
0.45	0.943	0.825	0.880
0.50	0.940	0.849	0.892
0.55	0.937	0.857	0.895
0.60	0.937	0.873	0.904
0.65	0.930	0.875	0.901
0.70	0.917	0.893	0.905
0.75	0.907	0.907	0.907
0.80	0.897	0.934	0.915
0.85	0.890	0.943	0.916
0.90	0.870	0.953	0.909
0.95	0.827	0.973	0.894

初期学習のカテゴリはスポーツカテゴリを選択する。本実験ではユーザの選好が野球に関するツイートである場合を仮定する。初期学習を行っただけの状態では、スポーツに関するツイートかどうかという基準でフィルタリングするため、野球に関するツイートだけをフィルタリングすることはできない。そこで、野球に関するツイートであるかという基準で判定されたツイートを各追加学習手法によって学習し、初期学習の状態から再現率を下げることなくどれだけ F 値を向上させたかという評価基準で比較を行う。

4.2.1 実験データ

スポーツカテゴリと非スポーツカテゴリで CQA データを 10,000 件ずつ初期学習データとして使用する。本実験では、追加学習によって followee に対するユーザの選好を学習することで、初期学習の状態からどれだけ F 値が向上するかを評価する。そのため実験データには、public_timeline から収集したものではなく、特定のユーザから収集したツイートをを用いる。実験データの条件として、選好ツイートである野球に関するツイートと、非選考ツイートである野球以外のツイートの両方が含まれている必要がある。また、非選考ツイートの中に野球以外のスポーツに関するツイートが含まれていることが必要である。なぜなら、野球に関するツイートとスポーツに関係のないツイートのどちらかしか存在しなければ、初期学習状態でも十分にフィルタリングを行うことができるからである。野球以外のスポーツに関するツイートであれば、初期学習状態のフィルタは選好ツイートとして出力するが、実際には非選考ツイートであるためフィルタリングが失敗する。このような場合に追加学習が必要となる。よって、追加学習によるフィルタの性能向上を比較するためには、初期学習状態のフィルタでは上手くフィルタリングできない野球以外のスポーツに関するツイート

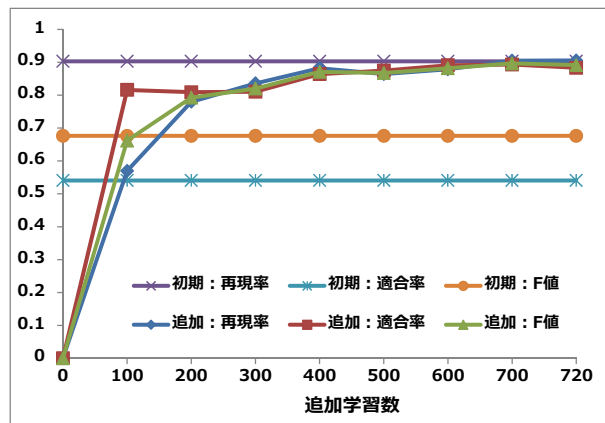


図 3 初期学習フィルタと追加学習フィルタによるフィルタリング結果

がなければならない。本実験では、上記の条件を満たすユーザから 800 件のツイートを取得して実験データとして用いる。実験データ 800 件中、スポーツに関するツイートは 507 件存在し、そのうち野球に関するツイートは 295 件であった。

4.2.2 実験方法

本実験では、800 件の実験データを用いた 10 分割交差検定を行う。720 件のツイートを追加学習を行い 80 件のツイートをフィルタリングする処理を 10 回行い、各値の平均によって評価を行う。また、各手法における最適なパラメータについて決定する。

本実験におけるベースラインとして、初期学習のみのフィルタと追加学習のみのフィルタによって、ツイートをフィルタリングした結果を図 3 に示す。初期学習フィルタは野球に限らずサッカーやゴルフのツイートも出力するため、再現率は高いが適合率が低くなり、F 値も低い値になっている。当然追加学習は行わないため、これ以上結果が向上することはない。一方、追加学習フィルタは初期学習を行っていないため、追加学習数が少ない間は再現率が非常に低い。追加学習数が多くなるにつれて再現率は向上するが、初期状態の再現率が低いと選好のツイートが出力されにくくなるだけでなく、追加学習の機会が少なくなりフィルタの性能向上の妨げとなる可能性ある。そのため、追加学習手法の条件として初期学習状態から再現率を下げることなく、追加学習数が増えるに連れて適合率を向上させることで F 値を向上させることが挙げられる。

4.2.3 実験結果

重み付き追加学習においては、追加学習の際に加算する重みが重要となる。そこで、加算する重みを σ として、 σ の値を変化させてフィルタリングを行った。結果をに図 4 に示す。実験の結果、 $\sigma = 1000$ の時に F 値が最大の 0.887 となった。

階層フィルタリングでは、初期学習フィルタと追加学習フィルタで二度のフィルタリングを行うため、再現率が低下することが考えられる。そこで、階層フィルタリングに適した θ の値を設定する。そこで、 θ を 0.1 から 0.9 まで変化させてフィルタリングを行った。結果を図 5 に示す。 $\theta = 0.7$ の時に F 値は 0.890 となった

並列フィルタリングでは、追加学習フィルタによって計算さ

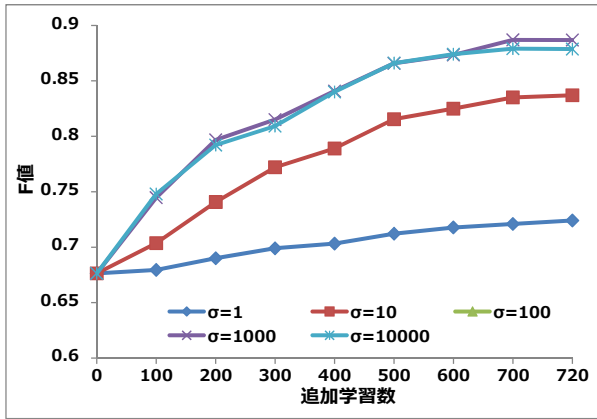


図 4 重み付き追加学習による F 値向上

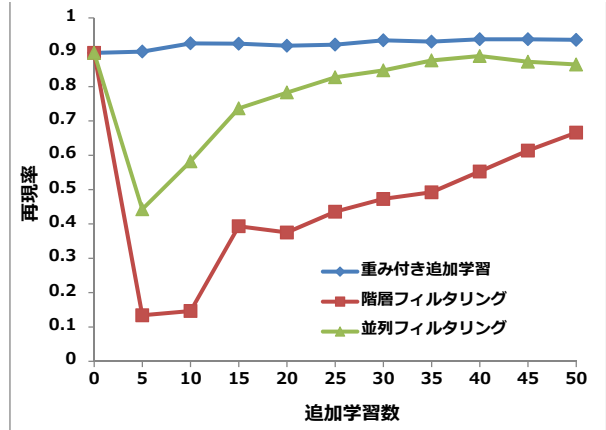


図 7 追加学習数が少ない場合の再現率

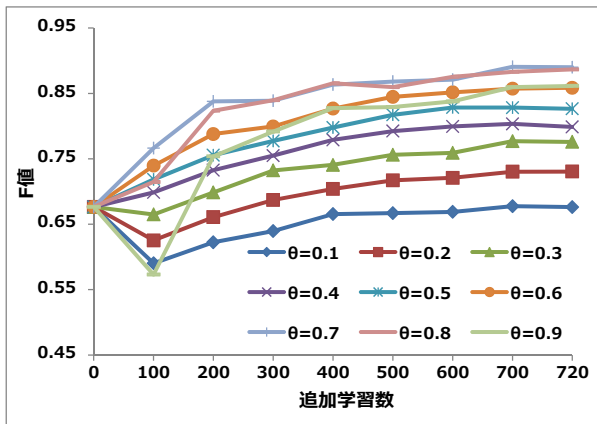


図 5 階層フィルタリングによる F 値向上

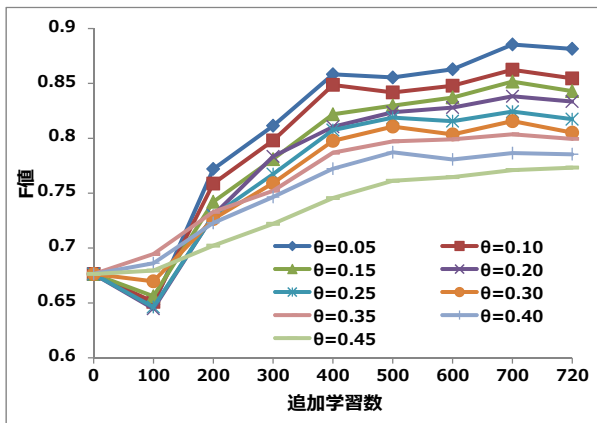


図 6 並列フィルタリングによる F 値向上

れたツイート T の選好確率 $P(T)$ の値によって、初期学習フィルタと追加学習のどちらでフィルタリングするかを決定する。具体的には、 $|0.5 - P(T)| > \delta$ の場合だけ追加学習フィルタを用いて、それ以外の場合には初期学習フィルタを用いる。 δ が 0.5 に近いほど初期学習フィルタを優先しやすく、 δ が 0 に近いほど追加学習フィルタを優先しやすくなる。 δ の値を 0.05 から 0.5 の間で変化させてフィルタリングを行った。結果を図 6 に示す。 $\delta = 0.05$ の時に F 値は 0.885 となった。

追加学習による F 値の向上については、階層フィルタリ

ングが最も高く 0.890 となったが、最低並列フィルタリングでも 0.885 と十分に高く、優位な差が見られなかった。そこで、追加学習数が少ない状態での再現率を比較するため、追加学習数を 0~50 の間で変化させてフィルタリングを行った。各手法のパラメータは前述の実験において F 値が最大となった時の値を用いる。結果を図 7 に示す。再現率の最小値について、階層フィルタリングでは 0.134、並列フィルタリングでは 0.443 と初期学習の状態から大きく低下しているが、重み付き追加学習では再現率が 0.901 と初期学習の状態より上がっている。追加学習によって再現率が低下すると本来見られるはずの選好のツイートが見られなくなり、追加学習の機会も少なくなってしまうため、この点においては重み付き追加学習に優位な差があると言える。

4.2.4 考 察

追加学習手法の比較の結果、F 値の向上においては大きな差が見られなかったが、追加学習数の少ない場合の再現率において、重み付き追加学習が他の手法を大きく上回る結果となった。これは、階層フィルタリングと並列フィルタリングが、追加学習数の少ない状態から追加学習データが結果に大きく影響するのに対し、重み付き追加学習では追加学習数が増えるごとに追加学習データの影響が徐々に大きくなるためであると考えられる。初期学習の状態からの再現率の変化において他の手法に対して優位性があり、再現率、適合率、F 値の全てがベースラインより大きく向上しているため、重み付き追加学習は提案法に適した手法であると言える。

5. ま と め

本研究では、ユーザの選好は followee によって異なるということに着目し、フォロー別にフィルタを作成することによって、各 followee に対するユーザの選好に基づいてツイートをフィルタリングする手法を提案した。followee のツイートだけを用いてフィルタの学習を行うと、フォロー直後のフィルタの性能が低いことが考えられる。そこで、分類済みの文書集合を用いて初期学習を行い、ユーザが判定したツイートによって逐次的にフィルタの追加学習を行うことで、フォロー直後から十分に機能し、利用するたびに性能が向上する学習手法を提案した。

フィルタリング手法には、追加学習が高速なベイジアンフィルタを採用し、初期学習に用いる分類済み文書集合として CQA のデータを採用した。実験の結果、選好がカテゴリとほとんど一致する場合には、初期学習のみでも F 値が 0.916 と高い値でツイートをフィルタリングすることができた。

フィルタの追加学習手法について、重み付き追加学習、階層フィルタリング、並列フィルタリングという 3 つの手法を考案し比較を行った。評価尺度として、追加学習による F 値の最大値に加えて、追加学習数の少ない状態の再現率にも着目して評価を行った。800 件のツイートをを用いて 10 分割交差検定を行った結果、F 値については階層フィルタリングが 0.890 と最も向上したが、他の手法と比べて優位な差は無かった。一方、追加学習数の少ない状態の再現率においては、階層フィルタリング、並列フィルタリングが初期学習の状態から大きく低下したのに対し、重み付き追加学習では初期学習の状態より向上した。よって、提案法におけるフィルタの追加学習には重み付き追加学習が適しているということが明らかとなった。

今後の課題として、重み付き追加学習の際に加算する重みを動的に決定する手法の検討が挙げられる。

謝 辞

本研究の一部は、筑波大学図書館情報メディア系プロジェクト研究による助成を受けたものである。本研究の実装・評価に際し、大学共同利用機関法人 国立情報学研究所から提供を受けた Yahoo!知恵袋のデータを利用している。ここに記して謝意を示す。

文 献

- [1] 田中淳史, 田島敬史. Twitter のフォロー関係のユーザの意図に基づく分類. 第 3 回データ工学と情報マネジメントに関するフォーラム (DEIM Forum 2011) F5-1, 2011.
- [2] 岩木祐輔, アダムヤトフト, 田中克己. マイクロブログにおける有用な記事の発見支援. 第 1 回データ工学と情報マネジメントに関するフォーラム (DEIM Forum 2009) A6-6, 2009.
- [3] 濱田翔吾, 黒澤義明, 目良和也, 竹澤寿幸. ユーザの知的欲求による選好に基づいたマイクロブログの記事分類. 情報処理学会研究報告. 自然言語処理研究会報告, Vol. 2011, No. 7, pp. 1–11, 2011.
- [4] 竹中姫子, 古宮嘉那子, 小谷善行. ベイジアンフィルターを用いた twitter におけるツイートのハッシュタグ分類. 情報処理学会研究報告. DD[デジタル・ドキュメント], Vol. 2011, No. 1, pp. 1–6, 2011.
- [5] 小坂龍一, 青野雅樹. 機械学習を用いた tweet の多カテゴリ分類. 第 4 回データ工学と情報マネジメントに関するフォーラム (DEIM Forum 2012) F5-1, 2012.
- [6] 西田京介, 坂野遼平, 藤村考, 星出高秀. データ圧縮によるツイート話題分類. 日本データベース学会論文誌, Vol. 10, No. 1, pp. 1–6, 2011.
- [7] B. Sriram, D. Fuhry, and M. Demirbas. Short text classification on twitter to improve information filtering. *Proceedings of 33rd International ACM SIGIR Conference on Re-*