

# 文体の違いを考慮したマイクロブログの極性判定

前田 裕      嶋田 和孝      遠藤 勉

九州工業大学大学院 情報工学府

{h\_maeda, shimada, endo}@pluto.ai.kyutech.ac.jp

## 1 はじめに

近年、評判分析の対象として、Twitterが注目されている。Goら[2]は、「:)」や「:(」などのemoticonを利用した訓練データの獲得と機械学習による極性判定(肯定・否定判定)の手法について提案した。Jiangら[3]は、Twitter特有の機能である「リプライ」などに着目し、極性判定する手法について提案している。Brodyら[1]は、Twitterによく見られる「Cooool」のような繰り返し表現に着目し、その繰り返し表現の正規化や感情との関連性などについて検証している。

ここで、我々はTwitter上で用いられる文体に着目する。Twitter上にはさまざまな文体が存在し、その特性が異なる。例えば、話し言葉に近い口語体では、特徴的な文末表現や記号などによる顔文字の多用によってその感情を表現する傾向があるのに対し、書き言葉的な文語体では、言語表現そのものによってその感情が表されることが多い。このような文体の違いを考慮することは、さまざまな場面に有効であると考えられる。

本論文では、このような文体の違いに着目した極性判定の手法を提案する。文体の判別については、人手によるルールと機械学習を併用する。極性判定のための訓練データの獲得については、Goらの手法と同様に、いくつかのemoticonや言語表現を利用する。文体を考慮した手法とそうでない手法を比較し、提案手法の有効性を検証する。

## 2 文体の定義

文中で利用される語や表現にはいくつかのバリエーションがある。機能語の出現頻度などの文体情報の利用は、blogの性別推定に関するタスク[5]などでその有効性が述べられている。

本稿では、極性判定というタスクの基、書き手の年齢や性別、どのようなネットコミュニティに属しているかでその文体が異なる点について着目する。例えば、「学校」という語の持つ印象は年代によって若干異な

表 1: 各文体の例

文体	例文
緩い文体	夜桜とってもキレイだったよ～
荒い文体	夜桜ガチでやばかったわ
硬い文体	夜桜の美しさはかつてないほどだった。

表 2: 各文体の特徴・印象

文体	特徴・印象
緩い文体	口語調, 女性的, 平和的, 緩やか
荒い文体	口語調, 男性的, 堂々, 2ch的
硬い文体	文語調, 長文が多い, 教科書的

ると考えられる。実際のTwitterでの使用例を観察すると、若年層の間では「学校の宿題メンドイ・・・」のような否定的文書が多いのに対して、社会人の間では回想として肯定的に語られることが多いといった傾向の違いが見られる。このような場合、その語の極性は書き手によって異なるため、一つのカテゴリでは適切に扱えない可能性が残る。別の例としては、「\(^o^)/」という顔文字の肯定性/否定性の解釈がある。この表現は一般に喜びを表す顔文字として使われることが多い。一方で、2ちゃんねるを中心としたネットコミュニティではしばしば否定的表現として用いられる。このような場合も、最終的な文の極性判別の際に誤識別の原因となる。一方で、これらの点を適切に区別して扱うことができれば、極性判定の精度は大きく向上すると考えられる。

以上のような考察から、本稿では、表1のような3種類の文体を導入する。それぞれの分類の基準は、表2に示される主観的な印象に基づいている。「硬い文体」は比較的高い年齢層で使われ、「緩い文体」や「荒い文体」は若年層で使われることが多いため、世代で異なる極性の切り分けに有効である。同様に、コミュニティが異なる「緩い文体」と「荒い文体」とでは、使われる表現も大きく異なっており、文体の切り分けは極性判定に有効に作用すると考えられる。

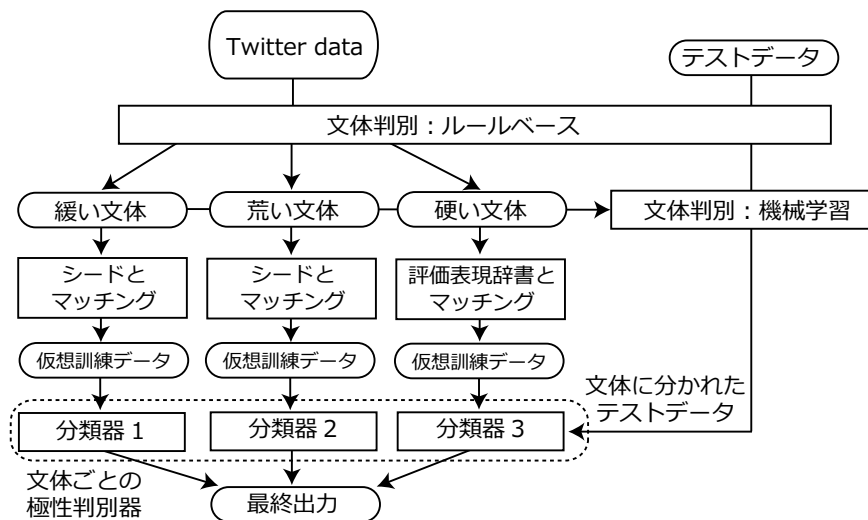


図 1: 提案手法の概要

### 3 提案手法

本節では、提案手法について述べる。処理の流れを図 1 に示す。提案手法では、まず、入力を前節で定義した 3 つの文体に分類する。その後、シード表現や評価表現辞書などを用いて、仮想的な訓練データを自動的に作成し、文体ごとに分類器を作成する。評価の際には、テストデータを同じような方法で文体ごとに分け、その文体に応じた分類器で極性判定を行う。以降、まず、文体判別の手法について述べ、次に訓練データの作成手法と極性判定のための分類器について述べる。

#### 3.1 文体判別

入力データを文体ごとに分類するために、文体判別を行う必要がある。本論文では、状況に応じてヒューリスティックなルールによる判別と機械学習による判別の 2 つの方法を利用する。それぞれの文体判別方法について具体的な説明を以下に行う。

##### 3.1.1 ルールによる文体判別

まず、ヒューリスティックなルールに基づいて文体を判別する。この手法は、再現率よりも精度を重視した文体判別を行うためである。ルールとしては、文体ごとに表 3 のようなパターンを数百ほど用意し、正規表現を用い、マッチングによって文体を判別する。文体を判別できる特徴は主に語尾や記号に表れる傾向があり、語尾一致と文中一致の 2 種類を用意した。いずれのパターンにもマッチしない場合には、どの文体にも分類されない。後述する極性判定のための分類器の作成に用いる訓練データを獲得する際には、このルールが適用されたもののみが利用される。

表 3: 文体判別の語尾パターン例

文体名	語尾パターン
緩い文体	だよね、だよお、だよ～、ね！、だもん
荒い文体	だな、だな w、だよな、だろ、かよ
硬い文体	だ。、である。、です。

表 4: 文体判別の文中パターン例

文体名	文中パターン
緩い文体	☆彡、≧▽≦、わーい
荒い文体	ざまあ、きめえ、お前、俺
硬い文体	(なし)

##### 3.1.2 機械学習による文体判別

ルールによる文体判別は精度重視の枠組みであり、再現率は低くなる。これは、次節で説明するように、高精度な文体ごとの仮想訓練データを得るためである。一方で、テストデータを評価する際に文体进行分类する場合、あまりに再現率が低いと、そもそも適切な分類器を適用できない可能性がある。そこで、再現率を重視したアプローチが必要となる。

ここでは、ルールベースで分類された各文体を訓練データとする機械学習を適用する。具体的にはナイーブベイズ分類器を利用し、素性としては文字 3-gram を利用する。実際には、各文体に対して、ルールベースの手法で得られた 10 万ツイートが学習データとして自動的に与えられる。

#### 3.2 訓練データの取得

文体ごとに極性判定を行うためには、機械学習に適用する訓練データが必要となる。しかしながら、十分な訓練データを人手で用意することは高コストである。本論文では、「緩い文体」と「荒い文体」のための訓練

表 5: シード表現の例

文体名	肯定シード	否定シード
緩い文体	♪	;)
荒い文体	きたこれ, おっしゃ	地獄, 糞すぎ

データ獲得にはシード表現を用い、「硬い文体」のための訓練データ獲得には既存の評価極性辞書に基づくスコアリング手法を適用することで対応する。

まず、シードを用いた方法について説明する。Turney[4]は、レビュー文書の極性判定において、語やフレーズの極性値を自動的に算出するために、“excellent”や“poor”といった肯定および否定の極性をよく表す語（シード）との共起度を利用する手法を提案した。Goら[2]は、emoticonを利用して訓練データを自動的に獲得している。本研究でもこれらに倣い、いくつかのシード表現をTwitterのデータに適用することで、仮想的な訓練データを獲得する。具体的には、表5に示す表現を、各文体のシード表現として利用する<sup>1</sup>。例えば「緩い文体」に関しては、ルールベースの文体判別手法で得られたツイート群で「♪」を含んでいるものを仮想的な肯定訓練データ、「;)」を含んでいるものを仮想的な否定訓練データとして扱うことになる。なお、各文体において、訓練データ数は肯定事例10万ツイート、否定事例10万ツイートを取得する。

次に、評価極性辞書を用いた方法について説明する。「硬い文体」の訓練データは、評価極性辞書を利用したスコアリングによって獲得する。本論文では、高村らの単語感情極性対応表[6]を利用する。この対応表には各単語に対して極性の強さを示す値が1から-1の範囲で付けられている。獲得処理では、対象となるツイート中に辞書のエントリと一致するものがあれば、その語に対応したスコアを加算する。最終的には、その合計値をそのツイートのスコアとする。本論文では、合計値が1.8以上のツイートを肯定事例、-1.8以下を否定事例として扱った。獲得される訓練データ数は「緩い文体」および「荒い文体」と同様である。

### 3.3 文体ごとの極性判別

本手法では、前節までに得られた仮想的な訓練データを用いて、各文体ごとに極性判別のためのPN分類器を作成する。すなわち、提案手法では、「緩い文体」のための分類器、「荒い文体」のための分類器、「硬い文体」のための分類器の3つが作成される。各分類器はすべてナイーブベイズ分類器であり、素性としては、

<sup>1</sup>表中の「;)」は泣いた顔文字や汗をかいた顔文字をターゲットとした意図がある。

文字 n-gram と単語 n-gram<sup>2</sup>について最適な組み合わせを用いる。テストデータを評価する際には、ルールベースおよび機械学習に基づく文体判別手法を適用し、同定された文体に対応する分類器へ渡され、最終的な極性判断が行われる。

## 4 実験

### 4.1 実験設定

訓練データなどの作成のために、Twitterから収集した約1億ツイートを利用した。これとは別に、評価用のテストデータとして、692ツイートを用意した。この692ツイートは人手によって肯定もしくは否定のラベルが付けられている。ここで、全テストデータ中、315ツイートが肯定文であり、377ツイートが否定文である。実験では、提案手法が出力した極性と的一致性を評価尺度として利用した。

### 4.2 結果と考察

実験結果を表6に示す。文体名の隣にある括弧内の数字は、ルールベースもしくは機械学習に基づく文体判別によってその文体へ分類されたツイート数を表している。例えば、「緩い文体」には692ツイート中307ツイートが分類されたことを表している。表中の「分類器(緩)」、「分類器(荒)」および「分類器(硬)」は各文体ごとの極性判定用の分類器を意味し、素性はその手法でもっとも精度がよかった場合の素性を表している。表で、例えば、「緩い文体」のための分類器である「分類器(緩)」で、「緩い文体」と判断されたツイートに対する極性判定の精度は88.3%であり、同様にもし「分類器(緩)」で異なる文体である「荒い文体」のツイートを極性判定した場合の精度は73.1%であることを表している。すなわち、表中の対角成分(88.3%, 82.8%, 78.6%)が、ある文体と判定されたツイートを適切な分類器で極性判定した精度を表しており、それ以外の場所は、文体と分類器にミスマッチが起きている状態の精度を表している。各文体で、もっとも精度が高いのは、提案手法が適切に適用されたその対角成分の数値であり、提案手法の有効性が確認された。

表中の「全体」と書かれた行は、ある特定の分類器のみで全データを分類した場合の精度を表している。これはすなわち、従来からある、文体を考慮していない分類器と等価であり、その精度と各文体ごとにカスタマイズされた分類器が正しい文体を分類した場合の精度比較することで、文体を考慮する提案手法の精度

<sup>2</sup>形態素解析にはMeCab(<http://mecab.sourceforge.net/>)を利用。

