

ツイートの長さや句読点に基づく年齢・性別の推定

江口 大賀†

明治大学総合数理学部 先端メディアサイエンス学科 菊池研究室†

1 はじめに

Twitter 上で投稿される文章は、身の回りの出来事や趣味に関する事を口語体で投稿される事が多い。そこで使われる言葉使いや単語は、ユーザの年代や性別によって変わる事が予想される。例えば長浜らは、ツイートから得られた単語の χ^2 乗値を用いるアルゴリズムを用いて、ユーザの性別の推定を行った [1]。男子は「僕、俺」などの名詞を多用し、女子は「*、ω」などの記号を多用する傾向があった。また、品詞の出現割合では、男女間で大きな偏りがなかったことを報告している。

そこで本研究では、自然言語処理の段階でストップワードに指定されて削除されがちである句読点等に注目し、Twitter に投稿された文章から、ユーザの年齢と性別の属性推定を行う。10 代のユーザは句読点の使用回数が 20 代以上のユーザより少なく、年代が上がるにつれて句読点の使用回数が多くなるという傾向が、予想されるからである。しかし、若い年代の句読点の使用回数が少ないのは、投稿する文章量が少ないからではないかという懸念点が挙げられる。そこで、文章の長さにも考慮し、句読点を使用する割合に注目した解析を行う。以上の本実験のシステム構成図を図 1 に示す。

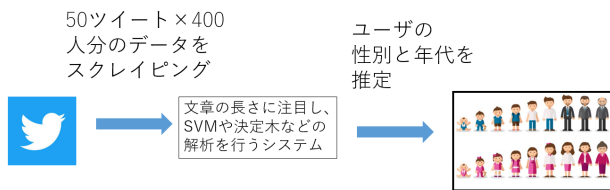


図 1 システム構成図

2 提案手法

2.1 データ収集

本研究のデータ収集では、プロフィールから Twitter のユーザを検索する「ツイプロ」[4]を用いる。このサービスは各アカウントのプロフィールの文章から、年齢・

性別・地域・職業・趣味などを分類している。「ツイプロ」を用いて、プロフィールから年代と性別が分かるユーザのアカウント名を集め、TwitterAPI 等を用いて、収集したユーザのツイートを取得する。

収集した合計 400 人分のユーザの属性と人数を表 1 に示す。1 人当たり 50 個のツイート、1 年で計 20,000 個のツイートを収集する。2019 年度と 2020 年度で合わせて 40,000 個のツイートを収集する。本研究のデータ取得時期を表 2 に示す。

表 1 収集したユーザの属性の統計値

性別\年代	10 代	20 代	30 代	40 代	合計
男性	50	50	50	50	200
女性	50	50	50	50	200
合計	100	100	100	100	400

表 2 データ取得時期

	2019 年度	2020 年度
期間	2019 年 11 月 4 日 ～2019 年 11 月 6 日	2020 年 11 月 30 日 ～2020 年 12 月 2 日

2.2 用いた手法

本実験では、3 つの手法による検証を行う。各手法の仮説、用いたデータ数、その手法による検証を行った年度を表 3 に示す。

表 3 手法の仮説とデータ数一覧

手法	立てた仮説	学習データ数	評価データ数	検証した年度
1	各属性で最も差が出る単語は、句読点である	200	200	2019 2020
2	若い年代は、句読点を使わない	100	100	2019 2020
3	若い年代は、文章が短い	100	100	2020

また、各手法の説明変数と目的変数を表 4 に示す。なお、手法 3 の説明変数の「生起確率」は、ユーザの全ツイートの文章における、句読点が含まれる割合を表した数値である。すなわち、以下の式で示される。

$$\text{生起確率} = \frac{\text{ユーザの句読点の使用件数}}{\text{ユーザの全ツイートの文字数の合計}}$$

†Department of Frontier Media Science, School of Interdisciplinary Mathematical Science, Meiji University, Kikuchi Laboratory.

表4 手法の説明変数と目的変数

手法	説明変数	目的変数
1	400人のユーザ内の誰かが3回以上使用した単語(2019年度は1,2052個, 2020年度は1,1147個)の使用件数	10代, 20代, 30代, 40代の4年代と性別の計8種類
2	「。」と「、」の2単語の使用件数	ある年代とそれ以外の年代の2種類
3	「。」と「、」の2単語の生起確率	

2.3 実験方法

様々な機械学習を用いて年代や性別の学習を行い、各属性を推定する。手法1ではランダムフォレスト、手法2・3ではサポートベクタマシン(以下、SVMと称する)を用いる。

本実験は、pythonを用いる。ツイートの収集では、urllib, pyquery, TwitterAPIを用いる。得られたツイートの自然言語処理にはjanomeを用いる。属性推定における決定木とサポートベクタマシンには、sklearnを用いる。

3 実験結果

3.1 手法1: 特徴語の抽出

手法1から抽出された特徴語の上位10語と、その重要度を表5に示す。なお、重要度の算出には、Random Forst Classifier内の関数であるfeature_importance[5]を用いた。重要度は、各々の説明変数の値が、目的変数を算出するのにどれ位重要かを示す。句読点の2単語の重要度は、1位と3位であった。また、特徴語の上位5語の各属性における平均使用件数を表6に示す。

表5 特徴語と重要度

順位	単語	重要度 [%]
1	。	0.84
2	私	0.74
3	、	0.62
4	を	0.61
5	まし	0.46
6	僕	0.42
7	ない	0.39
8	です	0.37
9	ある	0.37
10	!	0.36

表6 各属性の特徴語の平均使用件数

	10代		20代		30代		40代	
	男性	女性	男性	女性	男性	女性	男性	女性
。	7.1	7.0	33.7	21.4	45.4	41.2	65.7	37.3
私	0.4	2.0	0.3	2.3	1.5	4.3	2.0	6.4
、	10.6	20.4	34.4	20.8	39.3	39.3	60.9	41.5
を	7.8	7.9	21.3	13.7	27.9	23.2	38.2	26.7
まし	3.2	7.0	7.4	6.4	8.3	9.2	11.1	10.0

3.2 手法2: 句読点による年代ごとの分類

各年代の句読点の使用件数の平均と標準偏差を表7に示す。この結果から、年代が上がるにつれて、句読点の使用件数が多くなる事が分かる。

表7 各年代の句読点の使用件数の統計値

年代 \ 統計値	平均		標準偏差	
	。	、	。	、
10代	7.1	15.5	11.8	15.9
20代	27.5	27.6	43.1	29.4
30代	43.3	39.3	39.5	32.5
40代	52.0	51.2	40.6	37.9

句読点の使用件数から、SVMを用いて年代の推定を行う。その結果を表8に示す。この表での再現率は、100個の評価データの属性を推定し、正解したデータ数の割合とする。10代の推定の再現率が、特に高い事が分かった。4年代の中で最も再現率が高かった10代の推定の散布図を図2に示す。

表8 ある年代100人とそれ以外の年代100人の句読点による推定の再現率

	10代	20代	30代	40代
再現率 [%]	70.6	56.0	62.3	67.4

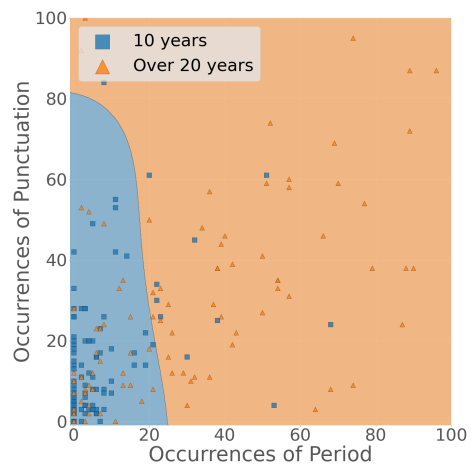


図2 10代100人と20,30,40代の100人の句読点によるSVMの分類の散布図

3.3 手法1・2:2019年度と2020年度の実験結果の比較

手法1・2における2年間の実験結果を比較した。表9より、手法1の結果では句点(「。」のこ)と読点(「、」のこ)ともに重要度の値が、2019年度より下がった。また、表10より、手法2の10代の推定の再現率が、2019年度の結果よりも5.6%低下した。

表9 手法1:“特徴語と重要度”の比較

順位	2019年度		2020年度	
	単語	重要度 [%]	単語	重要度 [%]
1	を	0.93	。	0.84
2	。	0.87	私	0.74
3	、	0.72	、	0.62
4	私	0.70	を	0.61
5	まし	0.57	まし	0.46

表10 手法2:“ある年代100人とそれ以外の年代100人の句読点による推定の再現率”の比較

	10代	20代	30代	40代
2019年度の再現率 [%]	76.2	54.5	61.8	62.9
2020年度の再現率 [%]	70.6	56.0	62.3	67.4

3.4 手法3:文章の長さを考慮した句読点による分類

各ユーザ群の文字数(1ツイート当たり)の平均を表11に示す。文字数をカウントする上で、URLやメンション(Twitterのユーザ名を載せる事)等は、正規表現を用いて削除した。

この結果から、年代が上がるにつれて、1ツイート当たりの文章量が多くなる事が分かった。中でも、10代男子の1ツイート当たりの平均文字数が32文字であり、特に少ない。

表11 各属性の文字数(1ツイート当たり)の平均

	男子	女子
10代	32.0	40.0
20代	50.7	42.8
30代	58.7	57.1
40代	68.6	58.2

句読点の生起確率の平均と標準偏差を、各年代に分けて表12に示す。10代の句点(「。」のこ)と読点(「、」)の生起確率は、0.36%と0.89%であり、4年代の中で最も低い値である。また、句読点の生起確率を用いたSVMによる年代の推定結果を表13に示す。

表12 各年代の句読点の生起確率 [%] の統計値

年代 \ 統計値	平均		標準偏差	
	。	、	。	、
10代	0.36	0.89	0.49	0.68
20代	1.06	1.08	1.20	0.84
30代	1.41	1.29	1.06	0.80
40代	1.54	1.51	1.06	0.86

表13 ある年代100人とそれ以外の年代100人の句読点の生起確率による推定の再現率

	10代	20代	30代	40代
再現率 [%]	73.0	55.9	59.2	66.0

3.5 手法2と手法3の実験結果の比較

各属性の句読点の平均使用件数(手法2)と平均生起確率(手法3)から、年代毎に偏差値を求め、表14に示す。

表14 句読点の平均使用件数と平均生起確率の偏差値

年代 \ 偏差値	使用件数 (平均)		生起確率 (平均)	
	。	、	。	、
10代	35.1	36.5	34.0	36.9
20代	47.0	45.6	49.2	45.1
30代	56.3	54.4	56.9	54.2
40代	61.4	63.4	59.7	63.7

表8と表13より、10代の推定の再現率においては、手法3が手法2よりも2.4%高い。一方で、20~40代の3年代の推定の再現率においては、手法3が手法2よりも低い。

4 考察

4.1 手法1の仮説「各属性で最も差が出る単語は、句読点である」

句読点の2単語の重要度は、表9より、2019年度と2020年度の両年で3位内の順位という結果であった。この結果は、手法1の仮説を支持するものである。

手法1の実験結果の問題点として、重要度の値が上位の単語でも低い事が挙げられる。例えば、2020年度の最も高い重要度は、「。」の0.84%であった。これは、説明変数の数が多い事が原因だと考えられる。

4.2 手法2の仮説「若い年代は、句読点を使わない」

表7から、年代が上がるにつれて、句読点の使用件数が多くなる傾向がある事が分かる。この結果は、手法2の仮説を支持するものである。

10代の推定の再現率は、表10より、2019年度が76.2%、2020年度が70.6%であり、他の年代の推定よりも高い結果であった。従って、句読点の使用回数に注目することで、10代か20代以上かの2択ならば、高い精度で推定する事ができると言える。

4.3 手法3の仮説「若い年代は、文章が短い」

表11より、10代の1ツイート当たりの文字数が、他の年代よりも短い結果となった。この結果より、手法3の仮説は支持された。

10代の句読点の生起確率は、表12と表14より、句読点の使用件数と同様に、他の年代よりも顕著に低い結果であった。

表13より、10代の推定の再現率は73.0%であり、他の年代の推定よりも高い結果であった。更に、これは表8の手法2の再現率よりも2.4%高い。よって、10代の推定には、句読点の使用回数よりも、句読点を使用する割合に注目する事が有効であると言える。

5 おわりに

本研究では、従来の実験ではストップワードに指定される事が多い句読点に注目し、性別や年代等の属性推定を行った。その結果、句読点の使用回数によって、10代か20代以上かの推定を行ったところ、再現率が70.6%であり、他年代の推定の再現率より10%ほど高いという結果が出た。

さらに、文章の長さを考慮し、句読点を使用する割合に注目した年代の推定も行った。その結果、「若い年代は、句読点を使用する頻度が低い」という特徴が示された。また、10代か20代以上かの推定においては、句読点を使用する割合に注目する事が有効である事が分かった。

参考文献

- [1] 長浜祐貴, 遠藤聡志, 當間愛晃, 赤嶺有平, 山田考治, “Twitterの投稿文章による人物像の推定”, 2012年度JSiSE学生研究発表会, 2013.
- [2] 岩朝史展, 松本和幸, 吉田稔, 北研二, “Twitterユーザの属性別感情推定の検討”, 言語処理学会第22回年次大会 発表論文集, pp.389-392, 2016.
- [3] 江口大賀, 菊池浩明, “ツイートの文章に使われている句読点に基づく属性推定”, 情報処理学会第82回全国大会, pp.3_431-3_432, 2020.
- [4] s21g Inc., “ツイプロ”, (<https://twpro.jp/>), 2021年1月参照)

- [5] “sklearn.ensemble.RandomForestClassifier — scikit-learn 0.24.1 documentation”, (<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>), 2021年1月参照)