

選択式天気情報を用いた ソーシャルメディアからの有用投稿抽出

萩行 正嗣

株式会社ウェザーニューズ

hangyo@wni.com

1 はじめに

近年、マイクロブログなどのソーシャルメディアを活用することで、実世界の動向を掴むことが広く行われている。人手で全ての投稿を確認することは困難なこと、比較可能な指標が求められることなどから、キーワードの出現頻度など動向の性質を示す何らかの数値的指標を用いて定量的に実世界の動向を把握することがおこなわれている。一方、投稿を数値量として表現することは、自然言語などで表現された個々の投稿の情報を捨象することである。ソーシャルメディアを活用した意思決定や動向の詳細な分析の際には、これらの捨象された情報が重要な役割を果たす場合があり、実際の投稿を人の目によって確認したいニーズが存在する。しかし、上述のように全ての投稿を人間が確認することは困難であるため、有益な情報が得られそうな投稿を優先して確認するための指標が必要となる。

ソーシャルメディアから得ることができる情報はソーシャルメディアの性質に大きく依存し、また何を有益とするかは情報の利用者によって異なる。本研究では、ウェザーリポートという気象情報を中心としたソーシャルメディアを対象として研究を行う。ウェザーリポートは話題が気象を中心という性質がはっきりとしており、そこから情報を得ようとする利用者も気象に関連する情報を求めていると考えられる。本研究では、ウェザーリポートに含まれる自然文で記述されたコメントから気象状況を知るにあたって有用な情報を得ることを目的とする。

一般に気象状況を知る手段としては、アメダスなどによる定点観測、気象レーダーや気象衛星による面的な観測がある。また、ウェザーリポートの各投稿(リポート)にはユーザーが入力した選択式の天気情報(例:パラバラ, ザーザー, 暑い, 寒い)が含まれており、これも気象状況を知るために活用できる。これらの観測

手段およびウェザーリポートのコメントの長所と短所をまとめたものを表1に示す。

ウェザーリポートによるユーザーからの投稿では、観測機による観測に比べると以下の2つの長所がある。1つ目は、観測機による観測が困難な現象を捉えることができることである。このような現象には、道路冠水や積雪など意志判断を行うにあたって重要な現象も多く含まれる。例えば下記の(1)では、ユーザーの周辺の地点で冠水の被害があることが分かる。

(1) 所々冠水してます

2つ目は、定性的ではあるものの、人間の感覚にあった表現が使用されており、情報の受け手に状況が伝わりやすいことである。数値による雨量表現よりも、下記の(2)の方が、多くの人にとって具体的な降雨の状況が想像しやすいと考えられる。

(2) バケツをひっくり返したような豪雨です

また、ウェザーリポート内の要素で比較すると、コメントによる記述は選択式の情報に比べて詳細な情報が得られるという長所がある。下記の(3)では、ただ寒いというだけでなく、ストーブをつけるほどの寒さであるという程度が分かる。

(3) 朝は肌寒くて、ストーブをつけた ($\geq \nabla \leq$)

コメントにはユーザーの感情が表現されることもあり、これも意思決定においては重要な手掛りとなる場合がある。下記の(4)は、大雨の際のリポートの例だが、防災上の対策等を行う意思決定において非常に有用といえる。

(4) 危険を感じます

これらをふまえ、本研究では、気象状況、被害状況への言及や感情表現が含まれるコメントを有用な情報とし、これらを含むリポートを抽出の対象とする。

*1 例えばアメダスの観測では雨と雪の判別が困難である。

表1 観測手段の長所と短所

観測手法	長所	短所
アメダス 気象レーダー	定量的な観測が可能 面的な観測が可能	観測点が限られる 降雨を直接観測しているわけではない 観測不能な場所(ブラインドエリア)が存在する
選択式天気情報 コメント	観測機で観測困難な現象も観測できる*1 報告の内容が多岐に渡る	報告内容が限られる 定量的解析が困難

本研究では、選択式の天気情報を学習に用いることで、半教師あり学習的にコメントにスコア付けを行い、このスコアをもとに有用なりポートを抽出する手法を提案する。本研究の目的は、選択式の天気情報を予測することではなく、有用なりポートを抽出することである。しかし、コメントから選択式の天気情報を予測する枠組みでレポートに与えられるスコアは、有用なりポート抽出において有用であると考えられる。例えば、下記の(5)のレポートでは、選択式の天気は「ザーザー」であった。しかし、このコメントから雨に関する意志決定に有用な情報は得られるとは言い難い。

(5) おはようございます。寒い朝でしたね

一方、コメントに対し機械学習的手法でスコアを付与する場合には、雨に関する情報が少ないため、低いスコアが与えられ、有用な情報はないと判断される。また、下記の(6)では、選択式の天気は「影うっすら*2」であったが、コメントからは雨について有用な重要が得られると言える。

(6) 南に凄い雨雲

このような場合でも、コメントに対して機械学習的にスコアを付与することで、有用なコメントと判定できる。

ソーシャルメディアからの情報抽出では、素性を抽出するための前処理として形態素解析や表現の正規化が行うことが多い。しかし、崩れた表現は著者の感情を表現している場合が多く、正規化を行ってしまうと、この情報を失ってしまうことになる。また、顔文字も著者の感情を表現するが、顔文字全体ではなく部分ごとの意味が重要となる場合がある。

(7) 今朝は寒い((((; ㄥ)))

((((; ㄥ)))は本来は恐怖などで震えることを表現する顔文字であるが、(1)では寒くて震えることを表現している。これは、顔文字両端の括弧の並びが震えを表現していることから判別できる。しかし、顔文字全体を一まとめにして扱ってしまうと、顔文字の部分ごとの情

*2 薄曇りを示す。

表2 ウェザーレポートの選択式の項目の例

天気情報名	選択肢の例
天気	影はっきり、影なし、 ポツポツ、バラバラ、ザーザー、 フワフワ、シンシン、ドカドカ
体感	ムシッと暑い、ジリジリ暑い、 暖かい、ちょうどいい、 涼しい、肌寒い、寒い

報を扱うことができない。そこで、本研究ではこれら
の問題に対して頑健な文字 n-gram を単位として扱う。

本研究では、ウェザーレポートに含まれる選択式の天気情報を活用し、文字 n-gram と気象状況(「雨」「雪」「暑さ」「寒さ」)との関連度を算出する。その関連度をもとにレポートのコメントと気象状況との関連度を算出し、有用なりポートを抽出する手法を提案する。

2 ウェザーレポート

ウェザーレポートは(株)ウェザーニューズの運営する気象情報に特化したソーシャルメディアである。1日平均約2万通のレポートが投稿される。投稿されるレポートには、写真やコメント、GPSの位置情報に加え、表2に示したような選択式の天気情報が付与されている。レポートの話題は天気のことが中心だが、紅葉など季節に関する話題や体感などに関連する生活の話題なども多く含まれる。

3 提案手法

各りポート r は、 $(wx, comment)$ で構成される。 wx は対象とする気象状況(「雨」「雪」「暑さ」「寒さ」)に対応する天気情報であり、「雨」の場合は表3の天気情報のいずれかの値を取る。 $comment$ はレポートの本文である。

提案手法では、各りポートのコメントと気象状況との関連度を半教師あり学習的に算出し、この関連度が高いものを有用なりポートとする。具体的な関連度の算出は以下で説明する。

表3 天気情報と雨の関連度

天気情報 (w_x)	w_{wx}
ポツポツ	1
バラバラ	2
サー	4
ザーザー	8
ゴォー	16
その他	0

表4 天気情報と雪の関連度

天気情報 (w_x)	w_{wx}
ベチャ(みぞれ)	1
チラ	2
フワフワ	4
バチバチ	4
(あられ)	4
シンシン	8
ドカドカ	16
ブワァー(吹雪)	24
その他	0

表5 天気情報と暑さの関連度

天気情報 (w_x)	w_{wx}
少し暑い	2
ムシッと暑い	3
カラッと暑い	3
暑くなった	4
蒸し暑くなった	4
暑い	5
ジリジリ暑い	5
耐えられない	10
暑さ	10
その他	0

表6 天気情報と寒さの関連度

天気情報 (w_x)	w_{wx}
涼しい	1
涼しくなった	1
涼しい気分になった	1
肌寒い	3
しっとり肌寒い	3
寒い	5
極寒	10
その他	0

3.1 n-gram 気象状況辞書の作成

提案手法では各文字 n-gram(cn) と気象状況の関連度 ($R(cn)$) を算出する。本研究では、文字 n-gram として、1-gram から 4-gram を利用する。まず、 cn と天気情報 (w_x) との $pmi(cn, w_x)$ を以下の式で計算する。

$$pmi(cn, w_x) = \log \left(\frac{\text{freq}(cn, w_x) \times N}{\text{freq}(cn) \times \text{freq}(w_x)} \right)$$

ここで $\text{freq}(x)$ は x を含むレポートの総数であり、 N は全レポート数である。

次に、 $pmi(cn, w_x)$ を重み付き和を以下の式で計算することで n-gram と気象状況の関連度 $R(cn)$ とする。

$$R(cn) = \sum_{\{w_x | pmi(cn, w_x) > 0\}} pmi(cn, w_x) \times w_{w_x}$$

ここで w_{w_x} は w_x との気象状況との関連度に対応する値だが、本研究では人手により設定することとし、「雨」の場合には表3のように設定した。

3.2 レポートと気象状況の関連度算出

レポート r に対する気象状況との関連度 $R(r)$ は、 $comment$ に含まれる n-gram に対する $R(cn)$ の合計をレポートの長さで正規化することで計算する。具体的には以下の式となる。

$$R(r) = \frac{\sum_{cn \in N\text{-gram}(comment)} r(cn)}{\text{length}(comment)}$$

ここで、 $N\text{-gram}(comment)$ は $comment$ に含まれる n-gram の集合である。

4 実験

4.1 実験設定

2014年のウェザーレポート 6,986,600通を利用して n-gram 気象状況辞書の作成とレポートと気象状況の関連度算出を行なった。「雨」「雪」「暑さ」「寒さ」に対応する w_{wx} は表3、表4、表5、表6のように設定した。

表7 $R(cn)$ の例(雨)

cn	$R(cn)$
雨風です	105.1
雨が窓	104.3
で冠水	97.0
停電しま	93.1
学校は休	92.9
水たまり	51.0
今はパラ	11.0
足湯	1.0

表8 $R(cn)$ の例(雪)

cn	$R(cn)$
開かなく	237.0
強く吹雪	233.3
0cm 以	223.6
雪まみれ	206.4
雪降ろし	188.3
積雪は無	51.45
路はツル	15.6
今夜も暖	2.0

表9 $R(cn)$ の例(暑さ)

cn	$R(cn)$
の暑さで	88.3
あぢい	74.3
あつーい	73.0
ゝ(´Д)	70.0
熱中症注	63.0
射しは強	34.3
一雨来る	13.0
また曇り	7.4

表10 $R(cn)$ の例(寒さ)

cn	$R(cn)$
して寒い	45.3
((+_	41.1
ストーブ	38.2
))}}	35.2
さぶっ	34.1
路凍結	28.9
月が綺麗	13.6
爽やかに	1.2

4.2 n-gram 気象状況辞書

作成された辞書の例を表7、表8、表9、表10に示す。「雨:雨風です」*3、「雪:強く吹雪」など気象状況を直接表現する n-gram だけでなく、「雨:停電しま」「雨:学校は休」「雪:開かなく」といった気象状況の影響を表現する n-gram に対しても高い関連度を与えていることが分かる。また、「暑い:あぢい」「寒い:さぶっ」のような崩れた表現や、「暑い:ゝ(´Д)」「寒い:((+_」「寒い:))}}」のような顔文字にも関連度も与えることができている。「))}}」は顔文字の中で震えを示す文字列であり、顔文字の構成性を扱うことができるといえる。

*3 気象状況「雨」と n-gram「雨風です」の関連度を示す。

