

集合注意の潜在的構造

Latent Structure of Collective Attention

笹原 和俊 *1

Kazutoshi Sasahara

*1 名古屋大学 大学院情報科学研究科

Graduate School of Information Science, Nagoya University

We propose “association networks” to visualize word-association patterns in a large dataset of tweets (short text messages) during collective attention events. First, we train the word2vec model to obtain vector representations of terms based on semantic similarities, constructing association networks: given some tokens as seeds, the associated tokens are linked with each other using the trained word2vec model, and considering the resulting tokens as new seeds, the same procedure is repeated. Using two sets of Twitter data—the 2011 Japan earthquake and the 2011 FIFA Women’s World Cup—we demonstrate how association networks visualize collective attention on these events. The results indicate the proposed method to be a useful tool for visualizing the implicit nature of collective attention that is otherwise invisible.

1. はじめに

インターネットのソーシャル化が急激に進行している現在、人々はオンライン上でこれまでにない規模で自発的に情報のやり取りを行い、その詳細は時々刻々とデジタルに記録・蓄積されるようになった。ソーシャルメディアは便利な情報発信のツールという域を超え、人々と情報を急速に結びつけ相互作用を促し、それによって実世界における人々の行動も変化するというサイクルが生じている。例えば、「アラブの春」と称された一連の民主化運動やバイト店員が悪ふざけた写真を投稿し炎上が多発した「バイトテロ」では、Twitter や Facebook などのソーシャルメディアが大量の情報や感情を媒介し、善かれ悪しかれこれらの事象の顛末に大きな影響を与えたことはよく知られた事実である。

このような大規模なソーシャルデータを利用して、「つながりすぎた世界」における人間行動や社会現象を探求しようとする試みが、計算社会科学の分野を中心として現在盛んに行われている [Lazer 09, Miller 11]。本論文ではオープンデータが利用可能な Twitter を対象として、「集合注意」(Collective Attention) と呼ばれるつながりすぎた世界に特有の創発現象に着目する。集合注意とは、実世界やネット上で生じた出来事が契機となって、ソーシャルメディアを介して人々の関心や注意がその出来事に集中する現象のことで、Twitter においては投稿数のパースト的増加や不安定な変動として捉えられることが知られている [Lehmann 12, Sasahara 13]。この現象は人間本性のマクロな現れと考えられ、ネット時代の人間行動を研究する上で重要である。集合注意の本質的理解のためには、投稿数のパーストのような顕在的性質だけでなく、その背後にある潜在的性質の探求が不可欠であり、そのためには集合注意の意味的構造にまで踏み込んだ理解が必要になる。しかしながら、ソーシャルデータから集合注意の意味的構造を定量化するための有効な手法がなかったことから、集合注意の理解はまだまだ表面的なレベルにとどまっている。そこで本論文では、集合注意の意味的構造を可視化するための新しい方法を提案し、実験を通じて集合注意の潜在的性質を観察する。

2. 方法

2.1 データの収集と前処理

実験には Twitter からスノーボール・サンプリングによって継続的に収集したユーザータイムラインのデータを使用した。ツイートの収集には Twitter REST API を用い、2011 年 4 月から約 1 年間をかけて約 40 万人のユーザーから約 5 億ツイートを取得した。そのうち、2011 年 3 月 11 日から 15 日 (東日本大震災) および 2011 年 7 月 16 日から 18 日 (FIFA 女子ワールドカップ準決勝と決勝) を実験用のデータセットとして準備した。データサイズは、震災関連のデータセットが 610,1930 ツイート、ワールドカップ関連が 7,497,877 ツイートである。

日本語形態素解析システム MeCab を用いて、ツイートのデータを分かち書き処理した。形態素解析の精度を上げるために、Wikipedia 日本語版の見出語データ、オープンソースの日本語入力システム Mozc の顔文字データと先行研究 [Sasahara 13] で使用された 5 種類の顔文字 (T, T, ^, ^, ^, ^, ^) を取り込んで日本語辞書 NAIST Japanese Dictionary をあらかじめアップデートした。したがって、形態素解析の結果は単語だけでなく複合語も含まれるため、ここでは分割されたものを「トークン」と呼ぶ。

2.2 word2vec による単語のベクトル化

Mikolov らが提案した word2vec は、コーパス内の単語を意味的演算が可能なたちで低次元のベクトルに高速に変換できる有効な手法である [Mikolov 13a, Mikolov 13b]。以下では Skip-gram モデルに絞って、word2vec がどのように単語ベクトルを作るのかを説明する。

Skip-gram モデルとは、単語 w_t が与えられたときにその周辺単語 (w_t の前後の c 個の単語) を推定する言語モデルである。図 1 のような三層のフィードフォワード型のニューラルネットワークにこれを学習させる。入力層に単語 w_t を入力し、正解の単語 w_{t+j} の出力確率が高くなるようにニューラルネットワークの重みを調整し、これを繰り返すことで学習が行われる。したがって、周辺単語の分布が似ていれば、単語ベクトルどうしも似た値をとるようになる。周辺単語の分布が似ているということは、単語の意味が似ていると捉えることができる。word2vec では隠れ層のサイズが入力層のそれと比べて著しく小さく設定されているため、例えばコーパスの語彙数が 100 万、隠れ層の

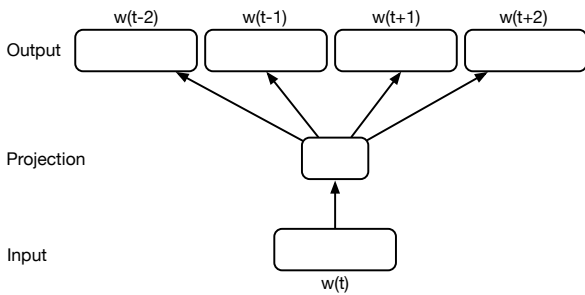


図 1: Skip-gram モデルの例 ([Mikolov 13b] を参考に作成).

サイズを 200 とすると, 100 万次元のベクトル (1-of-K 形式) は 200 次元のベクトルに圧縮される.

正確には, 入力文 (単語系列) を w_1, w_2, \dots, w_T とすると Skip-gram モデルの目的関数は次のようになる.

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t) \quad (1)$$

単語 w の入力ベクトルを v_w , 出力ベクトルを v'_w , コーパスに含まれる語彙数を W とすると, 式 (1) の $p(w_{t+j} | w_t)$ は次のソフトマックス関数で計算される.

$$p(w_o | w_I) = \frac{\exp(v'_{w_o} \cdot v_{w_I})}{\sum_{w=1}^W \exp(v'_{w_o} \cdot v_{w_I})} \quad (2)$$

W が大きくなると式 (2) の計算コストが大きくなるため, word2vec ではネガティブサンプリング等の手法を導入して計算量を削減している.

前述のように, word2vec は単語の並びの条件付き確率を学習するため, 使用文脈にもとづく単語の意味を反映したベクトルが構成可能である. その結果として, 単語の意味の演算ができたり, 単語の類推ができたりといった既存の言語モデルにはない特徴をもつ [Mikolov 13a, Mikolov 13b]. 例えば, Bag-of-Words モデルにもとづく LDA (Latent Dirichlet Allocation) [Blei 03] では単語の並びの情報を捨ててしまうため, 意味を反映した単語ベクトルを得ることは期待できない. したがって, 使用された単語の意味的つながりによって集合注意の潜在的特徴を可視化する場合, word2vec を用いた単語のベクトル化が適していると言える.

2.3 連想ネットワークの構成方法

word2vec の学習結果を可視化するための方法について述べる. 前述のように, ツイートの分かち書きから得られるのは単語と複合語の両方なので, これらを区別せずにトークンと呼ぶ. トークン 1 のベクトル v_{w_1} とトークン 2 のベクトル v_{w_2} の類似度を以下のコサイン類似度 s で測る.

$$s = \cos(v_{w_1}, v_{w_2}) = \frac{v_{w_1} \cdot v_{w_2}}{|v_{w_1}| |v_{w_2}|} \quad (3)$$

式 (3) を計算することによって, あるトークン w をシードとして指定した時に, w と周辺文脈が似たトークンを類似度順に取り出すことができる.

そこで, 類似度の閾値を s_{th} として, ツイートから作成したトークンベクトルから次の手順でネットワークを構成する. ま

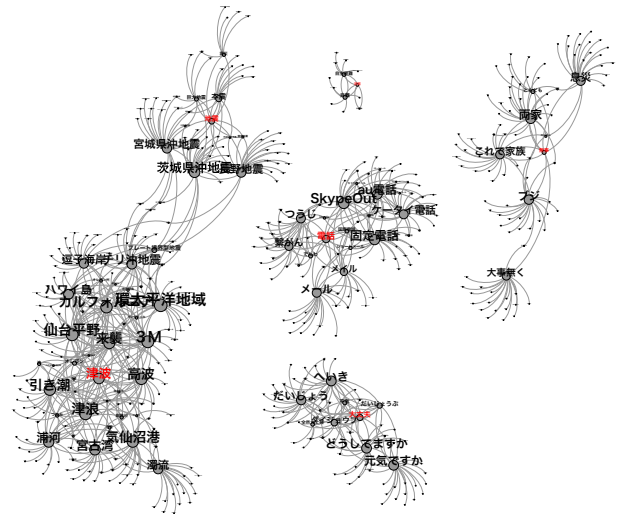


図 2: 頻出トークンから構成した連想ネットワーク (東日本大震災) 連想ネットワークのシードトークンは地震, 大丈夫, 無事, 津波, 電話, 避難でトークンの表示サイズは回数に比例.

ず, シードとなるトークン w とのコサイン類似度が s_{th} 以上のトークンを列挙し, 上位 N 個を採用する. シードとなるトークンが複数ある場合は, すべてのシードトークンに対して同様の作業をする. 次に, 採用したトークンすべてを新たなシードとして, 同様にコサイン類似度が s_{th} 以上のトークンを列挙し, 上位 N 個を採用する. このような手続きで得たトークンがネットワークのノードとなる. 次はリンクの貼り方であるが, トークン w_1 をシードとして w_2 が採用された場合, これらには意味的な類似関係があるとみなし, $w_1 \rightarrow w_2$ とリンクを貼る. トークン w_1 から採用されたトークンが複数ある場合は, それらすべてに対してリンクが貼られる. 連想ネットワークの配置には Force-directed レイアウトを用いる [Bastian 09]. このような一連の手続きは, word2vec の学習結果であるトークンの知識表現を使って連想ゲームをすることに相当するので, このネットワークを「連想ネットワーク」(Association Networks) と呼ぶことにする.

以下の実験では $s_{th} = 0.6$, $N = 20$ とした. ここで注意すべき点は, 上記手続きでトークンを列挙する際に $s_{th} = 0.6$ を満たすものが 20 個未満の場合もありうるということ, この条件を満たすトークンが 1 つもない場合はリンクが作成されない, ということである.

3. 結果

東日本大震災と FIFA 女子ワールドカップ 2011 に関するツイートを対象として提案した可視化手法に関する実験を行う. この実験の目的は, 集合注意の対象となるイベントの違いが連想ネットワークの構造にどのように反映するのかを観察することである.

3.1 頻出トークンの連想ネットワーク

まず, 震災関連のデータセットを用いて連想ネットワークによる可視化を行う. 先行研究において, 東北地方太平洋沖地震が発生した 2011 年 3 月 11 日 14 時 46 分以降に投稿されたツイートに頻出した単語は, 「地震」(1 位), 「大丈夫」(2 位), 「無事」(4 位), 「津波」(5 位), 「電話」(8 位), 「避難」(9 位) などであったことが報告されている [Sasahara 13]. そこで, こ

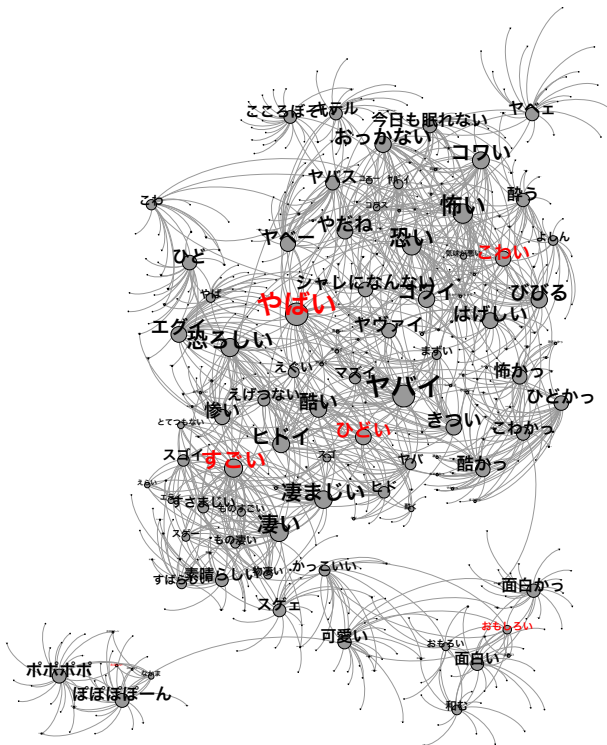


図 3: 感情形容詞の連想ネットワーク (東日本大震災): 連想ネットワークのシードトークンは「すごい」、「やばい」、「こわい」、「おもしろい」、「たのしい」、「ひどい」。トークンの表示サイズは度数に比例。

れら 6 つの震災関連の頻出トークンをシードとして構成した連想ネットワークが図 2 である。最頻出の「地震」は「茨城県沖地震」や「チリ沖地震」などの過去の大地震の名称を介して「津波」とつながり、地震や津波と関連するトークンを多数含みながら巨大なクラスターを形成している。このことから、東北地方太平洋沖地震が発生した直後に、人々は過去に起きた大地震や大津波を想起して直面している状況とそのイメージを重ねていたことが想像される。

それ以外の震災関連の頻出トークンは互いにつながることはなく、それぞれを起点とする特徴的なしかし狭域の意味的ネットワークを形成した。例えば、「電話」は「メール」や「Skype」とつながり、「大丈夫」は「どうですか」や「元気ですか」とつながった。また、「電話」と「でんわ」、「大丈夫」と「だいじょうぶ」のような表記揺れによるつながりや「避難」を「非難」とする頻出タイプによるつながりも確認された。興味深い連想の飛躍が見られたのは「無事」を起点とするネットワークで、「両家」や「これで家族」といった直感的でないつながりが観察された。連想ネットワークに震災以外の文脈が紛れ込んでいることから、震災当日は必ずしもネガティブなつづきのみがタイムラインを独占していたわけではないことがわかる。

3.2 感情形容詞の連想ネットワーク

同じデータを用いて、今度は感情に関わる 6 種類の形容詞「すごい」、「やばい」、「こわい」、「おもしろい」、「たのしい」、「ひどい」をシードとして連想ネットワークを構成して観察する。図 3 を見ると、6 種類の形容詞を起点とするクラスターどうしは連想によってすべてつながっていることがわかる。特に注目すべきは、ネガティブなトークン「こわい」や「ひど

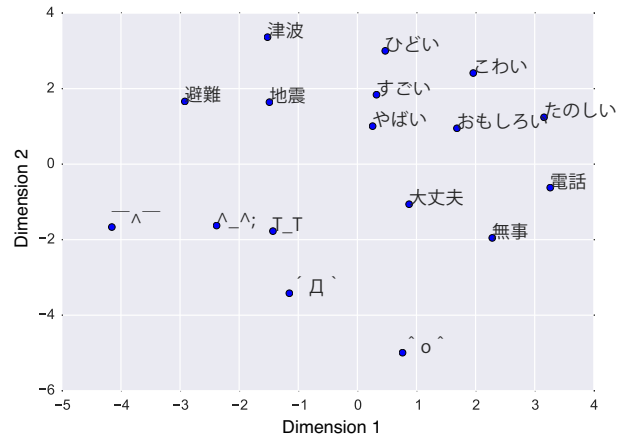


図 4: 震災関連トークンのベクトル空間における位置関係

い」と中立的なトークン「すごい」や「やばい」は連想が広範囲に広がって、巨大クラスターを形成していることである。その中には、「こわい」の同義語「恐ろしい」や「すごい」の同義語「凄まじい」などのトークンや、震災と関連する「今日も眠れない」や「よしん」などのトークンが密に相互連結していることがわかる。口語では「やばい」はポジティブな意味にもネガティブな意味にも使われるが、ここではもっぱらネガティブな意味で使用されていた様子が見える。このようにネガティブなトークンがネガティブなトークンと連鎖的に結びついているということは、震災によってマクロレベルで心的状態が落ち込んでいたことを視覚的に表現していると推測できる。

一方、ポジティブなトークン「たのしい」や「おもしろい」は連想が広がらず、小さいクラスターを形成するにとどまっている。しかし興味深いのは、「たのしい」は「ぼぼぼーん」(震災時に流れた CM のフレーズの 1 つ) と結びつき、「なかま」(この CM のキャラクター) や「可愛い」を経由して「おもしろい」につながっていることである。複数のクラスターをまたぐ連想の接続様式から、震災時に繰り返し流れたこの CM が人々の気分を和ませるのに少しなりとも寄与したのではないかと、ということを考えられる。災害時におけるメディア接触が人々の心的状態にどのように影響したのかは、社会心理学の重要な研究テーマであり、図 3 のような可視化はそのヒントになるかもしれない。

3.3 シードトークンの意味的位置関係

連想ネットワークのシードに使用したすべてのトークンのベクトル空間における位置関係を確認する(ただし、顔文字をシードとした連想ネットワークは紙面の都合で割愛)。図 4 は使用したすべてのシードトークンのベクトルを多次元尺度構成法で二次元に配置したものである。この空間の左上方向はネガティブさ、逆に右下方向はポジティブさに関係していることが読み取れる。そして、「地震」、「津波」、「避難」は近接し、これらとは空間的に離れた場所で「電話」、「無事」、「大丈夫」は近接している。感情形容詞は感情形容詞どうし、絵文字は絵文字どうしでまとまっているが、その中でも「すごい」と「やばい」や「T.T」と「^_^」など意味が近いものどうしが近接している様子がわかる。このように、震災時のツイートのデータセットでは、震災と関連した集合注意を反映したトークンの意味的關係性が生じていることが確認できる。

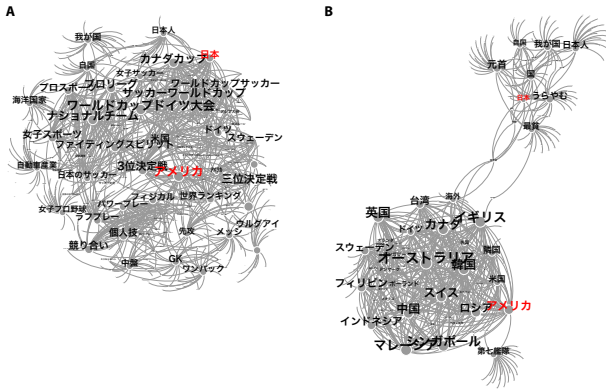


図 5: 「日本」と「アメリカ」をシードとした異なる文脈における連想ネットワーク。(A) はワールドカップのデータセットから構成した連想ネットワーク, (B) は震災のデータセットから構成した連想ネットワーク.

3.4 異なる文脈における連想ネットワーク

これまでの実験結果を踏まえると、集合注意の種類が違えば同じトークンをシードとしても異なる構造をもった連想ネットワークになることが考えられる。FIFA 女子ワールドカップのデータセットを用いてそれを確かめる。2011 年の女子ワールドカップはドイツで開催され、決勝戦では日本チームとアメリカチームが対決して、PK の末に日本チームが優勝した。そこで、「日本」と「アメリカ」をシードとして連想ネットワーク作成した。その結果が図 5(A) である。「日本」と「アメリカ」は国名なので、必ずしもサッカー関連のトークンばかりが連想でつながる必然性はないが、図 5(A) ではやはりワールドカップに関するトークンで単一のクラスターを形成している。「我が国」や「自動車産業」といったワールドカップ以外の文脈も混じってはいるが、「ファイティングスピリット」や「三位決定戦」などのサッカーを連想させるトークンがほとんどである。このことから、確かにワールドカップに関する集合注意が生じていることが確認できる。

次に、前節までの実験で使用した震災関係のデータセットを用いて、「日本」と「アメリカ」をシードとして連想ネットワークを構成したのが図 5(B) である。これを見ると、「日本」からはじまる連想と「アメリカ」からはじまる連想ががらうじて「世界中」というトークンでつながってはいるものの、ほぼ独立した 2 つのクラスターを形成していることがわかる。震災という文脈では「日本」は「うらやむ」や「最貧」と連想でつながり、一方、「アメリカ」は日本以外の外国名と連想でつながっている。したがって、予想通り、同じトークンをシードとしても集合注意の種類によって異なる構造の連想ネットワークが構成されることが確認できた。

4. まとめ

本論文では、ソーシャルメディアで生じる集合注意の潜在的構造を可視化する手法を提案した。連想ネットワークでは意味を反映したトークンベクトルを用いて、シードに指定したトークンの連想の連鎖によってネットワークを構成するため、有意義な結合様式や意外な関係性が表現される。

実験の結果、集合注意の意味的構造を反映したトークンのつながりやサブトピックに相当すると思われるクラスターが現れたり、支配的な感情がより誇張して表現されたりといった集

合注意の潜在的特徴を連想ネットワークに見出すことができた。このような結果は、Bag-of-Words にもとづくモデルや共起ネットワークによる可視化では得難い。

意味的連想による集合注意の潜在的構造の可視化は、実世界イベントを不特定多数のユーザーの自発的かつ主観的発言から特徴付けるといった行為に他ならない。連想ネットワークは不可視なものを可視化するための手法として有効であり、計算社会科学における新たな探索ツールとなる。

謝辞

本研究は公益財団法人堀科学芸術振興財団の研究助成を受けた。

参考文献

- [Bastian 09] Bastian, M., Heymann, S., and Jacomy, M.: Gephi: An Open Source Software for Exploring and Manipulating Networks, in *Proceedings of the Third International Conference on Weblogs and Social Media*, pp. 361–362 (2009)
- [Blei 03] Blei, D. M., Ng, A. Y., and Jordan, M. I.: Latent Dirichlet Allocation, *The Journal of Machine Learning Research*, Vol. 3, pp. 993–1022 (2003)
- [Lazer 09] Lazer, D., Pentland, A., Adamic, L., Aral, S., Barabási, A.-L. c. c., Brewer, D., Christakis, N., Contractor, N., Fowler, J., Gutmann, M., Jebara, T., King, G., Macy, M., Roy, D., and Van Alstyne, M.: Computational Social Science, *Science*, Vol. 323, No. 515, pp. 721–722 (2009)
- [Lehmann 12] Lehmann, J., Gonçalves, B., Ramasco, J., and Cattuto, C.: Dynamical Classes of Collective Attention in Twitter, in *Proceedings of the 21st International Conference on World Wide Web*, pp. 251–260 (2012)
- [Mikolov 13a] Mikolov, T., Chen, K., Corrado, G., and Dean, J.: Efficient Estimation of Word Representations in Vector Space, in *Proceedings of International Conference on Learning Representations (ICLR)* (2013)
- [Mikolov 13b] Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J.: Distributed Representations of Words and Phrases and their Compositionality, in *Advances in Neural Information Processing Systems 26*, pp. 3111–3119 (2013)
- [Miller 11] Miller, G.: Social Scientists Wade Into the Tweet Stream, *Science*, Vol. 333, No. 6051, pp. 1814–1815 (2011)
- [Sasahara 13] Sasahara, K., Hirata, Y., Toyoda, M., Kitsueregawa, M., and Aihara, K.: Quantifying Collective Attention from Tweet Stream, *PLoS ONE*, Vol. 8, No. 4: e61823 (2013)