

集合注意の創発とRTネットワークのダイナミクス

Emergence of Collective Attention and Dynamics of RT Networks

笹原和俊 *1 平田祥人 *2 豊田正史 *2 喜連川優 *2 合原一幸 *2
 Kazutoshi Sasahara Yoshito Hirata Masashi Toyoda Masaru Kitsuregawa Kazuyuki Aihara

*1名古屋大学大学院情報科学研究科 *2東京大学生産技術研究所
 Graduate School of Information Science, Nagoya University Institute of Industrial Science, The University of Tokyo

Quantifying distinct behavioral patterns using online social data is important for exploring collective social phenomena. Here we investigated “collective attention” on Twitter, a popular social networking service. To quantify it, we focused on the fact that tweet activity exhibits a burst-like increase and an irregular oscillation when a real-world event occurs; otherwise, it follows regular circadian rhythms. The difference between these states was measured using the Jensen-Shannon divergence, which corresponds to the intensity of collective attention. We associated irregular incidents with their corresponding events based on the popularity and the enhancement of key terms in tweets. The results of applying this method to 490 million Japanese tweets over 400,000 users revealed 60 cases of collective attentions. Retweet networks were also investigated to understand collective attention in terms of social interactions. Our findings provide new insights into social behavior in the digital era.

1. はじめに

ソーシャル・ネットワーキングサービス (SNS) は「今・ここ」を超えたコミュニケーションを促進し、時々刻々と生成されるソーシャルデータは電子化された行動の「化石」となる。大規模ソーシャルデータは人々の社会行動を定量的に探索する上で重要な手がかりを与えてくれる。SNS の 1 つである Twitter は「世の中の今を伝え合う」ツールで、ユーザは今どうしているのかを 140 文字以内でつぶやき (ツイート), 別のユーザがつぶやきで反応し (リプライやリツイート (RT)), その連鎖によってユーザネットワーク上を瞬間に情報が伝搬する。Twitter はリアルタイム性、ネットワーク性の高い社会行動を観察するのに適した系であるため、盛んに研究が行われており、これまでに様々な社会行動の特性が報告されている [Gonçalves 11, Wu 11, Golder 11, Weng 12].

本研究では、ツイートストリームにおける定常と非定常の差に着目して、Twitter 上で生じる集合注意 (Collective Attention) を定量的に捉えるための手法を提案し、集合注意の創発とその特性を探索する。また、RT ネットワークの構成と分析を通じて、集合注意の背後にある社会的相互作用の性質を明らかにする。

2. 方法

2.1 データセット

Twitter REST API*1 を利用してスノーボールサンプリングを行い、約 40 万人のユーザから約 5 億ツイートのデータ (ユーザタイムライン) を収集した。各データには、ツイートのテキストの他に、ユーザプロフィール、タイムスタンプや位置情報などのメタデータが含まれている。2010 年と 2011 年のデータを解析対象とし、ユーザ ID、ツイートのテキスト、タイムスタンプを解析に用いた。

2.2 集合注意の検出と定量化

通常、ツイートストリームは三相の概日リズムを示すが、実世界において大きなイベントが生じると、ツイートのバースト的な増加や不安定な振動が生じる。例えば、図 1 は 2011 年の 2 月と 3 月における 1 時間毎のツイート数の比較で、東北地方太平洋沖地震のあった 3 月 11 日の午後にはツイートが急増し、その後、概日リズムが乱れている様子がわかる。このようなツイートストリームの変化は、人々の注意の変化を集団レベルで反映していると考えられる。この仮定にもとづき、ツイートストリームの定常状態からのずれを集合注意の創発とみなし、ずれの程度を集合注意の強度とする。

ツイートストリームの定常状態と非定常状態の差を定量化するために、ツイートストリームを確率分布に直し、次式で定義される Jensen-Shannon ダイバージェンス (JS) を計算する。

$$JS(P, Q) = \frac{1}{2} \left(KL(P, \frac{P+Q}{2}) + KL(Q, \frac{P+Q}{2}) \right),$$

$$KL(P, Q) = \sum_i p_i \log_2 \frac{p_i}{q_i}.$$

JS は Kullback-Leibler ダイバージェンス (KL) を対称化したもので、確率分布 $P = \{p_i\}$ と $Q = \{q_i\}$ の差異を測るのに用いられる [Lin 91]. 定義から明らかなように、 JS は非負で常に有界のため、実データの解析に適している。データセットから 30 分毎にツイート数を集計し ($dt = 30$), 各日毎に求めたツイートの確率分布を P , 年平均のそれを Q として、上式より JS 値を計算する。

JS 値が大きくなるのは通常とは異なるツイートストリームを示した日で、何らかのイベントが実世界で生じ、それによって集団レベルでツイート行動が誘発されたと考えられる。ただし、図 2 が示すように、土日祝日は平日よりも JS が優位に大きくなり ($P < 0.001$; Mann-Whitney U test), 外的イベントがなくても通常と異なるツイートストリームを示す。そこで、土日祝日の JS 値のヒストグラム (図 2 の赤) をガウス関数でフィットして、平均 μ から 3σ 以上大きな JS 値を示す事象のみを集合注意の対象とした ($\mu = 0.0019$, $\sigma = 0.0009$, $JS_{\text{thresh}} = 0.005$).

連絡先: 笹原和俊, 名古屋大学大学院情報科学研究科, 〒464-8601 名古屋市千種区不老町, sasahara@nagoya-u.jp

*1 <https://dev.twitter.com/docs/api/>

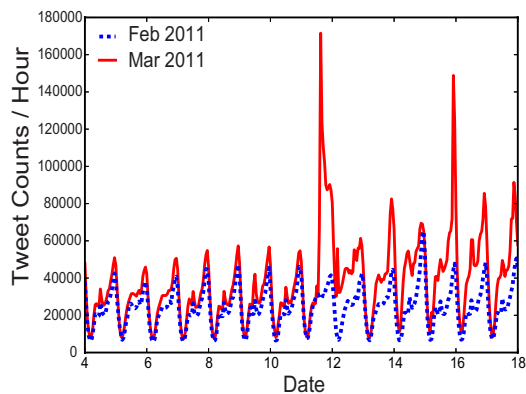


図 1: ツイートストリームの定常と非定常.

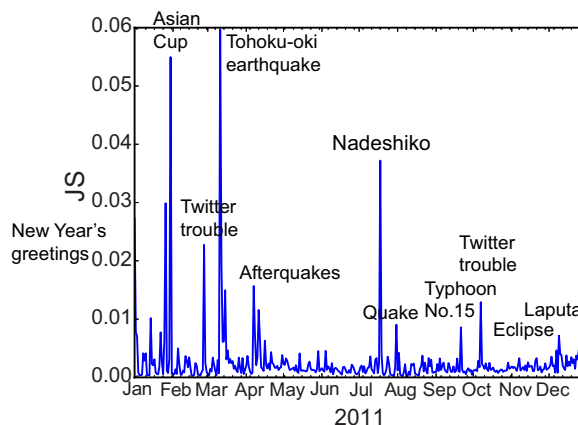


図 3: 集合注意の内容とその強度 (2011 年).

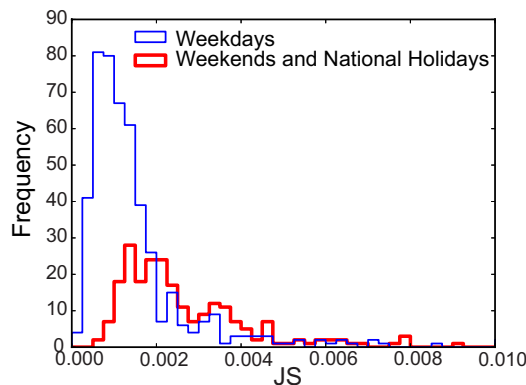


図 2: 平日と土日祝日の JS 値の分布.

その場合は上記処理を再帰的に行う.

3. 結果

JS 値が 0.005 以上を示した事象は、2010 年は 34 件、2011 年は 26 件あった. さらに JS 値が 0.01 以上の大きな値を示した事象は、2010 年は 12 件、2011 年では 6 件あった. 以下では、提案手法で検出されたこれらの集合注意の特性を見ていく. ただし紙面の都合で 2011 年の解析結果を中心に述べる. (詳細は [Sasahara 13] を参照)

3.1 集合注意の特性

2010 年と 2011 年のデータから検出された集合注意を同定して内訳を調べたところ、自然災害、スポーツイベント、文化、科学技術、政治、年間行事などに分類された. 中でも特に大きな JS 値を示したのが、自然災害に関するものだった. 大地震が発生した 3 月 11 日は 2011 年で最大強度の集合注意が生じ、この日から 4 日間連続で JS 値は 0.005 を超えた. このように連続した高強度の集合注意が観測されたのは、東日本大震災の時のみである. この際のツイートの内容は、M9.0 の大地震や東北地方沿岸部を襲った大津波、それらに起因した様々なトラブルに関するものがほとんどだった. 図 4C を見ると、「地震」、「津波」、「停電」など、日常では頻繁に使われることのないタームが大量にツイートされていたことがわかる. 人々の注意の移ろいやすさを考えると、いかに日本人が大震災から大きく、持続的に影響を受けたのかがわかる. 大震災以外にも、地震や台風などの自然災害に関連した集合注意の強度は総じて大きかった.

その他に高い JS 値を示したのはスポーツイベントに関するもので、日本女子サッカーチーム・なでしこジャパンのワールドカップ優勝や日本男子サッカーチームのアジアカップ優勝、バンクーバーオリンピックのフィギュアスケート決勝やプロ野球の日本シリーズに関するものがあつた. いずれもゲームの進行やテレビ中継と同期して、選手やチームへの応援や歓喜の声がツイートされた. これらの大強度な集合注意の場合、タームの頻度ランクのみから内容がそれとわかるものがほとんどだった.

中程度の JS 値を示した事象に目を向けてみると、文化や科学技術や政治などに関する興味深い集合注意が見られた. 例えば、小惑星探査機「はやぶさ」の帰還や皆既月食、人気アニメの「天空の城ラピュタ」(英語版タイトルは “Castle in the Sky”)

2.3 集合注意と関連するイベントの同定

次に、集合注意が生じている時間帯に発生したツイートのテキストを形態素解析してターム (t) を抽出し、得られたデータから頻度 (tf) を求めてランク付けした. 形態素解析には MeCab^{*2} と NAIST-jdic^{*3} を用いた. また、条件 $\{t \mid tf > tf_{before}, tf \text{ and } tf_{before} \geq \text{mean}\}$ を満たすタームに着目して、頻度増加率 (tf/tf_{before}) を求め、ランク付けした. 頻度ランクの場合、「今日」や「http」などのトリビアル・タームで上位が占められてしまうが、頻度増加率の場合には、トリビアル・タームは tf/tf_{before} がほぼ 1 になり、イベント当日の頻度 (tf) が前日の頻度 (tf_{before}) より増大した重要なタームは、 tf/tf_{before} が大きな値になると期待できる. したがって、得られたタームの頻度ランクだけでなく頻度増加率ランクを手がかりとして、検出された事象に対応する実世界のイベントを同定する. 本研究ではタームとして名詞のみを用いる.

2.4 集合注意と関連する RT ネットワーク

集合注意を社会的相互作用の観点から捉えるために、RT ネットワークを構成して、その構造を分析した. 例えば、ユーザ A があるイベントに関連するキーワード (名詞) を含むユーザ B のツイートをリツイートした場合 (“RT @user-B…” や “via @user-B…”), 各ユーザをノードとして A→B のようにリンクを結ぶ. ただし、本研究では非公式リツイートのみを解析対象とする. リツイートがさらにリツイートされることもあり、

*2 <http://mecab.googlecode.com/>

*3 <http://sourceforge.jp/projects/naist-jdic/>

や「ルパン三世—カリオストロの城」のテレビ放送や選挙速報に関連するものがあつた。特に、「天空の城ラピュタ」のクライマックス・シーンと同期したツイートは、アニメやネットの日本独特の文化や習慣を反映していて興味深い。明示的な指示があつたわけではないが、ユーザたちは主人公たちが呪文を唱えるのと時を同じくして一齐に「パルス」とツイートした(図4D)。この同期的ツイートは、当時の瞬間ツイート数(TPS)の最高記録を樹立した*4。この例のように、行動学的に興味深い集合注意は中強度のものが多い。その場合、タームの頻度ランクよりも頻度増加率の方がイベント同定のヒントになることが多い。

一方、正月や年越しなどの年間行事は、1日を通して全体的にツイートストリームが違うという非バースト・分散的なパターンが見られた。また、Twitterのサービスの不具合によって、つまりツイートできないことによってJS値が高い値を示すこともあり、2010年は7件、2011年は2件、そのような事象が確認された。1日のうちに異なる種類の集合注意が複数生じたケースもあつた。

3.2 集合注意と社会的相互作用

次に、非定常なツイートストリームを生み出すメカニズムを探るためにリツイート(RT)を解析した。図4のCとDに示されたキーワード(赤字)を含むツイートをリツイートしたユーザ同士のつながりの構造を調べた。図5は集合注意と関連したRTネットワークの例で、ノードがユーザを表し、リンクはRT関係、ノードの大きさはリツイートされた回数に比例する。どちらも通常とは違うツイートストリームの生成と関係しているのだが、明らかに両者のRTネットワークの構造が違う。東日本大震災のRTネットワークでは、巨大なノードを中心として大多数のノードが結合しており、これは大震災に関連したツイートが注目を集め、リツイートによってたくさんのユーザがつながったことを意味する。ちなみに、最大ノードは一般のユーザで、阪神淡路大震災の教訓をツイートしたことが多くのリツイートを生んだ理由だった。一方、天空の城ラピュタのRTネットワークには突出して大きなノードや大規模なノードの結合はなく、リツイートは局所的にのみ起こっていることがわかる。

これらのことを定量的に示したのが図6である。この図は、東日本大震災と天空の城ラピュタを含む4つの事象に関するRTネットワークに含まれる連結成分のサイズ分布を箱ひげ図で示したもので、縦軸は対数表示であることに注意されたい。この図を見ると、東日本大震災やなでしこジャパンの優勝などの大強度の集合注意の場合、RTネットワークには桁違いに巨大な連結成分(RTコア)があることがわかる。一方、皆既月食や天空の城ラピュタに関連した中規模の集合注意の場合は、RTコアは存在しない。

さらに、関連キーワードを含むツイートとリツイートの数をイベントの当日と前日で比較してみると、東日本大震災の場合、ツイート増加率は121倍でリツイート増加率は117倍、天空の城ラピュタの場合は、ツイート増加率が103倍でリツイート増加率が5倍となり、集合注意の創発と関連キーワードの増加が相関することが確認できた。

これらの結果は、集合注意下における社会的相互作用の違いに起因すると考えられる。つまり、集合注意の強度や特性にはツイートの量だけでなく、リツイートの連鎖の有無が関係していることが示唆される。

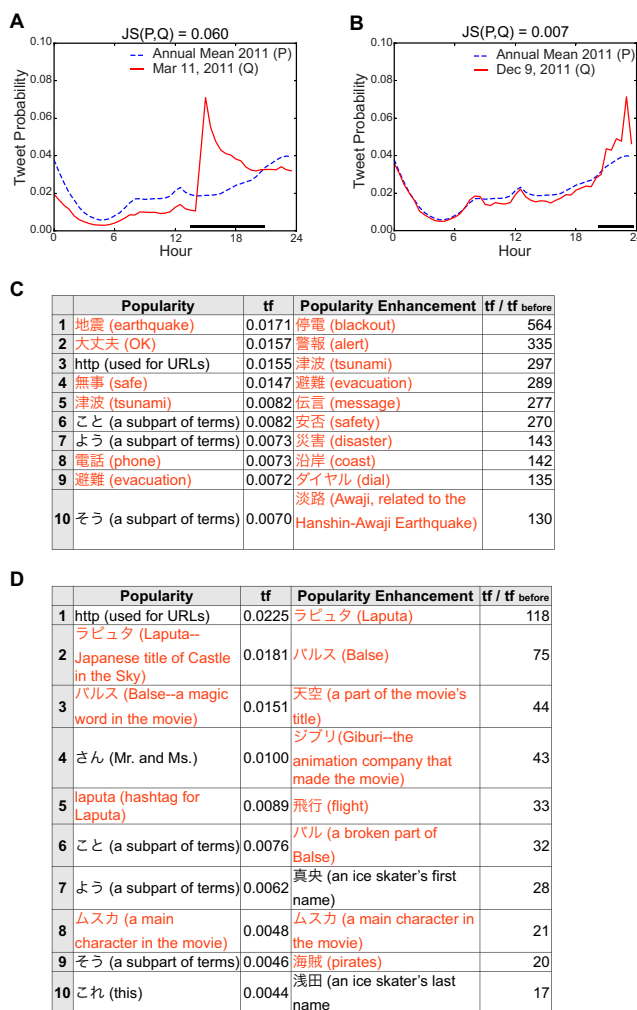


図4: 集合注意と関連したツイートストリームとタームランク。東日本大震災に関連したツイートの確率分布(A)とタームランク(C)、および、天空の城ラピュタに関連したツイートの確率分布(B)とタームランク(D)。各表の左列はタームの頻度ランク、右列は頻度増加率ランクで、赤字は集合注意と関連したキーワード。

4. 議論

本論文では、ツイートストリームの定常状態からの逸脱を集団注意の創発、その逸脱度合いを集団注意の強度として、実データを解析し、集団注意の特性を調べた。その結果、JS値が0.005を超える集合注意と関連する事象が、2010年と2011年を合わせて60件検出された。これらの事象は、自然災害、スポーツイベント、文化、科学技術、政治、年間行事などに分類され、特に自然災害とスポーツイベントは大強度、それ以外は中強度の集合注意と関連していた。さらに、JS値が大きかった事象に関してRTネットワークを構成して、集合注意の背後にある社会的相互作用を調べた。その結果、大強度の集合注意と中強度のそれは、RTネットワークに含まれる連結成分のサイズ分布が大きく異なることがわかった。これは、通常と違うツイートストリームを生み出すメカニズムには2つあり、大強度の集合注意はツイートとリツイートの大規模な連鎖に起因し、中強度のものは局所的な大量ツイートに起因することを示唆し

*4 <https://twitter.com/twittercomms/statuses/146751974904311808>

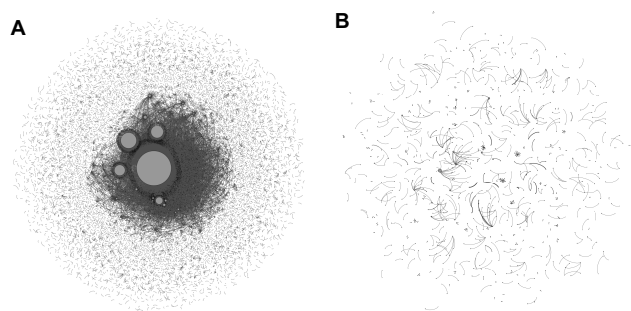


図 5: 集合注意と関連した RT ネットワークの例. (A) 東日本大震災 (ノード数は 27,340, リンク数は 27,709) (B) 天空の城ラピュタ (ノード数は 1,183, リンク数は 793).

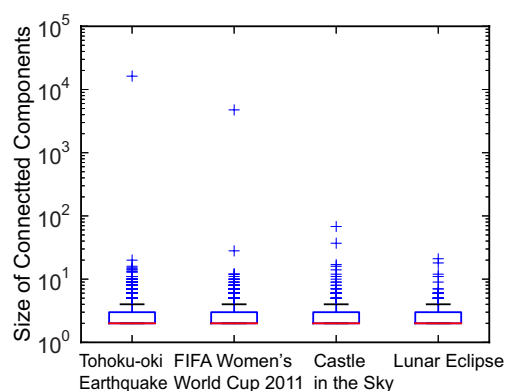


図 6: RT ネットワークの連結成分のサイズ分布. 左から東日本大震災, なでしこジャパンの優勝, 天空の城ラピュタ, 皆既月食.

ている.

類似した先行研究はツイートのバースト的变化に着目するものが主流だが [Lehmann 12], 提案手法では非バースト的・分散的なツイートパターンをもつ事象も捉えることができる. さらに, 集合注意を検出するだけでなく, その強度を定量化して比較することができるのが大きな特徴である. 本手法を利用する際の注意点として, dt や JS_{thresh} が小さすぎると偽イベントの検出確率が増えるため, データの性質に応じてこれらのパラメータを適切に設定する必要がある. また本研究では, 解析対象を日本語のツイートに限ったが, この方法自体は他言語のツイートにも応用可能である. 例えば英語の場合, すでにツイートは単語単位で区切られているため形態素解析の必要はなく, この手法をそのまま使うことができる. ただしその場合, 解析するデータの地理的局所性 (例えば, 国や地域, 位置情報) やネットワークの局所性 (例えば, ユーザのグループやクラス) を考慮した方が有意な結果が得られると予想される.

ソーシャルメディアの登場によって, 行動の「化石」がデジタルに蓄積されるようになると, 大規模ソーシャルデータから集団レベルの行動パターンやダイナミクスを定量化することが, 新しい行動学の遂行にとって重要になる. 提案手法は予測能力を有さないが, 集合注意を簡便に定量化・比較することができるという点で有効であり, 結果として得られる情報は, デジタル時代の社会行動の理解に新しい洞察をもたらす.

謝辞

本研究は, 総合科学技術会議により制度設計された最先端研究開発支援プログラム (FIRST 合原最先端数理モデルプロジェクト) により, 日本学術振興会を通して助成されたものです.

参考文献

- [Golder 11] Golder, S. A. and Macy, M. W.: Diurnal and Seasonal Mood Vary with Work, Sleep, and Daylength Across Diverse Cultures., *Science*, Vol. 333, No. 6051, pp. 1878–1881 (2011)
- [Gonçalves 11] Gonçalves, B., Perra, N., and Vespignani, A.: Modeling Users' Activity on Twitter Networks: Validation of Dunbar's Number, *PLoS ONE*, Vol. 6, No. 8, e22656 (2011)
- [Lehmann 12] Lehmann, J., Gonçalves, B., Ramasco, J. 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)
- [Lin 91] Lin, J.: Divergence Measures Based on the Shannon Entropy, *IEEE Transactions on Information Theory*, Vol. 37, No. 1, pp. 145–151 (1991)
- [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)
- [Weng 12] Weng, L., Flammini, A., Vespignani, A., and Menczer, F.: Competition Among Memes in a World With Limited Attention., *Scientific Reports*, Vol. 2, p. 335 (2012)
- [Wu 11] Wu, S., Hofman, J. M., Mason, W., and Watts, D. J.: Who Says What to Whom on Twitter, in *Proceedings of the 20th International Conference on World Wide Web*, pp. 705–714 (2011)