事象情報グラフによる複数文書中の時系列情報可視化 Graph-based Visualization of Timeseries Information in Multiple Texts

河田 裕成^{*1} Hironari Kawada 赤石 美奈^{*2} Mina Akaishi 細部 博史^{*2} Hiroshi Hosobe

*1 法政大学大学院情報科学研究科 Graduate School of Computer and Information Sciences, Hosei University *2 法政大学情報科学部 Faculty of Computer and Information Sciences, Hosei University

Newspapers and other media often describe particular topics over some periods by giving a series of texts that introduces new information and corrects previous information. By analyzing such texts, we can understand how a topic changed over time and how the media treated it. However, analyzing such texts manually is time-consuming. To solve this problem, we propose a method for visualizing texts that describe a topic changing over time. For this purpose, we first generate an event information graph from such texts by using an exiting method. Then we transform it into two different graphs. The first graph shows changes in the topic that we obtain by separating the initial graph and comparing the resulting time-dependent graphs. The second graph depicts a summary of the topic that we obtain by extracting high-frequency words from the initial graph. The results of our experiment show that our method can visualize changes and important information described in the texts.

1. はじめに

新聞等のメディアが報道する時事情報は、時間と共に刻々と 変化する.新たな情報が次々と報じられる中で、時には過去の 情報が間違いだったと判明することもある.こうした時系列的な 情報変化に関する文書群を分析することで、時事情報に関する 状況の推移と、メディアにおけるその扱いを知ることができる.し かし、長期間に渡って報じられた時事情報について、文書群中 の情報変化を人手により分析する作業には手間と時間がかかる.

そこで、本論文では、文書中の情報を表現する有向グラフに 対してグラフ操作を行うことによる、情報変化の可視化手法を提 案する.具体的には、先行研究の田代らの手法[田代 2013]を 用いて文書に含まれる実世界の出来事の情報をグラフ化した 後、それらの時系列的な比較と重要部分の抽出を行い、それぞ れを別のグラフとして可視化する.これにより、時系列情報に関 する情報変化や、時系列全体を通して出現する重要な情報の 分析支援を試みる.

2. 関連研究

情報の時系列的変化の可視化手法としては、伊藤らの研究 [伊藤 2013]がある.伊藤らは、係り受け解析により動詞句と名 詞句のペアをイベントとして抽出した.また、イベント群の時系列 的な構成や頻度の変化を3次元的に可視化した.伊藤らの手 法を用いることで、特定キーワードに関する話題の変化を可視 化することができる.しかし、キーワードとその関連単語の間に ある関係は明示されず、ユーザの推測に任せている.

田代ら[田代 2013]は、文書から事象を抽出し、単語間の関 係をネットワークとして可視化する手法を提案した.ここで、田代 らの定義する「事象」とは、実世界の出来事を表す情報のことで ある.具体的な手法は、まず、係り受け解析とCRFによって文書 中における 5W1H 等の事象属性句を推測する.その後、人手 で定義したルールによって事象属性句間の関係を推測し、事 象情報として LOD により表現する.田代らはこの手法を用いて 特定テーマに関する文書群中の事象情報を可視化し、メディア

連絡先:河田裕成, e-mail:17t0009@cis.k.hosei.ac.jp

による情報の差異を分析した.一方で,田代らは情報の時系列 的変化は研究の対象としていない.本論文では,田代らの手法 を用いて文書中の事象情報を可視化する.加えて,文書の公 開時期を考慮することで,文書の題材に関する情報の推移を分 析する.

3. 提案手法

本論文では、時系列的な情報変化を扱うために、特定テーマ に関する、情報公開日時の明示された文書を分析対象とする. その例としてニュース記事、歴史年表、SNS への投稿文等が挙 げられる.図1に示すような提案手法の流れを説明する.

まず,2章で述べた田代らの手法を用いて文書から事象情報 を抽出し,文書中の事象情報を有向グラフとして得る.この事象 情報グラフを2通りのアプローチで操作することで,時系列的な 情報変化と文書のテーマに関する重要情報をそれぞれ可視化 する.以下,3.1節において,事象情報グラフの概要を述べる. その後,事象情報グラフを用いた,情報変化の可視化手法を 3.2節で,重要情報の可視化手法を 3.3節でそれぞれ説明する.





3.1 事象情報グラフ

田代らの手法による事象情報グラフは、形態素からなる事象 属性句をノードとし、「Activity」や「Status」、「Location」など、13 種の関係をエッジとしている.これらのエッジは、文書内におけ る事象属性句の係り受け関係や、それを構成する形態素の品 詞情報、手がかり語を用いるルールによって設定される.なお、 ノードとなる語は形態素解析等により原形化することで、活用に よる表現揺れを可能な限り吸収している.例えば、「会社員の男 が本を買った」という文から生成される有向グラフでは、「会社員」 に Modifier 関係で修飾された「男」から、「買う」に動作を表す Activity 関係でエッジが張られる.さらに「買う」は、動作対象を 表す Object 関係で「本」に接続される.

3.2 情報変化グラフの生成

時間経過による新たな出来事の発生や,過去の情報に対す る補足や訂正がなされることを,時系列的な情報変化とする.こ れを可視化した情報変化グラフは,時系列中のある期間に出現 した事象情報群と,それらが直前の期間からどのように変化した かを表現する.その中で,事象情報は3.1節で述べたノードとエ ッジからなる有向グラフとして表現されるが,各ノードが情報変 化の種別を表すように着色されているという点で異なる.ここで の種別とは,直前の時期に出現したノードの①削除,②継続し た出現,対象時期における③新規ノードの出現の3つである.

ユーザは着色されたノードとそれらの接続関係を見ることで、 情報変化を分析する.例えば、②継続出現のノードに対して① 削除と③新規出現のノードが状態を表す「Status」関係で接続さ れている場合、①のノードが表す単語に関して、時間経過によ る状態の推移が発生したと考えることができる.図1左側のフロ ーに示す、情報変化の分析と可視化手法を説明する.

まず,文書の公開日情報を用い,事象情報グラフをユーザの 指定する一定期間毎に分割する.各グラフは,対応する期間に 出現した事象情報群を表す.次に,グラフ群を時系列順にソー トし,直前期間のグラフとのノード構成の差異を計算する.これ により,先に述べた①,②,③が得られる.可視化では,各グラ フ中の①のノードを青色,②を黒色,③を赤色に着色する.

3.3 重要情報グラフの抽出

本論文での文書の重要情報とは、文書の扱うテーマを説明 するものとする。例えば、自動車事故に関する文書であれば、 当事者の名前や事故の起きた状況等が重要情報にあたる。事 象情報グラフにおいて、重要情報は文書のテーマに強く関連 する単語ノードや、それらと近い距離にあるノードからなるサブ グラフとして現れると考えた。そして、文書のテーマに強く関連 する単語は、そのテーマを扱う文書群において時系列に関わら ず出現しやすいことに注目した。以上より、重要情報グラフは事 象情報グラフから抽出したサブグラフに対して、ノードが表す単 語の出現時期が判別できるように、ノードの着色処理を行うこと によって可視化することとした。

ユーザは着色されたノードによるサブグラフを見ることで、高 頻度で出現する単語にどのようなものがあるか知ることができる.



図2 単語の出現時期を示すノードの着色例

また,高頻度語間の関係や,それらを低頻度のノードが表す単 語とその出現時期を分析することができる.以下に,図1右側の フローで表される,重要情報の可視化手法を説明する.

始めに、文書の公開日情報を用いて、単語の時系列的な出 現頻度を計算する.これは、時系列の各期間での単語の出現 有無を2値で求め、それらを合計したものである.次に、出現頻 度が閾値以上であるノードを求める.その後、これらのノードを 始点とし、そこから指定距離以内で到達できるノードを探索する. 最後に、探索によって得られたノードと始点ノードから構成され るサブグラフを抽出する.なお、ノードの出現頻度に関する閾値 やグラフ探索の距離は、パラメータとして変更可能である.

可視化では、文書の公開時期ごとに色を割り当て、各ノード を出現時期に対応する色で着色している。図2に着色例を示す。 ここでは、ノード「証言する」は時系列的に初期の2文書のみに 出現し、ノード「海上自衛隊」は全ての文書において登場した。 よって、「証言する」は青色2色で、「海上自衛隊」は青色から緑 色、そして赤色までのグラデーションとなるように着色されている。

4. 適用実験

実際の文書群に対して本手法を適用し,情報の変化と重要 情報を可視化できているか評価する.本論文では,2014年1月 15日に発生した海上自衛隊の輸送艦「おおすみ」とプレジャー ボート「とびうお」の衝突事故を対象とし,実験を行った.運輸安 全委員会の船舶事故調査報告書*1によると,この事故では,左 舷前方から接近してきた「とびうお」と,それに対し右転して回避 しようとした「おおすみ」が衝突したとされる.その結果,転覆し た「とびうお」の乗員2名が死亡した.この事故に関して,2014 年1月から11月までに公開された毎日新聞社の新聞記事9日 分14記事を人手で収集し,事象情報を可視化した.図3に田 代らの手法による事象情報グラフの全体図を示す.このグラフ は,431個のノードと885のエッジからなる.提案した2つの手法 による可視化結果について,以下に述べる.

4.1 情報変化グラフ例

可視化における記事の分割期間は1日単位とし、可視化を 行った.提案手法では、時系列的に分割した期間毎に事象情 報グラフが出力される.本節では、出力例として2014年1月16 日に公開された記事に関する情報変化グラフを図4に示し、評 価を行う.このグラフは、1月15日から16日にかけて発生した、 記事中の事象情報変化を可視化したものである.

図 4 のグラフには, 124 の色付きノードが存在する. その情報 変化種別ごとの内訳は, 情報の削除が 43, 新規出現が 64, 継





*1 http://jtsb.mlit.go.jp/jtsb/ship/detail.php?id=6270



図5 情報変化グラフ中の新規情報出現例

続出現が 17 ノードである. なお, 継続出現として可視化された 17ノードのうち9つは事故の当事者名や「衝突事故」など, 事故 の基本的な情報を表すもので, 図で示した1月15日から同16 日までの期間に特有な情報は半数以下であった.

ここで、ノードの色に注目することで、より詳細な情報変化の 分析を行う.図4のグラフ中における情報の新規出現例を図5 に、特定単語に関する情報変化の例を図6に示す.ただし、図 6中のノードにおける人物名は黒塗りとして伏せてある.

図5では、「新聞社などの取材に対し、男性が15日夜に証言 した」という情報の新規出現が可視化されている.また、ノード 「男性」に向かう Modifier 関係のエッジは「釣り船とびうお」から 伸びており、「男性」が事故で転覆した「釣り船とびうお」の関係 者であると把握できる.同様に、「証言する」からの Object 関係 のエッジを辿ることで、証言の内容を知ることができる.

図6では、ノード「2人」に関して「15日に病院へ搬送された」 から「16日に死亡した」という時間経過による情報の推移が可視 化されている.この「2人」とは、とびうおの乗員のことを指してい る.しかし、グラフでは「2人」と「釣り船とびうお」の間にエッジは 存在しないため、グラフのみからそうした情報を得ることはできな い.実際の文書において、図6に示した2つの情報はそれぞれ、 1つの文として出現する.しかしながら、「釣り船とびうお」という 単語はこれらの文中に存在しなかった.そのため、事象情報グ ラフ構築時にノード間の関係設定が期待通りに行われず、情報 の不足が発生したと考えられる.

以上より,提案手法は,時間経過による文書中の単語の出現 や削除に注目した情報変化を可視化できる.グラフにおいて, 単語を表すノードは事象内の関係を基に接続される.ユーザは エッジに付与された関係ラベルやノードの色が示す情報変化種 別の組み合わせを見ることで,情報がどのように変化したか知る ことができる.一方で,提案手法は文中で明示的に書かれてい ない関係を可視化することはできない.また,情報変化の構造 や,グラフ中におけるその範囲はユーザが判断する必要がある.

4.2 重要情報グラフ例

時系列全体の文書からなるグラフに対し,全9日中8日以上の頻度で出現した単語ノードから,距離2で到達できるノードからなる重要情報グラフを抽出した.その結果を図7に示す.なお,時系列的な出現頻度によるノードの着色は,図2の右側に示すように,全9日のそれぞれに色を割り当てることとした.

図 7 において、閾値の 8 回を超える頻度で出現したノードは 「海上自衛隊」「輸送艦おおすみ」「釣り船とびうお」「衝突事故」 「瀬戸内海」の 5 つである.これらの高出現頻度ノードの描画さ れた部分を図 8 に示す.高頻度のノードの表す単語や、それら



図6 情報変化グラフ中の情報推移出現例

がどのように接続されているかを見ることで、文書の対象とする テーマの基本的な情報を知ることができる.

まず,高頻度のノード間の関係について分析する. 図 8 に点線四角で強調したように,「海上自衛隊」に修飾される「輸送艦 おおすみ」と「釣り船とびうお」が,Activity 関係で「衝突事故」に 接続されている.加えて,「衝突事故」は「瀬戸内海」と Location 関係をもっている.以上のことから,分析対象の文書が扱うテー マは,瀬戸内海で発生した海上自衛隊の輸送艦と釣り船の衝 突事故であることが分かる.

次に、高頻度ノードに接続される低頻度ノードにどのようなものがあるか分析する。例えば、「衝突事故」は「横切る形」で修飾され、Time関係で「1月」や「15日」と接続されている。また、「輸送艦おおすみ」は「右旋回している」を経由して「釣り船とびうお」や「衝突事故」と接続されている。こうした関係から、事故が発生したのは1月15日であることや、輸送艦は右旋回中に釣り船と衝突したという情報が出現したと分かる。

ここで、ノードの色によって表現される単語の出現時期に注 目する. 図9では、「釣り船とびうお」に修飾された「2人」は、「搬 送される」と「死亡する」にActivity関係のエッジをもつ. これらの ノードの色によれば、「搬送される」は時系列の最初期のみに出 現する.そして、「死亡する」は「搬送される」の出現直後と時系 列の末期周辺に出現している.これは、図6で示したような情報 変化が、図9での点線矢印のようなノード色の変化として現れて いると解釈できる.また、「とびうおの乗員が死亡した」という情報 は、始めに出現した後はしばらく文書中で取り上げられなかっ たことが分かる.実際,分析対象の記事では,この情報が出現 しなかった1月17日から28日までの記事では、1記事を除い て調査で判明した事実や関係者の証言のみを述べていた. 一 方で、例外である1記事には、それらに加えて「死亡した乗員2 人のうち1人の通夜が営まれた」という情報が含まれていた.こ の記事では「2人」という単語を用いずに以上の情報を述べてい たため、ノードの出現頻度計算に誤りが生じたと考えられる.

このように、提案手法による時系列的な出現頻度を用いたサ ブグラフでは、高頻度で出現する情報と、それに付随する低頻 度な情報を可視化することができる。特に、高頻度ノードが示す 単語や、それらの関係は、3.3節で定義した重要情報を表すも のであった.また、低頻度ノードは、重要情報に関連する情報 が、時系列の中でどのように出現したかを示した.ただし、低頻 度の情報にはある一時期のみ注目されたものや後に誤りだと判 明したものなど、様々な種類があることが推測される。従って、 ユーザは低頻度の情報が事実を正しく表しているとは限らない ことに注意する必要がある。



5. 課題

4 章での可視化結果より得られた,提案手法における課題に ついて述べる.提案システムで利用する事象抽出手法に関する 問題の1つとして,グラフ中におけるノードの多さが挙げられる. その原因としては、「こと」や「する」等、それ自体では多くの情報 を持たないノードの存在がある.こうしたノードは他のノードと接 続されやすく、グラフ操作による削除が難しい一方,表示が煩 雑になる.もう1つの問題として、文中における単語の省略に起 因する、エッジの設定ミスがある.先行研究と異なり、本論文で はグラフの分割を行うために1グラフ当たりの文書数が少なくな りやすい.そのため、グラフ中の情報欠損が比較的発生しやす く、分析の支障となる可能性がある.

これらへの対策としては、田代らの事象抽出手法を改良する ことによる、ノードとなる単語のフィルタリングや、照応解析による 情報の補完などが考えられる.

(1) 情報変化グラフにおける課題

情報変化の可視化では、連続する2期間における短期的な 変化を計算している.そのため、初期に錯綜していた情報が次 第に訂正され、統一されていくといった、長期的視点での情報 変化を可視化することができない.また、情報変化を検知する 以上の処理を行わないため、情報の矛盾や人物状態の推移等、 情報変化の特徴を判別することができない.

これらの改善策として,情報変化の対象や変化パターン等, 長期的変化の可視化に有用な情報の利用を検討中である.

(2) 重要情報グラフにおける課題

提案した重要情報可視化手法は、ノードがどのような単語を 表すかを考慮せず、出現頻度とノード間の距離によってサブグ ラフを抽出する. そのため、特にサブグラフの末端ノードにおい て、それがどのような事象に出現したのか解釈するのが難しい. 例えば、図 7 のグラフでは「衝突事故」が「に関し」と接続されて いるが、「に関し」の対象となるノードは抽出できていない.

この問題の解決策として、ノード間の関係やエッジの向きを考慮することで、ノードの削除や追加取得を行うことが考えられる.

6. おわりに

本論文では,時系列的に複数作製される文書群における情報の変化と,文書の対象とするテーマに関連する重要情報の



図8 高出現頻度ノードに関する部分



図9 時間経過に関する重要情報例

可視化手法を提案した.その後,実際の新聞記事に対して提 案手法を適用し,可視化結果を評価した.

実際の文書群に適用した結果,提案手法は短期間での情報 の新規出現や時間経過による変化を可視化できることを確認し た.また,重要情報として文書のテーマを表すような情報や,そ れに関する様々な情報を可視化できることを示した.

今後の予定としては、人物等に注目することによる長期的な 情報変化の可視化や、ノード間の関係などを考慮することによ る重要情報のより高精度な抽出などを考えている.

参考文献

- [田代 2013] 田沼和浩, 王冕, 越川兼地, 森田武史, 長野伸 一, 清雄一, 中川博之, 田原康之, 川村隆浩, 大須賀昭彦: "Linked Data を用いたソーシャルメディア×マスメディアの 比較実験", 人工知能学会全国大会論文集, 2013.
- [伊藤 2012] 伊藤正彦, 吉永直樹, 豊田正史, 喜連川優: "係 り受け解析を用いたブログユーザの行動・興味に関する時 系列推移 3 次元可視化システム", 電気情報通信学会論文 誌 D, Vol. J95-D, No.7, pp.1454-1466, 2012.