

イベントの社会的影響度を用いた 伝記生成のための内容選択

川添 正太郎 徳永 健伸
東京工業大学 情報理工学院
{kawazoe.s.aa@m,take@c}.titech.ac.jp

概要

本研究では、人物に関するイベント集合を入力とし、イベントの社会的影響度によって伝記生成のための内容選択を行う。実験では、人物に言及するニュース速報ツイートをイベントとし、その「いいね数」を社会的影響度として、重要なイベントの選択を行った。Wikipediaの記事の出典の集合に対して評価を行い、社会的影響度によって重要なイベントを選択できることを示す。

1 はじめに

自然言語生成 (Natural Language Generation; 以下 NLG と略す) は、自然言語で書かれたテキストを生成するコンピュータシステムを扱う、計算言語学の一分野である [1]。内容選択 (content selection) は、NLG の部分問題であり、入力から言及すべき重要な情報を選択する過程である。たとえば、株価の概況生成において株価のチャートが暴騰している部分を選択したり、ニュース記事の要約において重要な出来事が記述されている部分を選択する。

しかし、NLG の入力は何なのかは曖昧であり [2, 3]、また重要度とは何を意味するのかも曖昧である [4, 5]。本研究では、NLG の一種である伝記生成 (biography generation) において、これらの問いに部分的に答えることを目指す。伝記は人物を紹介するテキストである。たとえば、Wikipedia の「バラク・オバマ」の記事 [6] には、経歴、大統領時代の仕事、家族のことなどが記述されている。以下に記事に含まれている文を示す。

2009年10月9日にノルウェー・ノーベル委員会はオバマの「核無き世界」に向けた国際社会への働きかけを評価して2009年度のノーベル平和賞を彼に受賞させることを決定したと発表した。

この文は、人間によって、どのような入力から、どのような重要度にもとづいて、選択されたのだろうか。ノーベル平和賞を受賞することは、オバマにとって重要なイベントであることは明らかだろう。このように、人間が情報の内容を理解してそれを重要だと判断したとき、そのような重要度のことを**内在的重要度** [5, 7] と定義する。

2 問題

伝記生成の既存研究の問題は、内在的重要度による内容選択を行っていないことである。これは、裏を返せば、内在的重要度を用いなくても内容選択ができてしまうような入力を用いているということでもある。一方、オバマの例からわかるように、人間は内在的重要度による内容選択を行っている。本研究では、内在的重要度によって内容選択を行うことにより、人間の内容選択を再現することを目的とする。既存研究におけるこの問題は、以下で示すように、二つの場合に分けられる。

第一に、入力において情報の絞り込みがほとんど終わっている場合である。たとえば、WikiBio [8] は、Wikipediaの人物に関するインフォボックス (属性と値の組から構成された表) を入力としている。しかし、インフォボックスの時点でその人物に関する重要な情報がまとまっており、人物に関するある種の要約となっている。WikiSum [9] は、Wikipediaの記事 (ただし、人物に限らない) における出典の集合が入力として与えられ、その記事の最初のセクションを出力する。いずれの場合も、入力の時点で内容選択はほとんど終わっている。

第二に、表層的重要度によって内容選択を行っている場合である。表層的重要度とは、内容を見なくても計算できるような重要度のことである。たとえば、事前に指定された型にあてはまるような情報を重要とみなすトップダウンな手法は、内容

を見なくても重要度を計算できる。Zhou ら [10] は、DUC2004 [11] のタスク 5 において、ニュース記事の文に対し、10 クラス (fame, education, nationality, work など) への分類を行っている。これは、伝記にある種の型を想定し、伝記的かどうかを判断している。しかし、トップダウンな手法は人物個人の背景を考慮せず、人物一般にあてはまる形で重要度を計算している。たとえば、オバマがノーベル平和賞を受賞することの内在的重要度は、オバマ個人の背景 (核廃絶に対する働きかけへの評価、あまり成果を出していない中での早すぎる受賞、そのことへの批判など) に依存するだろう。

また、入力における情報の位置や頻度を用いるボトムアップな手法も、内容を見なくても重要度を計算できる。つまり、内在的重要度と相関する表層的な情報を用いることによって内容選択を行っている。これは、入力を作成した人間が重要だと考えた情報は、何かしらの表層的な特徴をもって入力に現れるからである。たとえば、WikiBio では、位置と内在的重要度が相関している [12]。つまり、入力であるインフォボックスはその作成者が重要な情報を上に配置する傾向にあるため、上の方にある情報ほど、正解テキストで言及されやすい。DUC2004 のタスク 5 を解いている代表的な手法 [10, 13] は、単語頻度を利用して重要度を計算する。つまり、ニュース記事に対して単語頻度ベースのスコアを計算し、スコアが高い単語が含まれる文を抽出している。しかし、入力に重要な情報が頻出するのはコーパスがそのような性質を偶然持っただけであり、また重要な単語が含まれている文が重要とは限らないため、内在的重要度を扱っているとはいえない。

3 仮説

以上の問題が生じた原因は、伝記生成の内容選択を、計算機で解くことができる簡単なタスクへと落とし込んだことによって、実際に人間によって行われている内容選択から乖離してしまったことだと思われる。したがって、人間によって書かれた伝記が、どのような入力から、どのような重要度にもとづいて、生成されているかを考察する必要がある。

伝記生成の内容選択では、社会的影響度のあるイベントが発生するたびに、そのイベントの内容が伝記に追記される。たとえば、オバマがノーベル平和賞を受賞したというイベントが起こると、それが重要だと判断した Wikipedia 編集者が、その内容を伝

記に追記したと考えられる。社会的影響度が高いイベントほど、記事に追記される確率が高くなる。このことから、入力を人物に関するイベントの集合とし、重要度をイベントの社会的影響度と考えるのが妥当だろう。イベント集合は、既存研究に比べ、情報の絞り込みがあまり終わっていない。また、イベントの社会的影響度は内在的重要度と直接関係すると考えられる。実験によって検証すべき仮説は以下の通りである。

イベント集合から伝記を生成する際の内容選択において、社会的影響度は、他の表層的な重要度よりも、重要な情報を選択する能力が高い。

4 実装

本研究で実装した内容選択システムは、ある人物 i の URI (たとえば、http://dbpedia.org/resource/Barack_Obama) がクエリとして与えられたとき、入力すなわち人物 i に関するイベント集合 X を収集し、そこから人物 i の伝記で言及されるべきイベントの集合 $\tilde{Y} (\subseteq X)$ を選択する。

ツイッターによるイベントの収集 入力すなわち人物 i に関するイベント集合 X には、ニュース速報をツイートするツイッターアカウント CNN Breaking News (@cnnbrk) のうち、人物 i に言及しているすべてのツイート (リツイートを除く) を用いる。ただし、DBpedia Spotlight [14] を用いてツイートのテキストに対してエンティティリンキングを行うことで、ツイートが人物へ言及しているかどうかを判定している。ツイート t が持つ情報のうち、本研究で利用するのは、(1) テキスト s_t , (2) 日付 d_t , (3) いいね数 n_t の三つである。社会的影響度にはいいね数を用いる。ここで、イベントの社会的影響度を計算するのは難しいため、いいね数が与えられているとしている。

スコアリング 各ツイート $t \in X$ に対してスコアを割り当て、スコアが大きい順に並べた上位 p (%) のツイートを取得して内容選択を行う。スコアリングの方法として以下の 7 種類を検討した。

- **Random** は、ランダムにスコアを割り当てる。
- **Newest** は、新しいツイートほど良いスコアを割り当てる。一般に、新しいツイートほどいいね数が多くなるため、時間の影響を差し引くためにこれをベースラインに加える。
- **Oldest** は、古いツイートほど良いスコアを割り

当てる。

- **TF-IDF** は、ツイートのテキストに含まれる各単語の TF-IDF 値の平均をスコアとする。TF-IDF 値の計算における文書は、各人物に対し、その人物が言及されているすべてのツイートのテキストを結合したものである。このとき、一文書が一人物に対応する。ある人物に対して、TF が高い単語は、ツイート中で人物名と頻繁に共起し、IDF が高い単語は、他の人物名とはあまり共起しない。したがって、TF-IDF 値が高い単語は、その人物に関するキーワードであると解釈できる (付録 A)。
- **LexRank** では、LexRank [15] で計算されたスコアを用いる。ツイートのテキストを文とみなし、ツイートのテキストを結合したものを文書とみなして、continuous LexRank を適用する。
- **Like** は、いいね数をスコアとする。
- **LikeRatio** は、いいね数を正規化したスコアを用いる。分母は、当該ツイートの月における CNN Breaking News のすべてのツイートの平均いいね数である。

5 評価

30 個以上のツイートで言及されている 72 人の人物を評価の対象とする (付録 B)。

人物 i に対し、正解のイベント集合 Y には、人物 i に関する Wikipedia の記事の出典の集合を用いる。出典には、出典となった記事のタイトルや URL、記事が書かれた日付などが含まれている (付録 C)。出典 c が持つ情報のうち、本研究で利用するのは、(1) 記事タイトルのテキスト s_c 、(2) 記事が書かれた日付 d_c の二つである。ただし、 $d_{oldest} \leq d_c \leq d_{newest}$ となる出典のみを用いる。ここで、 d_{oldest} および d_{newest} はそれぞれ人物 i に関するツイートのうち最も古い日付および最も新しい日付である。

評価指標には、精度 (precision) と再現率 (recall) を用いる。ランキング上位 p (%) のイベントを選択するとき、人物 i に対し、精度 P^i および再現率 R^i を計算する。そして、すべての人物にわたって精度および再現率の平均 $\frac{1}{N} \sum_{i=1}^N P^i$ および $\frac{1}{N} \sum_{i=1}^N R^i$ をとる。ただし、 $N (= 72)$ は人物の数である。上記の手続きを、 p を 10 から 100 まで 10 刻みで変化させて行い、平均の PR 曲線 (precision-recall curve) を描く。

精度と再現率の計算の際、ツイートと出典がイベ

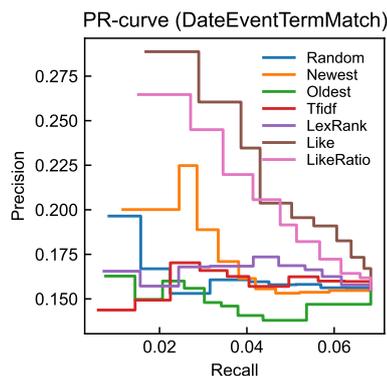


図 1 PR 曲線

ントとして一致しているかどうかの判定を実装する必要がある。イベントは「いつ・どこで・誰が・誰に・何をした」と表現できる [16, 17]。「誰が」あるいは「誰に」は、人物 i になる可能性が高い。「いつ」はツイートおよび出典の日付であると考えられる。「どこで」と「何をした」は、テキストに記載されている可能性がある。これらを踏まえ、一致判定方法 **DateEventTermMatch** を説明する。

DateEventTermMatch では、日付 (年月日) が一致し、かつイベント用語が一致したら、イベントも一致したと考える。イベント用語は、動作や事件の発生を特徴付ける単語である [16]。Li ら [16] にならない、イベント用語は動詞および動作名詞 (action noun) とする。この一致判定では、「いつ (日付)」「誰が (人物)」「何をした (イベント用語)」の一致に着目している。人物 i に関する精度 P^i は以下の式で計算する。

$$P^i = \frac{\sum_{t \in \tilde{Y}} \mathbb{1}(\exists c \in Y (d_t = d_c \wedge ET_t \cap ET_c \neq \phi))}{|\tilde{Y}|}$$

ただし、 ET_t はツイート t のテキストに含まれるイベント用語の集合 (ET_c も同様)、 ϕ は空集合、 $\mathbb{1}$ は指示関数 (引数の条件が真ならば 1、偽ならば 0 を返す関数) である。再現率も同様に計算する。

6 結果

PR 曲線を図 1 に示す。図から、いいね数を用いた手法 (**Like** と **LikeRatio**) が、他の手法を上回っていることがわかる。これは、イベント集合からの内容選択において、社会的影響度のほうが、他の表層的な重要度よりも重要な情報を選択する能力が高いことを示している。

表 1 に **Like** と **TF-IDF** の比較を示す。上のツイートは、 $p = 20$ (%) において **Like** で真陽性となるが

表1 Likeで取得できたツイート(上)とTF-IDFで取得できたツイート(下).

ツイート(いいね数)/出典

An independent autopsy into the **death** of **George Floyd** found that he died from "asphyxiation from sustained pressure," which contradicts the county medical examiner's preliminary report (17415) / Independent autopsy reveals George Floyd died from 'asphyxiation' as lawyers call for first-degree murder charges

Serena Williams wins 5th Wimbledon title. (73) / Wimbledon 2012 – Serena Williams stretched to three sets, wins 5th title

表2 Likeで偽陽性となったツイート.

ツイート(いいね数)

Trump appears to be skipping a side-event at the G20 virtual summit focused on pandemic preparedness. The President has just arrived at his golf course in Virginia. (45540)

TF-IDFで偽陰性となったものである。太字は人物に関してTF-IDF値が上位30位以内の単語であり、これは人物に関するキーワードとみなせる。キーワードとなっているのは人物名を除くと“death”のみである。これは、George Floydは白人警官に殺害されて死亡したことがたびたび注目されているためである。それ以外の単語は専門用語などが多く、キーワードにはならないため、TF-IDFで偽陰性となったと考えられる。しかし、Likeで真陽性となったのは、このイベントの社会的影響度が高いからだと考えられる。

一方、下のツイートは、TF-IDFで真陽性となるがLikeで偽陰性となったものである。Serena Williamsはテニス選手であり、“Wimbledon”はテニス大会の名前である。この大会は定期的に行われており、Serena Williamsは何度も優勝争いを行っているため、“wins”, “Wimbledon”, “title”は、Serena Williamsに関するツイートに何度も出現する。よって、これらの単語のTF-IDF値が高くなり、このツイートがTF-IDFで真陽性となったと考えられる。一方、これはLikeで(そしてLikeRatioでも)偽陰性となった。原因としては、このツイートの日付は2012年7月7日と古く、いいね数と時間には相関があるため、いいね数が小さくなったのだと考えられる。これは、いいね数による内容選択は、ツイッターの性質上、古いツイートを取得しにくいことを示している。

$p = 20$ (%)においてLikeで偽陽性になったツイートのうち、スコア(いいね数)が大きいものを表2に示す。このツイートは、Donald Trumpがコロナ対

表3 社会的影響度が高いイベント用語の例.

ツイート/出典

Israeli Prime Minister Benjamin Netanyahu will be indicted on corruption **charges**, pending a final hearing / Israeli Prime Minister Benjamin Netanyahu to be indicted on corruption charges

Former Minneapolis Police Officer Derek Chauvin is **found** guilty of murder and manslaughter in the death of George Floyd / Derek Chauvin found guilty of George Floyd's murder

British Prime Minister Boris Johnson has **tested** positive for coronavirus. He says on Twitter that he is isolating with "mild symptoms." / PM Boris Johnson tests positive for coronavirus

応を放棄してゴルフに行き遊んでいるという印象を与えるため、「炎上」していると考えられる。しかし、このツイートは偽陽性となった。これは、炎上しているツイートはいいね数が大きくなるが、それが必ずしも伝記にとって重要であることを意味しない、ということを示唆している。

$p = 20$ (%)において真陽性となったツイートのうち、NewestよりLikeに多く出現するイベント用語として、“charge”, “find”, “test”が挙げられる(付録D)。これらのイベント用語に対応する例を表3に示す。“charge”は「起訴」，“find”は“found guilty”つまり「有罪になる」，“test”は“has tested positive for coronavirus”つまり「コロナ陽性となる」という意味で、それぞれ用いられる。いずれも、炎上しやすいイベントを表しているため、いいね数が大きくなりやすい。そして、表に示したツイートはすべて真陽性である。これは、社会的影響度によって伝記生成の内容選択ができることを示している。

7 おわりに

本研究では人物に言及するCNN Breaking Newsのツイートをイベントとし、そのいいね数を社会的影響度として、伝記の内容選択を行った。実験により、社会的影響度は、他の表層的重要度よりも、重要なイベントを選択する能力が高いことを示した。

今後の課題として、(1)ツイッターによるイベント収集の網羅性を向上させること、(2)トップダウンな内容選択と社会的影響度によるボトムアップな内容選択を組み合わせること、(3)伝記生成にとって適切なイベントの単位(atomic event) [17, 18]を明確にすること、(4)ツイートのテキストが与えられたとき、そのいいね数を予測すること、(5)多数決による重要度の計算 [19]を回避することが挙げられる。

参考文献

- [1] Ehud Reiter and Robert Dale. **Building natural language generation systems**. Cambridge university press, 2000.
- [2] David D. McDonald. Issues in the choice of a source for natural language generation. **Computational Linguistics**, Vol. 19, No. 1, pp. 191–197, 1993.
- [3] Roger Evans, Paul Piwek, and Lynne Cahill. What is NLG? In **Proceedings of the International Natural Language Generation Conference**, pp. 144–151, Harriman, New York, USA, July 2002. Association for Computational Linguistics.
- [4] Maxime Peyrard. A simple theoretical model of importance for summarization. In **Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics**, pp. 1059–1073, Florence, Italy, July 2019. Association for Computational Linguistics.
- [5] Markus Zopf. **Towards Context-free Information Importance Estimation**. PhD thesis, Technische Universität, Darmstadt, August 2019.
- [6] バラク・オバマ - wikipedia. <https://ja.wikipedia.org/w/index.php?title=%E3%83%90%E3%83%A9%E3%82%AF%E3%83%BB%E3%82%AA%E3%83%90%E3%83%9E&oldid=87303693>.
- [7] Kai Hong and Ani Nenkova. Improving the estimation of word importance for news multi-document summarization. In **Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics**, pp. 712–721, Gothenburg, Sweden, April 2014. Association for Computational Linguistics.
- [8] Rémi Lebret, David Grangier, and Michael Auli. Neural text generation from structured data with application to the biography domain. In **Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing**, pp. 1203–1213, Austin, Texas, November 2016. Association for Computational Linguistics.
- [9] Peter J. Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. Generating wikipedia by summarizing long sequences. In **6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings**. OpenReview.net, 2018.
- [10] Liang Zhou, Miruna Ticea, and Eduard Hovy. Multi-document biography summarization. In **Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing**, pp. 434–441, Barcelona, Spain, July 2004. Association for Computational Linguistics.
- [11] Paul Over and James Yen. An introduction to duc-2004. **National Institute of Standards and Technology**, 2004.
- [12] Tianyu Liu, Kexiang Wang, Lei Sha, Baobao Chang, and Zhifang Sui. Table-to-text generation by structure-aware seq2seq learning. In Sheila A. McIlraith and Kilian Q. Weinberger, editors, **Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence**, pp. 4881–4888, New Orleans, Louisiana, USA, 2018. AAAI Press.
- [13] Fadi Biadisy, Julia Hirschberg, and Elena Filatova. An unsupervised approach to biography production using Wikipedia. In **Proceedings of ACL-08: HLT**, pp. 807–815, Columbus, Ohio, June 2008. Association for Computational Linguistics.
- [14] Joachim Daiber, Max Jakob, Chris Hokamp, and Pablo N. Mendes. Improving efficiency and accuracy in multilingual entity extraction. In **Proceedings of the 9th International Conference on Semantic Systems (I-Semantics)**, 2013.
- [15] Günes Erkan and Dragomir R. Radev. Lexrank: Graph-based lexical centrality as salience in text summarization. **J. Artif. Intell. Res.**, Vol. 22, pp. 457–479, 2004.
- [16] Wenjie Li, Mingli Wu, Qin Lu, Wei Xu, and Chunfa Yuan. Extractive summarization using inter- and intra- event relevance. In **Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics**, pp. 369–376, Sydney, Australia, July 2006. Association for Computational Linguistics.
- [17] Vasileios Hatzivassiloglou and Elena Filatova. Domain-independent detection, extraction, and labeling of atomic events. 2003.
- [18] Elena Filatova and John Prager. Tell me what you do and I’ll tell you what you are: Learning occupation-related activities for biographies. In **Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing**, pp. 113–120, Vancouver, British Columbia, Canada, October 2005. Association for Computational Linguistics.
- [19] Ani Nenkova and Rebecca Passonneau. Evaluating content selection in summarization: The pyramid method. In **Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004**, pp. 145–152, Boston, Massachusetts, USA, May 2 - May 7 2004. Association for Computational Linguistics.
- [20] Barack obama - wikipedia. https://en.wikipedia.org/w/index.php?title=Barack_Obama&oldid=1064944028.

A キーワードの例

表4 TF-IDF 値が上位 20 位となる単語. i は人物の URI (<http://dbpedia.org/resource/> を省略) である.

i	キーワード
Donald_Trump	trump, president, says, donald, house, white, campaign, follow live, russia, updates, north, cnn, michael, korea, impeachment mueller, cohen, new, order
Barack_Obama	obama, president, says, barack, watch, cnn, isis, iraq house, live, white, military, romney, make, ukraine, address jobs, syria, iran, signs
Mitt_Romney	romney, mitt, santorum, gingrich, paul, projects, primary, votes cnn, win, obama, electoral, gop, caucuses, counted, says delegates, poll, jobs, 11

B 人物の例

100 個以上のツイートで言及されている人物を表 5 に示す. i は人物の URI (<http://dbpedia.org/resource/> を省略), $|X|$ はツイート数, $|Y|$ は出典数, \bar{n}_t は平均いいね数, d_{oldest} は最古の日付, d_{newest} は最新の日付である.

表5 $|X|$ が 100 以上の人物.

i	$ X $	$ Y $	\bar{n}_t	d_{oldest}	d_{newest}
Donald_Trump	3036	809	1360	2012-02-02	2021-10-29
Barack_Obama	1769	427	331	2008-01-04	2021-01-11
Mitt_Romney	294	442	294	2008-01-16	2020-11-08
Hillary_Clinton	280	373	650	2007-11-30	2020-10-09
Vladimir_Putin	210	364	471	2011-09-24	2021-10-19
Joe_Biden	209	430	5671	2008-08-23	2021-11-01
Bill_Clinton	167	247	689	2008-02-06	2021-10-17
Pope_Francis	152	319	1025	2007-03-02	2021-10-14
Rick_Santorum	112	195	42	2011-06-14	2018-03-25
James_Comey	111	235	910	2015-12-16	2020-09-30
Kim_Jong-un	110	164	807	2010-09-28	2021-06-30

C 出典の例

On August 23, 2008, Obama announced his selection of Delaware Senator Joe Biden as his vice presidential running mate.

上の文 [20] の出典は以下ようになる.

```
{{cite news |access-date = September 20, 2008
|url=https://www.nytimes.com/2008/08/24/us/politics/24biden.html
|last1 = Nagourney |first1 = Adam |first2 = Jeff |last2 = Zeleny |work = The New York Times
|date = August 23, 2008 |title = Obama Chooses Biden as Running Mate
|archive-url= https://web.archive.org/web/20090401222653/http://
www.nytimes.com/2008/08/24/us/politics/24biden.html|archive-date=April 1, 2009 |url-status=live }}
```

D イベント用語の例

表6 $p = 20$ (%) において真陽性となったツイートのイベント用語. 括弧内は頻度.

順位	Newest	TF-IDF	Like
1	say (32)	win (34)	say (31)
2	win (15)	say (21)	charge (14)
3	sentence (9)	president (7)	win (12)
4	charge (8)	charge (6)	president (10)
5	suspend (6)	fire (5)	sentence (9)
6	president (6)	murder (4)	resign (6)
7	resign (6)	session (4)	find (6)
8	leave (5)	resign (3)	test (6)
9	call (5)	reopen (3)	suspend (5)
10	murder (5)	drop (3)	endorse (5)