

映画レビューの時系列による評価視点の抽出と可視化

岑 天霞^{†1} 渡邊 英徳^{†2}

本研究では、映画レビューの本文を投稿時間とともに分析することを目的とする。現在、映画情報サイトには、映画のレビューに関する大量のデータが蓄積されている。このデータからは、一般の人々による映画作品に対する評価視点を読み取ることができる。そして社会環境の変化に応じて、同一の作品に対する評価視点は変化し続ける。しかし、既存の手法では、時間軸に沿った評価視点の変化については考慮されていない。そこで、本研究では感染症・災害・戦争を題材とした映画のレビューを対象として、時間軸に沿った評価視点の変化を抽出し、可視化する手法を提案する。本稿では、パンデミックをテーマにした映画『コンテイジョン』のレビューを年別・波別の時期ごとに分割し、単語の特徴度尺度を用いて各時期のレビューにおける特徴的な名詞を評価視点として抽出・可視化する。さらに、抽出結果を当時の社会的な背景と照らし合わせて、提案手法の有効性を評価する。その結果、時系列に基づく映画レビューの評価視点を抽出する手法が有効であり、社会の環境の変化によって映画作品への評価も変化していることが示された。

Extraction and Visualization of Information from Movie Reviews Based on Temporal Analysis

TIANXIA CEN^{†1} HIDENORI WATANABE^{†1}

The purpose of this study is to analyze the main text of a movie review with its submission time. Currently, movie information sites accumulate a large amount of data on movie reviews. From this data, we can understand the viewpoint of the general public's evaluation of films. In response to changes in the social environment, the evaluation of the same work continues to change. However, the existing methods do not take into account the change of evaluation viewpoint along the time axis. In this study, we propose a method to extract and visualize changes in the evaluation viewpoint along the time axis for reviews of movies on infectious diseases, disasters, and wars. In this paper, we divide the review of the pandemic-themed movie "contagion" into periods by year and wave, and extract and visualize the characteristic nouns in the review at each period using the feature scale of words. In addition, the effectiveness of the proposed method is evaluated by comparing the extracted results with the social background of the time. As a result, it is shown that the method of extracting the evaluation viewpoint of film reviews based on time series is effective, and that the evaluation of film works also changes with changes in the social environment.

1. はじめに

本研究は、ネット上における時間情報を持つ映画レビューの変化を捉えるために、時系列による映画レビューの評価視点を抽出と可視化することである。

現在、ネットの普及によって、誰でも簡単にネットを通じて発信・受信できるようになっている。特に、映画情報サイトには、作品に対する大衆の議論を読み取れる莫大なレビューデータが長期に蓄積されている。こうした映画レビューは、レビューのテキスト以外に、レビューの投稿日時も付与されており、時系列に沿って分析することができるデータである。また、戦争・パンデミックなどのイベントにより、状況が激変している現代社会においては、フィクション作品のストーリーが現実の世界で再現されることもある。

また、社会環境の変化によって、同じ映画に対する議論が変化し続けると考えられる。例えば、2020年「COVID-19」の流行によって、パンデミックをテーマにしたフィクション作品を改めて見直すようなことが起きていた。

そこで本研究は、パンデミックをテーマにした映画『コンテイジョン (Contagion)』(2011年)のレビューを対象に実験を行い、映画レビューの変化を捉える可能性を考察する。

2. 先行研究

2.1 評価情報抽出

映画のレビューを分析する既存研究では、評価視点を抽出・要約する方法が多く検討されている。例えば、Huら[1]は、製品の評価視点を出現頻度によって抽出し、各視点に対して肯定的・否定的な顧客レビューの数を特定している。また、小林らは、デジタル商品レビューにおける評価視点と評価表現を組み合わせた共起パターンに基づいて、特徴的なユーザの意見とその根拠を抽出している[2]。しかし、これらの手法では、人間の判断に基づく評価視点・評価表現の辞書を作成する必要がある。本研究のような大規模データを扱う場合には適さない。

一方、内山ら(2004)は、8つの統計的尺度を比較検討し、特定分野の英語文書(TOEIC 試験模擬問題)から特徴的な単語を抽出する有効性を確かめられた[3]。乾ら(2013)は、宿泊

^{†1} 東京大学大学院 情報学環・学際情報学府
The University of Tokyo

施設のロコミから評価視点を抽出するために、内山らが提案した対数尤度比(Log-Likelihood Ratio, LLR)に基づいて、特徴的な評価視点の抽出手法を提案した[4].

LLRは、ある特定の評価対象において観測される評価視点の確率と、全ての評価対象において観測される評価視点の確率の比の対数として算出される。従って、レビュー群において、一般的な評価視点であればLLRの値が小さくなり、独特な評価視点であればLLRの値が大きくなる。つまり、LLRは評価視点の辞書を必要とせず、レビュー群から特徴的な評価視点を直接抽出することを可能にする。

3. 提案手法

2章で述べたように、映画レビューにおける評価視点のバリエーションは多岐に渡っており、作品そのもの以外について話題も含まれる。従って、人力による評価視点の辞書作成は困難である。そこで、有意な評価視点を抽出するために、特徴度尺度に沿った順位付けを行なう。

内山ら(2004)は特定分野(TOEIC 試験模擬問題)の英語文書から特徴的な単語を抽出するために、8つの統計的尺度を提案している[5]。乾ら(2013)は、宿泊施設レビューから評価視点を抽出するために、内山らの尺度のうち対数尤度比(Log-Likelihood Ratio, LLR)に基づいて、単語を特徴度尺度に沿ってランク付けしている。LLRでは、レビューに含まれる全ての単語をもとに特徴的な単語が抽出される。一方、映画レビューにおける評価視点は専ら名詞で構成されている。そこで本研究では、特徴度尺度LLR映画のレビューの評価情報抽出へ適用する際に、便宜的に名詞のみを評価視点として扱う。LLRの算出法を応用し、品詞を名詞に限定することにより、計算コストを節約する。本研究は、この特徴語尺度をRFRと定義する。

「特定の時期」及び「全ての時期」についての名詞(評価視点)の相対出現頻度RF(Relative frequency)をそれぞれ求め、次いでそれらをもとに、評価視点の相対出現頻度比RFRを算出する。

特定時期の映画 M_j のレビューにおける名詞 t_i の相対出現頻度RFは、

$$rf_{i,j} = \frac{n_{i,j}}{n_j} \quad (1)$$

算出される(1)。ここで t_i 、 i 、 $n_{i,j}$ 、 j と n_j は

- t_i : ある名詞
- i : 特定の名詞 t_i の番号
- $n_{i,j}$: 映画 M_j のレビューにおける t_i の出現回数
- j : 特定時期の映画 M_j の番号
- n_j : 映画 M_j の特定時期のレビュー件数である。

また、全ての時期のレビューにおける名詞の相対出現頻度RFは、

$$rf_{i,j}' = \frac{N_i}{N} \quad (2)$$

で算出される(2)。

ここでは、

- N_i : 「全ての時期」の映画レビューにおける t_i の出現回数
- N : 「全ての時期」のレビュー件数の総数である。

RFとRF'の比は、

$$\frac{rf_{i,j}}{rf_{i,j}'} = \frac{n_{i,j} * N}{n_i * n_j} \quad (3)$$

対数をとって正規化し、常に正とするために真数に1を加えると、評価視点の相対出現頻度比RFRが算出できる

$$rfr_{i,j} = \log\left(\frac{n_{i,j} * N}{n_i * n_j} + 1\right) \quad (4)$$

4. 応用実験

本章では、パンデミックを題材とした映画のレビューを対象として、3章で提案した特徴度尺度RFRを適用し、映画のレビューから各時期の特徴的な評価視点を抽出する。時間軸に沿った抽出結果を当時の社会的な背景と照らし合わせて、各時期の評価視点の変化を考察し、提案手法の有効性を評価する。

4.1 感染症を題材とした映画のレビュー分析

まず、実験結果を捉えやすいために、急激に変化している社会事件「COVID-19」に関連するパンデミックをテーマにした映画『コンテイジョン(Contagion)』(2011年)のレビューを実験データとする。『コンテイジョン』は、スティーブン・ソダーバーグ監督による高い死亡率をもつ感染症の脅威とパニックを描く映画である。「COVID-19」が流行している現実に酷似しているため、予言映画と言われている。そこで、本研究は、日本の映画情報サイト「Filmarks映画」[5]における『コンテイジョン』のレビューデータ(2022年2月)を取得し、時系列によるレビュー可視化分析を行う。

4.2 年ごとと波ごとの抽出

まずは、RFRを応用して、映画レビューを時期別に区切りし、RFRを用いて各時期のレビューの特徴的な評価視点の抽出する。適切な時間軸区切りを見つけるために、年別・波別に分けてそれぞれの特徴語抽出した。

本稿では、まず、2016年から2021年までの6年間で年ごとに特徴語を抽出した(表1)。

また、「COVID-19」が日本で流行する波の第1波から第5波まで、波ごとに特徴語を抽出した(表2)。

各波の時間帯は以下の通りである。

第1波: 2020 3~5月

第2波: 2020年7~9月

第3波: 2020年11月~2021 02月

第4波：2021年3~5月

第5波：2021年7~9月

4.3 結果の考察

表1は、年別による『コンテイジョン』レビューの特徴的な評価視点の抽出結果を示したものである。

表1から以下のことがいえる。

「COVID-19」が流行する前は、映画そのものに関連する単語が特徴語として抽出されている。例えば、2016年においては、「派手」「演技」「俳優」などが映画製作に関する単語が特徴であった。一方、2020年から「COVID-19」が流行したため、レビューには、映画そのものに関連する言葉よりも、社会の現状に関連する評価視点が増えている。例えば、2020年の「緊急事態宣言」により「自粛」という言葉が抽出された特徴語に登場し、2021年のワクチン接種により「接種」という言葉が登場し。

一方、2017年の「ゾンビ」、2018年の「ホラー」というワードが特徴語となった理由に関して、一見では手がかりがない。「ゾンビ」「ホラー」を含むレビューの一部を以下に抜粋すると、これらの単語の文脈がわかる。

「…日本で去年流行ったB級ゾンビ作品みたいにバアアと感染して終わりっ！てな感じじゃなく、しっかりしてたラスト1分が何とも言えない(;ω;)」

「なぜか完全にゾンビものだと勘違いしてて、「ゾンビ登場までずいぶん引っ張るなあ…」なんて思ってた映画が終わってました…」

「…こうゆうパンデミック系、わたしはホラーより怖いかも。」

2017年は、日本映画『カメラを止めるな!』が世界的なブームとなり、ゾンビを題材にしたホラー映画が話題となった時期でした。したがって、これらの話題はこの映画のレビューにも反映されている。抽出された特徴語は、各時期の話題の変化が反映されている。

表2は、波別によるレビューの特徴語の抽出結果を示したものである。表2から以下のことがわかる。

「COVID-19」が日本で流行した第1波では、映画『コンテイジョン』のレビューに特徴的な評価情報として、「外出」「自粛」といった当時のニュースに関する話題のほか、「Netflix」「話題」「酷似」などの単語から、この映画を見直すきっかけも示されている。第4波以降のレビューでは、「接種」という単語が特徴的な評価視点として上位にランキングされている。第4波の2021年3月から5月までは、日本でのワクチン接種が開始した時期である。第5波の抽出結果では、「陰謀論」というワードが出現し、当時インターネット上ではワクチンに関する陰謀論が話題となった。

以上のことから、各時期の抽出結果には、その時期の話

題となった出来事や人々が関心のあるものが反映されていることがわかる。本研究で提案する特徴語抽出手法は、時系列によって、映画レビューの各時期の特徴的な評価視点の抽出に有効といえる。

表1『コンテイジョン』のレビューの年ごとに抽出された結果

No.	2013	2014	2015	2016	2017	2018	2019	2020	2021
1	緊張	現代	派手	派手	ゾンビ	ホラー	期待	経済	禍
2	新種	系	期待	演技	マットデ イモン	期待	万	外出	接種
3	地味	期待	新種	ジャーナ リスト	静か	好き	マットデ イモン	タイムリー	論
4	豪華	間	無駄	免疫	地味	退屈	系	Netflix	頃
5	退屈	緊張	ドラマ	それぞれ	期待	系	静か	外	まんま
6	キャスト	群像	緊張	病	群像	風	死者	自粛	前
7	盛り上がり	派手	有名	音楽	役者	主役	退屈	タイミング	陰謀
8	出演	マリオン・コ ティヤール	人類	物語	新種	ゾンビ	キャスト	現状	予測
9	自体	音楽	音楽	役者	全体	新種	ドラマ	時期	予見
10	残念	免疫	風	妻	級	様々	結局	話題	年
11	立場	ウィンスレ ット	盛り上がり	アウトブ レイク	系	主人公	豪華	現場	驚き
12	ウィンスレ ット	静か	病気	俳優	監督	マット・ デイモン	パンデミ ック	方々	以前
13	全体	有名	ゾンビ	新種	ソダーバ ーグ	謎	好き	ご時世	科学
14	映像	好き	ベス	主役	派手	妻	側	今後	マスク
15	派手	監督	地味	人類	有名	役	ゾンビ	最前線	予言

表2『コンテイジョン』のレビューの波ごとに抽出された結果 (全体と比較)

	第一波	第二波	第三波	第四波	第五波
1	外出	頃	禍	禍	禍
2	自粛	現代	COVID	緊急	接種
3	現状	脅威	論	接種	論
4	Netflix	禍	騒動	予知	陰謀
5	外	月	はず	まんま	科学
6	今後	マスク	嫌	頃	万
7	危機	仕方	SARS	前	予測
8	情勢	オチ	予測	実感	順番
9	家	以前	時代	SARS	理解
10	いま	感覚	予見	調査	マスク
11	状況	実感	ネタ	身近	前
12	方々	期間	本当	以前	年
13	話題	WHO	当たり前	娘	頃
14	経済	当たり前	陰謀	年	驚き
15	酷似	とき	人事	スピード	感覚

表3『コンテイジョン』のレビューの波ごとに抽出された結果 (前と比較)

No.	第一波	第二波	第三波	第四波	第五波
1	コロナ	禍	はず	普通	理解
2	今	脅威	生活	まんま	コンテイジョン
3	状況	謎	論	病気	陰謀
4	医療	マスク	騒動	好き	登場
5	たち	オチ	風	リアリティ	対策
6	デマ	月	人事	自粛	コウモリ
7	日本	仕方	たくさん	スター	ドキュメンタリー
8	前	まんま	予見	ストーリー	マスク
9	開発	アウトブレイク	表現	危機	咳
10	年	全体	崩壊	調査	中国
11	こと	人類	マットデ イモン	評価	論
12	ワクチン	系	嫌	宣言	影響
13	新型	当たり前	一番	予知	ラストシーン
14	みんな	病気	陰謀	以上	部分
15	本	豚	ウィンスレ ット	頭	WHO

しかし、部分と全体を比較するとき、「コロナ禍」の「禍」が全体と比べと特徴のある「名詞」であるため、第三、四、五波では、「禍」が全て特徴的な評価視点の一位となっている。そこで、ある期間から突然増えた言葉を捉えるために、全期間と比較するのではなく、直前の時期と比較して、各時期の特徴語を抽出する実験を行った(表3)。2019年までのデータを第一波の比較対象として計算する。

表3から以下のことがわかる。

第一波の間、映画レビューには「コロナ」の他、「医療」「デマ」「日本」などの言葉が前より増えた。また、「コロナ禍」の「禍」が第二波から増えた。「陰謀論」という言葉は、ワクチン最初に接種し始めた第三波から出現し、第五波の頃がレビューに増え、当時インターネット上ではワクチンに関する陰謀論が話題となったことが原因だと考えられる。

以上のことから、各時期の抽出結果には、その時期の話題となった出来事や人々に関心のあるものが反映されていることがわかる。本研究で提案する特徴語抽出手法は、時系列によって、映画レビューの各時期の特徴的な評価視点の抽出に有効といえる。

5. 可視化

本章では、抽出された特徴的な評価視点の可視化について検討する。

図1は、一つの可視化の例である。背景の紫のグラフは、朝日新聞社提供している「新型コロナウイルスの日本全国におけるひごとの感染者数」のグラフである。その前面は、表3の抽出された一部の特徴的な評価視点の可視化ワードクラウドである。感染者数の時系列による変化を各時期の評価視点ワードクラウドと照らし合わせる。図1から分かるように、第一波の時、「コロナ」「医療」「デマ」などのキーワードが急に増えたと考える。それが現実社会の状況と繋がっているとわかる。例えば、「デマ」を言及したレビューは以下のようにになっている。

「デマに惑わされトレペや納豆まで買う。」

「またその略奪行為・デマが横行して荒稼ぎをする輩が出てくる...と今現実に行っていることのダイジェスト版を見せられるようです。」

当時日本国内でトレペの買い占めが起きていたため、このような話題が盛り上がっていた。

また、第五波は前より、ひごとの感染者数が増え、自分自身の経験と比べ、映画に対する「理解」を話していることが多くなった。

「新型コロナウイルス以前と以降で映画で描かれている状況への理解度が変わるので、この映画を見た感想も変わる

かと思います。」

図2では、タイムラインと評価視点ワードクラウドで可視化したシステムである。下には、同時期の関連記事です。

可視化システムの上側はタイムラインと評価視点ワードクラウドで、下には、同時期の関連記事です。まず月別で見えるように、2020年1月のワードクラウドに「中国」という言葉があります。そのニュース記事のキーワードも中国です。

ワードの相対出現頻度比 RFR に基づいて特徴評価視点の抽出し、ワードクラウドで可視化、時系列による映画レビューの変化を捉えることができる。映画作品のレビューは社会状況と時間の変化につながっていることが分かる。

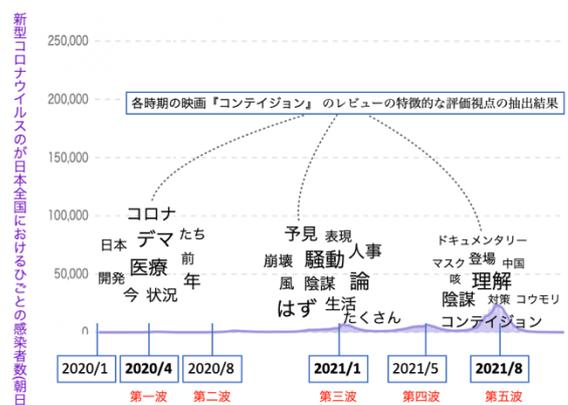


図1 感染者数と抽出結果の可視化



図2 ニュース記と抽出結果の可視化

6. おわりに

本稿では、映画レビューの時間による変化を捉えるために、名詞の相対出現頻度の比較する特徴度尺度 RFR を導入して、時系列によるワードクラウドの可視化を行った。映画レビューの投稿時間に従って、レビューをワードクラウドで特徴的なワードを可視化することで、各時期の特徴を捉えた。

実験では、同一作品のレビュー群を対象として、他の時

期と比べて、特定の時期に投稿されたレビューの特徴的な評価視点を抽出する。その結果、各時期における抽出結果が、社会におけるできごとに応じて変化していることが確認できた。これらのことから、提案する手法により、映画レビューの特徴と変化をよりよく捉えられることが確かめられた。

従って、本研究の目的は達成されたと考える。

本研究で提案した手法により、アーカイブされていた個々の映画作品のレビューにおける特徴的な評価視点と、社会状況に応じた評価視点の時系列変化が読み取れるようになった。本手法を用いて、アーカイブされていたレビューを分析することで、その映画作品に対する評価と、変化

していく社会状況との関わりを明らかにすることができる。このことにより、今後の映画作品制作や批評のあり方を改善していける可能性がある。

また、本研究の手法は、映画以外のレビューにおける特徴的な評価視点の抽出・時系列変化の読み取りにも適用することができる。本手法を発展させることで、例えば読書レビュー・SNSなど、Webにアーカイブされているコミュニケーションなどの自由記述文や、時間情報を含む大規模テキストデータの分析にも寄与するものと考えられる。

参考文献

1 Minqing Hu; Bing Liu: Mining Opinion Features in Customer Reviews, In Proceedings of the 19th National Conference on Artificial Intelligence, 2004, pp. 755-760.

2 小林のぞみ, 乾健太郎, 松本裕治, 立石健二, 福島俊一: 意見抽出のための評価表現の収集, 自然言語処理, 2005, Vol. 12, No. 2, pp. 203 -222.

3 内山将夫, 中條清美, 山本英子, 井佐原均: 英語教育のための分野特徴単語の選定尺度の比較, 自然言語処理, 2004, Vol. 11, No. 3, pp. 165-197.

4 乾孝司, 板谷悠人, 山本幹雄, 新里圭司, 平手勇宇, 山田薫: 意見集約における相対的特徴を考慮した評価視点の構造化, 自然言語処理, 2013, vol. 20, no. 1, pp. 3-25.

5 日本の映画情報 サービス. <https://filmmarks.com>