

# 大規模ティックデータを用いたニュース記事分類と文章トピック時系列変化の可視化 News Article Classification and Visualization of Text Topic Time-Series Changes Using Large-Scale Tick Data

西 良浩<sup>†</sup>      高橋 大志<sup>†</sup>  
Yoshihiro Nishi   Hiroshi Takahashi

## 1. はじめに

金融市場において配信されるニュース記事は資産価格評価に対し影響を与える重要な情報である。企業のファンダメンタルズや投資家のセンチメントなどが、資産価格に影響を与えるとの報告が行われており [4, 7], ニュース記事は時流に応じてそれらに関する情報を金融市場に配信し、株価に対して影響をもたらしていると考えられる。

自然言語データを機械学習モデルや深層学習モデルを用いて分析する際には、一般的にはデータ数が多いほど分析の精度は向上する傾向にある。しかしながら、金融市場において発信されたニュース記事には限りがあり、特に特定の期間内において指定のメディア内で発信された対象企業の分析を行う場合のように範囲が絞られる場合において獲得できるニュース記事の数には制約がある。近年では、自然言語生成技術により、分析に用いるデータを増やし、精度向上を試みる取り組みがなされているが、ニュース記事が主な話題としているトピックがなにかを時系列で考慮し自然言語生成を行い、ニュース記事のようなデータを拡充する取り組みはまだ十分に行われていない。

本研究では、ニュース記事が主として扱ったトピックの時系列変化を検証し、さらに変動があった単語分布の比較を行う。本論文の構成は次の通りである。まず第 2 章で関連研究について述べる。第 3 章では分析対象とするデータについて述べる。第 4 章では分析手法について述べ、第 5 章では分析結果について述べる。最後に第 6 章で結論を述べる。

## 2. 関連研究

ニュースと株価変動の関連性を分析した研究は数多くある。例えば、ニュースの文章を用いてテキストマイニングを行い株価変動の分析した研究では、ニュースに含まれるファンダメンタルズやセンチメントの情報が株価に反映される可能性があることが報告されている [2, 3]。

近年では自然言語処理技術の発展に伴い、分析に使用するデータを自然言語処理モデルにより生成し、分析データを増やす取り組みも行われている [1, 8, 9, 10]。ニュース記事が主なトピックとして扱う話題は時系列で変化することが考えられるため、そのトピックの移り変わりを考慮した自然言語生成を行うことでより精度の高い生成ができる可能性がある。しかしながら、ニュース記事が主な話題としているトピックがなにかを時系列で考慮し自然言語生成を

行い、ニュース記事の株価変動への影響を評価する取り組みはまだ十分に行われていない。

## 3. データ

本研究では、分析の対象期間を 2014 年 4 月から 2017 年 3 月までとし、マーケットデータとニュースデータを用いて分析を行った。2014 年時点で時価総額が最も高い上位 3 社(トヨタ自動車株式会社, 日産自動車株式会社, 本田技研工業株式会社)を主要な自動車企業とし、分析対象とした。

### 3.1 マーケットデータ

マーケットデータとして、国内上場企業のうち、自動車企業のなかから主要 3 社 (トヨタ自動車株式会社, 日産自動車株式会社, 本田技研工業株式会社) の株式価格を対象として分析を行った。分析対象期間は、2014 年 4 月から 2017 年 3 月とした。2014 年から 2017 年までの 4 年間分の市場取引データ約 19 億件を取得し、対象企業・対象期間に絞り、分析に使用した。本分析に用いたマーケットデータには、取引成立価格や取引量などの株式取引に関する情報が含まれており、各行にマイクロ秒単位のタイムスタンプが付されている。また、本分析では、企業属性に関するデータも日経 NEEDS, Datastream から取得し分析に用いた。

### 3.2 ニュースデータ

ニュースデータとして、トムソン・ロイター社が配信を行ったニュース記事を用いた。日本企業に関するニュースは主として英語もしくは日本語により配信されている。配信されたニュースのテキスト情報には、ヘッドラインと本文があり、ヘッドラインは本文内の重要な内容を要約したテキストデータである。ニュースには配信された日時のタイムスタンプが付されている。

本研究ではニュース配信の前後 1 分間に取引があった英語のヘッドラインを用いて分析を行う。2014 年 4 月から 2017 年 3 月までのトヨタ自動車株式会社, 日産自動車株式会社, 本田技研工業株式会社に関するニュース 2,326 件を取得した。

## 4. 分析手法

本研究では対象のニュース記事データに対し、トピック時系列変化の分析を行い、その後、文章トピックに変化が生じている年を対象とし年・ラベルごとにニュース記事の単語分布を比較する。

<sup>†</sup> 慶應義塾大学大学院 経営管理研究科  
Graduate School of Business Administration, Keio University

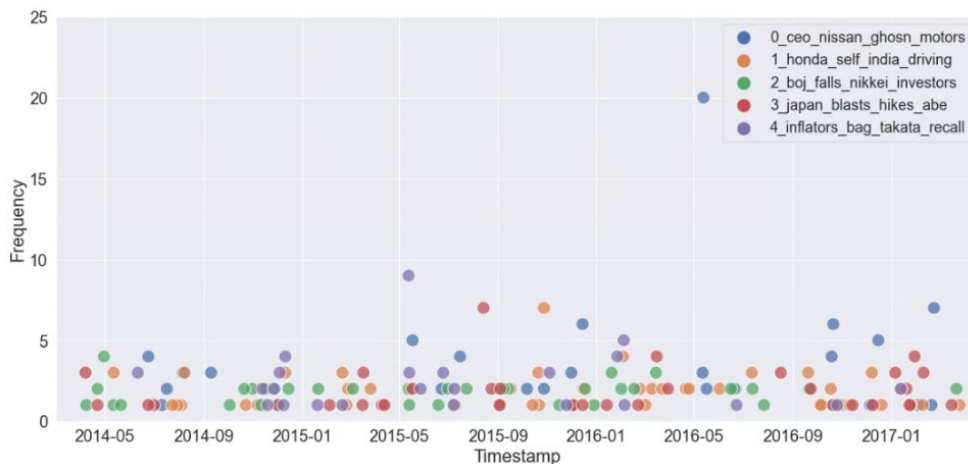


図 1 2014 年 4 月から 2017 年 3 月におけるニュース記事のトピックトレンド結果

#### 4.1 ニュース記事のトピックトレンド分析

トピックモデルとは、教師無しでテキスト中に扱われている共通のテーマを明らかにするために有効な手法である。本研究ではトピックモデルによるニュース記事トレンドの分析に BERTopic を用いる[5]。従来のトピックモデリング技術はクラスタの重心に近い単語がそのクラスタやトピックを代表とすると仮定しているが、実際にはクラスタが常にクラスタ重心を中心とした領域の中にあるとは限らないため、トピックモデルにより算出したトピック表現の結果が誤解を招く可能性がある。BERTopic はクラスタリング技術と TF-IDF のクラススペースでのバリエーションを活用して異なるトピック表現を生成するモデルであり、より柔軟なトピックモデリングを可能としている。

本研究では 2014 年 4 月から 2017 年 3 月のニュース記事に対して BERTopic を用い 5 つのトピックを抽出し、ニュース記事トピックトレンドの変化が生じているかを検証する。

#### 4.2 株価変動率に基づくラベル付け

次に、(1)の定義式により 2014 年 4 月から 2017 年までのニュース記事にラベル付けを行う。該当のニュース記事が配信された前後 1 分間の平均株価を取得し、ニュース記事ごとに株価変動率を算出する。本研究では大規模な株式価格ティックデータを用いて株価変動を算出し、株価変動状況の実態をもとにして各ニュース記事のラベルを算出する。

$$\text{Stock price fluctuation rate } (\alpha) = \frac{(\text{Average stock price 1 minute after news}) - (\text{Average stock price 1 minute before news})}{(\text{Average price 1 minute before news})} \times 100$$

Positive:  $\alpha > 0\%$   
Negative:  $\alpha < 0\%$

#### 4.3 年・ラベルごと単語分布の比較

各年のラベル別単語分布の可視化および比較には Scattertext を用いる[6]。Scattertext とはラベル間の単語分布の差異を可視化するライブラリである。テキスト中に含ま

れる単語を発生頻度に基づき 2 次元上に散布図で出力し、傾向に違いがあるかを捉えることができる。

本研究では株価変動に基づきラベル付けしたニュース記事を用い、Positive な影響を与えたニュース記事と Negative な影響を与えたニュース記事の単語分布を可視化し、年ごとに変化が生じているかを検証する。

### 5. 分析結果

本研究では対象データに対し、BERTopic を介したニュース記事のトピックトレンド分析、株価変動率によるラベル付け、Scattertext による年・ラベルごと単語分布の比較を行った。

#### 5.1 ニュース記事トピックトレンドの分析結果

BERTopic を介したニュース記事トピックトレンド分析の結果、抽出された 5 つのトピックは 0\_ceo\_nissan\_ghosn\_motors, 1\_honda\_self\_india\_driving, 2\_boj\_falls\_nikkei\_investors, 3\_japan\_blasts\_hikes\_abe, 4\_inflators\_bag\_takata\_recall であった。縦軸は期間内における該当トピックの発生回数、横軸は時間を示している。特にトピック 0 は日産自動車株式会社のカルロス・ゴーン氏に関するトピックであり、2016 年 5 月ごろに話題としてピークを迎えた。また、トピック 4 はタカタ株式会社のエアバッグ・リコール問題に関するトピックであり、2015 年の 5 月ごろに話題としてピークを迎えた。他いづれのトピックにおいても各トピックの出現頻度は時系列によって変化している。

#### 5.2 株価変動率に基づくラベル付け結果

株価変動率に基づく 2014 年 4 月から 2017 年 3 月までのニュース記事 2,326 件へのラベル付けの結果、2014 年 4 月から 2017 年 3 月にロイターニュースが配信したポジティブ・ネガティブなニュース記事は大きな偏りはなく、2014

