

# 日本中世古記録を対象としたトピック抽出自動化システムの構築

鳥居克哉<sup>1</sup> 中村覚<sup>2</sup> 山田太造<sup>2</sup> 稗方和夫<sup>1</sup>

**概要:** 本研究では、日本史学者の史料研究支援のために、史料群に対する可用性と有用性を高めるトピック抽出を自動で行うシステムの開発を行った。ルールベースにより抽出した人名及び N-gram や Sentencepiece によって分割した用語から Bag-of-Word を生成し、LDA(Latent Dirichlet Allocation)を適用することでトピック分析を行った。さらに、史料と人物索引表を入力としてこの一連の分析を行う Web システムをクラウド上に構築した。また、鎌倉時代の公卿である藤原(勘解由小路)経光が記した『民経記』を対象にこのシステムを利用し、トピック分析の結果が史実に整合していることが確認でき、有効性が示された。

**キーワード:** トピックモデル, LDA, 日本史, 古記録, オンライン分析

## Development of an Automated Topic Extraction System for Ancient Japanese Medieval Records

KATSUYA TORII<sup>†1</sup> SATORU NAKAMURA<sup>†2</sup>  
TAIZO YAMADA<sup>†2</sup> KAZUO HIEKATA<sup>†1</sup>

**Abstract:** In this study, we developed a system that automatically extracts topics to increase the availability and usefulness of historical documents to support Japanese historians in their research on historical documents. We generated a Bag-of-Words from the names of people extracted by the rule base and the terms divided by N-gram and Sentencepiece, and applied LDA (Latent Dirichlet Allocation) to analyze the topics. In addition, we constructed a web system on the cloud to perform this series of analysis using historical documents and a person index table as input. In addition, we used this system to analyze the "Minkeiki" written by Fujiwara (Kadenokoji) Tsunemitsu, a kuge of the Kamakura period, and confirmed that the results of the topic analysis were consistent with the historical facts, demonstrating its effectiveness.

**Keywords:** Topic Model, LDA, Japanese History, Old Diary, OLAP

### 1. はじめに

日本史学者は、歴史資料(以下、史料)を収集・読解・分析することで研究課題を解明していく[1]。特に分析の過程では史料を分析する際に大量の史料を読み込み、その中から研究課題に当てはまるトピックを見つける必要がある。本研究では、データ駆動型人文科学研究の一環として、自然言語処理技術により日本史学者の資料研究の分析過程を支援することを目的とする。本研究では史料から抽出した人物・語句を対象としてトピック抽出を行ない、その結果をトピックに属する語句と人物のリストや時系列関係で可視化した。また、ユーザーがアクセスし一連の分析を行うことができる Web Application として開発を行なった。ケーススタディでは、鎌倉時代の公卿である藤原(勘解由小路)経光が記した『民経記』の史料を対象に、システムを用い日本史学者へのインタビューを通じてトピック抽出の評価・検討を行った。本稿は[2]をベースとして、語句抽出手法を再検討し、システムを改良した。さらに、日本史学者へのインタビューを通じてシステム評価を行った。

### 2. 関連研究と本研究の位置付け

#### 2.1 トピックモデルのデータ分析への適用事例

本研究で扱うトピックモデルは大規模なデータを分析するための手法として広い分野で利用されているアプローチである。テキストデータだけではなく、顧客購買データ解析、画像処理、人の行動解析など様々な分野で用いられている。Sukhija ら[3]は、トピックモデルのうち、LDA(Latent Dirichlet Allocation)[4]を社会科学データに適用し、消費活動と人種・階級・ジェンダーとの関係を抽出、可視化し社会科学的な分析を行った。Heng ら[5]はアマゾンのレビューに LDA を適用し、消費者の食品購入に影響を与える要因を分析した。トピックモデルを実行する上で必要な単語の多重表現である Bag-of-Words に、どのデータを選択するかによって、多種多様な分析が可能である。

トピックモデルを歴史研究に応用した先行事例として山田ら[6]の研究が挙げられる。山田らは史料内の人物を対象に LDA を特定の史料に適用することで、史料中の人物の共起関係を基に潜在する意味関係を検出し、人物間の関係性を検出した。また、それを発展させ、時空間変化と人物の関係を抽出した。

<sup>†1</sup> 東京大学大学院新領域創成科学研究科  
Graduate School of Frontier Sciences, The University of Tokyo

<sup>†2</sup> 東京大学史料編纂所

Historiographical Institute, The University of Tokyo

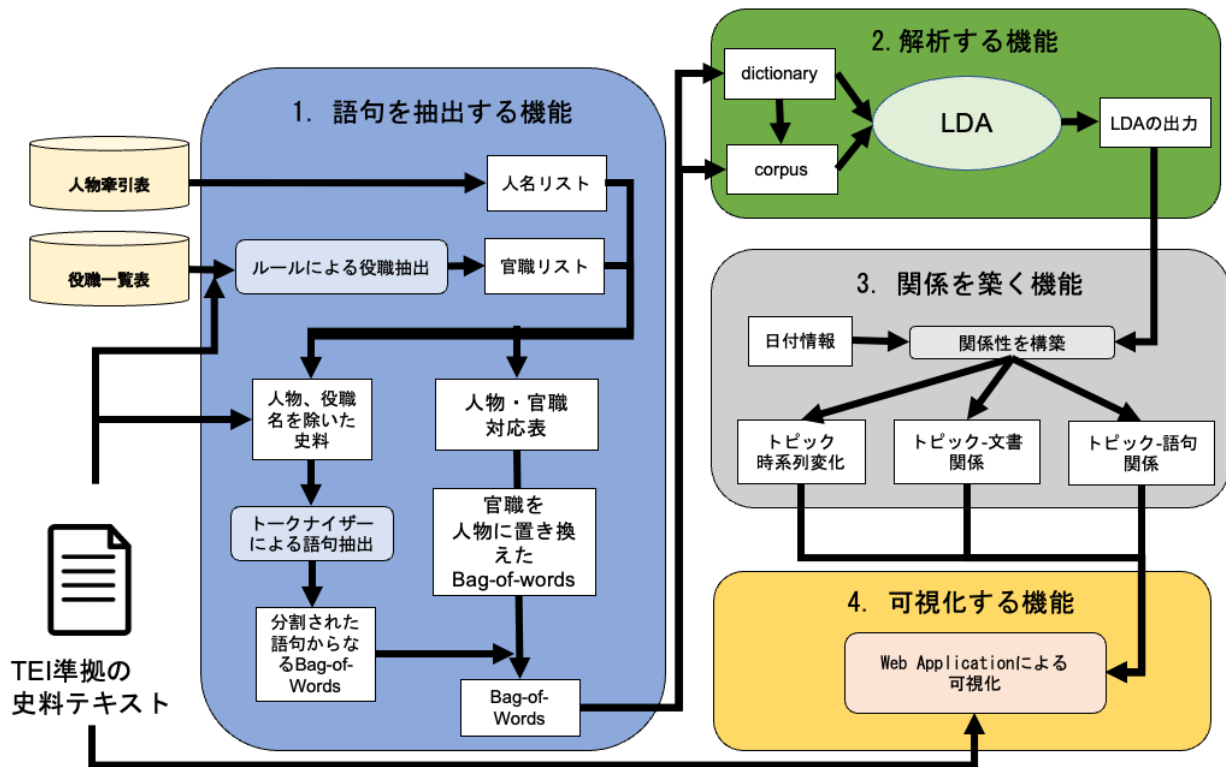


図1 システムの概要図

Figure1 Overview of the proposed System

## 2.2 本研究の位置付け

先行研究では Bag-of-Words に人物だけではなく、史料に出現する語句もトピックモデルに組み込むことで、各文書のトピックを推定する。

また、既存の支援システムとして、Voyant Tools や KHcoder などがあげられるが、これらの既存のシステムとの差別化を図りながら、UI の開発を行う。

また、本研究では平安中期～鎌倉中期の変体漢文の史料を対象にユーザーのユーズに即時対応可能な OLAP 分析システムの開発を目指す。

## 3. 開発したシステム

開発したシステムの概念図を図 1 に示す。ユーザーは TEI に準拠した史料と人物牽引表・役職一覧表をシステムの入力とすることで可視化された結果を取得することができる。なお、TEI に準拠した史料とは、構造化ルールを定める TEI ガイドラインに基づいてタグ付け及び構造化が行われた史料のことで、人物牽引表とは史料を編纂する際に作成される史料中に出現する人物名と出現場所を記したデータである。

以下、開発したシステムを 4 つの機能に分けて説明する。

### 3.1 語句を抽出する機能

トピック抽出において史料の解析を実行する上で必要な、単語の多重表現である Bag-of-Words を作成する。まず人物牽引表から人物リストを作成し、役職一覧表を元に、史料中の役職パターンに一致する表現を抽出し、史料中の官職を人物名に置き換えた人物のみからなる Bag-of-Words を作成する。なお、本ケーススタディでは <役職(家名+諱)> というパターンを採用した。一方、史料から役職名を除いた史料からトークナイザーによる単語抽出を行うことで抽出した語句のみからなる Bag-of-Words を作成する。これらを組み合わせ、人物・史料中の語句が含まれる Bag-of-Words を作成できる。

本稿では語句を抽出する際、トークナイザーとして教師なしの手法である N-gram と Sentencepiece[7] を比較検討した。

N-gram では、N=2, 3 の分割で抽出された、出現回数 50 以上の出現頻度の高い語句のうち、Web 上で公開されている百科事典(コトバンク)に存在する語句を選定した。

Sentencepiece は、教師なしトークナイザーで主にニューラルモデルの学習前で語彙の大きさがあらかじめ定められている場合に用いられる手法である。テキストをサブワードに分割するため、語彙数を小さくしつつ未知語をなくすることができるという特徴がある。また、文章から直接学習

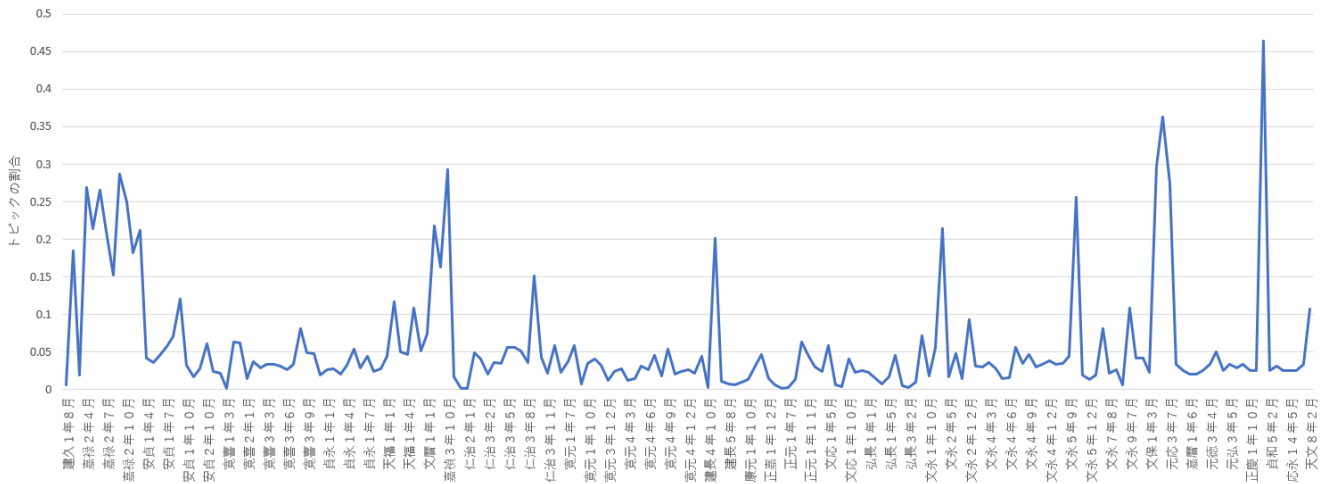


図4 トピック3の月別の時系列変化

Figure4 Monthly Time Series Change Graph for Topic3

することができ、言語固有の前処理・後処理に依存しない、エンドツーエンドのシステムを構築することができる。

『民経記』の嘉禄2年4月分のテキストデータについて人手で調査した正解件数227件に対する結果を下記に示す。なお、N-gramにおいては『民経記』に対応するよう後処理を行なったが、Sentencepieceはどの史料にも対応できるように使用した分、精度が大きく落ちている。

Sentencepieceの精度が上がるように今後改善する計画ではあるが、本稿のケーススタディではN-gramによる手法を用いた。

	適合率	再現率	F値
N-gram	0.92040	0.81858	0.86651
Setencepiece	0.43096	0.40517	0.41767

図2 語句抽出の結果

Figure2 Results of Word extraction

### 3.2 解析する機能

トピックモデルは対象とするデータから潜在する話題(トピック)を検出することができる教師なし学習の一つである。本研究では、トピックモデルとして、LDAを選定し、史料に出現する人物と語句を史料の特徴語として捉え、LDAの変数として組み込むことで、各史料のトピックを推定する。本研究におけるLDAによる史料の生成確率は次式のとおりである。

$$p(d|\alpha, \beta) = \int Dir(\theta|\alpha) \left( \prod_{n=1}^{|d|} \sum_{k=1}^C p(w_n|z_k, \beta) p(z_k|\theta) \right) d\theta \quad (1)$$

$\alpha$ と $\beta$ はパラメータ、 $z = z_1, z_2, \dots, z_c$ は潜在トピック、 $\theta = \theta_1, \theta_2, \dots, \theta_c$ は潜在トピックの生成確率、 $Dir(\theta|\alpha)$

はディレクレ分布、 $d = (w_1, w_2, \dots, w_{|d|})$ は史料、 $w_n$ は特徴語、 $|d|$ は史料dの単語数を表す。

LDAの計算にはオープンソースライブラリであるGensimを用いた。Gensimでは、トピックモデルの推定方法のうち、従来の変分ベイズを応用したオンライン変分ベイズ[7]を使用することで、高速な学習を実現している。

入力には語彙とID、出現回数に対応表であるdictionaryとdictionaryに登録されているIDからなるBag-of-Wordsであるcorpusが必要であり、3.1で作成したBag-of-Wordsからdictionaryとcorpusを作成した。また、passesについては、本研究ではdefault値としてperplexityが凡そ収束する50を設定した。本ケーススタディにおいてPerplexityが収束する様子を図3に図示する。

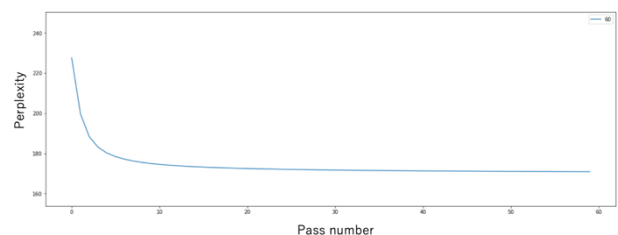


図3 Perplexityの収束

Figure3 Perplexity changes with the number of training passes

### 3.3 関係性を築く機能

LDAを実行後の出力を整備し、UIに表示すべき内容を作成する。ここではトピックに対して関連度が高い単語を再び語句と人物に区別した。また、日記である対象史料から日付情報を抽出し、トピックの時間変化抽出を行った。図4はトピック数10の場合のTopic3の月別の時系列変化グラフである。

## Topic3

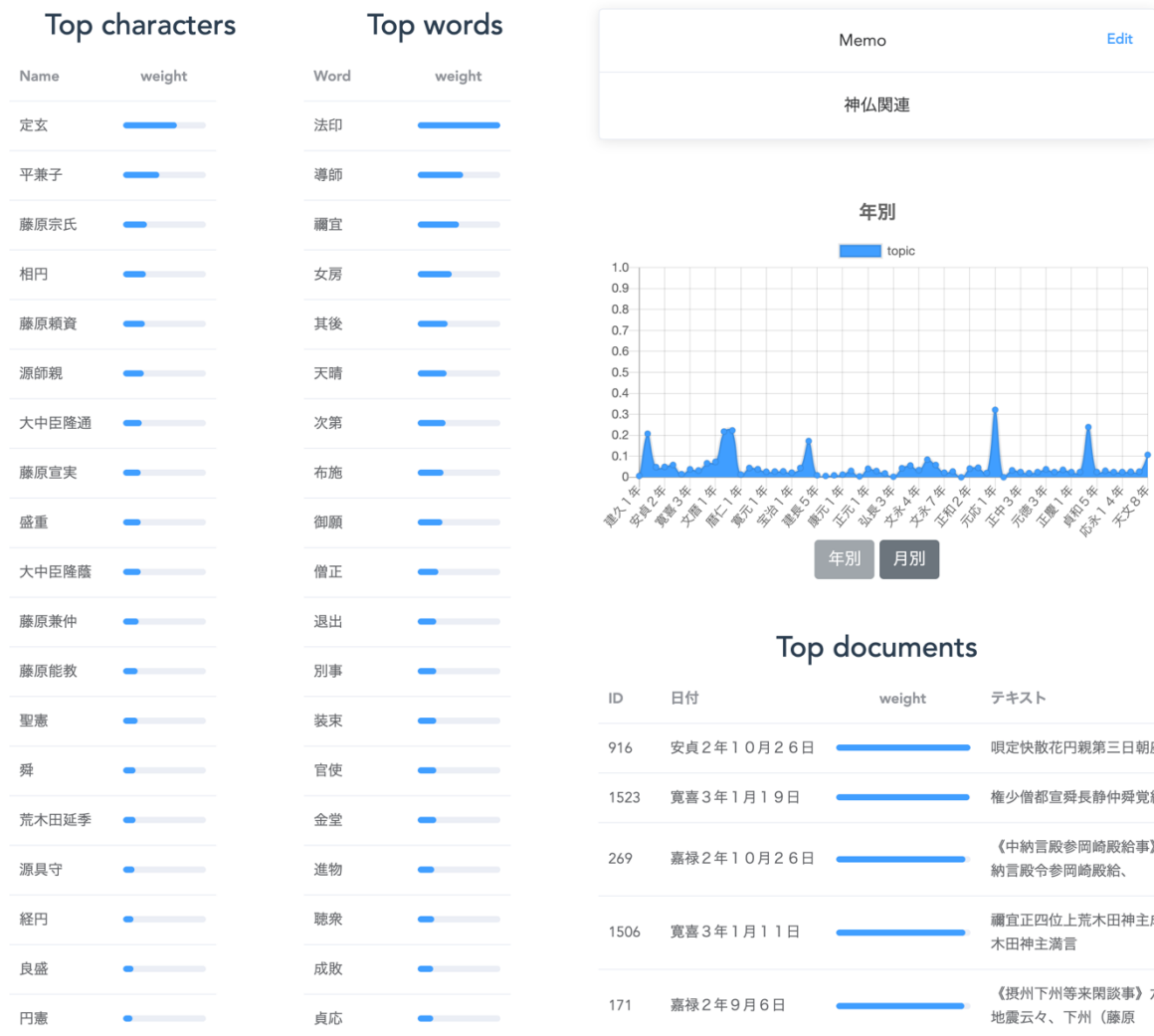


図5 UIの例. 語句と人物のリストとトピックの時系列変化

Figure5 Monthly Time Series Change Graph for Topic3

### 3.4 Web で可視化する機能

トピックごとの語句、語句ごとのトピックの割合、史料ごとのトピックの割合を表示する3つの画面を作成する。データをリストやグラフの形で表示し、史料中の語句をトピックで色付けし、リンク機能を実装することでユーザビリティの向上を目指した。また、Web Application上で入力フォームを用意し、ユーザーが入力後、同じサイト内で結果を確認することを可能とした。図5はWeb Applicationのトピック数10におけるTopic3におけるUIである。画面の左側にTopic3における重要な人物と語句のリスト、右上にトピックの時系列変化、右下に関連度が高い日付の史料を可視化している。

## 4. ケーススタディ

### 4.1 『民経記』

本研究のケーススタディとして、『民経記』を対象とする史料とする。『民経記』は鎌倉中期の公家である、藤原(広橋)経光(1212~1274)の日記であり、別に『経光卿記』、『経光卿曆記』、『経光御記』、『中光記』、『民経御記』等の称呼がある。自筆原本42巻は、国立歴史民俗資料館蔵史料に蔵する。『大日本古記録民経記』(全十一冊)として編纂史料集が発刊されている。ケーススタディでは、編纂所がサービスしている古記録フルテキストデータベースに格納されている民経記データを使用した。文字数は1,226,975文字、日条数は6630であった。



Topic1	Topic2	Topic3	Topic4	Topic5	Topic6	Topic7	Topic8	Topic9	Topic10
人名 語句	人名 語句	人名 語句	人名 語句	人名 語句	人名 語句	人名 語句	人名 語句	人名 語句	人名 語句
藤原経光 評定	藤原経光 幸清	藤原経光 神吉	藤原経光 定玄	藤原経光 法印	藤原経光 源通方	藤原経光 奉行	藤原経光 殿上	藤原経光 藤原氏	藤原経光 御方
藤原経光 参院	藤原経光 藤原兼頼	藤原経光 沐浴	藤原経光 平兼子	藤原経光 導師	藤原経光 藤原教実	藤原経光 延候	藤原経光 藤原隆親	藤原経光 此間	藤原経光 藤原隆親
平国朝 文書	平国朝 皇子内親王	平国朝 天恩	平国朝 藤原宗氏	平国朝 職宣	平国朝 参内	平国朝 平知宗	平国朝 退出	平国朝 藤原公経	平国朝 御奉
藤原頼資 候へ	藤原頼資 藤原公子	藤原頼資 天晴	藤原頼資 相丹	藤原頼資 女房	藤原頼資 藤原道家	藤原頼資 相触	藤原頼資 藤原忠高	藤原頼資 参進	藤原頼資 御所
藤原経長 御所	藤原経長 藤原頼資	藤原経長 月徳	藤原経長 藤原頼資	藤原経長 其後	藤原経長 藤原信盛	藤原経長 仰下	藤原経長 起座	藤原経長 藤原実有	藤原経長 雑色
藤原光国 正月	藤原光国 盛家	藤原光国 拜官	藤原光国 藤原親	藤原光国 天晴	藤原光国 藤原忠高	藤原光国 内覧	藤原光国 藤原良平	藤原光国 次第	藤原光国 源房
藤原宗雅 奉行	藤原宗雅 藤原経光	藤原宗雅 成徳	藤原宗雅 大中臣隆通	藤原宗雅 次第	藤原宗雅 平有親	藤原宗雅 行事	藤原宗雅 安倍貞季	藤原宗雅 出御	藤原宗雅 興茂在盛
藤原定朝 沙汰	藤原定朝 藤原信子	藤原定朝 母倉	藤原定朝 藤原宣実	藤原定朝 布施	藤原定朝 藤原兼高	藤原定朝 奏聞	藤原定朝 源頼定	藤原定朝 其後	藤原定朝 藤原実家
藤原公基 日至	藤原公基 藤原資定	藤原公基 大小	藤原公基 盛重	藤原公基 節願	藤原公基 家光	藤原公基 外記	藤原公基 源通忠	藤原公基 瓶子	藤原公基 中原俊職
藤原宗経 上皇	藤原宗経 大中臣公行	藤原宗経 土公	藤原宗経 大中臣隆隆	藤原宗経 僧正	藤原宗経 藤原家朝	藤原宗経 沙汰	藤原宗経 中原雅親	藤原宗経 信光	藤原宗経 遠御
藤原高朝 供養	藤原高朝 藤原兼仲	藤原高朝 奥書	藤原高朝 藤原兼仲	藤原高朝 退出	藤原高朝 藤原資朝	藤原高朝 宣命	藤原高朝 藤原経光	藤原高朝 吉書	藤原高朝 平繁俊
嵯峨天皇 下官	嵯峨天皇 藤原道家	嵯峨天皇 加冠	嵯峨天皇 藤原能教	嵯峨天皇 別事	嵯峨天皇 藤原親朝	嵯峨天皇 此間	嵯峨天皇 藤原信盛	嵯峨天皇 御簾	嵯峨天皇 平輔兼
紀実光 参候	紀実光 藤原氏	紀実光 以後	紀実光 聖書	紀実光 装束	紀実光 藤原実氏	紀実光 天晴	紀実光 大中臣永隆	紀実光 束帯	紀実光 藤原泰通
安楽 延引	安楽 藤原通雅	安楽 雨下	安楽 舜	安楽 官使	安楽 平経高	安楽 退出	安楽 藤原宗平	安楽 障子	安楽 永親
仁明天皇 御参	仁明天皇 藤原実基	仁明天皇 日出	仁明天皇 荒木田延季	仁明天皇 金堂	仁明天皇 平範輪	仁明天皇 下知	仁明天皇 藤原実世	仁明天皇 中門	仁明天皇 利子内親王
藤原兼仲 寛喜	藤原兼仲 藤原家実	藤原兼仲 日入	藤原兼仲 経門	藤原兼仲 経宗	藤原兼仲 参入	藤原兼仲 藤原実任	藤原兼仲 列立	藤原兼仲 藤原隆行	藤原兼仲 舍人
藤原宗親 自筆	藤原宗親 門助法親王	藤原宗親 不問	藤原宗親 良盛	藤原宗親 成敗	藤原宗親 中原師兼	藤原宗親 相催	藤原宗親 藤原伊平	藤原宗親 参入	藤原宗親 藤原家朝
宗尊親王 目録	宗尊親王 藤子内親王	宗尊親王 朝問	宗尊親王 門書	宗尊親王 貞心	宗尊親王 平範輪	宗尊親王 日時	宗尊親王 藤原公光	宗尊親王 気色	宗尊親王 藤原家季
藤原朝長 元年	藤原朝長 清原元久	藤原朝長 不祝病	藤原朝長 藤原任子	藤原朝長 停止	藤原朝長 源頼朝	藤原朝長 源有教	藤原朝長 西面	藤原朝長 藤原為家	藤原朝長 女院
高倉 法勝寺	高倉 薬師	高倉 禪林寺	高倉 公兼	高倉 宿所	高倉 中原親俊	高倉 数刻	高倉 橘親氏	高倉 御座	高倉 藤原長子
藤原宗雅 出御	藤原宗雅 鳥羽重久	藤原宗雅 還御	藤原宗雅 釈迦	藤原宗雅 釈久	藤原宗雅 藤原宗平	藤原宗雅 御覽	藤原宗雅 源頼平	藤原宗雅 下賜	藤原宗雅 季範
藤原隆親 自今	藤原隆親 嵯峨天皇	藤原隆親 下食	藤原隆親 尊海	藤原隆親 閑談	藤原隆親 安倍国通	藤原隆親 出納	藤原隆親 平有親	藤原隆親 北面	藤原隆親 藤原隆通
道深法親王 員数	道深法親王 仁明天皇	道深法親王 手足	道深法親王 良通	道深法親王 淨衣	道深法親王 藤原有長	道深法親王 勸文	道深法親王 平時高	道深法親王 妻戸	道深法親王 源維長
清原宗尚 御事	清原宗尚 龜山天皇	清原宗尚 狼藉	清原宗尚 荒木田成定	清原宗尚 供給	清原宗尚 藤原為家	清原宗尚 天皇	清原宗尚 宗基	清原宗尚 其儀	清原宗尚 藤原忠高
丹後 天福	丹後 藤原種子	丹後 解除	丹後 憲実	丹後 建久	丹後 藤原泰通	丹後 徒末	丹後 宗清	丹後 上儀	丹後 大江重房
藤原高朝 春日	藤原高朝 知仁親王	藤原高朝 入末	藤原高朝 大中臣能隆	藤原高朝 休息	藤原高朝 橘親氏	藤原高朝 每事	藤原高朝 定宗	藤原高朝 大進	藤原高朝 章親
順徳天皇 不可	順徳天皇 藤原家経	順徳天皇 殿対	順徳天皇 藤原実基	順徳天皇 散花	順徳天皇 和氣時定	順徳天皇 御所	順徳天皇 藤原兼経	順徳天皇 職事	順徳天皇 源通方
慈神 不快	慈神 藤子内親王	慈神 遙行	慈神 後鳥羽天皇	慈神 御分	慈神 佐伯康長	慈神 今度	慈神 源仲兼	慈神 相促	慈神 藤原基氏
四条天皇 終日	四条天皇 藤原経俊	四条天皇 結婚	四条天皇 盛家	四条天皇 鳥羽	四条天皇 源家清	四条天皇 神宝	四条天皇 源定通	四条天皇 長押	四条天皇 秀仁親王

図7 トピックに関連する語句と人物 (オレンジ: 筆者関連, 黄色: 政務関連, 青: 神仏関連, ピンク: 皇族関連, 緑: 暦注)

Figure7 Words and names of people related to the topic (orange: Author's relatives, yellow: Government affairs, blue: Religious affairs, pink: royal family, green: Calendar notes)

#### 4.2 語句抽出評価

『民経記』における役職と語句の抽出精度について述べる。パターンマッチング+N-gram の手法で役職を抽出した場合の嘉禄2年4月分の出力件数は53件、人手で作成した正解件数は64件だった。N-gram とコトバンクによって抽出した場合の嘉禄2年4月分の出力件数は201件、人手で調査した正解件数は227件だった。それぞれの適合率、再現率、F値を図6に記す。

	適合率	再現率	F 値
役職	0.6875	0.6000	0.6408
語句	0.9253	0.8193	0.8691

図6 役職・語句抽出の結果

Figure6. Results of data preprocessing

#### 4.3 トピックの内容に関する検証

システムで出力した結果を元に、評価・検討を行った。例としてトピック数を10とし、LDAの結果の検証を行う。図7は10のトピックの中から考察のために選択したトピックと、上位10語句と人物である。日本史学者のアドバイスを元にカテゴリ別で色分けを行った。

Topic1 は天皇関係のトピックである。父である頼資がいることから、経光の親族と天皇間のやりとりのトピックであることがわかる。

Topic2 は暦注と経光の親族、天皇関連のトピックであることがわかる。暦注が登場することからよく登場する人名

が並ぶ。Topic1 と Topic2 からは親族の情報を得ることができた。Topic3 と Topic10 は神仏関係の語句と人名が多いことから神仏のトピックであるとわかる。Topic4~Topic7 は政務関係のトピックである。これらのトピックの差は、専門家以外が Table7-6 のみで判断するのは難しい。Topic8 は書状に関するトピックであり、ここから「紙背文書」の情報が得られる。Topic9 は神仏関係のトピックとも捉えられるが、「朝儀典礼」の情報も抽出できる。「行事」「奉行」「御祭り」「大嘗祭」などの朝儀典礼関連の単語が見られるからである。朝講堂などの寺社の地名も、長講堂御八講などの行事に関係している。

また、各トピックに関して、時系列変化と対応させることで、トピックの割合が高くなる年代を特定し、史料の文章と照らし合わせ整合性を確認した。例えば、神仏関係のトピックである Topic3 の時系列変化である 図4では、2つ目の山である嘉禎元年10月と嘉禎3年10月に、仏教行事である「維摩会」の記述があるのが確認できた。

日本史学者へのインタビューによるシステムの検証を行った。Web Application を使ってトピックの妥当性について検証していただいた。システム全体に対してのフィードバックが得られ、分析結果が有用であり、UI がわかりやすいという評価が得られた。さらに、システムの改善すべき点として、日記数を反映させた時系列グラフの作成、紙背文書や裏書などを除外した状態で分析の必要性などがあげられた。

## 5. 考察

語句・役職抽出について、ある程度の精度は確認することができたが、精度向上の余地がある。今回はケーススタディにて **Sentencepiece** を用いることはなかったが、今後後処理に外部辞書を利用するなど、この技術を適切に用いることで精度向上を目指したい。

システム全体について、人物だけではなく語句も **LDA** に組み込むことにより、即座にトピックの理解ができるようになった。また、史料に含まれていることが確認されている史実を対象に、トピック分析を行った結果その史実に整合したトピックが割り当てられていることが確認され、有用性が示された。

## 6. おわりに

本稿では史料から **N-gram** を用いることでトピック抽出の対象とする人物・語句の抽出を行い、それらを **Bag-of-Words** としてトピック抽出を行なった。また、一連の動作を **Web Application** として自動で行うシステムとして開発し、『民経記』を対象に動作の検証とインタビューによる質的評価の結果から有用性を示し、資料研究の分析過程を支援することができた。

今後の研究計画としては、より柔軟性の高く、持続可能なシステムの設計及び開発を行う。具体的には、**Web Application** 上でユーザーが **Bag-of-Words** に入れる語彙の選定や、史料によって異なる役職パターンを入力を行い、その後情報学的分析結果の確認が自由にできるようなシステムの開発を行う予定である。さらに、システムを継続して利用し、恒常的に日本史学者の史料研究を支援するようなシステムの設計を行い、提案したい。

**謝辞** 本研究の成果の一部は **JSPS** 科研費 18H03576 および 19K20626 の助成を受けたものである

## 参考文献

- [1] 中村覚:デジタルアーカイブと **Linked Data** を用いた歴史学研究支援に関する研究, デジタル・ヒューマニティーズ 129-43 2019.
- [2] 鳥居克哉, 中村覚, 山田太造, 稗方和夫「日本史学者の要求分析に基づく歴史資料のトピック推定システムの開発」, 情報処理学会第 83 回全国大会. Vol.83, No.2ZC-05, pp.1-2, 2021 年.
- [3] N. Sukhija, M. Tatineni, N. Brown, M. V. Moer, P. Rodriguez and S. Callicott: Topic Modeling and Visualization for Big Data in Social Sciences, 2016 Intl IEEE Conferences on Ubiquitous Intelligence & Computing, Advanced and Trusted Computing, Scalable Computing and Communications, Cloud and Big Data Computing, Internet of People, and Smart World Congress, Toulouse, pp. 1198-1205, 2016.
- [4] D. M. Blei, A. Y. Ng, and M. I. Jordan: Latent Dirichlet Allocation, Journal of Machine Learning Research, vol. 3, pp. 993-1022(2003).
- [5] YanHeng, Zhifeng Gao, Yuan Jiang, Xuqi Chen: Exploring hidden factors behind online food shopping from Amazon reviews:

- A topic mining approach, Journal of Retailing and Consumer Services, Volume 42, May 2018, Pages 161-168, 2018.
- [6] Kudo, T., & Richardson, J. : Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808. 06226*, 2018
  - [7] Matthew D. Hoffman, David M. Blei, and Francis Bach. 2010. Online learning for Latent Dirichlet Allocation. In Proceedings of the 23rd International Conference on Neural Information Processing Systems - Volume 1 (NIPS'10). Curran Associates Inc., Red Hook, NY, USA, 856-864.