

# 潜在的なトピックの類似度に基づくトピック追跡

芹澤翠 (指導教員：小林一郎)

## 1 はじめに

ある事象を理解する際、時間的な内容の変化を把握することで、その全体像を掴み、深く理解することが可能になる。一方、単一の話題を取り上げていると思われる文書においても記載されている話題は複数のトピックにより構成されていることも多い。そのため、正確にその話題の経時変化を捉えるためには、その様な潜在的なトピックの追跡を行う必要がある。このことを考慮し、本研究では、確率的潜在意味解析によりトピックの抽出を行い、時系列データであるニュース記事を対象にトピック追跡を行うことを目的とする。

## 2 トピック抽出

本稿では、トピック抽出に潜在的ディリクレ配分法(LDA) [1] を用いる。これは、一文書に複数トピックが含まれることを表現できるトピックモデルであり、各文書は潜在トピックの混合分布として、トピックは語の確率分布として表現される。

### 2.1 トピック内の語の特徴量

LDA によって抽出したトピック内の語の特徴量として、以下の term-score が定義されている [2]。

$$\text{term-score}_{k,v} = \hat{\beta}_{k,v} \log \left( \frac{\hat{\beta}_{k,v}}{\left( \prod_{j=1}^K \hat{\beta}_{j,v} \right)^{\frac{1}{K}}} \right) \quad (1)$$

$\hat{\beta}_{k,v}$  : トピック  $k$  での語  $v$  の出現確率  
 $K$  : トピックの総数

これは、tf-idf 値の考え方に基づいた尺度であり、単語のトピック内の出現確率  $\hat{\beta}_{k,v}$  が tf 値に相当し、語の 1 トピック内での出現しやすさを表しており、残りの部分は、全トピックで頻繁に現れる語には値が低くなるため idf 値に相当している。

### 2.2 トピック間類似度

本稿では、対象文書群の持つトピック数の決定および対象期間でのトピック追跡をトピックの類似度に基づいて行う。トピックの類似度は、抽出された各トピックをそのトピック内の特徴語とその特徴量を各次元に対応付けたベクトルのコサイン類似度によって測る。

### 2.3 トピック数の決定

LDA は予め与えられたトピック数の下にトピックを抽出するが、トピック数は文書から陽に観測することはできない。そこで、本稿では、次のようにトピック数を決定する。まず、本来存在するであろうトピック数より大きめの値を意図的にトピック数として与えてトピックを抽出する。そして、抽出したトピックに対し、閾値<sup>1</sup>以上の類似度を持つトピック組を同じ内容を持つ‘類似トピック’、その中に含まれていないトピックを‘単独トピック’と見なし、類似トピックを 1 つのトピックとしてまとめることで、‘結合トピック’を生成する。ここで、本研究における先行研究 [3] に

<sup>1</sup> 閾値は、類似度の乖離に基づき決定される。

より、複数の結合トピックに含まれるようなトピック(‘重複トピック’と呼ぶ)を主張性の弱いトピックと捉え、「単独トピック数」と「重複トピックを除いた結合トピック数」の和をその文書に潜在するトピック数と判定する。

## 3 トピック追跡

抽出したトピック群について連続する 2 日間の各トピック間の類似度が閾値<sup>1</sup>以上ならば関連があるとし、トピックを追跡する。前日から関連の付かなかったトピックは、新たに生じたトピックと見なす。以下に、本手法におけるトピック追跡の流れを説明する。

### step 1. 対象文書の前処理

本稿では、名詞を複合処理した複合語と複合処理されなかった名詞を LDA によるトピック抽出の処理対象とした。複合語は新聞社や記者により同じ意味の語でも表現方法が異なる可能性があるため、複合語の統一を対象期間の全文書に対して次のルールに基づいて行う。

- サ変接続の名詞を含む場合は複合処理を行わない
- 構成する名詞に包含関係のある複合語は、構成する名詞数の少ない複合語へ置き換える

### step 2. トピック抽出

対象期間の各日に対して、以下の処理を行う。

#### 1. トピック抽出 (1 回目)

文書に本来存在するトピック数より多めと想定されるトピック数を指定し、LDA を用いてトピック抽出を行う。

#### 2. トピック数の決定

1. において抽出されたトピック群に対し、各トピック間の類似度に基づきトピックの結合を行う。結合後の「単独トピック数」と「重複トピックを除いた結合トピック数」の和を対象文書の持つ潜在トピック数とする。

#### 3. トピック抽出 (2 回目)

決定したトピック数を指定し、再度 LDA によりトピック抽出を行う。ここで得られたトピック群を対象文書の持つトピック群とする。

### step 3. トピック追跡

抽出したトピック群について、各トピック間の類似度に基づき連続する 2 日間のトピックの関連付けを行う。

## 4 実験

提案手法によるトピック追跡の実験を行い、本手法の妥当性を確認する。

### 4.1 実験仕様

対象とするニュース記事は、ニュースサイト「YOMIURI ONLINE (読売新聞)」, 「毎日 jp (毎日新聞)」からキーワード「尖閣」を与えて収集した 2010 年 11

