

## Web ニュースからの LDA に基づく興味学習手法

東原 智幸<sup>†</sup> 渥美 雅保<sup>†</sup>創価大学 大学院 工学研究科 情報システム工学専攻<sup>†</sup>

## 1. はじめに

文書の生成過程を確率的にモデル化したトピックモデルが、文書分類、推薦システムなど、さまざまな分野に適用されている。代表的な手法である LDA[1]では、文書を複数のトピックの混合により表現し、その混合確率を単語頻度データより学習している。本論文では、WEB ニュースに対するユーザの興味応答から、興味有集合、興味無集合に分割し、各集合に対して LDA を適用し、興味有トピック、無トピック間の距離からユーザの興味をモデル化する手法を提案する。そして、LDA に入力される特徴の違いによる興味の学習と推論への影響について評価する。

## 2. LDA に基づく興味学習

## 2.1. Latent Dirichlet Allocation (LDA)

LDA は、コーパスの確率的な生成モデルであり、基本的なアイデアは、 $k$  個の潜在トピック  $z$  の混合分布として表現される。各トピックは、単語の分布により特徴づけられる。

LDA では、各ドキュメントの単語集合  $w$  に対して、コーパス  $D$  でつぎのような生成モデルを仮定する。コーパス  $D$  は、 $M$  個のドキュメントの単語集合である。

- 1) ドキュメントの単語数  $N$  の選択
- 2)  $\theta \sim$  ディリクレ分布 ( $\alpha$ ) 選択
- 3) 各単語  $w_n$  について
  - (ア) トピック  $z_n \sim$  多項分布 ( $\theta$ ) の選択
  - (イ) トピック  $z_n$  で条件づけられた多項確率  $p(w_n | z_n, \beta)$  により単語  $w_n$  を選択

この生成過程のグラフィカルモデルは図 1 のようになる。

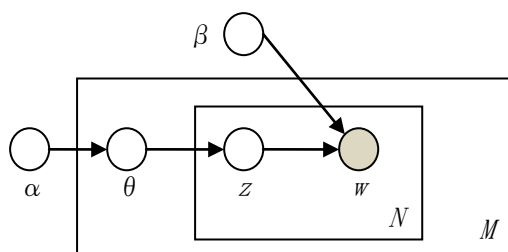


図 1 LDA のグラフィカルモデル

これらの生成過程と変分ベイズ法を用いてトピック混合分布  $\theta$  と各トピックにおける単語分布

$\beta = \{\beta^1, \dots, \beta^k, \dots, \beta^k\}$  が求まる。求めた  $\theta$  と  $\beta$  を用いて新規文書に対するトピック混合分布  $\theta^{new}$  を推論することができる。

## 2.2. LDA による興味の学習と推論

本研究では興味あり、興味なしニュースそれぞれに対して LDA を適用し、以下に示す手順に従って興味の学習と推論を行う。

## ● 学習部

- 1) 興味あり文書集合を用いて LDA を適用し、トピック混合分布  $\theta_{pos}$  と単語分布  $\beta_{pos}$  を求める
- 2) 興味なし文書集合についても同様の操作を行い  $\theta_{neg}$ ,  $\beta_{neg}$  を求める

## ● 推論部

- 1) 新規文書の単語頻度情報を作成
- 2) 単語頻度情報と興味ありデータからの学習分布を用いて LDA の推論計算により新規文書の  $\theta_{pos}^{new}$  を求め、新規文書の平均単語分布  $\beta_{pos}^{new}$  を以下の式により求める

$$\beta_{pos}^{new} = \sum_k \theta_{pos}^{new,k} \cdot \beta_{pos}^k$$

- 3) 同様に興味なしデータからの学習分布から  $\theta_{neg}^{new}$  を求め

$$\beta_{neg}^{new} = \sum_k \theta_{neg}^{new,k} \cdot \beta_{neg}^k$$

を計算

- 4)  $\beta_{pos}^{new}$  と各トピックの単語分布  $\beta_{pos}^k$  間の距離を Hellinger Distance により計算し、最短距離  $dist_{pos}$  を求める
- 5) 同様に興味なし  $\beta_{neg}^{new}$ ,  $\beta_{neg}^k$  用いて最短距離  $dist_{neg}$  を求める
- 6) 興味の判定

判定法 1 または 2 を用いて興味の判定を行う。

判定法 1: 興味あり最短距離  $dist_{pos}$  のみを用いて、 $dist_{pos}$  が閾値より小さい場合に興味ありと判定する。閾値は、0.2, 0.4, 0.6, 0.8, 1.0 の 5 つで判定する。

判定法 2: 興味ありトピックに近く、興味なしトピックからは遠いとき、興味ニュースであるという基準を用いて、 $dist_{pos} < dist_{neg}$  の場合は興味あり、条件を満たさない場合には、興味なしと判定する。

## 3. 実験

## 3.1. ニュース評価システム

ニュース評価システム(図 2)を作成し、2010

LDA-based Method for Learning User's Interests from Web News  
<sup>†</sup>Tomoyuki Higashihara, Masayasu Atsumi, Division of Information Systems Science, Graduate School of Engineering, Soka University

/12/16~2011/01/05/の期間、5名のユーザにニュース評価を行ってもらった。ユーザは、システムによって提示されたニュースリストから、閲覧したいニュースを選択し、ニュース本文を閲覧する。システムは、ニュースの更新、画面上の表示の切り替えや、評価履歴保存を行う。また、ユーザが閲覧したニュースを興味ありニュース、それ以外を興味なしと判定している。表1にユーザごとの評価ニュース総数と興味ありニュース数を掲載する。

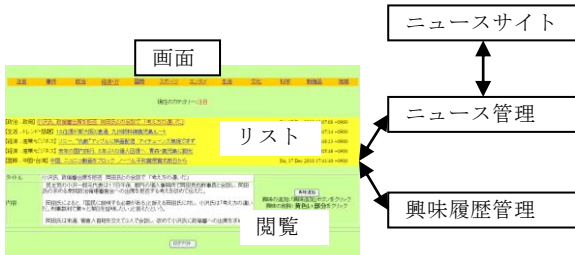


図2 評価システム概要

表1 各ユーザの評価総数と興味あり数

ユーザ	A	B	C	D	E
興味あり	8	58	19	53	49
総数	436	233	224	243	410

### 3.2. 実験データ

ユーザごとに得られた評価データを2つの集合に分割し、1つを学習用、一方を推論用として使用する。分割した集合を、興味ありニュース集合、興味なしニュース集合に分ける。興味なしニュース数は、興味なしデータが増加した場合の振舞について評価を行うため、興味ありニュース数の1倍、2倍の集合を用意する。用意された各集合に含まれるニュースタイトルと本文を形態素解析し、単語IDと頻度を求め、実験データを用意した。

#### 3.2.1. Supervised-LDA (SLDA) との比較

各ユーザの実験データに対して提案手法とSLDA[2]を適用し、推論精度比較を行った。SLDAは、文章の生成モデルはLDAと同じであり、文章のクラスの経験分布と回帰係数 $\eta$ の内積を平均とした正規分布から文書に対する応答を生成すると仮定している。

図3は、SLDA、判定法1、判定法2について、ユーザ5人の平均を求めたグラフである。横軸は、False Positive (FP) 率で興味なしニュースを興味ありと判定した比率、縦軸は、True Positive (TP) 率で、興味ありニュースを興味ありと判定した比率を表す。グラフ中のラベル1, 2は興味なしデータ数の倍率を表し、1.0, 0.8は判定法1での閾値を表す。SLDAは、興味なしデータ数を2倍にすると、TP, FP率ともに下がる傾向にある。これは、

興味なしデータが増加することでSLDA内の興味なし経験分布が平均化されたためではないかと考察される。判定法2は、興味なしデータ数が増加してもTP率は、上昇する傾向にある。また、判定法1よりも、TP率が高くなることから、興味なしデータを用いることが興味判定手法に有効であることが分かる。

図4は、ユーザC, DのTP, FP率のグラフである。ユーザCは、興味なしデータ数を2倍にすると、TP率, FP率ともに上昇する傾向があり、他2名でも同様の傾向が確認された。ユーザDでは、TP率の上昇, FP率の下降が確認され、他1名のユーザでも同様の傾向が確認された。

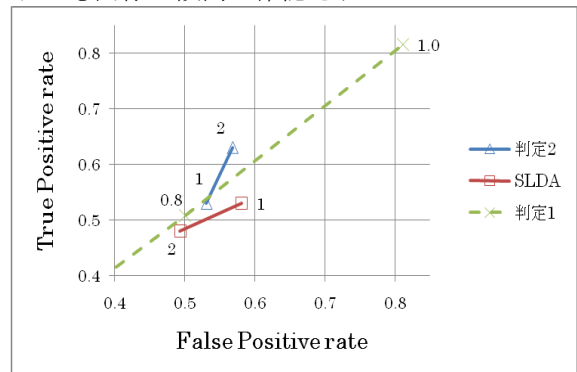


図3 5人の平均

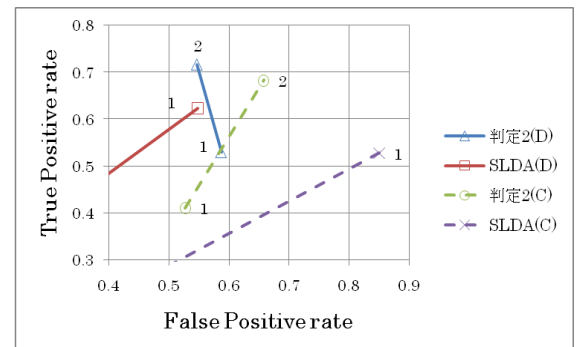


図4 ユーザC, Dの例

### 4. まとめ

LDAを興味ありデータ、興味なしデータそれぞれに適用し、興味の学習と推定を行う手法を提案した。また、実験により興味なしデータを用いることが精度向上に有効であることがわかった。今後の課題としては、長期的な評価データを用いた場合の精度比較、単純な単語の組み合わせだけではなく、格フレームのような意味的組み合わせを考慮した興味学習の手法の開発があげられる。

#### 参考文献

[1] Blei: Latent Dirichlet Allocation, JMLR, 3:993-1022, 2003.  
 [2] David Blei and Jon McAuliffe, Supervised topic models, NIPS, 2007.