

Twitterに適したオンライン学習可能なトピックモデルの検討 A Study on Online Topic Model for Twitter

佐々木 謙太郎[†] 吉川 大弘[†] 古橋 武[†]
Kentarō Sasaki Tomohiro Yoshikawa Takeshi Furuhashi

1. はじめに

近年急速に普及し、注目を集めている情報源として、Twitterを代表とするマイクロブログがある。一方、Latent Dirichlet Allocation (LDA)[1]は、様々な分野で応用されているトピックモデルであり、Twitterに対して適用した研究も数多く報告され始めている[2][3]。Twitterでは、ユーザはツイートと呼ばれる140文字以内の短いメッセージを投稿する。この140文字以内という制限により、ユーザが気軽に情報を発信できるため、Twitterは速報性、リアルタイム性が極めて高い情報源である。しかしその反面、ツイートの短さと、ノイズの多さのため、従来の自然言語処理技術が有効に働かないという問題点が存在する。そのため、上述のLDAを直接適用しても適切には機能せず、通常LDAをツイート集合に適用する場合、1ツイートを1文書とする、あるいは1ユーザの全ツイートを1文書とする方法がとられる。これに対して、Twitterの特徴を考慮し、1ツイートが1トピックから成るという仮定に基づいたTwitter-LDAが提案されている[4]。ただし、このモデルでは、通常のLDAと同様に時間発展を考慮しておらず、またデータが更新されるごとに全データを用いて再度学習を行う必要があるため、オンライン学習も困難である。そのため、Twitter上のトピックの時間発展を追跡することはできない。

本稿ではまず、Twitter-LDAを改良し、ユーザごとにトピック語と一般語との割合を推定する新しいモデルを考案する。パープレキシティを用いた実験により、その改良モデルの妥当性を評価する。さらに、考案したモデルを拡張し、Twitterに適したオンライン学習可能なトピックモデルを提案する。

2. Twitter-LDAの改良

図1(a)に、Twitter-LDAのグラフィカルモデルを示す。Twitter-LDAでは、一つのツイートを構成する単語は、どのトピックにおいても出現するような一般語の分布 θ_B と、そのツイートの持つトピック k の単語分布 θ_k のどちらか一方から生成されることを仮定している。単語がどちらの分布から生成されるかは潜在変数 y によって決まり、 $y=0$ ならば θ_B から、 $y=1$ ならば θ_k から生成される。

Zhaoらは、[4]においてLDAとの比較の際、Twitter-LDAは各トピックにおける出現確率の高い語にノイズ(一般語)が少ないことを実験により示している。しかし我々の行った実験では、Twitter-LDAは、言語モデルの性能に対する評価尺度として一般的に用いられるパープレキシティ[1]の上では、LDAよりも悪くなるという結果になった。これは、潜在変数 y の生成される分布 π が、すべてのユーザに共通であり、どのユーザもトピック語と一般語を同じ割合で含むツイートをするという仮定を置いているためだと考えられる。この仮定は実際のツイートの生成過程をうまく表現しておらず、ユーザによってトピック語と一般語の割合は変化すると仮定する方がより適切であると考えられる。そこで本稿では、図1(b)に示すように、潜在変数 y の分布 π がユーザごとに異なるという仮定に基づいたモデルを提案する。

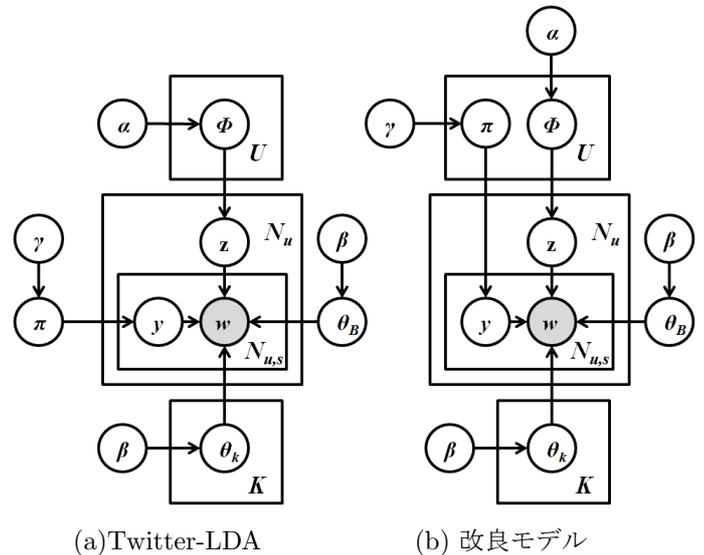


図1: グラフィカルモデル

3. パープレキシティ評価実験

2節で示したTwitter-LDAの改良モデルの妥当性を評価するために、パープレキシティに基づく評価実験を行った。パープレキシティは、学習によって得られたモデルが、実際に観測された単語をどれだけ予測出来るかを評価する指標である。パープレキシティが低いほど、モデルの予測性能が高いことを示している。実験には、2013年3月13日に収集したユーザ数409、ツイート数9176、語彙数2660の日本語ツイートデータを用いた。このツイートデータのうち、90%のツイートを学習に用い、残り10%をパープレキシティの算出に用いた。各モデルのハイパーパラメータはすべて0.1とした。また、学習にはギブスサンプリングを用い、反復回数は500とした。なお、LDAは各ユーザの全ツイートを1文書として学習を行った。LDA、Twitter-LDA、提案モデルに対して、トピック数を10~100まで10ずつ変化させ、それぞれ10試行学習を行った際のパープレキシティの平均値を図2に示す。

結果から、提案する改良モデルは、パープレキシティにおいてLDAとTwitter-LDAの両方に対して改善していることがわかる。このことから、提案するモデルがTwitterの特徴をより適切に反映しており、一般語とトピック語の割合がユーザごとに異なるという仮定が妥当であると考えられる。

4. 提案モデル

2節で示した改良モデルを、オンライン学習の可能なトピックモデルに拡張する。本稿では、時間発展を考慮し、オンライン学習可能なトピックモデルの一つである、Topic Tracking Model (TTM)[5]の機構に基づいてモデルの拡張を行う。

TTMは、時間変化するユーザの興味と、トピックの発展を追跡することができるトピックモデルである。TTMにおいて、時刻 t におけるユーザ u の興味分布 $\phi_{t,u}$ は、平均が一時刻前の興味分布の推定値 $\hat{\phi}_{t-1,u}$ であり、精度(分散の逆数)が

[†]名古屋大学大学院工学研究科計算理工学専攻

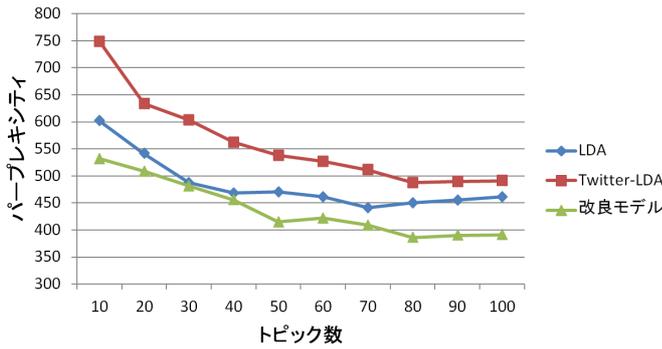


図 2: パープレキシティの比較

$\alpha_{t,u}$ である以下のディリクレ事前分布に従って生成される。

$$P(\phi_{t,u} | \hat{\phi}_{t-1,u}, \alpha_{t,u}) \propto \prod_k \phi_{t,u,k}^{\alpha_{t,u} \hat{\phi}_{t-1,u,k} - 1} \quad (1)$$

時刻 t におけるトピック k の単語分布 $\theta_{t,k}$ についても同様に、平均が一時刻前の単語分布の推定値 $\hat{\theta}_{t-1,k}$ 、精度が $\beta_{t,k}$ である以下のディリクレ事前分布から生成される。

$$P(\theta_{t,k} | \hat{\theta}_{t-1,k}, \beta_{t,k}) \propto \prod_v \theta_{t,k,v}^{\beta_{t,k} \hat{\theta}_{t-1,k,v} - 1} \quad (2)$$

本稿では、以上に述べた TTM の機構を、2 節で述べた改良モデルに取り入れたモデルを提案する。提案モデルの生成過程とグラフィカルモデルをそれぞれ図 3, 図 4 に示す。提案モデルにより、Twitter の特徴を考慮した上でユーザの興味およびトピックの流行の変化を逐次推定することが可能となる。

5. おわりに

本稿では、Twitter に適したオンライン学習が可能なトピックモデルを提案した。初めに、Twitter 向けに LDA を改良したトピックモデルである Twitter-LDA に対して、一般語とトピック語の割合はユーザごとに異なるという仮定を新たに加えた改良モデルを考案した。実験の結果から、改良モデルは Twitter-LDA や LDA と比較してパープレキシティが改善されることを示した。本稿ではさらに、この改良モデルに対して、時間発展を考慮したトピックモデルである TTM の機構を取り入れた新しいモデルを提案した。

今後は、提案モデルを Twitter におけるリアルタイムでの話題の検出、時間発展の追跡などに適用し、既存のモデルとの比較評価を行っていく予定である。

参考文献

- [1] David M. Blei, Andrew Y. Ng, and Michael I. Jordan: Latent dirichlet allocation, Machine Learning Research, Vol. 3, pp. 993-1022, 2003
- [2] Marco Pennacchiotti and Siva Gurumurthy: Investigating topic models for social media user recommendation, In WWW2011, pp. 101-102, 2011
- [3] Jianshu Weng, Ee-Peng Lim, Jing Jiang, and Qi He: Twiterrank: finding topic-sensitive influential twitterers, In WSDM 2010, 2010
- [4] Wayne Xin Zhao, Jing Jiang, Jianshu Weng, Jing He, Ee-Peng Lim, Hongfei Yan, and Xiaoming Li: Comparing twitter and traditional media using topic models,

1. Draw $\theta_{t,B} \sim \text{Dirichlet}(\lambda)$
2. For each topic $k = 1, \dots, K$,
 - (a) draw $\theta_{t,k} \sim \text{Dirichlet}(\beta_{t,k} \hat{\theta}_{t-1,k})$
3. For each user $u = 1, \dots, U$,
 - (a) draw $\phi_{t,u} \sim \text{Dirichlet}(\alpha_{t,u} \hat{\phi}_{t-1,u})$
 - (b) draw $\pi_{t,u} \sim \text{Beta}(\gamma)$
 - (c) for each tweet $s = 1, \dots, N_u$
 - i. draw $z_{t,u,s} \sim \text{Multinomial}(\phi_{t,u})$
 - ii. for each word $v = 1, \dots, N_{u,s}$
 - A. draw $y_{t,u,s,v} \sim \text{Bernoulli}(\pi_{t,u})$
 - B. draw $w_{t,u,s,v} \sim \text{Multinomial}(\theta_{t,B})$ if $y_{t,u,s,v} = 0$ or $\text{Multinomial}(\theta_{t,z_{t,u,s}})$ if $y_{t,u,s,v} = 1$

図 3: 提案モデルにおけるツイートの生成過程

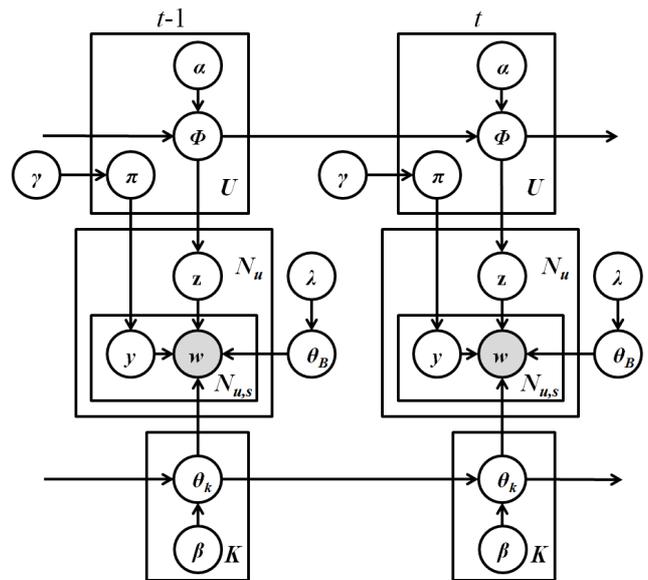


図 4: 提案モデル

In Proceedings of the 33rd European conference on Advances in information retrieval, pp. 338-349, 2011

- [5] T. Iwata, S. Watanabe, T. Yamada, and N. Ueda: Topic tracking model for analyzing consumer purchase behavior, in Proc. IJCAI, pp. 1427-1432, 2009