

大規模テキストを対象にした分散表現に基づくトピック抽出

尾崎花奈 (指導教員：小林一郎)

1 はじめに

トピックモデルは、文書の中に潜在的に存在するトピックを自動で抽出するためのモデルである。代表的な手法である LDA (Latent Dirichlet Allocation) [1] は、各文書に潜在トピックがあると仮定し、統計的に共起しやすい単語の集合が生成される要因を、潜在トピックという観測できない確率変数で定式化する。Dasら [2] によって提案された Gaussian LDA は、LDA の改良モデルであり、LDA に単語の分散表現 (Word embedding) を組み合わせている。従来の LDA ではトピックごとの単語分布をカテゴリ分布としていたのに対し、Gaussian LDA では embedding 空間上の多次元ガウス分布としている。単語の意味の関係性を事前知識として持つため、Gaussian LDA の方がトピックごとの自己相互情報量 (PMI) の値が高くなったと報告している。また、訓練データには出現しない未知語に対してもトピックを推定できるようになった。

Gaussian LDA における事後分布推定方法では、周辺化ギブスサンプリングを用いているが、本稿では SVI (Stochastic Variational Inference) [3][4] を用いることによって、計算時間の大幅な短縮を行い、大規模なコーパスに対しての効率的な処理を目指す。

2 Gaussian LDA

まず、LDA におけるトピックを生成する分布を多次元ガウス分布にするモデルが Hu ら [5] によって提案された。このモデルに、単語の分散表現を組み合わせたものが Gaussian LDA [2] である。Embedding のツールとしては、word2vec [6] を用いている。連続空間に Embedding された単語ベクトルに対し、トピック k を同空間上での多次元ガウス分布としている。

単語の分散表現を用いることによって、トピック内の意味的結束性が向上し、実験結果として従来の LDA と比較して PMI が上昇することが確認されている。また、トピックの分布に連続分布を用いることによって、従来の LDA では対応できていなかった未知語に対しても、もう一度モデルでの推定を行うことなしに潜在トピックを割り当てることが可能になっている。

Gaussian LDA の生成モデルは以下ようになる。

1. for $k = 1$ to K
 - (a) Draw topic covariance $\Sigma_k \sim \mathcal{W}^{-1}(\Psi, \mu)$
 - (b) Draw topic mean $\mu_k \sim \mathcal{N}(\mu, \frac{1}{\beta} \Sigma_k)$
2. for each document d in corpus D
 - (a) Draw topic distribution $\theta_d \sim \text{Dir}(\alpha)$
 - (b) for each word index n from 1 to N_d
 - (a) Draw a topic $z_n \sim \text{Categorical}(\theta_d)$
 - (b) Draw $v_{d,n} \sim \mathcal{N}(\mu_{z_n}, \Sigma_{z_n})$

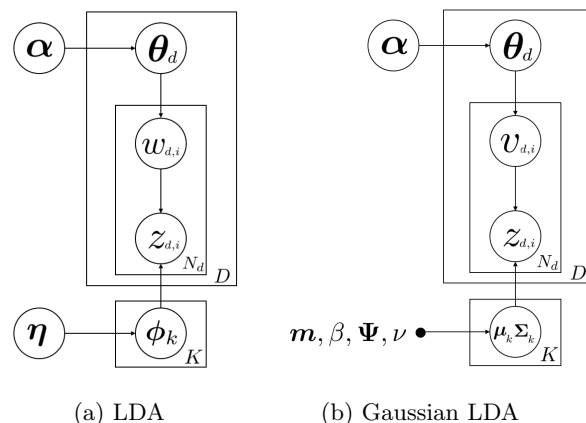


図 1: LDA と Gaussian LDA のグラフィカルモデル

ここで θ_d は従来の LDA と同じく文書 d のトピック分布を表すが、 μ_k, Σ_k はそれぞれトピック k における多次元ガウス分布の平均と分散を表している。また、 $v_{d,n}$ は単語ベクトルを表す。

LDA のグラフィカルモデルを図 1(a) に、Gaussian LDA のグラフィカルモデルを図 1(b) に示す。

3 SVI を用いたトピック推定

Gaussian LDA において、事後分布の推定に用いていたのは周辺化ギブスサンプリングであった。しかし、ギブスサンプリングは実装が簡潔である利点はあるが、計算時間が非常にかかる。

そこで本稿では、確率的変分近似法 (SVI: Stochastic Variational Inference) [4] を用いることによって、計算時間の大幅な減少を実現し、大規模なデータに対して効率的にトピック解析することを目指す。

変分ベイズにおいては、真の事後分布に対してより簡単な近似分布 $q(z, \theta, \mu, \Sigma)$ を考え、対数周辺尤度 $p(v|\alpha, \zeta)$ の変分下限を最大にする $q(z, \theta, \mu, \Sigma)$ を求める。

$$\begin{aligned} \log p(v|\alpha, \zeta) &\geq L(v, \phi, \gamma, \zeta) \\ &\triangleq \mathbb{E}_q[\log p(v, z, \theta, \mu, \Sigma|\alpha, \zeta)] \\ &\quad - \mathbb{E}_q[\log q(z, \theta, \mu, \Sigma)]. \end{aligned} \quad (1)$$

平均場近似に基づいて、近似分布 q に対して次のように各確率変数に独立性の仮定をおく。

$$q(z, \theta, \mu, \Sigma) = q(z)q(\theta)q(\mu, \Sigma). \quad (2)$$

単語ごとのトピック割り当て z のパラメータを ϕ 、文書ごとのトピック分布 θ のパラメータを γ 、トピックごとの単語分布の平均と分散 μ, Σ のパラメータを $\zeta = (\mathbf{m}, \beta, \Psi, \nu)$ とすると、近似分布 q はそれぞれ以下のように表される。

$$q(z_{di} = k) = \phi_{dw_{dik}}; \quad q(\theta_d) = \text{Dir}(\theta_d | \gamma_d),$$

$$q(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \text{NIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathbf{m}, \beta, \boldsymbol{\Psi}, \nu). \quad (3)$$

また、パラメータ ϕ, γ は以下のように定義される。

$$\phi_{dwk} \propto \exp\{\mathbb{E}_q[\log \theta_{dk}] + \mathbb{E}_q[\log N(\mathbf{v}_{dw} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]\},$$

$$\gamma_{dk} = \alpha + \sum_w n_{dw} \phi_{dwk}. \quad (4)$$

従来の変分近似法における LDA の学習では、文書データ全体に対して繰り返し学習が必要であったが、SVI は文書を逐次的に学習する。 $q(z_d), q(\theta_d)$ は各文書ごとに学習される近似事後分布であるので逐次学習を行う必要はなく、逐次学習の対象となるのは $q(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ である。よって、パラメータ $\zeta = (\mathbf{m}, \beta, \boldsymbol{\Psi}, \nu)$ の更新において、確率的自然勾配法を用いた最適化を行う。

n_t 個の単語を含む t 番目の文書において、 ζ は固定して ϕ_t と γ_t の最適化を行う。次に、 ζ の中間パラメータ $\zeta^* = (\mathbf{m}^*, \beta^*, \boldsymbol{\Psi}^*, \nu^*)$ を以下の式で求める。

$$\beta_k^* = \beta + D \sum_w n_{tw} \phi_{twk}; \quad \nu_k^* = \nu + D \sum_w n_{tw} \phi_{twk},$$

$$\mathbf{m}_k^* = \frac{\beta \mathbf{m} + D \sum_w n_{tw} \phi_{twk} \bar{\mathbf{v}}_k}{\beta_k^*},$$

$$\boldsymbol{\Psi}_k^* = \boldsymbol{\psi} + \mathbf{C}_k + \frac{\beta D \sum_w n_{tw} \phi_{twk} (\bar{\mathbf{v}}_k - \mathbf{m})(\bar{\mathbf{v}}_k - \mathbf{m})^T}{\beta_k^*}. \quad (5)$$

ここで、

$$\bar{\mathbf{v}}_k = \frac{\sum_w n_{tw} \phi_{twk} \mathbf{v}_{tw}}{\sum_w n_{tw} \phi_{twk}},$$

$$\mathbf{C}_k = D \sum_w n_{tw} \phi_{twk} (\mathbf{v}_{tw} - \bar{\mathbf{v}}_k)(\mathbf{v}_{tw} - \bar{\mathbf{v}}_k)^T. \quad (6)$$

D はコーパスの数を表しており、 ζ の計算を文書 t の複製 D 個に対して適用することを意味している。この操作によって、パラメータ ϕ, γ, ζ を更新する各イテレーションにおいてコーパス全体を必要とすることがなくなり、大規模なデータに対して逐次的な計算が可能になる。次のイテレーションに用いる ζ は、 $\rho_t \triangleq (\tau_0 + t)^{-\kappa}$, $\kappa \in (0.5, 1]$ で与えられるステップサイズによって、前回のイテレーションの ζ と更新された ζ^* に対して重みをかけることによって以下の式で求められる。

$$\zeta = (1 - \rho_t) \zeta + \rho_t \zeta^*. \quad (7)$$

また、 q のもとでの $\log \theta_{dk}$ と $\log N(\mathbf{v}_{dw} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ の期待値はそれぞれ

$$\mathbb{E}_q[\log \theta_{dk}] = \Psi(\gamma_{dk}) - \Psi\left(\sum_{i=1}^K \gamma_{di}\right),$$

$$\mathbb{E}_q[\log N(\mathbf{v}_{dw} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)] = -\frac{1}{2} \mathbf{v}_{dw}^T \langle \boldsymbol{\Sigma}_k^{-1} \rangle \mathbf{v}_{dw}$$

$$+ \mathbf{v}_{dw}^T \langle \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \rangle - \frac{1}{2} \langle \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \rangle$$

$$- \frac{1}{2} \langle \log |\boldsymbol{\Sigma}_k^{-1}| \rangle, \quad (8)$$

と表される。ただし、 Ψ はディガンマ関数を表し、 $\langle \rangle$ は期待値を表すものとする。

アルゴリズムは以下のようになる。

Algorithm 1 SVI for Gaussian LDA

Define $\rho_t \triangleq (\tau_0 + t)^{-\kappa}$
Initialize $\mathbf{m}, \beta, \boldsymbol{\Psi}, \nu$ randomly.
for $t = 0$ to ∞ **do**
 Estep:
 initialize $\gamma_{tk} = 1$ (The constant 1 is arbitrary.)
 repeat
 Set $\phi_{twk} \propto \exp\{\mathbb{E}_q[\log \theta_{tk}] + \mathbb{E}_q[\log N(\mathbf{v}_{tw} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]\}$
 Set $\gamma_{tk} = \alpha + \sum_w n_{tw} \phi_{twk}$
 until $\frac{1}{K} \sum_k |\text{change in } \gamma_{tk}| < 0.00001$
 Mstep:
 Compute ζ_k^* with Eq.(5)
 Set $\zeta = (1 - \rho_t) \zeta + \rho_t \zeta^*$
end for

4 まとめと今後の課題

本研究では、単語の分散表現を取り入れたトピックモデルにおいて、大規模テキストに対応できる効率の良い計算方法を導入する提案を行なった。今後は、実験を通じて提案手法の正当性を検証していく。

参考文献

- [1] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. (2003) Latent dirichlet allocation, *J. Mach. Learn. Res.*, 3:993-1022, March.
- [2] Rajarshi Das, Manzil Zaheer, and Chris Dyer. (2015) Gaussian LDA for topic models with word embeddings. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics*.
- [3] Matthew D. Hoffman, David M. Blei, Francis Bach. (2010) Online Learning for Latent Dirichlet Allocation. In *Advances in Neural Information Processing Systems*.
- [4] Matthew D. Hoffman, David M. Blei, Chong Wang, John Paisley. (2013) Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303-1347.
- [5] Pengfei Hu, Wenju Liu, Wei Jiang, and Zhanlei Yang. (2012) Latent topic model based on Gaussian-LDA for audio retrieval. In *Pattern Recognition*, volume 321 of *CCIS*, pages 556-563. Springer.
- [6] Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. (2013) Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 746-751, Atlanta, Georgia, June.