

LDA

(latent Dirichlet allocation)

正田 備也

masada@rikkyo.ac.jp

Contents

PLSA の復習

PLSA の問題点

PLSA のベイズ版である LDA

PLSA (probabilistic latent semantic analysis)

- ▶ 同じ文書内でも、異なる単語トークンは異なる単語多項分布から生成されうる（＝異なるトピックを表現しうる）
- ▶ 文書によって、各トピックの出現確率が異なる
- ▶ PLSA では、単語多項分布をトピック (topic) と呼ぶ
- ▶ PLSA は最もシンプルなトピックモデル
 - ▶ トピックモデルは、単語トークンの “クラスタリング”
 - ▶ 同一文書内の同一単語の異なるトークンは区別されない (bag-of-words)

Notations

- ▶ 語彙集合 $\mathcal{V} = \{1, \dots, W\}$
- ▶ トピック集合 $\mathcal{T} = \{1, \dots, K\}$
 - ▶ 語彙やトピックをその添字と同一視している。
- ▶ 文書集合 $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_D\}$
- ▶ 文書 \mathbf{x}_d の i 番目のトークンとして現れる単語を、 $x_{d,i}$ という確率変数で表す
- ▶ 文書 \mathbf{x}_d の i 番目の単語 $x_{d,i}$ が表現するトピックを、 $z_{d,i}$ という確率変数で表す
- ▶ $x_{d,i}$ の値は観測されているが、 $z_{d,i}$ の値は観測されていない
 - ▶ つまり、 $z_{d,i}$ は潜在変数。

PLSAにおける同時分布

- ▶ PLSA では、文書 x_d の i 番目のトークンがトピック k を表現し、かつそのトピックを表現するために単語 w が使われる同時確率、つまり $p(x_{d,i} = w, z_{d,i} = k)$ は

$$p(x_{d,i} = w, z_{d,i} = k) = p(z_{d,i} = k)p(x_{d,i} = w|z_{d,i} = k) \quad (1)$$

- ▶ $p(z_{d,i} = k)$ は、文書 x_d の i 番目のトークンがトピック k を表現する確率
- ▶ $p(x_{d,i} = w|z_{d,i} = k)$ は、文書 x_d の i 番目のトークンがトピック k を表現するとき、単語 w が使われる確率
- ▶ さらに PLSA では以下のように仮定する（次スライド） 5 / 19

PLSAにおいて仮定すること

- ▶ どの i, i' についても $p(z_{d,i} = k) = p(z_{d,i'} = k)$ と仮定
 - ▶ そこで、 $p(z_{d,\cdot} = k) = \theta_{d,k}$ とおく
 - ▶ 同じ文書内なら、どの単語トークンであれ、トピック k を表現する確率は、同じ（場所によってトピックの確率が違ったりしない）
- ▶ どの d, d' や i, i' についても、 $p(x_{d,i} = w | z_{d,i} = k) = p(x_{d',i'} = w | z_{d',i'} = k)$ と仮定
 - ▶ そこで、 $p(x_{\cdot,\cdot} = w | z_{\cdot,\cdot} = k) = \phi_{k,w}$ とおく
 - ▶ 同じコーパス内なら、どの文書のどの単語トークンであれ、それがトピック k を表現するために使われるならば（条件付き確率の条件の部分）、 k を表現するためにどの単語が使われるかの確率は、同じ
 - ▶ つまり、単語確率分布とトピックが一对一に対応している

PLSAにおける観測データの尤度

個々の単語トークンにおけるトピックと単語の同時分布は

$$p(x_{d,i} = w, z_{d,i} = k) = p(z_{d,i} = k)p(x_{d,i} = w|z_{d,i} = k) = \phi_{k,x_{d,i}}\theta_{d,k} \quad (2)$$

潜在変数である $z_{d,i}$ を周辺化

$$p(x_{d,i} = w) = \sum_{z_i=1}^K p(x_{d,i} = w, z_{d,i} = k) = \sum_{k=1}^K \phi_{k,x_{d,i}}\theta_{d,k} \quad (3)$$

各トークンの独立性の仮定より

$$p(\mathbf{x}_d) = \prod_{i=1}^{N_d} \left(\sum_{k=1}^K \phi_{k,x_{d,i}}\theta_{d,k} \right) \quad (4)$$

各文書の独立性の仮定より

$$p(\mathcal{D}) = \prod_{d=1}^D \prod_{i=1}^{N_d} \left(\sum_{k=1}^K \phi_{k,x_{d,i}}\theta_{d,k} \right) \quad (5)$$

Contents

PLSA の復習

PLSA の問題点

PLSA のベイズ版である LDA

PLSAの問題点とベイズ化による改良

- ▶ 各文書におけるトピック確率 $\theta_d = (\theta_{d,1}, \dots, \theta_{d,K})$ に関して、異なる文書の間で何の関係性も仮定されていない
 - ▶ θ_d と $\theta_{d'}$ の間に何の関係もない。
- ▶ このことが過学習をもたらすかもしれない
- ▶ そこで、全文書にわたって θ_d が同一のディリクレ事前分布 $\text{Dir}(\alpha)$ から draw されると仮定する
- ▶ 他はそのまま
 - ▶ 各トピックの単語確率 ϕ_k についても別のディリクレ分布 $\text{Dir}(\beta)$ を導入できるが、これはそうしなくてもよい。

Contents

PLSA の復習

PLSA の問題点

PLSA のベイズ版である LDA

PLSA と LDA の比較

PLSA における \mathbf{x}_d の尤度

$$\begin{aligned} p(\mathbf{x}_d; \boldsymbol{\theta}_d, \Phi) &= \sum_{\mathbf{z}_d} p(\mathbf{x}_d, \mathbf{z}_d; \boldsymbol{\theta}_d, \Phi) = \sum_{\mathbf{z}_d} p(\mathbf{z}_d; \boldsymbol{\theta}_d) p(\mathbf{x}_d | \mathbf{z}_d; \Phi) \\ &= \prod_{i=1}^{N_d} \left(\sum_{z_{d,i}=1}^K p(z_{d,i}; \boldsymbol{\theta}_d) p(x_{d,i} | z_{d,i}; \Phi) \right) = \prod_{i=1}^{N_d} \left(\sum_{k=1}^K \phi_{k,x_{d,i}} \theta_{d,k} \right) \end{aligned}$$

LDA における \mathbf{x}_d の尤度

$$\begin{aligned} p(\mathbf{x}_d; \Phi, \alpha) &= \int p(\boldsymbol{\theta}_d; \alpha) p(\mathbf{x}_d | \boldsymbol{\theta}_d; \Phi) d\boldsymbol{\theta}_d \\ &= \int \sum_{\mathbf{z}_d} p(\boldsymbol{\theta}_d; \alpha) p(\mathbf{z}_d | \boldsymbol{\theta}_d) p(\mathbf{x}_d | \mathbf{z}_d; \Phi) d\boldsymbol{\theta}_d \end{aligned} \tag{6}$$

LDAの変分ベイズ法

Jensen の不等式を適用して ELBO を求める

$$\begin{aligned}\ln p(\mathbf{x}_d; \Phi, \alpha) &= \ln \int \sum_{\mathbf{z}_d} p(\boldsymbol{\theta}_d; \alpha) p(\mathbf{z}_d | \boldsymbol{\theta}_d) p(\mathbf{x}_d | \mathbf{z}_d; \Phi) d\boldsymbol{\theta}_d \\ &= \ln \int \sum_{\mathbf{z}_d} q(\mathbf{z}_d, \boldsymbol{\theta}_d) \frac{p(\boldsymbol{\theta}_d; \alpha) p(\mathbf{z}_d | \boldsymbol{\theta}_d) p(\mathbf{x}_d | \mathbf{z}_d; \Phi)}{q(\mathbf{z}_d, \boldsymbol{\theta}_d)} d\boldsymbol{\theta}_d \\ &\geq \int \sum_{\mathbf{z}_d} q(\mathbf{z}_d, \boldsymbol{\theta}_d) \ln \frac{p(\boldsymbol{\theta}_d; \alpha) p(\mathbf{z}_d | \boldsymbol{\theta}_d) p(\mathbf{x}_d | \mathbf{z}_d; \Phi)}{q(\mathbf{z}_d, \boldsymbol{\theta}_d)} d\boldsymbol{\theta}_d\end{aligned}\tag{7}$$

$q(\boldsymbol{\theta}_d)$ を求める

$q(\mathbf{z}_d, \boldsymbol{\theta}_d)$ は $q(\mathbf{z}_d)q(\boldsymbol{\theta}_d)$ と factorize すると仮定する。そして、 $q(\mathbf{z}_d)$ を固定する。

$$\begin{aligned}\ln p(\mathbf{x}_d; \boldsymbol{\Phi}, \boldsymbol{\alpha}) &\geq \int \sum_{\mathbf{z}_d} q(\mathbf{z}_d)q(\boldsymbol{\theta}_d) \ln \frac{p(\boldsymbol{\theta}_d; \boldsymbol{\alpha})p(\mathbf{z}_d|\boldsymbol{\theta}_d)p(\mathbf{x}_d|\mathbf{z}_d; \boldsymbol{\Phi})}{q(\mathbf{z}_d)q(\boldsymbol{\theta}_d)} d\boldsymbol{\theta}_d \\ &= \int q(\boldsymbol{\theta}_d) \left[\sum_{\mathbf{z}_d} q(\mathbf{z}_d) \ln p(\boldsymbol{\theta}_d; \boldsymbol{\alpha})p(\mathbf{z}_d|\boldsymbol{\theta}_d) \right] d\boldsymbol{\theta}_d - \int q(\boldsymbol{\theta}_d) \ln q(\boldsymbol{\theta}_d) d\boldsymbol{\theta}_d + \text{const.} \\ &= -D_{\text{KL}}(q(\boldsymbol{\theta}_d) \parallel \frac{1}{Z} \exp \left[\sum_{\mathbf{z}_d} q(\mathbf{z}_d) \ln p(\boldsymbol{\theta}_d; \boldsymbol{\alpha})p(\mathbf{z}_d|\boldsymbol{\theta}_d) \right]) + \text{const.}\end{aligned}\tag{8}$$

以上より、 $q(\boldsymbol{\theta}_d) \propto \exp \left[\sum_{\mathbf{z}_d} q(\mathbf{z}_d) \ln p(\boldsymbol{\theta}_d; \boldsymbol{\alpha})p(\mathbf{z}_d|\boldsymbol{\theta}_d) \right]$ のとき、ELBO は最大。つまり、 $q(\boldsymbol{\theta}_d) \propto p(\boldsymbol{\theta}_d; \boldsymbol{\alpha}) \exp \left[\sum_{\mathbf{z}_d} q(\mathbf{z}_d) \ln p(\mathbf{z}_d|\boldsymbol{\theta}_d) \right]$ のとき、ELBO は最大。

$$\begin{aligned}
\sum_{\mathbf{z}_d} q(\mathbf{z}_d) \ln p(\mathbf{z}_d | \boldsymbol{\theta}_d) &= \sum_{\mathbf{z}_d} q(\mathbf{z}_d) \ln \prod_{i=1}^{N_d} \theta_{d, z_{d,i}} = \sum_{i=1}^{N_d} \sum_{\mathbf{z}_d} q(\mathbf{z}_d) \ln \theta_{d, z_{d,i}} \\
&= \sum_{i=1}^{N_d} \sum_{z_{d,i}=1}^K q(z_{d,i}) \ln \theta_{d, z_{d,i}} = \sum_{k=1}^K \left(\sum_{i=1}^{N_d} q(z_{d,i} = k) \right) \ln \theta_{d,k} = \sum_{k=1}^K N_{d,k} \ln \theta_{d,k} \quad (9)
\end{aligned}$$

ただし、 $N_{d,k} \equiv \sum_{i=1}^{N_d} q(z_{d,i} = k)$ と定義した。

よって

$$\begin{aligned}
q(\boldsymbol{\theta}_d) &\propto \prod_{k=1}^K \theta_{d,k}^{\alpha_k - 1} \times \exp \left[\sum_{k=1}^K N_{d,k} \ln \theta_{d,k} \right] \\
&= \prod_{k=1}^K \theta_{d,k}^{\alpha_k + N_{d,k} - 1} \quad (10)
\end{aligned}$$

これは、変分事後分布 $q(\boldsymbol{\theta}_d)$ がディリクレ分布であることを意味する。

また、 $q(\boldsymbol{\theta}_d)$ のパラメータを ζ_d とすると、 $\zeta_{d,k} = \alpha_k + N_{d,k}$ が成り立つ。

$q(\mathbf{z}_d)$ を求める

今度は $q(\boldsymbol{\theta}_d)$ を固定する。

$$\begin{aligned}\ln p(\mathbf{x}_d; \boldsymbol{\Phi}, \boldsymbol{\alpha}) &\geq \int \sum_{\mathbf{z}_d} q(\mathbf{z}_d) q(\boldsymbol{\theta}_d) \ln \frac{p(\boldsymbol{\theta}_d; \boldsymbol{\alpha}) p(\mathbf{z}_d | \boldsymbol{\theta}_d) p(\mathbf{x}_d | \mathbf{z}_d; \boldsymbol{\Phi})}{q(\mathbf{z}_d) q(\boldsymbol{\theta}_d)} d\boldsymbol{\theta}_d \\ &= \sum_{\mathbf{z}_d} q(\mathbf{z}_d) \left[\ln p(\mathbf{x}_d | \mathbf{z}_d; \boldsymbol{\Phi}) + \int q(\boldsymbol{\theta}_d) \ln p(\mathbf{z}_d | \boldsymbol{\theta}_d) d\boldsymbol{\theta}_d \right] - \sum_{\mathbf{z}_d} q(\mathbf{z}_d) \ln q(\mathbf{z}_d) + \text{const.} \\ &= -D_{\text{KL}}(q(\mathbf{z}_d) \parallel \frac{1}{Z} \exp \left[\ln p(\mathbf{x}_d | \mathbf{z}_d; \boldsymbol{\Phi}) + \int q(\boldsymbol{\theta}_d) \ln p(\mathbf{z}_d | \boldsymbol{\theta}_d) d\boldsymbol{\theta}_d \right]) + \text{const.} \quad (11)\end{aligned}$$

以上より、 $q(\mathbf{z}_d) \propto p(\mathbf{x}_d | \mathbf{z}_d; \boldsymbol{\Phi}) \exp \left[\int q(\boldsymbol{\theta}_d) \ln p(\mathbf{z}_d | \boldsymbol{\theta}_d) d\boldsymbol{\theta}_d \right]$ のとき、ELBO は最大。

$q(\boldsymbol{\theta}_d)$ がパラメータ ζ_d のディリクレ分布であることを使うと、

$$\begin{aligned}
\int q(\boldsymbol{\theta}_d; \zeta_d) \ln p(\mathbf{z}_d | \boldsymbol{\theta}_d) d\boldsymbol{\theta}_d &= \int q(\boldsymbol{\theta}_d; \zeta_d) \ln \prod_{i=1}^{N_d} \theta_{d,z_{d,i}} d\boldsymbol{\theta}_d = \sum_{i=1}^{N_d} \int q(\boldsymbol{\theta}_d; \zeta_d) \ln \theta_{d,z_{d,i}} d\boldsymbol{\theta}_d \\
&= \sum_{i=1}^{N_d} \left\{ \psi(\zeta_{d,z_{d,i}}) - \psi\left(\sum_k \zeta_{d,k}\right) \right\} = \sum_{i=1}^{N_d} \psi(\zeta_{d,z_{d,i}}) + \text{const.}
\end{aligned} \tag{12}$$

よって、

$$\begin{aligned}
q(\mathbf{z}_d) &\propto p(\mathbf{x}_d | \mathbf{z}_d; \boldsymbol{\Phi}) \exp \left[\int q(\boldsymbol{\theta}_d) \ln p(\mathbf{z}_d | \boldsymbol{\theta}_d) d\boldsymbol{\theta}_d \right] \\
&= \prod_{i=1}^{N_d} \phi_{z_{d,i}, x_{d,i}} \exp \left(\sum_{i=1}^{N_d} \psi(\zeta_{d,z_{d,i}}) \right) = \prod_{i=1}^{N_d} \phi_{z_{d,i}, x_{d,i}} \exp \left(\psi(\zeta_{d,z_{d,i}}) \right)
\end{aligned} \tag{13}$$

つまり、

$$q(z_{d,i} = k) = \frac{\phi_{k,x_{d,i}} \exp \left(\psi(\zeta_{d,k}) \right)}{\sum_{l=1}^K \phi_{l,x_{d,i}} \exp \left(\psi(\zeta_{d,l}) \right)} \tag{14}$$

変分事後分布を使って ELBO を書き下す

$$\begin{aligned}\ln p(\mathbf{x}_d; \Phi, \alpha) &\geq \int \sum_{\mathbf{z}_d} q(\mathbf{z}_d) q(\boldsymbol{\theta}_d) \ln \frac{p(\boldsymbol{\theta}_d; \alpha) p(\mathbf{z}_d | \boldsymbol{\theta}_d) p(\mathbf{x}_d | \mathbf{z}_d; \Phi)}{q(\mathbf{z}_d) q(\boldsymbol{\theta}_d)} d\boldsymbol{\theta}_d \\ &= \int \sum_{\mathbf{z}_d} q(\mathbf{z}_d) q(\boldsymbol{\theta}_d) \ln p(\mathbf{z}_d | \boldsymbol{\theta}_d) d\boldsymbol{\theta}_d + \sum_{\mathbf{z}_d} q(\mathbf{z}_d) \ln p(\mathbf{x}_d | \mathbf{z}_d; \Phi) \\ &\quad - D_{\text{KL}}(q(\boldsymbol{\theta}_d) \parallel p(\boldsymbol{\theta}_d; \alpha)) - \sum_{\mathbf{z}_d} q(\mathbf{z}_d) \ln q(\mathbf{z}_d)\end{aligned}\tag{15}$$

式 (15) の右辺の最初の項を計算してみる。

$$\begin{aligned}\int \sum_{\mathbf{z}_d} q(\mathbf{z}_d) q(\boldsymbol{\theta}_d) \ln p(\mathbf{z}_d | \boldsymbol{\theta}_d) d\boldsymbol{\theta}_d &= \sum_{i=1}^{N_d} \sum_{\mathbf{z}_{d,i}=1}^K q(\mathbf{z}_{d,i}) \int q(\boldsymbol{\theta}_d) \ln \theta_{d,\mathbf{z}_{d,i}} d\boldsymbol{\theta}_d \\ &= \sum_{k=1}^K \left(\sum_{i=1}^{N_d} q(\mathbf{z}_{d,i} = k) \right) \left(\psi(\zeta_{d,k}) - \psi\left(\sum_l \zeta_{d,l}\right) \right)\end{aligned}\tag{16}$$

式 (15) の右辺の 2 番目の項を計算してみる。

$$\sum_{\mathbf{z}_d} q(\mathbf{z}_d) \ln p(\mathbf{x}_d | \mathbf{z}_d; \Phi) = \sum_{i=1}^{N_d} \sum_{z_{d,i}=1}^K q(z_{d,i}) \ln \phi_{z_{d,i}, x_{d,i}} = \sum_{i=1}^{N_d} \sum_{k=1}^K q(z_{d,i} = k) \ln \phi_{k, x_{d,i}} \quad (17)$$

トピック単語確率 Φ は、この項の全文書についての和 $\sum_{d=1}^D \sum_{i=1}^{N_d} \sum_{k=1}^K q(z_{d,i} = k) \ln \phi_{k, x_{d,i}}$ を最大化することで求めることができる。(ELBO の中で Φ を含むのはこの項だけだから。) 全文書の ELBO の和を \mathcal{L} と書くことにする。

$\sum_{w=1}^W \phi_{k,w} = 1$ が満たされなければならないので、ラグランジュの未定乗数法を使えば、

$$\frac{\partial \mathcal{L}}{\partial \phi_{k,w}} + \frac{\partial}{\partial \phi_{k,w}} \lambda_k \left(1 - \sum_{w=1}^W \phi_{k,w} \right) = \frac{\sum_{d=1}^D \sum_{i=1}^{N_d} q(z_{d,i} = k) \delta(x_{d,i} = w)}{\phi_{k,w}} - \lambda_k \quad (18)$$

$$\therefore \phi_{k,w} = \frac{\sum_{d=1}^D \sum_{i=1}^{N_d} q(z_{d,i} = k) \delta(x_{d,i} = w)}{\sum_{d=1}^D \sum_{i=1}^{N_d} q(z_{d,i} = k)} \quad (19)$$

LDAの変分ベイズ法のまとめ

以下の更新を繰り返し実行する。

$$q(z_{d,i} = k) \leftarrow \frac{\phi_{k,x_{d,i}} \exp(\psi(\zeta_{d,k}))}{\sum_{l=1}^K \phi_{l,x_{d,i}} \exp(\psi(\zeta_{d,l}))} \quad (20)$$

$$\zeta_{d,k} \leftarrow \alpha_k + \sum_{i=1}^{N_d} q(z_{d,i} = k) \quad (21)$$

$$\phi_{k,w} \leftarrow \frac{\sum_{d=1}^D \sum_{i=1}^{N_d} q(z_{d,i} = k) \delta(x_{d,i} = w)}{\sum_{d=1}^D \sum_{i=1}^{N_d} q(z_{d,i} = k)} \quad (22)$$