

変分ベイズ推論 (1)

正田 備也

masada@rikkyo.ac.jp

Contents

変分ベイズ推論とは

変分ベイズ推論の実例

前回のEMアルゴリズムでの議論のパターン

- ▶ 潜在変数 $\mathcal{Z} = \{z_1, \dots, z_N\}$ を含むモデリングを行いたい
- ▶ 確率モデルを指定することで同時分布

$p(\mathcal{X}, \mathcal{Z}) = p(\mathcal{Z})p(\mathcal{X}|\mathcal{Z}) = \prod_{i=1}^N p(z_i)p(x_i|z_i)$ が得られる

- ▶ 潜在変数 \mathcal{Z} の周辺化 $\sum_{\mathcal{Z}} p(\mathcal{X}, \mathcal{Z})$ により観測データの尤度 $p(\mathcal{X})$ は得られるのだが、大抵この尤度は計算できない
- ▶ Jensen の不等式を使い、対数尤度 $\ln p(\mathcal{X})$ の下界を得る

$$\ln p(\mathcal{X}) \geq \sum_{i=1}^N \sum_{z_i} q_{i,z_i} \ln \frac{p(z_i)p(x_i|z_i)}{q_{i,z_i}}$$

- ▶ この下界を最大化することで、様々な未知量を推定する 3 / 20

ベイズ的なモデリング

- ▶ 観測データを表す確率変数を $\mathcal{X} \equiv \{x_1, \dots, x_N\}$ とする
- ▶ 確率モデルのパラメータを表す確率変数を Θ とする
- ▶ ベイズ的なモデリングで知りたいのは、事後分布 $p(\Theta|\mathcal{X})$

$$p(\Theta|\mathcal{X}) = \frac{p(\mathcal{X}|\Theta)p(\Theta)}{p(\mathcal{X})} \quad (1)$$

- ▶ 変分ベイズ推論は $p(\Theta|\mathcal{X})$ を近似する分布 $q(\Theta)$ を求める
 - ▶ $q(\Theta)$ を変分事後分布 (variational posterior distribution) と呼ぶ
 - ▶ $q(\Theta)$ には比較的扱いやすい分布を選ぶ

EM アルゴリズムでの議論のパターンを適用

- ▶ 潜在変数 Θ を含むモデリングを行いたい
- ▶ 確率モデルを指定することで観測データと潜在変数の同時分布 $p(\mathcal{X}, \Theta) = p(\Theta)p(\mathcal{X}|\Theta) = p(\Theta) \prod_{i=1}^N p(x_i|\Theta)$ が得られる
- ▶ 潜在変数 Θ の周辺化 $\int p(\mathcal{X}, \Theta)d\Theta$ により観測データの周辺尤度 $p(\mathcal{X})$ は得られるのだが、大抵この尤度は計算できない
- ▶ Jensen の不等式を使い、対数周辺尤度 $\ln p(\mathcal{X})$ の下界を得る

$$\ln p(\mathcal{X}) \geq \int q(\Theta) \ln \frac{p(\Theta)p(\mathcal{X}|\Theta)}{q(\Theta)} d\Theta$$

- ▶ この下界を最大化することで、様々な未知量を推定する
 - ▶ この下界を ELBO(Evidence Lower BOund; 変分下限) と呼ぶ

変分ベイズ推論 (variational inference) とは

- ▶ Jensen の不等式を適用することで、ELBO を次のように得た

$$\ln p(\mathcal{X}) \geq \int q(\Theta) \ln \frac{p(\Theta)p(\mathcal{X}|\Theta)}{q(\Theta)} d\Theta$$

- ▶ 実は、ELBO を大きくすればするほど、 Θ が従う確率分布である $q(\Theta)$ が、事後分布 $p(\Theta|\mathcal{X})$ に近くなっていく
- ▶ つまり、この $q(\Theta)$ は、事後分布を近似する分布とみなせるような分布になっている
- ▶ $q(\Theta)$ を、変分事後分布 (variational posterior) と呼ぶ

対数周辺尤度と ELBO の差

- Jensen の不等式の左辺から右辺を引いたものを求めてみる

$$\begin{aligned} & \ln p(\mathcal{X}) - \int q(\Theta) \ln \frac{p(\Theta|\mathcal{X})p(\mathcal{X})}{q(\Theta)} d\Theta \\ &= \ln p(\mathcal{X}) - \int q(\Theta) \ln \frac{p(\Theta|\mathcal{X})}{q(\Theta)} d\Theta - \int q(\Theta) \ln p(\mathcal{X}) d\Theta \\ &= \ln p(\mathcal{X}) - \int q(\Theta) \ln \frac{p(\Theta|\mathcal{X})}{q(\Theta)} d\Theta - \ln p(\mathcal{X}) \int q(\Theta) d\Theta \\ &= \int q(\Theta) \ln \frac{q(\Theta)}{p(\Theta|\mathcal{X})} d\Theta = D_{\text{KL}}(q(\Theta) \parallel p(\Theta|\mathcal{X})) \end{aligned} \quad (2)$$

\therefore ELBO を $\ln p(\mathcal{X})$ に近づける $\Leftrightarrow q(\Theta)$ を $p(\Theta|\mathcal{X})$ に近づける

「変分 (variational)」の意味

- ▶ ELBO の最大化は、 $q(\Theta)$ を変化させることでおこなう
- ▶ このとき、 $q(\Theta)$ の密度関数そのものを変化させる
- ▶ 逆に言うと、 $q(\Theta)$ の密度関数が特定のかたちを持つと仮定した上で、その関数のパラメータを動かすのではない
 - ▶ パラメータについて微分することで最大化問題を解くのではなく、いわば “関数について微分する” ことで最大化問題を解いている
- ▶ 関数のかたちを決めてそのパラメータを動かすのではなく、関数のかたち自体を動かすことで問題を解く方法を、変分法と呼ぶ (cf. 汎関数微分)

実際の変分推論

- ▶ 実際には、 $q(\Theta)$ の密度関数が特定のかたちを持つと仮定した上で、その関数のパラメータを動かすことによって、ELBO を最大化することも多い

$$\ln p(\mathcal{X}) \geq \int q(\Theta) \ln \frac{p(\Theta)p(\mathcal{X}|\Theta)}{q(\Theta)} d\Theta$$

- ▶ 例えば、 $q(\Theta)$ が多変量正規分布だと仮定して、ELBO を最大化するような平均パラメータと共分散行列パラメータを求める、など
- ▶ 変分オートエンコーダでは、 $q(\Theta)$ が多変量正規分布だと仮定し、さらにその共分散行列が対角行列だと仮定する

Contents

変分ベイズ推論とは

変分ベイズ推論の実例

例：メッセージ受信数の変化点の検知

- ▶ この授業の最初に採り上げた例

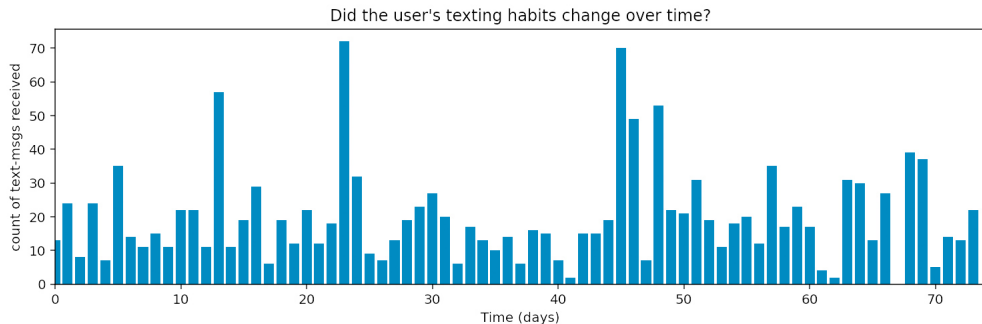


Figure: メッセージの受信数

モデルを指定する

- ▶ c_n が n 日目の受信数、 τ が受信数の変化点、 λ_1, λ_2 がそれぞれ $n < \tau, n \geq \tau$ の場合のポアソン分布のパラメータとする

$$\tau \sim \text{Uniform}(1, N)$$

$$\lambda_1 \sim \text{Gam}(a, b)$$

$$\lambda_2 \sim \text{Gam}(a, b)$$

$$c_n \sim \text{Poi}(\lambda_1) \quad \text{for } n < \tau$$

$$c_n \sim \text{Poi}(\lambda_2) \quad \text{for } n \geq \tau$$

ELBO を求める

同時分布は、 $\mathbf{c} = \{c_1, \dots, c_N\}$ として

$$\begin{aligned} p(\mathbf{c}, \lambda_1, \lambda_2, \tau) &= p(\mathbf{c} | \lambda_1, \lambda_2, \tau) p(\lambda_1; a, b) p(\lambda_2; a, b) p(\tau) \\ &= p(\lambda_1; a, b) p(\lambda_2; a, b) p(\tau) \prod_{n=1}^N p(c_n | \lambda_1)^{\delta(n < \tau)} p(c_n | \lambda_2)^{\delta(n \geq \tau)} \end{aligned} \quad (3)$$

ELBO は

$$\begin{aligned} \ln p(\mathbf{c}) &= \ln \int \sum_{\tau} p(\mathbf{c}, \lambda_1, \lambda_2, \tau) d\lambda_1 d\lambda_2 \\ &\geq \int \sum_{\tau} q(\lambda_1, \lambda_2, \tau) \ln \frac{p(\mathbf{c}, \lambda_1, \lambda_2, \tau)}{q(\lambda_1, \lambda_2, \tau)} d\lambda_1 d\lambda_2 \end{aligned} \quad (4)$$

平均場近似

- ▶ 変分事後分布 $q(\lambda_1, \lambda_2, \tau)$ について、
 $q(\lambda_1, \lambda_2, \tau) = q(\lambda_1)q(\lambda_2)q(\tau)$ と分解できることを仮定する
 - ▶ このような仮定を平均場近似 (mean field approximation) と呼ぶ

$$\begin{aligned}\ln p(\mathbf{c}) &\geq \int \sum_{\tau} q(\lambda_1, \lambda_2, \tau) \ln \frac{p(\mathbf{c}, \lambda_1, \lambda_2, \tau)}{q(\lambda_1, \lambda_2, \tau)} d\lambda_1 d\lambda_2 \\ &= \int \sum_{\tau} q(\lambda_1)q(\lambda_2)q(\tau) \ln \frac{p(\mathbf{c}, \lambda_1, \lambda_2, \tau)}{q(\lambda_1)q(\lambda_2)q(\tau)} d\lambda_1 d\lambda_2 \quad (5)\end{aligned}$$

- ▶ ELBO の最大化 $\Leftrightarrow D_{\text{KL}}(q(\lambda_1)q(\lambda_2)q(\tau) \parallel p(\lambda_1, \lambda_2, \tau|\mathbf{c}))$ の最小化 (cf. 式 (2))

$q(\lambda_1)$ のかたちを求める

$q(\lambda_2)$ と $q(\tau)$ を固定し、 $D_{\text{KL}}(q(\lambda_1)q(\lambda_2)q(\tau) \parallel p(\lambda_1, \lambda_2, \tau|\mathbf{c}))$ を最小にする $q(\lambda_1)$ を求める。

$$\begin{aligned} & D_{\text{KL}}(q(\lambda_1)q(\lambda_2)q(\tau) \parallel p(\lambda_1, \lambda_2, \tau|\mathbf{c})) \\ &= \int \sum_{\tau} q(\lambda_1)q(\lambda_2)q(\tau) \ln \frac{q(\lambda_1)q(\lambda_2)q(\tau)}{p(\lambda_1, \lambda_2, \tau|\mathbf{c})} d\lambda_1 d\lambda_2 \\ &= \int q(\lambda_1) \left\{ \ln q(\lambda_1) - \int \sum_{\tau} q(\lambda_2)q(\tau) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_2 \right\} d\lambda_1 + \text{const.} \\ &= \int q(\lambda_1) \ln \frac{q(\lambda_1)}{\exp \int \sum_{\tau} q(\lambda_2)q(\tau) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_2} d\lambda_1 + \text{const.} \\ &= D_{\text{KL}}(q(\lambda_1) \parallel \frac{1}{Z} \exp \int q(\lambda_2)q(\tau) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_2 d\tau) + \text{const.} \end{aligned} \tag{6}$$

$q(\lambda_1) = \frac{1}{Z} \exp \int \sum_{\tau} q(\lambda_2)q(\tau) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_2$ のとき、上の KL 情報量は最小。つまり、
 $\ln q(\lambda_1) = \int \sum_{\tau} q(\lambda_2)q(\tau) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_2 - \ln Z$ のとき、上の KL 情報量は最小。

$$\begin{aligned}
\ln q(\lambda_1) &= \int \sum_{\tau} q(\lambda_2) q(\tau) \ln \left\{ p(\lambda_1; a, b) p(\lambda_2; a, b) p(\tau) \prod_{n=1}^N p(c_n | \lambda_1)^{\delta(n < \tau)} p(c_n | \lambda_2)^{\delta(n \geq \tau)} \right\} d\lambda_2 - \ln Z \\
&= \ln p(\lambda_1; a, b) + \int q(\lambda_2) \ln p(\lambda_2; a, b) d\lambda_2 + \sum_{\tau} q(\tau) \ln p(\tau) \\
&\quad + \sum_{n=1}^N \sum_{\tau} q(\tau) \delta(n < \tau) \ln p(c_n | \lambda_1) + \sum_{n=1}^N \int \sum_{\tau} q(\lambda_2) q(\tau) \delta(n \geq \tau) \ln p(c_n | \lambda_2) d\lambda_2 - \ln Z \\
&= \ln \frac{b^a}{\Gamma(a)} \lambda_1^{a-1} e^{-b\lambda_1} + \sum_{n=1}^N \left(\sum_{\tau} q(\tau) \delta(n < \tau) \right) \ln \frac{\lambda_1^{c_n} e^{-\lambda_1}}{c_n!} + \text{const.} \\
&= \left(a - 1 + \sum_{n=1}^N \left(\sum_{\tau} q(\tau) \delta(n < \tau) \right) c_n \right) \ln \lambda_1 - \left(b + \sum_{n=1}^N \left(\sum_{\tau} q(\tau) \delta(n < \tau) \right) \right) \lambda_1 + \text{const.}
\end{aligned}$$

よって、 $q(\lambda_1)$ は、shape パラメータが $a + \sum_{n=1}^N \left(\sum_{\tau} q(\tau) \delta(n < \tau) \right) c_n$ で、rate パラメータが $b + \sum_{n=1}^N \left(\sum_{\tau} q(\tau) \delta(n < \tau) \right)$ のガンマ分布となる。
 $q(\lambda_2)$ についても同様に計算すると、やはりガンマ分布であることが分かる。

$q(\tau)$ のかたちを求める

$q(\lambda_1)$ と $q(\lambda_2)$ を固定し、 $D_{\text{KL}}(q(\lambda_1)q(\lambda_2)q(\tau) \parallel p(\lambda_1, \lambda_2, \tau|\mathbf{c}))$ を最小にする $q(\tau)$ を求める。

$$\begin{aligned} & D_{\text{KL}}(q(\lambda_1)q(\lambda_2)q(\tau) \parallel p(\lambda_1, \lambda_2, \tau|\mathbf{c})) \\ &= \int \sum_{\tau} q(\lambda_1)q(\lambda_2)q(\tau) \ln \frac{q(\lambda_1)q(\lambda_2)q(\tau)}{p(\lambda_1, \lambda_2, \tau|\mathbf{c})} d\lambda_1 d\lambda_2 \\ &= \sum_{\tau} q(\tau) \left\{ \ln q(\tau) - \int q(\lambda_1)q(\lambda_2) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_1 d\lambda_2 \right\} + \text{const.} \\ &= \sum_{\tau} q(\tau) \ln \frac{q(\tau)}{\exp \int q(\lambda_1)q(\lambda_2) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_1 d\lambda_2} + \text{const.} \\ &= D_{\text{KL}}(q(\tau) \parallel \frac{1}{Z} \exp \int q(\lambda_1)q(\lambda_2) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_1 d\lambda_2) + \text{const.} \end{aligned} \tag{7}$$

$q(\tau) = \frac{1}{Z} \exp \int q(\lambda_1)q(\lambda_2) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_1 d\lambda_2$ のとき、上の KL 情報量は最小。つまり、
 $\ln q(\tau) = \int q(\lambda_1)q(\lambda_2) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_1 d\lambda_2 - \ln Z$ のとき、上の KL 情報量は最小。

$$\begin{aligned}
\ln q(\tau) &= \int q(\lambda_1)q(\lambda_2) \ln \left\{ p(\lambda_1; a, b)p(\lambda_2; a, b)p(\tau) \prod_{n=1}^N p(c_n|\lambda_1)^{\delta(n<\tau)} p(c_n|\lambda_2)^{\delta(n\geq\tau)} \right\} d\lambda_1 d\lambda_2 - \ln Z \\
&= \int q(\lambda_1) \ln p(\lambda_1; a, b) d\lambda_1 + \int q(\lambda_2) \ln p(\lambda_2; a, b) d\lambda_2 + \ln p(\tau) \\
&\quad + \sum_{n=1}^N \delta(n < \tau) \int q(\lambda_1) \ln p(c_n|\lambda_1) d\lambda_1 + \sum_{n=1}^N \delta(n \geq \tau) \int q(\lambda_2) \ln p(c_n|\lambda_2) d\lambda_2 + \text{const.}
\end{aligned} \tag{8}$$

よって、 $q(\tau)$ はカテゴリカル分布であり、

$$q(\tau) \propto \exp \left[\sum_{n=1}^N \delta(n < \tau) \int q(\lambda_1) \ln p(c_n|\lambda_1) d\lambda_1 + \sum_{n=1}^N \delta(n \geq \tau) \int q(\lambda_2) \ln p(c_n|\lambda_2) d\lambda_2 \right]$$

となる。ただし、 $\sum_{\tau=1}^T q(\tau) = 1$ を満たす。

まとめ

- ▶ メッセージ受信数の変化点を検知するため、ベイズ的なモデルを立てた
- ▶ 事後分布を近似するために、変分ベイズ推論を行った
- ▶ その際、変分事後分布 $q(\lambda_1, \lambda_2, \tau)$ について、
 $q(\lambda_1, \lambda_2, \tau) = q(\lambda_1)q(\lambda_2)q(\tau)$ と分解できることを仮定した
- ▶ このように仮定すると、 $q(\lambda_1)$ と $q(\lambda_2)$ はガンマ分布となり、 $q(\tau)$ はカテゴリカル分布となった

課題9

- ▶ メッセージ受信数の変化点検知の例を考える。
- ▶ λ_1 の値が従う変分事後分布 $q(\lambda_1)$ は、ガンマ分布であることが分かった。
- ▶ そこで、 $q(\lambda_1)$ の shape パラメータを α_1 とし、rate パラメータを β_1 とする。
- ▶ このとき、 $\int q(\lambda_1) \ln p(\lambda_1; a, b) d\lambda_1$ を計算せよ。
 - ▶ これは ELBO の算出に必要な計算。