

10.3.3 変分下界

これまでの議論で、ベイズ線形回帰における事後分布と予測分布を変分推論で求めることができた。混合ガウスモデルと同様に、ベイズ線形回帰においても共役事前分布を設定することで、事前分布と同一の分布が更新式で得られることを確かめた。

この節では、次元数の異なる線形回帰モデルを比較するために、**変分下界**（エビデンス下界）を求めることとする。変分下界の導出を行い、これを実際のデータに当てはめてみた場合に、複雑なモデルが選択されなくなるケースを確かめる。

変分下界の導出に入る前に、変分下界を求めるモチベーションについて改めて整理したい。混合ガウス分布で変分下界を求めた節（10.2.2節）では、事後分布の収束判定のために変分下界を求めていた。今回のベイズ線形回帰でも同じく収束判定のために利用できるが、ここでは特にモデル比較のために利用する。

「10.1.4 モデル比較」では、下記の式より、モデルの事後分布をモデルの事前分布と変分下界で表すことができた。

$$q(m) \propto p(m) \exp \{ \mathcal{L}_m \}$$

$$\text{したがって } \mathcal{L}_m = \sum_{\mathbf{z}} q(\mathbf{z}|m) \ln \left\{ \frac{p(\mathbf{z}, \mathbf{y} | m)}{q(\mathbf{z} | m)} \right\}$$

~~~~~  
モデル  $m$  における変分下界

モデルの事前分布に一律な分布を仮定すると、モデルの事後分布  $q(m)$  は  $\mathcal{L}_m$  のみに比例することとなり、事後分布の大小は変分下界のみに依存する。事後分布最大のモデルを選択するには、次元数の異なる各のモデルで変分下界を求め、最大のモデルを選択すれば良い。

次に変分下界の導出に入る。変分下界はこれまで得られた分解の式を使って、次のように表現できる。

$$\begin{aligned} \mathcal{L}(q) &= \iint q(w, \alpha) \ln \left\{ \frac{p(w, \alpha, \mathbf{t})}{q(w, \alpha)} \right\} dw d\alpha \\ &= \iint q(w, \alpha) \ln p(w, \alpha, \mathbf{t}) dw d\alpha - \iint q(w, \alpha) \ln q(w, \alpha) dw d\alpha \\ &= E[\ln p(w, \alpha, \mathbf{t})] - E[\ln q(w, \alpha)] \\ &= E[\ln p(\mathbf{t} | w) p(w | \alpha) p(\alpha)] - E[\ln q(w) q(\alpha)] \\ &= E_w[\ln p(\mathbf{t} | w)] + E_w[\ln p(w | \alpha)] + E_\alpha[\ln p(\alpha)] \\ &\quad - E_w[\ln q(w)] - E_w[\ln q(\alpha)] \end{aligned}$$

ここで分解された期待値は分布の性質を使って次のように求めることができる。

$$\begin{aligned}
 E_w[\ln p(t|w)] &= E_w\left[\sum_{n=1}^N \ln \mathcal{N}(t_n | w' \Phi_n, \beta^{-1})\right] \\
 &= \frac{N}{2} \ln\left(\frac{\beta}{2\pi}\right) - \frac{\beta}{2} E_w[(t - \Phi w)'(t - \Phi w)] \\
 &= \frac{N}{2} \ln\left(\frac{\beta}{2\pi}\right) - \frac{\beta}{2} t't + \beta E_w[w' \Phi'] t - \frac{\beta}{2} E_w[w' \Phi \Phi w] \\
 &= \frac{N}{2} \ln\left(\frac{\beta}{2\pi}\right) - \frac{\beta}{2} t't + \beta m_N' \Phi' t - \frac{\beta}{2} E_w[\text{Tr}(\Phi \Phi w w')] \\
 &= \frac{N}{2} \ln\left(\frac{\beta}{2\pi}\right) - \frac{\beta}{2} t't + \beta m_N' \Phi' t - \frac{\beta}{2} \text{Tr}(\Phi \Phi (m_N m_N' + S_N))
 \end{aligned}$$

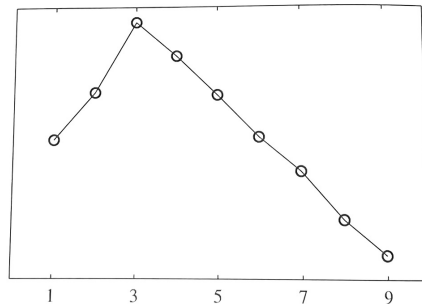
$$\begin{aligned}
 E_{w, \alpha}[\ln p(w|\alpha)] &= E_{w, \alpha}[\ln \mathcal{N}(w | 0, \alpha^{-1} \mathbf{I})] \\
 &= E_{w, \alpha}\left[-\frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln|\alpha^{-1} \mathbf{I}| - \frac{\alpha}{2} w' w\right] \\
 &= -\frac{M}{2} \ln(2\pi) - \frac{1}{2} E[M \ln(\alpha^{-1})] - \frac{1}{2} \underbrace{E[\alpha] E[w' w]}_{\text{独立だから、OK}} \\
 &= -\frac{M}{2} \ln(2\pi) + \frac{M}{2} E[\ln \alpha] - \frac{\alpha_N}{2b_N} E[\text{Tr}(w w')] \quad (\because (10.102)) \\
 &= -\frac{M}{2} \ln(2\pi) + \frac{M}{2} (\psi(\alpha_N) - \ln b_N) - \frac{\alpha_N}{2b_N} \text{Tr}[m_N m_N' + S_N] \\
 &\quad \left( \because \begin{pmatrix} (10.103) \\ (B.30) \end{pmatrix} \right) \\
 &= -\frac{M}{2} \ln(2\pi) + \frac{M}{2} (\psi(\alpha_N) - \ln b_N) - \frac{\alpha_N}{2b_N} [m_N' m_N + \text{Tr}(S_N)]
 \end{aligned}$$

$$\begin{aligned}
 E[\ln p(\alpha)] &= E[\ln \text{Gam}(\alpha | a_0, b_0)] \\
 &= E\left[\ln\left(\frac{1}{\Gamma(a_0)} b_0^{a_0} \alpha^{a_0-1} e^{-b_0 \alpha}\right)\right] \\
 &= -\ln \Gamma(a_0) + a_0 \ln b_0 + (a_0 - 1) E[\ln \alpha] - b_0 E[\alpha] \\
 &= -\ln \Gamma(a_0) + a_0 \ln b_0 + (a_0 - 1) [\psi(a_N) - \ln b_N] - b_0 \frac{\alpha_N}{b_N} \quad (10.110)
 \end{aligned}$$

$$\begin{aligned}
E_{w, \alpha} [\ln q(w)] &= E_{w, \alpha} [\ln \mathcal{N}(w | m_n, S_n)] \\
&= E_{w, \alpha} \left[ -\frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |S_n| - \frac{1}{2} (w - m_n)' S_n^{-1} (w - m_n) \right] \\
&= -\frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |S_n| - \frac{1}{2} E \left[ \text{Tr} \left\{ (w - m_n)(w - m_n)' S_n^{-1} \right\} \right] \\
&= -\frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |S_n| - \frac{1}{2} \text{Tr} \left[ E \left\{ (w - m_n)(w - m_n)' \right\} S_n^{-1} \right] \\
&= -\frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |S_n| - \frac{1}{2} \text{Tr} [S_n S_n^{-1}] \\
&= -\frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |S_n| - \frac{1}{2} \text{Tr} [I] \\
&= -\frac{1}{2} \ln |S_n| - \frac{M}{2} \ln(2\pi) - \frac{M}{2} \\
&= -\frac{1}{2} \ln |S_n| - \frac{M}{2} \{1 + \ln(2\pi)\} \quad (10.11)
\end{aligned}$$

$$\begin{aligned}
E[\ln q(\alpha)] &= E[\ln \text{Gam}(\alpha | a_n, b_n)] \\
&= E \left[ \ln \frac{1}{\Gamma(a_n)} b_n^{a_n} \alpha^{a_n-1} e^{-b_n \alpha} \right] \\
&= -\ln \Gamma(a_n) + a_n \ln b_n + (a_n - 1) E[\ln \alpha] - b_n E(\alpha) \\
&= -\ln \Gamma(a_n) + \underline{a_n \ln b_n} + \underline{(a_n - 1)} \{ \psi(a_n) - \underline{\ln b_n} \} - b_n \frac{\partial \ln \Gamma(a_n)}{\partial a_n} \\
&= -\ln \Gamma(a_n) + (a_n - 1) \psi(a_n) + \ln b_n - a_n \quad (10.12)
\end{aligned}$$

以上より、あるモデルにおける変分下界を求めることができた。下図は、人口データについて次数の異なる線形モデルの変分下界をプロットしたものである。図より、次数3において変分下界が最も大きくなり、モデルの事後分布も次数3で最大確率が割り当てられていると推察できる。これは複雑なモデルが選択されにくい性質を示している。



## 10.4 指数型分布族

混合ガウス分布の例とベイズ線形回帰の例を通じて、我々は次の2点の洞察を得た。

1. 共役事前分布を設定したパラメータは、変分推論後の近似分布も同一の分布となる
  2. 変分推論における独立性の仮定は、潜在変数とパラメータについての独立のみで良い
- 実はこれらの性質は指数型分布族についてより一般的に成り立つ。

本節では、一般的な性質として次のことを示す。「潜在変数が指数型分布で、パラメータがこれの共役事前分布を持つ場合、潜在変数とパラメータの独立性を仮定した変分推論により共役な事後分布が得られる」ここで、パラメータは**内包の変数**、潜在変数は**外延の変数**と呼ばれる。

さらに後半では、一般的なグラフィカルモデルの事後分布が局所的なメッセージパッシングを繰り返すことにより求まることを示す。これに必要な条件は、前半で示すように全ての確率変数が共役事前分布をもつ指数型分布族であることである。

まずは、一般的な指数型分布族についての変分推論により、共役な事後分布が得られることを示す。前提として、観測変数  $\mathcal{X}$  と潜在変数  $\mathcal{Z}$ 、パラメータ  $\eta$  は次のような分布をもつとする。

$$p(\mathcal{X}, \mathcal{Z} | \eta) = \prod_n h(\mathbf{x}_n, \mathbf{z}_n) g(\eta) \exp \left\{ \eta^T \mathbf{u}(\mathbf{x}_n, \mathbf{z}_n) \right\}$$

$$p(\eta | \nu_0, \chi_0) = f(\nu_0, \chi_0) g(\eta)^{\nu_0} \exp \left\{ \nu_0 \eta^T \chi_0 \right\}$$

同時分布  $p(\mathcal{X}, \mathcal{Z} | \eta)$  は指数型分布族であり、パラメータの事前分布  $p(\eta | \nu_0, \chi_0)$  は共役となるような指数型分布族として表している。(ref. 上巻 114頁 2.4.2節)

変分推論を行うため、 $q(\mathbf{z}, \eta) = q(\mathbf{z}) q(\eta)$  という独立性を仮定する。この仮定により、一般的な変分の更新式 (10.9) を適用することができる。

潜在変数の更新式は次のようになる。

$$\begin{aligned} \ln q^*(\mathbf{z}) &= E_{\eta} \left[ \ln \left\{ p(\mathcal{X}, \mathcal{Z} | \eta) p(\eta) \right\} \right] + \text{const} \\ &= E_{\eta} \left[ \ln p(\mathcal{X}, \mathcal{Z} | \eta) \right] + \text{const} \\ &= \sum_n E_{\eta} \left[ \ln h(\mathbf{x}_n, \mathbf{z}_n) + \ln g(\eta) + \eta^T \mathbf{u}(\mathbf{x}_n, \mathbf{z}_n) \right] + \text{const} \\ &= \sum_n \left\{ \ln h(\mathbf{x}_n, \mathbf{z}_n) + E[\eta^T] \mathbf{u}(\mathbf{x}_n, \mathbf{z}_n) \right\} + \text{const} \end{aligned}$$

両辺に指数をとると、指数型分布族として更新式を得る。

$$\begin{aligned} q^*(\mathbf{z}) &\propto \prod_n h(\mathbf{x}_n, \mathbf{z}_n) \exp \left\{ E[\eta^T] \mathbf{u}(\mathbf{x}_n, \mathbf{z}_n) \right\} \\ \therefore q^*(\mathbf{z}) &= \prod_n h(\mathbf{x}_n, \mathbf{z}_n) \underbrace{g(E(\eta))}_{\text{指数型分布の一般的な形から求まる}} \exp \left\{ E[\eta^T] \mathbf{u}(\mathbf{x}_n, \mathbf{z}_n) \right\} \end{aligned}$$

指数型分布の一般的な形から求まる。

次にパラメータの更新式は次のように求まる。

$$\begin{aligned}
 \ln q^*(\eta) &= E_z [\ln p(x, z | \eta) + \ln p(\eta | \nu_0, x_0)] + \text{const.} \\
 &= \sum_n^N \left\{ \underbrace{\ln E(h(x_n, z_n)) + \ln g(\eta)}_{\text{const.}} + \eta^T E(u(x_n, z_n)) \right\} \\
 &\quad + \underbrace{\ln f(\nu_0, x_0)}_{\text{const.}} + \nu_0 \ln g(\eta) + \nu_0 \eta^T x_0 + \text{const.} \\
 &= N \ln g(\eta) + \eta^T \sum_n^N E(u(x_n, z_n)) + \nu_0 \ln g(\eta) + \nu_0 \eta^T x_0 + \text{const.} \\
 &= (\nu_0 + N) \ln g(\eta) + \eta^T [\nu_0 x_0 + \sum_n^N E\{u(x_n, z_n)\}] + \text{const.} \\
 &= \nu_N \ln g(\eta) + \nu_N \eta^T x_N + \text{const.}
 \end{aligned}$$

(ここで,  $\nu_N = \nu_0 + N$   
 $\nu_N x_N = \nu_0 x_0 + \sum_n^N E\{u(x_n, z_n)\}$ )

両辺に指数をとり正規化すると、次の指数型分布族が得られる。

$$q^*(\eta) = f(\nu_N, x_N) g(\eta)^{\nu_N} \exp\{\nu_N \eta^T x_N\}$$

以上で、より一般的な仮定のもとで変分推論の更新式を得られることを確かめた。また、得られた二つの更新式は互いに依存関係があるため、これまで見てきた例と同様に繰り返し法によって分布を更新していく必要がある。

### 10.4.1 変分メッセージパッシング

10.4 指数型分布族の節では、 $z$  と  $\eta$  の2変数間のみの依存関係について考えてきた。より多くの変数の依存関係を表現する場合有向グラフで表すことができる。

この節では有向グラフで表現された確率分布について変分推論を実施できることを示す。特に、変分推論の一般的な結果 (10.9) を利用してノード  $\mathcal{X}_i$  の更新式を得る場合、ノード  $\mathcal{X}_i$  のマルコフブランケット内の期待値に関する局所的な計算のみが必要なことを示す。

まず、有向グラフにおける同時分布は次のように分解できる。

$$p(x) = \prod_i p(x_i | \text{pa}_i)$$

また、変分近似による分布は次のように分解できるとする。

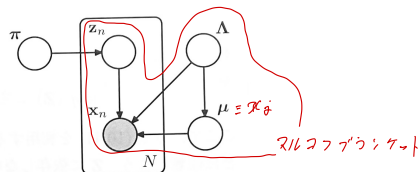
$$q(x) = \prod_i q(x_i)$$

ノード  $\mathcal{X}_i$  の変分推論による更新式は、(10.9) より次のように求まる。

$$\begin{aligned} \ln \mathcal{R}_i^*(x_i) &= E_{i \sim \tilde{\mathcal{D}}} [\ln p(x_i)] + \text{const.} \\ &= E_{i \sim \tilde{\mathcal{D}}} [\sum_j \ln p(x_i | p a_i)] + \text{const.} \\ &\quad \text{ } x_i \text{ の } p a_i \text{ への } p \text{ は } x_i \text{ による} \\ &\quad \text{項の } p \text{ だけ} \\ &= E_{i \sim \tilde{\mathcal{D}}} [\ln p(x_i | p a_i) + \sum_{i \in \mathcal{C}_i} \ln p(x_i | p a_i)] + \text{const.} \\ &\quad \text{ } i \text{ の子集合} \\ &= E_{i \in (\mathcal{MB}_{\tilde{\mathcal{D}}} / i)} [\ln p(x_i | p a_i) + \sum_{i \in \mathcal{C}_i} \ln p(x_i | p a_i)] + \text{const.} \\ &\quad \text{ } x_i \text{ の } \mathcal{MB}_{\tilde{\mathcal{D}}} \text{ からの } i \\ &\quad \text{ } x_i \text{ を除いた集合} \end{aligned}$$

この更新式から、ノード  $i$  の変分事後分布は、ノード  $i$  のマルコフブランケットに閉じた局所的な計算によって得られることがわかる。

(感想) ベイズ混合ガウス分布に当てはめて考えるとわかりやすかった。(10.55) から  $\pi$  と  $\Lambda$  ,  $\mu$  をそれぞれ別々に更新式を得ているのは、それぞれの変数がマルコフブランケットの外側にあること成り立っていると理解できる。



## 10.5 局所変分推論法

(ここでやろうとしていることをほぼ理解できなかったので、読み合わせとさせていただきます)

この節では特に、ロジスティックシグモイド関数の下界を新たなパラメータを用いて表現する。ここで求めた下界は、10.6.1 節で分布関数の近似として代用され、事後分布が求まるらしい (10.152 - 10.156 あたり)

また、今回の内容の中では、209頁最後の 10.131 の後あたりから何が言いたいかわからなくなった。図10.10 の解釈などがわからなかったので相談したい。