

変分ベイズ法とは

正田 備也

masada@rikkyo.ac.jp

Contents

変分ベイズ法とは

変分ベイズ法の実例

ベイズ的モデリングにおける変分法

- ▶ 観測データを表す確率変数を $\mathcal{X} \equiv \{x_1, \dots, x_N\}$ とする
- ▶ データモデルのパラメータを Θ とする
- ▶ ベイズ的なモデリングでは、 \mathcal{X} だけでなく Θ も確率変数
- ▶ 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
 - ▶ $q(\Theta)$ を変分法 (variational methods) で求める (後述)
 - ▶ $q(\Theta)$ を変分事後分布 (variational posterior distribution) と呼ぶ

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 4 / 25

この議論のパターンを事後分布の推論へ適用

- ▶ 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

変分ベイズ法 (variational Bayesian methods) 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
- ▶ $q(\Theta)$ は変分法 (variational method) で求められるので、変分事後分布 (variational posterior) と呼ばれる

「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、 $q(\Theta)$ の密度関数のかたち自体を変化させる
- ▶ 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
- ▶ とても直感的に言うと、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)

ELBO を最大化する根拠

- ▶ 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

変分事後分布に関する factorization の仮定

- ▶ モデルパラメータ Θ は、多数のパラメータからなっている
- ▶ Θ について、排他的なグループ分け $\Theta = \Theta_1 \cup \dots \cup \Theta_m$ を行い、これに対応して変分事後分布が $q(\Theta) = q(\Theta_1) \cdots q(\Theta_m)$ という積に分解されると仮定することが、よくある
 - ▶ 「posterior factorizes」でググってみよう
- ▶ 最も極端な場合、個々のパラメータが従う確率分布の積へ分解されると仮定することも、わりとある
 - ▶ つまり、 Θ が d 個のパラメータ $\theta_1, \dots, \theta_d$ からなるとすると、 $q(\Theta) = q(\theta_1) \cdots q(\theta_d)$ という積へ分解されると仮定する
 - ▶ これを平均場近似 (mean field approximation) と呼ぶ

より実地的な変分ベイズ法

- ▶ 上述の factorization の仮定をおくと、それだけで、変分事後分布の密度関数のかたちが決まってしまうこともある
 - ▶ この後示す例が、そうになっている
- ▶ しかし実際には、 $q(\Theta)$ の密度関数が特定のかたちを持つと仮定してしまった上で、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
 - ▶ 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

変分ベイズ法とは

変分ベイズ法の実例

変分ベイズ法によるデータ分析の手順

- ▶ データモデル $p(\mathcal{X}|\Theta)$ とモデルパラメータの事前分布 $p(\Theta; \Xi)$ を指定する
 - ▶ Ξ は事前分布のパラメータ、つまり、ハイパーパラメータ
- ▶ 同時分布 $p(\mathcal{X}, \Theta; \Xi)$ の式を書き下す
- ▶ ELBO の式を書き下す
- ▶ 変分事後分布 $q(\Theta; \Psi)$ について何らかの仮定を行う
- ▶ その仮定を利用して、変分事後分布のパラメータ Ψ を求めるための式（多くの場合、更新式）を得る
- ▶ この式を実装して計算機で動かす

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

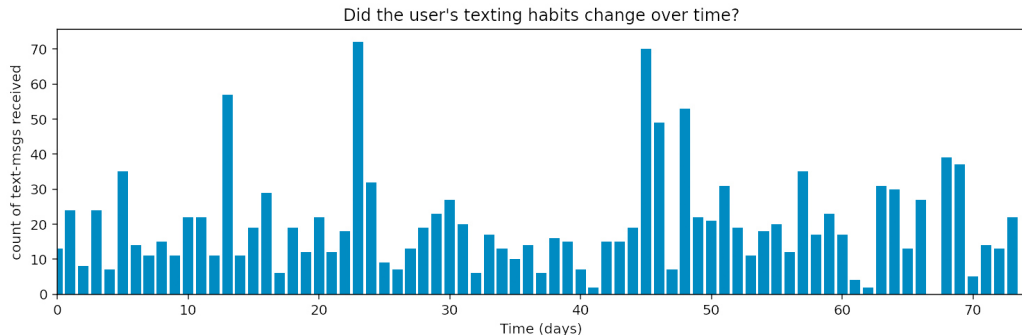


Figure: メッセージの受信数

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

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

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

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

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

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

同時分布を書き下す

同時分布は、観測データを $\mathbf{c} = \{c_1, \dots, c_N\}$ とすると

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

- ▶ $p(\lambda_i; a, b) \equiv \frac{b^a}{\Gamma(a)} \lambda_i^{a-1} e^{-b\lambda_i}$ for $i = 1, 2$
- ▶ $p(\tau) \equiv \frac{1}{N}$
- ▶ $\delta(\cdot)$ は、カッコ内の命題が真ならば 1、偽ならば 0
- ▶ $p(c_n | \lambda_i) \equiv \frac{\lambda_i^{c_n} e^{-\lambda_i}}{c_n!}$ for $i = 1, 2$

ELBO を書き下す

ELBO は

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

- ▶ このままではこれ以上議論を進められないので、変分事後分布 $q(\lambda_1, \lambda_2, \tau)$ について、それを単純化するような、何らかの仮定を行う

平均場近似の仮定

ここでは、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

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

▶ 同時分布の式 (3) を使って、ELBO をさらに詳しく書き下す

ハイパーパラメータ a, b は省略する。

$$\begin{aligned}
\ln p(\mathbf{c}) &\geq \int \sum_{\tau} q(\lambda_1)q(\lambda_2)q(\tau) \ln \frac{p(\mathbf{c}, \lambda_1, \lambda_2, \tau)}{q(\lambda_1)q(\lambda_2)q(\tau)} d\lambda_1 d\lambda_2 \\
&= \int \sum_{\tau} q(\lambda_1)q(\lambda_2)q(\tau) \ln \frac{p(\lambda_1)p(\lambda_2)p(\tau) \prod_{n=1}^N p(c_n|\lambda_1)^{\delta(n<\tau)} p(c_n|\lambda_2)^{\delta(n\geq\tau)}}{q(\lambda_1)q(\lambda_2)q(\tau)} d\lambda_1 d\lambda_2 \\
&= \int q(\lambda_1) \ln \frac{p(\lambda_1)}{q(\lambda_1)} d\lambda_1 + \int q(\lambda_2) \ln \frac{p(\lambda_2)}{q(\lambda_2)} d\lambda_2 + \sum_{\tau} q(\tau) \ln \frac{p(\tau)}{q(\tau)} \\
&\quad + \sum_{\tau} \sum_{n=1}^N \delta(n < \tau) \int q(\lambda_1) \ln p(c_n|\lambda_1) d\lambda_1 + \sum_{\tau} \sum_{n=1}^N \delta(n \geq \tau) \int q(\lambda_2) \ln p(c_n|\lambda_2) d\lambda_2 \\
&= -D_{\text{KL}}(q(\lambda_1) \parallel p(\lambda_1)) - D_{\text{KL}}(q(\lambda_2) \parallel p(\lambda_2)) - D_{\text{KL}}(q(\tau) \parallel p(\tau)) \\
&\quad + \sum_{\tau} \sum_{n=1}^N \delta(n < \tau) \int q(\lambda_1) \ln p(c_n|\lambda_1) d\lambda_1 + \sum_{\tau} \sum_{n=1}^N \delta(n \geq \tau) \int q(\lambda_2) \ln p(c_n|\lambda_2) d\lambda_2 \quad (6)
\end{aligned}$$

- ▶ 上記の ELBO の値を求めるコードも書いて、変分事後分布のパラメータ更新によって ELBO が徐々に大きくなっているかを、チェックする。

変分事後分布を求める

- ▶ この例の場合、平均場近似の仮定を置くと、変分事後分布の密度関数の式のかたちが決まってしまう
 - ▶ 結論を先取りすると、 $q(\lambda_1)$ と $q(\lambda_2)$ の密度関数の式はガンマ分布の密度関数の式のかたちに一致し、 $q(\tau)$ の質量関数の式はカテゴリカル分布の質量関数の式のかたちに一致する
- ▶ 具体的には、 $q(\lambda_1)$ と $q(\lambda_2)$ と $q(\tau)$ のうち2つを固定し、残りの1つについて、変分事後分布の事後分布に対するKL情報量 $D_{\text{KL}}(q \parallel p)$ がゼロになる条件を明らかにする
- ▶ すると、その変分事後分布の密度関数のかたちがおのずと決まってくる

$q(\lambda_1)$ の密度関数のかたちを求める

$q(\lambda_2)$ 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。

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

$q(\lambda_1) = \frac{1}{Z} \exp \int \sum_{\tau} q(\lambda_2)q(\tau) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_2$ 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。

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

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。

$$\begin{aligned} & D_{\text{KL}}(q(\lambda_1)q(\lambda_2)q(\tau) \parallel p(\lambda_1, \lambda_2, \tau|\mathbf{c})) \\ &= \int \sum_{\tau} q(\lambda_1)q(\lambda_2)q(\tau) \ln \frac{q(\lambda_1)q(\lambda_2)q(\tau)}{p(\lambda_1, \lambda_2, \tau|\mathbf{c})} d\lambda_1 d\lambda_2 \\ &= \sum_{\tau} q(\tau) \left\{ \ln q(\tau) - \int q(\lambda_1)q(\lambda_2) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_1 d\lambda_2 \right\} + \text{const.} \\ &= \sum_{\tau} q(\tau) \ln \frac{q(\tau)}{\exp \int q(\lambda_1)q(\lambda_2) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_1 d\lambda_2} + \text{const.} \\ &= D_{\text{KL}}(q(\tau) \parallel \frac{1}{Z} \exp \int q(\lambda_1)q(\lambda_2) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_1 d\lambda_2) + \text{const.} \end{aligned} \tag{8}$$

$q(\tau) = \frac{1}{Z} \exp \int q(\lambda_1)q(\lambda_2) \ln p(\lambda_1, \lambda_2, \tau|\mathbf{c}) d\lambda_1 d\lambda_2$ 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。

$$\begin{aligned}
\ln q(\tau) &= \int q(\lambda_1)q(\lambda_2) \ln \left\{ p(\lambda_1; a, b)p(\lambda_2; a, b)p(\tau) \prod_{n=1}^N p(c_n|\lambda_1)^{\delta(n<\tau)}p(c_n|\lambda_2)^{\delta(n\geq\tau)} \right\} d\lambda_1 d\lambda_2 - \ln Z \\
&= \int q(\lambda_1) \ln p(\lambda_1; a, b) d\lambda_1 + \int q(\lambda_2) \ln p(\lambda_2; a, b) d\lambda_2 + \ln p(\tau) \\
&\quad + \sum_{n=1}^N \delta(n < \tau) \int q(\lambda_1) \ln p(c_n|\lambda_1) d\lambda_1 + \sum_{n=1}^N \delta(n \geq \tau) \int q(\lambda_2) \ln p(c_n|\lambda_2) d\lambda_2 + \text{const.}
\end{aligned} \tag{9}$$

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、 $\sum_{\tau=1}^N q(\tau) = 1$ 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。