

Ch.4 近似 Bayes 推論

• \mathcal{X} : 観測データ.

\mathcal{Q} : 非観測変数 (パラメータ, 潜在変数など)

< Bayes 推論の流れ >

1 確率モデル $p(\mathcal{X}, \mathcal{Q})$ の設計.

2. 事後分布 $p(\mathcal{Q}|\mathcal{X})$ の計算, パラメータ更新 (学習)

3. 予測分布の計算 (予測)

• Ch.3 で扱ったモデルは単純だったため, 事後分布 $p(\mathcal{Q}|\mathcal{X})$ が厳密に得られた.

→ 複雑なモデルでは, 厳密に $p(\mathcal{Q}|\mathcal{X})$ を得るのは困難!

→ 近似手法 を用いて何とかしたい.

§4.1 サンプルリングに基づく推論手法.

• $p(\mathcal{Q}|\mathcal{X})$ を求める代わりに, $p(\mathcal{Q}|\mathcal{X})$ から複数のサンプルを得ることで事後分布の特性を調べることにする.

4.1.1 単純モンテカルロ法.

< 問題設定 >

$$\mathbb{E}_{p(\mathcal{Z})}[f(\mathcal{Z})] = \int f(\mathcal{Z})p(\mathcal{Z})d\mathcal{Z} \text{ を求める.}$$

[仮定] • 積分計算は困難.

サンプルリングが容易なことはない.

• 分布 $p(\mathcal{Z})$ からのサンプルリングが \checkmark できる.

Th. (大数の強法則) 独立に同分布に従う

$\{X_i\}$: d -次元 i.i.d. 確率変数列. $E[X_i] = \mu$.

$\Rightarrow P\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mu\right) = 1$. ← 確率1で標本平均は期待値に収束. □

• この定理を利用して積分を計算するのは

単純モンテカルロ法 (simple Monte Carlo method):

$\mathcal{X}^{(1)}, \dots, \mathcal{X}^{(T)} \stackrel{i.i.d.}{\sim} p(\mathcal{X})$ と T 個サンプリングして,

$\int f(\mathcal{X}) p(\mathcal{X}) d\mathcal{X} \approx \frac{1}{T} \sum_{t=1}^T f(\mathcal{X}^{(t)})$ と近似する.

→ 大数の法則により 標本平均 $\frac{1}{T} \sum_{t=1}^T f(\mathcal{X}^{(t)})$ は ← 確率1で 期待値 $E_{p(\mathcal{X})}[f(\mathcal{X})] = \int f(\mathcal{X}) p(\mathcal{X}) d\mathcal{X}$ に収束するので,

T を十分大きくとるようにする.

ex) N 個の観測データ $\mathcal{X} = \{\mathcal{X}_1, \dots, \mathcal{X}_N\}$: given.

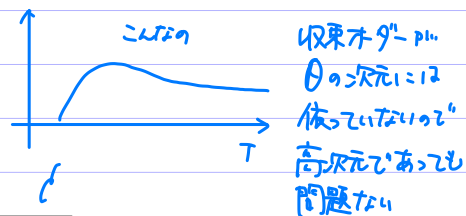
パラメータ θ をもつモデル $p(\mathcal{X}, \theta) = p(\mathcal{X}|\theta)p(\theta)$.

周辺尤度 $p(\mathcal{X}) = E_{p(\theta)}[p(\mathcal{X}|\theta)] = \int p(\mathcal{X}|\theta)p(\theta) d\theta$ を評価.

$p(\mathcal{X}) = \int p(\mathcal{X}|\theta)p(\theta) d\theta \approx \frac{1}{T} \sum_{t=1}^T p(\mathcal{X}|\theta^{(t)})$. ($\theta^{(t)} \stackrel{i.i.d.}{\sim} p(\theta)$)

と近似すればよい.

→ たゞし、この近似は収束が遅い.



(重複対数の法則より、収束のオーダーは $O\left(\sqrt{\frac{\log \log T}{T}}\right)$.)

$X_i \stackrel{i.i.d.}{\sim} p, E[X_i] = \mu, V[X_i] = \sigma^2 < \infty, S_n = \sum_{i=1}^n X_i \Rightarrow P\left(\limsup_{n \rightarrow \infty} \frac{S_n - n\mu}{\sqrt{2n \log \log n}} = \sigma\right) = 1, P\left(\liminf_{n \rightarrow \infty} \frac{S_n - n\mu}{\sqrt{2n \log \log n}} = -\sigma\right) = 1$.

4.1.2 棄却サンプリング

棄却サンプリング (rejection sampling):

単純モンテカルロ法の
サンプリングに利用できる。

密度の計算が困難な目標分布 $p(x)$ からのサンプルを得る方法

次の定理に基づく。

Th. (棄却サンプリング)

$p(x), q(x) : \text{pdf}$. $kq(x)$ を $p(x)$ の包絡函数 といふ

$k > 0$: 定数と $kq(x) \geq p(x)$ ($\forall x$) とするよう定める。

$x \sim q(x)$ と $u \sim \text{Uni}(0,1)$ は独立とする。

$\Rightarrow u \leq \frac{p(x)}{kq(x)}$ のもとでの x の条件付き pdf は $p(x)$ となる。 \square

pf. Bayesの定理より。

$$q(x | u \leq \frac{p(x)}{kq(x)}) = \frac{P(u \leq \frac{p(x)}{kq(x)} | x) q(x)}{P(u \leq \frac{p(x)}{kq(x)})}$$

ここで。

$$P(u \leq \frac{p(x)}{kq(x)} | x) = \int_0^{\frac{p(x)}{kq(x)}} 1 \, du = \frac{p(x)}{kq(x)}.$$

$$\begin{aligned} P(u \leq \frac{p(x)}{kq(x)}) &= \int P(u \leq \frac{p(x)}{kq(x)} | x) q(x) \, dx = \int \frac{p(x)}{kq(x)} q(x) \, dx \\ &= \frac{1}{k} \int p(x) \, dx = \frac{1}{k}. \end{aligned}$$

$$\therefore q(x | u \leq \frac{p(x)}{kq(x)}) = \frac{\frac{p(x)}{k}}{1/k} = p(x). \quad \blacksquare$$

\rightarrow 提案分布 (proposal distribution) $q(x)$ をサンプリングしやすい分布に

しておくことで、 $p(x)$ のサンプルを次のように得ることになる。

Algo. (棄却サンプリング)

Input: 目標分布 $p(x)$, 提案分布 $q(x)$, 定数 $k: kq(x) \geq p(x)$.

$T \in \mathbb{Z}_{>0}$

Output: $\{x^{(1)}, \dots, x^{(T)}\}$: $p(x)$ からのサンプル.

1 for $t = 1, 2, \dots, T$ do:

2 repeat

3 $x \sim q(x)$; $u \sim \text{Uni}(0,1)$

4 until $u \leq \frac{p(x)}{kq(x)}$; この条件が満たされるまで,
← サンプル x を棄却しリサンプリング.

5 $x^{(t)} \leftarrow x$ ← サンプル x を受容.

6 end for;

7 return $\{x^{(1)}, \dots, x^{(T)}\}$. □

- 目標分布 $p(x)$ の正規化定数 k が分からなくてもサンプリングできる.

$$p(x) = \frac{1}{Z_p} \tilde{p}(x), \quad Z_p = \int \tilde{p}(x) dx$$

← 計算が難しいことがある

のとき, $kq(x) \geq \frac{1}{Z_p} \tilde{p}(x) \Leftrightarrow kZ_p q(x) \geq \tilde{p}(x)$ となる. 改めて

$k \leftarrow kZ_p$ とすれば, $p(x)$ の代わりに $\tilde{p}(x)$ を使っても $p(x)$ からの

サンプリングが可能.

- 棄却サンプリングの効率: 2~4行目の繰り返し回数による.

受容確率は $P(u \leq \frac{p(x)}{kq(x)}) = \frac{1}{k}$. 各回の試行は独立.

→ 受容されるまでの繰り返り回数 は 幾何分布 $Ge(\frac{1}{k})$ に従う.

$$(Ge(p) \text{ の pmf: } p(n) = p(1-p)^{n-1}.)$$

繰り返り回数の期待値は k .

$p(x)$ を近似できるようなもの

• k はできるだけ小さい方がよい. そうなるように, $q(x)$ をうまく選ぶとよいから

そういった $q(x)$ を見つけるのは容易ではない.

• x の次元が大きくなると, k が大きくなる傾向にあり 効率が悪くなる.

(次元の呪い)

cf.) 効率を良くするために, $p(x)$ を区分的に近似する $q(x)$ を使った

適応的棄却サンプリング" といったものもある.

4.1.3 自己正規化重点サンプリング

• 期待値 $E_{p(x)}[f(x)] = \int f(x)p(x)dx$ を計算するのに,

$$\text{単純MC では } \int f(x)p(x)dx \approx \frac{1}{T} \sum_{t=1}^T f(x^{(t)}), \quad x^{(t)} \sim p(x)$$

と $p(x)$ からのサンプリング" ができる必要があった.

• $p(x)$ からのサンプリング" が容易でないときは, 上述の棄却サンプリング" を

利用できた. サンプリングされた $x \sim q(x)$ に対して

• $p(x)$ が小さいとき x は棄却されやすいので, $p(x)$ の大きいサンプルが得られやすい.

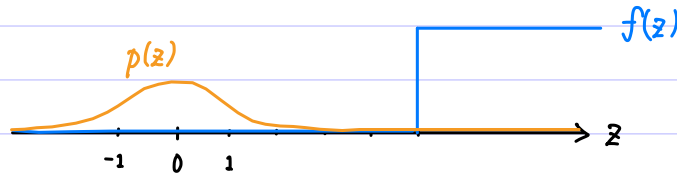
• $|f(x)|$ が小さいとき 和 $\sum_{t=1}^T f(x^{(t)})$ の寄与が小さいため,

$E_{p(x)}[f(x)]$ の収束に T をかなり大きくする必要もある.

ex) 極端 T="1000, $p(z) = \mathcal{N}(z | 0, 1)$,

$$f(z) = I(\{z \geq 5\}) = \begin{cases} 1 & (z \geq 5) \\ 0 & (z < 5) \end{cases}$$

とすると下図のような状況になる.



積分 $\int f(z)p(z)dz \neq 0$ T="1000, ↖ 2.8665×10^{-7} c"34

$z \geq 5$ のサンプルを得る確率 p (ほぼ 0 の T にあ, $z^{(t)} \sim p(z)$ としても

$$\frac{1}{T} \sum_{t=1}^T f(z^{(t)}) \text{ は } \text{ほぼ } T=1) 0 \text{ になってしまふこと } p \text{ (ほとんど).}$$

(近似はできていない p ...)

↖ $p(z)$ が大きくなる z は 高次元空間のごく一部.
そこで $|f(z)|$ が小さいと ほど 計算効率が悪くなる.

以上のような問題は, 高次元だとさらに悪化する (次元の呪い).

→ $|f(z)|p(z)$ の大きい領域を重点的にサンプリングすることで,

計算効率を上げる (T をなるべく小さくする) ことができれば良い.

重点サンプリング (importance sampling).

サンプリングが容易にできる提案分布 $q(z)$ を考える.

$f(z)p(z) \neq 0 \Rightarrow q(z) > 0$ と仮定. 計算したい期待値は,

$$\mathbb{E}_{p(z)}[f(z)]$$

↖ $\text{supp}(q) \supseteq \text{supp}(f \cdot p)$.
($\text{supp}(f) := \{z | f(z) \neq 0\}$: f の台 (support))

$$= \int f(z)p(z) dz \quad \leftarrow \text{積分は } f(z)q(z) \neq 0 \text{ のところを考慮すれば + 分.}$$

$$= \int f(z) \frac{p(z)}{q(z)} q(z) dz. \quad \leftarrow \text{このとき 仮定より } q(z) > 0.$$

∴ $w(\mathbf{z}) := \frac{p(\mathbf{z})}{q(\mathbf{z})}$ (重要度重み, importance weight) とおくと,

$$\mathbb{E}_{p(\mathbf{z})}[f(\mathbf{z})] = \int w(\mathbf{z}) f(\mathbf{z}) q(\mathbf{z}) d\mathbf{z} = \mathbb{E}_{q(\mathbf{z})}[w(\mathbf{z}) f(\mathbf{z})].$$

よって, 単純MCのように

$$\mathbb{E}_{p(\mathbf{z})}[f(\mathbf{z})] \approx \frac{1}{T} \sum_{t=1}^T w(\mathbf{z}^{(t)}) f(\mathbf{z}^{(t)}) \quad (\mathbf{z}^{(t)} \stackrel{i.i.d.}{\sim} q(\mathbf{z}))$$

と近似することと $T \rightarrow \infty$ のとき右辺は $\mathbb{E}_{q(\mathbf{z})}[w(\mathbf{z}) f(\mathbf{z})] = \mathbb{E}_{p(\mathbf{z})}[f(\mathbf{z})]$ に収束する.
← almost surely.

・ 近似値の分散について考えてみる.

$$I := \mathbb{E}_{p(\mathbf{z})}[f(\mathbf{z})], \quad \hat{I} := \frac{1}{T} \sum_{t=1}^T w(\mathbf{z}^{(t)}) f(\mathbf{z}^{(t)}) \quad \text{とおく}$$

$$\mathbb{E}_{q(\mathbf{z})}[\hat{I}] = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{q(\mathbf{z})}[w(\mathbf{z}) f(\mathbf{z})] = I \quad \text{なので.}$$

$$V_{q(\mathbf{z})}[\hat{I}]$$

$$= \mathbb{E}_{q(\mathbf{z})}[\hat{I}^2] - I^2$$

$$= \frac{1}{T^2} \mathbb{E}_{q(\mathbf{z})} \left[\sum_{t=1}^T w(\mathbf{z}^{(t)})^2 f(\mathbf{z}^{(t)})^2 + 2 \sum_{s < t} w(\mathbf{z}^{(s)}) f(\mathbf{z}^{(s)}) w(\mathbf{z}^{(t)}) f(\mathbf{z}^{(t)}) \right] - I^2$$

$$= \frac{1}{T^2} \sum_{t=1}^T \mathbb{E}_{q(\mathbf{z})}[w(\mathbf{z})^2 f(\mathbf{z})^2] + \frac{2}{T^2} \sum_{s < t} \mathbb{E}_{q(\mathbf{z})}[w(\mathbf{z}) f(\mathbf{z})]^2 - I^2$$

$$= \frac{1}{T} \mathbb{E}_{q(\mathbf{z})}[w(\mathbf{z})^2 f(\mathbf{z})^2] + \frac{2}{T^2} \cdot \frac{T(T-1)}{2} I^2 - I^2$$

$$= \frac{1}{T} \mathbb{E}_{q(\mathbf{z})}[w(\mathbf{z})^2 f(\mathbf{z})^2] - \frac{1}{T} I^2$$

$$= \frac{1}{T} \left(\mathbb{E}_{q(\mathbf{z})}[w(\mathbf{z})^2 f(\mathbf{z})^2] - I^2 \right). \quad \left(= \frac{1}{T} V_{q(\mathbf{z})}[w(\mathbf{z}) f(\mathbf{z})] \right)$$

・ $\mathbb{E}_{q(\mathbf{z})}[w(\mathbf{z})^2 f(\mathbf{z})^2] < \infty$ とならば, 近似値 \hat{I} の分散も ∞ にはならず.

$w(\mathbf{z})$ の分母は $q(\mathbf{z})$ なので, $q(\mathbf{z})$ が $p(\mathbf{z})$ より 0 に近づくほど $w(\mathbf{z})$ が

大きくなり分散が発散するものがあある。

→ $q(\mathbb{R})$ は $p(\mathbb{R})$ より "裾の厚い" ものを "選ぶ" とよい。

• $\mathbb{E}_{q(\mathbb{R})}[w(\mathbb{R})^2 f(\mathbb{R})^2] < \infty$ のとき, Jensen の不等式より,

$$\mathbb{E}_{q(\mathbb{R})}[w(\mathbb{R})^2 f(\mathbb{R})^2] \\ \geq \left(\mathbb{E}_{q(\mathbb{R})}[w(\mathbb{R}) |f(\mathbb{R})|] \right)^2.$$

$I \subseteq \mathbb{R}$: 区間, $h: I \rightarrow \mathbb{R}$: 凸函数.

X : I -値確率変数, $\mathbb{E}[X] < \infty$.

$$\Rightarrow \mathbb{E}[h(X)] \geq h(\mathbb{E}[X]).$$

等号成立は $X = \mathbb{E}[X]$, a.s. のとき

等号成立は, $|f(\mathbb{R})| w(\mathbb{R}) = \mathbb{E}_{q(\mathbb{R})}[|f(\mathbb{R})| w(\mathbb{R})]$, a.s.

$$\Leftrightarrow q(\mathbb{R}) \propto |f(\mathbb{R})| w(\mathbb{R}), \text{ a.s. のとき.}$$

→ q で $|f(\mathbb{R})| w(\mathbb{R})$ の大きい範囲を重点的にサンプリングすると良い

ことが正確にめされた

• 自己正規化重点サンプリング (self-normalized importance sampling).

上で p, q と pdf とした。こゝでは,

$$p(\mathbb{R}) = \frac{1}{Z_p} \tilde{p}(\mathbb{R}), \quad q(\mathbb{R}) = \frac{1}{Z_q} \tilde{q}(\mathbb{R})$$

として, $\tilde{p}(\mathbb{R}), \tilde{q}(\mathbb{R})$ のみで計算できる状況を考える。

(もちろん $q(\mathbb{R})$ からのサンプリングはできるものとする)

$$\mathbb{E}_{p(\mathbb{R})}[f(\mathbb{R})] \\ = \int f(\mathbb{R}) \frac{p(\mathbb{R})}{q(\mathbb{R})} q(\mathbb{R}) d\mathbb{R}. \\ = \frac{Z_q}{Z_p} \int f(\mathbb{R}) \frac{\tilde{p}(\mathbb{R})}{\tilde{q}(\mathbb{R})} q(\mathbb{R}) d\mathbb{R}.$$

こゝで, $\tilde{w}(\mathbb{R}) := \frac{\tilde{p}(\mathbb{R})}{\tilde{q}(\mathbb{R})}$ とすると,

$$\begin{aligned} \frac{z_p}{z_q} &= \frac{1}{z_q} \int \tilde{p}(z) dz = \frac{1}{z_q} \int \frac{\tilde{p}(z)}{\tilde{q}(z)} \tilde{q}(z) dz \\ &= \int \tilde{w}(z) q(z) dz = \mathbb{E}_{q(z)}[\tilde{w}(z)] \end{aligned}$$

なので、

$$\mathbb{E}_{p(z)}[f(z)] = \frac{\mathbb{E}_{q(z)}[\tilde{w}(z)f(z)]}{\mathbb{E}_{q(z)}[\tilde{w}(z)]}$$

よって、 T 個のサンプル $z^{(1)}, \dots, z^{(T)} \stackrel{i.i.d.}{\sim} q(z)$ を用いて

$$\mathbb{E}_{q(z)}[\tilde{w}(z)f(z)] \approx \frac{1}{T} \sum_{t=1}^T \tilde{w}(z^{(t)})f(z^{(t)}),$$

$$\mathbb{E}_{q(z)}[\tilde{w}(z)] \approx \frac{1}{T} \sum_{t=1}^T \tilde{w}(z^{(t)})$$

と近似することはできる, i.e.,

$$\mathbb{E}_{p(z)}[f(z)] \approx \frac{\frac{1}{T} \sum_{t=1}^T \tilde{w}(z^{(t)})f(z^{(t)})}{\frac{1}{T} \sum_{t=1}^T \tilde{w}(z^{(t)})} = \sum_{t=1}^T \frac{\tilde{w}(z^{(t)})}{\sum_{t=1}^T \tilde{w}(z^{(t)})} f(z^{(t)})$$

と近似する.

$\therefore w^*(z^{(t)})$ とおく

$$\cdot \sum_{t=1}^T w^*(z^{(t)}) = 1 \text{ なので "自己正規化"}$$

Algo. (自己正規化重点サンプリング)

Input: 目標分布 p の主要項 $\tilde{p}(z)$, 提案分布 $q(z) = \frac{1}{z_q} \tilde{q}(z)$,

関数 $f(z)$, $T \in \mathbb{Z}_{>0}$.

Output: 期待値 $I = \mathbb{E}_{p(z)}[f(z)]$ の近似値 \hat{I} .

1. for $t = 1, \dots, T$ do:

2. $z^{(t)} \sim q(z)$; $\tilde{w}(z^{(t)}) \leftarrow \frac{\tilde{p}(z^{(t)})}{\tilde{q}(z^{(t)})}$

3. end for;

$$4. \mathcal{J} \leftarrow \sum_{t=1}^T \tilde{w}(z^{(t)});$$

5. for $t = 1, \dots, T$ do :

$$6. w^*(z^{(t)}) \leftarrow \frac{1}{\mathcal{J}} \tilde{w}(z^{(t)})$$

7. end for ;

8. return $\sum_{t=1}^T w^*(z^{(t)}) f(z^{(t)})$ as \hat{I} .

実装の際は $\tilde{w}(z^{(t)})$ や $f(z^{(t)})$ の値を保存していなくても、
 $\tilde{w}(z^{(t)})$ や $f(z^{(t)}) \tilde{w}(z^{(t)})$ をサンプル生成の都度足しこんでいけばよい。

Remark (1) $\mathbb{E}_{q(z)} \left[\sum_{t=1}^T w^*(z^{(t)}) f(z^{(t)}) \right] \neq I$ であり、 $\sum_{t=1}^T w^*(z^{(t)}) f(z^{(t)})$ は I の

不偏推定量ではない。しかし、 $\mathbb{E}_{q(z)} \left[\sum_{t=1}^T w^*(z^{(t)}) f(z^{(t)}) \right] \xrightarrow{T \rightarrow \infty} I$, a.s.

であり、漸近不偏推定量になっている

(2) \hat{I} の分散が最小になるのは $q(z) \propto |f(z) - I| p(z)$ のとき

(Hesterberg, 1988)

(3) 棄却サンプリングは重点サンプリングの特別な場合とみなせる。

cf.) 古澄『ハイス計算統計学』

(4) 重点サンプリングでは、 $z^{(t)} \stackrel{i.i.d.}{\sim} q(z)$ としていたのを、 $p(z)$ からのサンプルは

得られないから、サンプリング-重点リサンプリング (sampling/importance resampling: SIR)

と利用すると近似的に $p(z)$ に従ったサンプルを得ることはできる。

cf.) Bishop: PRML, 古澄『ハイス計算統計学』など。

4.1.4. Markov連鎖モンテカルロ法

• Markov連鎖モンテカルロ法 (Markov Chain Monte Carlo: MCMC):

高次元の分布 $p(\mathbf{z})$ から効率的にサンプリングを行うための方法.

Def. (Markov連鎖)

$\{\mathbf{z}^{(i)}\}_{i=1}^{\infty}$: 確率変数列.

→ 前の"状態"にしか影響を及ぼさない.

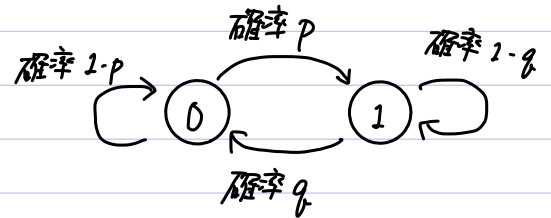
$$p(\mathbf{z}^{(t)} | \mathbf{z}^{(t-1)}, \dots, \mathbf{z}^{(1)}) = p(\mathbf{z}^{(t)} | \mathbf{z}^{(t-1)})$$

が成り立つとき, $\{\mathbf{z}^{(i)}\}_{i=1}^{\infty}$ を (1次の) Markov連鎖 という.

• $\mathcal{J}(\xi, \mathbf{z}) := p(\mathbf{z} | \xi)$: 遷移確率密度 という.

ex) (有限状態の Markov連鎖)

状態 0, 1 があって,



右図のように状態が遷移していく.

$$\mathbf{z}^{(i)} = 0 \text{ のときは } p(\mathbf{z}^{(i+1)} | \mathbf{z}^{(i)}) = \text{Bern}(\mathbf{z}^{(i+1)} | p).$$

$$\mathbf{z}^{(i)} = 1 \text{ のときは } p(\mathbf{z}^{(i+1)} | \mathbf{z}^{(i)}) = \text{Bern}(\mathbf{z}^{(i+1)} | 1-q).$$

$\mathbf{z}^{(i+1)}$ は $\mathbf{z}^{(i)}$ のみに依存. \Rightarrow 1次の Markov連鎖.

$$\mathcal{J}(\mathbf{z}^{(i)}, \mathbf{z}^{(i+1)}) = \text{Bern}(\mathbf{z}^{(i+1)} | p^{1-\mathbf{z}^{(i)}} (1-q)^{\mathbf{z}^{(i)}}) \text{ と表せる.}$$

このように行列の形で

$$\begin{pmatrix} 1-p & p \\ q & 1-q \end{pmatrix} \text{ と表す}$$

遷移確率行列

ex) $\mathbf{z}^{(i+1)} = \mathbf{z}^{(i)} + \mathbf{E}_i$, $\mathbf{E}_i \sim \mathcal{N}(0, \mathbf{I}_2)$

$\{\mathbf{z}^{(i)}\}_{i=1}^{\infty}$ は Markov連鎖.

$$\mathcal{J}(\mathbf{z}^{(i)}, \mathbf{z}^{(i+1)}) = \mathcal{N}(\mathbf{z}^{(i+1)} | \mathbf{z}^{(i)}, \mathbf{I}_2).$$



< やりたいこと >

$\mathcal{J}(\xi, \mathbb{X})$ をうまく設計して、サンプリングしたい分布 $p(\mathbb{X})$ に従うような

Markov連鎖 $\{\mathbb{X}^{(i)}\}_{i=1}^{\infty}$, $\mathbb{X}^{(i)} \sim p(\mathbb{X})$ を得たい.

$\mathbb{X}^{(i)}$ たちは同分布に従うが、独立ではないことに注意.

- $\mathbb{X}^{(i)} \sim p(\mathbb{X}) \Rightarrow \mathbb{X}^{(m)} \sim p(\mathbb{X})$ とおけるような分布を“定常分布”という:

Def. (定常分布)

確率分布 p に対し, 全ての \mathbb{X} に対して

$$p(\mathbb{X}) = \int p(\xi) \mathcal{J}(\xi, \mathbb{X}) d\xi$$

が成り立つとき, p を **定常分布** (stationary distribution) という.

- 定常分布はいつでも存在するとは限らないし, 存在しても一意とは限らない.

定常分布が存在するための十分条件として, 次のとおり.

Th. (詳細釣り合い条件)

分布 p と遷移確率密度 \mathcal{J} に対し, 全ての \mathbb{X}, ξ に対して

$$p(\mathbb{X}) \mathcal{J}(\mathbb{X}, \xi) = p(\xi) \mathcal{J}(\xi, \mathbb{X})$$

マスター方程式から出てくる式.
← \mathbb{X} から ξ に移る“確率”と
 ξ から \mathbb{X} に移る“確率”が等しいということ.

が満たされると, p は Markov連鎖の定常分布.

pf. 両辺を ξ について積分.

$$\int p(\mathbb{X}) \mathcal{J}(\mathbb{X}, \xi) d\xi = p(\mathbb{X}) \underbrace{\int \mathcal{J}(\mathbb{X}, \xi) d\xi}_{=1} = p(\mathbb{X}).$$

$$\therefore p(\mathbb{X}) = \int p(\xi) \mathcal{J}(\xi, \mathbb{X}) d\xi \quad \square$$

- MCMCでは \mathcal{J} が詳細釣り合い条件を満たすようになっている.

これかできているか？ 苦労したかい ...

・ p が定常になる \mathcal{X} が設計できたとしても、初期値 $x^{(1)}$ が $p(x)$ に従うように

サンプリング できなければ、 p に従う変数からなる Markov 連鎖 $\{x^{(i)}\}_{i=1}^{\infty}$

は得られない。

どの x の後かはものによる

遷移をくり返した後の“極限分布”が
定常分布と一致する、ということ。

↓
どんな $x^{(1)}$ から始めても、十分ステップ後の変数 $x^{(T)}$ が $p(x)$ に従うように

なっていくければよい。 $(x^{(T)}$ を改めて $x^{(1)}$ とおけば p に従う変数からなる

Markov 連鎖 $\{x^{(i)}\}_{i=1}^{\infty}$ が得られることになる (burn-in))

↑ $x^{(i)}$ は独立な変数をサンプリングしたければ、
これらから間引きすればよい。

→ エルゴード性 (ergodicity)

{	既約	: どの変数 x も、任意の x' へ有限回で遷移可能。	}	⇒ 定常分布の存在と一意性
	正再帰的	: どの変数 x も、有限回の遷移で戻ってくる。		
	非同期的	: どの変数 x も、周期的に現れない		

な Markov 連鎖のこととエルゴード的という。

Fact. (1) エルゴード的な Markov 連鎖 $\{x^{(i)}\}_{i=1}^{\infty}$ は 唯一 の定常分布 p をもつ。

(2) 任意の初期値 $x^{(1)}$ に対し、 $x^{(i)}$ の従う分布は

$i \rightarrow \infty$ の極限で p に法則収束する。

(3) $\mathbb{E}_{p(x)}[f(x)] < \infty$ のとき、

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n f(x^{(i)}) = \mathbb{E}_{p(x)}[f(x)] \leftarrow \text{この計算をする分には burn-in とおいた方がいい}$$

が成立する (Markov 連鎖に対するエルゴード定理)

→ MCMC では基本的にエルゴード的な Markov 連鎖にしておくべき!

4.1.5 Metropolis-Hastings法.

• \mathcal{J} をうまく設計するために, 次のように仮定してみる:

$$\mathcal{J}(\xi, \mathbb{R}) = A(\xi, \mathbb{R}) q(\mathbb{R} | \xi)$$

• $q(\mathbb{R} | \xi)$: 遷移確率密度, 提案分布.

• $A(\xi, \mathbb{R}) \in [0, 1]$: 採択率

→ ξ から \mathbb{R} へと仮に遷移したうす, 確率 $A(\xi, \mathbb{R})$ で \mathbb{R} を採択し,

確率 $1 - A(\xi, \mathbb{R})$ で \mathbb{R} を棄却して ξ に戻る. ということ.

このとき, 詳細釣り合い式は, $\forall \mathbb{R}, \xi$ で

$$p(\mathbb{R}) A(\mathbb{R}, \xi) q(\xi | \mathbb{R}) = p(\xi) A(\xi, \mathbb{R}) q(\mathbb{R} | \xi)$$

$$\Leftrightarrow \frac{A(\xi, \mathbb{R})}{A(\mathbb{R}, \xi)} = \frac{p(\mathbb{R}) q(\xi | \mathbb{R})}{p(\xi) q(\mathbb{R} | \xi)} =: \rho(\xi, \mathbb{R}).$$

Lem. 任意の $p(\mathbb{R}), q(\mathbb{R} | \xi)$ に対し, 次の $A(\xi, \mathbb{R})$ は 詳細釣り合い式を満たす:

$$A(\xi, \mathbb{R}) = \min(1, \rho(\xi, \mathbb{R}))$$

pf. $\rho(\mathbb{R}, \xi) = \rho(\xi, \mathbb{R})^{-1}$ に注意する.

$$\frac{A(\xi, \mathbb{R})}{A(\mathbb{R}, \xi)} = \frac{\min(1, \rho(\xi, \mathbb{R}))}{\min(1, \rho(\xi, \mathbb{R})^{-1})} = \frac{\min(\rho(\xi, \mathbb{R}), \rho(\xi, \mathbb{R})^2)}{\min(\rho(\xi, \mathbb{R}), 1)}$$

$$= \begin{cases} \frac{\rho(\xi, \mathbb{R})}{1} & (\rho(\xi, \mathbb{R}) \geq 1) \\ \frac{1}{\frac{\rho(\xi, \mathbb{R})^2}{\rho(\xi, \mathbb{R})}} & (\rho(\xi, \mathbb{R}) < 1) \end{cases}$$

$$= \rho(\xi, \mathbb{R})$$



→ $q(z|\xi)$ をうまく選ぶならば, $T(\xi, z) = q(z|\xi) \min(1, \rho(\xi, z))$

で Markov 連鎖のエルゴディク性が保証される.

$q(z|\xi) = \mathcal{N}(z|\xi, I_d)$ などとしておけば OK.

cf.) Robert and Casella "Monte Carlo Statistical Methods".

Algo. (Metropolis-Hastings)

Input: 目標分布 $p(z)$, 提案分布 $q(z|\xi)$. 繰り返し回数 $N \in \mathbb{Z}_+$.

Output: Markov 連鎖 $\{z^{(i)}\}_{i=1}^N$

1. $z^{(1)}$ をランダムに初期化;
2. for $i = 1, \dots, N-1$ do:
3. $z^* \sim q(z|z^{(i)}); u \sim \text{Unif}(0, 1);$
4. $r \leftarrow \rho(z^{(i)}, z^*);$
5. if $u \leq \min(1, r):$
6. $z^{(i+1)} \leftarrow z^* \quad \leftarrow$ 採択
7. else:
8. $z^{(i+1)} \leftarrow z^{(i)} \quad \leftarrow$ 棄却
9. return $\{z^{(i)}\}_{i=1}^N$.

Remark p の代わりに, 正規化されていない \tilde{p} を用いても良い ($p(\xi, z)$ は等しい)

$q(z|\xi) = q(\xi|z)$ (対称) のとき **Metropolis 法** という.

4.1.6 ハミルトニアン・モンテカルロ法

• $q \in \text{Gauss}$ 分布とした Metropolis-Hastings 法:

◦ 分散を大きくする

✓ $x^{(i)}$ と $x^{(i+1)}$ の相関 低.

1回の遷移での平均移動幅が大きくなるから

✗ 目標分布から外れて棄却率大.

◦ 分散を小さくする

✓ 採択率大

✗ $x^{(i)}$ と $x^{(i+1)}$ の相関 高 \Rightarrow 目標分布からのサンプルを得るまで時間がかかる.

\rightarrow 高次元だと分散の調整は困難. 特に採択率の低下が著しくなる.

高次元空間で平均移動距離を長く取ると、採択率を高くすることは

できるのか...?

• ハミルトニアン・モンテカルロ法 / ハイブリッド・モンテカルロ法

(Hamiltonian / Hybrid Monte Carlo: HMC):

解析力学的な物体の軌道のシミュレーション + Metropolis-Hastings.

\rightarrow 効率のよいサンプリングが可能.

4.1.6.1 Hamiltonian のシミュレーション.

• 少しだけ物理の復習.

- 質量 $m=1$ の物体と考える.

$\mathbf{x} = \mathbf{x}(\tau)$: 時刻 τ での位置ベクトル

$\mathbf{p} = \mathbf{p}(\tau) = \frac{d}{d\tau} \mathbf{x}$: 時刻 τ での運動量ベクトル

- $\mathcal{U}(\mathbf{x})$: 位置 \mathbf{x} における物体の位置エネルギー (ポテンシャル)

$\mathcal{K}(\mathbf{p}) = \frac{1}{2} \mathbf{p}^T \mathbf{p}$: 運動量 \mathbf{p} のときの物体の運動エネルギー.

- 今考える運動では, ポテンシャルによる力のみに働くとする.

$\rightarrow \frac{d}{d\tau} \mathbf{p} = - \frac{\partial}{\partial \mathbf{x}} \mathcal{U}(\mathbf{x})$: ポテンシャルの下での運動方程式.

$\mathcal{H}(\mathbf{x}, \mathbf{p}) := \mathcal{U}(\mathbf{x}) + \mathcal{K}(\mathbf{p})$: 系の全エネルギー (Hamiltonian)

- $(\mathbf{x}, \mathbf{p}) \in \mathbb{R}^d \times \mathbb{R}^d$ と考え, 空間 $\mathbb{R}^d \times \mathbb{R}^d$ を相空間 (phase space) という.

Lem. (Hamiltonの運動方程式)

$$\frac{d}{d\tau} \mathbf{x} = \frac{\partial}{\partial \mathbf{p}} \mathcal{H}(\mathbf{x}, \mathbf{p}), \quad \frac{d}{d\tau} \mathbf{p} = - \frac{\partial}{\partial \mathbf{x}} \mathcal{H}(\mathbf{x}, \mathbf{p}).$$

pf. $\frac{\partial}{\partial \mathbf{p}} \mathcal{H}(\mathbf{x}, \mathbf{p}) = \frac{\partial}{\partial \mathbf{p}} \mathcal{K}(\mathbf{p}) = \mathbf{p} = \frac{d}{d\tau} \mathbf{x}.$

$$\frac{\partial}{\partial \mathbf{x}} \mathcal{H}(\mathbf{x}, \mathbf{p}) = \frac{\partial}{\partial \mathbf{x}} \mathcal{U}(\mathbf{x}) = - \frac{d}{d\tau} \mathbf{p}. \quad \square$$

(ポイント)

- Hamiltonian は保存する. 実際,

$$\frac{d\mathcal{H}}{d\tau} = \sum_{i=1}^d \left(\frac{\partial \mathcal{H}}{\partial z_i} \frac{dz_i}{d\tau} + \frac{\partial \mathcal{H}}{\partial p_i} \frac{dp_i}{d\tau} \right) = \sum_{i=1}^d \left(\frac{\partial \mathcal{H}}{\partial z_i} \frac{\partial \mathcal{H}}{\partial p_i} - \frac{\partial \mathcal{H}}{\partial p_i} \frac{\partial \mathcal{H}}{\partial z_i} \right) = 0.$$

\rightarrow 物体は \mathcal{H} を一定に保ちながら運動する.

相空間内の運動の軌道は, Hamiltonの運動方程式を解けば求まる.

2. 可逆性 (reversibility):

$$T_\sigma(\mathbb{z}(\tau), p(\tau)) := (\mathbb{z}(\tau+\sigma), p(\tau+\sigma)) \text{ とすると,}$$

変換 T_σ は 逆変換 $T_\sigma^{-1} = T_{-\sigma}$ とも。

3. 体積保存 (volume preservation):

相空間内の領域 $D \subseteq \mathbb{R}^d \times \mathbb{R}^d$ について, D の体積と

$T_\sigma(D) = \{ T_\sigma(\mathbb{z}, p) \mid (\mathbb{z}, p) \in D \}$ の体積は等しい (Liouville の定理)。

(i.e., $\nu = \left(\frac{d\mathbb{z}}{d\tau}, \frac{dp}{d\tau} \right)$ とすると $\operatorname{div} \nu = 0$.)

→ T_σ を変数変換しても, 変換の Jacobian は 1 になる。
領域の体積拡大率

• T_σ による変換を求めるには, Hamilton の運動方程式

$$\frac{d}{d\tau} \mathbb{z} = \frac{\partial}{\partial p} \mathcal{H}(\mathbb{z}, p), \quad \frac{d}{d\tau} p = - \frac{\partial}{\partial \mathbb{z}} \mathcal{H}(\mathbb{z}, p)$$

を解くのはよい。数値的に解く。

$\varepsilon > 0$: 微小時間幅. $\mathbb{z}^\tau = (\mathbb{z}(\tau) \text{ の計算値})$ などこのように書く

• (前進) Euler 法 としては $\frac{d}{d\tau} \mathbb{z} \approx \frac{\mathbb{z}(\tau+\varepsilon) - \mathbb{z}(\tau)}{\varepsilon}$ などを近似して,

$$\begin{aligned} \mathbb{z}^{\tau+\varepsilon} &= \mathbb{z}^\tau + \varepsilon \frac{\partial}{\partial p} \mathcal{H}(\mathbb{z}, p) \Big|_{\mathbb{z}=\mathbb{z}^\tau, p=p^\tau} = \mathbb{z}^\tau + \varepsilon \frac{d}{d\tau} \mathbb{z} \Big|_{\mathbb{z}=\mathbb{z}^\tau, p=p^\tau} \\ &= \mathbb{z}^\tau + \varepsilon p^\tau \end{aligned}$$

$$\begin{aligned} p^{\tau+\varepsilon} &= p^\tau - \varepsilon \frac{\partial}{\partial \mathbb{z}} \mathcal{H}(\mathbb{z}, p) \Big|_{\mathbb{z}=\mathbb{z}^\tau, p=p^\tau} = p^\tau + \varepsilon \frac{d}{d\tau} p \Big|_{\mathbb{z}=\mathbb{z}^\tau, p=p^\tau} \\ &= p^\tau - \varepsilon \frac{\partial}{\partial \mathbb{z}} \mathcal{U}(\mathbb{z}) \Big|_{\mathbb{z}=\mathbb{z}^\tau} \end{aligned}$$

というスキームで ε 時刻先の相空間での座標を計算。

誤差 $\sim O(\epsilon)$

→ Euler法は離散化の数値誤差が大きい (1次精度)

• リ-プフロック法では、運動量の離散化を $\frac{\epsilon}{2}$ 刻みにして置く。

$$\begin{cases}
 1. & p^{\tau+\frac{\epsilon}{2}} = p^{\tau} - \frac{\epsilon}{2} \frac{\partial}{\partial z} \mathcal{U}(z) \Big|_{z=z^{\tau}} \\
 2. & z^{\tau+\epsilon} = z^{\tau} + \epsilon p^{\tau+\frac{\epsilon}{2}} \quad \leftarrow z \text{の計算で } p^{\tau} \\
 & \quad \quad \quad \tau, \tau+\epsilon \text{の平均を用いる。} \\
 3. & p^{\tau+\epsilon} = p^{\tau+\frac{\epsilon}{2}} - \frac{\epsilon}{2} \frac{\partial}{\partial z} \mathcal{U}(z) \Big|_{z=z^{\tau+\epsilon}}
 \end{cases}$$

$L \in \mathbb{Z}_{>0}$: ステップ数として、上の更新を L 回行い、時刻 ϵL 先の

物体の位置と運動量 $(z^{\tau+\epsilon L}, p^{\tau+\epsilon L})$ を計算する。

→ 数値誤差が Euler法よりも小さくなる (2次精度)

また、このスキームは可逆性があり、 $(z^{\tau+\epsilon L}, p^{\tau+\epsilon L})$ から $-\tau$ 方向に

リ-プフロック法を L 回行くと、もとの (z^{τ}, p^{τ}) が得られる。

さらに、シンプティック性を有しており、 $(z^{\tau}, p^{\tau}) \mapsto (z^{\tau+\epsilon L}, p^{\tau+\epsilon L})$ の

変換で Hamiltonian が保存される。 ← もちろん誤差はゼロか、振動するだけ

cf.) シンプティック積分法 (Symplectic Integrator)

• T_{ϵ} のシミュレーションをする上で、リ-プフロック法を利用するのが良い。

4.1.6.2 サンプリングアルゴリズムへの適用

• サンプリングしたい分布を $p(z) = \frac{1}{Z} \tilde{p}(z)$ とする。 ($z \in \mathbb{R}^d$)

$p \in \mathbb{R}^d$ を導入して、同時分布を $p(z, p) = p(z)p(p)$ とする。

(p と z は独立だと仮定している)

$(z, p) \in p(z, p)$ からサンプリングすること、 z は p と無関係に $p(z)$ からサンプリングされることになる。

• $p(p) = \mathcal{N}(p | 0, I_d) = \frac{1}{\sqrt{(2\pi)^d}} \exp(-\frac{1}{2} p^T p)$ とする。

$p(z) = \frac{1}{Z} \tilde{p}(z) = \frac{1}{Z} \exp(-\mathcal{U}(z))$ と表す。 ← $\tilde{p}(z) > 0$ なのでこのように表す。

同時分布は、

$p(z, p) = \frac{1}{Z} \frac{1}{\sqrt{(2\pi)^d}} \exp(-\mathcal{U}(z) - \frac{1}{2} p^T p) \propto \exp(-\mathcal{H}(z, p))$
= $\mathcal{K}(p)$

となる。

• Metropolis-Hastings法では、提案分布 $q(z|\xi)$ からのサンプリングをしていた。

Hamiltonian Monte Carlo では、代わりに $\phi \in \mathcal{N}(\phi | 0, I_d)$ からサンプリングし、

(ξ, ϕ) の位置から σ たりの時間発展させて得られる

$(z, p) = T_\sigma(\xi, \phi)$ への遷移を考える。
本当はこの leapfrog 離散近似を考えている。

T_σ の体積保存性から、推移確率密度が

$\mathcal{T}((\xi, \phi), (z, p)) = \exp(-\mathcal{H}(z, p) + \mathcal{H}(\xi, \phi))$

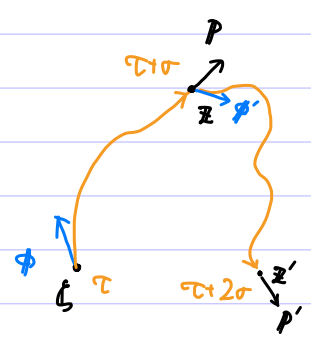
と表す。更に可逆性より同様に

$\mathcal{T}((z, p), (\xi, \phi)) = \exp(-\mathcal{H}(\xi, \phi) + \mathcal{H}(z, p))$

も成り立つ。すると、

$p(z, p) \mathcal{T}((z, p), (\xi, \phi)) = \exp(-\mathcal{H}(\xi, \phi))$

$p(\xi, \phi) \mathcal{T}((\xi, \phi), (z, p)) = \exp(-\mathcal{H}(z, p))$



Hamiltonian は時間発展で不変なので,

$$p(\mathbf{z}, p) \mathcal{T}((\mathbf{z}, p), (\xi, \phi)) = p(\xi, \phi) \mathcal{T}((\xi, \phi), (\mathbf{z}, p))$$

と詳細釣り合い条件が成立し, $p(\mathbf{z}, p)$ はこの Markov 連鎖の定常分布.

- 実際は T_0 を適用する代わりに leapfrog によるシミュレーションを行うので, 数値誤差のため Hamiltonian は完全に保存されるわけではない.

→ Metropolis 法での採択率

$$A((\xi, \phi), (\mathbf{z}, p)) := \min\left(1, \frac{p(\mathbf{z}, p)}{p(\xi, \phi)}\right) = \min\left(1, \exp(-\mathcal{H}(\mathbf{z}, p) + \mathcal{H}(\xi, \phi))\right)$$

を利用する. 数値誤差があっても, leapfrog 法では $\mathcal{H}(\mathbf{z}, p) \approx \mathcal{H}(\xi, \phi)$

となるので, 採択率は 1 に近くなる!

Algo. (Hamiltonian Monte Carlo)

Input: $\tilde{p}(\mathbf{z})$, $\varepsilon > 0$, $L \in \mathbb{Z}_{>0}$, $N \in \mathbb{Z}_{>0}$

Output: Markov 連鎖 $\{\mathbf{z}^{(i)}\}_{i=1}^N$

1. $\mathbf{z}^{(1)}$ を ランダムに初期化;
2. for $i = 1, \dots, N-1$ do:
3. $p \sim \mathcal{N}(p | 0, I_d)$; $u \sim \text{Unif}(0, 1)$;
4. $(\mathbf{z}^*, p^*) \leftarrow \text{leapfrog}_{\varepsilon, L}(\mathbf{z}^{(i)}, p)$; $r \leftarrow \frac{p(\mathbf{z}^*, p^*)}{p(\mathbf{z}^{(i)}, p)}$;
5. $\mathbf{z}^{(i+1)} \leftarrow \mathbf{z}^*$ if $u \leq \min(1, r)$ else $\mathbf{z}^{(i)}$
6. return $\{\mathbf{z}^{(i)}\}_{i=1}^N$.

4.1.6.3 Langevin 動力学法.

- Langevin 動力学法 (Langevin dynamics method):

$L=1$ と $L_T = \text{HMC}$.

leapfrog の式をまとめ,

$$\begin{aligned} \mathbb{Z}^{T+\varepsilon} &= \mathbb{Z}^T + \varepsilon p^{T+\frac{\varepsilon}{2}} = \mathbb{Z}^T + \varepsilon \left(p^T - \frac{\varepsilon}{2} \frac{\partial}{\partial \mathbb{Z}} \mathcal{U}(\mathbb{Z}) \Big|_{\mathbb{Z}=\mathbb{Z}^T} \right) \\ &= \mathbb{Z}^T - \frac{\varepsilon^2}{2} \frac{\partial}{\partial \mathbb{Z}} \mathcal{U}(\mathbb{Z}) \Big|_{\mathbb{Z}=\mathbb{Z}^T} + \varepsilon p^T. \end{aligned}$$

により \mathbb{Z}^* を計算する. \mathbb{Z}^* をそのまま採択してしまうことで計算を簡便にする.

- 少しランダムウォーク的な挙動がある.

4.1.7 Gibbs サンプルング

- 分布 $p(\mathbb{Z})$ から直接 \mathbb{Z} 全体を サンプルングするのが難しいとき,

$\mathbb{Z} = (\mathbb{Z}_1, \dots, \mathbb{Z}_M)$ と M 個に分けて, 次のように逐次的に
サンプルングする (Gibbs サンプルング).
Gibbs 分布からサンプルングするに便利だ.

Algo. (Gibbs サンプルング)

Input: $p(\mathbb{Z}) = p(\mathbb{Z}_1, \dots, \mathbb{Z}_M)$, $N \in \mathbb{Z}_0$

Output: サンプル $\{\mathbb{Z}^{(i)}\}_{i=1}^N$

1 $\mathbb{Z}^{(1)}$: ランダムに初期化;

2 for $i = 1, \dots, N-1$ do:

3 for $m = 1, \dots, M$ do:

4.
$$Z_m^{(i+1)} \sim p(Z_m | Z_1^{(i+1)}, \dots, Z_{m-1}^{(i+1)}, Z_{m+1}^{(i)}, \dots, Z_n^{(i)})$$

return $\{(Z_1^{(i)}, \dots, Z_n^{(i)})\}_{i=1}^N$

- ・ サンプルを得たい変数の数が膨大なとき、複数の確率モデルを組み合わせた巨大な確率モデルからのサンプリングに利用するとよい。
- ・ Gibbsサンプリングは、採択率が常に1のMetropolis法。

§4.2 最適化に基づく推論手法

・ サンプリング法

・ 確率的な近似

✓ サンプルサイズを増やせば

厳密解に近づく

✗ 必要なサンプルサイズが不明

✗ 計算コストが膨大

最適化に基づく手法

・ 決定的な近似

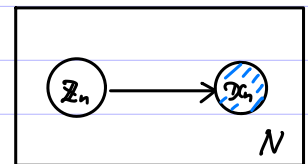
✓ 収束性能が良い

✗ 厳密解は決して求まらない

< 問題設定 (再確認) >

$\mathcal{X} = \{x_1, \dots, x_N\}$: 観測データ

$\mathcal{Z} = \{z_1, \dots, z_N\}$: 潜在変数 (非観測)



$p(\mathcal{X}, \mathcal{Z}) = p(\mathcal{X} | \mathcal{Z}) p(\mathcal{Z}) = \prod_{n=1}^N p(x_n | z_n) p(z_n)$: 確率モデル

$z_n \stackrel{i.i.d.}{\sim} p(z)$

事後分布 $p(\mathcal{Z} | \mathcal{X})$ を求めたい

4.2.1 変分推論法

・ **変分推論** (variational inference) :

分布の集合 \mathcal{D} に属する分布のうち、分布 $p(z)$ と "うまく近似する" 分布 $q^*(z)$ と

(何らかの) 変分問題 $q^*(z) = \underset{q \in \mathcal{D}}{\operatorname{argmin}} \mathcal{L}[q]$ を解くことで求める手法.

$\mathcal{L} : \mathcal{D} \rightarrow \mathbb{R}$ と汎関数 (functional) という.

"函数の函数"

Remark したがって変分問題を解く際には変分法を用いるが、ここで扱うものは

変分法を明示的に使わなくても解ける。

・ 事後分布について、Bayesの定理より

$$p(\theta | \mathcal{X}) = \frac{p(\mathcal{X} | \theta) p(\theta)}{p(\mathcal{X})} = p(\mathcal{X}, \theta) p(\mathcal{X})^{-1}$$

$$\Leftrightarrow \log p(\theta | \mathcal{X}) = -\log p(\mathcal{X}) + \log p(\mathcal{X}, \theta).$$

$p(\mathcal{X}, \theta)$ は確率モデル (これは自分で設定したものになっている)。LML.

周辺尤度 $p(\mathcal{X}) = \int p(\mathcal{X}, \theta) d\theta$ は、モデルが複雑になると

解析的に計算できない。

・ θ についての適当な分布 $q(\theta)$ を考える

$$\log p(\mathcal{X}) = \log p(\mathcal{X}, \theta) - \log p(\theta | \mathcal{X})$$

の両辺 $q(\theta)$ について期待値をとると、

$$\int q(\theta) \log p(\mathcal{X}) d\theta = \int q(\theta) \log p(\mathcal{X}, \theta) d\theta - \int q(\theta) \log p(\theta | \mathcal{X}) d\theta$$

$$(\text{l.h.s.}) = \log p(\mathcal{X}) \underbrace{\int q(\theta) d\theta}_{=1} = \log p(\mathcal{X})$$

(r.h.s.)

$$= \int q(\theta) \log \frac{p(\mathcal{X}, \theta)}{q(\theta)} \cdot \frac{q(\theta)}{p(\theta | \mathcal{X})} d\theta$$

$$= \mathbb{E}_{q(\theta)} \left[\log \frac{p(\mathcal{X}, \theta)}{q(\theta)} \right] + \mathbb{E}_{q(\theta)} \left[\log \frac{q(\theta)}{p(\theta | \mathcal{X})} \right]$$

$$= \mathbb{E}_{q(\theta)} \left[\log \frac{p(\mathcal{X}, \theta)}{q(\theta)} \right] + D_{\text{KL}} [q(\theta) \| p(\theta | \mathcal{X})].$$

$$\mathcal{L}[q] := \mathbb{E}_{q(\theta)} \left[\log \frac{p(\mathcal{X}, \theta)}{q(\theta)} \right] \text{ とおくと、}$$

$$\log p(x) = \mathcal{L}[q] + D_{KL}[q(z) \| p(z|x)] \quad \text{を得る.}$$

Kullback-Leibler divergence は「非負なもの」,

$$\log p(x) \geq \mathcal{L}[q].$$

$\mathcal{L}[q]$ は 対数周辺尤度 の下界を与える. これは
エビデンス モデルから与えられたときのデータxの出現する尤もらしさ.

エビデンス下界 (evidence lower bound: ELBO), あるいは

変分下界 (variational lower bound) という. cf.) $\log p(x|\theta)$ を θ について最大化するのは
最適法の種類が.

$\log p(x)$ が求まらなくても, $\mathcal{L}[q]$ を最大化すること

$\log p(x)$ の (下からの) よい評価が得られる.

Remark • $\mathcal{F}[q] := -\mathcal{L}[q]$ を **変分エネルギー** (variational energy) という.

• $\mathcal{L}[q]$ の設計の仕方はい他にもある.

cf.) Rényi Lower Bound.

$$\begin{aligned} \bullet \quad q^*(z) &:= \operatorname{argmax}_q \mathcal{L}[q] \quad \text{qによらず一定.} \\ &= \operatorname{argmax}_q \left(\log p(x) - D_{KL}[q(z) \| p(z|x)] \right) \\ &= \operatorname{argmin}_q D_{KL}[q(z) \| p(z|x)]. \end{aligned}$$

→ $q(z)$ が $p(z|x)$ をうまく近似するように最適化している.

• $D_{KL}[q(z) \| p(z|x)] \geq 0$ で, 等号成立は $q(z) = p(z|x)$ のとき.

$p(z|x)$ は $p(x)$ のために求まらなかったもので, これでは意味がない...

→ 考える分布 q をある範囲 \mathcal{D} に制限しよう.

• q を扱いやすい分布にできる p にも、代わりに表現力 p にも制限されてしまうため、

近似精度には限界 p がある。

• 制限の仕方。

1. q がパラメトリックな分布 $q(\mathcal{Z}|\xi)$ に制限する。

• パラメータ ξ が **変分パラメータ** (variational parameter) という

→ このとき ELBO は ξ の関数 $L(\xi)$ になる。最適なパラメータ ξ^* を決めるには、

$$\xi^* = \operatorname{argmax}_{\xi} L(\xi)$$

を適当な最適化手法を用いて解けばよい。

2 **平均場近似** (mean field approximation)。

$\mathcal{Z} = (\mathcal{Z}_1, \dots, \mathcal{Z}_M)$ と M に分割して、独立性を仮定：

$$q(\mathcal{Z}) = \prod_{i=1}^M q_i(\mathcal{Z}_i)$$

→ 各近似分布 $q_1(\mathcal{Z}_1), \dots, q_M(\mathcal{Z}_M)$ を順にそれぞれ最適化する

• $q_j(\mathcal{Z}_j)$ についての最適化をする 他 $q_i(\mathcal{Z}_i)$ ($i \neq j$) は固定する。

$$L[q]$$

$$= \mathbb{E}_{q(\mathcal{Z})} \left[\log \frac{p(\mathcal{X}, \mathcal{Z})}{q(\mathcal{Z})} \right]$$

$$= \int q_j(\mathcal{Z}_j) \prod_{i \neq j} q_i(\mathcal{Z}_i) \left(\log p(\mathcal{X}, \mathcal{Z}) - \sum_{i \neq j} \log q_i(\mathcal{Z}_i) - \log q_j(\mathcal{Z}_j) \right) d\mathcal{Z}$$

$$= \int q_j(\mathcal{Z}_j) \left(\int \prod_{i \neq j} q_i(\mathcal{Z}_i) \log p(\mathcal{X}, \mathcal{Z}) d\mathcal{Z}_{\substack{i \\ j \text{ 以外}}} \right) d\mathcal{Z}_j$$

$$- \sum_{i \neq j} \int q_i(\mathcal{Z}_i) \log q_i(\mathcal{Z}_i) d\mathcal{Z}_i - \int q_j(\mathcal{Z}_j) \log q_j(\mathcal{Z}_j) d\mathcal{Z}_j.$$

$\therefore \mathcal{L}[q] = \int \prod_{i \neq j} q_i(\theta_i) \log p(\mathcal{X}, \theta) d\theta_j =: \mathbb{E}_{i \neq j} [\log p(\mathcal{X}, \theta)]$ と書ける.

$\sum_{i \neq j} \int q_i(\theta_i) \log q_i(\theta_i) d\theta_i = \text{const.}$ に注意すると,

$$\mathcal{L}[q] = \int q_j(\theta_j) \left(\mathbb{E}_{i \neq j} [\log p(\mathcal{X}, \theta)] - \log q_j(\theta_j) \right) d\theta_j + \text{const.}$$

さらに, $\log \tilde{p}(\mathcal{X}, \theta_j) = \mathbb{E}_{i \neq j} [\log p(\mathcal{X}, \theta)] + \text{const.}$ とする分布 $\tilde{p}(\mathcal{X}, \theta_j)$

を考える. すると,

$$\mathcal{L}[q] = - \int q_j(\theta_j) \log \frac{q_j(\theta_j)}{\tilde{p}(\mathcal{X}, \theta_j)} d\theta_j + \text{const.}$$

$$= - D_{\text{KL}} [q_j(\theta_j) \parallel \tilde{p}(\mathcal{X}, \theta_j)] + \text{const.}$$

よって,

$$q_j^*(\theta_j)$$

$$= \operatorname{argmax}_{q_j} \mathcal{L}[q]$$

$$= \operatorname{argmin}_{q_j} D_{\text{KL}} [q_j(\theta_j) \parallel \tilde{p}(\mathcal{X}, \theta_j)]$$

← 一般には $\mathcal{L}[q]$ を変分法を用いて最適解を求め、
 $\mathcal{L}[q]$ は KL divergence の最小化なので簡単.

$$= \tilde{p}(\mathcal{X}, \theta_j)$$

$$\propto \exp \left(\mathbb{E}_{i \neq j} [\log p(\mathcal{X}, \theta)] \right).$$

$q_j^*(\theta_j)$ が分布になるように, 正規化定数を計算すると,

$$q_j^*(\theta_j) = \frac{\exp \left(\mathbb{E}_{i \neq j} [\log p(\mathcal{X}, \theta)] \right)}{\int \exp \left(\mathbb{E}_{i \neq j} [\log p(\mathcal{X}, \theta)] \right) d\theta_j}$$

座標降下法.

↓

→ 全ての因子 $q_i(\theta_i)$ を適当に初期化して, 各 $q_i(\theta_i)$ を上式で逐次最適化する.

$\mathcal{L}[q]$ は各 q_i に関して凸なので, このスケールで q^* に収束する.

Remark $\mathcal{L}[q]$

$$= \int q(\mathcal{Z}) \log \frac{p(\mathcal{X}, \mathcal{Z})}{q(\mathcal{Z})} d\mathcal{Z}$$

$$= \int q(\mathcal{Z}) \left(\log p(\mathcal{X}|\mathcal{Z}) + \log \frac{p(\mathcal{Z})}{q(\mathcal{Z})} \right) d\mathcal{Z}$$

$$= \mathbb{E}_{q(\mathcal{Z})} [\log p(\mathcal{X}|\mathcal{Z})] - D_{KL}[q(\mathcal{Z}) \| p(\mathcal{Z})].$$

$\mathcal{L}[q] \rightarrow \max$ するとき

第1項: 観測データの対数尤度の q による期待値

\rightarrow なるべく大きくなる.

$\log p(\mathcal{X}|\mathcal{Z})$ を最大にする \mathcal{Z} を $\hat{\mathcal{Z}}$ とし, $q(\mathcal{Z}) = \delta(\mathcal{Z} - \hat{\mathcal{Z}})$ (デルタ分布) のとき最大

第2項: q の p に対する KL-divergence.

\rightarrow なるべく小さくなる.

q が p から大きく離れた分布になるのを防ぐ. (正則化の効果)

変分推論により得られる q^* は, 2項のバランスのとれた解になっている.

• 変分 EM アルゴリズム (variational expectation maximization algo.)

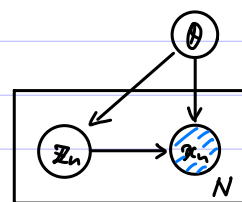
• 潜在変数 $\mathcal{Z} = \{z_1, \dots, z_N\}$ の他に, パラメータ θ が存在する

状況と考える.

• \mathcal{Z} : データ数が増えるにつれ次元も増える (外延的)

• θ : 固定次元 (内延的)

• モデル: $p(\mathcal{X}, \mathcal{Z}, \theta) = p(\mathcal{X}|\mathcal{Z}, \theta) p(\mathcal{Z}|\theta) p(\theta)$



- 変分EMアルゴリズムでは,

$$p(\mathcal{Z}, \theta) = p(\mathcal{Z} | \theta) p(\theta) \approx q(\mathcal{Z}, \theta) = q(\mathcal{Z}) q(\theta)$$

と平均場近似をして変分推論を行う.

- 変分E-step: パラメータの近似分布 $q(\theta)$ を固定し,

潜在変数の近似分布 $q(\mathcal{Z})$ を更新する.

$$q^*(\mathcal{Z}) = \frac{\exp(\mathbb{E}_{q(\theta)}[\log p(\mathcal{Z}, \mathcal{Z}, \theta)])}{\int \exp(\mathbb{E}_{q(\theta)}[\log p(\mathcal{Z}, \mathcal{Z}, \theta)]) d\mathcal{Z}}.$$

ここで,

$$\begin{aligned} & \exp(\mathbb{E}_{q(\theta)}[\log p(\mathcal{Z}, \mathcal{Z}, \theta)]) \\ &= \exp(\mathbb{E}_{q(\theta)}[\log p(\mathcal{Z}, \mathcal{Z} | \theta) + \log p(\theta)]) \\ &= \exp(\mathbb{E}_{q(\theta)}[\log p(\mathcal{Z}, \mathcal{Z} | \theta)]) \underbrace{\exp(\mathbb{E}_{q(\theta)}[\log p(\theta)])}_{\mathcal{Z} \text{ に依らない}} \end{aligned}$$

なので,

$$\begin{aligned} & \int \exp(\mathbb{E}_{q(\theta)}[\log p(\mathcal{Z}, \mathcal{Z} | \theta)]) d\mathcal{Z} \\ &= \exp(\mathbb{E}_{q(\theta)}[\log p(\theta)]) \int \exp(\mathbb{E}_{q(\theta)}[\log p(\mathcal{Z}, \mathcal{Z} | \theta)]) d\mathcal{Z}. \end{aligned}$$

結局,

$$q^*(\mathcal{Z}) = \frac{\exp(\mathbb{E}_{q(\theta)}[\log p(\mathcal{Z}, \mathcal{Z} | \theta)])}{\int \exp(\mathbb{E}_{q(\theta)}[\log p(\mathcal{Z}, \mathcal{Z} | \theta)]) d\mathcal{Z}}.$$

と更新する.

・変分 M-step: 潜在変数の近似分布 $q(\mathcal{Z})$ を固定し,

パラメータの近似分布 $q(\theta)$ を更新する.

$$q^*(\theta) = \frac{\exp(\mathbb{E}_{q(\mathcal{Z})}[\log p(\mathcal{X}, \mathcal{Z}, \theta)])}{\int \exp(\mathbb{E}_{q(\mathcal{Z})}[\log p(\mathcal{X}, \mathcal{Z}, \theta)]) d\theta}$$

先と同様に,

$$\begin{aligned} & \exp(\mathbb{E}_{q(\mathcal{Z})}[\log p(\mathcal{X}, \mathcal{Z}, \theta)]) \\ = & \exp(\mathbb{E}_{q(\mathcal{Z})}[\log p(\mathcal{X}, \mathcal{Z} | \theta)]) \exp(\mathbb{E}_{q(\mathcal{Z})}[\log p(\theta)]) \\ = & \exp(\mathbb{E}_{q(\mathcal{Z})}[\log p(\mathcal{X}, \mathcal{Z} | \theta)]) \exp(\log p(\theta)) \\ = & p(\theta) \exp(\mathbb{E}_{q(\mathcal{Z})}[\log p(\mathcal{X}, \mathcal{Z} | \theta)]). \end{aligned}$$

よって

$$q^*(\theta) = \frac{p(\theta) \exp(\mathbb{E}_{q(\mathcal{Z})}[\log p(\mathcal{X}, \mathcal{Z} | \theta)])}{\int p(\theta) \exp(\mathbb{E}_{q(\mathcal{Z})}[\log p(\mathcal{X}, \mathcal{Z} | \theta)]) d\theta}$$

と更新する.

cf.) (通常の) EM アルゴリズム.

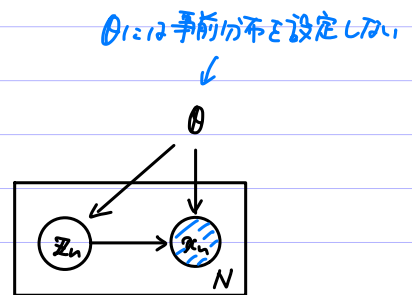
$$\text{モデル: } p(\mathcal{X}, \mathcal{Z}; \theta) = p(\mathcal{X} | \mathcal{Z}; \theta) p(\mathcal{Z}; \theta).$$

θ はパラメータ.

θ について, 対数尤度

$$\log p(\mathcal{X}; \theta) = \log \int p(\mathcal{X}, \mathcal{Z}; \theta) d\mathcal{Z}$$

を最大化する (θ の最尤推定).



q についての適当な分布 $q(z)$ を考えよ.

$$\begin{aligned}\log p(x; \theta) &= \log p(x; \theta) \int q(z) dz \\ &= \int q(z) \log p(x; \theta) dz \\ &= \int q(z) \log \frac{p(x, z; \theta)}{p(z|x; \theta)} dz \quad \left. \begin{array}{l} f=g \\ \Rightarrow \int_A f(x) dx = \int_A g(x) dx \end{array} \right\} \\ &= \int q(z) \log \frac{p(x, z; \theta)}{q(z)} dz + \int q(z) \log \frac{q(z)}{p(z|x; \theta)} dz \\ &= \int q(z) \log \frac{p(x, z; \theta)}{q(z)} dz + D_{KL}[q(z) \| p(z|x; \theta)].\end{aligned}$$

これを $\mathcal{L}_\theta[q] := \int q(z) \log \frac{p(x, z; \theta)}{q(z)} dz$ とおくと.

$$\begin{aligned}\log p(x; \theta) &= \mathcal{L}_\theta[q] + \underbrace{D_{KL}[q(z) \| p(z|x; \theta)]}_{\geq 0} \\ &\geq \mathcal{L}_\theta[q].\end{aligned}$$

下界 $\mathcal{L}_\theta[q]$ を θ に関して最大化することと尤度を最大化する θ を求める.

このために, q と θ に関して交互に最適化していく.

• E-step: θ を固定して $q(z)$ を更新する.

$$\mathcal{L}_\theta[q] = \underbrace{\log p(x; \theta)}_{q \text{ に依らず}} - D_{KL}[q(z) \| p(z|x; \theta)] \quad \text{となる.}$$

$$\begin{aligned}q^*(z) &= \operatorname{argmax}_{q(z)} \mathcal{L}_\theta[q] \\ &= \operatorname{argmin}_{q(z)} D_{KL}[q(z) \| p(z|x; \theta)] \\ &= p(z|x; \theta). \quad (\text{この事後分布})\end{aligned}$$

このとき,

この期待値をとるためのステップ "E-step".

$$\mathcal{L}[q^*] = \int q^*(z) \log p(x, z; \theta) dz - \int q^*(z) \log q^*(z) dz = \mathbb{E}_{q^*(z)}[\log p(x, z; \theta)] + \text{const}$$

• M-step: $q(z)$ を固定して θ を更新する.

$$\theta^* = \operatorname{argmax}_{\theta} \mathcal{L}_{\theta}[q]$$

$$= \operatorname{argmax}_{\theta} \int q(z) \log p(x, z; \theta) dz$$

$$= \operatorname{argmax}_{\theta} \mathbb{E}_{q(z)} [\log p(x, z; \theta)].$$

← θ によって最大化するステップを "M-step" と呼ぶ。
"最大化しよ" のこと。

• 以上の E-step, M-step を繰り返して $\mathcal{L}_{\theta}[q]$ が減少することは無い。

E-step を行う前のパラメータを θ^{old} とする。

E-step 終了後の q を q^{new} とすると, $\mathcal{L}_{\theta^{\text{old}}}[q^{\text{new}}] = \log p(x; \theta^{\text{old}})$ が成立。

M-step 終了後のパラメータを θ^{new} とすると, $q^{\text{new}}(z) = p(z|x; \theta^{\text{old}})$ とする。

$$\log p(x; \theta^{\text{new}}) = \mathcal{L}_{\theta^{\text{new}}}[q^{\text{new}}] + D_{\text{KL}}[q^{\text{new}}(z) \| p(z|x; \theta^{\text{new}})]$$

$$\geq \mathcal{L}_{\theta^{\text{old}}}[q^{\text{new}}] + D_{\text{KL}}[q^{\text{new}}(z) \| p(z|x; \theta^{\text{new}})]$$

$$= \log p(x; \theta^{\text{old}}) + D_{\text{KL}}[q^{\text{new}}(z) \| p(z|x; \theta^{\text{new}})].$$

$$= \log p(x; \theta^{\text{old}}) + D_{\text{KL}}[p(z|x; \theta^{\text{old}}) \| p(z|x; \theta^{\text{new}})].$$

$$\geq \log p(x; \theta^{\text{old}}).$$

M-step での
 θ によって最大化
していきよ" のこと
 $\mathcal{L}_{\theta^{\text{new}}}[q] \geq \mathcal{L}_{\theta^{\text{old}}}[q]$
が $\forall q$ によって成立

よって EM algo. 1 ステップで対数尤度は増加する。

$\theta^{\text{new}} = \theta^{\text{old}}$ のとき, またはそのときに限り等号成立。

対数尤度が収束したときのパラメータ θ が最大推定量になる。

4.2.2 例: 平均場近似による潜在変数モデルの学習.

・ 潜在変数モデル (latent variable model) による次元削減を扱う.

・ $\mathcal{X} = \{x_1, \dots, x_N\} \subseteq \mathbb{R}^d$: 観測データ

$\mathcal{Z} = \{z_1, \dots, z_N\} \subseteq \mathbb{R}^l$: 潜在変数 ($l \ll d$).

x_i と z_i で表現し, データの圧縮, 可視化, 特徴量抽出などに使う.

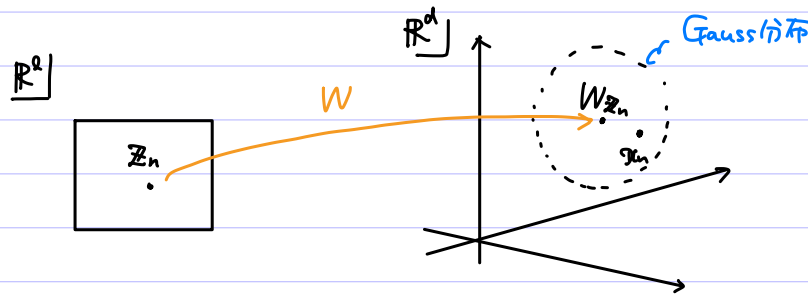
・ PCA, ICA, 行列分解, t-means などはこのモデルで表現できる.

4.2.2.1 線形次元削減への適用.

・ 仮定.

$W \in M_{d \times l}(\mathbb{R})$ と $l < d$, $\sigma_x^2 > 0$ と定数とす

$$p(\mathcal{X} | \mathcal{Z}, W) = \prod_{n=1}^N p(x_n | z_n, W) = \prod_{n=1}^N \mathcal{N}(x_n | Wz_n, \sigma_x^2 I_d)$$



\mathcal{Z} の事前分布: $p(\mathcal{Z}) = \prod_{n=1}^N \mathcal{N}(z_n | 0, I_l)$

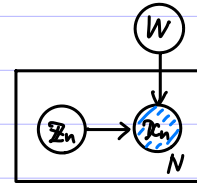
W の事前分布:

$$W = \begin{pmatrix} w_1^T \\ \vdots \\ w_d^T \end{pmatrix} \quad (w_j \in \mathbb{R}^l), \quad \sigma_w^2 > 0 : \text{定数} \quad \text{とす.}$$

$$p(W) = \prod_{j=1}^d \mathcal{N}(w_j | 0, \sigma_w^2 I_l).$$

• モデル:

$$\begin{aligned} p(\mathcal{X}, \theta, W) &= p(\mathcal{X} | \theta, W) p(\theta, W) \\ &= p(\mathcal{X} | \theta, W) p(\theta) p(W). \end{aligned}$$



• 事後分布 $p(\theta, W | \mathcal{X})$ に次々近似:

$$p(\theta, W | \mathcal{X}) \approx q(\theta) q(W).$$

変分推論 (変分EM algo.) を行う.

i 回目の更新後の近似分布を $q_i(\theta)$, $q_i(W)$ と書く.

• 変分 E-step. θ の近似分布を更新.

$$\begin{aligned} q_{i+1}(\theta) &\propto \exp(\mathbb{E}_{q_i(W)}[\log p(\mathcal{X}, \theta | W)]). && W \text{ に依存} \\ &= \exp(\mathbb{E}_{q_i(W)}[\log p(\mathcal{X} | \theta, W)] + \mathbb{E}_{q_i(W)}[\log p(\theta)]) \\ &= \exp(\mathbb{E}_{q_i(W)}[\log p(\mathcal{X} | \theta, W)]) \exp(\log p(\theta)) \\ &= p(\theta) \exp(\mathbb{E}_{q_i(W)}[\log p(\mathcal{X} | \theta, W)]). \end{aligned}$$

• 変分 M-step. W の近似分布を更新.

$$\begin{aligned} q_{i+1}(W) &\propto p(W) \exp(\mathbb{E}_{q_{i+1}(\theta)}[\log p(\mathcal{X}, \theta | W)]) \\ &\propto p(W) \exp(\mathbb{E}_{q_{i+1}(\theta)}[\log p(\mathcal{X} | \theta, W)]). \end{aligned}$$

• 各 step の計算をもう少し進めよう.

$$\begin{aligned} &\log p(\mathcal{X} | \theta, W) \\ &= \log \left(\prod_{n=1}^N \mathcal{N}(x_n | W z_n, \sigma_x^2 \mathbf{I}_d) \right) \end{aligned}$$

$$\begin{aligned}
&= \sum_{n=1}^N \log \mathcal{N}(\mathcal{X}_n | W \mathcal{Z}_n, \sigma_x^2 I_d) \quad \downarrow \mathcal{Z}_n, W \text{ に依存しない項は } \\
&\quad \text{Const. にまとめます.} \\
&= \sum_{n=1}^N \left(-\frac{1}{2\sigma_x^2} (\mathcal{X}_n - W \mathcal{Z}_n)^T (\mathcal{X}_n - W \mathcal{Z}_n) \right) + \text{Const.} \\
&= -\frac{1}{2\sigma_x^2} \sum_{n=1}^N \left(\underbrace{\mathcal{X}_n^T \mathcal{X}_n}_{\text{Const.}} - \mathcal{X}_n^T W \mathcal{Z}_n - \underbrace{\mathcal{Z}_n^T W^T \mathcal{X}_n}_{\text{スカラー-ベクトル積置換 E と E^T も同じ.}} + \mathcal{Z}_n^T W^T W \mathcal{Z}_n \right) + \text{Const.} \\
&= -\frac{1}{2\sigma_x^2} \sum_{n=1}^N \left(\mathcal{Z}_n^T W^T W \mathcal{Z}_n - 2 \mathcal{X}_n^T W \mathcal{Z}_n \right) + \text{Const.}
\end{aligned}$$

よって, $\mathcal{Z}_n^T W^T \mathcal{X}_n$ スカラー-ベクトル積置換 E と E^T も同じ.

$$\mathbb{E}_{q_i(w)} [\log p(\mathcal{X} | \mathcal{Q}, W)]$$

$$\begin{aligned}
&= -\frac{1}{2\sigma_x^2} \sum_{n=1}^N \left(\mathcal{Z}_n^T \mathbb{E}_{q_i(w)} [W^T W] \mathcal{Z}_n - 2 \mathcal{Z}_n^T \mathbb{E}_{q_i(w)} [W] \mathcal{X}_n \right) + \text{Const.} \\
&= -\frac{1}{2\sigma_x^2} \sum_{n=1}^N \left(\mathcal{Z}_n^T \left(\sum_{j=1}^d \mathbb{E}_{q_i(w)} [w_j w_j^T] \right) \mathcal{Z}_n - 2 \mathcal{Z}_n^T \left(\sum_{j=1}^d \mathcal{X}_{n,j} \mathbb{E}_{q_i(w)} [w_j] \right) \right) + \text{Const.}
\end{aligned}$$

$$\mathbb{E}_{q_{int}(\mathcal{Q})} [\log p(\mathcal{X} | \mathcal{Q}, W)]$$

$$\begin{aligned}
&= -\frac{1}{2\sigma_x^2} \sum_{n=1}^N \left(\mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n^T W^T W \mathcal{Z}_n] - 2 \mathcal{X}_n^T W \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n] \right) + \text{Const.} \quad \downarrow \text{スカラー-ベクトル積置換 E と E^T も同じ.} \\
&= -\frac{1}{2\sigma_x^2} \sum_{n=1}^N \left(\mathbb{E}_{q_{int}(\mathcal{Q})} [\text{tr}(\mathcal{Z}_n^T W^T W \mathcal{Z}_n)] - 2 \mathcal{X}_n^T W \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n] \right) + \text{Const.} \quad \downarrow \text{tr}(AB) = \text{tr}(BA) \\
&= -\frac{1}{2\sigma_x^2} \sum_{n=1}^N \left(\mathbb{E}_{q_{int}(\mathcal{Q})} [\text{tr}(W \mathcal{Z}_n \mathcal{Z}_n^T W^T)] - 2 \mathcal{X}_n^T W \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n] \right) + \text{Const.} \quad \downarrow \text{tr の計算} \\
&= -\frac{1}{2\sigma_x^2} \sum_{n=1}^N \left(\mathbb{E}_{q_{int}(\mathcal{Q})} \left[\sum_{j=1}^d w_j^T \mathcal{Z}_n \mathcal{Z}_n^T w_j \right] - 2 \mathcal{X}_n^T W \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n] \right) + \text{Const.} \\
&= -\frac{1}{2\sigma_x^2} \sum_{n=1}^N \left(\sum_{j=1}^d w_j^T \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n \mathcal{Z}_n^T] w_j - 2 \mathcal{X}_n^T W \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n] \right) + \text{Const.} \quad \downarrow \text{スカラー-ベクトル積置換} \\
&= -\frac{1}{2\sigma_x^2} \sum_{n=1}^N \left(\sum_{j=1}^d w_j^T \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n \mathcal{Z}_n^T] w_j - 2 \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n]^T W^T \mathcal{X}_n \right) + \text{Const.} \\
&= -\frac{1}{2\sigma_x^2} \sum_{n=1}^N \left(\sum_{j=1}^d w_j^T \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n \mathcal{Z}_n^T] w_j - 2 \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n]^T \left(\sum_{j=1}^d \mathcal{X}_{n,j} w_j \right) \right) + \text{Const.} \\
&= -\frac{1}{2\sigma_x^2} \sum_{j=1}^d \sum_{n=1}^N \left(w_j^T \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n \mathcal{Z}_n^T] w_j - 2 \mathcal{X}_{n,j} \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n]^T w_j \right) + \text{Const.} \\
&= -\frac{1}{2\sigma_x^2} \sum_{j=1}^d \left(w_j^T \left(\sum_{n=1}^N \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n \mathcal{Z}_n^T] \right) w_j - 2 w_j^T \left(\sum_{n=1}^N \mathcal{X}_{n,j} \mathbb{E}_{q_{int}(\mathcal{Q})} [\mathcal{Z}_n] \right) \right) + \text{Const.}
\end{aligned}$$

変分 M-step での更新は

$$\begin{aligned}
 & \log q_{i+1}(W) \\
 &= \log p(W) + \mathbb{E}_{q_{i+1}(q)}[\log p(\mathcal{X}|q, W)] + \text{const.} \\
 &= -\frac{1}{2} \sum_{j=1}^d \frac{1}{\sigma_w^2} w_j^T w_j - \frac{1}{2} \sum_{j=1}^d \left(w_j^T \left(\frac{1}{\sigma_x^2} \sum_{n=1}^N \mathbb{E}_{q_{i+1}(q)}[z_n z_n^T] \right) w_j - 2 w_j^T \left(\frac{1}{\sigma_x^2} \sum_{n=1}^N x_{n,j} \mathbb{E}_{q_{i+1}(q)}[z_n] \right) \right) + \text{const.} \\
 &= -\frac{1}{2} \sum_{j=1}^d \left(w_j^T \left(\frac{1}{\sigma_w^2} \mathbf{I}_d + \frac{1}{\sigma_x^2} \sum_{n=1}^N \mathbb{E}_{q_{i+1}(q)}[z_n z_n^T] \right) w_j - 2 w_j^T \left(\frac{1}{\sigma_x^2} \sum_{n=1}^N x_{n,j} \mathbb{E}_{q_{i+1}(q)}[z_n] \right) \right) + \text{const.} \\
 & \qquad \qquad \qquad =: \hat{\Sigma}_{w, i+1}^{-1} \qquad \qquad \qquad =: \hat{\Sigma}_{w, i+1}^{-1} \hat{\mu}_{w_j, i+1} \\
 &= -\frac{1}{2} \sum_{j=1}^d \left(w_j^T \hat{\Sigma}_{w, i+1}^{-1} w_j - 2 w_j^T \hat{\Sigma}_{w, i+1}^{-1} \hat{\mu}_{w_j, i+1} \right) + \text{const.} \\
 &= \sum_{j=1}^d \log \mathcal{N}(w_j | \hat{\mu}_{w_j, i+1}, \hat{\Sigma}_{w, i+1})
 \end{aligned}$$

つまり, $q_{i+1}(W) = \prod_{n=1}^d \mathcal{N}(w_j | \hat{\mu}_{w_j, i+1}, \hat{\Sigma}_{w, i+1})$

右辺は各 w_j の分布の積の形になっているので, $q_{i+1}(W) = \prod_{n=1}^d q_{i+1}(w_j)$ と独立な形に

書いて, 各 $j = 1, \dots, d$ について

$$\begin{aligned}
 q_{i+1}(w_j) &= \mathcal{N}(w_j | \hat{\mu}_{w_j, i+1}, \hat{\Sigma}_{w, i+1}) \\
 \hat{\mu}_{w_j, i+1} &= \frac{1}{\sigma_x^2} \hat{\Sigma}_{w, i+1} \sum_{n=1}^N x_{n,j} \mathbb{E}_{q_{i+1}(q)}[z_n] \\
 \hat{\Sigma}_{w, i+1} &= \left(\frac{1}{\sigma_w^2} \mathbf{I}_d + \frac{1}{\sigma_x^2} \sum_{n=1}^N \mathbb{E}_{q_{i+1}(q)}[z_n z_n^T] \right)^{-1}
 \end{aligned}$$

と更新すればよい.

結局, $q_i(z_n)$ も $q_i(w_j)$ も Gauss分布なので, 期待値のみ計算でき,

$$\begin{aligned}
 \mathbb{E}_{q_i(w)}[w_j] &= \mathbb{E}_{q_i(w_j)}[w_j] = \hat{\mu}_{w_j, i}, \\
 \mathbb{E}_{q_i(w)}[w_j w_j^T] &= \mathbb{E}_{q_i(w_j)}[w_j w_j^T] = \hat{\Sigma}_{w, i} + \hat{\mu}_{w_j, i} \hat{\mu}_{w_j, i}^T,
 \end{aligned}$$

$$\mathbb{E}_{q_{z_{n,l}}(z_n)}[z_n] = \mathbb{E}_{q_{z_{n,l}}(z_n)}[z_n] = \hat{\mu}_{z_{n,l}},$$

$$\mathbb{E}_{q_{z_{n,l}}(z_n)}[z_n z_n^T] = \mathbb{E}_{q_{z_{n,l}}(z_n)}[z_n z_n^T] = \hat{\Sigma}_{z_{n,l}} + \hat{\mu}_{z_{n,l}} \hat{\mu}_{z_{n,l}}^T.$$

最終的に、以下のように Gauss 分布のパラメータを更新すればよいことが分かった:

[変分 E-step]

$$\hat{\Sigma}_{z_{n,l}} = \left(I_d + \frac{d}{\sigma_x^2} \hat{\Sigma}_{w,l} + \frac{1}{\sigma_x^2} \sum_{j=1}^d \hat{\mu}_{w_{j,l}} \hat{\mu}_{w_{j,l}}^T \right)^{-1},$$

$$\hat{\mu}_{z_{n,l+1}} = \frac{1}{\sigma_x^2} \hat{\Sigma}_{z_{n,l}} \sum_{j=1}^d x_{n,j} \hat{\mu}_{w_{j,l}}.$$

[変分 M-step]

$$\hat{\Sigma}_{w,l+1} = \left(\frac{1}{\sigma_w^2} I_d + \frac{N}{\sigma_x^2} \hat{\Sigma}_{z_{n,l}} + \frac{1}{\sigma_x^2} \sum_{n=1}^N \hat{\mu}_{z_{n,l+1}} \hat{\mu}_{z_{n,l+1}}^T \right)^{-1},$$

$$\hat{\mu}_{w_{j,l+1}} = \frac{1}{\sigma_x^2} \hat{\Sigma}_{w,l+1} \sum_{n=1}^N x_{n,j} \hat{\mu}_{z_{n,l+1}}.$$

4.2.2.2 混合 Gauss 分布への適用.

- 連続な潜在変数 z の代わりに、離散の潜在変数

$$\mathcal{S} = \{s_1, \dots, s_N\} \in \{0,1\}^K, \quad \sum_{k=1}^K s_{n,k} = 1 \quad (n=1, \dots, N)$$

を割り当てることで、クラスタリングが可能である。

- 仮定.

$W = (w_1, \dots, w_K) \in M_{d,K}(\mathbb{R})$ ($w_j \in \mathbb{R}^d$) はパラメータ, $\sigma_x^2 > 0$ は定数とす.

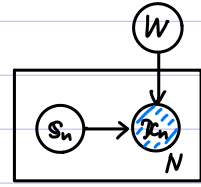
$$p(x_n | \mathcal{S}, W) = \prod_{n=1}^N p(x_n | s_n, W) = \prod_{n=1}^N \mathcal{N}(x_n | W s_n, \sigma_x^2 I_d).$$

\mathcal{S} の事前分布: $p(\mathcal{S}) = \prod_{n=1}^N \text{Cat}(s_n | \pi)$ (π : 定ベクトル)

W の事前分布: $\sigma_w^2 > 0$ とす $p(W) = \prod_{k=1}^K \mathcal{N}(w_k | 0, \sigma_w^2 I_d).$

モデル:

$$\begin{aligned}
 p(\mathcal{X}, \mathcal{S}, W) &= p(\mathcal{X} | \mathcal{S}, W) p(\mathcal{S}, W) \\
 &= p(\mathcal{X} | \mathcal{S}, W) p(\mathcal{S}) p(W)
 \end{aligned}$$



事後分布 $p(\mathcal{S}, W | \mathcal{X})$ に 2 次元近似:

$$p(\mathcal{S}, W | \mathcal{X}) \approx q(\mathcal{S}) q(W).$$

変分推論 (変分 EM algo.) を行う.

i 回目の更新後の近似分布を $q_i(\mathcal{S}), q_i(W)$ と書く.

変分 E-step. \mathcal{S} の近似分布を更新.

$$\begin{aligned}
 q_{i+1}(\mathcal{S}) &\propto \exp(\mathbb{E}_{q_i(W)}[\log p(\mathcal{X}, \mathcal{S} | W)]) \\
 &= p(\mathcal{S}) \exp(\mathbb{E}_{q_i(W)}[\log p(\mathcal{X} | \mathcal{S}, W)])
 \end{aligned}$$

変分 M-step. W の近似分布を更新.

$$\begin{aligned}
 q_{i+1}(W) &\propto p(W) \exp(\mathbb{E}_{q_{i+1}(\mathcal{S})}[\log p(\mathcal{X}, \mathcal{S} | W)]) \\
 &\propto p(W) \exp(\mathbb{E}_{q_{i+1}(\mathcal{S})}[\log p(\mathcal{X} | \mathcal{S}, W)])
 \end{aligned}$$

各 step の計算をもう少し進める.

$$\begin{aligned}
 &\log p(\mathcal{X} | \mathcal{S}, W) \\
 &= -\frac{1}{2} \sum_{n=1}^N \left(\frac{1}{\sigma_x^2} \mathcal{S}_n^T W^T W \mathcal{S}_n - \frac{2}{\sigma_x^2} \mathcal{X}_n^T W \mathcal{S}_n \right) + \text{const.} \\
 &= -\frac{1}{2} \sum_{n=1}^N \left(\frac{1}{\sigma_x^2} \left(\sum_{k=1}^K \mathcal{S}_{n,k} \mathcal{W}_k \right)^T \left(\sum_{k=1}^K \mathcal{S}_{n,k} \mathcal{W}_k \right) - \frac{2}{\sigma_x^2} \mathcal{X}_n^T \left(\sum_{k=1}^K \mathcal{S}_{n,k} \mathcal{W}_k \right) \right) + \text{const.} \\
 &= -\frac{1}{2} \sum_{n=1}^N \left(\frac{1}{\sigma_x^2} \sum_{k=1}^K \mathcal{S}_{n,k} \mathcal{W}_k^T \mathcal{W}_k - \frac{2}{\sigma_x^2} \mathcal{X}_n^T \left(\sum_{k=1}^K \mathcal{S}_{n,k} \mathcal{W}_k \right) \right) + \text{const.}
 \end{aligned}$$

$\downarrow \mathcal{S}_{n,k} \mathcal{S}_{n,k'} = \begin{cases} 0 & (k \neq k') \\ \mathcal{S}_{n,k} & (k = k') \end{cases}$

$$= -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \left(\frac{1}{\sigma_x^2} S_{n,k} \mathbf{w}_k^T \mathbf{w}_k - \frac{2}{\sigma_x^2} S_{n,k} \mathbf{g}_n^T \mathbf{w}_k \right) + \text{CONST}$$

∴

$$\begin{aligned} & \mathbb{E}_{q_i(w)} [\log p(\mathcal{X} | \mathcal{S}, w)] \\ &= -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \left(\frac{1}{\sigma_x^2} S_{n,k} \mathbb{E}_{q_i(w)} [\mathbf{w}_k^T \mathbf{w}_k] - \frac{2}{\sigma_x^2} S_{n,k} \mathbf{g}_n^T \mathbb{E}_{q_i(w)} [\mathbf{w}_k] \right) + \text{CONST.} \end{aligned}$$

$$\begin{aligned} & \mathbb{E}_{q_{i+1}(w)} [\log p(\mathcal{X} | \mathcal{S}, w)] \\ &= -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \left(\frac{1}{\sigma_x^2} \mathbb{E}_{q_{i+1}(w)} [S_{n,k}] \mathbf{w}_k^T \mathbf{w}_k - \frac{2}{\sigma_x^2} \mathbb{E}_{q_{i+1}(w)} [S_{n,k}] \mathbf{g}_n^T \mathbf{w}_k \right) + \text{CONST} \end{aligned}$$

変分 E-step 2'17,

$$\begin{aligned} & \log q_{i+1}(\mathcal{S}) \\ &= \log p(\mathcal{S}) + \mathbb{E}_{q_i(w)} [\log p(\mathcal{X} | \mathcal{S}, w)] + \text{CONST.} \\ &= \sum_{n=1}^N \log \text{Cat}(S_n | \boldsymbol{\pi}) - \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \left(\frac{1}{\sigma_x^2} S_{n,k} \mathbb{E}_{q_i(w)} [\mathbf{w}_k^T \mathbf{w}_k] - \frac{2}{\sigma_x^2} S_{n,k} \mathbf{g}_n^T \mathbb{E}_{q_i(w)} [\mathbf{w}_k] \right) \\ & \quad + \text{CONST.} \end{aligned}$$

$$= \sum_{n=1}^N \sum_{k=1}^K S_{n,k} \left(\log \pi_k - \frac{1}{2\sigma_x^2} \mathbb{E}_{q_i(w)} [\mathbf{w}_k^T \mathbf{w}_k] + \frac{1}{\sigma_x^2} \mathbf{g}_n^T \mathbb{E}_{q_i(w)} [\mathbf{w}_k] \right) + \text{CONST.}$$

$$= \sum_{n=1}^N \sum_{k=1}^K S_{n,k} \log \alpha_{n,k,i+1} + \text{CONST}$$

$$\hat{\pi}_{n,k,i+1} := \frac{\alpha_{n,k,i+1}}{\sum_{k'=1}^K \alpha_{n,k',i+1}}, \quad \hat{\boldsymbol{\pi}}_{n,i+1} := \begin{pmatrix} \hat{\pi}_{n,1,i+1} \\ \vdots \\ \hat{\pi}_{n,K,i+1} \end{pmatrix} \text{ と } \delta' < \epsilon,$$

$$\log q_{i+1}(\mathcal{S}) = \sum_{n=1}^N \log \text{Cat}(S_n | \hat{\boldsymbol{\pi}}_{n,i+1}) \Leftrightarrow q_{i+1}(\mathcal{S}) = \prod_{n=1}^N \text{Cat}(S_n | \hat{\boldsymbol{\pi}}_{n,i+1}).$$

右辺は各 S_n の分布の積の形なので, $q_{i+1}(\mathcal{S}) = \prod_{n=1}^N q_{i+1}(S_n)$ と独立な形に分解.

各 S_n ($n=1, \dots, N$) について,

$$\alpha_{n,k,i+1} = \exp\left(\log \pi_k - \frac{1}{2\sigma_x^2} \mathbb{E}_{q_i(w)}[\mathbf{w}_k^T \mathbf{w}_k] + \frac{1}{\sigma_x^2} \mathbf{g}_n^T \mathbb{E}_{q_i(w)}[\mathbf{w}_k]\right)$$

$$= \pi_k \exp\left(-\frac{1}{2\sigma_x^2} \mathbb{E}_{q_i(w)}[\mathbf{w}_k^T \mathbf{w}_k] + \frac{1}{\sigma_x^2} \mathbf{g}_n^T \mathbb{E}_{q_i(w)}[\mathbf{w}_k]\right) \quad (k=1, \dots, K)$$

$$\hat{\pi}_{n,k,i+1} = \frac{\alpha_{n,k,i+1}}{\sum_{k=1}^K \alpha_{n,k,i+1}} \quad (k=1, \dots, K), \quad \hat{\pi}_{n,i+1} = \begin{pmatrix} \hat{\pi}_{n,1,i+1} \\ \vdots \\ \hat{\pi}_{n,K,i+1} \end{pmatrix}$$

$$q_{i+1}(S_n) = \text{Cat}(S_n | \hat{\pi}_{n,i+1})$$

と更新すればよい。

変分 M-step では,

$$\log q_{i+1}(w)$$

$$= \log p(w) + \mathbb{E}_{q_{i+1}(s)}[\log p(s | \delta, w)] + \text{const.}$$

$$= -\frac{1}{2\sigma_w^2} \sum_{k=1}^K \mathbf{w}_k^T \mathbf{w}_k - \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \left(\frac{1}{\sigma_x^2} \mathbb{E}_{q_{i+1}(s)}[S_{n,k}] \mathbf{w}_k^T \mathbf{w}_k - \frac{2}{\sigma_x^2} \mathbb{E}_{q_{i+1}(s)}[S_{n,k}] \mathbf{g}_n^T \mathbf{w}_k \right)$$

+ const

$$= -\frac{1}{2} \sum_{k=1}^K \left(\left(\frac{1}{\sigma_w^2} + \frac{1}{\sigma_x^2} \sum_{n=1}^N \mathbb{E}_{q_{i+1}(s)}[S_{n,k}] \right) \mathbf{w}_k^T \mathbf{w}_k - 2 \mathbf{w}_k^T \left(\frac{1}{\sigma_x^2} \sum_{n=1}^N \mathbb{E}_{q_{i+1}(s)}[S_{n,k}] \mathbf{g}_n \right) \right)$$

+ const. $=: \hat{\sigma}_{k,i+1}^{-2}$ $=: \hat{\sigma}_{k,i+1}^{-2} \hat{\mu}_{k,i+1}$

$$= -\frac{1}{2} \sum_{k=1}^K \left(\mathbf{w}_k^T \left(\hat{\sigma}_{k,i+1}^2 \mathbf{I}_d \right)^{-1} \mathbf{w}_k - 2 \mathbf{w}_k^T \left(\hat{\sigma}_{k,i+1}^2 \mathbf{I}_d \right)^{-1} \hat{\mu}_{k,i+1} \right) + \text{const.}$$

$$= \sum_{k=1}^K \mathcal{N}(\mathbf{w}_k | \hat{\mu}_{k,i+1}, \hat{\sigma}_{k,i+1}^2 \mathbf{I}_d).$$

$$\Leftrightarrow q_{i+1}(w) = \prod_{k=1}^K \mathcal{N}(\mathbf{w}_k | \hat{\mu}_{k,i+1}, \hat{\sigma}_{k,i+1}^2 \mathbf{I}_d).$$

右辺は各 \mathbf{w}_k の分布の積の形なので, $q_{i+1}(w) = \prod_{k=1}^K q_{i+1}(\mathbf{w}_k)$ と独立な形に PVTC.

各 \mathbf{w}_k ($k=1, \dots, K$) に対して,

$$q_{i+1}(w_k) = \mathcal{N}(w_k | \hat{\mu}_{k,i+1}, \hat{\sigma}_{k,i+1}^2 I_d),$$

$$\hat{\sigma}_{k,i+1}^2 = \left(\frac{1}{\sigma_w^2} + \frac{1}{\sigma_x^2} \sum_{n=1}^N \mathbb{E}_{q_{i+1}(w)} [S_{n,k}] \right)^{-1}$$

$$\hat{\mu}_{k,i+1} = \frac{\hat{\sigma}_{k,i+1}^2}{\sigma_x^2} \sum_{n=1}^N \mathbb{E}_{q_{i+1}(w)} [S_{n,k}] q_n$$

と更新すればよい.

分布の形も分かたないので期待値を計算する.

$$\mathbb{E}_{q_i(w)} [w_k] = \mathbb{E}_{q_i(w_k)} [w_k] = \hat{\mu}_{k,i}$$

$\mathcal{N}(w_k | \hat{\mu}_{k,i}, \hat{\sigma}_{k,i}^2 I_d)$ のとき

w_k の各成分 $w_{k,j}$ はそれぞれ

$\mathcal{N}(\hat{\mu}_{k,i,j}, \hat{\sigma}_{k,i}^2)$ に従う.

$$\begin{aligned} \mathbb{E}_{q_i(w)} [w_k^T w_k] &= \mathbb{E}_{q_i(w_k)} [w_k^T w_k] = \sum_{j=1}^d \mathbb{E}_{q_i(w_k)} [w_{k,j}^2] \\ &= \sum_{j=1}^d (\hat{\sigma}_{k,i}^2 + \hat{\mu}_{k,i,j}^2) = d \hat{\sigma}_{k,i}^2 + \hat{\mu}_{k,i}^T \hat{\mu}_{k,i}. \end{aligned}$$

$$\mathbb{E}_{q_{i+1}(w)} [S_{n,k}] = \mathbb{E}_{q_{i+1}(S_n)} [S_{n,k}] = \hat{\pi}_{n,k,i+1}.$$

以上より、パラメータの更新式は次のとおり.

[変分 E-step]

$$\alpha_{n,k,i+1} = \pi_k \exp \left(-\frac{1}{2\sigma_x^2} (d \hat{\sigma}_{k,i}^2 + \hat{\mu}_{k,i}^T \hat{\mu}_{k,i}) + \frac{1}{\sigma_x^2} q_n^T \hat{\mu}_{k,i} \right)$$

$$\hat{\pi}_{n,k,i+1} = \frac{\alpha_{n,k,i+1}}{\sum_{k=1}^K \alpha_{n,k,i+1}}, \quad \hat{\pi}_{n,i+1} = \begin{pmatrix} \hat{\pi}_{n,1,i+1} \\ \vdots \\ \hat{\pi}_{n,K,i+1} \end{pmatrix}.$$

[変分 M-step]

$$\hat{\sigma}_{k,i+1}^2 = \left(\frac{1}{\sigma_w^2} + \frac{1}{\sigma_x^2} \sum_{n=1}^N \hat{\pi}_{n,k,i+1} \right)^{-1},$$

$$\hat{\mu}_{k,i+1} = \frac{\hat{\sigma}_{k,i+1}^2}{\sigma_x^2} \sum_{n=1}^N \hat{\pi}_{n,k,i+1} q_n.$$

4.2.3 Laplace近似

・ Laplace近似 (Laplace approximation):
事後分布 $p(\mathbf{z}|\mathbf{y})$ 近似

複雑な分布をより簡単な Gauss 分布を使って近似的に表現する手法.

・ 分布 $p(\mathbf{z})$ のモード (mode): $p(\mathbf{z})$ の極大値を与える点.

・ \mathbf{z}_0 : 近似したい分布 $p(\mathbf{z})$ の任意のモード.

$$\text{Gauss 分布では, } \log \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{1}{2}(\mathbf{z}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{z}-\boldsymbol{\mu}) + \text{const.}$$

の形にならているので, $\log p(\mathbf{z})$ を \mathbf{z}_0 あたりで Taylor 展開してみる.

$$\begin{aligned} & \log p(\mathbf{z}) \\ &= \underbrace{\log p(\mathbf{z}_0)}_{\text{const.}} + \left(\nabla_{\mathbf{z}} \log p(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} \right)^T (\mathbf{z}-\mathbf{z}_0) + \frac{1}{2} (\mathbf{z}-\mathbf{z}_0)^T \left(\nabla_{\mathbf{z}}^2 \log p(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} \right) (\mathbf{z}-\mathbf{z}_0) \\ & \quad + O(\|\mathbf{z}-\mathbf{z}_0\|^3) \end{aligned}$$

$$\approx \left(\nabla_{\mathbf{z}} \log p(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} \right)^T (\mathbf{z}-\mathbf{z}_0) + \frac{1}{2} (\mathbf{z}-\mathbf{z}_0)^T \left(\nabla_{\mathbf{z}}^2 \log p(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} \right) (\mathbf{z}-\mathbf{z}_0) + \text{const.}$$

\mathbf{z}_0 は $p(\mathbf{z})$ のモードであり, $\log p(\mathbf{z})$ の極大値でもあるので,

$$\nabla_{\mathbf{z}} \log p(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} = \mathbf{0}.$$

ゆえに,

$$\log p(\mathbf{z}) \approx \frac{1}{2} (\mathbf{z}-\mathbf{z}_0)^T \left(\nabla_{\mathbf{z}}^2 \log p(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} \right) (\mathbf{z}-\mathbf{z}_0) + \text{const.}$$

よって, $\Lambda(\mathbf{z}_0) := -\nabla_{\mathbf{z}}^2 \log p(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0}$ とおけば,
 ← \mathbf{z}_0 がモードならば, $\Lambda(\mathbf{z}_0)$ は正定値行列.

$$\log p(\mathbf{z}) = -\frac{1}{2} (\mathbf{z}-\mathbf{z}_0)^T \Lambda(\mathbf{z}_0) (\mathbf{z}-\mathbf{z}_0) + \text{const.}$$

$\Leftrightarrow p(\mathbf{z}) \approx C \exp\left(-\frac{1}{2} (\mathbf{z}-\mathbf{z}_0)^T \Lambda(\mathbf{z}_0) (\mathbf{z}-\mathbf{z}_0)\right) = \mathcal{N}(\mathbf{z}|\mathbf{z}_0, \Lambda(\mathbf{z}_0)^{-1})$ と近似できる.

- ・ 適当な最適化手法で $p(\mathbf{x})$ のモード \mathbf{x}_0 を求め, $\Lambda(\mathbf{x}_0)$ を計算することで近似 $\mathcal{N}(\mathbf{x} | \mathbf{x}_0, \Lambda(\mathbf{x}_0)^{-1})$ を得る.

Remark ・ 多峰的 (multimodal) な分布では, モード \mathbf{x}_0 の選び方により得られる近似が異なる.

- ・ モード \mathbf{x}_0 近傍のみを考慮した近似なので, 分布全体を近似できてはいない.
- ・ Hesse 行列の計算をする必要がある.
- ・ $p(\mathbf{x})$ が実変数の分布でないと, 直接適用できない.

4.2.4 モーメントマッチングによる近似.

4.2.4.1 モーメントマッチング.

- ・ 分布 $p(\mathbf{x})$ を, より簡単な分布 $q(\mathbf{x})$ で近似する.

[仮定] $q(\mathbf{x})$ は指数型分布族, i.e.,

$$q(\mathbf{x}; \eta) = h(\mathbf{x}) \exp(\eta^T t(\mathbf{x}) - a(\eta)).$$

- ・ $p \in q$ で近似するにあ, p の q に対する KL-divergence $D_{KL}[p(\mathbf{x}) \| q(\mathbf{x}; \eta)]$ (forward KL-divergence) を最小化する.

Remark ・ 変分推論では reverse KL-divergence $D_{KL}[q(\mathbf{x}; \eta) \| p(\mathbf{x})]$ を最小化した.

今日は, 「 p を真の分布とみたときの, q の p からの解離度」を最小化する.

- ・ 解く問題は $\eta^* = \underset{\eta \in \mathcal{H}}{\operatorname{argmin}} D_{KL}[p(\mathbf{x}) \| q(\mathbf{x}; \eta)]$.

但し, $\mathcal{H} = \{ \eta \mid \exp(a(\eta)) = \int \exp(\eta^T t(\mathbf{x})) h(\mathbf{x}) d\mathbf{x} < \infty \}$

← q が分布となるような範囲.

$D_{KL}[p(\mathbf{z}) \parallel q(\mathbf{z}; \eta)]$ を整理してみよう.

$$\begin{aligned} D_{KL}[p(\mathbf{z}) \parallel q(\mathbf{z}; \eta)] &= \int p(\mathbf{z}) \log \frac{p(\mathbf{z})}{q(\mathbf{z}; \eta)} d\mathbf{z} \\ &= \int p(\mathbf{z}) \log p(\mathbf{z}) d\mathbf{z} - \int p(\mathbf{z}) \log q(\mathbf{z}; \eta) d\mathbf{z} \\ &= -\mathbb{E}_p[\log q(\mathbf{z}; \eta)] + \underbrace{\mathbb{E}_p[\log p(\mathbf{z})]}_{\text{const.}} \\ &= -\mathbb{E}_p[\log h(\mathbf{z}) + \eta^T \mathbf{t}(\mathbf{z}) - a(\eta)] + \text{const.} \\ &= \underbrace{-\mathbb{E}_p[\log h(\mathbf{z})]}_{\text{const.}} - \eta^T \mathbb{E}_p[\mathbf{t}(\mathbf{z})] + a(\eta) + \text{const.} \\ &= -\eta^T \mathbb{E}_p[\mathbf{t}(\mathbf{z})] + a(\eta) + \text{const.} \end{aligned}$$

よって, $\eta^* = \underset{\eta}{\operatorname{argmin}} \left(-\eta^T \mathbb{E}_p[\mathbf{t}(\mathbf{z})] + a(\eta) \right)$.

[Claim] $-\eta^T \mathbb{E}_p[\mathbf{t}(\mathbf{z})] + a(\eta)$ は η により凸関数.

pf. $-\eta^T \mathbb{E}_p[\mathbf{t}(\mathbf{z})]$ は η により線形なため凸関数.

$\nabla_{\eta}^2 a(\eta) = \nabla_{\eta}^2 [\log Z(\eta)] > 0$ なるため凸関数. 凸関数の和はやはり凸関数. \square

[Claim] 自然パラメータ η の動く空間

$$\mathcal{H} = \left\{ \eta \mid \exp(a(\eta)) = \int \exp(\eta^T \mathbf{t}(\mathbf{z})) h(\mathbf{z}) d\mathbf{z} < \infty \right\}$$

は凸集合.

pf. $\forall \eta_1, \eta_2 \in \mathcal{H}, \forall \alpha \in (0, 1)$ とし, $\alpha \eta_1 + (1-\alpha) \eta_2 \in \mathcal{H}$ を示す.

$$\begin{aligned} &\int \exp((\alpha \eta_1 + (1-\alpha) \eta_2)^T \mathbf{t}(\mathbf{z})) h(\mathbf{z}) d\mathbf{z} \\ &= \int \exp(\alpha \eta_1^T \mathbf{t}(\mathbf{z})) \exp((1-\alpha) \eta_2^T \mathbf{t}(\mathbf{z})) h(\mathbf{z}) d\mathbf{z} \\ &= \int \left(\exp(\eta_1^T \mathbf{t}(\mathbf{z})) \right)^\alpha \left(\exp(\eta_2^T \mathbf{t}(\mathbf{z})) \right)^{1-\alpha} h(\mathbf{z}) d\mathbf{z} \end{aligned}$$

$$\leq \left(\int \exp(\eta_1^T \tau(\mathbf{z})) h(\mathbf{z}) d\mathbf{z} \right)^\alpha \left(\int \exp(\eta_2^T \tau(\mathbf{z})) h(\mathbf{z}) d\mathbf{z} \right)^{1-\alpha} \quad (\because \text{Hölderの不等式})$$

$< \infty$.



以上より、この問題は凸最適化問題なので、極値が最小値を与える。

$$\therefore -\mathbb{E}_{p(\mathbf{z})}[\tau(\mathbf{z})] + \nabla_{\eta} a(\eta) \Big|_{\eta=\eta^*} = 0. \quad \left. \begin{array}{l} \\ \end{array} \right\} \nabla_{\eta} a(\eta) = \mathbb{E}_q[\tau(\mathbf{z})]$$

$$\Leftrightarrow \mathbb{E}_{q(\mathbf{z}; \eta^*)}[\tau(\mathbf{z})] = \mathbb{E}_{p(\mathbf{z})}[\tau(\mathbf{z})]$$

結局、 q を指数型分布族に制限した際の最適な q は、十分統計量 $\tau(\mathbf{z})$ の

p による平均と q による平均が一致するようなパラメータ η^* を選んだときの $q(\mathbf{z}; \eta^*)$ になる。

このようにして近似する手法を **モーメントマッチング** (moment matching) という。

Remark $\cdot p(\mathbf{z})$ のモーメントとは $\mathbf{z} \sim p(\mathbf{z})$ のときの $\mathbb{E}_p[\mathbf{z}]$, $\mathbb{E}_p[\mathbf{z}\mathbf{z}^T]$ などの量という

$\mathbb{E}_{p(\mathbf{z})}[\tau(\mathbf{z})]$ は、 $p(\mathbf{z})$ のモーメントではないのだが、指数型分布族 $q(\mathbf{z}; \eta)$ に対し、

$\mu = \mathbb{E}_{q(\mathbf{z}; \eta^*)}[\tau(\mathbf{z})]$ のことを「モーメントパラメータ」と呼ぶ。実際、 $\tau = \tau(\mathbf{z})$ の分布は

$q(\tau; \eta) = \tilde{h}(\tau) \exp(\eta^T \tau - a(\eta))$ であり、 $\mathbb{E}_{q(\tau; \eta)}[\tau] = \nabla_{\eta} a(\eta) = \mathbb{E}_{q(\mathbf{z}; \eta)}[\tau(\mathbf{z})]$ となる。
(cf.) 百田『数理統計学』

μ は $q(\tau; \eta)$ の1次のモーメント。「モーメントマッチング」とはモーメントパラメータをあわせること、と考えられる。

4.2.4.2 仮定密度フィルタリング

制御分野での呼び方もと。

"online Bayesian learning", "weak marginalization" など。

仮定密度フィルタリング (assumed density filtering)

データ $\mathcal{D}_1, \mathcal{D}_2, \dots$ を逐次的に学習する。

パラメータ θ の事前分布を $p(\theta)$ とすると、データ \mathcal{D}_1 を観測後の事後分布は、

$p(\theta | \mathcal{D}_1) \propto p(\mathcal{D}_1 | \theta) p(\theta)$ となる。尤度関数 $p(\mathcal{D}_1 | \theta)$ と $p(\theta)$ が共役なら、

$$p(\theta | \mathcal{D}_1), p(\theta | \mathcal{D}_1, \mathcal{D}_2), p(\theta | \mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3), \dots$$

は全て $p(\theta)$ と同じ形の分布で解析的に計算できる。

- ・ 共役でない状況を考える。 $p(\theta | \mathcal{D}_1)$ に対する近似分布として、 $p(\theta)$ と同じ分布になるよう $q_1(\theta)$ を設定し、

$$q_1(\theta) \approx r_1(\theta) = \frac{1}{Z_1} p(\mathcal{D}_1 | \theta) p(\theta), \quad Z_1 = \int p(\mathcal{D}_1 | \theta) p(\theta) d\theta.$$

と近似する。 q_1 は $q_1(\theta) = \operatorname{argmin}_{q_1} D_{KL}[r_1(\theta) \| q_1(\theta)]$ により決定する。

p が指数型分布族のとき、これはモーメントマッチングでできる。

以降も、

$$q_i(\theta) \approx r_i(\theta) = \frac{1}{Z_i} p(\mathcal{D}_i | \theta) q_{i-1}(\theta), \quad Z_i = \int p(\mathcal{D}_i | \theta) q_{i-1}(\theta) d\theta,$$

$$q_i(\theta) = \operatorname{argmin}_{q_i} D_{KL}[r_i(\theta) \| q_i(\theta)]$$

と近似していく。こうすることで、 $q_i(\theta)$ は同じ分布を保ったまま逐次的に

事後分布を更新していく。

以降 $f_i(\theta) = p(\mathcal{D}_i | \theta)$ と書く。

4.2.4.3 1次元 Gauss 分布の例。 $\mathcal{N}(\theta | \mu_i, \nu_i) = \exp\left(-\frac{1}{2\nu_i} \theta^2 + \frac{\mu_i}{\nu_i} \theta - \frac{\mu_i^2}{2\nu_i} - \frac{1}{2} \log 2\pi\nu_i\right)$

- ・ $q_i(\theta) = \mathcal{N}(\theta | \mu_i, \nu_i)$ の場合。これは指数型分布族なのでモーメントマッチングでできる。

十分統計量は $\pi(\theta) = \begin{pmatrix} \theta \\ \theta^2 \end{pmatrix}$ 。 $r_{i+1}(\theta)$ の正規化定数は、

$$Z_{i+1} = \int f_{i+1}(\theta) q_i(\theta) d\theta.$$

$$= \int f_{i+1}(\theta) \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1}{2\nu_i} (\theta - \mu_i)^2\right) d\theta.$$

まず,

$$\begin{aligned} & \frac{\partial}{\partial \mu_i} \log Z_{i+1} \\ &= \frac{1}{Z_{i+1}} \int f_{i+1}(\theta) \frac{1}{\sqrt{2\pi v_i}} \exp\left(-\frac{1}{2v_i}(\theta - \mu_i)^2\right) \left(\frac{\theta - \mu_i}{v_i}\right) d\theta \\ &= \frac{1}{v_i} \int r_{i+1}(\theta) (\theta - \mu_i) d\theta \\ &= \frac{1}{v_i} \left(\mathbb{E}_{r_{i+1}(\theta)}[\theta] - \mu_i \right). \\ \therefore \mathbb{E}_{r_{i+1}(\theta)}[\theta] &= \mu_i + v_i \frac{\partial}{\partial \mu_i} \log Z_{i+1}. \end{aligned}$$

次に,

$$\begin{aligned} & \frac{\partial}{\partial v_i} \log Z_{i+1} \\ &= \frac{1}{Z_{i+1}} \int f_{i+1}(\theta) \frac{1}{\sqrt{2\pi}} \left(-\frac{1}{2} \frac{1}{v_i \sqrt{v_i}} + \frac{1}{\sqrt{v_i}} \frac{(\theta - \mu_i)^2}{2v_i^2} \right) \exp\left(-\frac{1}{2v_i}(\theta - \mu_i)^2\right) d\theta \\ &= \int r_{i+1}(\theta) \left(-\frac{1}{2v_i} + \frac{(\theta - \mu_i)^2}{2v_i^2} \right) d\theta \\ &= \frac{1}{2v_i^2} \int r_{i+1}(\theta) (-v_i + \theta^2 - 2\mu_i \theta + \mu_i^2) d\theta \\ &= -\frac{1}{2v_i} + \frac{1}{2v_i^2} \left(\mathbb{E}_{r_{i+1}(\theta)}[\theta^2] - 2\mu_i \mathbb{E}_{r_{i+1}(\theta)}[\theta] + \mu_i^2 \right). \\ \therefore \mathbb{E}_{r_{i+1}(\theta)}[\theta^2] &= v_i - \mu_i^2 + 2\mu_i \mathbb{E}_{r_{i+1}(\theta)}[\theta] + 2v_i^2 \frac{\partial}{\partial v_i} \log Z_{i+1}. \end{aligned}$$

モメントマップにより,

$$\mathbb{E}_{q_{i+1}(\theta)}[\theta] = \mathbb{E}_{r_{i+1}(\theta)}[\theta], \quad \mathbb{E}_{q_{i+1}(\theta)}[\theta^2] = \mathbb{E}_{r_{i+1}(\theta)}[\theta^2]$$

とわかる。よって,

$$\mathbb{E}_{q_{i+1}(\theta)}[\theta] = \mu_{i+1} \text{ かつ } \mu_{i+1} = \mu_i + v_i \frac{\partial}{\partial \mu_i} \log Z_{i+1}.$$

$$\mathbb{E}_{q_{i+1}(\theta)}[\theta^2] = v_{i+1} + \mu_{i+1}^2 \text{ かつ } v_{i+1} = v_i - \mu_i^2 + 2\mu_i \mu_{i+1} + 2v_i^2 \frac{\partial}{\partial v_i} \log Z_{i+1}.$$

$$\begin{aligned}
 V_{i+1} &= V_i - (\mu_{i+1} - \mu_i)^2 + 2V_i^2 \frac{\partial}{\partial V_i} \log Z_{i+1} \\
 &= V_i - V_i^2 \left(\left(\frac{\partial}{\partial \mu_i} \log Z_{i+1} \right)^2 - 2 \frac{\partial}{\partial V_i} \log Z_{i+1} \right).
 \end{aligned}$$

各正規化定数 Z_i が計算できれば、このように逐次更新ができる。

4.2.4.4 ガンマ分布の例.

$$\begin{aligned}
 q_i(\theta) &= \text{Gam}(\theta | a_i, b_i) = \frac{b_i^{a_i}}{\Gamma(a_i)} \theta^{a_i-1} e^{-b_i \theta} \\
 &= \exp((a_i-1) \log \theta - b_i \theta + a_i \log b_i - \log \Gamma(a_i))
 \end{aligned}$$

の場合. 指数型分布族. 十分統計量は $\eta(\theta) = \begin{pmatrix} \log \theta \\ \theta \end{pmatrix}$.

$f_{i+1}(\theta)$ の正規化定数は,

$$Z_{i+1}(a_i, b_i) = \int f_{i+1}(\theta) \text{Gam}(\theta | a_i, b_i) d\theta.$$

$$\mathbb{E}_{f_{i+1}(\theta)}[\theta]$$

$$\begin{aligned}
 &= \frac{1}{Z_{i+1}(a_i, b_i)} \int \theta f_{i+1}(\theta) \frac{b_i^{a_i}}{\Gamma(a_i)} \theta^{a_i-1} e^{-b_i \theta} d\theta \\
 &= \frac{1}{Z_{i+1}(a_i, b_i)} \int f_{i+1}(\theta) \frac{b_i^{a_i}}{\Gamma(a_i)} \theta^{a_i} e^{-b_i \theta} d\theta \\
 &= \frac{1}{Z_{i+1}(a_i, b_i)} \frac{\Gamma(a_i+1)}{b_i \Gamma(a_i)} Z_{i+1}(a_i+1, b_i) \quad \left. \begin{array}{l} \Gamma(a_i+1) \\ = a_i \Gamma(a_i) \end{array} \right\} \\
 &= \frac{a_i Z_{i+1}(a_i+1, b_i)}{b_i Z_{i+1}(a_i, b_i)}.
 \end{aligned}$$

$$\mathbb{E}_{f_{i+1}(\theta)}[\log \theta]$$

$$\begin{aligned}
 &= \frac{1}{Z_{i+1}(a_i, b_i)} \int \log \theta f_{i+1}(\theta) \frac{b_i^{a_i}}{\Gamma(a_i)} \theta^{a_i-1} e^{-b_i \theta} d\theta \\
 &= \frac{1}{Z_{i+1}(a_i, b_i)} \int f_{i+1}(\theta) \frac{b_i^{a_i}}{\Gamma(a_i)} \left(\frac{\partial}{\partial a_i} \theta^{a_i-1} \right) e^{-b_i \theta} d\theta \\
 &= \frac{1}{Z_{i+1}(a_i, b_i)} \frac{b_i^{a_i}}{\Gamma(a_i)} \frac{\partial}{\partial a_i} \int f_{i+1}(\theta) \theta^{a_i-1} e^{-b_i \theta} d\theta
 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{Z_{i+1}(a_i, b_i)} \frac{b_i^{a_i}}{\Gamma(a_i)} \frac{\partial}{\partial a_i} \left(Z_{i+1}(a_i, b_i) \frac{\Gamma(a_i)}{b_i^{a_i}} \right) \\
&= \frac{\partial}{\partial a_i} \log \left(Z_{i+1}(a_i, b_i) \frac{\Gamma(a_i)}{b_i^{a_i}} \right) \\
&= \frac{\partial}{\partial a_i} \log Z_{i+1}(a_i, b_i) + \frac{d}{da_i} \log \Gamma(a_i) - \log b_i.
\end{aligned}$$

ここで、ディガンマ関数 $\psi(a) = \frac{d}{da} \log \Gamma(a)$ を使うと、

$$E_{q_{i+1}(\theta)}[\log \theta] = \frac{\partial}{\partial a_i} \log Z_{i+1}(a_i, b_i) + \psi(a_i) - \log b_i.$$

モーメントマッチングをすると、

$$E_{q_{i+1}(\theta)}[\theta] = \frac{a_{i+1}}{b_{i+1}} = \frac{a_i Z_{i+1}(a_i+1, b_i)}{b_i Z_{i+1}(a_i, b_i)}$$

$$E_{q_{i+1}(\theta)}[\log \theta] = \psi(a_{i+1}) - \log b_{i+1} = \frac{\partial}{\partial a_i} \log Z_{i+1}(a_i, b_i) + \psi(a_i) - \log b_i.$$

となる。この更新は数値的に行うしかない。

この妥当性は不明。

→ 必ず KL-divergence は最小にはならない。

解析的な更新式を得るため、 $\log \theta$ の代わりに θ^2 の平均をあわせることにする。

$$\begin{aligned}
&E_{q_{i+1}(\theta)}[\theta^2] \\
&= \frac{1}{Z_{i+1}(a_i, b_i)} \int \theta^2 f_{i+1}(\theta) \frac{b_i^{a_i}}{\Gamma(a_i)} \theta^{a_i-1} e^{-b_i \theta} d\theta \\
&= \frac{1}{Z_{i+1}(a_i, b_i)} \int f_{i+1}(\theta) \frac{b_i^{a_i}}{\Gamma(a_i)} \theta^{a_i+1} e^{-b_i \theta} d\theta \\
&= \frac{1}{Z_{i+1}(a_i, b_i)} \frac{\Gamma(a_i+2)}{b_i^2 \Gamma(a_i)} Z_{i+1}(a_i+2, b_i) \\
&= \frac{(a_i+1) a_i Z_{i+1}(a_i+2, b_i)}{b_i^2 Z_{i+1}(a_i, b_i)}.
\end{aligned}$$

更新のためには、以下のようになる。

$$E_{q_{i+1}(\theta)}[\theta] = \frac{a_{i+1}}{b_{i+1}} = \frac{a_i Z_{i+1}(a_i+1, b_i)}{b_i Z_{i+1}(a_i, b_i)}.$$

$$E_{q_{i+1}(\theta)}[\theta^2] = \frac{a_{i+1}^2}{b_{i+1}^2} = \frac{(a_i+1) a_i Z_{i+1}(a_i+2, b_i)}{b_i^2 Z_{i+1}(a_i, b_i)} - \left(\frac{a_i Z_{i+1}(a_i+1, b_i)}{b_i Z_{i+1}(a_i, b_i)} \right)^2.$$

$$a_{i+1} = \frac{a_i Z_{i+1}(a_i+1, b_i)}{b_i Z_{i+1}(a_i, b_i)} \quad b_{i+1} \in \text{第2式に代入},$$

$$\frac{1}{b_{i+1}} \frac{a_i Z_{i+1}(a_i+1, b_i)}{b_i Z_{i+1}(a_i, b_i)} = \frac{(a_i+1) a_i Z_{i+1}(a_i+2, b_i)}{b_i^2 Z_{i+1}(a_i, b_i)} - \left(\frac{a_i Z_{i+1}(a_i+1, b_i)}{b_i Z_{i+1}(a_i, b_i)} \right)^2$$

$$\frac{1}{b_{i+1}} a_i Z_{i+1}(a_i+1, b_i) = \frac{(a_i+1) a_i Z_{i+1}(a_i+2, b_i)}{b_i} - \frac{a_i^2 Z_{i+1}(a_i+1, b_i)^2}{b_i Z_{i+1}(a_i, b_i)}$$

$$\frac{1}{b_{i+1}} = \frac{(a_i+1) Z_{i+1}(a_i+2, b_i)}{b_i Z_{i+1}(a_i+1, b_i)} - \frac{a_i Z_{i+1}(a_i+1, b_i)}{b_i Z_{i+1}(a_i, b_i)}$$

$$= \frac{1}{b_i} \left(\frac{(a_i+1) Z_{i+1}(a_i+2, b_i)}{Z_{i+1}(a_i+1, b_i)} - \frac{a_i Z_{i+1}(a_i+1, b_i)}{Z_{i+1}(a_i, b_i)} \right)$$

$$\therefore b_{i+1} = b_i \left(\frac{(a_i+1) Z_{i+1}(a_i+2, b_i)}{Z_{i+1}(a_i+1, b_i)} - \frac{a_i Z_{i+1}(a_i+1, b_i)}{Z_{i+1}(a_i, b_i)} \right)^{-1}$$

$$a_{i+1} = a_i \frac{Z_{i+1}(a_i+1, b_i)}{Z_{i+1}(a_i, b_i)} \left(\frac{(a_i+1) Z_{i+1}(a_i+2, b_i)}{Z_{i+1}(a_i+1, b_i)} - \frac{a_i Z_{i+1}(a_i+1, b_i)}{Z_{i+1}(a_i, b_i)} \right)^{-1}$$

$$= a_i \left(\frac{(a_i+1) Z_{i+1}(a_i+2, b_i) Z_{i+1}(a_i, b_i)}{Z_{i+1}(a_i+1, b_i)^2} - a_i \right)^{-1}$$

各正規化定数 $Z(a_i, b_i)$ が計算できれば, このような更新ができる.

4.2.5 例: モーメントマッチングによるプロビット回帰モデルの学習.

← probability unit に関連

プロビット回帰 (probit regression) モデル:

2値分類のための一般化線形モデルの一つ. モデルによる出力値は分類確率.

リンク関数は **プロビット関数** (probit func.):

$$\text{probit}(p) = \Phi^{-1}(p), \quad p \in (0, 1) \quad \leftarrow \Phi \text{ は標準正規分布の cdf.}$$

cf.) ロジスティック回帰も同様のモデル. こちらはリンク関数が **ロジット関数** (logit func.)

← ロジスティック関数の逆関数 「対数オッズ比」ともいう.

$$\text{logit}(p) = \log \frac{p}{1-p}, \quad p \in (0, 1)$$

の一般化線形モデル.

ロジスティック回帰よりもプロビット回帰の方が外れ値に敏感らしい. (PRML Ch 4.3)

入力: $\mathcal{X} = \{x_1, \dots, x_N\} \in \mathcal{R}$, ← 多次元でもできる

出力: $\mathcal{Y} = \{y_1, \dots, y_N\}$

プロビットモデル: $p(y|x, w) = \Phi(ywx)$. ($w \in \mathcal{R}$ はパラメータ)

尤度関数: $p(\mathcal{Y}|\mathcal{X}, w) = \prod_{n=1}^N p(y_n|x_n, w) = \prod_{n=1}^N \Phi(y_n w x_n)$.

w の事前分布 $p(w) = \mathcal{N}(w|0, v_0)$ ($v_0 > 0$ はパラメータ)

モデルの周辺尤度は

$$p(\mathcal{Y}|\mathcal{X}) = \int p(\mathcal{Y}|\mathcal{X}, w) p(w) dw.$$

これは解析的に計算できないので、 $\epsilon \rightarrow 0$ まで追加していく

仮定密度フィルタリングをする。

$f_i(w) = p(y_i|x_i, w)$ とし、パラメータ w の事後分布を

$$q_i(w) = \mathcal{N}(w|\mu_i, v_i)$$

で近似する。 $\mu_0 = 0$, $q_0(w) = p(w) = \mathcal{N}(w|\mu_0, v_0)$ とする。

i ステップ目の更新は

$$q_{i+1}(w) \approx \frac{1}{Z_{i+1}} f_{i+1}(w) q_i(w)$$

をモーメントマッチングで行えばよい。

正規化定数

$$Z_{i+1} = \int p(y_{i+1}|x_{i+1}, w) \mathcal{N}(w|\mu_i, v_i) dw$$

は解析的に計算することからできる。

• $x_{i+1} = 0$ のとき, $p(y_{i+1} | x_{i+1}, w) = \Phi(0) = \frac{1}{2}$.

このとき $Z_{i+1} = \int \Phi(0) \mathcal{N}(w | \mu_i, \nu_i) dw = \Phi(0) = \frac{1}{2}$.

• $x_{i+1} \neq 0$ のとき, $y_{i+1} \in \{\pm 1\}$ と仮定し $x_{i+1} y_{i+1} \neq 0$.

$$Z_{i+1}$$

$$= \int_{-\infty}^{\infty} \Phi(y_{i+1} w x_{i+1}) \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1}{2\nu_i} (w - \mu_i)^2\right) dw$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1}{2\nu_i} (w - \mu_i)^2\right) \left(\int_{-\infty}^{y_{i+1} w x_{i+1}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} t^2\right) dt \right) dw$$

$s = \frac{t}{x_{i+1} y_{i+1}}$ とおく.

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1}{2\nu_i} (w - \mu_i)^2\right) \left(\int_{-\infty}^w \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} x_{i+1}^2 y_{i+1}^2 s^2\right) x_{i+1} y_{i+1} ds \right) dw$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1}{2\nu_i} u^2\right) \left(\int_{-\infty}^{u+\mu_i} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} x_{i+1}^2 y_{i+1}^2 s^2\right) x_{i+1} y_{i+1} ds \right) du$$

$u = w - \mu_i$ とおく.

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1}{2\nu_i} u^2\right) \left(\int_{-\infty}^{\mu_i} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} x_{i+1}^2 y_{i+1}^2 (u+v)^2\right) x_{i+1} y_{i+1} dv \right) du$$

積分順序の変更

$$= \int_{-\infty}^{\mu_i} \frac{1}{\sqrt{2\pi}} x_{i+1} y_{i+1} \left(\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1}{2\nu_i} u^2\right) \exp\left(-\frac{1}{2} x_{i+1}^2 y_{i+1}^2 (u+v)^2\right) du \right) dv$$

$$= \int_{-\infty}^{\mu_i} \frac{x_{i+1} y_{i+1}}{\sqrt{2\pi}} \left(\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{u^2}{2\nu_i}\right) \exp\left(-\frac{1}{2} x_{i+1}^2 y_{i+1}^2 (u^2 + 2uv + v^2)\right) du \right) dv$$

展開

$$= \int_{-\infty}^{\mu_i} \frac{x_{i+1} y_{i+1}}{\sqrt{2\pi}} \exp\left(-\frac{x_{i+1}^2 y_{i+1}^2 v^2}{2}\right) \left(\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1}{2\nu_i} \left((1 + \nu_i x_{i+1}^2 y_{i+1}^2) u^2 + 2\nu_i x_{i+1}^2 y_{i+1}^2 uv \right)\right) du \right) dv$$

展開

$$= \int_{-\infty}^{\mu_i} \frac{x_{i+1} y_{i+1}}{\sqrt{2\pi}} \exp\left(-\frac{x_{i+1}^2 y_{i+1}^2 v^2}{2}\right) \left(\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1+C_{i+1}}{2\nu_i} \left(u + \frac{C_{i+1} v}{1+C_{i+1}}\right)^2 + \frac{C_{i+1}^2 v^2}{2\nu_i(1+C_{i+1})}\right) du \right) dv$$

$C_{i+1} = \nu_i x_{i+1}^2 y_{i+1}^2$ とおく. 平方完成

$$= \int_{-\infty}^{\mu_i} \frac{x_{i+1} y_{i+1}}{\sqrt{2\pi}} \exp\left(-\frac{x_{i+1}^2 y_{i+1}^2 v^2}{2} + \frac{C_{i+1}^2 v^2}{2\nu_i(1+C_{i+1})}\right) \left(\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1+C_{i+1}}{2\nu_i} \left(u + \frac{C_{i+1} v}{1+C_{i+1}}\right)^2\right) du \right) dv$$

$z = z''$,

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1+C_{i+1}}{2\nu_i} \left(u + \frac{C_{i+1} v}{1+C_{i+1}}\right)^2\right) du$$

$$= \frac{1}{\sqrt{1+C_{i+1}}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\nu_i}} \exp\left(-\frac{1+C_{i+1}}{2\nu_i} \left(u + \frac{C_{i+1} v}{1+C_{i+1}}\right)^2\right) du$$

$$= \frac{1}{\sqrt{1+C_{i+1}}} \int_{-\infty}^{\infty} \mathcal{N}\left(u \mid -\frac{C_{i+1} v}{1+C_{i+1}}, \frac{\nu_i}{1+C_{i+1}}\right) du = \frac{1}{\sqrt{1+C_{i+1}}}$$

よって,

$$\begin{aligned} Z_{i+1} &= \int_{-\infty}^{\mu_i} \frac{x_{i+1} y_{i+1}}{\sqrt{2\pi}} \exp\left(-\frac{x_{i+1}^2 y_{i+1}^2 v^2}{2} + \frac{c_{i+1}^2 v^2}{2v_i(1+c_{i+1})}\right) \frac{1}{\sqrt{1+c_{i+1}}} dv \\ &= \int_{-\infty}^{\mu_i} \frac{x_{i+1} y_{i+1}}{\sqrt{2\pi(1+c_{i+1})}} \exp\left(-\frac{1}{2} \cdot \frac{1}{v_i} \left(c_{i+1} - \frac{c_{i+1}^2}{1+c_{i+1}}\right) v^2\right) dv \\ &= \int_{-\infty}^{\mu_i} \frac{x_{i+1} y_{i+1}}{\sqrt{2\pi(1+c_{i+1})}} \exp\left(-\frac{1}{2} \frac{c_{i+1}}{v_i(1+c_{i+1})} v^2\right) dv \\ &= \int_{-\infty}^{\mu_i} \frac{x_{i+1} y_{i+1}}{\sqrt{2\pi(1+c_{i+1})}} \exp\left(-\frac{1}{2} \frac{x_{i+1}^2 y_{i+1}^2}{1+c_{i+1}} v^2\right) dv \\ z &= \frac{x_{i+1} y_{i+1}}{\sqrt{1+c_{i+1}}} v = d_{i+1} v \quad \text{とおくと,} \end{aligned}$$

$$\begin{aligned} Z_{i+1} &= \int_{-\infty}^{d_{i+1} \mu_i} \frac{dz}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} z^2\right) \cdot \frac{1}{d_{i+1}} dz \\ &= \Phi(d_{i+1} \mu_i) \\ &= \Phi\left(\frac{x_{i+1} y_{i+1} \mu_i}{\sqrt{1+v_i x_{i+1}^2 y_{i+1}^2}}\right). \end{aligned}$$

これは $x_{i+1} = 0$ のときも定義され値が一一致するから、任意の $x_{i+1} \in \mathbb{R}$ で成り立つ。

以上より, $a_{i+1} := \frac{x_{i+1} y_{i+1} \mu_i}{\sqrt{1+v_i x_{i+1}^2 y_{i+1}^2}}$ とおくと $Z_{i+1} = \Phi(a_{i+1})$.

モーメントマッチングを行おう。

$$\begin{aligned} &\frac{\partial}{\partial \mu_i} \log Z_{i+1} \\ &= Z_{i+1}^{-1} \frac{\partial}{\partial \mu_i} Z_{i+1} \\ &= \Phi(a_{i+1})^{-1} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} a_{i+1}^2\right) \frac{\partial}{\partial \mu_i} a_{i+1} \end{aligned}$$

$$\begin{aligned}
&= \Phi(a_{i1})^{-1} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} a_{i1}^2\right) \frac{a_{i1}}{\mu_i} \\
&\quad \frac{\partial}{\partial \mu_i} \log Z_{i1} \\
&= Z_{i1}^{-1} \frac{\partial}{\partial \mu_i} Z_{i1} \\
&= \Phi(a_{i1})^{-1} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} a_{i1}^2\right) \frac{\partial}{\partial \mu_i} a_{i1} \\
&= -\frac{1}{2} \Phi(a_{i1})^{-1} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} a_{i1}^2\right) \frac{x_{i1}^3 y_{i1}^3 \mu_i}{\left(\sqrt{1 + \nu_i x_{i1}^2 y_{i1}^2}\right)^3} \\
&= -\frac{1}{2} \Phi(a_{i1})^{-1} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} a_{i1}^2\right) \frac{a_{i1}^3}{\mu_i^2}.
\end{aligned}$$

よって更新式は,

$$\begin{aligned}
\mu_{i1} &= \mu_i + \nu_i \frac{\partial}{\partial \mu_i} \log Z_{i1} = \mu_i + \frac{1}{\sqrt{2\pi}} \frac{\nu_i a_{i1}}{\mu_i} \Phi(a_{i1})^{-1} \exp\left(-\frac{1}{2} a_{i1}^2\right). \\
\nu_{i1} &= \nu_i - \nu_i^2 \left(\left(\frac{\partial}{\partial \mu_i} \log Z_{i1} \right)^2 - 2 \frac{\partial}{\partial \nu_i} \log Z_{i1} \right) \\
&= \nu_i - \nu_i^2 \left(\frac{1}{2\pi} \Phi(a_{i1})^{-2} \exp(-a_{i1}^2) \frac{a_{i1}^2}{\mu_i^2} + \Phi(a_{i1})^{-1} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} a_{i1}^2\right) \frac{a_{i1}^3}{\mu_i^2} \right) \\
&= \nu_i - \frac{\nu_i^2 a_{i1}^2}{\mu_i^2} \Phi(a_{i1})^{-2} \left(\frac{1}{2\pi} \exp(-a_{i1}^2) + \frac{a_{i1}}{\sqrt{2\pi}} \Phi(a_{i1}) \exp\left(-\frac{1}{2} a_{i1}^2\right) \right)
\end{aligned}$$

4.2.6 期待値伝播法.

・ 仮定密度フィルタリング: 一度学習したパラメータは捨ててしまい, 再び学習に便利.

✓ メモリ効率が良い

✗ 逐次的に入ってくるパラメータの順序に近似結果が強く依存.

・ 期待値伝播法 (expectation propagation)

仮定密度フィルタリングをバッチ学習もできるように一般化した手法.

最適化の過程で同じパラメータを何度も用いることで, 精度の高い近似が可能.

- パラメータ θ に事前分布 $p(\theta)$ を設定し、各 x_n に対し $p(x_n|\theta)$ に従うとする。

考えるモデルは、

$$p(\mathcal{X}, \theta) = p(\theta) \prod_{n=1}^N p(x_n|\theta).$$

これを、

$$f_n(\theta) := \begin{cases} p(\theta) & (n=0) \\ p(x_n|\theta) & (n=1, \dots, N) \end{cases}$$

とすると、 $p(\mathcal{X}, \theta) = \prod_{n=0}^N f_n(\theta)$ と表せる。 f_n を **因子** (factor) といい。

- 事後分布は、

$$p(\theta|\mathcal{X}) = \frac{1}{p(\mathcal{X})} p(\mathcal{X}, \theta) = \frac{1}{p(\mathcal{X})} \prod_{n=0}^N f_n(\theta).$$

この事後分布の近似を、次のように近似因子 \tilde{f}_n の積で表す:

$$p(\theta|\mathcal{X}) \approx q(\theta) = \frac{1}{Z} \prod_{n=0}^N \tilde{f}_n(\theta). \quad (Z: \text{正規化定数})$$

\tilde{f}_n を指数型分布族の分布にしておけば、 q も指数型分布族の分布になる。

- $q(\theta)$ が $p(\theta|\mathcal{X})$ をよく近似するよう、各 \tilde{f}_n を逐次的に更新する。

$q(\theta)$ のパラメータを適当に初期化しておく。

現在の $q(\theta)$ を $q_{old}(\theta)$ と書く。 \tilde{f}_n の更新を行う。

まず、 $q_{old}(\theta)$ から n 番目の因子を取り除く:

$$q_{\setminus n}(\theta) = \frac{q_{old}(\theta)}{\tilde{f}_n(\theta)}$$

\tilde{f}_n の代わりに f_n を用いる分布を $r(\theta) = \frac{1}{Z_n} f_n(\theta) q_{\setminus n}(\theta)$ とし、

新たな近似分布 $q_{\text{new}}(\theta)$ を次をみたすように更新:

$$q_{\text{new}}(\theta) = \underset{q}{\operatorname{argmin}} D_{\text{KL}}[r(\theta) \parallel q(\theta)].$$

これは、 q が指数型分布族なのでモーメントマッチングでできる。

最後に、 \tilde{f}_n を次のように更新する。

$$\tilde{f}_n(\theta) \leftarrow \sum_n \frac{q_{\text{new}}(\theta)}{q_n(\theta)}.$$

以上の更新を $n = 0, \dots, N$ に対して繰り返し行う。

- 以上の計算は、各 \sum_n が計算できるならば実行可能
- 収束性については理論的に保証されてはいないが、実験的には良い性能を示す。