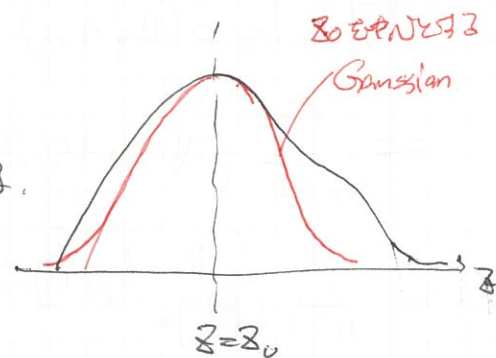


4.4 Laplace 近似

○ 目的

確率密度分布 $p(z)$ を Gaussian に近似する。



○ 1変数の case

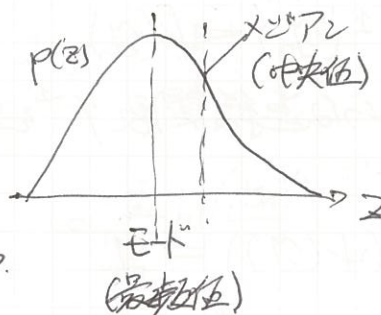
仮定: $p(z) = \frac{1}{Z} f(z)$ (4.125)

正規化係数 $Z = \int dz f(z)$

・モード ... 最頻値

$p(z)$ は $z = z_0$ でモードになる

$\Leftrightarrow p(z)$ は $z = z_0$ で最大値をとる。



$\Rightarrow \frac{dp}{dz} \Big|_{z=z_0} = 0$ \downarrow z は定数。

$\Leftrightarrow \frac{df(z)}{dz} \Big|_{z=z_0} = 0$

・ Gaussian の対数関数は 2 次関数。

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) \quad (2.42)$$

$$\begin{aligned} \log \mathcal{N}(x|\mu, \sigma^2) &= -\frac{1}{2\sigma^2}(x-\mu)^2 + \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) \\ &= \left(-\frac{1}{2\sigma^2}\right)x^2 + \frac{\mu}{\sigma^2}x + \left(-\frac{\mu^2}{2\sigma^2} + \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)\right) \\ &= (x \text{ の 2 次関数}). \end{aligned}$$

7.17. $\log f(z)$ を $z=z_0$ 近傍で z の 2 次関数 とおくと
 ✓ Taylor 展開

$$\log f(z) = \log f(z_0) + \underbrace{\frac{d}{dz} \log f(z) \Big|_{z=z_0}}_{\parallel 0} (z-z_0) + \frac{1}{2} \underbrace{\frac{d^2}{dz^2} \log f(z) \Big|_{z=z_0}}_{\parallel A} (z-z_0)^2 + O((z-z_0)^3)$$

$$\frac{1}{f(z)} \frac{df(z)}{dz} \Big|_{z=z_0} \parallel 0$$

$$A = -\frac{d^2}{dz^2} \log f(z) \Big|_{z=z_0} \quad (4.128)$$

$$\approx \log f(z_0) - \frac{1}{2} A (z-z_0)^2$$

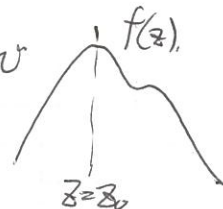
↓ 両辺の対数と比較

$$f(z) \approx f(z_0) \exp\left(-\frac{A}{2} (z-z_0)^2\right) \quad (4.129)$$

• $f(z_0)$ を求めよ (正規化)

$$\text{仮定: } A = -\frac{d^2}{dz^2} \log f(z) \Big|_{z=z_0} > 0$$

$f(z)$ は $z=z_0$ で
 局所極大
 (上凸)
~~(4.128)~~



$$1 = \int_{-\infty}^{\infty} dz f(z) = f(z_0) \int_{-\infty}^{\infty} dz \exp\left(-\frac{A}{2} (z-z_0)^2\right)$$

$\parallel \sqrt{\frac{A}{2}} (z-z_0) = x \text{ とおく}$

$$= \sqrt{\frac{2\pi}{A}} f(z_0) \underbrace{\sqrt{\frac{2}{A}} \int_{-\infty}^{\infty} dx e^{-x^2}}_{\parallel \sqrt{\pi}}$$

$$\therefore f(z_0) = \frac{\sqrt{A}}{\sqrt{2\pi}}$$

8.2. 近似正規分布 $g(z)$ は

$$g(z) = \underbrace{\sqrt{\frac{A}{2\pi}}}_{\parallel \frac{1}{\sqrt{2\pi A}}} \exp\left(-\frac{A}{2} (z-z_0)^2\right) \quad (4.130)$$

$$= \mathcal{N}(z|z_0, A^{-1})$$

$\parallel \frac{1}{2A^{-1}}$

0M次元

$$\text{分布: } p(\mathbf{z}) = \frac{1}{Z} f(\mathbf{z})$$

仮定: $p(\mathbf{z})$ は $\mathbf{z} = \mathbf{z}_0$ 付近で $\nabla \log f(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0} = 0$

$$\nabla f(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0} = 0.$$

$\log f(\mathbf{z})$ は $\mathbf{z} = \mathbf{z}_0$ 付近で Taylor 展開可能.

$$\begin{aligned} \log f(\mathbf{z}) &\stackrel{\text{Taylor}}{\simeq} \log f(\mathbf{z}_0) + \underbrace{\nabla \log f(\mathbf{z}_0)}_{=0} + \frac{1}{2} \underbrace{\nabla \nabla \log f(\mathbf{z}_0)}_{\substack{A = -\nabla \nabla \log f(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0} \\ \text{ヘッセ行列}}} (\mathbf{z} - \mathbf{z}_0)^T (\mathbf{z} - \mathbf{z}_0) \\ &= \log f(\mathbf{z}_0) - \frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T A (\mathbf{z} - \mathbf{z}_0) \end{aligned} \quad (4.132)$$

↓ 両辺の指数をとり

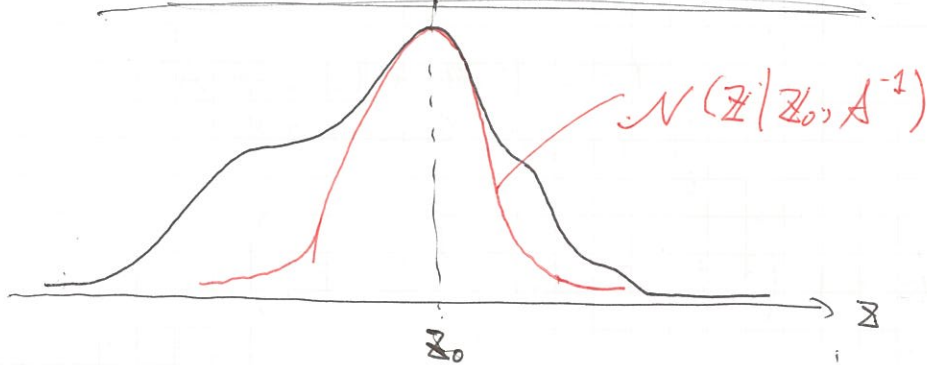
$$f(\mathbf{z}) \simeq f(\mathbf{z}_0) \exp\left(-\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T A (\mathbf{z} - \mathbf{z}_0)\right) \quad (4.133)$$

$$\downarrow \int d\mathbf{z} f(\mathbf{z}) = 1 \Rightarrow f(\mathbf{z}_0) = \frac{\sqrt{|A|}}{\sqrt{2\pi}^m}$$

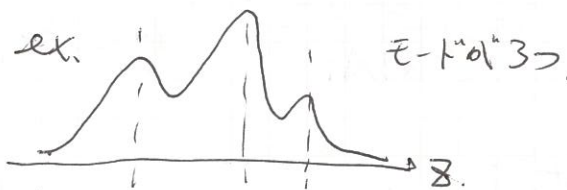
$$\begin{aligned} q(\mathbf{z}) &= \frac{\sqrt{|A|}}{\sqrt{2\pi}^m} \exp\left(-\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T A (\mathbf{z} - \mathbf{z}_0)\right) \\ &= \mathcal{N}(\mathbf{z} | \mathbf{z}_0, A^{-1}) \end{aligned} \quad (4.134)$$

つまり...

モード \bar{x}_0 を見つければ近似できる.



○注意



- 一般的に分布は多峰的
→ このモードを使うと Laplace 近似が異なる.
- Laplace 近似を使う際、真の分布の正規化係数 ~~は必要~~ ~~知る必要は~~
- 中心極限定理より、観測データが増えるほど
事後確率分布は Gaussian に近づく.
→ Laplace 近似が有効
- Gaussian 故 実数変数にのみ適用できる. 局所的
- 真の分布のある1点 (モード付近) の特性しか捉えられない.
cf. 10章 ... 全体的な特性を捉える方法.

4.4.1 モデルの比較とBIC

○ 正規化係数 Z の近似

(\mathbf{A})⁻¹ の行列

$$Z = \int d\mathbf{z} f(\mathbf{z})$$

$f(\mathbf{z})$ は 平均値: モード \mathbf{z}_0 , 分散 \mathbf{A}^{-1} の

(4.133) Gaussian 近似

$$\simeq f(\mathbf{z}_0) \int d\mathbf{z} \exp\left(-\frac{1}{2}(\mathbf{z}-\mathbf{z}_0)^T \mathbf{A}(\mathbf{z}-\mathbf{z}_0)\right)$$

$$= f(\mathbf{z}_0) \frac{(2\pi)^{\frac{M}{2}}}{\sqrt{|\mathbf{A}|}} \quad (4.135) \quad \mathbf{z} \in \mathbb{R}^M \quad \text{空間の次元の数}$$

○ モデルエビデンス

L 個のモデル $\mathcal{M} = \{M_i\}_{i=1, \dots, L}$ の中から
訓練データ \mathcal{D} を最もよく説明するモデルを選ぶ。

訓練データ \mathcal{D} が与えられたときのモデルの事後確率

$$p(M_i | \mathcal{D}) \propto p(\mathcal{D} | M_i) p(M_i)$$

仮定:

$p(M_i) = \text{const.}$ (離散均一分布 $\propto \frac{1}{L}$ のように考える)

どのモデルに選ばれたい?

モデルの不確実性

モデルの事後確率は $p(\mathcal{D} | M_i)$ で決まる。 \leftarrow どのモデルに選ばれたいか
 $p(\mathcal{D} | M_i)$ で決まる。

< model evidence >

“モデル M_i の下で \mathcal{D} が
どれくらいよく説明されるか?”

モデル M_i のパラメータ $\{\theta_i\}$ をもつ。ここで、モデルエビデンスは

$$p(\mathcal{D} | M_i) = \int d\theta_i p(\mathcal{D} | \theta_i, M_i) p(\theta_i | M_i)$$

(3.68)

~~(4.136)~~

と表される。

$$\underbrace{p(\mathcal{D}|M_i)}_{\parallel \sum} = \int d\theta_i \underbrace{p(\mathcal{D}|\theta_i, M_i)}_{\parallel f(\theta_i)} p(\theta_i|M_i) \quad (4.136)$$

$$\mathcal{Z} \simeq f(\mathbf{z}_0) \frac{(2\pi)^{\frac{M}{2}}}{|A|} \quad (4.135)$$

$\downarrow \log$

$$\log \mathcal{Z} \simeq \log f(\mathbf{z}_0) + \frac{M}{2} \log(2\pi) - \frac{1}{2} \log |A|$$

(4.135)
通用後
log する。
(演習 4.22) \downarrow

$$\log p(\mathcal{D}|M_i) \simeq \log p(\mathcal{D}|\theta_{\text{MAP}}, M_i) \quad (4.137)$$

事後分布があるモードでの θ の値

最適パラメータ
を使って評価した対数尤度

$$+ \log p(\theta_{\text{MAP}}|M_i) + \frac{M}{2} \log(2\pi) - \frac{1}{2} \log |A|$$

Occam 係数

where

モデルの複雑さにはペナルティを掛ける

$$A = -\nabla \nabla \log f(\mathbf{z}_0) \quad (4.132)$$

$$= -\nabla \nabla \log p(\mathcal{D}|\theta_{\text{MAP}}, M_i) p(\theta_{\text{MAP}}|M_i)$$

$$= -\nabla \nabla \log p(\mathcal{D}|\theta_{\text{MAP}}, M_i) \quad (4.138)$$

$\frac{p(\theta_{\text{MAP}}|M_i)}{p(\theta_{\text{MAP}}|M_i)}$ (Hessian matrix) この誤極

typo
 ~~$\nabla \nabla p(\theta_{\text{MAP}}|M_i)$~~
 $\neq 0$
?

$$= \nabla \nabla (-\log p(\mathcal{D}|\theta_{\text{MAP}}, M_i)) - \nabla \nabla \log p(\theta_{\text{MAP}}|M_i) \quad (4.138)$$

「正規化係数 \mathcal{Z} を近似すること」 \parallel (演習 4.23)
モデルエビデンスを計算してる。」
H

⑦ 式 (4.139) の導出 (演習 4.23)

~~パラメータの事前分布 $p(\theta)$~~

仮定:

$p(\theta | M_i)$ は Gaussian $\mathcal{N}(\theta | m, V_0)$.

$$\mathcal{N}(\theta | m, V_0) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{\sqrt{|V_0|}} \exp\left(-\frac{1}{2}(\theta - m)^T V_0^{-1}(\theta - m)\right) \quad (243)$$

$$\log \mathcal{N}(\theta | m, V_0) = -\frac{1}{2}(\theta - m)^T V_0^{-1}(\theta - m) + \log\left(\frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{\sqrt{|V_0|}}\right)$$

$$\begin{aligned} (\langle \theta | - \langle m |) \hat{V}_0 (|\theta\rangle - |m\rangle) &= \langle \theta | \hat{V}_0 | \theta \rangle + \langle m | \hat{V}_0 | m \rangle - \langle \theta | \hat{V}_0 | m \rangle - \langle m | \hat{V}_0 | \theta \rangle \\ \hat{V}_0^T &= \hat{V}_0. \\ &= \langle \theta | \hat{V}_0 | \theta \rangle - 2 \langle \theta | \hat{V}_0 | m \rangle + \langle m | \hat{V}_0 | m \rangle \\ &\rightarrow \theta^T V_0^{-1} \theta - 2 \theta^T V_0^{-1} m + m^T V_0^{-1} m. \end{aligned}$$

$$= \underbrace{-\frac{1}{2} \theta^T V_0^{-1} \theta}_{\theta \text{ の二次項}} + \underbrace{\theta^T V_0^{-1} m}_{\theta \text{ の一次項}} - \underbrace{\frac{1}{2} m^T V_0^{-1} m}_{\text{定数項}} + \log\left(\frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{\sqrt{|V_0|}}\right)$$

$$\nabla \log \mathcal{N}(\theta | m, V_0) = -V_0^{-1} \theta + V_0^{-1} m.$$

$$\nabla \nabla \log \mathcal{N}(\theta | m, V_0) = -V_0^{-1}.$$

52. (4.138) の

$$A = H - \nabla \nabla \log p(\theta_{\text{MAP}} | M_i)$$

$$= H + V_0^{-1}.$$

Chap. 2. (4.137) is

$$\log p(\mathcal{D} | M_i) \quad (4.137)$$

$$\simeq \underbrace{\log p(\mathcal{D} | \theta_{\text{MAP}}, M_i)}_{\textcircled{1}} + \underbrace{\log p(\theta_{\text{MAP}} | M_i)}_{\substack{\textcircled{2} \\ \mathcal{N}(\theta_{\text{MAP}} | \mu, V_0)}} + \underbrace{\frac{M}{2} \log(2\pi)}_{\textcircled{3}} - \underbrace{\frac{1}{2} \log |A|}_{\textcircled{4}}$$

$$= \log p(\mathcal{D} | \theta_{\text{MAP}}, M_i) \textcircled{1}$$

$$- \frac{1}{2} (\theta_{\text{MAP}} - \mu)^T V_0^{-1} (\theta_{\text{MAP}} - \mu) + \underbrace{\log \left(\frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{\sqrt{|V_0|}} \right)}_{\text{Const.}} \textcircled{2}$$

$$+ \underbrace{\frac{M}{2} \log(2\pi)}_{\text{Const.}} \textcircled{3}$$

$$- \frac{1}{2} \log |H + V_0^{-1}| \textcircled{4}$$

$$= |H(1 + H^{-1}V_0^{-1})|$$

$$= |H| |1 + H^{-1}V_0^{-1}| \quad \text{依定}$$

$$\simeq |H| \quad \checkmark \quad \text{Gauss 分布 } \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2\sigma^2}}$$

$$\rightarrow V_0 \text{ 很大}$$

$$\rightarrow V_0^{-1} \text{ 微小}$$

$$\simeq \log p(\mathcal{D} | \theta_{\text{MAP}}, M_i) - \frac{1}{2} (\theta_{\text{MAP}} - \mu)^T V_0^{-1} (\theta_{\text{MAP}} - \mu) - \frac{1}{2} \log |H| + \text{const.}$$

✓

$$H = -\nabla \nabla \log p(\mathcal{D} | \theta_{\text{MAP}}) \quad \text{独立同分布}$$

$$= \sum_{n=1}^N H_n$$

△ 仮定: i.i.d data.

$$p(\mathcal{D} | \theta_{\text{MAP}}) = \prod_{n=1}^N p_n$$

$$= \mathcal{N}\left(\frac{1}{N} \sum_{n=1}^N H_n\right)$$

n番目のデータ点への寄与

$$= \mathcal{N}(\hat{H})$$

データ数

よ),

$$\log |H| = \log |N \hat{H}|$$

Hの次元数... M

θのパラメータ数

$$\nabla \nabla \log p(\mathcal{D} | \theta) = \left(\frac{\partial^2 \log p(\mathcal{D} | \theta)}{\partial \theta_i \partial \theta_j} \right)_{i,j}$$

$$= \log (N^M |\hat{H}|)$$

$$= M \log N + \log |\hat{H}|$$

O(N) vs O(1)

N >> 1

$$\simeq M \log N$$

よ),

$$\log p(\mathcal{D} | M_i)$$

$$\simeq \log p(\mathcal{D} | \theta_{\text{MAP}}, M_i) - \frac{1}{2} (\theta_{\text{MAP}} - m) V_0^{-1} (\theta_{\text{MAP}} - m) - \frac{1}{2} \log |H| + \text{const.}$$

V_0^{-1} は微小 → 0

$$\simeq \log p(\mathcal{D} | \theta_{\text{MAP}}, M_i) - \frac{1}{2} M \log N \quad (4.139)$$

< ベイズ情報量規準 (BIC)
Bayesian Information Criterion >

or

< Schwarz criterion >

□.

cf. 赤池情報量規準 (AIC)

$$\log p(\mathcal{D} | w_{\text{MLE}}) - M \quad (4.73)$$

BICの導出

→ BICは赤池情報量規準の近似

