

スパースベイズ推定

株式会社 シグナルアナリシス 関原謙介

1 データモデル

本当は x_1, x_2, \dots, x_N を観測したいのだがこれらは直接観測できず、代わりに一群の観測データ y_1, y_2, \dots, y_M が得られる状況を考えよう。知りたい量 x_1, x_2, \dots, x_N と観測データ y_1, y_2, \dots, y_M を、次の列ベクトル x と y で表現する。すなわち、

$$\text{未知量: } x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \quad \text{観測データ: } y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{bmatrix} \quad (1)$$

と定義する。この x は未知量ベクトルあるいは解ベクトル、 y は観測ベクトルあるいはデータベクトルと呼ばれる。ベクトル x の集合は解空間 (solution space)、ベクトル y の集合は観測空間 (observation space) と呼ばれることもある。

ベクトル x とベクトル y は線形な関係、

$$y = Hx \quad (2)$$

で結ばれているとする。ここで H は未知量 x と観測結果 y を結びつける $M \times N$ の行列である。本書では未知量 x と観測データ y は共に実数であり、 H も実数行列であるとする。さらに、観測結果には、信号成分に加法的にノイズ ε が重畳すると仮定すれば、観測データのモデルは

$$y = Hx + \varepsilon \quad (3)$$

となる。ノイズベクトル ε は M 次元の列ベクトルで、

$$\varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_M \end{bmatrix} \quad (4)$$

であり、 j 番目の要素 ε_j は j 番目の観測データ y_j に重畳するノイズを表す。式 (3) で表される観測モデルを線形離散モデルと呼ぶ。

2 スパースベイズ推定：確率モデル

線形離散モデル

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \varepsilon$$

を仮定し，観測データ \mathbf{y} から未知な量 \mathbf{x} を求める問題を引き続き考えよう．事前分布として， \mathbf{x} の各要素に対し独立で同一な正規分布，

$$x_j \sim \mathcal{N}(x_j|0, \alpha^{-1}) \quad (5)$$

を仮定して，この事前分布のもとで L_2 ノルム正則化ミニマムノルム解が導かれる．それでは，この事前分布を少しだけ変えて，分散（あるいは精度）が各要素 x_j に固有の値を取る，すなわち，

$$x_j \sim \mathcal{N}(x_j|0, \alpha_j^{-1}) \quad (6)$$

としたらどのような解が得られるであろうか．実は， L_2 ノルム正則化ミニマムノルム解とは全く異なるスパースな解が求まるのである．事前分布として式 (6) を用いた \mathbf{x} の推定法はスパースベイズ推定と呼ばれる．以下にこのスパースベイズ推定を説明する．

精度行列を用いて議論を進めることにして，事前確率分布とデータ尤度が，

$$p(\mathbf{x}|\Phi) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \Phi^{-1}) \quad (7)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{H}\mathbf{x}, \beta^{-1}\mathbf{I}) \quad (8)$$

の確率モデルに従うとする．説明を簡単にするため，ノイズの精度行列は $\beta\mathbf{I}$ として，精度 β は既知であるとするとする．また， Φ が事前分布の精度行列で，対角行列：

$$\Phi = \begin{bmatrix} \alpha_1 & \cdots & 0 \\ 0 & \ddots & 0 \\ 0 & \cdots & \alpha_N \end{bmatrix}$$

と仮定する．したがって，

$$p(\mathbf{x}|\Phi) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \Phi^{-1}) = \prod_{j=1}^N \mathcal{N}(x_j|0, \alpha_j^{-1}) \quad (9)$$

と表記することもできる．ここで，各 α_j がみな等しい，つまり， $\alpha_1 = \cdots = \alpha_N$ とすれば， L_2 正則化ミニマムノルム推定に等しくなる．ここで，各 α_j が皆等しいとの仮定を置かないと，ミニマムノルム解とは全く異なるスパースな解が求まる．

3 推定の定式化

まず，ベクトル α を $\alpha = [\alpha_1, \dots, \alpha_N]^T$ と定義して導入する．ここで， $\Phi = \text{diag}(\alpha)$ である．ここで， $\text{diag}(\cdot)$ は括弧内に記載されたベクトルを対角成分に持つ行列を意味する．式 (9) は，あらためて，

$$p(\mathbf{x}|\alpha) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \text{diag}(\alpha)^{-1}) = \prod_{j=1}^N \mathcal{N}(x_j|0, \alpha_j^{-1}) \quad (10)$$

とも書ける．ここで，推定すべき未知パラメータは，そもそもの未知量 x と，事前分布に現れるハイパーパラメータ α である．観測データ y に重畳するノイズの分散 β^{-1} も一般的には未知の量であるが，ここでは既知と仮定する．

ここで，真にベイズ的な取り扱いをしようとするれば両方の未知パラメータに対する結合事後分布 $p(x, \alpha|y)$ を

$$p(x, \alpha|y) = \frac{p(y|x, \alpha)p(x, \alpha)}{\iint p(y|x, \alpha)p(x, \alpha)d\alpha dx} \quad (11)$$

として求めたいのだが，分母の積分が計算できず，結合事後分布 $p(x, \alpha|y)$ を求めることができない．したがって，代わりに，

$$p(x|y) = \int p(x|y, \alpha)p(\alpha|y)d\alpha \approx p(x|y, \hat{\alpha}) \quad (12)$$

から，事後分布 $p(x|y)$ の近似解 $p(x|y, \hat{\alpha})$ を求め，これを最大とする x の MAP 推定解を求める．

ここで，事前分布と尤度が正規分布（式 (7) および (8)）の場合，事後分布 $p(x|y, \alpha)$ も正規分布，

$$p(x|y, \alpha) = \mathcal{N}(x|\bar{x}, \Gamma^{-1}) \quad (13)$$

で与えられる．ここで，精度行列と平均は

$$\Gamma = \Phi + \beta H^T H \quad (14)$$

$$\bar{x} = \beta \Gamma^{-1} H^T y \quad (15)$$

で与えられる．したがって， $\hat{\alpha}$ が求まっていれば，式 (14) および (15) の式において $\Phi = \text{diag}(\hat{\alpha})$ を用いて計算した事後分布が $p(x|y, \hat{\alpha})$ である．したがって， $p(x|y, \hat{\alpha})$ を最大とする x ，すなわち， $\Phi = \text{diag}(\hat{\alpha})$ として，式 (15) で与えられる \bar{x} が未知量 x の MAP（および MMSE）推定解である．

それでは， $\hat{\alpha}$ はどうやって求めることができるであろうか．もちろん，EM アルゴリズムを用いて求めることができるが，EM アルゴリズムは推定が劣決定 ($M \ll N$) の場合，収束が遅いという問題が知られている．したがって，ここでは周辺尤度 $p(y|\alpha)$ を，EM アルゴリズムのように間接的にではなく，直接最大とする方法を説明する．

4 周辺尤度関数の導出

α の推定解 $\hat{\alpha}$ は周辺尤度 $\log p(y|\alpha)$ を最大とする α として求める．まず周辺尤度関数を導出するために，

$$p(y|\alpha) = \int p(y, x|\alpha)dx = \int p(y|x)p(x|\alpha)dx$$

に

$$p(x|\alpha) = \frac{|\Phi|^{1/2}}{(2\pi)^{N/2}} \exp\left[-\frac{1}{2}\mathbf{x}^T \Phi \mathbf{x}\right] \quad (16)$$

$$p(y|x) = \left(\frac{\beta}{2\pi}\right)^{M/2} \exp\left[-\frac{\beta}{2}\|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2\right] \quad (17)$$

を代入して整理する．すると，

$$p(\mathbf{y}|\boldsymbol{\alpha}) = \frac{|\boldsymbol{\Phi}|^{1/2}}{(2\pi)^{N/2}} \left(\frac{\beta}{2\pi}\right)^{M/2} \int \exp[-D(\mathbf{x})] d\mathbf{x} \quad (18)$$

を得る．ここで，

$$D(\mathbf{x}) = \frac{\beta}{2} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 + \frac{1}{2} \mathbf{x}^T \boldsymbol{\Phi} \mathbf{x} \quad (19)$$

である．

上式の $D(\mathbf{x})$ を書き直してみると，

$$D(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T [\boldsymbol{\Phi} + \beta \mathbf{H}^T \mathbf{H}] \mathbf{x} - \beta \mathbf{x}^T \mathbf{H}^T \mathbf{y} + \mathcal{C} \quad (20)$$

と表すことができる [問題 5.1]¹．なお上式では \mathbf{x} を含まない項は \mathcal{C} で表した．この $D(\mathbf{x})$ は平方完成された

$$D(\mathbf{x}) = \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}})^T \boldsymbol{\Gamma} (\mathbf{x} - \bar{\mathbf{x}}) + \xi \quad (21)$$

の形で表すことができる [問題 5.2]．ここで，右辺の ξ は \mathbf{x} を含まない残りの項を表す．上式の $D(\mathbf{x})$ は $\mathbf{x} = \bar{\mathbf{x}}$ で最小値を取り，その最小値が ξ である．この ξ の値は式 (19) に $\mathbf{x} = \bar{\mathbf{x}}$ を代入することにより求まり， $\xi = D(\bar{\mathbf{x}})$ である．したがって結局， $D(\mathbf{x})$ は，

$$D(\mathbf{x}) = \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}})^T \boldsymbol{\Gamma}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) + D(\bar{\mathbf{x}}) \quad (22)$$

と求まる．

この結果を用いて式 (18) 右辺の積分を実行してみよう．まず，

$$\int \exp[-D(\mathbf{x})] d\mathbf{x} = \exp[-D(\bar{\mathbf{x}})] \int \exp\left[-\frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}})^T \boldsymbol{\Gamma} (\mathbf{x} - \bar{\mathbf{x}})\right] d\mathbf{x}$$

が成り立つ．上式右辺の積分は，平均 $\bar{\mathbf{x}}$ で精度行列 $\boldsymbol{\Gamma}$ である正規分布の全積分を考慮すれば，

$$\int \exp\left[-\frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}})^T \boldsymbol{\Gamma} (\mathbf{x} - \bar{\mathbf{x}})\right] d\mathbf{x} = \frac{(2\pi)^{N/2}}{|\boldsymbol{\Gamma}|^{1/2}}$$

と表されるので，式 (18) にこれらの結果を代入して，

$$p(\mathbf{y}|\boldsymbol{\alpha}) = \frac{|\boldsymbol{\Phi}|^{1/2}}{(2\pi)^{N/2}} \left(\frac{\beta}{2\pi}\right)^{M/2} \exp[-D(\bar{\mathbf{x}})] \frac{(2\pi)^{N/2}}{|\boldsymbol{\Gamma}|^{1/2}} \quad (23)$$

を得る．さらに，上式両辺の対数を取って，定数項を無視すれば，

$$\log p(\mathbf{y}|\boldsymbol{\alpha}) = -\frac{1}{2} \log |\boldsymbol{\Gamma}| + \frac{1}{2} \log |\boldsymbol{\Phi}| + \frac{M}{2} \log \beta - D(\bar{\mathbf{x}}) \quad (24)$$

を得る．

¹問題番号は「ベイズ信号処理」(共立出版)の問題番号を指す．

ここで, $D(\bar{x})$ をさらに以下のように変形する. まず,

$$\begin{aligned}
 D(\bar{x}) &= \frac{\beta}{2} \|\mathbf{y} - \mathbf{H}\bar{x}\|^2 + \frac{1}{2} \bar{x}^T \Phi \bar{x} \\
 &= \frac{\beta}{2} \left(\mathbf{y}^T \mathbf{y} - 2\bar{x}^T \mathbf{H}^T \mathbf{y} + \bar{x}^T \mathbf{H}^T \mathbf{H} \bar{x} \right) + \frac{1}{2} \bar{x}^T \Phi \bar{x} \\
 &= \frac{1}{2} \left[\beta \mathbf{y}^T \mathbf{y} - 2\bar{x}^T \beta \mathbf{H}^T \mathbf{y} + \bar{x}^T (\Phi + \beta \mathbf{H}^T \mathbf{H}) \bar{x} \right] \\
 &= \frac{1}{2} \left[\beta \mathbf{y}^T \mathbf{y} - 2\bar{x}^T \beta \mathbf{H}^T \mathbf{y} + \bar{x}^T \Gamma \bar{x} \right] \quad (25)
 \end{aligned}$$

である. さらに, $\beta \mathbf{H}^T \mathbf{y} = \Gamma \bar{x}$ の関係を用いれば, 最終的に

$$D(\bar{x}) = \frac{1}{2} \mathbf{y}^T \left[\beta^{-1} \mathbf{I} + \mathbf{H} \Phi^{-1} \mathbf{H}^T \right]^{-1} \mathbf{y} = \frac{1}{2} \mathbf{y}^T \Sigma_y^{-1} \mathbf{y} \quad (26)$$

となることを示すことができる [問題 5.3]. ここで, Σ_y はモデルデータ共分散行列と呼ばれ

$$\Sigma_y = \beta^{-1} \mathbf{I} + \mathbf{H} \Phi^{-1} \mathbf{H}^T \quad (27)$$

である.

ふたたび式 (24) に戻って $D(\bar{x})$ 以外の部分を計算する. 行列式に関して,

$$|\Phi| |\beta^{-1} \mathbf{I} + \mathbf{H} \Phi^{-1} \mathbf{H}^T| = |\beta^{-1} \mathbf{I}| |\Phi + \beta \mathbf{H}^T \mathbf{H}| \quad (28)$$

が成り立つ. したがって, 式 (27) および (14) を代入すれば,

$$|\Phi| |\Sigma_y| = |\beta^{-1} \mathbf{I}| |\Gamma| \quad (29)$$

より,

$$\log |\Sigma_y| = \log |\Gamma| - M \log(\beta) - \log |\Phi| \quad (30)$$

を得る. 式 (26) および (30) を式 (24) に代入すれば, 結局, 周辺尤度 $\log p(\mathbf{y}|\alpha)$ として

$$\log p(\mathbf{y}|\alpha) = -\frac{1}{2} \log |\Sigma_y| - \frac{1}{2} \mathbf{y}^T \Sigma_y^{-1} \mathbf{y} \quad (31)$$

を得る. したがって, $\log p(\mathbf{y}|\alpha)$ を最大とする α , あるいは,

$$\mathcal{F}(\alpha) = \log |\Sigma_y| + \mathbf{y}^T \Sigma_y^{-1} \mathbf{y} \quad (32)$$

として, $\mathcal{F}(\alpha)$ を最小とする α として最適推定値 $\hat{\alpha}$ を求める. すなわち, この $\mathcal{F}(\alpha)$ が α を推定するためのコスト関数である.

5 周辺尤度 (31) の自由エネルギーを用いた導出

周辺尤度 (式 (31)) は自由エネルギーを用いて導出でき, 前節に述べたものより若干簡便に導出できる². 評価関数 $\mathcal{F}[q, \alpha]$:

$$\mathcal{F}[q, \alpha] = \int dx q(x) [\log p(x, \mathbf{y}|\alpha) - \log q(x)] \quad (33)$$

²自由エネルギーについては「ベイズ信号処理」(共立出版)を参照のこと.

は自由エネルギーと呼ばれる．この $\mathcal{F}[q, \alpha]$ は2つの量，ハイパーパラメータ α と任意の確率分布 $q(x)$ の関数である．第4.6節の議論によれば，自由エネルギー $\mathcal{F}[q, \alpha]$ を最大にする確率分布は未知量 x の (α をこの時点での値に固定した) 事後分布 $p(x|\mathbf{y}, \alpha)$ であることを示している．すなわち，自由エネルギーを $q(x)$ について最大にすることはEMアルゴリズムのEステップに対応する．Eステップ終了時点での自由エネルギーは

$$\begin{aligned} \mathcal{F}[p(x|\mathbf{y}, \alpha), \alpha] &= \int dx p(x|\mathbf{y}, \alpha) [\log p(x, \mathbf{y}|\alpha) - \log p(x|\mathbf{y}, \alpha)] \\ &= \int dx p(x|\mathbf{y}, \alpha) \log \frac{p(x, \mathbf{y}|\alpha)}{p(x|\mathbf{y}, \alpha)} = \int dx p(x|\mathbf{y}, \alpha) \log p(\mathbf{y}|\alpha) \\ &= \log p(\mathbf{y}|\alpha) \int dx p(x|\mathbf{y}, \alpha) = \log p(\mathbf{y}|\alpha) \end{aligned} \quad (34)$$

となり，自由エネルギーの値は周辺尤度 $\log p(\mathbf{y}|\alpha)$ に等しい．したがって，

$$\begin{aligned} \log p(\mathbf{y}|\alpha) &= \int dx p(x|\mathbf{y}, \alpha) [\log p(x, \mathbf{y}|\alpha) - \log p(x|\mathbf{y}, \alpha)] \\ &= E[\log p(x, \mathbf{y}|\alpha) - \log p(x|\mathbf{y}, \alpha)] \\ &= E[\log p(\mathbf{y}|x) + \log p(x|\alpha)] + \mathcal{H}[p(x|\mathbf{y}, \alpha)] \end{aligned} \quad (35)$$

が成立する．ここで， $E[\cdot]$ は事後分布 $p(x|\mathbf{y}, \alpha)$ で平均を取ることを意味し， $\mathcal{H}[p(x|\mathbf{y}, \alpha)]$ は事後分布のエントロピーである．

式(35)に

$$\begin{aligned} p(x|\alpha) &= \mathcal{N}(x|\mathbf{0}, \Phi^{-1}) \\ p(\mathbf{y}|x) &= \mathcal{N}(\mathbf{y}|\mathbf{H}x, \beta^{-1}\mathbf{I}) \\ p(x|\mathbf{y}, \alpha) &= \mathcal{N}(x|\bar{x}, \Gamma^{-1}) \end{aligned}$$

を代入して整理する．ただしここで，

$$\Phi = \text{diag}[\alpha_1, \dots, \alpha_N]$$

である．すると，定数を見捨て， $\mathcal{H}[p(x|\mathbf{y}, \alpha)] = -\log |\Gamma|$ を用いて，

$$\log p(\mathbf{y}|\alpha) = \frac{1}{2} [\|\Phi\| + \|\beta\mathbf{I}\| - \|\Gamma\|] - \frac{1}{2} E[\beta\|\mathbf{y} - \mathbf{H}x\|^2] - \frac{1}{2} E[\mathbf{x}^T \Phi \mathbf{x}] \quad (36)$$

を得る．さらに，

$$E[\beta\|\mathbf{y} - \mathbf{H}x\|^2] + E[\mathbf{x}^T \Phi \mathbf{x}] = \beta [\mathbf{y}^T \mathbf{y} - 2\bar{x}^T \mathbf{H}^T \mathbf{y}] + E[\mathbf{x}^T (\beta \mathbf{H}^T \mathbf{H} + \Phi) \mathbf{x}] \quad (37)$$

であり，ここで，

$$\begin{aligned} E[\mathbf{x}^T (\beta \mathbf{H}^T \mathbf{H} + \Phi) \mathbf{x}] &= E[\mathbf{x}^T \Gamma \mathbf{x}] = E[\text{tr}(\mathbf{x} \mathbf{x}^T \Gamma)] \\ &= E[\text{tr}(\mathbf{x} \mathbf{x}^T \Gamma)] = \text{tr}[E(\mathbf{x} \mathbf{x}^T) \Gamma] \\ &= \text{tr}[(\bar{x} \bar{x}^T + \Gamma^{-1}) \Gamma] = \bar{x}^T (\beta \mathbf{H}^T \mathbf{H} + \Phi) \bar{x} \end{aligned} \quad (38)$$

であるので，結局，

$$\begin{aligned} E[\beta\|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2] + E[\mathbf{x}^T \Phi \mathbf{x}] \\ = \beta \left[\mathbf{y}^T \mathbf{y} - 2\bar{\mathbf{x}}^T \mathbf{H}^T \mathbf{y} + \bar{\mathbf{x}}^T \mathbf{H}^T \mathbf{H} \bar{\mathbf{x}} \right] + \bar{\mathbf{x}}^T \Phi \bar{\mathbf{x}} \\ = \beta\|\mathbf{y} - \mathbf{H}\bar{\mathbf{x}}\|^2 + \bar{\mathbf{x}}^T \Phi \bar{\mathbf{x}} \quad (39) \end{aligned}$$

を得る．したがって，

$$\log|\Sigma_y| = -(|\Phi| + |\beta\mathbf{I}| - |\Gamma|)$$

を式 (36) に代入すれば，結局，周辺尤度として式 (31)：

$$\log p(\mathbf{y}|\alpha) = -\frac{1}{2} \log|\Sigma_y| - \frac{1}{2} [\beta\|\mathbf{y} - \mathbf{H}\bar{\mathbf{x}}\|^2 + \bar{\mathbf{x}}^T \Phi \bar{\mathbf{x}}]$$

を得る．ここで，

$$\Sigma_y = \beta^{-1}\mathbf{I} + \mathbf{H}\Upsilon\mathbf{H}^T$$

である．

6 ハイパーパラメータ α の更新式

ハイパーパラメータ α の推定値 $\hat{\alpha}$ は，式 (32) に示されるコスト関数 $\mathcal{F}(\alpha)$ を用いて，

$$\hat{\alpha} = \underset{\alpha}{\operatorname{argmin}} \mathcal{F}(\alpha)$$

として求める．この最小値を計算するため，

$$\frac{\partial \mathcal{F}(\alpha)}{\partial \alpha_j} = \frac{\partial}{\partial \alpha_j} \log|\Sigma_y| + \frac{\partial}{\partial \alpha_j} \mathbf{y}^T \Sigma_y^{-1} \mathbf{y} \quad (40)$$

を計算する．まず，右辺第一項は，再び式 (30) を用いて，

$$\frac{\partial}{\partial \alpha_j} \log|\Sigma_y| = \frac{\partial}{\partial \alpha_j} [-M \log \beta - \log|\Phi| + \log|\Gamma|] = -\frac{\partial}{\partial \alpha_j} \log|\Phi| + \frac{\partial}{\partial \alpha_j} \log|\Gamma| \quad (41)$$

である．ここで，まず，

$$\frac{\partial}{\partial \alpha_j} \log|\Phi| = \frac{\partial}{\partial \alpha_j} \sum_{j=1}^N \log \alpha_j = \alpha_j^{-1} \quad (42)$$

である．また， $\Omega_{j,j}$ を (j, j) 要素のみ 1，他の要素はすべてゼロであるような $N \times N$ の行列として，

$$\begin{aligned} \frac{\partial}{\partial \alpha_j} \log|\Gamma| &= \operatorname{tr} \left[\Gamma^{-1} \frac{\partial}{\partial \alpha_j} \Gamma \right] = \operatorname{tr} \left[\Gamma^{-1} \frac{\partial}{\partial \alpha_j} [\Phi + \beta \mathbf{H}^T \mathbf{H}] \right] \\ &= \operatorname{tr} \left[\Gamma^{-1} \frac{\partial}{\partial \alpha_j} \Phi \right] = \operatorname{tr} [\Gamma^{-1} \Omega_{j,j}] = \Sigma_{j,j} \quad (43) \end{aligned}$$

となる．ここで， $\Sigma_{j,j}$ は事後分布の共分散行列 $\Sigma (= \Gamma^{-1})$ の (j, j) 成分である．

次に，式 (40) の右辺第 2 項を微分する．まず，式 (25) および式 (26) を用いれば，

$$\frac{\partial}{\partial \alpha_j} \mathbf{y}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y} = \frac{\partial}{\partial \alpha_j} [\beta \|\mathbf{y} - \mathbf{H}\bar{\mathbf{x}}\|^2 + \bar{\mathbf{x}}^T \boldsymbol{\Phi} \bar{\mathbf{x}}] = \bar{\mathbf{x}}^T \left[\frac{\partial}{\partial \alpha_j} \boldsymbol{\Phi} \right] \bar{\mathbf{x}} = \bar{\mathbf{x}}^T \boldsymbol{\Omega}_{j,j} \bar{\mathbf{x}} = \bar{x}_j^2 \quad (44)$$

したがって，

$$\frac{\partial \mathcal{F}(\boldsymbol{\alpha})}{\partial \alpha_j} = \Sigma_{j,j} - \alpha_j^{-1} + \bar{x}_j^2 = 0 \quad (45)$$

であるので，結局，

$$\hat{\alpha}_j^{-1} = \Sigma_{j,j} + \bar{x}_j^2 \quad (46)$$

を得る．実は，上式は EM アルゴリズムによる更新式に等しい [問題 5.5]．

一方，式 (45) は

$$1 - \alpha_j \Sigma_{j,j} = \alpha_j \bar{x}_j^2 \quad (47)$$

と変形できる．一方，式 (14) は

$$\mathbf{I} - \boldsymbol{\Gamma}^{-1} \boldsymbol{\Phi} = \beta \boldsymbol{\Gamma}^{-1} \mathbf{H}^T \mathbf{H}$$

と変形できる．ここで，上式左辺の (j, j) 成分は $1 - \alpha_j \Sigma_{j,j}$ に等しいことに注目すれば，式 (47) は

$$\hat{\alpha}_j = \frac{[\beta \boldsymbol{\Gamma}^{-1} \mathbf{H}^T \mathbf{H}]_{j,j}}{\bar{x}_j^2} \quad (48)$$

と表記することができる．ここで，上式右辺の $[\cdot]_{j,j}$ は，括弧内の行列の (j, j) 成分を取ることを意味する．式 (48) に示す更新式は MacKay によって提案されたもので，EM アルゴリズムによる更新式 (46) よりも， $M \ll N$ なる劣決定の状況において，収束が早いことが知られている．

式 (48) の右辺において， $\boldsymbol{\Gamma}$ にも α_j は含まれているのに，この α_j を無視しているように感じられるため，読者はここで妙に思うかもしれない．しかし，式 (46) において， $\Sigma_{j,j}$ も \bar{x}_j^2 も α_j に依存している，単に，その依存性が明示的に示されていないだけである．したがって，式 (46) と (48) との違いは，式 (46) においては， $\Sigma_{j,j}$ と \bar{x}_j^2 を古い α_j の値を用いて計算し，式 (46) にしたがって α_j の値を更新するのに対して，式 (48) においては， $\boldsymbol{\Gamma}$ と \bar{x}_j^2 を古い α_j の値を用いて計算し， α_j の値を更新することである．上の議論から示唆されるように，更新式の導出にはいくつかの任意性があり，結果として得られた更新式のどれがより優れたものであるかを理論的にあらかじめ予測するのは難しく，実際の問題に適用してみて初めて分かるのが普通である．式 (48) に示す MacKay の更新式はこのように heuristic (発見的) に見いだされ，提案されたものである．式 (46) に示す EM アルゴリズムの更新式と異なり，更新ごとに周辺尤度を増加させる理論的な保障はない．

式 (48) に示すように， $\boldsymbol{\alpha}$ の更新には事後分布のパラメータ $\boldsymbol{\Gamma}$ と $\bar{\mathbf{x}}$ が必要である．したがって，このアルゴリズムはやはり再帰的なものとなる．まず， $\boldsymbol{\alpha}$ に適当な初期値を仮定して，事後分布のパラメータを式 (14) および (15) から計算する．これは EM アルゴリズムの E ステップと全く同じである．このステップから得られた事後分布のパラメータを用いて式 (48) よりハイパーパラメータ $\boldsymbol{\alpha}$ を更新する．このステップは更新式が異なっているが，EM アルゴリズムの M ステップに相当する．上記のステップを収束条件が満たされるまで繰り返す．収束の判定は， $\hat{\boldsymbol{\Phi}} = \text{diag}(\hat{\boldsymbol{\alpha}})$ として，コスト関数 $\mathcal{F}(\hat{\boldsymbol{\alpha}})$ を計算し，更新を繰り返してもほとんど減少しなくなったときに計算を打ち切る．

7 凹関数の性質を用いた更新式の導出

周辺尤度の式 (31) を増加させる別の更新式を導こう。この導出では、凹関数の性質を用いるため、精度 α_j に対応した分散を ν_j ($\nu_j = \alpha_j^{-1}$) とし、列ベクトル $\boldsymbol{\nu}$ を $\boldsymbol{\nu} = [\nu_1, \dots, \nu_N]$ と定義する。精度ベクトル $\boldsymbol{\alpha}$ の代わりにこの分散ベクトル $\boldsymbol{\nu}$ を、精度行列 $\boldsymbol{\Phi}$ の代わりに共分散行列 $\boldsymbol{\Upsilon}$ ($\boldsymbol{\Upsilon} = \boldsymbol{\Phi}^{-1}$) を用いる。したがって、式 (31) の周辺尤度は

$$\log p(\mathbf{y}|\boldsymbol{\nu}) = -\frac{1}{2} \log |\boldsymbol{\Sigma}_y| - \frac{1}{2} [\beta \|\mathbf{y} - \mathbf{H}\bar{\mathbf{x}}\|^2 + \bar{\mathbf{x}}^T \boldsymbol{\Upsilon}^{-1} \bar{\mathbf{x}}] \quad (49)$$

と表される。

ここで、 $\log |\boldsymbol{\Sigma}_y|$ は $\boldsymbol{\nu}$ に関して凹関数であり、補助変数 z を用いて、

$$z^T \boldsymbol{\nu} - z_0 \geq \log |\boldsymbol{\Sigma}_y|$$

が常に成立する³。したがって、

$$\log p(\mathbf{y}|\boldsymbol{\nu}) \geq -\frac{1}{2} (z^T \boldsymbol{\nu} - z_0) - \frac{1}{2} [\beta \|\mathbf{y} - \mathbf{H}\bar{\mathbf{x}}\|^2 + \bar{\mathbf{x}}^T \boldsymbol{\Upsilon}^{-1} \bar{\mathbf{x}}] \quad (50)$$

であり、補助的コスト関数 $\tilde{\mathcal{F}}(\boldsymbol{\nu}, z)$ を

$$\tilde{\mathcal{F}}(\boldsymbol{\nu}, z) = -\frac{1}{2} (z^T \boldsymbol{\nu} - z_0) - \frac{1}{2} [\beta \|\mathbf{y} - \mathbf{H}\bar{\mathbf{x}}\|^2 + \bar{\mathbf{x}}^T \boldsymbol{\Upsilon}^{-1} \bar{\mathbf{x}}]$$

と定義すれば、必ず、

$$\log p(\mathbf{y}|\boldsymbol{\nu}) \geq \tilde{\mathcal{F}}(\boldsymbol{\nu}, z)$$

が成立するため、 $\tilde{\mathcal{F}}(\boldsymbol{\nu}, z)$ を増加させる $\boldsymbol{\nu}$ は周辺尤度 $\log p(\mathbf{y}|\boldsymbol{\nu})$ を増加させる。

補助変数 z の更新値 \hat{z} は、

$$\hat{z} = \frac{\partial}{\partial \boldsymbol{\nu}} \log |\boldsymbol{\Sigma}_y|$$

から求め、³「ベイズ信号処理」5.5.2 項に示すように

$$\hat{z}_j = \mathbf{h}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{h}_j \quad (51)$$

となる。

$\boldsymbol{\nu}$ の更新解は

$$\hat{\boldsymbol{\nu}} = \underset{\boldsymbol{\nu}}{\operatorname{argmax}} \tilde{\mathcal{F}}(\boldsymbol{\nu}, z) \quad (52)$$

として求めることができる。すなわち、 $\boldsymbol{\nu}$ の更新値 $\hat{\boldsymbol{\nu}}$ は

$$\underset{\boldsymbol{\nu}}{\operatorname{argmax}} \tilde{\mathcal{F}}(\boldsymbol{\nu}, z) = \underset{\boldsymbol{\nu}}{\operatorname{argmin}} [\bar{\mathbf{x}}^T \boldsymbol{\Upsilon}^{-1} \bar{\mathbf{x}} + z^T \boldsymbol{\nu}] = \underset{\boldsymbol{\nu}}{\operatorname{argmin}} \sum_{j=1}^N \left[z_j \nu_j + \frac{\bar{x}_j^2}{\nu_j} \right] \quad (53)$$

から求まる。したがって、

$$\frac{\partial}{\partial \nu_j} \sum_{j=1}^N \left[z_j \nu_j + \frac{\bar{x}_j^2}{\nu_j} \right] = z_j - \frac{\bar{x}_j^2}{\nu_j^2} = 0 \quad (54)$$

³説明については「ベイズ信号処理」(共立出版)を参照のこと。

とにおいて、式 (51) を用いることにより

$$\hat{\nu}_j = \frac{|\bar{x}_j|}{\sqrt{z_j}} = \frac{|\bar{x}_j|}{\sqrt{\mathbf{h}_j^T \boldsymbol{\Sigma}_y^{-1} \mathbf{h}_j}} \quad (55)$$

を得る。 \bar{x} は事後分布の平均であり、式 (14), (15) から求まる。

式 (55) に示すように、 ν の更新には事後分布のパラメータ \bar{x} が必要である。したがって、このアルゴリズムはやはり再帰的なものとなる。まず、 ν に適当な初期値を仮定して、事後分布のパラメータを式 (14) および (15) から計算する。これは EM アルゴリズムの E ステップと全く同じである。このステップから得られた事後分布のパラメータを用いて式 (55) よりハイパーパラメータ ν を更新する。このステップは EM アルゴリズムの M ステップに相当する。上記のステップを収束条件が満たされるまで繰り返す。収束の判定はコスト関数 $\mathcal{F}(\hat{\nu})$ を計算し、更新を繰り返してもほとんど減少しなくなったときに計算を打ち切る。

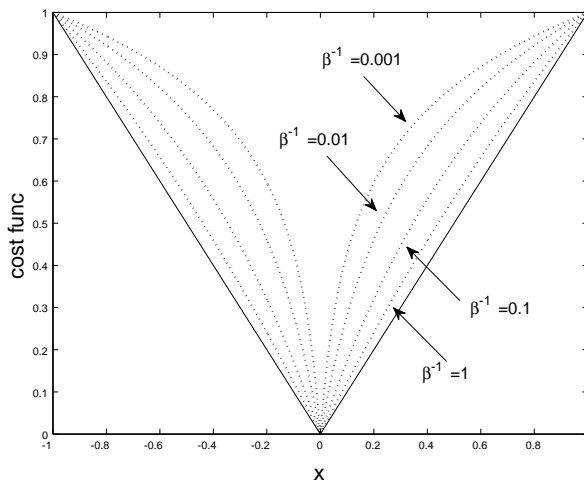


図 1: 式 (64) に示すコスト関数 $\varphi(x)$ の計算結果。 $\varphi(x)$ のプロットを、4 種類の β^{-1} の値について点線で示す。また、比較のため L_1 ノルム解の制約条件 $|x|$ を実線で示す。それぞれのコスト関数は、 $|x| = 1$ で値 1 を取るよう規格化して示す。

8 どうしてスパースな解が得られるのかについての考察

線形離散モデル $\mathbf{y} = \mathbf{H}\mathbf{x} + \boldsymbol{\varepsilon}$ もとで、事前確率分布を以下の正規分布、

$$p(\mathbf{x}|\boldsymbol{\alpha}) = \prod_{j=1}^N \mathcal{N}(x_j|0, \alpha^{-1}) \quad (56)$$

を仮定すると、ベイズ推定解 \bar{x} はスパースな解、すなわち、 $\bar{x}_1, \dots, \bar{x}_N$ のなかで、ほとんどがのものがゼロに非常に近く、少数が明らかにノンゼロな値を持つという解が得られる。本節では、どうしてスパースな解が得られるのかについてコスト関数を用いて考察してみよう。

ここ精度 α ではなく分散 ν を用いて議論しよう⁴。コスト関数を変数 ν を用いてもう 1 度表記してみ

⁴ $\alpha = (\alpha_1, \dots, \alpha_N)$ および $\nu = (\nu_1, \dots, \nu_N)$ であり、 $\nu_j = 1/\alpha_j$ である。

ると,

$$\mathcal{F}(\boldsymbol{\nu}) = \log |\boldsymbol{\Sigma}_y| + \mathbf{y}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y} : \quad \boldsymbol{\Sigma}_y = \beta^{-1} \mathbf{I} + \mathbf{H} \boldsymbol{\Upsilon} \mathbf{H}^T \quad (57)$$

ただし, $\boldsymbol{\Upsilon} = \text{diag}(\boldsymbol{\nu})$ である. したがって, $\boldsymbol{\nu}$ の更新値は

$$\hat{\boldsymbol{\nu}} = \underset{\boldsymbol{\nu}}{\text{argmin}} (\log |\boldsymbol{\Sigma}_y| + \mathbf{y}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{y})$$

と表すことができる. ここで, 上式右辺の第 2 項は

$$\mathbf{y}^T \boldsymbol{\Sigma}_y \mathbf{y} = \min_{\mathbf{x}} [\beta \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 + \mathbf{x}^T \boldsymbol{\Upsilon}^{-1} \mathbf{x}] \quad (58)$$

と表すことができる. [問題 5.6]. したがって, これらの式を組み合わせると, \mathbf{x} を求めるためのコスト関数として,

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\text{argmin}} \mathcal{F} : \quad \mathcal{F} = \beta \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 + \phi(\mathbf{x}) \quad (59)$$

を導くことができる. ここで, 制約条件 $\phi(\mathbf{x})$ は

$$\phi(\mathbf{x}) = \min_{\boldsymbol{\nu}} (\mathbf{x}^T \boldsymbol{\Upsilon}^{-1} \mathbf{x} + \log |\boldsymbol{\Sigma}_y|) = \min_{\boldsymbol{\nu}} \left(\sum_{j=1}^N \frac{x_j^2}{\nu_j} + \log |\boldsymbol{\Sigma}_y| \right) \quad (60)$$

と与えられる.

ここで, なぜスパースな解が得られるかを考察するため, この制約条件 $\phi(\mathbf{x})$ の特徴を調べてみよう. と言っても, この制約条件の式には $\boldsymbol{\nu}$ を含む $\log |\boldsymbol{\Sigma}_y|$ の項が存在するため, この $\phi(\mathbf{x})$ を計算するのは簡単ではない. そこで, 少し話を簡単にするため, 行列 \mathbf{H} の列ベクトルが正規直交系を成す, すなわち, $\mathbf{h}_i^T \mathbf{h}_j = I_{i,j}$ が成り立つと仮定してみよう. この場合

$$\log |\boldsymbol{\Sigma}_y| = \sum_{j=1}^N \log(\beta^{-1} + \nu_j)$$

であるので, 式 (60) を用いれば,

$$\phi(\mathbf{x}) = \min_{\boldsymbol{\nu}} \sum_{j=1}^N \left(\frac{x_j^2}{\nu_j} + \log(\beta^{-1} + \nu_j) \right) \quad (61)$$

を得る. したがって,

$$\phi(\mathbf{x}) = \sum_{j=1}^N \varphi(x_j) \quad (62)$$

と表すことができ, ここで,

$$\varphi(x_j) = \min_{\nu_j} \left(\frac{x_j^2}{\nu_j} + \log(\beta^{-1} + \nu_j) \right) \quad (63)$$

である. この $\varphi(x_j)$ は

$$\varphi(x_j) = \frac{2|x_j|}{\sqrt{x_j^2 + 4\beta^{-1}} + |x_j|} + \log \left(\beta^{-1} + \frac{x_j^2}{2} + \frac{1}{2}|x_j| \sqrt{x_j^2 + 4\beta^{-1}} \right) \quad (64)$$

と表すことができる [問題 5.7] .

式 (64) に示すコスト関数 $\varphi(x)$ の計算結果を図 1 に示す . 同図に $\varphi(x)$ のプロットを , $\beta^{-1} = 0.001$, $\beta^{-1} = 0.01$, $\beta^{-1} = 0.1$ および $\beta^{-1} = 1$ の 4 種類の場合について , 点線で示す . 比較のため , L_1 ノルムの場合の制約条件 $|x|$ のプロットを実線で示す . 図 1 に示されるように , $\varphi(x)$ は L_1 ノルムにおける制約条件 $|x|$ に非常に類似した形をしており , スパースな解を生じることが理解できよう . さらに , $\varphi(x)$ が , β^{-1} すなわちノイズの分散に依存することも見て取れる . すなわち , ノイズ分散が小さい場合には , $\varphi(x)$ は , L_1 ノルム制約より , さらに急峻となり , L_0 ノルム制約により近づく . すなわち , スパースベイズ推定のコスト関数における制約条件は , 固定のものではなく , ノイズの大きさによって急峻さを調節し , ノイズが小さな時には , よりスパースネスを強調した L_0 ノルム制約に近いものになる . 反対にノイズが大きい場合には , $\varphi(x)$ は L_1 ノルム制約に近いものとなり , したがって , スパースネスはより穏やかなものに調節されることになる .