

### 6.4.6 ラプラス近似 (続)

前回はガウス過程の分類問題を扱った。その中で、予測分布を積分の形式で求め、積分を使わない表現を得るためにラプラス近似を導入するところまで行った。今回はラプラス近似を引き続き行い、予測分布の最終形を導く。その後、超パラメータの最適化と、ニューラルネットワークとの関係を考察する。

分類問題のガウス過程の予測分布は次のような形になった。

$$\begin{aligned}
 P(t_{N+1} | t_N) &= \int P(t_{N+1} = 1 | a_{N+1}) P(a_{N+1} | t_N) da_{N+1} \quad \left[ \begin{array}{l} \text{重畳積分の近似公式} \\ (4.153) \text{ を使う。} \end{array} \right] \\
 &\quad \underbrace{\left( \frac{1}{\sigma(a_{N+1})} \right)}_{\text{ヘルムス-イ分布のσ}} \underbrace{\left( \frac{1}{\sigma(a_{N+1})} \right)}_{\text{ガウス分布で近似する。}} \quad (6.76) \\
 &\quad // \\
 &\int P(a_{N+1} | a_N) P(a_N | t_N) da_N \quad \left[ \begin{array}{l} \text{ガウス分布の周辺化の} \\ \text{公式を使う (2.115)} \end{array} \right] \\
 &\quad \underbrace{\left( \frac{1}{\sigma(a_N)} \right)}_{\text{ガウス分布 (4.44) のガウス過程 (6.66) (6.67)}} \underbrace{\left( \frac{1}{\sigma(a_N)} \right)}_{\text{ガウス分布 (4.44) のガウス過程 (6.66) (6.67)}} \quad (6.77) \\
 &\quad // \\
 &\log P(a_N | t_N) \text{ を モーダル付近で 2 次近似する。} \\
 &\quad // \\
 &\log P(a_N) + \log P(t_N | a_N) + \text{const.} \\
 &\quad // \\
 &\underbrace{-\frac{1}{2} a_N^T C_N^{-1} a_N - \frac{1}{2} \ln |C_N| + t_N^T a_N - \sum_{n=1}^N \ln(1 + e^{a_n})}_{\Psi(a_N) \text{ とおく。}} + \text{const.}
 \end{aligned}$$

つまり、 $\Psi(a_N)$  のモードの値とモード付近の2次微分の値が求まると、芋づる式に予測分布の形が得られるのである。まずモードから求める。

$$\begin{aligned}
 \nabla \Psi(a_N) &= -C_N^{-1} a_N + t_N - \left( \sum_{n=1}^N \frac{e^{a_n}}{1 + e^{a_n}} \delta_{1n} \quad \sum_{n=1}^N \frac{e^{a_n}}{1 + e^{a_n}} \delta_{2n} \quad \dots \right)^T \\
 &= -C_N^{-1} a_N + t_N - \left( \frac{e^{a_1}}{1 + e^{a_1}} \quad \frac{e^{a_2}}{1 + e^{a_2}} \quad \dots \quad \frac{e^{a_N}}{1 + e^{a_N}} \right)^T \\
 &= -C_N^{-1} a_N + t_N - \left( \frac{1}{1 + e^{-a_1}} \quad \frac{1}{1 + e^{-a_2}} \quad \dots \quad \frac{1}{1 + e^{-a_N}} \right)^T \\
 &= -C_N^{-1} a_N + t_N - \left( \sigma(a_1) \quad \sigma(a_2) \quad \dots \quad \sigma(a_N) \right)^T \\
 &= -C_N^{-1} a_N + t_N - \Phi(a_N)
 \end{aligned}$$

新しく定義する。

|

$\nabla \Psi(\alpha_N) = 0$  を満たす  $\alpha_N$  を求めたいが、 $\sigma(\alpha_N)$  が  $\alpha_N$  に依存するため陽に求めることができない。そこで反復再重み付け最小二乗法を用いる。

～おさらい：反復再重み付け最小二乗法～

ニュートン・ラフソン法 (a.k.a. ニュートン法) において、ヘッセ行列が更新のたびに更新する場合の更新手順のことを言う。ニュートン法の更新式は次の形である。

$$w^{(new)} = w^{(old)} - H^{-1} \nabla E(w^{old})$$

ここで、ヘッセ行列  $H$  が  $w$  に依存する場合には、更新のたびに  $H$  も更新する必要がある。 $H$  を「重み」と考えると、都度重みを更新していることに相当するため、反復再重み付け最小二乗法と呼ばれる。(ref. Chapter.4/210228輪講.pdf)

ニュートン法に必要な2次微分  $\nabla \nabla \Psi(\alpha_N)$  を求める。これは既に記載したように、ラプラス近似での分散にも再利用される。

$$\begin{aligned} \nabla \nabla \Psi &= -C_N^{-1} - \nabla \sigma(\alpha_N) \\ &= -C_N^{-1} - \left( \frac{\partial}{\partial \alpha_N} \sigma(\alpha_1)^T \quad \frac{\partial}{\partial \alpha_N} \sigma(\alpha_2)^T \quad \dots \quad \frac{\partial}{\partial \alpha_N} \sigma(\alpha_N)^T \right)^T \\ &= -C_N^{-1} - \begin{pmatrix} \frac{d}{d\alpha_N} \sigma(\alpha_1) & 0 & \dots & 0 \\ 0 & \ddots & \ddots & \frac{d}{d\alpha_N} \sigma(\alpha_N) \end{pmatrix} \\ &= -C_N^{-1} - \begin{pmatrix} \sigma(\alpha_1)(1-\sigma(\alpha_1)) & 0 & \dots & 0 \\ 0 & \ddots & \ddots & \sigma(\alpha_N)(1-\sigma(\alpha_N)) \end{pmatrix} \quad \left( \because \frac{d}{dx} \sigma(x) = \sigma(x)(1-\sigma(x)) \right) \\ &= -C_N^{-1} - W_N \end{aligned}$$

得られた勾配とヘッセ行列をニュートン法の更新式に代入すると次の更新式を得る。(演習6.25)

$$\begin{aligned} \alpha_N^{new} &= \alpha_N - (\nabla \nabla \Psi(\alpha_N))^{-1} \nabla \Psi(\alpha_N) \\ &= \alpha_N + (C_N^{-1} + W_N)^{-1} (\mathbb{I}_N - \sigma(\alpha_N) - C_N^{-1} \alpha_N) \\ &= (C_N^{-1} + W_N)^{-1} \{ (C_N^{-1} + W_N) \alpha_N + \mathbb{I}_N - \sigma(\alpha_N) - C_N^{-1} \alpha_N \} \end{aligned}$$

$$\begin{aligned}
 \alpha_n^{new} &= (\underbrace{\mathbb{I} \mathcal{C}_n^{-1} + W_n \mathcal{C}_n \mathcal{C}_n^{-1}}_{\text{作り出す}})^{-1} \{ \mathcal{C}_n^{-1} \alpha_n + W_n \alpha_n + \mathbf{t}_n - \sigma(\alpha_n) - \mathcal{C}_n^{-1} \alpha_n \} \\
 &= \{ (\mathbb{I} + W_n \mathcal{C}_n) \mathcal{C}_n^{-1} \}^{-1} \{ \mathbf{t}_n - \sigma(\alpha_n) + W_n \alpha_n \} \\
 &\quad \downarrow (AB)^{-1} = B^{-1}A^{-1} \\
 &= \mathcal{C}_n (\mathbb{I} + W_n \mathcal{C}_n)^{-1} \{ \mathbf{t}_n - \sigma(\alpha_n) + W_n \alpha_n \} \quad (6.83)
 \end{aligned}$$

これより、ニュートン法（今回は反復再重み付け最小二乗法でもある）の更新式を得られたので繰り返し計算によって、モード  $\alpha_n^*$  と負のヘッセ行列  $H = -\nabla \nabla \Psi = W_n + \mathcal{C}_n^{-1}$  が算出される。これより、目的の確率分布  $P(\alpha_n | \mathbf{t}_n)$  は次のガウス分布となる。（上巻 p.215, (4.132)式）

$$P(\alpha_n | \mathbf{t}_n) = \mathcal{N}(\alpha_n | \alpha_n^*, H^{-1})$$

最後に、予測分布が得られるまでのロードマップを振り返る。我々は(6.77) 式の  $p(\alpha_n | \mathbf{t}_n)$  を求めることができたので、予測分布  $p(\mathbf{t}_{n+1} | \mathbf{t}_n)$  を得ることができた。

$$\begin{aligned}
 p(\mathbf{t}_{n+1} | \mathbf{t}_n) &= \int \underbrace{P(\mathbf{t}_{n+1} = 1 | \alpha_{n+1})}_{\substack{\text{ベルヌーイ分布の密度} \\ \sigma(\alpha_{n+1})}} \underbrace{P(\alpha_{n+1} | \mathbf{t}_n)}_{\substack{\text{ガウス分布に近似する} \\ (6.76)}} d\alpha_{n+1} \quad \left. \begin{array}{l} \text{重畳積分の近似公式} \\ (4.153) \text{を使う} \end{array} \right\} \\
 &\quad // \\
 &= \int \underbrace{P(\alpha_{n+1} | \alpha_n)}_{\substack{\text{ガウス分布} \\ (\because \text{回帰のガウス過程}) \\ (6.66) (6.67)}} \underbrace{P(\alpha_n | \mathbf{t}_n)}_{\substack{\text{ガウス分布} \\ (\because \text{ラプラス近似})}} d\alpha_n \quad (6.77) \quad \left. \begin{array}{l} \text{ガウス分布の周辺化の} \\ \text{公式を使う} (2.115) \end{array} \right\}
 \end{aligned}$$

※ 最終形を露わな形で書きたい気持ちはあったけど、答え合わせができないのでここで止めます…最後にハイパーパラメータ  $\theta$  の最適化を扱う。基本的なアイデアとしては「6.4.3 超パラメータの学習」同様に、尤度関数を最大にする最尤推定量を求める。尤度関数は次のように表される。

$$p(\mathbf{t}_N | \theta) = \int \underbrace{P(\mathbf{t}_N | \alpha_N)}_{\text{ベルヌーイ}} \underbrace{P(\alpha_N | \theta)}_{\text{ガウス}} d\alpha_N$$

(感想)

6.90 式の導出できませんでした。困った点は次の3つ

1. (6.90) に  $|I| = |W_N + C_N^{-1}|$  や  $\sigma_N^*$  が現れているので  $p(\sigma_N^* | \epsilon_N)$  を使いそうだが、積分形式以外に  $p(\sigma_N^* | \epsilon)$  の分布が現れる形に変形できない。
2. 積分の代わりにラプラス近似を行うとこのとだが、 $\epsilon_N$  の分布はベルヌーイ分布なのでガウス分布に近似して良いのか怪しい。
3. 積分の中身の  $p(\sigma_N^* | \theta)$  にのみラプラス近似を適用しように思えたが、重畳積分の近似公式 (4.151 - 4.155) の形が出てこないで違うっぽい。

6.90式が次のように導出できたとする。

$$\ln p(\epsilon_N | \theta) = \Psi(\sigma_N^*) - \frac{1}{2} \ln |W_N + C_N^{-1}| + \frac{N}{2} \ln(2\pi)$$

$$\text{ここで } \Psi(\sigma_N^*) = \ln p(\sigma_N^* | \theta) + \ln p(\epsilon_N | \sigma_N^*)$$

この  $\theta$  に関する停留点を求めることで、超パラメータの最適化を行う。

$$\left( \frac{\partial \ln p(\epsilon_N | \theta)}{\partial \theta_i} \right) = \sum_{n=1}^N \underbrace{\frac{\partial \Psi(\sigma_N^*)}{\partial \sigma_n^*} \cdot \frac{\partial \sigma_n^*}{\partial \theta_i}}_{=0 \text{ に似る?}} - \frac{1}{2} \sum_{n=1}^N \underbrace{\frac{\partial \ln |W_N + C_N^{-1}|}{\partial \sigma_n^*} \cdot \frac{\partial \sigma_n^*}{\partial \theta_i}}_{(6.93) \text{ 式で扱う。}} + 0$$

(2理由)  $\sigma_n^*$  は  $\frac{\partial \Psi(\sigma_n^*)}{\partial \sigma_n^*} = 0$  より得た

$$= -\frac{1}{2} \sum_{n=1}^N \underbrace{\frac{\partial \ln |W_N + C_N^{-1}|}{\partial \sigma_n^*}}_{(C.22) \frac{\partial}{\partial x} \ln |A| = \text{Tr}(A^{-1} \frac{\partial A}{\partial x}) \text{ を使う}}$$

$$= -\frac{1}{2} \sum_{n=1}^N \text{Tr} \left\{ (W_N + C_N^{-1})^{-1} \frac{\partial (W_N + C_N^{-1})}{\partial \sigma_n^*} \right\} \frac{\partial \sigma_n^*}{\partial \theta_i}$$

$C_N$  は  $\{C_n\}_{i,j} = \delta(x_i, x_j) + \ell^{-1} \delta_{ij}$  により、 $\sigma_n^*$  に依存しない。

$$= -\frac{1}{2} \sum_{n=1}^N \text{Tr} \left\{ (W_N + C_N^{-1})^{-1} \frac{\partial W_N}{\partial \sigma_n^*} \right\} \frac{\partial \sigma_n^*}{\partial \theta_i}$$

$$= \underbrace{\{C_N^{-1} (I + C_N W_N)\}^{-1} C_N}_{= (I + C_N W_N)^{-1} C_N} \frac{\partial}{\partial \sigma_n^*} \begin{pmatrix} \sigma_n^*(1 - \sigma_n^*) & & 0 \\ 0 & \dots & 0 \\ & & \sigma_n^*(1 - \sigma_n^*) \end{pmatrix}$$

ただし、 $\sigma_n^* = \sigma(\sigma_n^*)$  とおく。

$$= \begin{pmatrix} 0 & \dots & \frac{\partial}{\partial \sigma_n^*} \sigma_n^*(1 - \sigma_n^*) & 0 \\ & & & \ddots \\ 0 & & & 0 \end{pmatrix} \quad (n, n) \text{ の成分だけ}$$

$\text{Tr}()$  とすると  $(n, n)$  成分のみ残る。

$$= -\frac{1}{2} \sum_{n=1}^N \left[ (I + C_N W_N)^{-1} C_N \right]_{nn} \frac{\partial}{\partial \sigma_n^*} \sigma_n^*(1 - \sigma_n^*) \frac{\partial \sigma_n^*}{\partial \theta_i}$$

$$\begin{aligned}
 \left( \frac{\partial \ln P(t_n | \theta)}{\partial \theta_i} \right) &= -\frac{1}{2} \sum_{n=1}^N \left[ (I + C_n W_n)' C_n \right]_{nn} \frac{\partial}{\partial a_n^*} (\sigma_n^* - \sigma_n^{*2}) \frac{\partial a_n^*}{\partial \theta_i} \\
 &= -\frac{1}{2} \sum_{n=1}^N \left[ (I + C_n W_n)' C_n \right]_{nn} \sigma_n^* (1 - \sigma_n^*) (1 - 2\sigma_n^*) \frac{\partial a_n^*}{\partial \theta_i} \quad (6.92)
 \end{aligned}$$

(6.92)より示す。

(感想)

勾配を (6.92) で示してしまっただけ、教科書だと (6.91) 式を挟んでいる点で異なる。(6.91) をどのように導出したかわからなかったなので相談したい。(6.92) で現れた  $\frac{\partial a_n^*}{\partial \theta_i}$  を更に計算する。

$$\begin{aligned}
 \frac{\partial a_n^*}{\partial \theta_i} &= \frac{\partial}{\partial \theta_i} \{ C_n (t_n - \sigma_n) \} & (\because 6.84; a_n^* = C_n (t_n - \sigma_n)) \\
 &= \frac{\partial C_n}{\partial \theta_i} (t_n - \sigma_n) + C_n W_n \frac{\partial a_n^*}{\partial \theta_i} & (6.93)
 \end{aligned}$$

(6.93)より示す。

$$\Leftrightarrow \frac{\partial a_n^*}{\partial \theta_i} + C_n W_n \frac{\partial a_n^*}{\partial \theta_i} = \frac{\partial C_n}{\partial \theta_i} (t_n - \sigma_n)$$

$$\Leftrightarrow (I + C_n W_n) \frac{\partial a_n^*}{\partial \theta_i} = \frac{\partial C_n}{\partial \theta_i} (t_n - \sigma_n)$$

$$\Leftrightarrow \frac{\partial a_n^*}{\partial \theta_i} = (I + C_n W_n)^{-1} \frac{\partial C_n}{\partial \theta_i} (t_n - \sigma_n) \quad (6.94)$$

ここ教科書と異なす

全く教科書通りにうまく行かなかったが、とりあえず勾配が得られたのでこれを元に何らかの最適化法を適用して最適なハイパーパラメータを取得する。

## 6.4.7 ニューラルネットワーク (NN) との関係

カラム的な内容なので、各々読むこととする。主張としては次の通り。

- NNはデータサイズに応じてユニット数を変える必要があるが、ガウス過程にはない。
- NNの事前分布を広いクラス (?) に設定すると、予測分布がガウス分布に近くなる。