

パターン認識と機械学習

第3章 線形回帰モデル

橋本龍二

January 6, 2021

1 線形基底関数モデル

入力変数 $\begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix}$ が与えられた時、その線形結合を取ることで、目標変数 t の予測値 $y(\mathbf{x}, \mathbf{w})$ を

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D$$

と構成する。(線形回帰)

次に、入力変数 \mathbf{x} に対して非線形な基底関数 $\phi(\mathbf{x}), \dots, \phi(\mathbf{x})$ を用いて、

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$

なるモデルを構成する。 $\phi_0(\mathbf{x}) = 1$ とすれば、

$$\begin{aligned} y(\mathbf{x}, \mathbf{w}) &= \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) \\ &= {}^T \mathbf{w} \phi(\mathbf{x}) \end{aligned}$$

と書ける。ただし、 $\mathbf{w} = \begin{pmatrix} w_0 \\ \vdots \\ w_{M-1} \end{pmatrix}$ 、 $\phi = \begin{pmatrix} \phi_0 \\ \vdots \\ \phi_{M-1} \end{pmatrix}$

1. ガウス基底関数

$$\phi_j(x) = \exp\left\{-\frac{(x - \mu_j)^2}{2s^2}\right\}$$

2. シグモイド基底関数

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right), \quad \sigma(a) = \frac{1}{1 + \exp(-a)}$$

以下、目標変数 t は一次元とする。

1.1 最尤推定と最小二乗法

条件付きガウスノイズ分布のもとでの線形モデルを考える。すなわち、

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon, \quad \epsilon \sim N(0, \frac{1}{\beta})$$

として、

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = N(t|y(\mathbf{x}, \mathbf{w}), \frac{1}{\beta})$$

を考える。 N 個の入力データ集合 $\mathbf{X} = (\mathbf{x}_1 \dots \mathbf{x}_N)$ とその目標値のデータ集合 $\mathbf{t} = \begin{pmatrix} t_1 \\ \vdots \\ t_N \end{pmatrix}$ を考える。 $y(\mathbf{x}, \mathbf{w}) = {}^T \mathbf{w} \phi(\mathbf{x})$ とすると、尤度関数は

$$\begin{aligned} p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) &= \prod_{n=1}^N N(t_n|y(\mathbf{x}, \mathbf{w}), \frac{1}{\beta}) \\ &= \prod_{n=1}^N N(t_n|{}^T \mathbf{w} \phi(\mathbf{x}), \frac{1}{\beta}) \\ &= (2\pi)^{-\frac{N}{2}} \times \beta^{\frac{N}{2}} \times \exp\left[-\frac{\beta}{2} \sum_{n=1}^N \{t_n - {}^T \mathbf{w} \phi(\mathbf{x}_n)\}^2\right] \\ \log p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) &= \frac{N}{2} \log \beta - \frac{N}{2} \log(2\pi) - \frac{\beta}{2} \sum_{n=1}^N \{t_n - {}^T \mathbf{w} \phi(\mathbf{x}_n)\}^2 \end{aligned}$$

最尤推定によって \mathbf{w} を求める。対数尤度の最大化は二乗和誤差関数の最小化である。以下、 \mathbf{x} を省略する。 $(\mathbf{x}$ の分布のモデル化は目指していない)
対数尤度関数を微分して 0 と置く。

$$\begin{aligned} 0 = \frac{\partial}{\partial \mathbf{w}} \log p(\mathbf{t}|\mathbf{w}, \beta) &= \beta \sum_{n=1}^N \{t_n - {}^T \mathbf{w} \phi(\mathbf{x}_n)\} {}^T \phi(\mathbf{x}_n) \\ \therefore 0 &= \sum_{n=1}^N t_n {}^T \phi(\mathbf{x}_n) - {}^T \mathbf{w} \left\{ \sum_{n=1}^N \phi(\mathbf{x}_n) {}^T \phi(\mathbf{x}_n) \right\} \end{aligned}$$

ここで、

$$\begin{aligned}
\sum_{n=1}^N t_n {}^T \phi(\mathbf{x}_n) &= \sum_{n=1}^N t_n (\phi_0(\mathbf{x}_n) \ \dots \ \phi_{M-1}(\mathbf{x}_n)) \\
&= (t_1 \ \dots \ t_N) \begin{pmatrix} \phi_0(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix} \\
\sum_{n=1}^N \phi(\mathbf{x}_n) {}^T \phi(\mathbf{x}_n) &= \sum_{n=1}^N \begin{pmatrix} \phi_0(\mathbf{x}_n) \\ \vdots \\ \phi_{M-1}(\mathbf{x}_n) \end{pmatrix} (\phi_0(\mathbf{x}_n) \ \dots \ \phi_{M-1}(\mathbf{x}_n)) \\
&= \begin{pmatrix} \phi_0(\mathbf{x}_1) & \dots & \phi_0(\mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ \phi_{M-1}(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix} \begin{pmatrix} \phi_0(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}
\end{aligned}$$

であるから、

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

と置けば(計画行列)、

$$\begin{aligned}
0 &= {}^T \mathbf{t} \Phi - {}^T \mathbf{w} {}^T \Phi \Phi \\
\therefore {}^T \mathbf{w} &= {}^T \mathbf{t} \Phi ({}^T \Phi \Phi)^{-1} \\
\therefore \mathbf{w}_{ML} &= {}^T \{ {}^T \mathbf{t} \Phi ({}^T \Phi \Phi)^{-1} \} \\
&= \{ {}^T ({}^T \Phi \Phi) \}^{-1} {}^T ({}^T (\mathbf{t} \Phi)) \\
&= ({}^T \Phi \Phi)^{-1} {}^T \Phi \mathbf{t}
\end{aligned}$$

1.2 最小二乗法の幾何学

$\mathbf{t} = \begin{pmatrix} t_1 \\ \vdots \\ t_N \end{pmatrix}$ は N 次元数ベクトル空間の要素。 Φ の j 番目の列ベクトルを φ_j と置く。

$$\varphi_j = \begin{pmatrix} \phi_j(\mathbf{x}_1) \\ \vdots \\ \phi_j(\mathbf{x}_N) \end{pmatrix} \in \mathbb{R}^N, \quad j = 1, \dots, M$$

$M < N$ ならば、 M 本のベクトル φ_j は M 次元の線形部分空間 S を張る。次に、推定値 $y(\mathbf{x}_j, \mathbf{w})$ を並べたベクトルを \mathbf{y} と置くと、

$$\begin{aligned}\mathbf{y} &= \begin{pmatrix} y(\mathbf{x}_1, \mathbf{w}) \\ \vdots \\ y(\mathbf{x}_j, \mathbf{w}) \end{pmatrix} \\ &= \begin{pmatrix} w_0\phi_0(\mathbf{x}_1) + \dots + w_{M-1}\phi_{M-1}(\mathbf{x}_1) \\ \vdots \\ w_0\phi_0(\mathbf{x}_N) + \dots + w_{M-1}\phi_{M-1}(\mathbf{x}_N) \end{pmatrix} \\ &= \mathbf{w}_o\varphi_0 + \dots + \mathbf{w}_{M-1}\varphi_{M-1}\end{aligned}$$

より、 \mathbf{y} はベクトル φ_j の任意の線形結合より、 S 内の要素である。今、二乗誤差；

$$\sum_{n=1}^N \{t_n - y(\mathbf{x}_n, \mathbf{w})\}^2$$

とは \mathbf{y} と \mathbf{t} の二乗ユークリッド距離である。すなわち、二乗和誤差の最小化とは $\mathbf{t} \in \mathbb{R}^N$ の \mathbb{R}^N の部分空間 $S = \mathbb{R}^M$ への正射影 \mathbf{y} を選ぶことである。

1.3 逐次学習

データ点を一度に一つだけ用いてモデルのパラメータを順次学習していく方法を逐次学習という。パターン n が与えられた時、確率的勾配降下法ではパラメータベクトル \mathbf{w} を

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n$$

を用いて更新する。

1.4 正則化最小二乗法

誤差関数に正則化項を加えて、

$$\frac{1}{2} \sum_{n=1}^N \{t_n - {}^T \mathbf{w} \boldsymbol{\phi}(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

を最小化することで、過学習を防ぐ。 $q = 1$ のとき lasso、 $q = 2$ のとき二次正則化項という。lasso では、いくつかの係数 w_j が 0 になり、疎な解が得られる。

2 バイアス-バリアンス分解

入力 \mathbf{x} に対して t の値に対する特定の推定値 $y(\mathbf{x})$ を選ぶことで損失 $L(t, y(\mathbf{x}))$ を被るとしよう。期待損失 $E[L]$ は、

$$E[L] = \int \int L(t, y(\mathbf{x})) p(\mathbf{x}, t) d\mathbf{x} dt$$

今、 $L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$ とすれば、

$$\begin{aligned} \{y(\mathbf{x}) - t\}^2 &= \{y(\mathbf{x}) - E[t|\mathbf{x}] + E[t|\mathbf{x}] - t\}^2 \\ &= \{y(\mathbf{x}) - E[t|\mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - E[t|\mathbf{x}]\}(E[t|\mathbf{x}] - t) + (E[t|\mathbf{x}] - t)^2 \end{aligned}$$

なので、期待二乗損失 $E[L]$ は

$$\begin{aligned} E[L] &= \int \int \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt \\ &= \int \{y(\mathbf{x}) - E[t|\mathbf{x}]\}^2 \left\{ \int p(\mathbf{x}, t) dt \right\} d\mathbf{x} \\ &\quad + 2 \int \{y(\mathbf{x}) - E[t|\mathbf{x}]\} \left\{ \int E[t|\mathbf{x}] \int p(\mathbf{x}, t) dt - \int t p(t|\mathbf{x}) dt \right\} p(\mathbf{x}) d\mathbf{x} \\ &\quad + \int \int (E[t|\mathbf{x}] - t)^2 p(t, \mathbf{x}) dt d\mathbf{x} \\ &= \int \{y(\mathbf{x}) - E[t|\mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} \\ &\quad + 2 \int \{y(\mathbf{x}) - E[t|\mathbf{x}]\} \{E[t|\mathbf{x}] p(\mathbf{x}) - E[t|\mathbf{x}] p(\mathbf{x})\} d\mathbf{x} \\ &\quad + \int \left\{ \int (t - E[t|\mathbf{x}])^2 p(t|\mathbf{x}) dt \right\} p(\mathbf{x}) d\mathbf{x} \\ &= \int \{y(\mathbf{x}) - E[t|\mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int var[t|\mathbf{x}] p(\mathbf{x}) d\mathbf{x} \end{aligned}$$

の形で書くことができる。第二項は $y(\mathbf{x})$ に依存しないので、 $h(\mathbf{x}) = E[t|\mathbf{x}]$ が期待二乗損失を最小化する回帰関数である。第二項は t の分布の分散を \mathbf{x} に関して平均したもので、目標データが本質的に持つ変動(ノイズ)である。

<回帰の考え方>

入力 \mathbf{x} に対して t の推定値を選ぶ上で、期待損失を最小化するのは $E[t|\mathbf{x}]$ である。

→ $E[t|\mathbf{x}]$ を求めたい。

→ $p(t|\mathbf{x}) = N(t|y(\mathbf{x}, \mathbf{w}), \frac{1}{\beta})$ とすることで、 $E[t|\mathbf{x}]$ を $y(\mathbf{x}, \mathbf{w})$ によってモデル

化

ここで、頻度主義的には \mathbf{w} を点推定する。ベイズの立場では、 \mathbf{w} の事後分布を求める。

さて、今無限のデータ点を取れるのであれば推定値 $y(\mathbf{x})$ は理想の精度で $h(\mathbf{x})$ に近づけることができる。しかし、現実には有限のデータ集合 \mathcal{D} を得ることが出来ないので、それに対する予測関数 $y(\mathbf{x}; \mathcal{D})$ を構成する。 $y(\mathbf{x}; \mathcal{D}) = h(\mathbf{x})$ は厳密には成り立たない。ではその誤差 $\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2$ はどのように変動するだろうか。

$$\begin{aligned}\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 &= \{y(\mathbf{x}; \mathcal{D}) - E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &\quad + 2\{y(\mathbf{x}; \mathcal{D}) - E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}\{E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}\end{aligned}$$

これ全体に \mathcal{D} に関して期待値を取る。 $y(\mathbf{x}; \mathcal{D})$ でない項は定数であるから、

$$E_{\mathcal{D}}[\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] = \{E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 + E_{\mathcal{D}}[\{y(\mathbf{x}; \mathcal{D}) - E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]$$

右辺第一項は(バイアス)²、第二項はバリアンスである。

- バイアス： $E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})$ 全てのデータ集合 \mathcal{D} の取り方に関する予測値の平均が理想的な回帰関数 $h(\mathbf{x})$ からどのくらい離れているか
- バリアンス： $E_{\mathcal{D}}[\{y(\mathbf{x}; \mathcal{D}) - E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]$ 各のデータ集合 \mathcal{D} に対する解 $y(\mathbf{x}; \mathcal{D})$ の分散

以上より、期待損失は二乗バイアスとバリアンス、ノイズに分解できる。すなわち、

$$\begin{aligned}E[L] &= \int \{E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} \\ &\quad + \int E_{\mathcal{D}}[\{y(\mathbf{x}; \mathcal{D}) - E_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x} \\ &\quad + \int var[t|\mathbf{x}] p(\mathbf{x}) d\mathbf{x}\end{aligned}$$

バイアスとバリアンスはモデルの複雑さに依存する。単純なモデルではバリアンスは小さくなるがバイアスは大きくなる。複雑なモデルではバイアスは小さくなるがバリアンスは大きくなる。

3 ベイズ線形回帰

線形回帰モデルをベイズ的に取り扱うことによって、最尤推定の過学習を回避するとともに訓練データだけからモデルの複雑さを自動的に決定する。

3.1 パラメータの分布

$y(\mathbf{x}, \mathbf{w}) = {}^T \mathbf{w} \boldsymbol{\phi}(\mathbf{x})$ とすると、尤度関数は

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N N(t_n | {}^T \mathbf{w} \boldsymbol{\phi}(\mathbf{x}_n), \frac{1}{\beta})$$

と、 \mathbf{w} の二次関数の指數なので、共役事前分布は

$$p(\mathbf{w}) = N(\mathbf{w} | \mathbf{m}_0, S_0)$$

で与えるのが適切である。 $\mathbf{a} = \begin{pmatrix} \mathbf{t} \\ \mathbf{w} \end{pmatrix}$ とおく。

$$\begin{aligned} \log p(\mathbf{a}) &= \sum_{n=1}^N {}^T (t_n - {}^T \mathbf{w} \boldsymbol{\phi}(\mathbf{x}_n)) \beta (t_n - {}^T \mathbf{w} \boldsymbol{\phi}(\mathbf{x}_n)) \\ &\quad + {}^T (\mathbf{w} - \mathbf{m}_0) S_0^{-1} (\mathbf{w} - \mathbf{m}_0) \\ &= {}^T (\mathbf{t} - \Phi \mathbf{w}) \beta (\mathbf{t} - \Phi \mathbf{w}) + {}^T (\mathbf{w} - \mathbf{m}_0) S_0^{-1} (\mathbf{w} - \mathbf{m}_0) \\ &= {}^T \mathbf{w} {}^T \Phi \beta \Phi \mathbf{w} + {}^T \mathbf{w} S_0^{-1} \mathbf{w} \\ &\quad - {}^T \mathbf{w} {}^T \Phi \beta \mathbf{t} - {}^T \mathbf{t} \beta \Phi \mathbf{w} + {}^T \mathbf{t} \beta \mathbf{t} \\ &\quad + {}^T \mathbf{w} S_0^{-1} \mathbf{m}_0 - {}^T \mathbf{m}_0 S_0^{-1} \mathbf{w} \\ &= {}^T \begin{pmatrix} \mathbf{w} \\ \mathbf{t} \end{pmatrix} \begin{pmatrix} {}^T \Phi \beta \Phi + S_0^{-1} & - {}^T \Phi \beta \\ - \beta \Phi & \beta \end{pmatrix} \begin{pmatrix} \mathbf{w} \\ \mathbf{t} \end{pmatrix} - {}^T \begin{pmatrix} \mathbf{w} \\ \mathbf{t} \end{pmatrix} \begin{pmatrix} S_0^{-1} \mathbf{m}_0 \\ \mathbf{0} \end{pmatrix} - {}^T \begin{pmatrix} S_0^{-1} \mathbf{m}_0 \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{w} \\ \mathbf{t} \end{pmatrix} \\ &= {}^T \mathbf{a} \Lambda \mathbf{a} - {}^T \mathbf{a} B - {}^T B \mathbf{a} \\ &= {}^T (\mathbf{a} - \Lambda^{-1} B) \Lambda (\mathbf{a} - \Lambda^{-1} B) \end{aligned}$$

ただし、

$$\Lambda = \begin{pmatrix} {}^T \Phi \beta \Phi + S_0^{-1} & - {}^T \Phi \beta \\ - \beta \Phi & \beta \end{pmatrix}, \quad B = \begin{pmatrix} S_0^{-1} \mathbf{m}_0 \\ \mathbf{0} \end{pmatrix}$$

と置いた。よって、

$$\Sigma_{\mathbf{w}|\mathbf{t}} = (S_0^{-1} + {}^T\Phi\beta\Phi)^{-1}$$

条件付き平均を求める。

$$\begin{aligned}\Lambda^{-1} &= \begin{pmatrix} S_0 & S_0 {}^T\Phi \\ \Phi S_0 & \beta^{-1} + \Phi S_0 {}^T\Phi \end{pmatrix} \\ \therefore \boldsymbol{\mu}_a &= \begin{pmatrix} S_0 & S_0 {}^T\Phi \\ \Phi S_0 & \beta^{-1} + \Phi S_0 {}^T\Phi \end{pmatrix} \begin{pmatrix} S_0^{-1} \mathbf{m}_0 \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{m}_0 \\ \Phi \mathbf{m}_0 \end{pmatrix}\end{aligned}$$

より、

$$\begin{aligned}\boldsymbol{\mu}_{\mathbf{w}|\mathbf{t}} &= \Sigma_{\mathbf{w}|\mathbf{t}} \{ ({}^T\Phi\beta\Phi + S_0^{-1}) \mathbf{m}_0 + {}^T\Phi\beta(\mathbf{t} - \Phi \mathbf{m}_0) \\ &= \Sigma_{\mathbf{w}|\mathbf{t}} (S_0^{-1} \mathbf{m}_0 + {}^T\Phi\beta \mathbf{t})\end{aligned}$$

事前分布を $p(\mathbf{w}|\alpha) = N(\mathbf{w}|\mathbf{0}, \alpha^{-1}I)$ と置けば、

$$\begin{aligned}\Sigma_{\mathbf{w}|\mathbf{t}} &= (\alpha I + \beta {}^T\Phi\Phi)^{-1} \\ \boldsymbol{\mu}_{\mathbf{w}|\mathbf{t}} &= \beta \Sigma_{\mathbf{w}|\mathbf{t}} {}^T\Phi \mathbf{t}\end{aligned}$$

3.2 予測分布

訓練データが得られた時、新しい t の値を予測する。

$$\begin{aligned}p(t|\mathbf{t}, \alpha, \beta) &= \int p(t, \mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \int p(\mathbf{w}|\mathbf{t}, \alpha, \beta) p(t|\mathbf{w}, \beta) d\mathbf{w}\end{aligned}$$

より、「 \mathbf{w} の分布」と「平均が \mathbf{w} に依存する t の条件付き分布」が与えられたときの t の周辺分布なので、 μ_t, Σ_t は、

$$\begin{aligned}\mu_t &= {}^T \boldsymbol{\mu}_{\mathbf{w}|\mathbf{t}} \phi(\mathbf{x}) \\ \Sigma_t &= \frac{1}{\beta} + \phi(\mathbf{x}) \Sigma_{\mathbf{w}|\mathbf{t}} {}^T \phi(\mathbf{x})\end{aligned}$$

である。予測分布の平均を t の推定値として採用する。

3.3 等価カーネル

事後平均 $\mu_{w|t}$ を w の推定値として採用する。

$$\begin{aligned} y(\mathbf{x}, \mu_{w|t}) &= {}^T \boldsymbol{\mu}_{w|t} \boldsymbol{\phi}(\mathbf{x}) \\ &= {}^T (\beta \Sigma_{w|t} {}^T \Phi \mathbf{t}) \boldsymbol{\phi}(\mathbf{x}) \\ &= \beta {}^T \mathbf{t} \Phi \Sigma_{w|t} \boldsymbol{\phi}(\mathbf{x}) \\ &= \beta {}^T \boldsymbol{\phi}(\mathbf{x}) \Sigma_{w|t} {}^T \Phi \mathbf{t} \end{aligned}$$

(最終行について、 y はスカラーなので全体に転置をとって良い)

${}^T \Phi \mathbf{t} = \sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) t_n$ (「行列 × ベクトル」は列ベクトルの一次結合で表せる。) なので、

$$\begin{aligned} y(\mathbf{x}, \mu_{w|t}) &= \sum_{n=1}^N \beta {}^T \boldsymbol{\phi}(\mathbf{x}) \Sigma_{w|t} \boldsymbol{\phi}(\mathbf{x}_n) t_n \\ &= \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n \end{aligned}$$

(\mathbf{x} はテストデータのあるベクトル。 \mathbf{x}_n は N 個の訓練データの n 番目のベクトル)

このように、点 \mathbf{x} に対する予測分布の平均は訓練データの目標変数 t_n の線形結合で与えられる。

ここで、二本の N 次元ベクトル \mathbf{x}, \mathbf{x}' を引数にもつ関数

$$k(\mathbf{x}, \mathbf{x}') = \beta {}^T \boldsymbol{\phi}(\mathbf{x}) \Sigma_{w|t} \boldsymbol{\phi}(\mathbf{x}')$$

を平滑化行列または等価カーネルという。等価カーネルは重みの役割を果たす。例えば p158 の図 3.11 では $x = 0$ としたときの x' の関数としての等価カーネルがプロットされている。 $x' = 0$ 付近の重みが重要視されていることがわかる。ベイズ線形回帰では、訓練データの目標値の線形結合によって予測を行う。この際、訓練データのうち予測したいデータ点 \mathbf{x} に近い点の目標データが重視される、ということ。

ガウス過程を用いることで、基底関数を用いずに、近傍点での相関を強く、より離れた点への相関は小さい局所的なカーネルを定義することで予測ができる。

4 ガウスモデル比較

データがモデル $\{\mathcal{M}_i\}$, ($i = 1, \dots, L$) のどれかにしたがって生成されているとする (モデルは \mathcal{D} 上の確率分布)。訓練データ集合 \mathcal{D} が与えられた時、モ

ルの事後分布は次で与えられる。

$$p(\mathcal{M}_i|\mathcal{D}) \propto p(\mathcal{M}_i)p(\mathcal{D}|\mathcal{M}_i)$$

事前分布 $p(\mathcal{M}_i)$ は大体分からないので、全てのモデルの事前確率は等しいとする。この時、モデルエビデンス $p(\mathcal{D}|\mathcal{M}_i)$ が重要な働きをする。モデルエビデンスは、

$$p(\mathcal{D}|\mathcal{M}_i) = \int p(\mathcal{D}, \theta|\mathcal{M}_i)d\theta$$

と、モデル \mathcal{M}_i の空間で θ を周辺化した尤度関数と見ることができるため、周辺尤度とも呼ばれる。また、モデル \mathcal{M}_i の \mathcal{M}_i の \mathcal{M}_j に対するベイズ因子を

$$\frac{p(\mathcal{D}|\mathcal{M}_i)}{p(\mathcal{D}|\mathcal{M}_j)}$$

と定義する。

一旦モデルの事後分布がわかると、予測分布は

$$\begin{aligned} p(t|\mathbf{x}, \mathcal{D}) &= \sum_{i=1}^L p(t, \mathcal{M}_i|\mathbf{x}, \mathcal{D}) \\ &= \sum_{i=1}^L p(\mathcal{M}_i|\mathbf{x}, \mathcal{D})p(t|\mathbf{x}, \mathcal{M}_i, \mathcal{D}) \end{aligned}$$

と書ける(混合分布)。

w を一次元とする。事後分布は $p(w|\mathcal{D}, \mathcal{M}_i) \propto p(w|\mathcal{M}_i)p(\mathcal{D}|w, \mathcal{M}_i)$ なので、周辺尤度は

$$\begin{aligned} p(\mathcal{D}|\mathcal{M}_i) &= \int p(\mathcal{D}, w|\mathcal{M}_i)dw \\ &= \int p(\mathcal{D}|w, \mathcal{M}_i)p(w|\mathcal{M}_i)dw \\ &\simeq \frac{1}{\Delta w_{prior}} \int p(\mathcal{D}|w, \mathcal{M}_i)dw \\ &\simeq \frac{1}{\Delta w_{prior}} \times p(\mathcal{D}|w_{MAP})\Delta w_{posterior} \\ \log p(\mathcal{D}|\mathcal{M}_i) &\simeq \log p(\mathcal{D}|w_{MAP}) + \log \frac{\Delta w_{posterior}}{\Delta w_{prior}} \end{aligned}$$

第一項は尤もな w によるデータへのフィッティング度を表し、第二項はモデルの複雑さに対する罰則項(負)である。

では、パラメータ M 個を幾つにするかを選ぶ問題を考える。 $w_{posterior}/w_{prior}$ が全てのパラメータで等しいとすると、

$$\log p(\mathcal{D}|\mathcal{M}_i) \simeq \log p(\mathcal{D}|w_{MAP}) + M \log\left(\frac{\Delta w_{posterior}}{\Delta w_{prior}}\right)$$

となる。実際に、周辺尤度を最大化することで中くらいの複雑さのモデルが選ばれる。単純すぎるモデルでは \mathcal{D}_0 を生成する確率 $p(\mathcal{D}_0|\mathcal{M}_i)$ は低くなる(限られたデータ集合しか作り出せない)し、複雑すぎるモデルではありうるデータの組み合わせが多すぎてピンポイントで \mathcal{D}_0 を生成する確率はやはり小さくなる。このように、ベイズ線形回帰では(事前分布が単純、など一定の条件のもとで)パラメータの事後分布をパラメータで積分することで、解析的にモデル選択を行うことができる。

5 エビデンス近似

訓練データ \mathbf{t} がえられたとき、新しい入力 \mathbf{x} に対して目標変数 t の値を予測する。事前分布の精度 α とガウスノイズ β に事前分布を導入する。予測分布は、

$$\begin{aligned} p(t|\mathbf{t}) &= \int \int \int p(t, \mathbf{w}, \alpha, \beta | \mathbf{t}) d\mathbf{w} d\alpha d\beta \\ &= \int \int \int p(\alpha, \beta | \mathbf{t}) p(\mathbf{w} | \alpha, \beta, \mathbf{t}) p(t | \mathbf{w}, \beta, \mathbf{t}) d\mathbf{w} d\alpha d\beta \end{aligned}$$

5.1 エビデンス関数の評価

$p(\mathbf{w} | \mathbf{t}, \alpha, \beta) \propto p(\mathbf{w} | \alpha) p(\mathbf{t} | \mathbf{w}, \beta)$ より、事後分布をパラメータで周辺化することで α, β の周辺尤度 $p(\mathbf{t} | \alpha, \beta)$ が次のように得られる。

$$\begin{aligned} p(\mathbf{t} | \alpha, \beta) &= \int p(\mathbf{t}, \mathbf{w} | \alpha, \beta) d\mathbf{w} \\ &= \int p(\mathbf{w} | \alpha) p(\mathbf{t} | \mathbf{w}, \beta) d\mathbf{w} \end{aligned}$$

ここで、

$$\begin{aligned}
 p(\mathbf{t}|\mathbf{w}, \beta) &= \left(\frac{\beta}{2\pi}\right)^{N/2} \exp\left[-\frac{\beta}{2} \sum_{n=1}^N \{t_n - {}^T \mathbf{w} \boldsymbol{\phi}(\mathbf{x}_n)\}^2\right] \\
 &= \left(\frac{\beta}{2\pi}\right)^{N/2} \exp\left\{-\frac{\beta}{2} {}^T (\mathbf{t} - \Phi \mathbf{w}) (\mathbf{t} - \Phi \mathbf{w})\right\} \\
 p(\mathbf{w}|\alpha) &= \frac{\alpha^{M/2}}{(2\pi)^{M/2}} \exp\left(-\frac{1}{2} \alpha {}^T \mathbf{w} \mathbf{w}\right)
 \end{aligned}$$

より、

$$p(\mathbf{t}|\alpha, \beta) = \left(\frac{\beta}{2\pi}\right)^{N/2} \left(\frac{\alpha}{2\pi}\right)^{M/2} \int \exp(-E[\mathbf{w}]) d\mathbf{w}$$

ただし、

$$\begin{aligned}
 E[\mathbf{w}] &= \beta E_D[\mathbf{w}] + \alpha E_W[\mathbf{w}] \\
 &= \frac{\beta}{2} {}^T (\mathbf{t} - \Phi \mathbf{w}) (\mathbf{t} - \Phi \mathbf{w}) + \frac{\alpha}{2} {}^T \mathbf{w} \mathbf{w}
 \end{aligned}$$

$E[\mathbf{w}]$ を \mathbf{w} について平方完成して積分を解くことを考える。

$$\begin{aligned}
 E[\mathbf{w}] &= \frac{\beta}{2} {}^T (\mathbf{t} - \Phi \mathbf{w}) (\mathbf{t} - \Phi \mathbf{w}) + \frac{\alpha}{2} {}^T \mathbf{w} \mathbf{w} \\
 &= \frac{\beta}{2} ({}^T \mathbf{t} \mathbf{t} - {}^T \mathbf{t} \Phi \mathbf{w} - {}^T \mathbf{w} {}^T \Phi \mathbf{t} + {}^T \mathbf{w} {}^T \Phi \Phi \mathbf{w}) + \frac{\alpha}{2} {}^T \mathbf{w} \mathbf{w} \\
 &= \frac{1}{2} {}^T \mathbf{w} (\beta {}^T \Phi \Phi + \alpha I) \mathbf{w} - \frac{\beta}{2} {}^T \mathbf{w} {}^T \Phi \mathbf{t} - \frac{\beta}{2} {}^T \mathbf{t} \Phi \mathbf{w} + \frac{\beta}{2} {}^T \mathbf{t} \mathbf{t} \\
 &= \frac{1}{2} {}^T (\mathbf{w} - A^{-1} B) A (\mathbf{w} - A^{-1} B) + \frac{\beta}{2} {}^T \mathbf{t} \mathbf{t} - \frac{1}{2} {}^B A^{-1} B
 \end{aligned}$$

ただし、

$$A = \alpha I + \beta {}^T \Phi \Phi, \quad B = \beta {}^T \mathbf{w} {}^T \Phi \mathbf{t}$$

とおいた。

$$\begin{aligned}
 \mathbf{m}_N &= A^{-1} B \\
 &= \beta A^{-1} {}^T \Phi \mathbf{t}
 \end{aligned}$$

と置くと、 $B = A\mathbf{m}_N$ なので、

$$\begin{aligned}
 \frac{\beta}{2} {}^T \mathbf{t} \mathbf{t} - {}^T B A^{-1} B &= \frac{\beta}{2} {}^T \mathbf{t} \mathbf{t} - \frac{1}{2} {}^T \mathbf{m}_N {}^T A \mathbf{m}_N \\
 &= \frac{\beta}{2} {}^T \mathbf{t} \mathbf{t} - \frac{1}{2} {}^T \mathbf{m}_N (\alpha I + \beta {}^T \Phi \Phi) \mathbf{m}_N \\
 &= \frac{\beta}{2} {}^T (\mathbf{t} - \Phi \mathbf{m}_N) (\mathbf{t} - \Phi \mathbf{m}_N) + \frac{\alpha}{2} {}^T \mathbf{m}_N \mathbf{m}_N \\
 &= E[\mathbf{m}_N]
 \end{aligned}$$

より、

$$E[\mathbf{w}] = \frac{1}{2} {}^T (\mathbf{w} - \mathbf{m}_N) A (\mathbf{w} - \mathbf{m}_N) + E[\mathbf{m}_N]$$

したがって、

$$\begin{aligned}
 \int \exp(-E[\mathbf{w}]) d\mathbf{w} &= \exp(-E[\mathbf{m}_N]) \int \exp\left\{-\frac{1}{2} {}^T (\mathbf{w} - \mathbf{m}_N) A (\mathbf{w} - \mathbf{m}_N)\right\} d\mathbf{w} \\
 &= \exp\{-E[\mathbf{m}_N]\} \frac{(2\pi)^{M/2}}{|A|^{1/2}}
 \end{aligned}$$

となるので、周辺尤度は

$$\begin{aligned}
 p(\mathbf{t}|\alpha, \beta) &= \left(\frac{\beta}{2\pi}\right)^{N/2} \left(\frac{\alpha}{2\pi}\right)^{M/2} \exp\{-E[\mathbf{m}_N]\} \frac{(2\pi)^{M/2}}{|A|^{1/2}} \\
 \log p(\mathbf{t}|\alpha, \beta) &= \frac{M}{2} \log \alpha + \frac{N}{2} \log \beta - E[\mathbf{m}_N] - \frac{1}{2} \log |A| - \frac{N}{2} \log(2\pi)
 \end{aligned}$$