

PRML 輪読会

線形回帰モデル (3.1.1 節 – 3.1.2 節)

鈴木拓己

November 20, 2024

線形基底関数モデル

最尤推定と最小二乗法

最小二乗法の幾何学

まとめ

演習問題 (Appendix)

線形基底関数モデル

最尤推定と最小二乗法

最小二乗法の幾何学

まとめ

演習問題 (Appendix)

回帰分析の問題設定および線形回帰モデルの概要

訓練データ: $(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)$, $\mathbf{x}_i \in \mathbb{R}^D$, $t_i \in \mathbb{R}$.

$$t = \underbrace{y(\mathbf{x})}_{\text{function}} + \underbrace{\varepsilon}_{\text{r.v.: error}} \longrightarrow \text{関数 } y \text{ を推定}$$

- 線形回帰モデル: 入力変数 (に対する非線形変換) の線形結合により t を予測
 - 線形結合のパラメータ \mathbf{w} に関しては線形 \longrightarrow 解析が容易
 - 入力変数に関しては非線形 \longrightarrow モデルの表現力が高い
- 回帰モデルの構成法
 - 関数 $y: \mathbb{R}^D \ni \mathbf{x} \longmapsto t \in \mathbb{R}$ を直接構成
 - * 最も単純なアプローチ
 - 予測分布 $p(t \mid \mathbf{x})$ をモデル化
 - * 実変数に対しては、損失関数として二乗損失関数を選ぶのが一般的であり、このとき最適解は $\mathbb{E}[t \mid \mathbf{x}]$ となる (1.5.5 節)

線形回帰モデルは解析的に扱いやすく、より洗練されたモデルの基礎として重要

最も単純な線形回帰モデル

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1x_1 + \cdots + w_Dx_D, \quad \text{where } \mathbf{x} = (x_1, \dots, x_D)^\top \in \mathbb{R}^D.$$

- 単に線形回帰 (linear regression) と呼ばれる
- パラメータ w_0, \dots, w_D に関して線形
- 入力変数 x_i に関しても線形であるため、表現力に乏しい

→ 入力変数に関して非線形な関数の線形結合を考え、モデルのクラスを拡張

線形基底関数¹モデル

$$y(\mathbf{x}, \mathbf{w}) = \underbrace{w_0}_{\text{bias 項}} + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}), \quad \text{where } \mathbf{x} = (x_1, \dots, x_D)^\top \in \mathbb{R}^D. \quad (1.1)$$

- M : モデルの表現力を制御するハイパーパラメータ
- ϕ_j : 基底関数 (basis function)
- 基底関数からなる特徴写像 (feature map) ϕ により, 入力 \mathbf{x} を特徴空間 (feature space) へ写像:

$$\phi: \mathbb{R}^D \ni \mathbf{x} \mapsto \phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^\top \in \mathbb{R}^{M-1}$$

- $\phi_0(\mathbf{x}) := 1$ とすれば, 式 (1.1) は $y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$.

¹<https://ja.wikipedia.org/wiki/%E5%9F%BA%E5%BA%95%E9%96%A2%E6%95%B0>

Example 1 (1 変数多項式回帰)

データ $(x, t) \in \mathbb{R} \times \mathbb{R}$ に対して, $\phi_j(x) := x^j$ とおくと,

$$\begin{aligned}\phi: \mathbb{R} \ni x &\longmapsto (1, x, x^2, \dots, x^{M-1})^\top \in \mathbb{R}^M, \\ y(x, \mathbf{w}) &= w_0 + w_1 x + w_2 x^2 + \dots + w_{M-1} x^{M-1}.\end{aligned}$$

Example 2 (多変数多項式回帰 (e.g. $D = 2$))

データ $(\mathbf{x}, t) = ((x_1, x_2), t) \in \mathbb{R}^2 \times \mathbb{R}$ に対して, $\phi_{k,\ell} := x_1^k x_2^\ell$ とおくと,

$$\begin{aligned}\phi: \mathbb{R}^2 \ni \mathbf{x} &\longmapsto (1, x_1, x_2, x_1 x_2, \dots, x_1^{M_1-1} x_2^{M_2-1})^\top \in \mathbb{R}^{M_1 M_2}, \\ y(\mathbf{x}, \mathbf{w}) &= w_{0,0} + w_{1,0} x_1 + w_{0,1} x_2 + w_{1,1} x_1 x_2 + \dots + w_{M_1-1, M_2-1} x_1^{M_1-1} x_2^{M_2-1}.\end{aligned}$$

その他の基底関数の例 (1 変数の場合)

- スプライン関数²
 - 多項式では入力空間のある領域の変化が他のすべての領域に及んでしまう
 - 入力空間を分割し、各領域で区分的に多項式をあてはめることにより解決できる
- ガウス型基底関数 (動径基底関数, RBF³ の例)

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\} \propto \mathcal{N}(\mu_j, s^2).$$

- シグモイド基底関数

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

logistic sigmoid function

²<https://ja.wikipedia.org/wiki/%E3%82%B9%E3%83%97%E3%83%A9%E3%82%A4%E3%83%B3%E6%9B%B2%E7%B7%9A>

³<https://ja.wikipedia.org/wiki/%E6%94%BE%E5%B0%84%E5%9F%BA%E5%BA%95%E9%96%A2%E6%95%B0>

その他の基底関数の例 (1 変数の場合)

(左) 多項式関数 (中央) ガウス型基底関数 (右) シグモイド基底関数

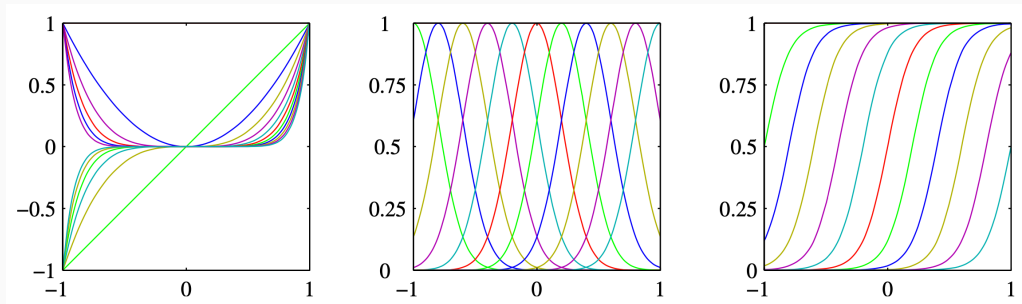


Figure 1: 基底関数の例

その他の基底関数の例 (1 変数の場合)

- tanh 関数

- 下記の式変形より、ロジスティックシグモイド関数の線形結合と等価 (→ 演習 3.1)

$$\begin{aligned}\tanh(a) &= \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)} = \frac{1 - \exp(-2a)}{1 + \exp(-2a)} \\ &= \frac{2}{1 + \exp(-2a)} - \frac{1 + \exp(-2a)}{1 + \exp(-2a)} = \underline{2\sigma(2a)} - 1.\end{aligned}\tag{1.2}$$

- フーリエ基底関数

- 様々な周波数の sin 関数, cos 関数による級数展開 (cf. フーリエ級数展開)

- ウェーブレット基底関数⁴

- フーリエ基底では周波数ごとに定まる基底関数は入力空間において無限に広がる
- ウェーブレット基底は入力空間でも周波数領域でも局所的
- 詳細は長くなりそうなので、定義と簡単な例のみ確認 (次ページ以降)

⁴<https://ja.wikipedia.org/wiki/%E3%82%A6%E3%82%A7%E3%83%BC%E3%83%96%E3%83%AC%E3%83%83%E3%83%88>

その他の基底関数の例 (1 変数の場合)

Definition 1 (ウェーブレット基底)

関数 $\psi \in L^2(\mathbb{R}) = \{f: \mathbb{R} \rightarrow \mathbb{C} \mid \int_{\mathbb{R}} |f(x)|^2 dx < \infty\}$ に対し, 関数 $\psi_{m,n}$ を

$$\psi_{m,n}(t) = 2^{m/2} \psi(2^m t - n), \quad m \in \mathbb{Z}, \quad n \in \mathbb{Z}$$

で定義する. 関数列 $\{\psi_{m,n}\}_{m \in \mathbb{Z}, n \in \mathbb{Z}}$ が $L^2(\mathbb{R})$ において基底⁵をなすとき, $\psi_{m,n}$ を ウェーブレット, ψ をその マザーウェーブレット, $\{\psi_{m,n}\}_{m \in \mathbb{Z}, n \in \mathbb{Z}}$ を ウェーブレット基底という. ウェーブレット基底 $\{\psi_{m,n}\}_{m \in \mathbb{Z}, n \in \mathbb{Z}}$ が正規直交条件

$$\langle \psi_{m,n}, \psi_{k,\ell} \rangle = \int_{\mathbb{R}} \psi_{m,n}(t) \overline{\psi_{k,\ell}(t)} dt = \delta_{mk} \delta_{n\ell}$$

を満たすとき, $\{\psi_{m,n}\}_{m \in \mathbb{Z}, n \in \mathbb{Z}}$ を 正規直交ウェーブレット基底という.

⁵すなわち, $\forall f \in L^2(\mathbb{R})$ に対して, $f = \sum_{m \in \mathbb{Z}} \sum_{n \in \mathbb{Z}} \langle f, \psi_{m,n} \rangle \psi_{m,n}$.

その他の基底関数の例 (1 変数の場合)

Example 3 (直交ウェーブレットの例: Haar ウェーブレット⁶)

マザーウェーブレット ψ を

$$\psi(t) = \begin{cases} 1, & 0 \leq t < 1/2, \\ -1, & 1/2 \leq t < 1, \\ 0, & \text{otherwise,} \end{cases}$$

で定めると, $\{\psi_{m,n}\}_{m \in \mathbb{Z}, n \in \mathbb{Z}}$ は直交ウェーブレットとなる.

Haar ウェーブレットのグラフは次ページ.

他にもいろいろな直交ウェーブレットがある (気になる人は pdf⁷ を参照してください).

⁶<https://ja.wikipedia.org/wiki/%E3%83%8F%E3%83%BC%E3%83%AB%E3%82%A6%E3%82%A7%E3%83%BC%E3%83%96%E3%83%AC%E3%83%83%E3%83%88>

⁷<http://wwwcs.ce.nihon-u.ac.jp/lab/moritaleb224w.pdf>

その他の基底関数の例 (1 変数の場合)

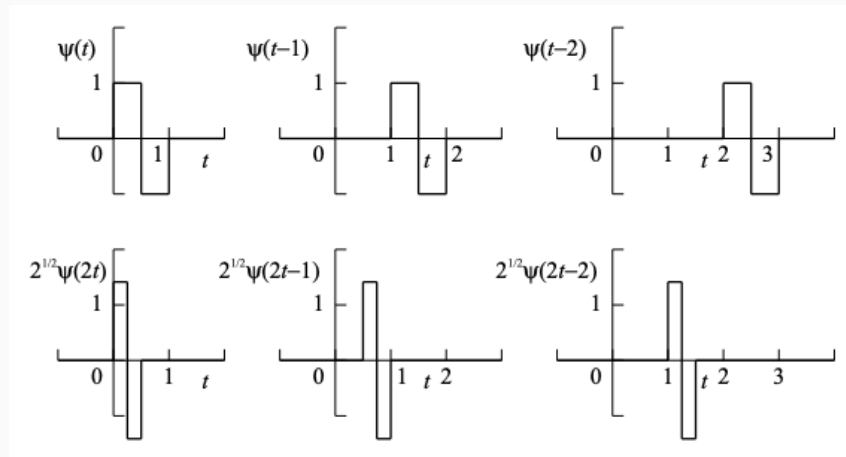


Figure 2: Haar ウェーブレットの場合の関数 $\psi_{0,0} = \psi, \psi_{0,1}, \psi_{0,2}, \psi_{1,0}, \psi_{1,1}, \psi_{1,2}$

線形基底関数モデル

最尤推定と最小二乗法

最小二乗法の幾何学

まとめ

演習問題 (Appendix)

以後の議論は基底関数の形に依らないため、基底関数は特に限定しないことにする。
また、簡単のため、目標変数 t が 1 次元の場合を扱う (多次元の場合は 3.1.5 節)。

ここでは、1 章の多項式フィッティングで扱った、最小二乗法と最尤推定との関係をより詳細に議論する。

1 章で見たように、多項式回帰において二乗和誤差関数を最小化するパラメータ w は、ガウスノイズモデルの下での最尤推定解と一致する：

$$t = \underbrace{y(\mathbf{x}, \mathbf{w})}_{\text{deterministic}} + \underbrace{\varepsilon}_{\text{stochastic}}, \quad \text{where } \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

すなわち、目標変数 t の予測分布は、精度パラメータ β を用いて、

$$p(t \mid \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t \mid y(\mathbf{x}, \mathbf{w}), \beta^{-1}), \quad \text{where } \beta^{-1} = \sigma^2$$

で与えられる。

1.5.5 節で見たように、二乗損失関数を仮定すれば、入力 \boldsymbol{x} の値に対する最適な予測値は目標変数の条件付き期待値

$$\mathbb{E}[t \mid \boldsymbol{x}] = \int_{\mathbb{R}} t p(t \mid \boldsymbol{x}) dt = \int_{\mathbb{R}} t \mathcal{N}(t \mid y(\boldsymbol{x}, \boldsymbol{w}), \beta^{-1}) dt = y(\boldsymbol{x}, \boldsymbol{w})$$

で与えられる。

Remark 1

ガウスノイズの仮定は、条件付き分布 $p(t \mid \boldsymbol{x})$ が単峰性のため、応用場面によっては不適切である可能性がある。

→ 14.5.1 節で条件付きガウス混合分布への拡張を議論 (多峰性の条件付き分布)

最尤推定により、予測分布 $p(t \mid \mathbf{x}, \mathbf{w}, \beta)$ のパラメータ \mathbf{w}, β を推定する.

入力 $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ と対応する目標値 t_1, \dots, t_N からなるデータ集合を考え、 $\mathbf{t} = (t_1, \dots, t_N)^\top$ とする. 各データ点が分布 $\mathcal{N}(t \mid y(\mathbf{x}, \mathbf{w}), \beta^{-1})$ から独立に生成されたと仮定すると、尤度関数は、

$$p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n \mid y(\mathbf{x}_n, \mathbf{w}), \beta^{-1}) = \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

で与えられる.

教師あり学習問題では、入力変数の分布をモデル化しないため、 \mathbf{x} は条件としてしか現れないので、以後は尤度関数の表記から \mathbf{x} を省略して $p(\mathbf{t} \mid \mathbf{w}, \beta)$ と書く.

このとき、対数尤度関数は、

$$\begin{aligned}\ln p(\mathbf{t} \mid \mathbf{w}, \beta) &= \ln \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) = \sum_{n=1}^N \ln \mathcal{N}(t_n \mid \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) \\&= \sum_{n=1}^N \ln \left(\frac{1}{(2\pi\beta^{-1})^{1/2}} \exp \left(-\frac{(t_n - \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n))^2}{2\beta^{-1}} \right) \right) \\&= \sum_{n=1}^N \left(\frac{1}{2} \ln \beta - \frac{1}{2} \ln(2\pi) - \beta \frac{1}{2} (t_n - \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n))^2 \right) \\&= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \underbrace{\beta \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n))^2}_{E_D(\mathbf{w}): \text{sum-of-squares error}}\end{aligned}\tag{2.1}$$

となる．最尤推定解を得るには、これを \mathbf{w} と β について最大化すればよい．

式 (2.1) より，パラメータ \boldsymbol{w} に関する対数尤度関数 $\ln p(\mathbf{t} \mid \boldsymbol{w}, \beta)$ の最大化は，二乗和誤差関数 $E_D(\boldsymbol{w})$ の最小化と等価である．

- ガウスノイズモデル下の最尤推定と最小二乗法との等価性

ここで，以下で与えられる計画行列 (design matrix)

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix} \in \mathbb{R}^{N \times M}$$

を考えると，二乗和誤差関数 $E_D(\boldsymbol{w})$ は

$$E_D(\boldsymbol{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \boldsymbol{w}^\top \boldsymbol{\phi}(\mathbf{x}_n))^2 = \frac{1}{2} (\mathbf{t} - \Phi \boldsymbol{w})^\top (\mathbf{t} - \Phi \boldsymbol{w}).$$

パラメータ w に関する対数尤度関数 $\ln p(\mathbf{t} \mid w, \beta)$ の最大化

$$\begin{aligned} E_D(w) &= \frac{1}{2}(\mathbf{t} - \Phi w)^\top (\mathbf{t} - \Phi w) = \frac{1}{2} \left(\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \Phi w - w^\top \Phi^\top \mathbf{t} + w^\top \Phi^\top \Phi w \right) \\ &= \frac{1}{2} \left(w^\top \Phi^\top \Phi w - 2\mathbf{t}^\top \Phi w + \mathbf{t}^\top \mathbf{t} \right) \end{aligned}$$

であり,

(1) $\Phi^\top \Phi$ は対称 ($\because (\Phi^\top \Phi)^\top = \Phi^\top (\Phi^\top)^\top = \Phi^\top \Phi$),

(2) $\Phi^\top \Phi$ はグラム行列⁸ $\implies \Phi^\top \Phi \succeq 0$,

ゆえ, $E_D(w)$ の w に関する最小化は制約なし凸 2 次計画問題⁹ なので,

$$w_{\text{ML}} = \operatorname{argmin}_{w \in \mathbb{R}^M} E_D(w) \iff \left. \frac{\partial}{\partial w} E_D(w) \right|_{w=w_{\text{ML}}} = 0.$$

⁸<https://mathlandscape.com/gram-matrix/>

⁹http://www.me.titech.ac.jp/~mizu_lab/lib/pdf/kougisiryousuuti/handout/11/suuti_kougi11-6.pdf

パラメータ \boldsymbol{w} に関する対数尤度関数 $\ln p(\mathbf{t} \mid \boldsymbol{w}, \beta)$ の最大化

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{w}} \ln p(\mathbf{t} \mid \boldsymbol{w}, \beta) &= \frac{\partial}{\partial \boldsymbol{w}} (-\beta E_D(\boldsymbol{w})) = -\frac{\beta}{2} \frac{\partial}{\partial \boldsymbol{w}} (\mathbf{t} - \Phi \boldsymbol{w})^\top (\mathbf{t} - \Phi \boldsymbol{w}) \\ &= -\frac{\beta}{2} \frac{\partial}{\partial \boldsymbol{w}} \left(\boldsymbol{w}^\top \Phi^\top \Phi \boldsymbol{w} - 2\mathbf{t}^\top \Phi \boldsymbol{w} + \mathbf{t}^\top \mathbf{t} \right) \\ &= -\frac{\beta}{2} \left(2\Phi^\top \Phi \boldsymbol{w} - 2\Phi^\top \mathbf{t} \right) = -\beta \left(\Phi^\top \Phi \boldsymbol{w} - \Phi^\top \mathbf{t} \right) \\ &\quad \left(\because \frac{\partial}{\partial \boldsymbol{w}} \boldsymbol{w}^\top \mathbf{A} \boldsymbol{w} = (\mathbf{A} + \mathbf{A}^\top) \boldsymbol{w}, \quad \frac{\partial}{\partial \boldsymbol{w}} \boldsymbol{a}^\top \boldsymbol{w} = \boldsymbol{a} \right)\end{aligned}$$

ゆえ、 $\Phi^\top \Phi$ が正則ならば、最尤推定解 $\boldsymbol{w}_{\text{ML}}$ は、

$$\left. \frac{\partial}{\partial \boldsymbol{w}} \ln p(\mathbf{t} \mid \boldsymbol{w}, \beta) \right|_{\boldsymbol{w}=\boldsymbol{w}_{\text{ML}}} = 0 \iff \Phi^\top \Phi \boldsymbol{w}_{\text{ML}} - \Phi^\top \mathbf{t} = 0 \quad (2.2)$$

$$\iff \boldsymbol{w}_{\text{ML}} = \left(\Phi^\top \Phi \right)^{-1} \Phi^\top \mathbf{t}. \quad (2.3)$$

パラメータ w に関する対数尤度関数 $\ln p(\mathbf{t} | \mathbf{w}, \beta)$ の最大化

- 方程式 (2.3) を最小二乗問題の正規方程式 (normal equation) という.
- 行列 $\Phi^\dagger \equiv (\Phi^\top \Phi)^{-1} \Phi^\top \in \mathbb{R}^{M \times N}$ を行列 Φ のムーア-ペンローズの擬似逆行列¹⁰ (Moore-Penrose pseudo-inverse matrix) という.
 - 通常の逆行列の概念の非正方行列への拡張
 - 実際, Φ が正則ならば, Φ^{-1} が存在し,

$$\Phi^\dagger \equiv (\Phi^\top \Phi)^{-1} \Phi^\top = \Phi^{-1} (\Phi^\top)^{-1} \Phi^\top = \Phi^{-1}.$$

- 任意の行列 M に対して, 定義の 4 条件を満たす擬似逆行列 M^\dagger が一意に存在
- 式 (2.2) において $\Phi^\top \Phi$ が正則でない場合については後述.
 - 計画行列 Φ がフルランクでない場合に発生する (多重共線性¹¹ (Multicollinearity)).

¹⁰<https://ja.wikipedia.org/wiki/%E3%83%A0%E3%83%BC%E3%82%A2%E3%83%BB%E3%83%9A%E3%83%B3%E3%83%AD%E3%83%BC%E3%82%BA%E9%80%86%E8%A1%8C%E5%88%97>

¹¹<https://ja.wikipedia.org/wiki/%E5%A4%9A%E9%87%8D%E5%85%B1%E7%B7%9A%E6%80%A7>

パラメータ β に関する対数尤度関数 $\ln p(\mathbf{t} \mid \mathbf{w}_{\text{ML}}, \beta)$ の最大化

次に、パラメータ β に関して対数尤度関数 $\ln p(\mathbf{t} \mid \mathbf{w}_{\text{ML}}, \beta)$ を最大化する．

$$\begin{aligned}\frac{\partial}{\partial \beta} \ln p(\mathbf{t} \mid \mathbf{w}_{\text{ML}}, \beta) &= \frac{\partial}{\partial \beta} \left(\frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}_{\text{ML}}^\top \phi(\mathbf{x}_n))^2 \right) \\ &= \frac{N}{2} \frac{1}{\beta} - \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}_{\text{ML}}^\top \phi(\mathbf{x}_n))^2\end{aligned}$$

ゆえ、最尤推定解 β_{ML} は、

$$\left. \frac{\partial}{\partial \beta} \ln p(\mathbf{t} \mid \mathbf{w}_{\text{ML}}, \beta) \right|_{\beta=\beta_{\text{ML}}} = 0 \iff \frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \underbrace{(t_n - \mathbf{w}_{\text{ML}}^\top \phi(\mathbf{x}_n))^2}_{\text{residual}}$$

として得られ、 $\beta_{\text{ML}}^{-1} = \sigma_{\text{ML}}^2$ は回帰関数 $y(\mathbf{x}, \mathbf{w}_{\text{ML}})$ 周りでの目標値の残差分散で与えられることがわかる．

バイアスパラメータ w_0 の役割

二乗和誤差関数 $E_D(\mathbf{w})$ において、バイアスパラメータ w_0 を明示的に書くと、

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \phi(\mathbf{x}_n))^2 = \frac{1}{2} \sum_{n=1}^N \left(t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n) \right)^2$$
$$\implies \frac{\partial}{\partial w_0} E_D(\mathbf{w}) = - \sum_{n=1}^N \left(t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n) \right)$$

ゆえ、 $E_D(\mathbf{w})$ を最小化するパラメータ \hat{w}_0 は、

$$\left. \frac{\partial}{\partial w_0} E_D(\mathbf{w}) \right|_{w_0=\hat{w}_0} = 0 \iff \hat{w}_0 = \frac{1}{N} \sum_{n=1}^N t_n - \sum_{j=1}^{M-1} w_j \left(\frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n) \right) \equiv \bar{t} - \sum_{j=1}^{M-1} w_j \overline{\phi_j}.$$

バイアスパラメータ w_0 は、訓練データ $\{(\mathbf{x}_n, t_n)\}_{n=1}^N$ に関する、目標値の平均 \bar{t} と、基底関数の値の平均の重み付き和 $\sum_{j=1}^{M-1} w_j \overline{\phi_j}$ との差を埋め合わせる役割をしている。

線形基底関数モデル

最尤推定と最小二乗法

最小二乗法の幾何学

まとめ

演習問題 (Appendix)

計画行列 Φ の第 j 列を $\varphi_j = (\phi_0(\mathbf{x}_1), \dots, \phi_{j-1}(\mathbf{x}_N))^T \in \mathbb{R}^N$ とおくと,

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix} = \begin{pmatrix} \varphi_1 & \varphi_2 & \cdots & \varphi_M \end{pmatrix} \in \mathbb{R}^{N \times M}.$$

このとき, 基底関数の数 M がデータ点数 N よりも小さければ¹², M 個のベクトル $\{\varphi_j\}_{j=1}^M$ は M 次元線形部分空間 $\mathcal{S} \subsetneq \mathbb{R}^N$ を張る:

$$\mathcal{S} = \text{Im } \Phi = \text{span}_{\mathbb{R}}(\varphi_1, \dots, \varphi_M) = \left\{ \sum_{j=1}^M w_j \varphi_j \mid \mathbf{w} \in \mathbb{R}^M \right\} \subsetneq \mathbb{R}^N.$$

¹² $M \geq N$ かつ Φ がフルランクならば $\mathcal{S} = \mathbb{R}^N$.

最小二乗解の幾何学的解釈

ベクトル $\mathbf{y} \in \mathbb{R}^N$ を

$$\mathbf{y} := (y(\mathbf{x}_1, \mathbf{w}), \dots, y(\mathbf{x}_N, \mathbf{w}))^\top$$

で定義すると, $y(\mathbf{x}_n, \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x}_n)$ より,

$$\mathbf{y} = (\mathbf{w}^\top \phi(\mathbf{x}_1), \dots, \mathbf{w}^\top \phi(\mathbf{x}_N))^\top = \Phi \mathbf{w} \in \mathcal{S}.$$

ここで, 最小二乗解 \mathbf{w}_{ML} は, $\|\cdot\|_2$ をユークリッドノルムとして,

$$\mathbf{w}_{\text{ML}} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^M} E_D(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^M} \frac{1}{2} \|\mathbf{t} - \underbrace{\Phi \mathbf{w}}_{\mathbf{y}}\|_2^2$$

で与えられるので, 最小二乗解 $\mathbf{w}_{\text{ML}} \in \mathbb{R}^M$ を求めることは, ユークリッド距離で $\mathbf{t} \in \mathbb{R}^N$ に最も近い $\Phi \mathbf{w} = \mathbf{y} \in \mathcal{S}$ を選ぶことに相当する. この $\mathbf{y} \in \mathcal{S}$ は, $\mathbf{t} \in \mathbb{R}^N$ の部分空間 \mathcal{S} 上への正射影に対応する (次ページ参照).

最小二乗解の幾何学的解釈

直観的な説明

- $N = 3, M = 2$ の場合は下図

$$\begin{array}{ll}\text{minimize} & \|\mathbf{t} - \mathbf{y}\|_2^2 \\ \text{subject to} & \mathbf{y} \in \mathcal{S}\end{array}$$

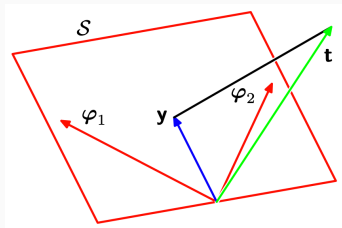


Figure 3: 最小二乗解の幾何学的解釈

理論的な説明

行列 $\Phi (\Phi^\top \Phi)^{-1} \Phi^\top$ は任意のベクトル $\mathbf{v} \in \mathbb{R}^N$ を $\mathcal{S} = \text{span}_{\mathbb{R}}(\varphi_1, \dots, \varphi_M)$ に正射影する (\because 演習 3.2). したがって,

$$\mathbf{w}_{\text{ML}} = \left(\Phi^\top \Phi \right)^{-1} \Phi^\top \mathbf{t}$$

より,

$$\hat{\mathbf{y}} = \Phi \mathbf{w}_{\text{ML}} = \Phi \left(\Phi^\top \Phi \right)^{-1} \Phi^\top \mathbf{t} \quad (3.1)$$

は $\mathbf{t} \in \mathbb{R}^N$ の部分空間 \mathcal{S} 上への正射影となっている.

$\Phi^\top \Phi$ が正則でない (非正則に近い) 場合の対処法

式 (2.2) において $\Phi^\top \Phi$ が正則でない場合には、最尤推定解 w_{ML} を求められない。特に、2 つ以上の基底ベクトル φ_j が線形従属の場合、 $\text{rank } \Phi < \min(M, N)$ となり、このような問題が発生する。また、2 つ以上の基底ベクトル φ_j が線形従属に近い場合にも同様の問題が発生し、求まるパラメータ w_{ML} が数値的に不安定になる。この問題の対処法として以下のような手法がある。

- 特異値分解¹³ (singular value decomposition; SVD)
 - 陽に逆行列 $(\Phi^\top \Phi)^{-1}$ を計算せずとも、擬似逆行列 Φ^\dagger を計算することができる。
- 正則化最小二乗法 (→ 3.1.4 節)
 - 正則化項を加えることにより、 $\Phi^\top \Phi$ が正則でない場合でも最適解 w_{ML} が得られる。
 - 特に、 $q = 2$ (Ridge 回帰) の場合は解析的に最適解が求められる。

$$\text{minimize } \|\mathbf{t} - \Phi \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_q^q \quad \text{subject to } \mathbf{w} \in \mathbb{R}^M, \quad \text{where } \|\cdot\|_q \text{ is } L^q\text{-norm}^{14}.$$

¹³<https://ja.wikipedia.org/wiki/%E7%89%B9%E7%95%B0%E5%80%A4%E5%88%86%E8%A7%A3>

¹⁴<https://ja.wikipedia.org/wiki/Lp%E7%A9%BA%E9%96%93>

Theorem 1 (特異値分解定理¹⁵)

$M \in \mathbb{C}^{m \times n}$ を $\text{rank } M = r$ の行列とする. このとき, $\Sigma \in \mathbb{C}^{m \times n}$ が一意に存在し,

$$M = U \Sigma V^* \quad (3.2)$$

の分解が成り立つ. ただし, U, V はそれぞれ $m \times m, n \times n$ のユニタリ行列であり, Σ は半正定値行列 MM^* あるいは M^*M の正の固有値¹⁶ の平方根 $\sigma_1 \geq \dots \geq \sigma_r > 0$ および $q = \min(m, n)$, $\sigma_{r+1} = \dots = \sigma_q = 0$ を用いて

$$\Sigma = \begin{cases} \begin{pmatrix} \Delta & \mathbf{0} \end{pmatrix}, & (m < n), \\ \Delta, & (m = n), \\ \begin{pmatrix} \Delta \\ \mathbf{0} \end{pmatrix}, & (m > n), \end{cases} \quad \text{where } \Delta = \text{diag}(\sigma_1, \dots, \sigma_q) = \begin{pmatrix} \sigma_1 & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \sigma_q \end{pmatrix}$$

で与えられる. 分解 (3.2) を M の特異値分解, $\sigma_1, \dots, \sigma_q$ を M の特異値と呼ぶ.

¹⁵証明: <https://qiita.com/gyu-don/items/2e4fd7e945d8c349afb5>

¹⁶ MM^* と M^*M の非零固有値は一致する: <http://teagis.ip.is.saga-u.ac.jp/svd.pdf>

特異値分解の計算と幾何学的解釈

ユニタリ行列による変換は等長変換なので、特異値分解 $M = U\Sigma V^*$ におけるそれぞれの線形変換は、

- V^* : 回転 (\because ユニタリ行列)
- Σ : 拡大・縮小
- U : 回転 (\because ユニタリ行列)

を表す．また、 $M = U\Sigma V^*$ より

$$MM^* = U\Sigma\Sigma^T U^*, \quad M^*M = V\Sigma^T\Sigma V^*,$$

ゆえ、特異値分解における U, V はそれぞれ MM^*, M^*M の固有値分解によって得られる．

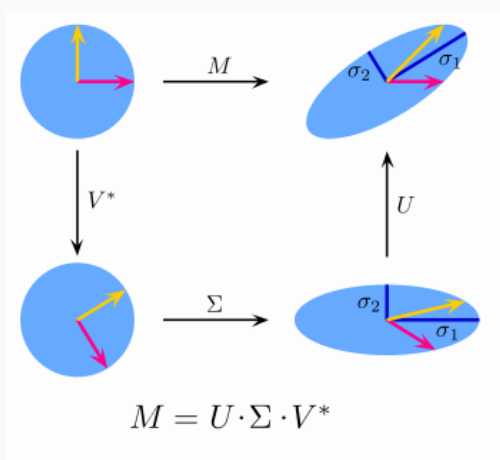


Figure 4: 特異値分解の図示

特異値分解による擬似逆行列 Φ^\dagger の計算

特異値分解定理により、計画行列 $\Phi \in \mathbb{R}^{N \times M}$ に対して、直交行列 $\mathbf{U} \in \mathbb{R}^{N \times N}$, $\mathbf{V} \in \mathbb{R}^{M \times M}$ および行列 $\Sigma \in \mathbb{R}^{N \times M}$ が存在し、 $\Phi = \mathbf{U}\Sigma\mathbf{V}^\top$ の分解が成り立つ。いま、 \mathbf{U}, \mathbf{V} は直交行列ゆえ、

$$\mathbf{U}\mathbf{U}^\top = \mathbf{U}^\top\mathbf{U} = \mathbf{I}, \quad \mathbf{V}\mathbf{V}^\top = \mathbf{V}^\top\mathbf{V} = \mathbf{I}, \quad \mathbf{U}^{-1} = \mathbf{U}^\top, \quad \mathbf{V}^{-1} = \mathbf{V}^\top$$

が成り立つことに注意すると、

$$\begin{aligned} \Phi^\dagger &= (\Phi^\top \Phi)^{-1} \Phi^\top = (\mathbf{V}\Sigma^\top \mathbf{U}^\top \mathbf{U}\Sigma\mathbf{V}^\top)^{-1} \mathbf{V}\Sigma^\top \mathbf{U}^\top \\ &= \mathbf{V} (\Sigma^\top \Sigma)^{-1} \mathbf{V}^\top \mathbf{V}\Sigma^\top \mathbf{U}^\top = \mathbf{V} (\Sigma^\top \Sigma)^{-1} \Sigma^\top \mathbf{U}^\top \\ &\equiv \mathbf{V}\Sigma^\dagger \mathbf{U}^\top. \quad (\Sigma^\dagger \in \mathbb{R}^{M \times N}) \end{aligned}$$

ただし、 $\Sigma^\dagger \equiv (\Sigma^\top \Sigma)^{-1} \Sigma^\top$ は、 Σ の非零成分¹⁷の逆数を成分とする行列の転置。

¹⁷数値計算では、許容誤差よりも大きい要素のみが非零と見なされ、他の要素は0に置き換えられる。

線形基底関数モデル

最尤推定と最小二乗法

最小二乗法の幾何学

まとめ

演習問題 (Appendix)

- 線形基底関数モデル
 - 基底関数 ϕ_j からなる特徴写像 ϕ によって線形回帰モデルを拡張
 - 様々な基底関数が使われる
- 最尤推定と最小二乗法
 - 最小二乗解とガウスノイズモデル下の最尤推定解は一致
 - 計画行列 Φ がフルランクならば，最小二乗解 w_{ML} は解析的に求められる
- 最小二乗法の幾何学
 - 最小二乗解 w_{ML} の求解は，訓練データの目標値からなるベクトル $\mathbf{t} \in \mathbb{R}^N$ のモデル空間 $\mathcal{S} = \text{Im } \Phi = \text{span}_{\mathbb{R}}(\varphi_1, \dots, \varphi_M)$ 上への正射影を求めることに対応する
 - $\Phi^{\top} \Phi$ が非正則に近い場合の対処法
 - * 特異値分解
 - * 正則化最小二乗法
 - * 正則化の一種として，主成分回帰¹⁸ (principal component regression; PCR) などもある

¹⁸https://en.wikipedia.org/wiki/Principal_component_regression

線形基底関数モデル

最尤推定と最小二乗法

最小二乗法の幾何学

まとめ

演習問題 (Appendix)

演習 3.1 (基本) www

\tanh 関数とロジスティックシグモイド関数は次のように関係付けられることを示せ.

$$\tanh(a) = 2\sigma(2a) - 1. \quad (5.1)$$

さらに, 次の形のロジスティックシグモイド関数の線形結合

$$y(x, \mathbf{w}) = w_0 + \sum_{j=1}^M w_j \sigma\left(\frac{x - \mu_j}{s}\right)$$

は次の形の \tanh 関数の線形結合

$$y(x, \mathbf{u}) = u_0 + \sum_{j=1}^M u_j \tanh\left(\frac{x - \mu_j}{2s}\right)$$

と等価であることを示し, 新しいパラメータ $\{u_0, \dots, u_M\}$ ともとのパラメータ $\{w_0, \dots, w_M\}$ を関係付ける式を求めよ.

(解答) 式 (1.2) より, 式 (5.1) が成り立つ. これを用いると,

$$y(x, \mathbf{u}) = u_0 + \sum_{j=1}^M u_j \tanh\left(\frac{x - \mu_j}{2s}\right) = u_0 + \sum_{j=1}^M u_j \left\{ 2\sigma\left(\frac{x - \mu_j}{s}\right) - 1 \right\} = u_0 - \sum_{j=1}^M u_j + \sum_{j=1}^M 2u_j \sigma\left(\frac{x - \mu_j}{s}\right)$$

ゆえ, パラメータの関係式は $w_0 = u_0 - \sum_{j=1}^M u_j$, $w_j = 2u_j$ ($j = 1, \dots, M$) で与えられる.

演習 3.2 (標準)

行列

$$\Phi \left(\Phi^\top \Phi \right)^{-1} \Phi^\top$$

は任意のベクトル \boldsymbol{v} を Φ の列ベクトルで張られる空間の上に正射影することを示せ．そしてこの結果を使って，最小二乗解 (2.3) は図 3 で示した多様体 S の上にベクトル \boldsymbol{t} を正射影することに対応していることを示せ．

(解答) 行列 $\Phi = (\varphi_1, \dots, \varphi_M) \in \mathbb{R}^{N \times M}$ に対して，部分空間 $S \in \mathbb{R}^N$ を $S = \text{span}_{\mathbb{R}}(\varphi_1, \dots, \varphi_M)$ とおく．任意のベクトル $\boldsymbol{v} \in \mathbb{R}^N$ の S 上への正射影 \boldsymbol{p} は， $\boldsymbol{w} = (w_1, \dots, w_M) \in \mathbb{R}^M$ を用いて， $\boldsymbol{p} = \sum_{j=1}^M w_j \varphi_j = \Phi \boldsymbol{w}$ と表せる．このとき，ベクトル $\boldsymbol{v} - \boldsymbol{p}$ と S の基底 $\varphi_1, \dots, \varphi_M$ は直交するので，

$$\begin{aligned} \langle \boldsymbol{v} - \boldsymbol{p}, \varphi_j \rangle &= 0, \quad \forall j \in \{1, \dots, M\} \iff \Phi^\top (\boldsymbol{v} - \boldsymbol{p}) = 0 \\ &\iff \Phi^\top (\boldsymbol{v} - \Phi \boldsymbol{w}) = 0 \\ &\iff \boldsymbol{w} = \left(\Phi^\top \Phi \right)^{-1} \Phi^\top \boldsymbol{v} \end{aligned}$$

が成り立つ．したがって， \boldsymbol{v} の正射影 \boldsymbol{p} は，

$$\boldsymbol{p} = \Phi \boldsymbol{w} = \Phi \left(\Phi^\top \Phi \right)^{-1} \Phi^\top \boldsymbol{v}$$

で表されるので，行列 $\Phi \left(\Phi^\top \Phi \right)^{-1} \Phi^\top$ は任意のベクトル \boldsymbol{v} を Φ の列ベクトルで張られる空間の上に正射影する．最小二乗解 (2.3) が多様体 S の上にベクトル \boldsymbol{t} を正射影することに対応していることは，式 (3.1) より成り立つ．

□

演習 3.3 (基本)

それぞれのデータ点 t_n に重み要素 $r_n > 0$ が割り当てられており、二乗和誤差関数が

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N r_n \left\{ t_n - \mathbf{w}^\top \phi(\mathbf{x}_n) \right\}^2$$

となるデータ集合を考える．このとき、この誤差関数を最小にする解 \mathbf{w}^* についての式を求めよ．また、(i) ノイズの分散がデータに依存する場合、(ii) データ点に重複がある場合に照らして、それぞれ重み付き二乗和誤差関数の解釈を与えよ．

(解答) 行列 \mathbf{R} を $\mathbf{R} = \text{diag}(r_1, \dots, r_N) \in \mathbb{R}^{N \times N}$ とすると、

$$\begin{aligned} E_D(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N r_n \left\{ t_n - \mathbf{w}^\top \phi(\mathbf{x}_n) \right\}^2 = \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^\top \mathbf{R} (\mathbf{t} - \Phi \mathbf{w}) \\ &= \frac{1}{2} (\mathbf{t}^\top - \mathbf{w}^\top \Phi^\top) \mathbf{R} (\mathbf{t} - \Phi \mathbf{w}) = \frac{1}{2} (\mathbf{w}^\top \Phi^\top \mathbf{R} \Phi \mathbf{w} - 2 \mathbf{t}^\top \mathbf{R} \Phi \mathbf{w} + \mathbf{t}^\top \mathbf{R} \mathbf{t}) \end{aligned}$$

ゆえ、

$$\frac{\partial}{\partial \mathbf{w}} E_D(\mathbf{w}) = \Phi^\top \mathbf{R} \Phi \mathbf{w} - \Phi^\top \mathbf{R} \mathbf{t} \implies \left. \frac{\partial}{\partial \mathbf{w}} E_D(\mathbf{w}) \right|_{\mathbf{w}=\mathbf{w}^*} = 0 \iff \mathbf{w}^* = \left(\Phi^\top \mathbf{R} \Phi \right)^{-1} \Phi^\top \mathbf{R} \mathbf{t}.$$

演習 3.3 (基本)

それぞれのデータ点 t_n に重み要素 $r_n > 0$ が割り当てられており、二乗和誤差関数が

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N r_n \left\{ t_n - \mathbf{w}^\top \phi(\mathbf{x}_n) \right\}^2$$

となるデータ集合を考える．このとき、この誤差関数を最小にする解 \mathbf{w}^* についての式を求めよ．また、(i) ノイズの分散がデータに依存する場合、(ii) データ点に重複がある場合に照らして、それぞれ重み付き二乗和誤差関数の解釈を与えよ．

(解答続き) (i) ノイズの分散がデータに依存する場合、式 (2.1) より、重み付き二乗和誤差関数はデータ点 t_n のノイズの分散が r_n^{-1} で与えられている場合の二乗和誤差関数と解釈できる．(ii) データ点に重複がある場合、重み付き二乗和誤差関数はデータ点 (\mathbf{x}_n, t_n) が r_n 個重複している場合の二乗和誤差関数と解釈できる．