

カーネル法

中田和秀

東京科学大学 工学院 経営工学系

機械学習入門

<https://www.nakatalab.iee.e.titech.ac.jp/text/nakata.html>

概要

ここではカーネルを使うことによって、単純な計算によって複雑な分析を行う一般的な方法論について説明を行う。その後で、カーネル回帰、カーネル Ridge 回帰、ガウス過程回帰について説明する。なお、サポートベクターマシンやカーネル主成分分析は別スライドで説明を行っている。

目次：

1. カーネルトリック
2. カーネル回帰
2. カーネル Ridge 回帰
3. ガウス過程回帰

記号の使い方：

- $A := B$ は、 B で A を定義する、 B を A に代入することを意味する
- $[n]$ は n までのインデックスの集合を表し $[n] := \{1, 2, \dots, n\}$

高次元空間への写像

- 入力空間を高次元の特徴空間に写像して、その特徴空間上で単純な分析
- 写像 $\phi: \mathbb{R}^n \rightarrow V$ を使って、 x_d の代わりに $\phi(x_d)$ を使う
- 元の空間 \mathbb{R}^n でみると複雑な分析となる

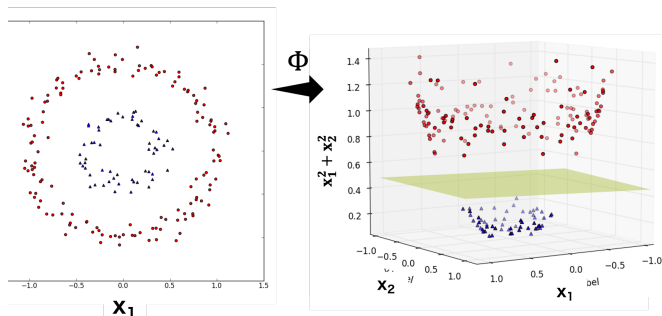


Figure: <https://axa.biopapyrus.jp/machine-learning/svm/kernel-svm.html>

カーネルトリック

写像 $\phi: \mathbb{R}^n \rightarrow V$ をどう決めるかが難しい

もし

- (1) x_i と x_j から直接 $\phi(x_i)^T \phi(x_j)$ の値を計算できる
- (2) 高次元空間上での内積 $\phi(x_i)^T \phi(x_j)$ にのみ基づいて分析

\implies 写像 $\phi(x)$ を陽に定めることなく分析が可能

まず、(1) について議論する。

$\phi(x_a)^T \phi(x_b)$ に相当するカーネル関数 ($\mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$) を導入

$$\text{Ker}(x_a, x_b)$$

論点：

写像 ϕ の存在が担保されるように、カーネル関数を構築する必要がある

マーサーの条件

$K_{ij} := \text{Ker}(\mathbf{x}_i, \mathbf{x}_j) \quad (i, j \in [D])$ とする。

$$\mathbf{K} = \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T$$

$$\widetilde{\mathbf{X}} := \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_D)^T \end{pmatrix} \in \mathbb{R}^{D \times |V|}$$

必要条件： カーネル行列 $\mathbf{K} = \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T$ は半正定値行列

実は必要十分条件である

マーサーの条件

任意の $i, j \in [D]$ に対し、 $K_{ij} := \text{Ker}(\mathbf{x}_i, \mathbf{x}_j)$ で作られる カーネル行列 $\mathbf{K} \in \mathbb{R}^{D \times D}$ が半正定値となるとき、写像 ϕ の存在が保証される

カーネルの例

- $\text{Ker}(\mathbf{x}_a, \mathbf{x}_b) := \mathbf{x}_a^T \mathbf{x}_b$ 線形カーネル
- $\text{Ker}(\mathbf{x}_a, \mathbf{x}_b) := (\mathbf{x}_a^T \mathbf{x}_b + 1)^d \quad (d \in \mathbb{N})$ 多項式カーネル
- $\text{Ker}(\mathbf{x}_a, \mathbf{x}_b) := \exp\{-\gamma \|\mathbf{x}_a - \mathbf{x}_b\|^2\}$ RBF カーネル、ガウシアンカーネル

カーネルの証明

線形カーネル：

$K_{ij} := \mathbf{x}_i^T \mathbf{x}_j$ のとき、 $\mathbf{K} = \mathbf{X}\mathbf{X}^T$ より半正定値行列

多項式カーネル：

$K_{ij} := (\mathbf{x}_i^T \mathbf{x}_j + 1)^d$ から作られるカーネル行列を \mathbf{K}_d とする。

- $\mathbf{K}_1 = \mathbf{X}\mathbf{X}^T + \mathbf{e}\mathbf{e}^T$ より半正定値行列
- $\mathbf{K}_d, \mathbf{K}_1$ が半正定値のとき、 $\mathbf{K}_{d+1} = \mathbf{K}_d \circ \mathbf{K}_1$ より、 \mathbf{K}_{d+1} も半正定値

(証明は次スライド)

$\mathbf{C} = \mathbf{A} \circ \mathbf{B}$ は Hadamard (アダマール) 積を意味し $C_{ij} := A_{ij}B_{ij}$

RBF カーネル：

後述

Hadamard 積の半正定値性

$\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times n}$ が半正定値のとき、 $\mathbf{Z} := \mathbf{X} \circ \mathbf{Y} \in \mathbb{R}^{n \times n}$ も半正定値

$$\text{Hadamard 積} \quad Z_{ij} := X_{ij}Y_{ij} \quad (i, j \in [n])$$

証明：

$\mathbf{X} = \mathbf{X}^{1/2} \mathbf{X}^{1/2}$, $\mathbf{Y} = \mathbf{Y}^{1/2} \mathbf{Y}^{1/2}$ なので、 $\mathbf{X}^{1/2}, \mathbf{Y}^{1/2}$ の i 列目を $\mathbf{x}_i, \mathbf{y}_i$ とする。

$$\begin{aligned} \mathbf{v}^T (\mathbf{X} \circ \mathbf{Y}) \mathbf{v} &= \sum_{i \in [n]} \sum_{j \in [n]} v_i v_j (\mathbf{x}_i^T \mathbf{x}_j) (\mathbf{y}_i^T \mathbf{y}_j) = \sum_{i \in [n]} \sum_{j \in [n]} v_i v_j (\mathbf{x}_i^T \mathbf{x}_j) (\mathbf{y}_j^T \mathbf{y}_i) \\ &= \sum_{i \in [n]} \sum_{j \in [n]} v_i v_j \text{Tr}(\mathbf{x}_i^T \mathbf{x}_j \mathbf{y}_j^T \mathbf{y}_i) = \sum_{i \in [n]} \sum_{j \in [n]} v_i v_j \text{Tr}(\mathbf{y}_i \mathbf{x}_i^T \mathbf{x}_j \mathbf{y}_j^T) \\ &= \text{Tr} \left(\sum_{i \in [n]} \sum_{j \in [n]} v_i v_j \mathbf{y}_i \mathbf{x}_i^T \mathbf{x}_j \mathbf{y}_j^T \right) = \text{Tr} \left(\left(\sum_{i \in [n]} v_i \mathbf{y}_i \mathbf{x}_i^T \right) \left(\sum_{j \in [n]} v_j \mathbf{x}_j \mathbf{y}_j^T \right) \right) \\ &= \text{Tr} (\mathbf{M} \mathbf{M}^T) = \sum_{i \in [n]} \sum_{j \in [n]} (M_{ij})^2 \geq 0. \end{aligned}$$

写像 ϕ の例：多項式カーネル

$\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$ に対する 3 次多項式カーネル $\text{Ker}(\mathbf{x}, \mathbf{y}) := (\mathbf{x}^T \mathbf{y} + 1)^3$ を考える。

$$\begin{aligned}\text{Ker}(\mathbf{x}, \mathbf{y}) &:= (\mathbf{x}^T \mathbf{y} + 1)^3 \\ &= (x_1 y_1 + x_2 y_2 + 1)^3 \\ &= x_1^3 y_1^3 + 3x_1^2 y_1^2 x_2 y_2 + 3x_1 y_1 x_2^2 y_2^2 + x_2^3 y_2^3 \\ &\quad + 3x_1^2 y_1^2 + 6x_1 y_1 x_2 y_2 + 3x_2^2 y_2^2 + 3x_1 y_1 + 3x_2 y_2 + 1\end{aligned}$$

よって、 $\phi(\mathbf{x}) : \mathbb{R}^2 \rightarrow \mathbb{R}^{10}$ を

$$\phi(\mathbf{x}) := \left(x_1^3, \sqrt{3}x_1^2 x_2, \sqrt{3}x_1 x_2^2, x_2^3, \sqrt{3}x_1^2, \sqrt{6}x_1 x_2, \sqrt{3}x_2^2, \sqrt{3}x_1, \sqrt{3}x_2, 1 \right)^T$$

と定めると、次の関係が成り立つ。

$$\phi(\mathbf{x})^T \phi(\mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^3 = \text{Ker}(\mathbf{x}, \mathbf{y})$$

$\phi(\mathbf{x}) \in \mathbb{R}^{10}$ に対する線形な判別は、 $\mathbf{x} \in \mathbb{R}^2$ に対する非線形な判別（3 次多項式）

無限次元ベクトルと関数

(1) 有限次元のベクトル： $\boldsymbol{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix} \in \mathbb{R}^m$

これは次のように定義した f で表すことができる。

- $f(1) = x_1, f(2) = x_2, \dots, f(m) = x_m$

関数の定義域は $1, 2, \dots, m$

(2) 無限次元のベクトル： $\boldsymbol{x} \in V$

これは次のように定義した f で表すことができる（2例挙げる）。

- $f(s) = x_s$ 関数の定義域は $s \in \mathbb{R}$

- $f(s) = x_s$ 関数の定義域は $s \in \mathbb{R}^n$

無限次元のベクトルは関数で表すことができる。

無限次元ベクトルと関数

- 有限次元のベクトルの和とスカラー倍を自然に拡張した形で、関数の和とスカラー倍が定義できる。…… 関数空間は線形空間
- 内積も自然な形で拡張すると、 $\mathbf{x}_a^T \mathbf{x}_b \rightarrow \int f_a(s) f_b(s) ds, \int f_a(s) f_b(s) ds$

内積の説明：

M 次元の場合、 $\mathbf{x}_a^T \mathbf{x}_b = \sum_{k \in [M]} f_a(k) f_b(k)$ であるが、無限次元に拡張するとき、無限個の和とすると発散する可能性があるため、要素数で割ることにする。

$$\text{内積} = \frac{1}{M} \sum_{k \in [M]} f_a(k) f_b(k)$$

関数の定義域 \mathcal{S}_k を $-k$ から k までは $\frac{1}{k}$ の間隔で取ったものとする。

$$\mathcal{S}_k := \left\{ -\frac{k^2}{k}, -\frac{k^2-1}{k}, \dots, \frac{0}{k}, \dots, \frac{k^2-1}{k}, \frac{k^2}{k} \right\}$$

このとき、 $\lim_{k \rightarrow \infty} \frac{1}{|\mathcal{S}_k|} \sum_{k \in \mathcal{S}_k} f_a(k) f_b(k) = \int_{\mathbb{R}} f_a(s) f_b(s) ds$

写像 ϕ の例： RBF カーネル

$\phi: \mathbb{R}^n \rightarrow V$

ベクトル \mathbf{x} が入力されると、関数 $f(\mathbf{s})$ を出力

$$f(\mathbf{s}) := \beta \exp\{-2\gamma\|\mathbf{x} - \mathbf{s}\|^2\} \quad (\alpha, \beta \text{ は定数})$$

$$\begin{aligned} \phi(\mathbf{x}_a) \text{ と } \phi(\mathbf{x}_b) \text{ の内積: } & \int \beta \exp\{-2\gamma\|\mathbf{x}_a - \mathbf{s}\|^2\} \beta \exp\{-2\gamma\|\mathbf{x}_b - \mathbf{s}\|^2\} d\mathbf{s} \\ & = \beta^2 \left(\frac{\pi}{4\gamma}\right)^{n/2} \exp\{-\gamma\|\mathbf{x}_a - \mathbf{x}_b\|^2\} \end{aligned}$$

α, β を適切に設定すると、RBF カーネル

$$\text{Ker}(\mathbf{x}_a, \mathbf{x}_b) := \exp\{-\gamma\|\mathbf{x}_a - \mathbf{x}_b\|^2\}$$

に対応する高次元写像となる

※ なお、無限次元は扱いにくいので、カーネル回帰以降は有限次元のイメージで話を進める

式変形

$$\begin{aligned}\phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) &:= \int \beta \exp\{-2\gamma\|\mathbf{s} - \mathbf{x}_1\|^2\} \beta \exp\{-2\gamma\|\mathbf{s} - \mathbf{x}_2\|^2\} d\mathbf{s} \\&= \beta^2 \int \exp\{-2\gamma(\|\mathbf{s} - \mathbf{x}_1\|^2 + \|\mathbf{s} - \mathbf{x}_2\|^2)\} d\mathbf{s} \\&= \beta^2 \int \exp\left\{-\gamma\left(4\left\|\mathbf{s} - \frac{\mathbf{x}_1 + \mathbf{x}_2}{2}\right\|^2 + \|\mathbf{x}_1 - \mathbf{x}_2\|^2\right)\right\} d\mathbf{s} \\&= \beta^2 \exp\{-\gamma\|\mathbf{x}_1 - \mathbf{x}_2\|^2\} \int \exp\left\{-4\gamma\left\|\mathbf{s} - \frac{\mathbf{x}_1 + \mathbf{x}_2}{2}\right\|^2\right\} d\mathbf{s} \\&= \beta^2 \left(\frac{\pi}{4\gamma}\right)^{n/2} \exp\{-\gamma\|\mathbf{x}_1 - \mathbf{x}_2\|^2\}\end{aligned}$$

$$\begin{aligned}2\|\mathbf{s} - \mathbf{x}_1\|^2 + 2\|\mathbf{s} - \mathbf{x}_2\|^2 &= 4\|\mathbf{s}\|^2 - 4\mathbf{s}^T(\mathbf{x}_1 + \mathbf{x}_2) + 2\|\mathbf{x}_1\|^2 + 2\|\mathbf{x}_2\|^2 \\&= 4\left\|\mathbf{s} - \frac{\mathbf{x}_1 + \mathbf{x}_2}{2}\right\|^2 - \|\mathbf{x}_1 + \mathbf{x}_2\|^2 + 2\|\mathbf{x}_1\|^2 + 2\|\mathbf{x}_2\|^2 \\&= 4\left\|\mathbf{s} - \frac{\mathbf{x}_1 + \mathbf{x}_2}{2}\right\|^2 + \|\mathbf{x}_1 - \mathbf{x}_2\|^2\end{aligned}$$

$$4\gamma = \frac{1}{2\sigma^2} \text{ より、 } \frac{1}{2\pi\sigma^2} = \frac{4\gamma}{\pi}$$

RBF カーネルのカーネル行列

RBF カーネル： $K_{ij} = \text{Ker}(\mathbf{x}_i, \mathbf{x}_j) := \exp\{-\gamma\|\mathbf{x}_i - \mathbf{x}_j\|^2\}$

すでに写像 $\phi(\mathbf{x})$ は構築できているため、 \mathbf{K} が半正定値行列であることを示す必要はないが、以下のように写像 $\phi(\mathbf{x})$ を使うと示すことができる。

任意の $\mathbf{v} \in \mathbb{R}^D$ に対し、

$$\begin{aligned}\mathbf{v}^T \mathbf{K} \mathbf{v} &= \sum_{i,j \in [D]} v_i v_j \exp\{-\gamma\|\mathbf{x}_i - \mathbf{x}_j\|^2\} \\&= \sum_{i,j \in [D]} v_i v_j \left(\frac{4\gamma}{\pi}\right)^{n/2} \int \exp\{-2\gamma\|\mathbf{x}_i - \mathbf{s}\|^2\} \exp\{-2\gamma\|\mathbf{x}_j - \mathbf{s}\|^2\} d\mathbf{s} \\&= \left(\frac{4\gamma}{\pi}\right)^{n/2} \int \sum_{i,j \in [D]} v_i v_j \exp\{-2\gamma\|\mathbf{x}_i - \mathbf{s}\|^2\} \exp\{-2\gamma\|\mathbf{x}_j - \mathbf{s}\|^2\} d\mathbf{s} \\&= \left(\frac{4\gamma}{\pi}\right)^{n/2} \int \left(\sum_{i \in [D]} v_i \exp\{-2\gamma\|\mathbf{x}_i - \mathbf{s}\|^2\} \right)^2 d\mathbf{s} \\&\geq 0.\end{aligned}$$

カーネル法

再掲 もし

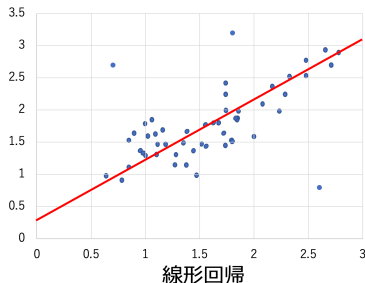
- (1) x_i と x_j から直接 $\phi(x_i)^T \phi(x_j)$ の値を計算できる
- (2) 高次元空間上での内積 $\phi(x_i)^T \phi(x_j)$ にのみ基づいて分析

⇒ 写像 $\phi(x)$ を陽に定めることなく分析が可能

(2) の主な分析法

- カーネル回帰、カーネル Ridge 回帰
- ガウス過程回帰
- サポートベクターマシン (サポートベクターマシンのスライド参照)
- カーネル主成分分析 (特徴抽出のスライド参照)

線形回帰



データ： $\{(\mathbf{x}_d, y_d)\}_{d \in [D]}$, $\mathbf{x}_d \in \mathbb{R}^n$, $y_d \in \mathbb{R}$

線形関数（アフィン関数）： $y = \mathbf{w}^T \mathbf{x} + b$

$$\tilde{\mathbf{x}} := \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix} \in \mathbb{R}^{1+n}, \quad \tilde{\mathbf{w}} := \begin{pmatrix} b \\ \mathbf{w} \end{pmatrix} \in \mathbb{R}^{1+n} \quad \text{とすると、} \quad \mathbf{w}^T \mathbf{x} + b = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$$

以下では、1次元を増やした $\tilde{\mathbf{x}}_d, \tilde{\mathbf{w}}$ を $\mathbf{x}_d, \mathbf{w} \in \mathbb{R}^n$ として考える。

線形回帰

予測器

$$y \simeq \hat{y} = \boldsymbol{w}^T \boldsymbol{x} \quad (\text{パラメタは } \boldsymbol{w})$$

誤差関数

予測誤差の2乗

$$L(y, \hat{y}) := \frac{1}{2}(y - \hat{y})^2 \quad y: \text{目標} \quad \hat{y}: \text{予測値}$$

学習

$$\min_{\boldsymbol{w}} \sum_{d \in [D]} (y_d - \boldsymbol{w}^T \boldsymbol{x}_d)^2$$

目的関数の $1/2D$ は省略

線形回帰は最小2乗法や重回帰とも呼ばれる

学習アルゴリズム

$$\mathbf{X} := \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_D^T \end{pmatrix} \in \mathbb{R}^{D \times n}, \quad \mathbf{y} := \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_D \end{pmatrix} \in \mathbb{R}^D \quad \text{とする。}$$

学習の行列表現

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

\mathbf{X} が列フルランク $\text{rank}(\mathbf{X}) = n$ と仮定

学習アルゴリズム

最適解は解析的に与えられる

$$\mathbf{w}^* := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\begin{aligned} F(\boldsymbol{w}) &:= \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w}\|^2 = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w}) \\ &= \boldsymbol{y}^T \boldsymbol{y} - 2\boldsymbol{y}^T \boldsymbol{X}\boldsymbol{w} + \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X}\boldsymbol{w} \end{aligned}$$

最適解の必要条件 $\nabla F(\boldsymbol{w}^*) = \mathbf{0}$ を考える。

$$\begin{aligned} \nabla F(\boldsymbol{w}^*) &= -2\boldsymbol{X}^T \boldsymbol{y} + 2\boldsymbol{X}^T \boldsymbol{X}\boldsymbol{w}^* = \mathbf{0} \\ \iff \boldsymbol{X}^T \boldsymbol{X}\boldsymbol{w}^* &= \boldsymbol{X}^T \boldsymbol{y} \\ \iff \boldsymbol{w}^* &= (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y} \quad (\because \text{rank}(\boldsymbol{X}) = n \text{ より逆行列は存在}) \end{aligned}$$

また、 $\nabla^2 F(\boldsymbol{w})$ は半正定値である。なぜなら、任意の $\boldsymbol{v} \in \mathbb{R}^n$ に対して

$$\boldsymbol{v}^T \nabla^2 F(\boldsymbol{w}) \boldsymbol{v} = 2\boldsymbol{v}^T \boldsymbol{X}^T \boldsymbol{X}\boldsymbol{v} = 2(\boldsymbol{X}\boldsymbol{v})^T (\boldsymbol{X}\boldsymbol{v}) \geq 0$$

よって、 $F(\boldsymbol{w})$ は凸関数となり、 $\nabla F(\boldsymbol{w}^*) = \mathbf{0}$ は最適解の必要十分条件

以上より、 $\boldsymbol{w}^* = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y}$ が最適解。

カーネル回帰 1

$\mathbf{X} := \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_D^T \end{pmatrix}$ の代わりに $\widetilde{\mathbf{X}} := \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_D)^T \end{pmatrix}$ を使う

カーネル回帰

$$\min_{\widetilde{\mathbf{w}}} \|\mathbf{y} - \widetilde{\mathbf{X}}\widetilde{\mathbf{w}}\|^2$$

$\widetilde{\mathbf{w}}^* := (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T \mathbf{y}$ としたいが、問題がある。

- 高次元空間の次元 $|V|$ がデータ数 D より大きいと、解が一意に定まらない。
($\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}}$ が正則にならない)
- 我々が使えるのはカーネル行列 $\mathbf{K} := \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T$
 $K_{ij} := \text{Ker}(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$

カーネル回帰 2

以下では、 $D < |V|$ で、 \mathbf{K} は正則を仮定する (RBF カーネルなど)

カーネル回帰

多数ある最適解の中で、 $\tilde{\mathbf{w}}$ のノルムが小さなものを採用する。

$$\begin{aligned} \min_{\tilde{\mathbf{w}}} \quad & \frac{1}{2} \|\tilde{\mathbf{w}}\|^2 \\ \text{s.t.} \quad & \mathbf{y} = \widetilde{\mathbf{X}} \tilde{\mathbf{w}} \end{aligned}$$

ラグランジュの未定乗数法

$$L(\tilde{\mathbf{w}}, \boldsymbol{\lambda}) := \frac{1}{2} \tilde{\mathbf{w}}^T \tilde{\mathbf{w}} + \boldsymbol{\lambda}^T (\mathbf{y} - \widetilde{\mathbf{X}} \tilde{\mathbf{w}})$$

最適解の必要十分条件

$$\nabla_{\tilde{\mathbf{w}}} L(\tilde{\mathbf{w}}, \boldsymbol{\lambda}) := \tilde{\mathbf{w}} - \widetilde{\mathbf{X}}^T \boldsymbol{\lambda} = \mathbf{0}$$

$$\nabla_{\boldsymbol{\lambda}} L(\tilde{\mathbf{w}}, \boldsymbol{\lambda}) := \mathbf{y} - \widetilde{\mathbf{X}} \tilde{\mathbf{w}} = \mathbf{0}$$

これを解くと、

$$\tilde{\mathbf{w}}^* = \widetilde{\mathbf{X}}^T \mathbf{K}^{-1} \mathbf{y}$$

新しい \mathbf{x} に対する $f(\mathbf{x}) = (\tilde{\mathbf{w}}^*)^T \tilde{\mathbf{x}}$ の値は次のように計算できる。

$$\begin{aligned} f(\mathbf{x}) &= (\tilde{\mathbf{w}}^*)^T \tilde{\mathbf{x}} \\ &= \tilde{\mathbf{x}}^T \widetilde{\mathbf{X}}^T \mathbf{K}^{-1} \mathbf{y} \end{aligned}$$

$\boldsymbol{\alpha} := \mathbf{K}^{-1} \mathbf{y}$ とすると

$$\begin{aligned} &= \tilde{\mathbf{x}}^T \widetilde{\mathbf{X}}^T \boldsymbol{\alpha} \\ &= \tilde{\mathbf{x}}^T \sum_{d \in [D]} \alpha_d \tilde{\mathbf{x}}_d \\ &= \sum_{d \in [D]} \alpha_d \text{Ker}(\mathbf{x}, \mathbf{x}_d) \end{aligned}$$

基底関数による意味付け

基底関数： $f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_M(\mathbf{x})$

予測器 ($\mathbf{w} \in \mathbb{R}^M$ はパラメタ)： $f(\mathbf{x}) := \sum_{k \in [M]} w_k f_k(\mathbf{x})$

$$\phi(\mathbf{x}) := \begin{pmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{pmatrix} \text{ とすると、 } f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

2乗損失を使ったパラメタ \mathbf{w} を学習は線形回帰となる。

$$\min_{\mathbf{w}} \sum_{d \in [D]} (y_d - \mathbf{w}^T \phi(\mathbf{x}_d))^2 \quad \Longleftrightarrow \quad \min_{\mathbf{w}} \|\mathbf{y} - \widetilde{\mathbf{X}}\mathbf{w}\|^2$$

$\mathcal{F} = \{f(\mathbf{x}) : \exists \mathbf{w} \in \mathbb{R}^M, f(\mathbf{x}) = \sum_{k \in [M]} w_k f_k(\mathbf{x})\}$ とすると、

$$\min_{f \in \mathcal{F}} \sum_{d \in [D]} (y_d - f(\mathbf{x}_d))^2$$

基底関数の例

例 1 : 線形回帰

特徴量 $\boldsymbol{x} \in \mathbb{R}^n$

$$f_k(\boldsymbol{x}) = \boldsymbol{e}_k \quad (k = 1, 2, \dots, n)$$

$\boldsymbol{e}_k \in \mathbb{R}^n$ は、 k 番目の要素が 1 で、残りの要素が 0 のベクトル

例 2 : 多項式回帰

特徴量 $x \in \mathbb{R}$

$$f_0(x) = 1, f_1(x) = x, f_2(x) = x^2, \dots, f_M(x) = x^M$$

例 3 : RBF カーネルを用いたカーネル回帰

特徴量 $\boldsymbol{x} \in \mathbb{R}^n$

$$f_{\boldsymbol{s}}(\boldsymbol{x}) := \exp\{-\gamma\|\boldsymbol{x} - \boldsymbol{s}\|^2\} \quad (\boldsymbol{s} \in \mathbb{R}^n)$$

無限個の基底関数

カーネル Ridge 回帰

正則化項を加えた線形回帰

Ridge 回帰

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$$

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

カーネル Ridge 回帰

$$\min_{\tilde{\mathbf{w}}} \|\mathbf{y} - \tilde{\mathbf{X}}\tilde{\mathbf{w}}\|^2 + \lambda \|\tilde{\mathbf{w}}\|^2$$

$$\tilde{\mathbf{w}}^* = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

- $\lambda > 0$ で、 $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I}$ は正則 \rightarrow 最適解は一つ
- $\mathbf{K} = \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T$ を使って計算できない

カーネル Ridge 回帰

$K = \widetilde{X}\widetilde{X}^T$ に対して、次の性質が成り立つ ($\lambda > 0$)。

$$\begin{aligned}\tilde{w}^* &= (\widetilde{X}^T \widetilde{X} + \lambda I)^{-1} \widetilde{X}^T y \\ &= \widetilde{X}^T (K + \lambda I)^{-1} y\end{aligned}$$

※ 上式では $I \in \mathbb{R}^{|V| \times |V|}$, 下式では $I \in \mathbb{R}^{D \times D}$ であることに注意

新しい x に対する $f(x) = (\tilde{w}^*)^T \tilde{x}$ の値は次のように計算できる。

$$\begin{aligned}f(x) &= (\tilde{w}^*)^T \tilde{x} \\ &= \tilde{x}^T \widetilde{X}^T (K + \lambda I)^{-1} y\end{aligned}$$

$\alpha := (K + \lambda I)^{-1} y$ とすると

$$\begin{aligned}&= \tilde{x}^T \sum_{d \in [D]} \alpha_d \tilde{x}_d \\ &= \sum_{d \in [D]} \alpha_d \text{Ker}(x, x_d)\end{aligned}$$

式変形

Sherman–Morrison–Woodbury の公式より、

$$\begin{aligned}(\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} &= \frac{1}{\lambda} \mathbf{I} - \frac{1}{\lambda} \mathbf{I} \widetilde{\mathbf{X}}^T \left(\mathbf{I} + \widetilde{\mathbf{X}} \frac{1}{\lambda} \mathbf{I} \widetilde{\mathbf{X}}^T \right)^{-1} \widetilde{\mathbf{X}} \frac{1}{\lambda} \mathbf{I} \\ &= \frac{1}{\lambda} \mathbf{I} - \frac{1}{\lambda} \widetilde{\mathbf{X}}^T \left(\lambda \mathbf{I} + \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T \right)^{-1} \widetilde{\mathbf{X}}\end{aligned}$$

よって、

$$\begin{aligned}(\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} \widetilde{\mathbf{X}}^T \mathbf{y} &= \left(\frac{1}{\lambda} \mathbf{I} - \frac{1}{\lambda} \widetilde{\mathbf{X}}^T \left(\lambda \mathbf{I} + \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T \right)^{-1} \widetilde{\mathbf{X}} \right) \widetilde{\mathbf{X}}^T \mathbf{y} \\ &= \frac{1}{\lambda} \widetilde{\mathbf{X}}^T \mathbf{y} - \frac{1}{\lambda} \widetilde{\mathbf{X}}^T (\lambda \mathbf{I} + \mathbf{K})^{-1} \mathbf{K} \mathbf{y} \\ &= \frac{1}{\lambda} \widetilde{\mathbf{X}}^T (\lambda \mathbf{I} + \mathbf{K})^{-1} (\lambda \mathbf{I} + \mathbf{K}) \mathbf{y} - \frac{1}{\lambda} \widetilde{\mathbf{X}}^T (\lambda \mathbf{I} + \mathbf{K})^{-1} \mathbf{K} \mathbf{y} \\ &= \frac{1}{\lambda} \widetilde{\mathbf{X}}^T (\lambda \mathbf{I} + \mathbf{K})^{-1} \lambda \mathbf{I} \mathbf{y} \\ &= \widetilde{\mathbf{X}}^T (\lambda \mathbf{I} + \mathbf{K})^{-1} \mathbf{y}\end{aligned}$$

カーネル回帰とカーネル Ridge 回帰

カーネル回帰で得られる重み：

$$\begin{aligned}\tilde{\mathbf{w}}^* &:= \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \quad \frac{1}{2} \|\tilde{\mathbf{w}}\|^2 \\ \text{s.t.} \quad & \mathbf{y} = \widetilde{\mathbf{X}} \tilde{\mathbf{w}}\end{aligned}$$

カーネル Ridge 回帰で得られる重み：

$$\tilde{\mathbf{w}}^*(\lambda) := \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \|\mathbf{y} - \widetilde{\mathbf{X}} \tilde{\mathbf{w}}\|^2 + \lambda \|\tilde{\mathbf{w}}\|^2$$

このとき、次の関係が成り立つ。

$$\tilde{\mathbf{w}}^* = \lim_{\lambda \rightarrow +0} \tilde{\mathbf{w}}^*(\lambda)$$

説明：

$$\begin{aligned}\tilde{\mathbf{w}}^* &= \widetilde{\mathbf{X}}^T \mathbf{K}^{-1} \mathbf{y} \\ \tilde{\mathbf{w}}^*(\lambda) &= \widetilde{\mathbf{X}}^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y} \quad (\lambda > 0)\end{aligned}$$

表現定理

- カーネル回帰もカーネル Ridge 回帰も、 $\mathbf{w}^* = \widetilde{\mathbf{X}}^T \boldsymbol{\alpha}$ という形で表現できる。
- カーネル法では多くの場合同様の性質が成り立ち、表現定理という。

$\mathbf{w}^* = \widetilde{\mathbf{X}}^T \boldsymbol{\alpha}$ を前提とすれば、

カーネル回帰：

$$\min_{\boldsymbol{\alpha}} \|\mathbf{y} - \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T \boldsymbol{\alpha}\|^2$$

カーネル Ridge 回帰：

$$\min_{\boldsymbol{\alpha}} \|\mathbf{y} - \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T \boldsymbol{\alpha}\|^2 + \lambda \|\widetilde{\mathbf{X}}^T \boldsymbol{\alpha}\|^2$$

カーネル行列 \mathbf{K} が正則とすれば、最適解は簡単に導出できる。

$$\begin{aligned}\boldsymbol{\alpha}^* &= \mathbf{K}^{-1} \mathbf{y} \\ \boldsymbol{\alpha}^*(\lambda) &= (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}\end{aligned}$$

確率モデルに基づいた解析

データ $\{(\mathbf{x}_d, y_d)\}_{d \in [D]}$ が次の確率モデルから生成されたサンプル（標本）とする。

$$y_d = \mathbf{w}^T \mathbf{x}_d + \epsilon_d \quad (d \in [D])$$

$$\iff \mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$

ϵ はノイズ（測定誤差など）を表す確率変数

ここで確率変数 $\epsilon \in \mathbb{R}^D$ が従う仮定として、次のものを考える。

誤差項に対する仮定

- $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

線形回帰・カーネル回帰と最尤推定

このページの話は、線形回帰とカーネル回帰で共通
(チルダがついているかどうかの違い)

仮定：

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon, \quad \epsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

このとき、

$$\mathbf{y} \sim N(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I})$$

尤度関数：
$$l(\mathbf{y}; \mathbf{w}) = \frac{1}{(2\pi\sigma^2)^{D/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \right\}$$

対数尤度関数：
$$\log l(\mathbf{y}; \mathbf{w}) = -\frac{D}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

対数尤度関数の最大化は線形回帰と等価

$$\max_{\mathbf{w}} \log l(\mathbf{y}; \mathbf{w})$$

つまり、線形回帰やカーネル回帰は最尤推定法とみなせる

MAP 推定

ベイズの定理を使うと、データ $\mathcal{D} = \{(\mathbf{x}_d, y_d)\}_{d \in [D]}$ のもとでのパラメタ \mathbf{w} の確率は、次のようになる。

$$P(\mathbf{w}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathbf{w})P(\mathbf{w})}{P(\mathcal{D})}$$

- $P(\mathbf{w})$: パラメタ \mathbf{w} の事前確率
- $P(\mathbf{w}|\mathcal{D})$: データ \mathcal{D} が与えられた状況でのパラメタ \mathbf{w} の事後確率

MAP(maximum a posteriori) 推定

事後確率 $P(\mathbf{w}|\mathcal{D})$ が最大となるパラメタ \mathbf{w} を採用する

$$\max_{\mathbf{w}} \log P(\mathbf{w}|\mathcal{D}) \iff \max_{\mathbf{w}} \log P(\mathcal{D}|\mathbf{w}) + \log P(\mathbf{w})$$

カーネル回帰と MAP 推定

仮定

- パラメタの事前分布： $\tilde{\mathbf{w}} \sim N(\mathbf{0}, \tau^2 \mathbf{I})$
- $\mathbf{y} = \tilde{\mathbf{X}}\tilde{\mathbf{w}}$ 誤差は無し

$$\log P(\mathcal{D}|\tilde{\mathbf{w}}) + \log P(\tilde{\mathbf{w}}) = \begin{cases} -\frac{n}{2} \log 2\pi\tau^2 - \frac{1}{2\tau^2} \tilde{\mathbf{w}}^T \tilde{\mathbf{w}} & (\mathbf{y} = \tilde{\mathbf{X}}\tilde{\mathbf{w}}) \\ -\infty & (\mathbf{y} \neq \tilde{\mathbf{X}}\tilde{\mathbf{w}}) \end{cases}$$

この関数の最大化は、 $\mathbf{y} = \tilde{\mathbf{X}}\tilde{\mathbf{w}}$ のもとでの $\|\tilde{\mathbf{w}}\|^2$ の最小化と等価

つまり、カーネル回帰は上記の仮定での MAP 推定とみなせる

Ridge 回帰・カーネル Ridge 回帰と MAP 推定

このページの話は、Ridge 回帰とカーネル Ridge 回帰で共通
(チルダがついているかどうかの違い)

仮定

- パラメタの事前分布： $\mathbf{w} \sim N(\mathbf{0}, \tau^2 \mathbf{I})$
- $\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon$, $\epsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

$$\begin{aligned}\log P(\mathcal{D}|\mathbf{w}) + \log P(\mathbf{w}) &= -\frac{D}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) \\ &\quad - \frac{n}{2} \log 2\pi\tau^2 - \frac{1}{2\tau^2}\mathbf{w}^T\mathbf{w} \\ &= -\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) - \frac{1}{2\tau^2}\mathbf{w}^T\mathbf{w} + \text{定数}\end{aligned}$$

この関数の最大化は、 λ を適当に定めた $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda\|\mathbf{w}\|^2$ の最小化と等価
つまり、(カーネル) Ridge 回帰は上記の仮定の元での MAP 推定とみなせる

ガウス過程回帰

訓練データ： $\{(\boldsymbol{x}_d, y_d)\}_{d \in [D]}$

- y_d の平均は 0 となるように定数を引いておく。

このデータで学習を行い、新しい \boldsymbol{x} に対する y の「分布」を構築したい

仮定：

- \boldsymbol{x} によって決まる y はガウス過程である

ガウス過程

$y \in \mathbb{R}$ は $\boldsymbol{x} \in \mathbb{R}^n$ によって定まる確率変数とする。

任意に有限個の $\boldsymbol{x}_1, \boldsymbol{x}_2, \dots, \boldsymbol{x}_k$ を選んだとき、 (y_1, y_2, \dots, y_k) が多変量正規分布に従うとき、ガウス過程という。

ガウス過程回帰

ガウス過程回帰

データ $\{x_d, y_d\}_{d \in [D]}$ と任意の x に対する y において、次の関係が成り立つとする。

$$\begin{pmatrix} \mathbf{y} \\ y \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{0} \\ 0 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right)$$

$\mathbf{y}, \mathbf{0} \in \mathbb{R}^D$, $y, 0 \in \mathbb{R}$, $\Sigma_{11} \in \mathbb{R}^{D \times D}$, $\Sigma_{12} \in \mathbb{R}^{D \times 1}$, $\Sigma_{21} \in \mathbb{R}^{1 \times D}$, $\Sigma_{22} \in \mathbb{R}$

\mathbf{y} が観測された状態での y の条件付き分布は次のようになる。

$$y | \mathbf{y} \sim N(\Sigma_{21} \Sigma_{11}^{-1} \mathbf{y}, \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})$$

これを、 x における y の予測分布とする。

分散・共分散行列

論点：

- $\{\mathbf{x}_d\}_{d \in [D]}$ から分散・共分散行列 Σ をどう作るのか
- \mathbf{x} が「近い」とき共分散が大きい（より関係している）

カーネル関数を使う（半正定値性が保証されることに注意）。

$$\Sigma_{ij} := \text{Ker}(\mathbf{x}_i, \mathbf{x}_j)$$

カーネルの例

- $\text{Ker}(\mathbf{x}_a, \mathbf{x}_b) := \mathbf{x}_a^T \mathbf{x}_b$ 線形カーネル
- $\text{Ker}(\mathbf{x}_a, \mathbf{x}_b) := (\mathbf{x}_a^T \mathbf{x}_b + 1)^d \quad (d \in \mathbb{N})$ 多項式カーネル
- $\text{Ker}(\mathbf{x}_a, \mathbf{x}_b) := \exp\{-\gamma \|\mathbf{x}_a - \mathbf{x}_b\|^2\}$ RBF カーネル、ガウシアンカーネル

- 線形カーネルだと自由度が足りず \mathbf{y} を再現できない（ Σ_{11} が非正則）
- 多くの場合 RBF カーネルを使う

計算手順

データ： $\{(\mathbf{x}_d, y_d)\}_{d \in [D]}$

予測したい点： \mathbf{x}

$\mathbf{K} \in \mathbb{R}^{D \times D}$, $\mathbf{k} \in \mathbb{R}^D$, $k \in \mathbb{R}$ を次のように定義する。

$$K_{ij} := \text{Ker}(\mathbf{x}_i, \mathbf{x}_j) \quad (i, j \in [D])$$

$$k_i := \text{Ker}(\mathbf{x}_i, \mathbf{x}) \quad (i \in [D])$$

$$k := \text{Ker}(\mathbf{x}, \mathbf{x})$$

多変量正規分布の分散・共分散行列は次のようになる。

$$\begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} = \begin{pmatrix} \mathbf{K} & \mathbf{k} \\ \mathbf{k}^T & k \end{pmatrix}$$

多変量正規分布の条件付き確率の式より、データ \mathbf{y} のもとでの y の分布の予測は

$$y|\mathbf{y} \sim N(\mathbf{k}\mathbf{K}^{-1}\mathbf{y}, k - \mathbf{k}^T\mathbf{K}^{-1}\mathbf{k})$$

となる。

ガウス過程回帰の意味づけ

$\phi(\mathbf{x}) \in V$ とする。

モデル

- $\tilde{w}_i \sim N(0, 1) \quad (i \in |V|)$
- $y = \tilde{\mathbf{w}}^T \phi(\mathbf{x}), \quad \mathbf{y} = \widetilde{\mathbf{X}} \tilde{\mathbf{w}}$

$$\widetilde{\mathbf{X}} := \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_D)^T \end{pmatrix} \in \mathbb{R}^{D \times |V|},$$

$$\Rightarrow \mathbf{y} \sim N(\mathbf{0}, \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T)$$

よって、

$$\boldsymbol{\Sigma} = \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T \in \mathbb{R}^{D \times D}$$

$$\Sigma_{ij} := \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = \text{Ker}(\mathbf{x}_i, \mathbf{x}_j)$$

誤差を含むモデル

モデル

- $\tilde{w}_i \sim N(0, 1) \quad (i \in |V|), \quad \epsilon \sim N(0, \sigma^2),$
- $y = \tilde{w}^T \phi(x) + \epsilon, \quad \mathbf{y} = \widetilde{\mathbf{X}} \tilde{\mathbf{w}} + \boldsymbol{\epsilon}$

$$\widetilde{\mathbf{X}} := \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \phi(\mathbf{x}_2)^T \\ \vdots \\ \phi(\mathbf{x}_D)^T \end{pmatrix} \in \mathbb{R}^{D \times |V|}, \quad \boldsymbol{\epsilon} := \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_D \end{pmatrix} \in \mathbb{R}^D,$$

$$\Rightarrow \mathbf{y} \sim N(\mathbf{0}, \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T + \sigma^2 \mathbf{I})$$

よって、

$$\begin{aligned} \boldsymbol{\Sigma} &= \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T + \sigma^2 \mathbf{I} \\ \Sigma_{ij} &= \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) + \sigma^2 \delta_{ij} \\ &= \text{Ker}(\mathbf{x}_i, \mathbf{x}_j) + \sigma^2 \delta_{ij} \end{aligned}$$

ベイズ推定 1

ガウス過程回帰は、パラメータ w をベイズ推定したうえで、新しい x に対する y の分布を求めていると捉えることができる。

誤差付きの場合で説明：

まず、前述した条件付き確率の計算手順を \tilde{X} を使って書き直す。
($\phi(x)$ を \tilde{x} と表記する)

$$\begin{pmatrix} y \\ y \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{0} \\ 0 \end{pmatrix}, \begin{pmatrix} \tilde{X}\tilde{X}^T + \sigma^2 I & \tilde{X}\tilde{x} \\ \tilde{x}^T \tilde{X}^T & \tilde{x}^T \tilde{x} + \sigma^2 \end{pmatrix} \right)$$

このとき、多変量正規分布の条件付き分布は、

$$\begin{aligned} y|y &\sim N(\Sigma_{21}\Sigma_{11}^{-1}y, \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}) \\ &= N(\tilde{x}^T \tilde{X}^T (\tilde{X}\tilde{X}^T + \sigma^2 I)^{-1}y, \tilde{x}^T \tilde{x} + \sigma^2 - \tilde{x}^T \tilde{X}^T (\tilde{X}\tilde{X}^T + \sigma^2 I)^{-1} \tilde{X}\tilde{x}) \end{aligned}$$

である。

ベイズ推定 2

一方、ベイズ推定より、パラメタ w の事後分布は次のようになる。

$$\tilde{w}|\mathbf{y} \sim N((\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \sigma^2 \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}, \sigma^2 (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \sigma^2 \mathbf{I})^{-1})$$

このとき、 y の確率は

$$y|\mathbf{y} \sim N(\tilde{\mathbf{x}}^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \sigma^2 \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}, \sigma^2 \tilde{\mathbf{x}}^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \sigma^2 \mathbf{I})^{-1} \tilde{\mathbf{x}} + \sigma^2)$$

この2つの正規分布は平均も分散も等しいので、ガウス過程回帰のパラメタをベイズ推定しているといえる。

平均：

$$\tilde{\mathbf{x}}^T \tilde{\mathbf{X}}^T (\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T + \sigma^2 \mathbf{I})^{-1} \mathbf{y} = \tilde{\mathbf{x}}^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \sigma^2 \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

分散：

$$\tilde{\mathbf{x}}^T \tilde{\mathbf{x}} + \sigma^2 - \tilde{\mathbf{x}}^T \tilde{\mathbf{X}}^T (\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T + \sigma^2 \mathbf{I})^{-1} \tilde{\mathbf{X}} \tilde{\mathbf{x}} = \sigma^2 \tilde{\mathbf{x}}^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \sigma^2 \mathbf{I})^{-1} \tilde{\mathbf{x}} + \sigma^2$$

式変形

平均：

すでに、 $\widetilde{\mathbf{X}}^T (\widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T + \sigma^2 \mathbf{I})^{-1} \mathbf{y} = (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}} + \sigma^2 \mathbf{I})^{-1} \widetilde{\mathbf{X}}^T \mathbf{y}$ は説明しているので等式は成り立つ。

分散：

Sherman–Morrison–Woodbury の公式より、

$$(\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}} + \sigma^2 \mathbf{I})^{-1} = \frac{1}{\sigma^2} \mathbf{I} - \frac{1}{\sigma^2} \widetilde{\mathbf{X}}^T \left(\sigma^2 \mathbf{I} + \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T \right)^{-1} \widetilde{\mathbf{X}}$$

よって、

$$\begin{aligned} & \sigma^2 \widetilde{\mathbf{x}}^T (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}} + \sigma^2 \mathbf{I})^{-1} \widetilde{\mathbf{x}} + \sigma^2 \\ &= \sigma^2 \widetilde{\mathbf{x}}^T \left(\frac{1}{\sigma^2} \mathbf{I} - \frac{1}{\sigma^2} \widetilde{\mathbf{X}}^T \left(\sigma^2 \mathbf{I} + \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T \right)^{-1} \widetilde{\mathbf{X}} \right) \widetilde{\mathbf{x}} + \sigma^2 \\ &= \sigma^2 \widetilde{\mathbf{x}}^T \widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^T \widetilde{\mathbf{X}}^T \left(\sigma^2 \mathbf{I} + \widetilde{\mathbf{X}} \widetilde{\mathbf{X}}^T \right)^{-1} \widetilde{\mathbf{X}} \widetilde{\mathbf{x}} + \sigma^2 \end{aligned}$$

MAP 推定とベイズ推定

カーネル回帰・カーネル Ridge 回帰とガウス過程回帰の違いは、パラメタを MAP 推定しているかベイズ推定しているかの違いと捉えることができる

	カーネル	w	誤差	推定
線形回帰・カーネル回帰	有・無	確定的	あり	最尤推定 ¹
カーネル回帰	有	確率的	なし	MAP 推定
(カーネル) Ridge 回帰	有・無	確率的	あり	MAP 推定
ガウス過程回帰	有	確率的	なし	ベイズ推定
ガウス過程回帰＋誤差	有・無	確率的	あり	ベイズ推定

※ $n < D < |V|$ を想定

¹カーネル有りのときはノルム最小化も行う

参考：多変量正規分布

$N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ：平均ベクトル $\boldsymbol{\mu}$ 、分散共分散行列 $\boldsymbol{\Sigma}$ の正規分布

多変量正規分布

$\boldsymbol{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ のとき、確率密度関数は

$$f(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) := \frac{1}{(2\pi)^{n/2}} \frac{1}{\det(\boldsymbol{\Sigma})^{1/2}} \exp \left\{ -\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}) \right\}$$

$\boldsymbol{x} \sim N(\boldsymbol{\mu}, \sigma^2 \boldsymbol{I})$ のとき、確率密度関数は

$$f(\boldsymbol{x}; \boldsymbol{\mu}, \sigma^2 \boldsymbol{I}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2}(\boldsymbol{x} - \boldsymbol{\mu})^T(\boldsymbol{x} - \boldsymbol{\mu}) \right\}$$

参考：多変量正規分布

条件付き確率

$\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ のとき、 $\mathbf{x} := \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$, $\boldsymbol{\mu} := \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$, $\boldsymbol{\Sigma} := \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$

x_1 のもとでの x_2 の確率分布：

$$x_2 | x_1 \sim N(\mu_2 + \Sigma_{21} \Sigma_{11}^{-1} (x_1 - \mu_1), \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})$$

ベイズ推定

$\mathbf{w} \sim N(\mathbf{0}, \mathbf{I})$, $\mathbf{y} | \mathbf{w} \sim N(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I})$ のとき、

$$\mathbf{y} \sim N(\mathbf{0}, \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I})$$

$$\mathbf{w} | \mathbf{y} \sim N((\mathbf{X}^T \mathbf{X} + \sigma^2 \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}, \sigma^2 (\mathbf{X}^T \mathbf{X} + \sigma^2 \mathbf{I})^{-1})$$

参考：多変量正規分布

ベイズ推定

N は平均ベクトルと共分散行列と引数にとる正規分布

Σ ：分散・共分散行列、 Λ ：精度行列

$$\boldsymbol{w} \sim N(\boldsymbol{\mu}, \Sigma_1)$$

$$\boldsymbol{y}|\boldsymbol{w} \sim N(\boldsymbol{X}\boldsymbol{w} + \boldsymbol{b}, \Sigma_2) \quad \text{のとき、}$$

$$\boldsymbol{y} \sim N(\boldsymbol{X}\boldsymbol{\mu} + \boldsymbol{b}, \boldsymbol{X}\Sigma_1\boldsymbol{X}^T + \Sigma_2)$$

$$\boldsymbol{w}|\boldsymbol{y} \sim N((\Lambda_1 + \boldsymbol{X}^T\Lambda_2\boldsymbol{X})^{-1}(\boldsymbol{X}^T\Lambda_2(\boldsymbol{y} - \boldsymbol{b}) + \Lambda_1\boldsymbol{\mu}), (\Lambda_1 + \boldsymbol{X}^T\Lambda_2\boldsymbol{X})^{-1})$$

$$\text{ただし、} \Lambda_1 := \Sigma_1^{-1}, \Lambda_2 := \Sigma_2^{-1}$$