

カーネル多変量解析

第2章

カーネル多変量解析の仕組み(2)

サ ミンソン

2. 3確率モデルからの導入

- 多変量解析の目的・・・与えられたデータに基づいて、背後にある構造を推論
- ただし・・・データにノイズがのる、関数のうちの有限個の点でしか関数値が与えられない・・・などの問題



データの生成過程を確率分布を用いてモデル化する

(a)線形モデルのベイズ推論

▪まず、 $y = \mathbf{w}^T \mathbf{x}$ (式(1.1)) による関数近似の確率モデルを考える。

出力 $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ に、ランダムノイズ ε がのった、

$$y = f(\mathbf{x}) + \varepsilon \quad (2.28)$$

ε が独立な正規分布であるとする、 \mathbf{x} 、 f が与えられたもとで、分散を σ^2 として、

▪ y の条件付き確率精度

$$p(y | \mathbf{x}; f) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - f(\mathbf{x}))^2}{2\sigma^2}\right) \quad (2.29)$$

続き

- 線形モデル $f(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x}$
 - パラメータ \boldsymbol{w} を決めることと f を決めることは等価
 - \boldsymbol{w} の各成分が独立に平均 0、分散が $1/\lambda$ である正規分布から生成されたとする。

\boldsymbol{w} の分布を

$$p(\boldsymbol{w}) = \left(\frac{\lambda}{2\pi}\right)^{d/2} \exp\left(-\frac{\lambda}{2}\|\boldsymbol{w}\|^2\right) \quad (2.30)$$

とする。

サンプルの生成過程

- まず、 $p(\mathbf{w})$ に従ってパラメータ \mathbf{w} がランダムに決められる。
- \mathbf{w} で決まる関数 $f(\mathbf{x})$ を n 個の点で計算し、(2.29)の分布に従ってサンプル出力 y が観測される。

- パラメータとデータの同時確率分布

$$p(y^{(1)}, \dots, y^{(n)}, \mathbf{w}) = p(\mathbf{w}) \prod_{j=1}^n p(y^{(j)} | \mathbf{x}^{(j)}; f) \quad (2.31)$$

- データやパラメータの生成過程を確率分布で表す
→生成モデル

続き

- 生成モデル関数のパラメータが先に決まる
サンプルはその後に決まる
サンプルから関数を推定・・・ベイズの公式を利用

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)} \quad (2.32)$$

Bを観測した事後におけるAの分布・・・事後分布
最も大きくするAの値を見つける・・・MAP推定
—右辺の分子を最大

続き

- パラメータ w がA, 出力 y がBに相当する。
- ベイズの公式の右辺の分子は式(2.31)になり、対数を取ると、

$$\begin{aligned} & \sum_{i=1}^n \log p(y^{(i)} | \mathbf{x}^{(i)}; f) + \log p(\mathbf{w}) \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - f(\mathbf{x}^{(i)}))^2 - \frac{\lambda}{2} \|\mathbf{w}\|^2 + \text{定数} \end{aligned} \quad (2.33)$$

- MAP推定は、線形モデルで正規化付きの二乗誤差を最小にすること

(b) 正規過程からカーネルへ

- 入力値 $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$ に対して、 $f(\mathbf{x}^{(1)}), f(\mathbf{x}^{(2)}), \dots, f(\mathbf{x}^{(n)})$ が平均0、分散共分散行列 $V = (V_{ij})_{i,j=1,\dots,n}$ をもつ多次元正規分布に従うとする。
- $f(\mathbf{x}^{(i)})$ は正規確率変数、 $f(\mathbf{x}^{(i)})$ と $f(\mathbf{x}^{(j)})$ の共分散が

$$V_{ij} = \mathbb{E}_{f(\mathbf{x}^{(i)}), f(\mathbf{x}^{(j)})} [f(\mathbf{x}^{(i)}) f(\mathbf{x}^{(j)})] \quad (2.34)$$

となる。(Eはxに関する期待値)

- この確率過程をx上の正規過程と呼ぶ

つづき

- 任意の入力点 \mathbf{x}^{new} を固定し、その点での関数値の値の分布を調べる。サンプルに対するグラム行列を $K=(K_{ij})_{i,j=1,\dots,n}$ とおくと、サンプル $\mathbf{x}^{(i)}$ と \mathbf{x}^{new} に関する成分は $k(\mathbf{x}^{(i)}, \mathbf{x}^{\text{new}})$ で与えられる。式2.29と式2.31に相当する関数値とデータは

$$\begin{aligned} & p(y^{(1)}, \dots, y^{(n)}, f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(n)}), f(\mathbf{x}^{\text{new}})) \\ &= p(f(\mathbf{x}^{(1)}), f(\mathbf{x}^{(2)}), \dots, f(\mathbf{x}^{(n)}), f(\mathbf{x}^{\text{new}})) \prod_{i=1}^n p(y^{(i)} | \mathbf{x}^{(i)}; f) \quad (2.36) \end{aligned}$$

と書ける。

つづき

- 分布は全て正規分布→サンプルが与えられたもとでの関数値fの事後分布も正規分布。関数値fの平均は

$$\begin{aligned} E_{f(\mathbf{x}^{\text{new}})} [f(\mathbf{x}^{\text{new}}) | \mathcal{D}] \\ = \sum_{i=1}^n \{ (K + \sigma^2 I_n)^{-1} \mathbf{y} \}_i k(\mathbf{x}^{(i)}, \mathbf{x}^{\text{new}}) \end{aligned} \quad (2.37)$$

となる。(Dは、与えられたサンプル。)

- 関数近似の場合の $\boldsymbol{\alpha} = (K + \lambda I_n)^{-1} \mathbf{y}$ (1.17)

の $\boldsymbol{\alpha}$ を $f(\mathbf{x}) = \sum \alpha_i k(\mathbf{x}^{(i)}, \mathbf{x})$ に代入したものになっている。

- カーネル関数の重みつき和のモデルを正規化付きで求めた関数近似の結果 = 正規過程を事前分布として用いた場合のMAP推定の結果

2. 4 汎化能力の評価とモデル選択

- データ解析においては、汎化能力が重要
— サンプルにフィットするだけでなく、背後にある構造を正しく抽出することが求められる。
- モデル選択・・・汎化能力を高めるためにモデルの複雑度を調整すること
- 汎化能力の理論的な話・・・7章
汎化能力を評価する方法を説明

(a) クロスバリデーション

- 汎化能力は学習に使ったサンプル以外のデータに対する性能
 - ① 学習用のサンプルで学習
 - ② 残しておいたテスト用のサンプルで性能評価を行う
- 問題点
 - ・テストデータを残しすぎ→学習データが少ない
 - ・テストデータが少なすぎ→評価結果が不安定
- 学習用サンプルとテスト用サンプルの分け方をいろいろ変えて得たテスト誤差を平均する・・・クロスバリデーション

k-foldクロスバリデーション法 (k-fold CV)

- 「1」サンプルをk個のデータに分ける
- 「2」 $i=1, \dots, k$ に対し以下を繰り返す
 - (1) i 番目のグループを除いたデータで学習を行う
 - (2) i 番目のグループでテスト誤差を評価し r_i とおく
- 「3」 $\sum_{i=1}^k r_i / k$ をテスト誤差の推定値とする

(b)線形モデルのleave-one-out クロスバリデーション

・k-fold CV法は、一つだけのサンプルをテストデータとして除いておく方法・・・leave-one-outクロスバリデーションと呼ぶ

・ $y=f(x)$ を学習したとする。

i番目のサンプルの入力 $x^i \rightarrow \tilde{y}^{(i)}=f(x^{(i)})$

= サンプルの出力 $y^{(i)}$ のノイズ成分を除去した推定値

・カーネル回帰の場合、 $\tilde{y}^{(j)} = \sum_{i=1}^n \alpha_i k(x^{(i)}, x^{(j)})$ と

$$\alpha = (K + \lambda I_n)^{-1} y \quad (1.17) \quad \text{より、}$$

$$\tilde{y} = K \alpha = (K + \lambda I_n)^{-1} K y \quad (2.38) \quad \text{が得られる}$$

つづき

- 一般に $\tilde{y} = Hy$ という線形関係があるとする。
カーネル回帰の場合は $H = (K + \lambda I_n)^{-1} K$ である。
- leave-one-out CV誤差は、学習サンプルとテストサンプルを分ける手続きなしで、 $y^{(i)}$ と $\tilde{y}^{(i)}$ の重み付きの誤差平均

$$CV = \frac{1}{n} \sum_{i=1}^n \left(\frac{y^{(i)} - \tilde{y}^{(i)}}{1 - H_{ii}} \right)^2 \quad (2.39)$$

- ただし、 H_{ii} はHの第i対角成分である。

(c) 具体例

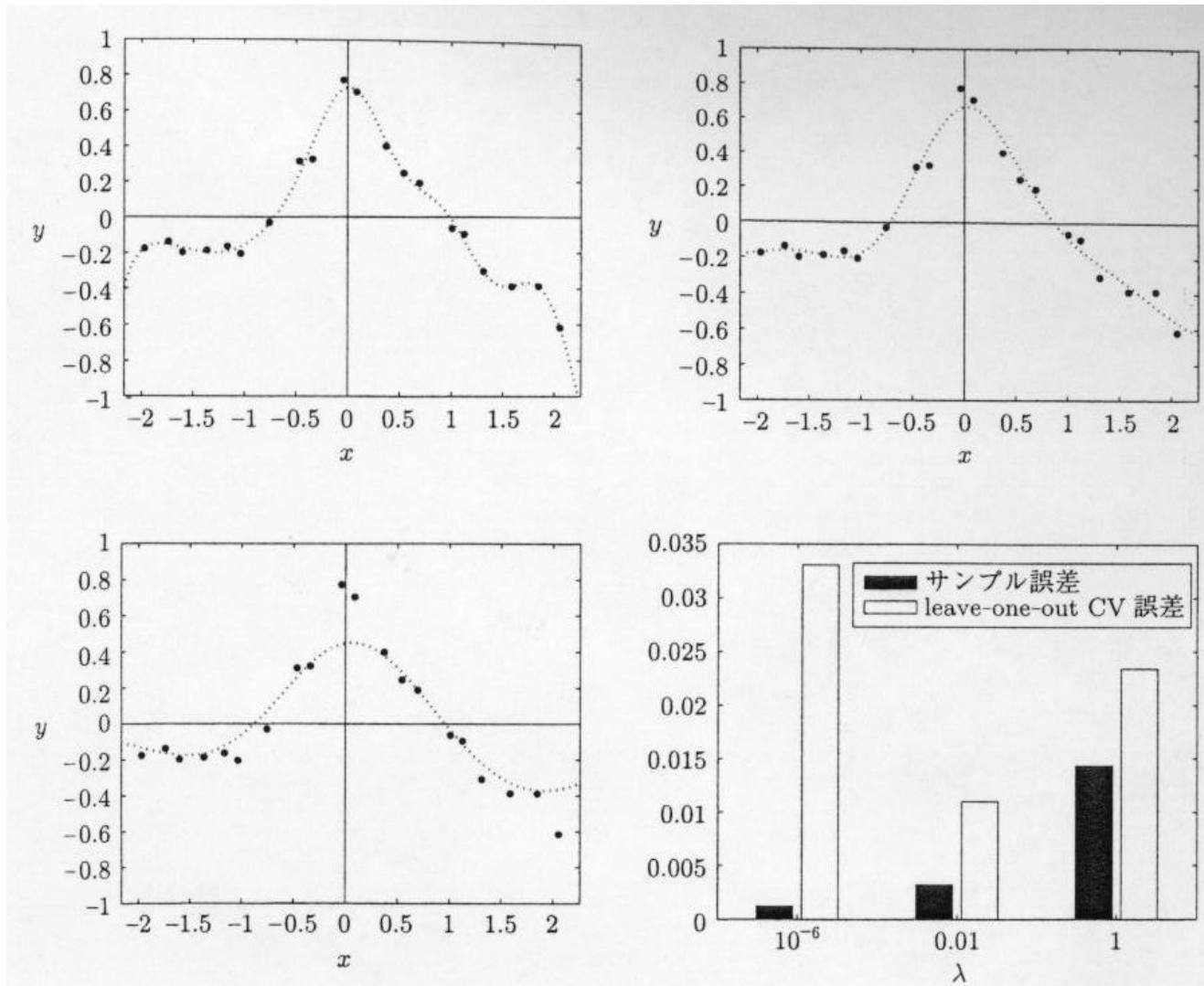


図 2.3 正則化パラメータ λ の値を変えた例. 左上・右上・左下の順に $\lambda=10^{-6}$, 0.01, 1. 右下は, それぞれのプロットのサンプル誤差と leave-one-out CV 誤差.

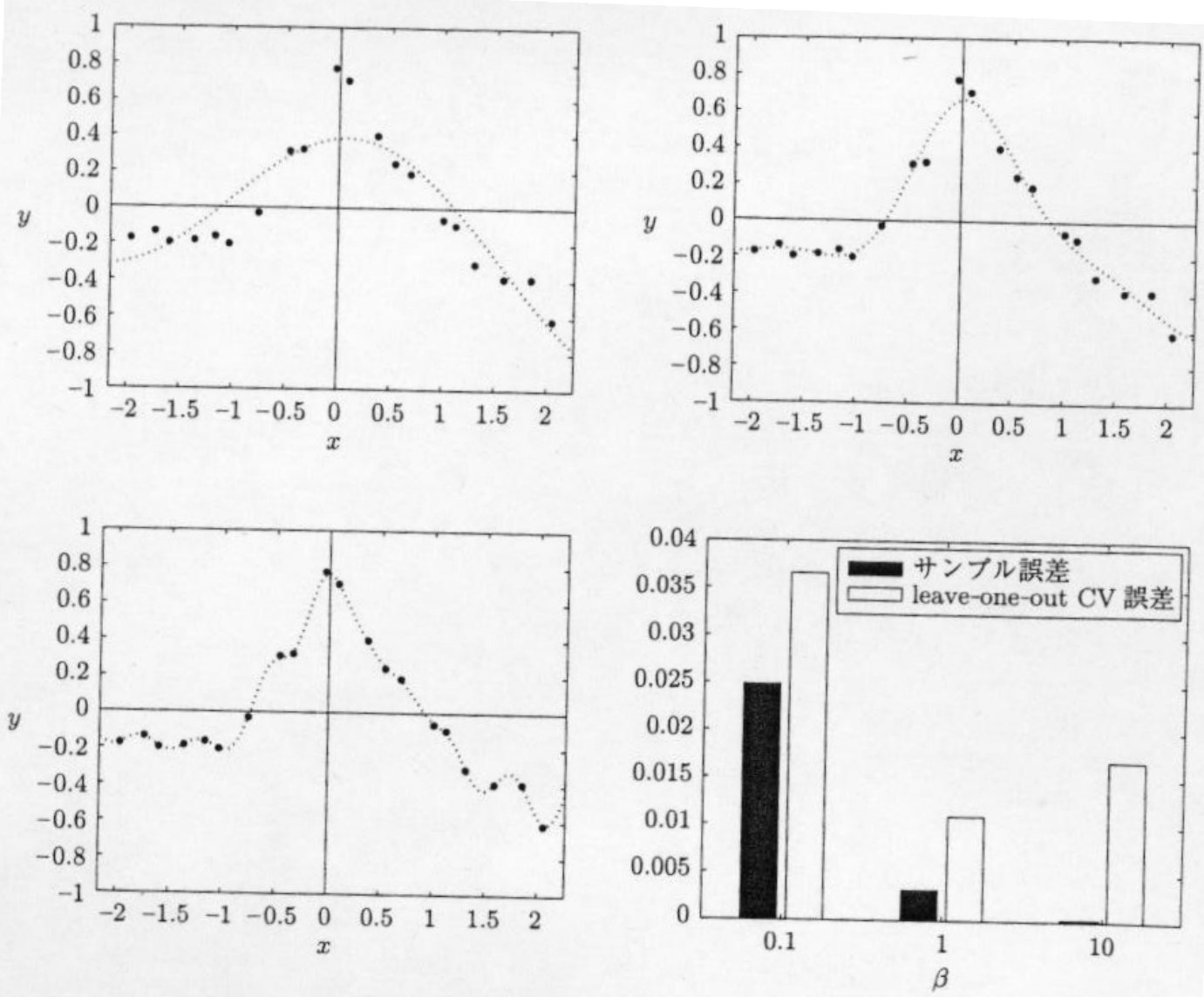


図 2.4 ガウスカーネルのパラメータ β の値を変えた例。左上・右上・左下の順に $\beta=0.1, 1, 10$ 。右下は、それぞれのプロットのサンプル誤差と leave-one-out CV 誤差。

つづき

- leave-one-out CV誤差は汎化能力を評価する一つの尺度と考えることができるので、モデル選択を行うことができる。
- いくつかの β と λ に対して leave-one-out CV誤差を計算し、最も小さな値となる β と λ を選ぶ