

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

2.3 確率モデルからの導入

2.4 汎化能力の評価とモデル選択

茨城大学工学部情報工学科
國井慎也

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すること
 - ただし、データには、ノイズがのっていたり、関数のうち有限個の点のデータしか関数値が与えられていなかったり、不確定な要素が多い
- このときの解決策として有効なのは、データの生成過程を確率分布を用いてモデル化すること

線形モデルによる関数近似の確率モデル

- ここで考える確率モデルは、出力 y は関数値 $f(x)=w^t x$ そのものではなく、ランダムノイズがのつた

$$y=f(x)+\varepsilon$$

- x, f が与えられたもとでの y の条件付き確率

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

ここで、 ε は独立な正規分布
この正規分布の分散を σ^2 とする

関数fの確率モデル

- 線形モデル $y=w^T x$ ではパラメータ w を決めることと f を決めることは等価である。
- W の各成分が独立に平均0で分散が $\frac{1}{\lambda}$ であるような正規分布から生成されたとすると、 w の分布は

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

生成モデル(1)

- サンプルは以下のような過程で生成されるものと仮定する。
 - $p(w)$ に従ってパラメータ w がランダムに決められる。
 - その w で決まる関数 $f(x)$ を n 個の点 $x^{(1)}, \dots, x^{(n)}$ で計算し、 $p(y|x;f)$ の分布に従ってノイズがのってサンプル出力 $y^{(1)}, y^{(2)}, \dots, y^{(n)}$ が観測される。

生成モデル(2)

- パラメータとデータの同時確率分布は

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

- このように、データやパラメータの生成過程を確率分布で表したものを生成モデルという

ベイズの定理(1)

- 生成モデルでは関数のパラメータが予め決まって、それからサンプルが決められた。
-しかし、サンプルから関数を推定したいときは、それを逆向きに行う必要がある。それをするのがベイズの定理である。

確率変数A,Bがあるとき

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

これをベイズの定理という

ベイズの定理(2)

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

- ベイズの定理は、AからBが生成される生成モデルがあるときに、逆にBを観測した上で、Aを推論する方法である。
- この確率はBを観測した事後におけるAの分布なので事後分布という

MAP推定

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

- $p(A|B)$ を最大にするようなAの値を見つけるにはMAP推定という手法を使う。
 - 右辺の分母はAの値に依存しないため、右辺の分子を最大にすることだけを考えればよい。
 - 線形モデルでは、パラメータ w がA、出力値 $y^{(1)}, y^{(2)}, \dots, y^{(n)}$ がBに相当する。

線形モデルのMAP推定(1)

- 線形モデルのベイズの公式の右辺の分子は式(2.31)と同じになる

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

- 式(2.31)の対数をとると、

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

線形モデルのMAP推定(2)

$$-\frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - f(x^{(i)}))^2 - \frac{\lambda}{2} \|w\|^2 + \text{定} \quad \square$$

- 第一項は二乗誤差にマイナスを付けたもので、第二項は正則化項に対応している。
- つまり、MAP推定は線形モデルでの正則化付きの二乗誤差を最小にすることと等価になっている。

正規過程

- パラメータ w ではなく、関数 $f(x)$ が正規分布に従うと考える。
- $f(x)$ が正規分布に従うとは、
任意の個数 n 個の入力値 $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ に対して、
 $f(x^{(1)}), \dots, f(x^{(n)})$ 平均0、分散共分散行列 $V = (V_{ij})_{i,j=0, \dots, n}$ をもつ
多次元正規分布 $N[0, V]$ に従う。
- $f(x^{(i)})$ は正規確率変数で、 $f(x^{(i)})$ と $f(x^{(j)})$ の共分散が

$$V_{ij} = E_{f(x^{(i)}), f(x^{(j)})} [f(x^{(i)}) f(x^{(j)})]$$

ここで $E[]$ は確率変数の期待値を表す。この確率過程を X 上の正規過程(GP=Gaussian Process)と呼ぶ

正規過程からカーネルへ

- V_{ij} は $x^{(i)}$ と $x^{(j)}$ で決まるので、関数 $k(x^{(i)}, x^{(j)})$ とみなすことができ、関数 k が正定値であって、カーネル関数とみなすことができる。
- カーネル関数 k は、そのグラム行列を分散共分散行列とみなすことによって、正規過程と等価になる。

正規過程のMAP推定

- サンプルが与えられたときの関数の推定をMAP推定で考える。
 - パラメータ w の分布を求めるのではなく、関数 $f(x)$ の分布を求める。
 - サンプル点の f に $f(x^{new})$ を加えた $n+1$ 個の関数値の事前分布は正規過程でモデル化できる。
- サンプルに対するグラム行列を $K=(K_{ij})_{i,j=0,\dots,n}$ とおくと f 同時確率分布は

$$p(f(x^{(1)}), f(x^{(2)}), \dots, f(x^{(n)}), f(x^{new}))$$

正規過程のMAP推定

$$p(f(x^{(1)}), f(x^{(2)}), \dots, f(x^{(n)}), f(x^{new}))$$

- 上の式は平均0の正規分布
- 分散共分散行列は、サンプル $x^{(i)}$ と $x^{(j)}$ に関する成分が K_{ij} で、サンプル $x^{(i)}$ と x^{new} に関する成分は $k(x^{(i)}, x^{new})$ で与えられる。
- $f(x)$ を固定したもとの y の分布は線形モデルと同じく式(2.29)となる

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

- 関数値とデータの同時分布は

$$p(f(y^{(1)}), \dots, f(y^{(n)}), f(x^{(1)}), f(x^{(2)}), \dots, f(x^{(n)}), f(x^{new}))$$
$$= p(f(x^{(1)}), f(x^{(2)}), \dots, f(x^{(n)}), f(x^{new})) \prod_{i=1}^n p(y^{(i)} | x^{(i)}; f)$$

- ここに出てくる分布はすべて正規分布

サンプルが与えられたもとでの $f(x^{(new)})$ の事後確率も正規分布になる。

$f(x^{new})$ の事後確率の平均

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

- ただし、与えられたサンプルをDと書いた
 $\{ \}_i$ はかっこの中のベクトルの第*i*成分を表す
- 事後分布が正規分布なら、MAP推定はその平均に一致する。

正規過程とカーネルの等価性

$$\sum_{i=1}^n \left\{ (K + \sigma^2 I_n)^{-1} y \right\}_i k(x^{(i)}, x^{new})$$

- 上の式は、関数近似の場合の式(1.17)の α を

$$f(x) = \sum \alpha_i k(x^{(i)}, x)$$

に代入したものである。ⁱ(ノイズの分散 σ^2 が正則化パラメータ λ に対応)

- つまり、カーネル関数の重みつき和のモデルを正則化付きで求めた関数近似の結果が、そのカーネル関数と等価な正規過程を事前分布として用いた場合のMAP推定と一致する

汎化能力の評価とモデル選択

- データ解析においては、汎化能力が高いあてはめを行うことが重要
 - カーネル法では、サンプル数と同じだけの自由度があるため、過学習を引き起こす。
 - モデルの形があまり複雑になりすぎないようにすることが必要
- 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をモデル選択という

クロスバリデーション

- 汎化能力とは、学習に使ったサンプル以外のデータに対する性能
 - サンプルを学習用とテスト用に分け、初めに学習用のサンプルで学習した後に、テスト用のサンプルで性能評価を行う。
 - 問題点として、テストデータを残しすぎると学習データが少なくて学習に十分な性能が出せず、一方、テストデータが少なすぎると、評価結果が不十分になる。

k-foldクロスバリデーション

- K-foldクロスバリデーション

[1]まず、サンプルをk個のグループに分ける。

[2] $i=1, \dots, k$ に対して以下を繰り返す。

(i)i番目のグループを除いたデータで学習を行う

(ii)i番目のグループでのテスト誤差を r_i とする

[3] $\sum_{i=1}^k \frac{r_i}{k}$ をテスト誤差の推定値(クロスバリデーション誤差)とする

線形モデルの leave-one-outクロスバリデーション(1)

- サンプルが n 個あるとき、 n -foldクロスバリデーションは1つだけのサンプルデータとしてのぞいておく方法。特にleave-one-outクロスバリデーションと呼ばれる。

- 関数近似で $y=f(x)$ を学習したとする。

このとき、入力 $x^{(i)}$ を学習した関数に入れると
 $\tilde{y}^{(i)} = f(x^{(i)})$ という結果が得られる。

これは、サンプル出力 $y^{(i)}$ のノイズを除去した推定値としてみなすことができる

線形モデルの leave-one-out クロスバリデーション (2)

- 線形回帰やカーネル回帰では、 $\tilde{y} = (\tilde{y}^{(1)}, \dots, \tilde{y}^{(1)})^T$ が $y = (y^{(1)}, \dots, y^{(1)})^T$ の線形変換でかける。
- 実際、カーネル回帰の場合、 $\tilde{y}^{(i)} = \sum_{i=1}^n \alpha_i k(x^{(i)}, x^{(j)})$
- と $\alpha = (K + \lambda I_n)^{-1} y$ より、

$$\tilde{y} = K \alpha = (K + \lambda I_n)^{-1} K y$$

- 一般に $\tilde{y} = H y$ という線形関係があるときを考える。カーネル回帰の場合は、 $H = (K + \lambda I_n)^{-1} K$ である。

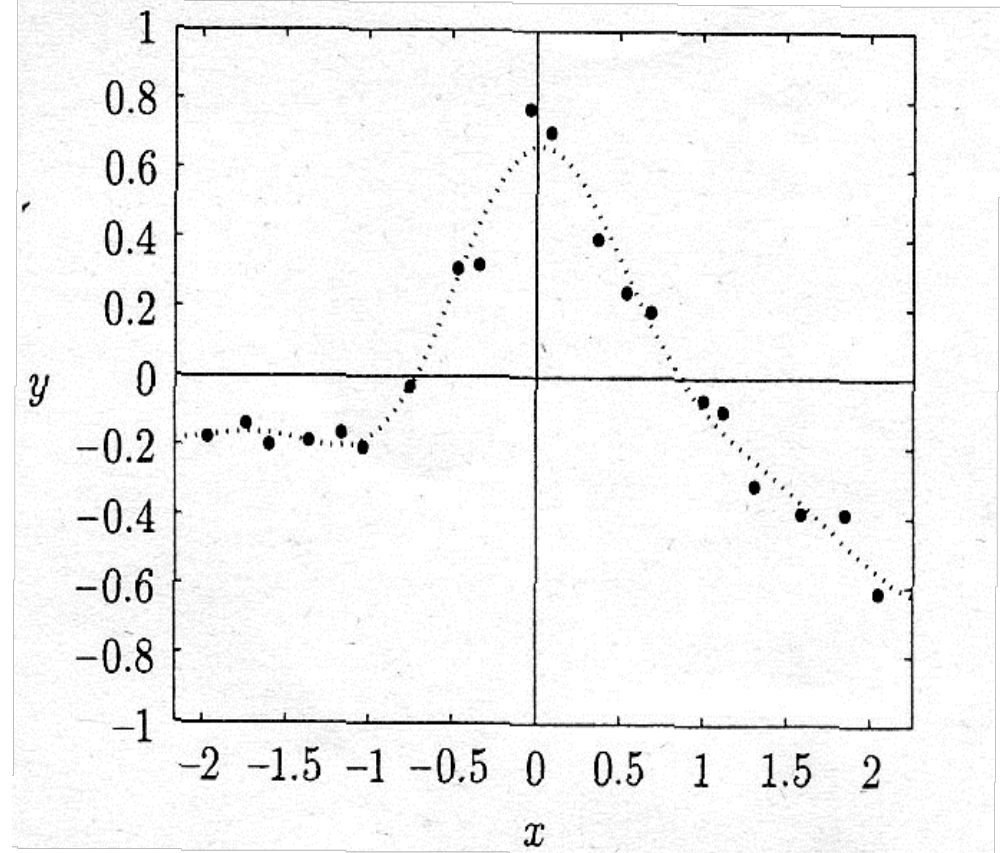
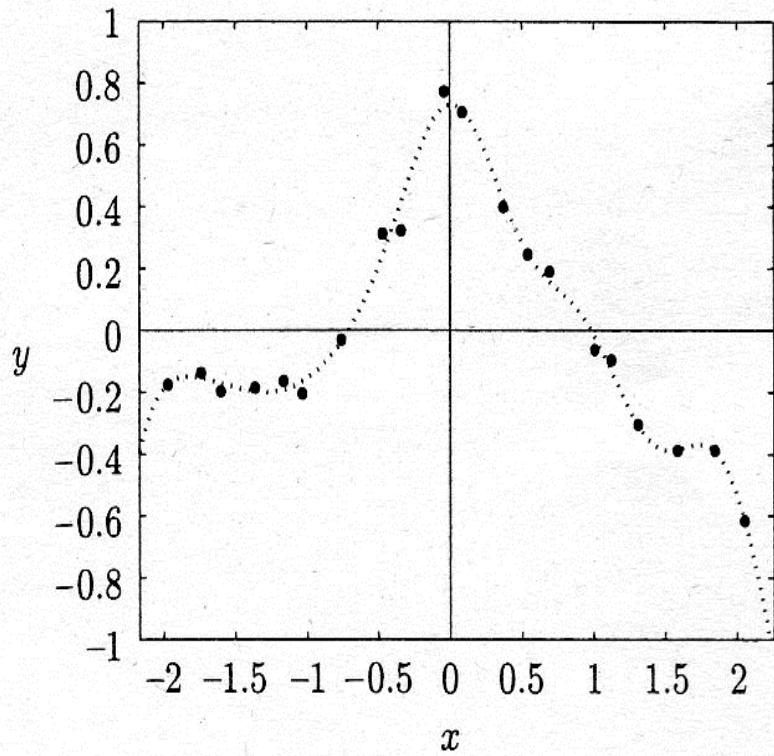
線形モデルの leave-one-outクロスバリデーション(3)

- このときのleave-one-outクロスバリデーション誤差(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、サンプル出力 $y^{(i)}$ とノイズ成分を除去した推定値 $\tilde{y}^{(i)}$ の重み付きの平均誤差として、

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

- ただし、 H_{ii} は H の第 i 対角成分である。

具体例(1)

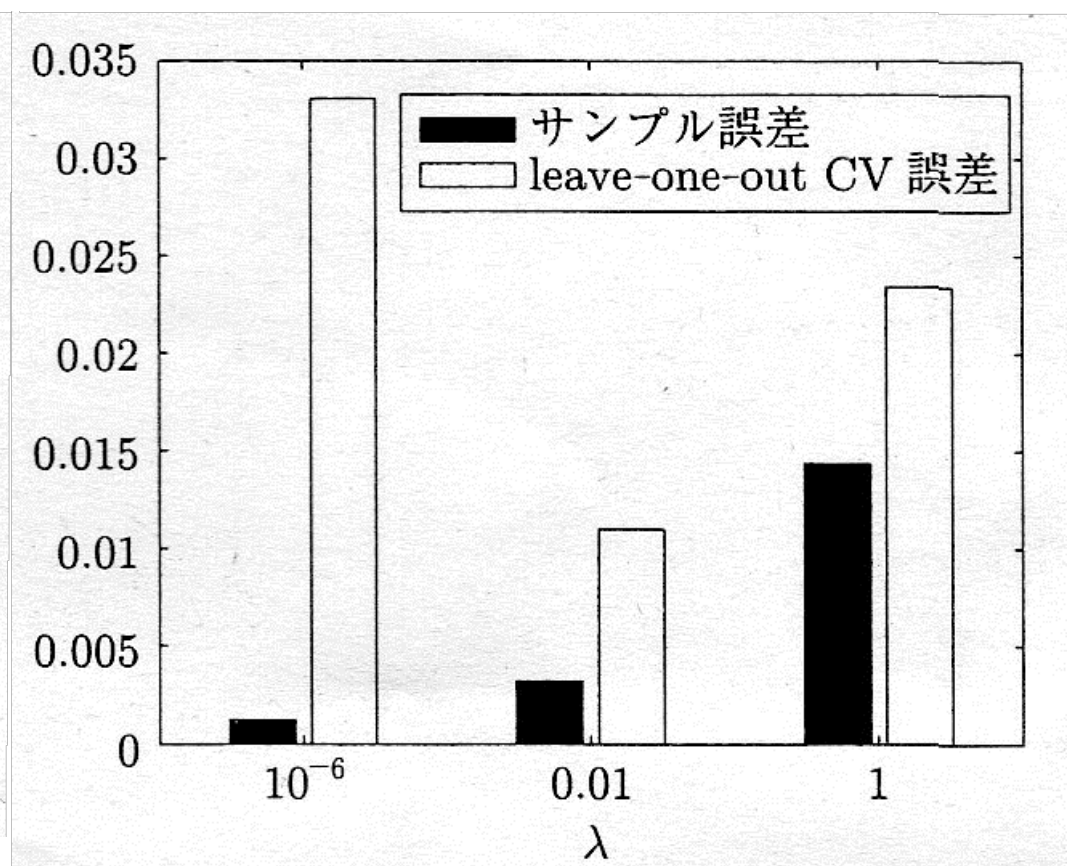
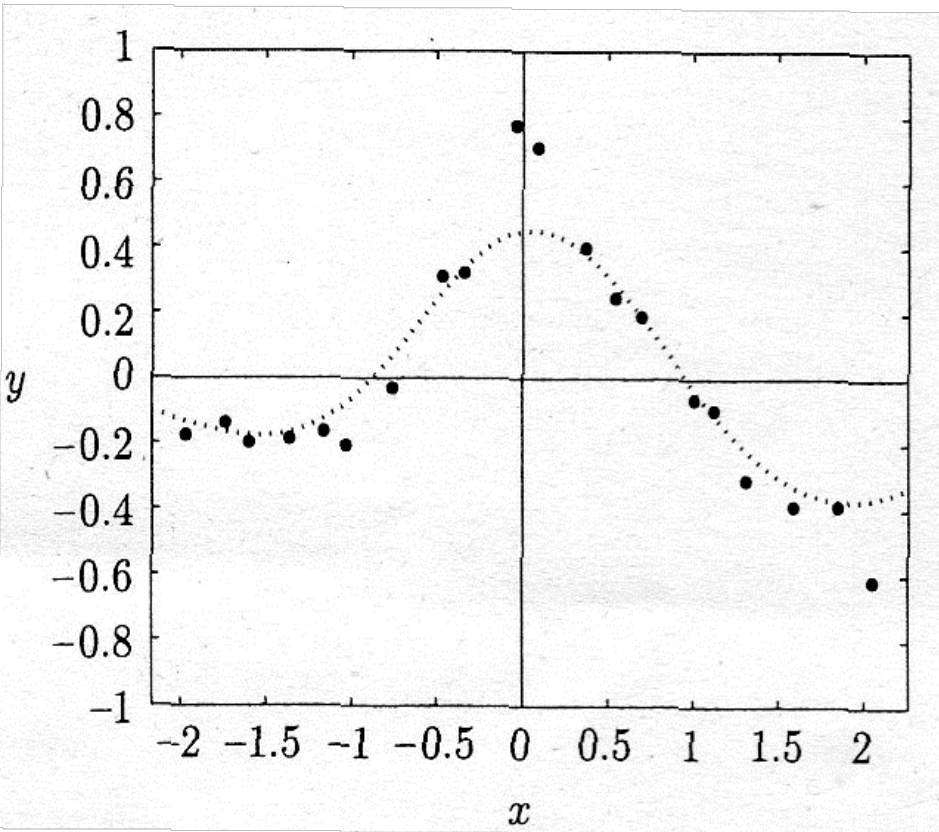


- $\beta = 1$ を固定

左図: $\lambda = 10^{-6}$

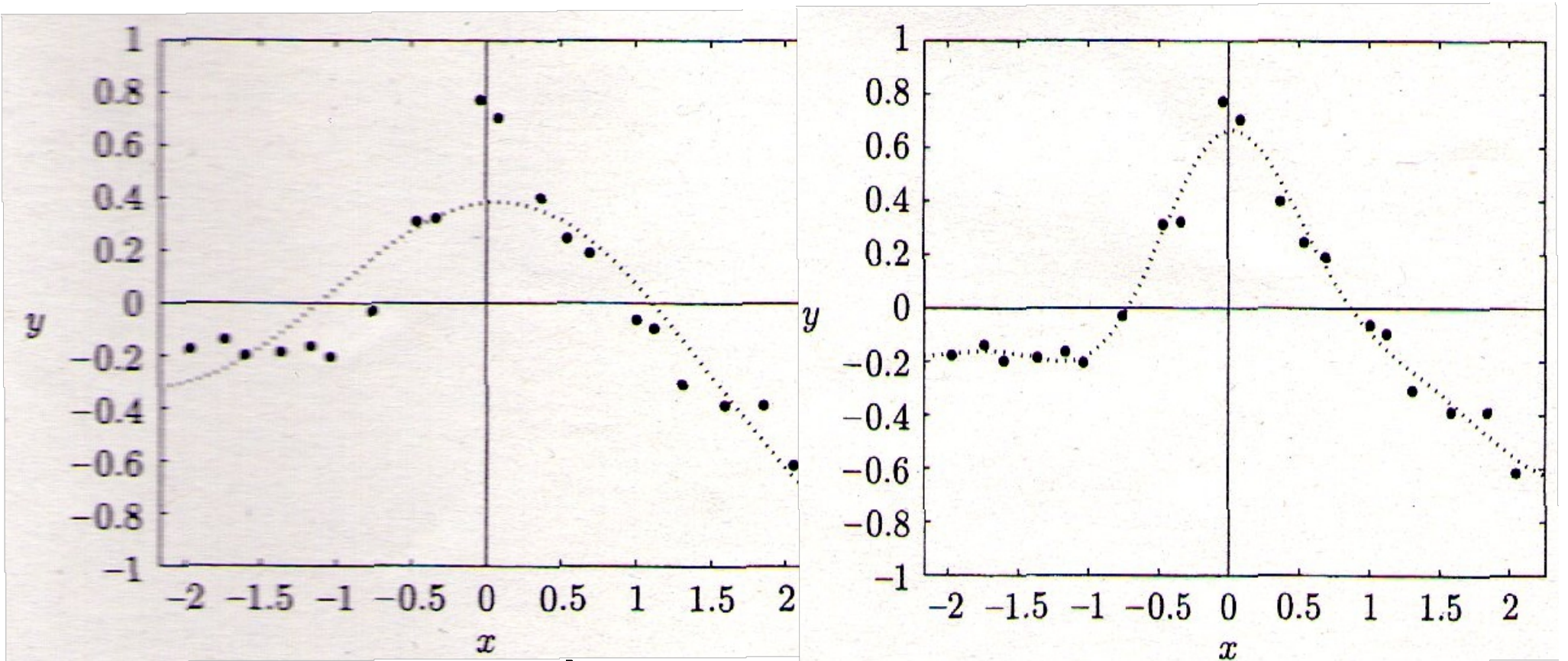
右図: $\lambda = 0.01$

具体例(2)



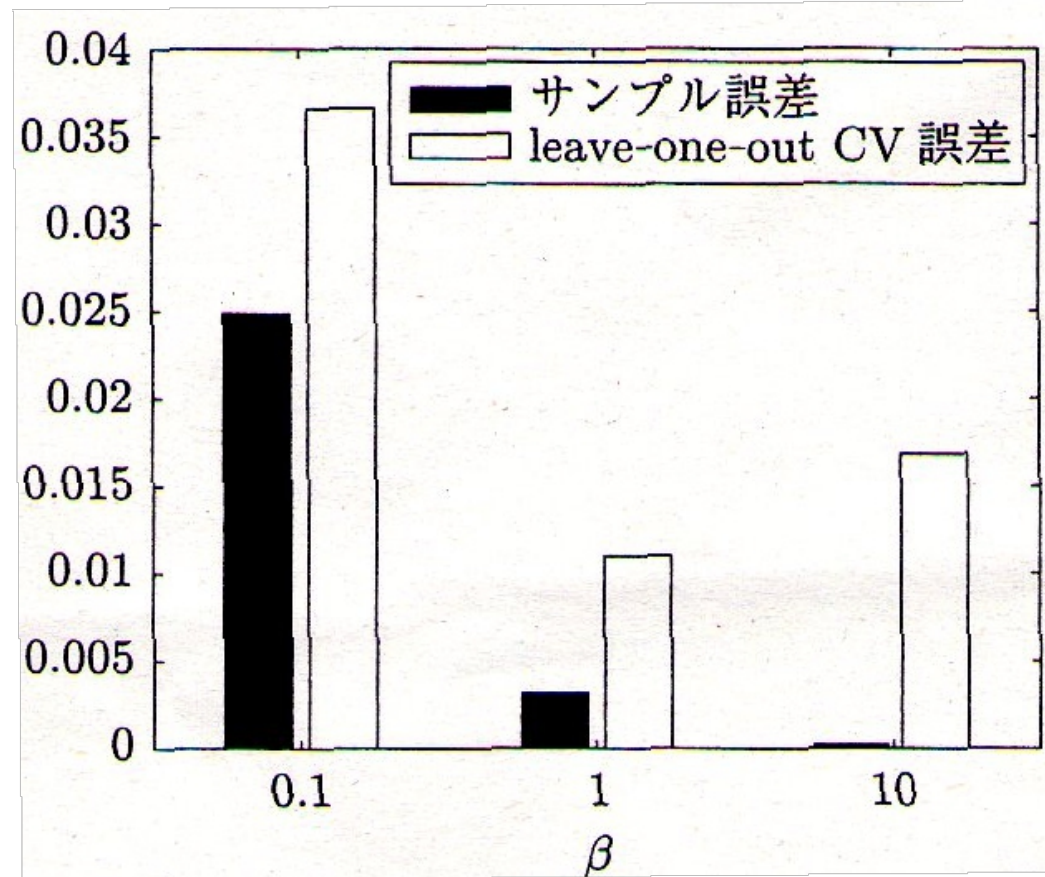
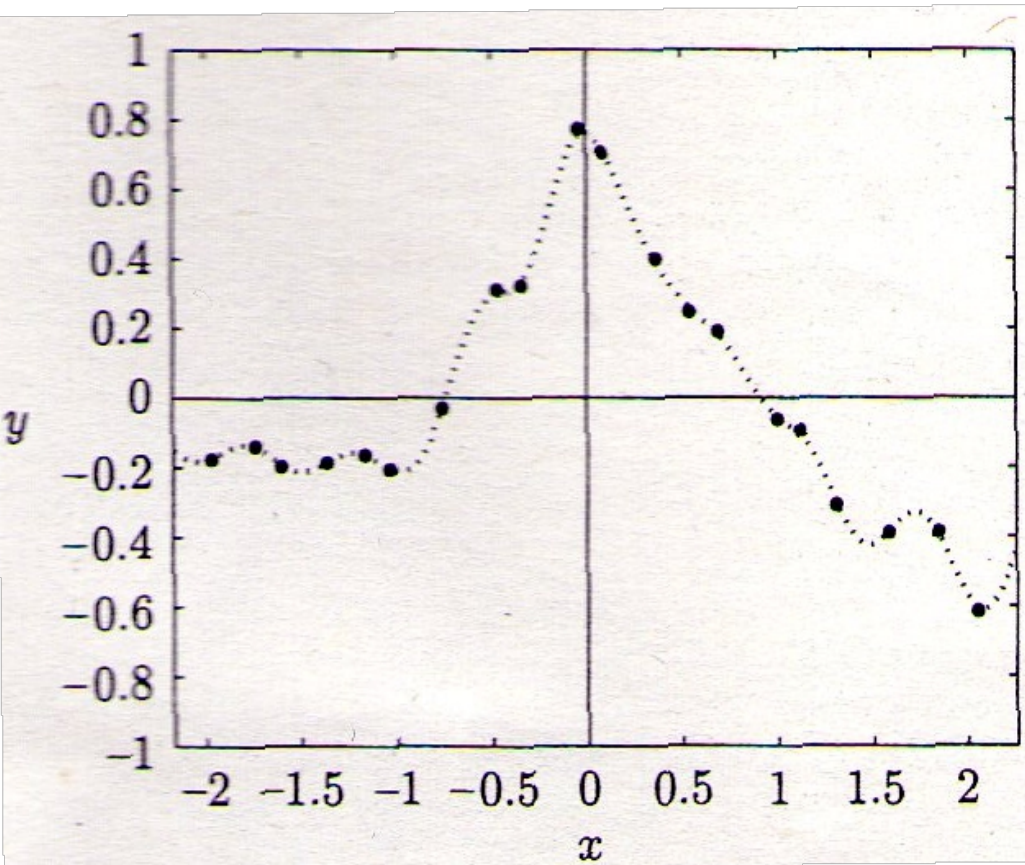
- 左図: $\beta=1, \lambda=1$
- 右図: サンプル誤差とCV誤差の比較

具体例(3)



- $\lambda = 0.01$ に固定、
- 左図: $\beta = 0.1$
- 右図: $\beta = 1$

具体例(4)



- 左図: $\lambda=0.01$, $\beta=10$
- 右図: サンプル誤差とCV誤差の比較

β と λ の選択

- β を増やすことと、 λ を減らすことが同じような振る舞いを示す。
- 結果的に、いくつかの β と λ に対してleave-one-out CV誤差を計算し、そのうち最も小さな値をとるような β と λ を選べばよい。