

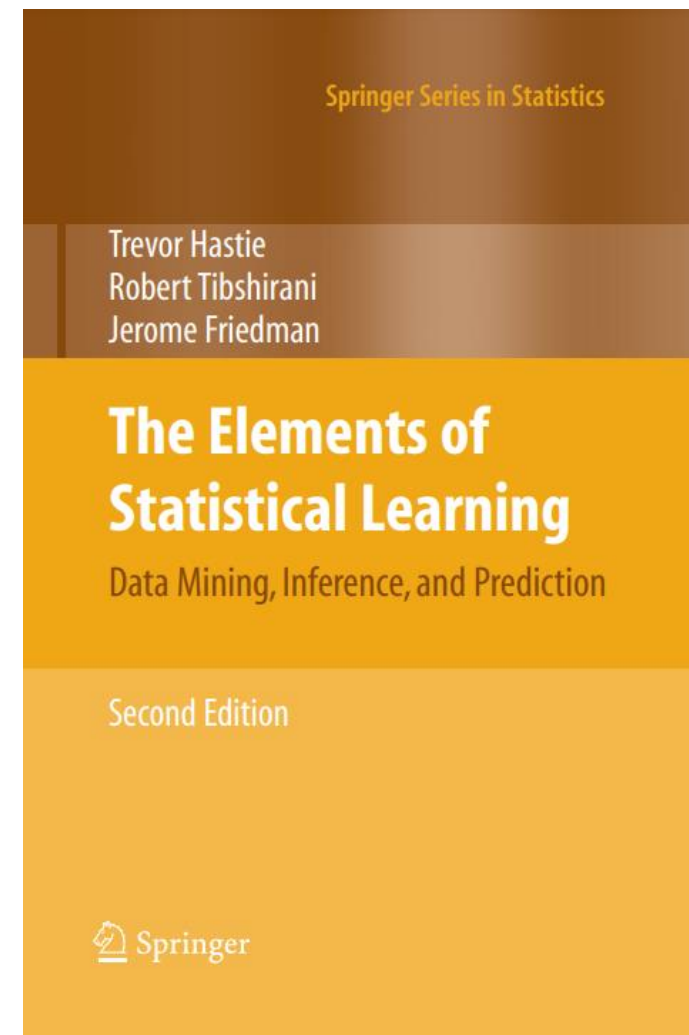


# THE ELEMENTS OF STATICAL LEARNING



## 第二章 教師あり学習の概要

- 最小 2 乗法
- 最近傍法
- 高次元での局所的手法



## ■ 線形モデル

線形モデルでは、入力ベクトル  $X^T = (X_1, X_2, \dots, X_p)$  が与えられたとき出力  $Y$  をこのように予測する

$$\hat{Y} = \underbrace{\hat{\beta}_0}_{\text{切片}} + \sum_{j=1}^p X_j \hat{\beta}_j.$$

$\hat{\beta}_0$  を  $\hat{\beta}$  に加えたら  
表記がカンタンになる。

$$\hat{Y} = X^T \hat{\beta},$$

線形モデルにどのように訓練データを当てはめればよいだろうか？  
⇒ 頻繁に使われるのは **最小 2 乗法**

## ■ 最小 2 乗法

$$\text{残差} = \text{実際の観測値} - \text{予測値}$$

残差 2 乗和 (RSS) : 予測モデルの精度を評価し、モデルの良さを評価する指標

$$\text{RSS}(\beta) = \sum_{i=1}^N (y_i - x_i^T \beta)^2.$$

最小 2 乗法 : 残差 2 乗和(RSS)を最小化するパラメータを見つけるための手法

## ・ 最小 2 乗法

RSS は  $\beta$  の 2 次関数

$$\text{RSS}(\beta) = \sum_{i=1}^N (y_i - x_i^T \beta)^2.$$

行列表記に

$$\text{RSS}(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$$



$$\frac{\partial \text{RSS}}{\partial \beta}$$

$$= 2(\mathbf{X}^T \mathbf{X})\beta - 2\mathbf{X}^T \mathbf{y}$$

$\beta$  で偏微分

$$2(\mathbf{X}^T \mathbf{X})\beta - 2\mathbf{X}^T \mathbf{y} = 0$$

最小なとき  
傾きが 0

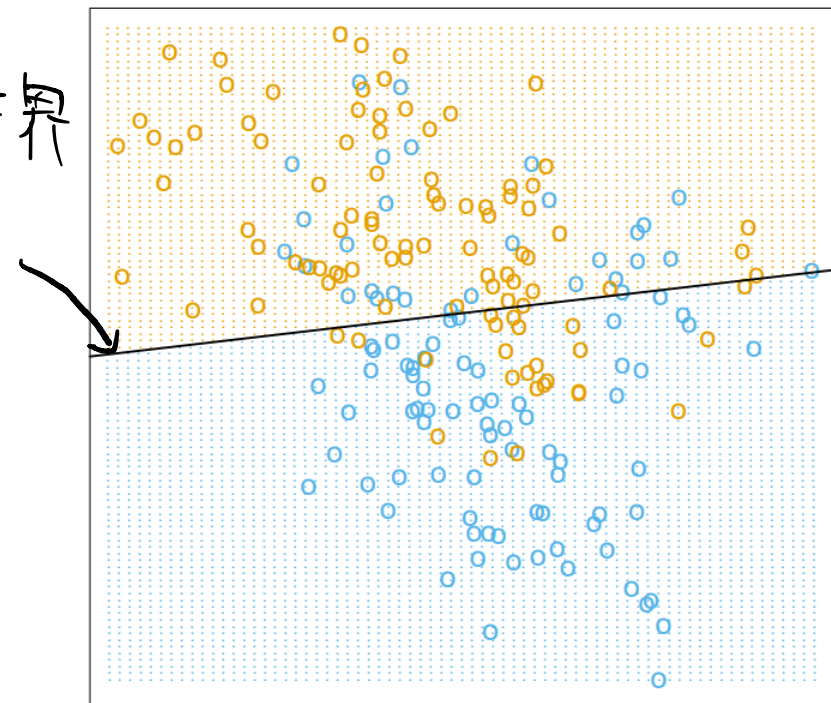
$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

## ■ 分類問題の線形モデル

The line is the decision boundary defined by  $x^T \bar{\beta} = 0.5$ .

応答変数 $Y$ を青色が0, オレンジ色が1として  
線形モデルを当てはめると...

決定境界



決定境界の両側に誤分類された点が存在している

⇒線形モデルの柔軟性の欠如? or 避けられないノイズ?

$$\hat{G} = \begin{cases} \text{ORANGE} & \text{if } \hat{Y} > 0.5, \\ \text{BLUE} & \text{if } \hat{Y} \leq 0.5. \end{cases}$$

## ■ k最近傍法

$\hat{Y}$ を予測する際に，入力 $x$ に最も近い訓練データ集合の観測値を利用する

⇒ $k$ 個の近くの点が 青色 or オレンジ色 で多数決を取る

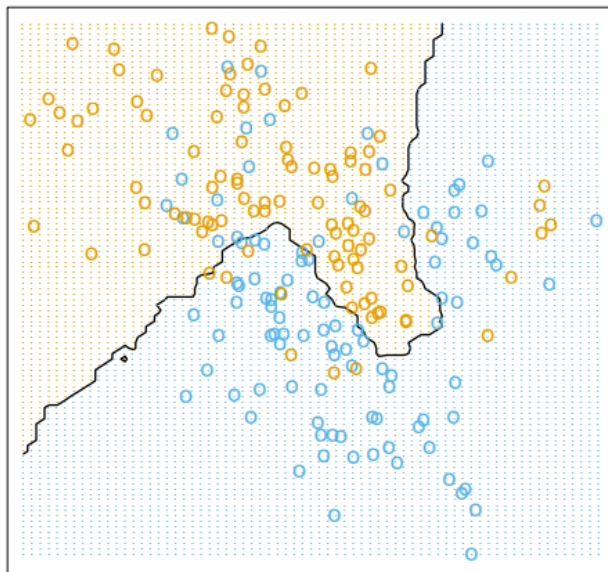
$$\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i,$$



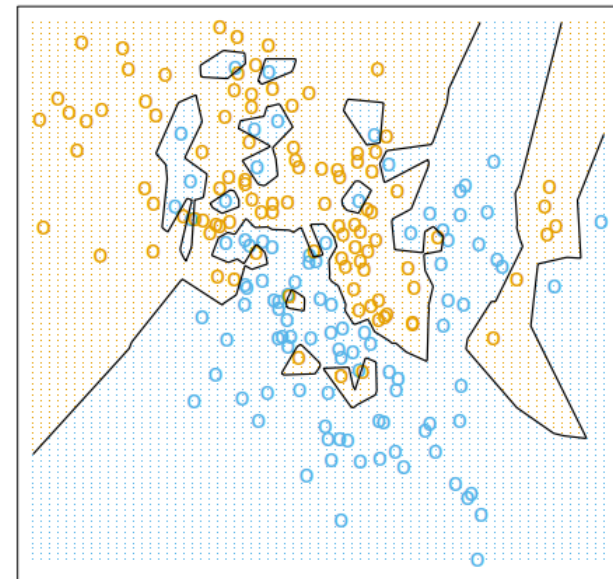
$N_k(x)$

⇒ 訓練データのうち  $x$  に近い  $k$  個の点

15-Nearest Neighbor Classifier



1-Nearest Neighbor Classifier



## ■ 損失関数

入力 $X$ が与えられたとき、 $Y$ を予測する関数 $f(X)$ を見つけるための予測に罰則を設ける

⇒予測に対する罰則を定義するための**損失関数**を導入

**損失関数**のうち最も頻繁に利用される、2乗誤差損失(squared error loss)は

$$L(Y, f(X)) = (Y - f(X))^2 \quad \text{で定義され,}$$

$f$ を選ぶ基準として、期待値予測誤差(expected prediction error)は

$$\text{EPE}(f) = E(Y - f(X))^2 \quad \text{で定義される}$$

⇒EPEを最小化することで良い予測を構築することが目標



# 期待予測誤差

$$E[X] = \int_{-\infty}^{\infty} x p(x) dx$$

$$= \int x p(dx)$$

$$\begin{aligned} \text{EPE}(f) &= E(Y - f(X))^2 \\ &= \int [y - f(x)]^2 \text{Pr}(dx, dy), \end{aligned}$$

同時分布  $\text{Pr}(dx, dy)$  を  $X$  条件づけて

$$\text{EPE}(f) = E_X E_{Y|X} ([Y - f(X)]^2 | X)$$

EPE を最小化するにはそれぞれの  $x$  について最小化を行う

$$f(x) = \text{argmin}_c E_{Y|X} ([Y - c]^2 | X = x)$$

$$f(x) = E(Y | X = x),$$

最小化するのとき  
条件つき期待値

$$\begin{aligned} &E[(Y - c)^2] \\ &= E[Y^2 - 2cY + c^2] \\ &= E[Y^2] - 2cE[Y] + c^2 \end{aligned}$$

## ■ ベイズ分類器

条件付確率を用いて最も確からしいクラスに分類する  
⇒理論的に最適な決定境界を示す

出力がカテゴリ型変数 $G$ である場合に  
損失関数と期待予測誤差を考えると...

$$\text{EPE} = E[L(G, \hat{G}(X))],$$

同時分布  $\Pr(X, Y)$  を  $X$  で条件づけ

$$\text{EPE} = E_X \sum_{k=1}^K L(\mathcal{G}_k, \hat{G}(X)) \Pr(\mathcal{G}_k | X)$$

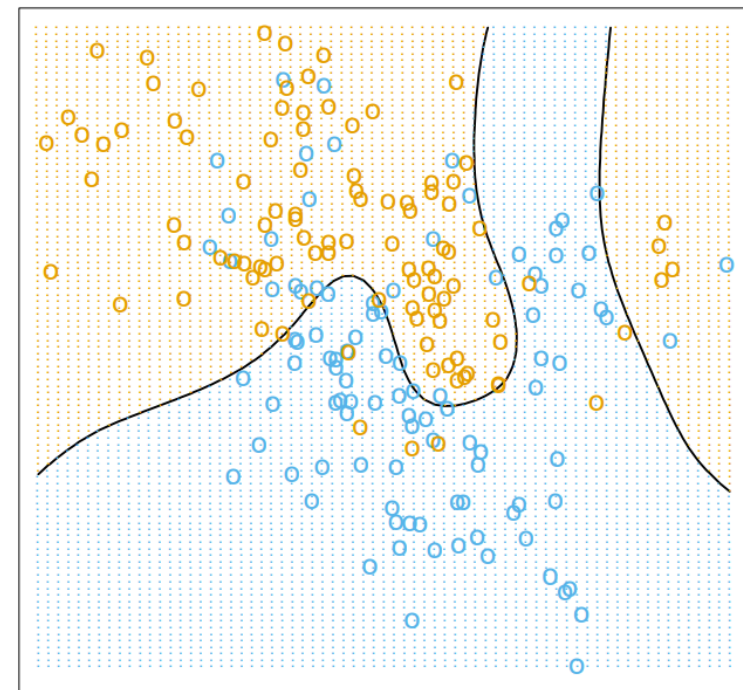
$$\hat{G}(x) = \operatorname{argmin}_{g \in \mathcal{G}} \sum_{k=1}^K L(\mathcal{G}_k, g) \Pr(\mathcal{G}_k | X = x).$$

または

$$\hat{G}(x) = \operatorname{argmin}_{g \in \mathcal{G}} [1 - \Pr(g | X = x)]$$

と表わすことができる

最小化したい



$$E[G] = C$$

## ■ 高次元での局所的手法

訓練データが多ければ多いほど,  $k$ 最近傍法を使えばより多くの近傍データを利用し最適な予測ができる...?

⇒高次元では上手く近似できなくなってしまう!

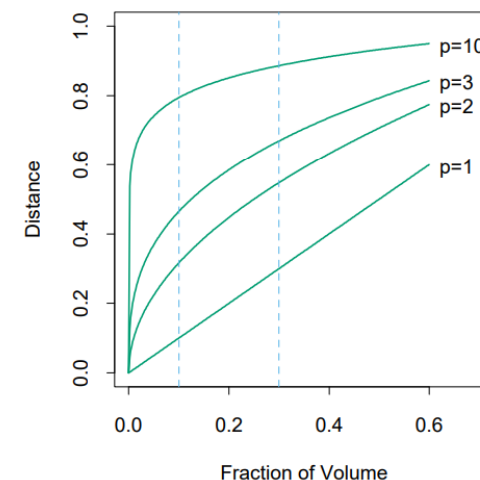
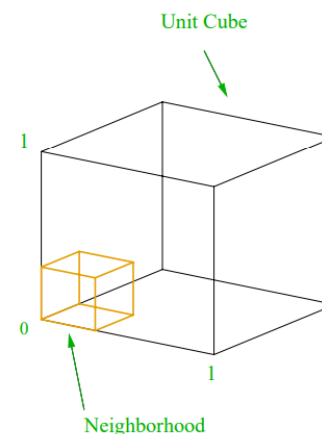
⇒この現象を次元の呪いとよぶ

例) 右の  $p$  次元単位立方体

近傍のデータ数が全データの割合  $r$  であるとする  
立方体の1辺の長さの期待値は  $e_p(r) = r^{1/p}$  で表される

$p = 10$  の場合,  $e_{10}(0.01) = 0.63$  となる

⇒ 全データの 0.01% を近傍に  
するためには 入力変数の 63%

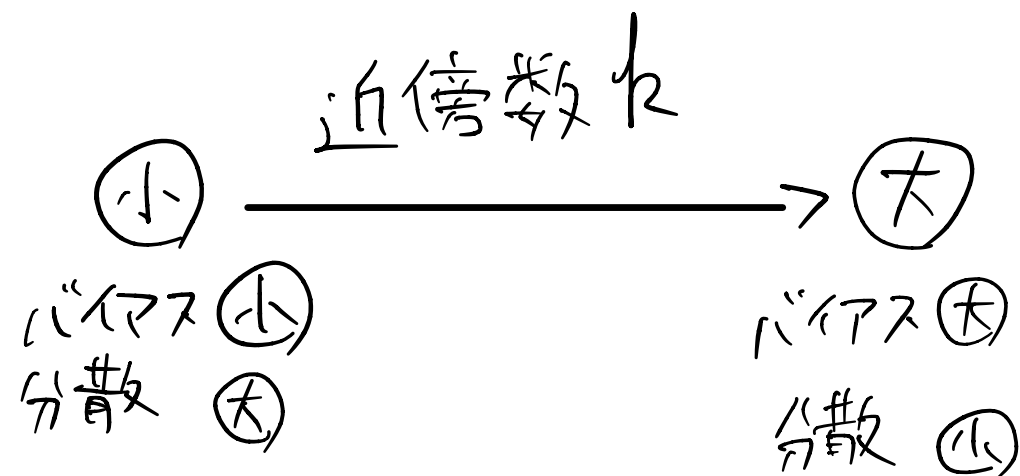


## ・ バイアスと分散のトレードオフ

バイアス：極めて複雑な実際の事象を単純なモデルで近似したために生じる誤差  
分散：異なる訓練データを使ったときにどの程度 $\hat{f}$ が変化するかを表す量

$Y = f(X) + \varepsilon$  からデータが生成されるとして,

$$\begin{cases} E(\varepsilon) = 0 \\ \text{Var}(\varepsilon) = \sigma^2 \end{cases}$$



$$\begin{aligned} \text{EPE}_k(x_0) &= E[(Y - \hat{f}_k(x_0))^2 | X = x_0] \\ &= \sigma^2 + [\text{Bias}^2(\hat{f}_k(x_0)) + \text{Var}_{\mathcal{T}}(\hat{f}_k(x_0))] \\ &= \underbrace{\sigma^2}_{\text{削減できない}} + \underbrace{\left[ f(x_0) - \frac{1}{k} \sum_{\ell=1}^k f(x_{(\ell)}) \right]^2 + \frac{\sigma^2}{k}}_{\text{削減可能}} \end{aligned}$$

→ 新しいテスト点による誤差

削減可能

## ・ バイアスと分散のトレードオフ

訓練誤差はモデルの複雑度を増やすほど（データに強く適合させるほど）減少する  
⇒過度に適合させると訓練データに特化したモデルになり  
大きな分散をもち、汎用性が悪化

逆にモデルが十分の複雑度をもっていない  
⇒過小適合となって大きなバイアスをもち、汎用性が悪化

