

GAMとその周辺

○田中 祐輔* 伊庭 克拓** 辻谷 将明***

*大阪電気通信大学大学院 工学研究科 情報工学専攻

**東京CRO株式会社 DM統計本部 統計解析部

***大阪電気通信大学 情報通信工学部 情報工学科

GAM and Its Related Problems

Yusuke Tanaka* Katsuhiko Iba** Masaaki Tsujitani***

*Osaka Electro-Communication University, Graduate School of Engineering

**TOKYO CRO, Inc.

***Osaka Electro-Communication University

発表内容

- はじめに
- 一般化加法モデル(GAM)
 - 平滑化スプライン
 - 薄板平滑化スプライン
- 適用例
 - 脊柱後湾症データ
 - 糖尿病網膜症データ
- シミュレーション実験
- まとめ

はじめに

- 非線形性をもつ統計的多変量解析の発展
 - 薄板平滑化スプライン ⇒ PROC TPSPLINE
 - 局所回帰 ⇒ PROC LOESS
 - 一般化加法モデル ⇒ PROC GAM
(Generalized Additive Models : GAM)



平滑化スプラインによるロジスティック判別の紹介と適用

2

ロジット変換

- 2値応答 : Y

$$\Pr(Y = 1) = \pi, \Pr(Y = 0) = 1 - \pi$$

ロジット変換 $(-\infty, +\infty)$

$$\ln\left(\frac{\pi}{1-\pi}\right) = c_0 + c_1x_1 + \cdots + c_kx_k$$

$Y \square Bin(1, \pi)$: ベルヌーイ分布

$$E[Y] = \pi, V[Y] = \pi(1-\pi)$$

3

ロジスティック判別

- ロジット逆変換

$$\pi = \frac{1}{1 + \exp\left\{-\left(c_0 + c_1x_1 + \cdots + c_kx_k\right)\right\}}$$

- 判別方法

$$\text{2値応答 } y = \begin{cases} 1: \text{第1群} \\ 0: \text{第2群} \end{cases}$$

観測値が**第1群**に属する($y = 1$)確率: π

観測値が**第2群**に属する($y = 0$)確率: $1 - \pi$

4

ロジスティック判別

- ロジット逆変換

$$\pi = \frac{1}{1 + \exp\left\{-\left(c_0 + c_1x_1 + \cdots + c_kx_k\right)\right\}}$$

- 判別方法

$$\pi \geq 0.5 \Rightarrow \text{第1群 } (y = 1)$$

$$\pi < 0.5 \Rightarrow \text{第2群 } (y = 0)$$

5

一般化加法モデル

(Generalized Additive Models : GAM)

- 応答変数に指数分布族を仮定
- 平滑化関数 $s(x)$ の加法モデル

GAM :

$$\text{logit}(\pi) = \ln\left(\frac{\pi}{1-\pi}\right) = \alpha + s(x_1) + \dots + s(x_k)$$

ロジット変換

平滑化スプライン

6

ペナルティ付き残差平方和

平滑化パラメータ

$$\sum_{i=1}^n [y_i - s(x_i)]^2 + \lambda \int \{s''(x)\}^2 dx$$

小さいほどモデルの当てはまりは良い

小さいほど滑らかな曲線 (曲げ弾性エネルギー)

➡ 最小にするスプライン関数

平滑化スプライン

$$y = s(x) = c_0 + c_1x + \frac{1}{12} \sum_{d=1}^n \theta_d |x - x_d|^3$$

7

自由度と平滑化パラメータ

- ハット行列

$$\text{応答 } \mathbf{y} \text{ の予測値: } \hat{\mathbf{y}} = \mathbf{H}_\lambda \mathbf{y}$$

- モデルの自由度

- 実効自由度(=有効パラメータ数)

$$df = \text{tr}(\mathbf{H}_\lambda)$$

- 平滑化パラメータ λ の決定

↔ 自由度 df の決定

8

クロス・バリデーション(CV)

- 1例消去(leaving-one-out)CV法

初期標本: $\mathbf{X} = \{\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(n)}\}$

↓ d 番目を除去: $\mathbf{X}^{(d)} = \{x_1^{(d)}, \dots, x_I^{(d)}; y^{(d)}\}$

$$\mathbf{X}_{[d]} = \{\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(d-1)}, \mathbf{X}^{(d+1)}, \dots, \mathbf{X}^{(n)}\}$$

$$CV = \frac{1}{n} \sum_{d=1}^n \left\{ y^{(d)} - \hat{y}_{[d]}^{(d)} \right\}$$

$\mathbf{X}_{[d]}$ で構築したモデルの、 $\mathbf{X}^{(d)}$ の予測値

➡ 最小にする平滑化パラメータ λ が最適。

クロス・バリデーション(CV)

● 1例消去(leaving-one-out)CV法

初期標本: $X = \{X^{(1)}, X^{(2)}, \dots, X^{(n)}\}$

↓ d 番目を除去: $X^{(d)} = \{x_1^{(d)}, \dots, x_l^{(d)}; y^{(d)}\}$

$X_{[d]} = \{X^{(1)}, X^{(2)}, \dots, X^{(d-1)}, X^{(d+1)}, \dots, X^{(n)}\}$

$$CV = -2 \sum_{d=1}^n \left\{ y^{(d)} \ln \hat{\pi}_{[d]}^{(d)} + (1 - y^{(d)}) \ln (1 - \hat{\pi}_{[d]}^{(d)}) \right\}$$

$X_{[d]}$ で構築したモデルの、 $X^{(d)}$ の予測値

➡ 最小にする平滑化パラメータ λ が最適 10

適用例1

● 脊柱後湾症データ

応答変数

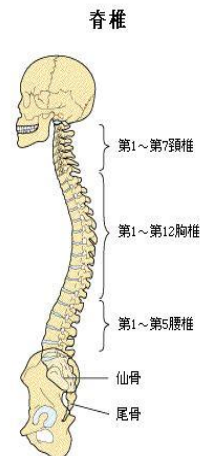
y : 術後の症状の有無

予測変数 $\begin{cases} 1: \text{症状が残る} \\ 0: \text{症状なし} \end{cases}$

x_1 : age 手術時の月齢

x_2 : start 何番目の脊椎から
先を手術したか

x_3 : num 手術した脊椎の個数



PROC GAMの文法

```
proc gam data=kyphosis( where=(id^=&d) ) ;
model kyp = param( start num ) spline( age , df=&df )
①          ②          / link = logit dist = binomial ;
run;
```

- ① model文の左辺は応答変数, 右辺はモデル式
- ② 線形項 `param(variable ...)`
 LOESS `loess(variable ,df=number)`
 平滑化スプライン `spline(variable ,df=number)`
 薄板平滑化スプライン `spline2(variable1,variable2 ,df=number)`

12

PROC GAMの文法

```
proc gam data=kyphosis( where=(id^=&d) ) ; ③
model kyp = param( start num ) spline( age , df=&df )
①          ②          / link = logit dist = binomial ;
run; ④
```

- ③ df値を指定しないとき、デフォルトdf=4
- ④ ③の場合, /の後に`method=gcv`を指定するとGCV規
準に基づいて最適な平滑化パラメータの探索

13

自由度の最適選択

- PROC GAMにおける最適自由度の決定
 - グリッド検索による平滑化パラメータの探索
 - method=gcvを指定した場合、モデルが収束しないことが起こり得る
- 提案法
主効果スプラインモデル：

$$\ln\left(\frac{\pi}{1-\pi}\right) = s(\text{age}) + \text{start} + \text{num}$$

- クロス・バリデーション(CV)規準に基づいて、各項の最適な自由度を探索

14

自由度の最適選択

- 最適選択の結果

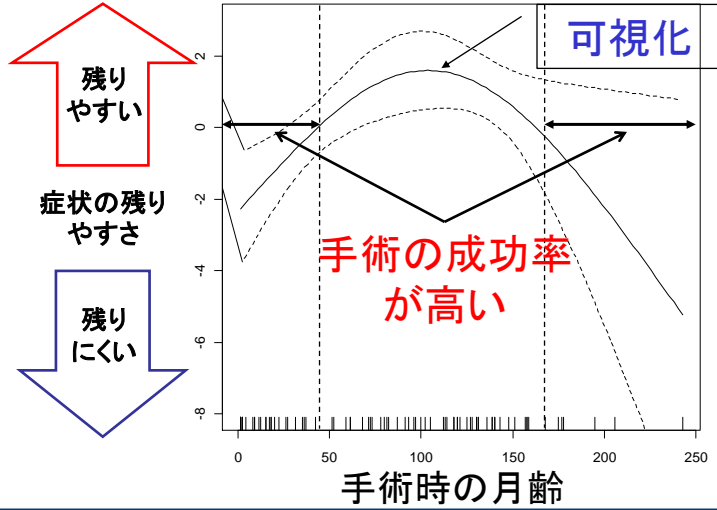
	CV	逸脱度	df(age)	誤判別率
1例消去CV法	65.92	56.39	2	0.133

- 他の判別手法との比較

$$\text{誤判別率} = \begin{cases} 0.133 : \text{GAM} \\ 0.205 : \text{ロジスティック判別} \\ 0.193 : \text{線形判別} \end{cases}$$

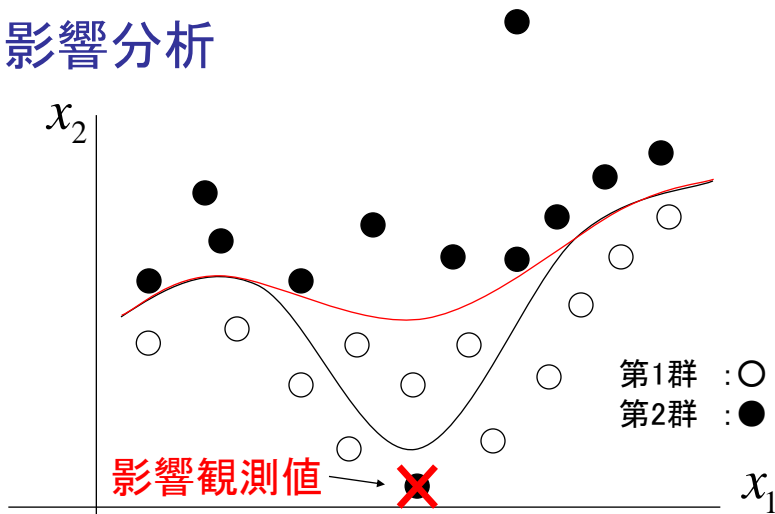
15

平滑化スプラインの予測値



16

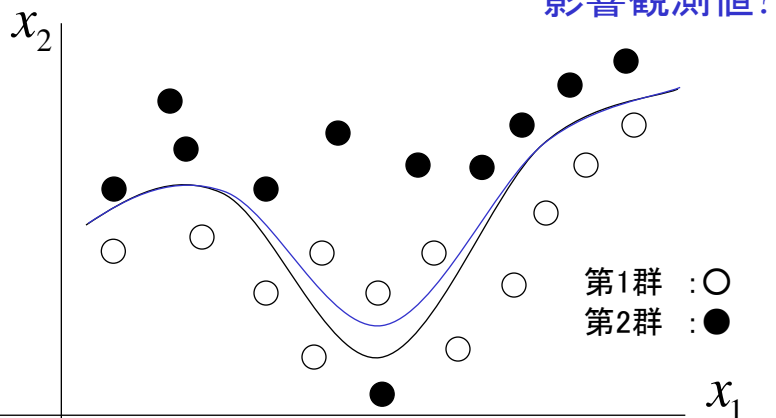
影響分析



曲線は簡単になり、誤判別は少なくなる

17

影響分析



曲線はあまり変わらず、誤判別も変わらない

DIFDEV

$$\Delta Dev_{[d]} = Dev - Dev_{[d]} \geq 0$$

Dev: すべての個体を用いたときの逸脱度

*Dev*_[d]: *d*番目の個体を取り除いたときの逸脱度

- 検定の方法

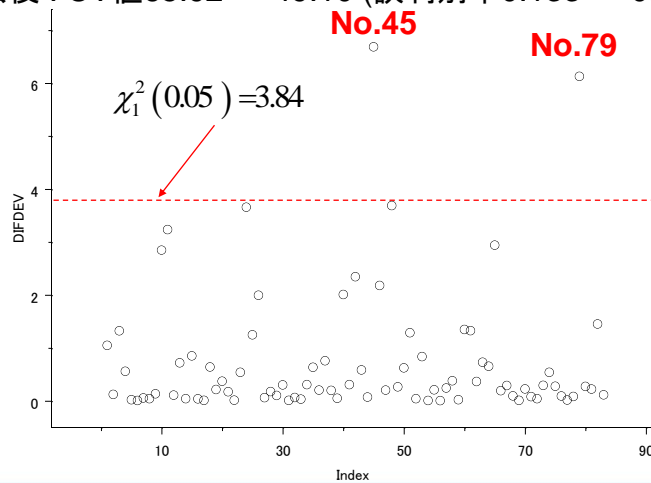
自由度1の χ^2 分布に基づき、有意水準0.05で検定

$$\Delta Dev_{[d]} \geq \chi^2(0.05) = 3.84$$

ならば、*d*番目の個体を除去

DIFDEV

除去後 : CV値65.92 ⇒ 49.10 (誤判別率0.133 ⇒ 0.111)



20

適用例2

● 糖尿病網膜症

- 網膜内の血管障害に伴う病気
- 669例の糖尿病患者における網膜症の進行の有無を調査

y : ret 糖尿病網膜症の進行 $\left\{ \begin{array}{l} 1 : \text{進行あり} \\ 0 : \text{進行なし} \end{array} \right.$

x_1 : dur 糖尿病の罹病期間 ; year

x_2 : gly 糖化ヘモグロビン ; %

x_3 : bmi 肥満度 ; kg/m²

交互作用

21

薄板平滑化スプライン

- 罹病期間とBMIとには交互作用あり

➡ 薄板平滑化スプラインの利用

交互作用スプラインモデル：

$$\ln\left(\frac{\pi}{1-\pi}\right) = gly + s(dur, bmi)$$

```
proc gam data=WESDR( where=(id^=&d) );
model ret = param( gly ) spline2( dur , bmi , df=&df )
/ link = logit dist = binomial ;
run;
```

22

自由度の最適選択

- 最適選択の結果

	CV	逸脱度	df	誤判別率
1例消去CV法	763.49	743.47	9	0.272
SAS自動選択	764.25	744.81	8.3	0.274

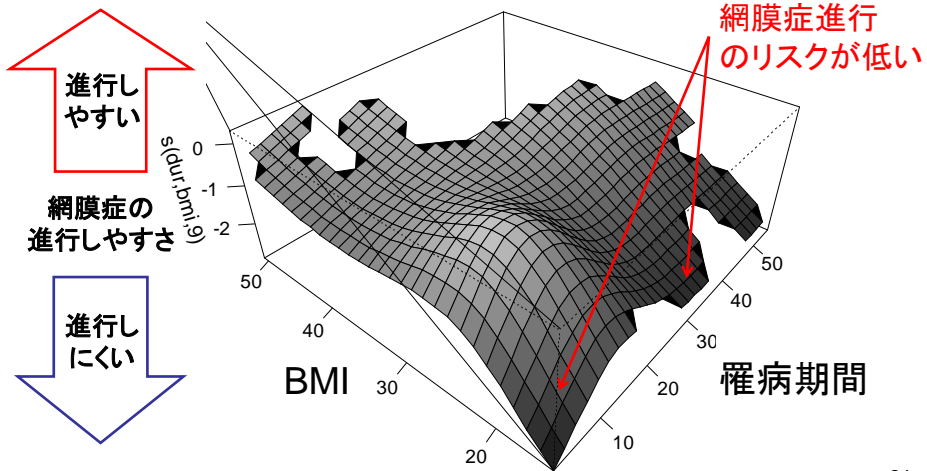
- 他の判別手法との比較

誤判別率 = {

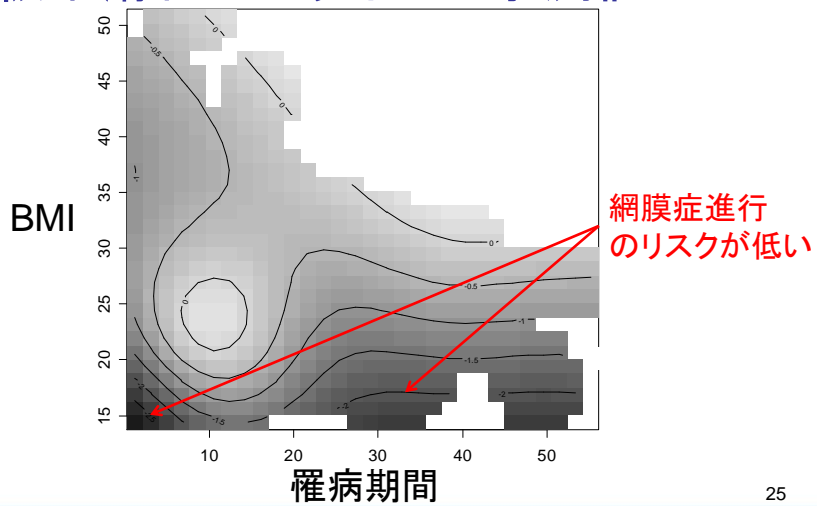
- 0.272 : GAM
- 0.306 : ニューラルネット
- 0.293 : ロジスティック判別
- 0.318 : 線形判別

23

薄板平滑化スプラインの予測値



薄板平滑化スプラインの予測値



シミュレーション・データの生成

手順1: 予測変数として, 一様乱数

$$x_1, x_2 \in [-1, +1]$$

を生成

手順2: x_1, x_2 を

$$f(x_1, x_2) = \sin(2 \times 3.14 \times x_1) + \alpha \times x_1 \times x_2 + \sin(2 \times 3.14 \times x_2)$$

へ代入

交互作用項

交互作用の調整用

$$\alpha = 1, 5$$

26

シミュレーション・データの生成

手順3: p を求める

$$p = \frac{1}{1 + \exp\{-f(x_1, x_2)\}}, 0 \leq p \leq 1$$

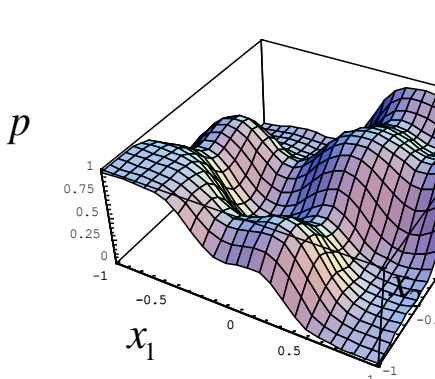
手順4: 手順3の p を用いて, ベルヌーイ乱数

$$y \sim \text{Bin}(1, p) \Rightarrow y = \begin{cases} 1 & \text{(第1群)} \\ 0 & \text{(第2群)} \end{cases}$$

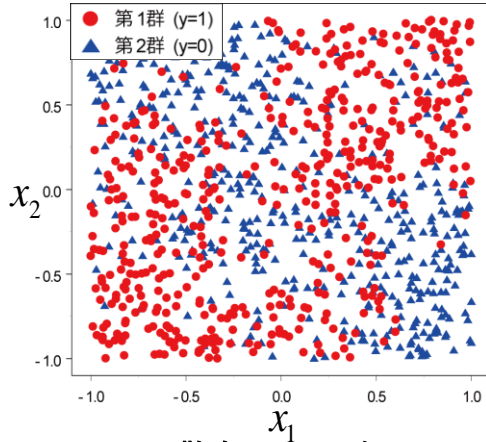
を発生

27

シミュレーション・データ($\alpha=5$)



(a) 3次元プロット

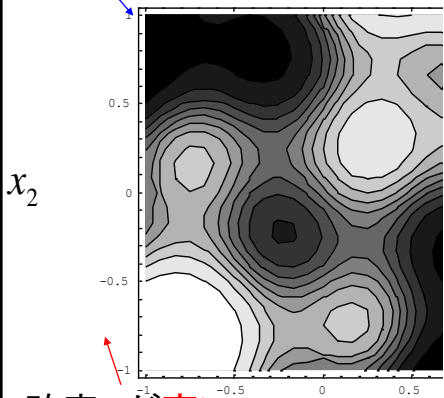


(b) 散布図(1000例)

28

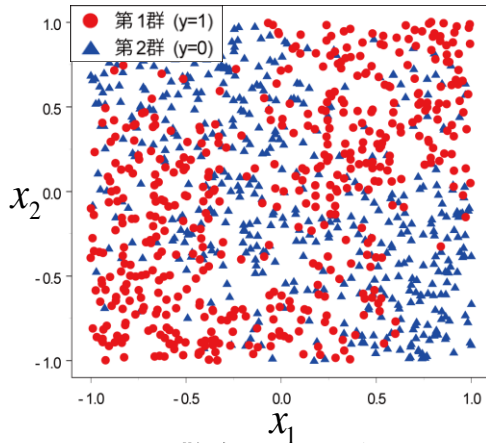
シミュレーション・データ($\alpha=5$)

確率 p が低い



確率 p が高い x_1

(a) 等高線



(b) 散布図(1000例)

29

シミュレーションの設定

- 観測値の構成
 - 予測変数 x_1, x_2 、2値応答 y
- シミュレーション・データ
 - 訓練標本(モデル構築用)
 - 検証標本(未知のデータ)
- サンプル・サイズ
 - 100, 200, 400, 600, 800, 1000例

30

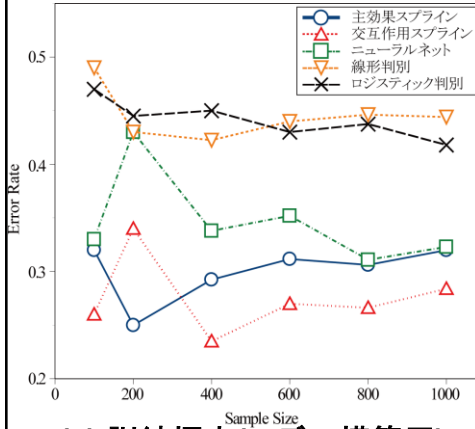
シミュレーションの設定

- 判別モデル
 - 主効果スプライン $\ln\left(\frac{p}{1-p}\right) = s(x_1) + s(x_2)$
 - 交互作用(薄板)スプライン $\ln\left(\frac{p}{1-p}\right) = s(x_1, x_2)$
 - ニューラルネットワーク
 - 隠れユニット数：100,200例=5個,
400,600,800,1000例=6個
 - 線形判別
 - ロジスティック判別

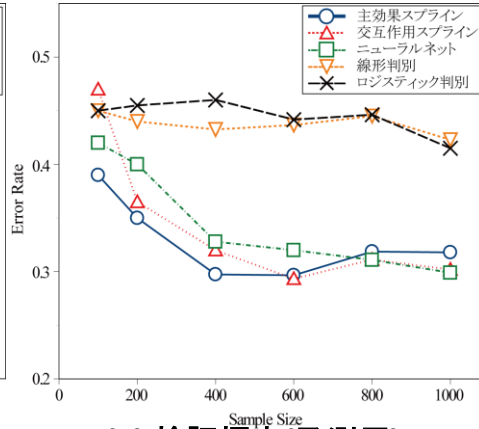
31

性能評価($\alpha=1$)

● 誤判別率の比較



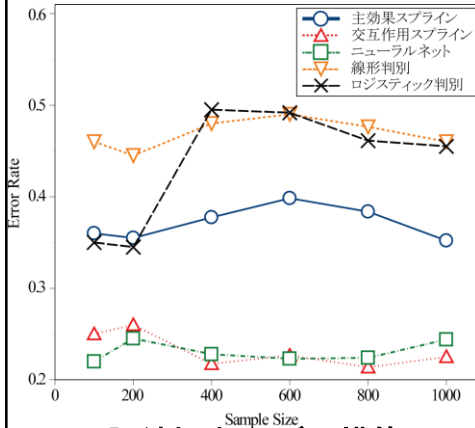
(a) 訓練標本(モデル構築用)



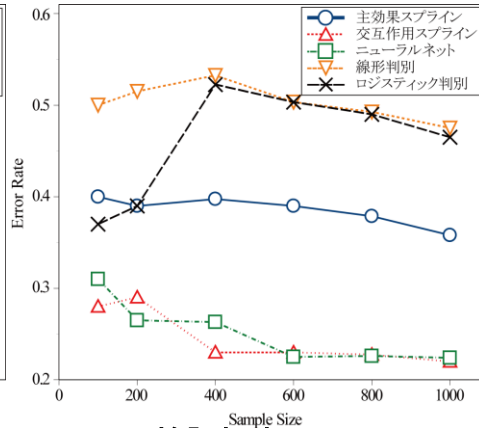
(b) 検証標本(予測用)

性能評価($\alpha=5$)

● 誤判別率の比較



(a) 訓練標本(モデル構築用)



(b) 検証標本(予測用)

まとめ

- 最適な自由度の決定が必要
- 共変量の非線形性の**可視化**
 - ➡ ニューラルネットでは不可能
- 交互作用効果の考慮