

# GAM(一般化加法モデル)による生存時間解析

大阪電気通信大学 情報通信工学部 辻谷将明

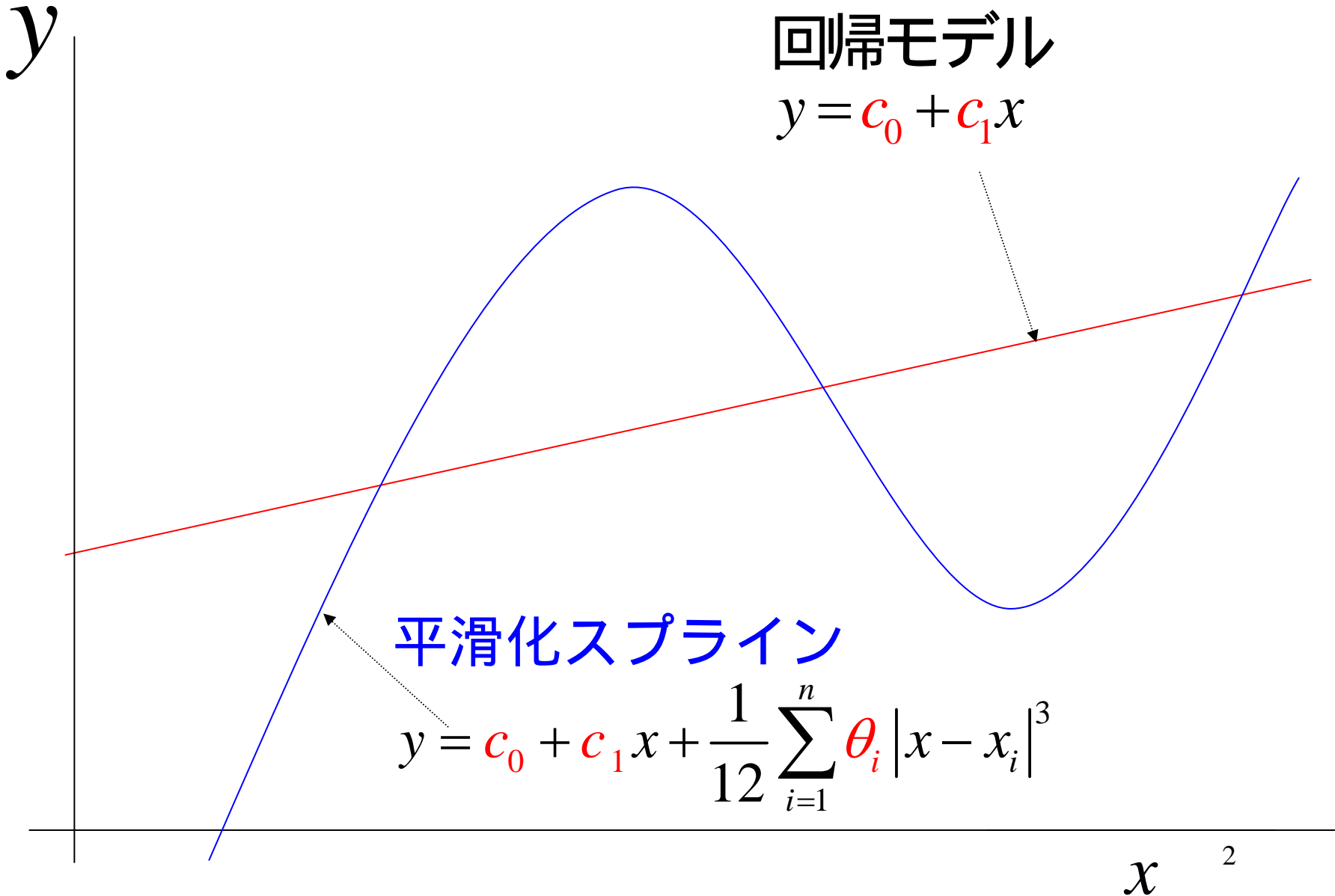
# 平滑化スプラインとは

回帰モデル

$$y = c_0 + c_1 x$$

平滑化スプライン

$$y = c_0 + c_1 x + \frac{1}{12} \sum_{i=1}^n \theta_i |x - x_i|^3$$



$x^2$

データ :  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$



$$y = c_0 + c_1 x + \frac{1}{12} \sum_{i=1}^n \theta_i |x - x_i|^3$$

平滑化パラメータ

$$\begin{pmatrix} \hat{\theta} \\ \hat{c} \end{pmatrix} = \begin{pmatrix} R + \lambda I_n & Q^t \\ Q & 0 \end{pmatrix}^{-1} \begin{pmatrix} y \\ 0 \end{pmatrix}$$

$$\mathbf{y} = (y_1, y_2, \dots, y_n)^t, \mathbf{Q} = \begin{pmatrix} 1 & 1 & \cdot & \cdot & 1 \\ x_1 & x_2 & \cdot & \cdot & x_n \end{pmatrix}$$

$$\mathbf{R} = \begin{pmatrix} 0 & \frac{|x_1 - x_2|^3}{12} & \cdot & \frac{|x_1 - x_n|^3}{12} \\ \frac{|x_2 - x_1|^3}{12} & 0 & \cdot & \frac{|x_2 - x_n|^3}{12} \\ \cdot & \cdot & \cdot & \cdot \\ \frac{|x_n - x_1|^3}{12} & \frac{|x_n - x_2|^3}{12} & \cdot & 0 \end{pmatrix} \leftarrow N\text{-ポイント}$$

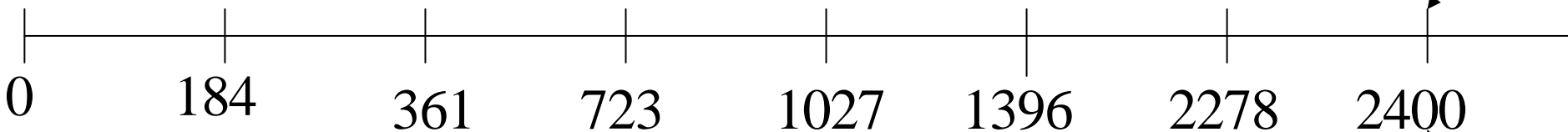
$$\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n)^t, \mathbf{c} = (c_0, c_1)^t$$

# PBC(原発性胆汁性肝硬変)データ(辻谷ら,2005)

時間依存性変数

No.	生存時間	打ち切り (=0)	共 変 量 (初診時の値)		
			年齢	プロトロンビン時間	ビリルビン値
1	400	0	58.8	12.2	14.5
.	.	.	.	.	.
9	2400	1	42.5	11.0	3.2
.	.	.	.	.	.
312	788	0	33.2	10.8	6.4

# 患者#9の入力データ(死亡例)

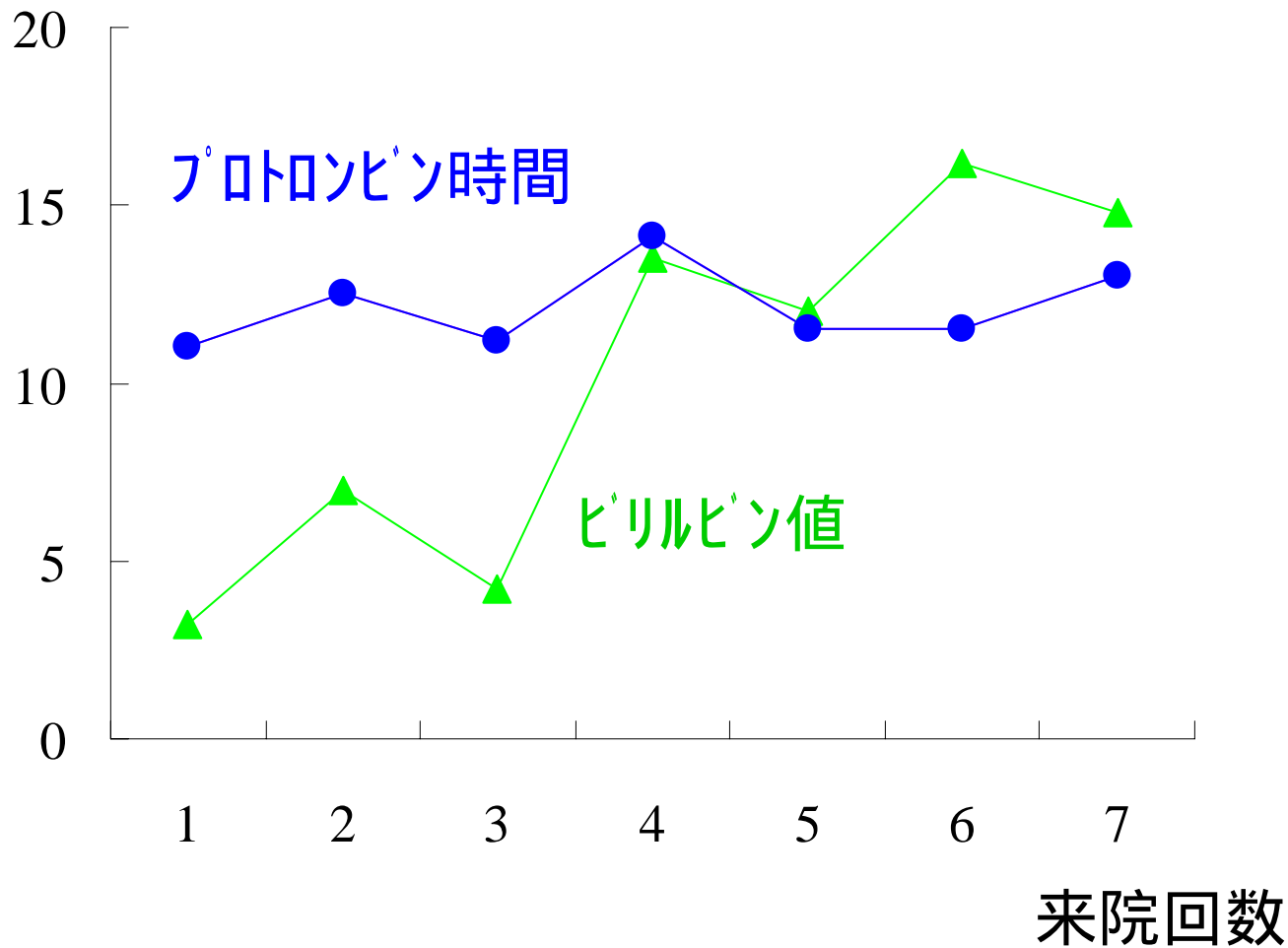


時間 区間 $l_9$	生存時間 $a_l^{<9>}$	年齢 $x_{l1}^{<9>}$	プロトロン ビン時間 $x_{l2}^{<9>}$	ビリルビン 値 $x_{l3}^{<9>}$	応答 $\delta^{<9>}$
1	92.0	42.5	11.0	3.2	0
2	272.5	43.0	12.5	7.0	0
3	542.0	43.5	11.2	4.2	0
4	875.0	44.5	14.1	13.5	0
5	1211.5	45.4	11.5	12.0	0
6	1837.0	46.4	11.5	16.2	0
7	2339.0	48.8	13.0	14.8	1

来院日

死亡

プロトンビン時間  
ビリルビン値



患者#9の時間依存型データ

# 企業倒産マトリックス(森平ら, 2002)

業種	90年 企業数	倒産企業数									
		91	92	93	94	95	96	97	98	99	00
建設	2379	15	25	26	20	16	22	44	65	35	40
製造	4194	12	21	27	19	23	19	19	42	28	29
卸し小売	1991	36	34	34	40	38	32	37	75	56	51
その他	1499	13	23	14	12	5	7	14	12	16	14
全体	13848	76	103	101	91	82	80	114	194	135	134



# 建設

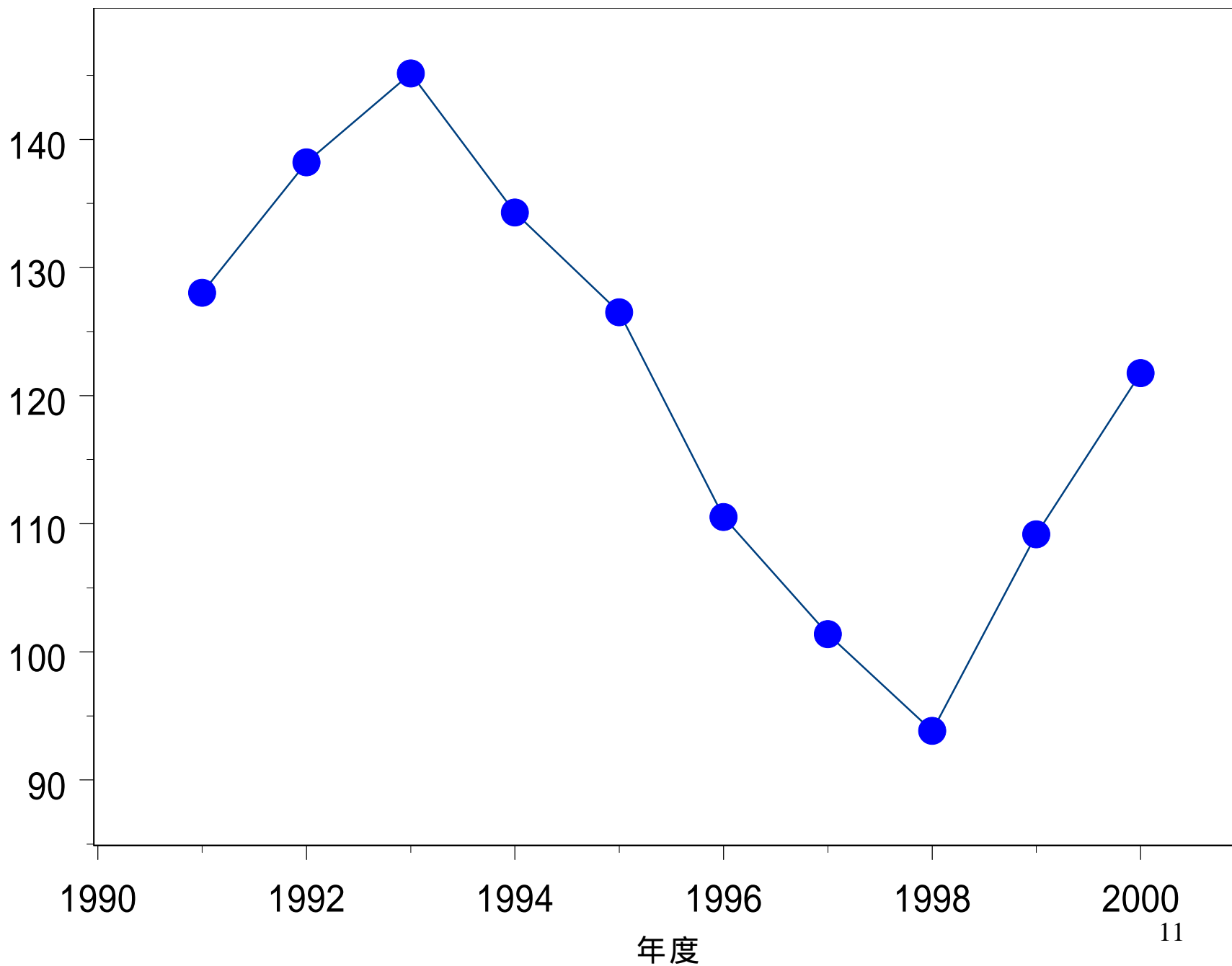
年度	企業数	倒産企業数
91	2379	15
92	2364	25
93	2339	26
94	2313	20
95	2293	16
96	2277	22
97	2255	44
98	2211	65
99	2146	35
00	2111	40

# 建設

時間依存型共変量

年度	企業数	倒産企業数	時点	為替レ-ト
91	2379	15	0.5	128.014
92	2364	25	1.5	138.20
93	2339	26	2.5	145.14
94	2313	20	3.5	134.29
95	2293	16	4.5	126.51
96	2277	22	5.5	110.53
97	2255	44	6.5	101.39
98	2211	65	7.5	93.83
99	2146	35	8.5	109.18
00	2111	40	9.5	121.76

為替レイト



# 二項型モデル

$N_i$  = 第  $t_i$  時点の直前まで生存していた企業数

$r_i$  = 第  $t_i$  時点の間に倒産した企業数

$r'_i$  = 第  $t_i$  時点の間に打切られた企業数  
( $N_{i+1} = N_i - r_i - r'_i$ )

時点	$N_i$	$r_i$	$r'_i$
$t_1$	2379	15	0
$t_2$	2364	25	0
$t_3$	2339	26	0
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$t_{10}$	2111	40	2071

$$r_i | N_i \sim B_i(N_i, h_i), i = 1 \sim n$$

i.e.,

$$f(r_i) = \binom{N_i}{r_i} h_i^{r_i} (1 - h_i)^{N_i - r_i}, r_i = 0, 1, 2, \dots, N_i$$

$h_i = \Pr \{ \text{第 } i \text{ 時点で倒産} \mid \text{第 } i \text{ 時点の直前まで生存} \}$   
: 離散 **ハザード**

打切りがなければ、 $r_i = N_i - N_{i+1}$  となり、 $r_1, r_2, \dots$  は、  
 $N_1, N_2, \dots$  によって完全に決まる



独立性のもとでの全尤度は使えない

$$\text{全对数尤度: } \ln L(\alpha, \zeta) \equiv \ln(\alpha) + \ln(\zeta)$$

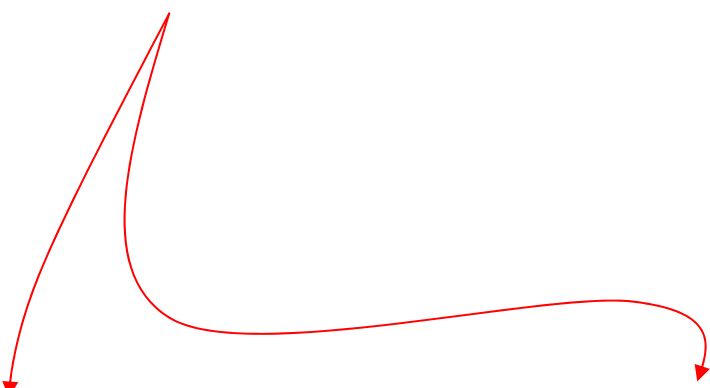
部分对数尤度

$$\ln(\alpha) = \ln \prod_{i=1}^n \binom{N_i}{r_i} h_i^{r_i} (1 - h_i)^{n_i - r_i}$$

$$r_i | N_i \overset{\text{ind.}}{\sim} \sim B_i(N_i, h_i)$$

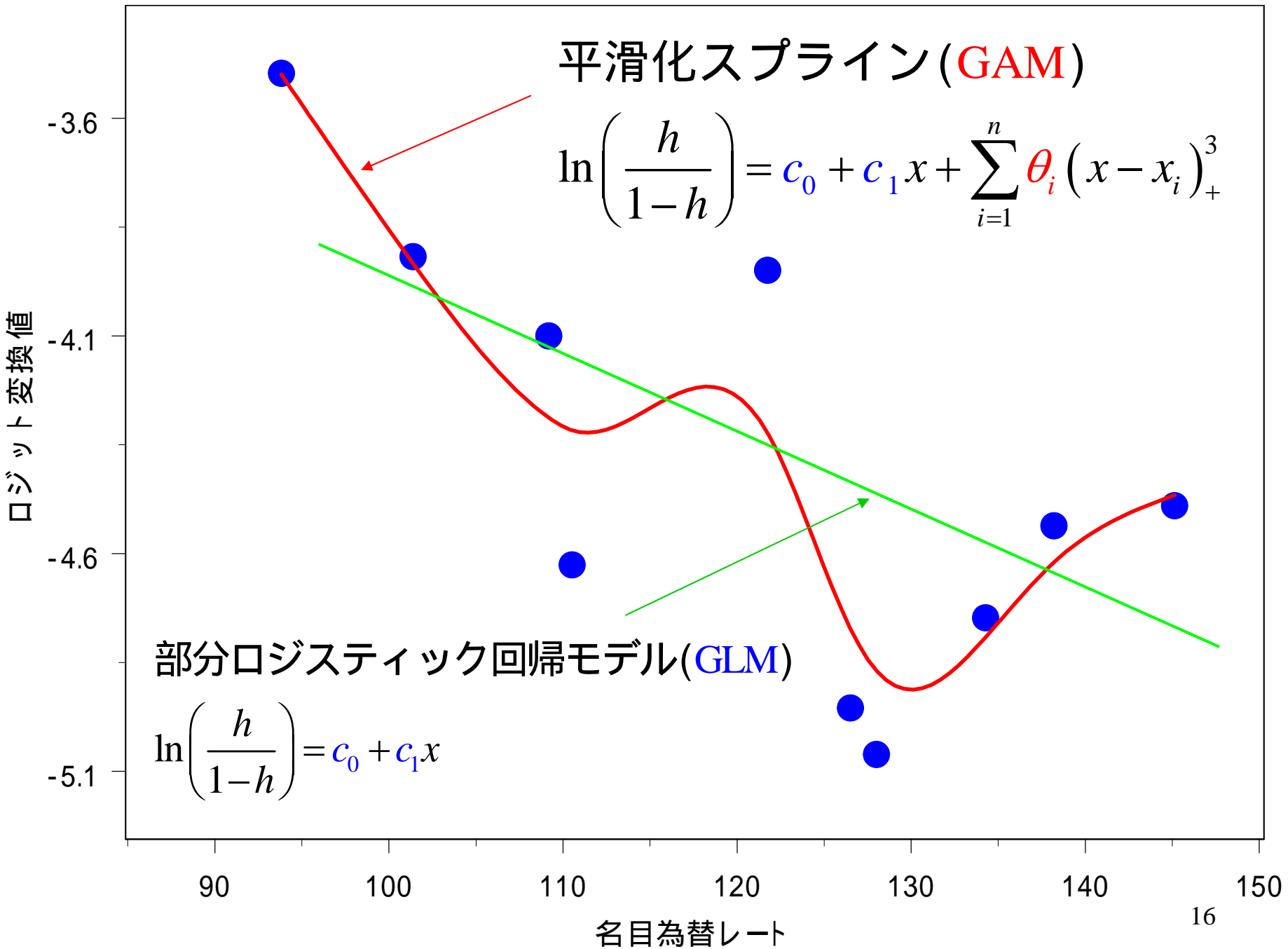
# 対数尤度最大

$$h_i = \begin{cases} \text{ロジスティック回帰モデル} \\ \text{平滑化スプライン} \\ \text{ニューラルネット} \end{cases}$$

$$\ln L(\mathcal{X} : \boldsymbol{\alpha}) = \sum_{i=1}^n \{ r_i \ln h_i + (N_i - r_i) \ln(1 - h_i) \}$$


$$\boldsymbol{\alpha} = \{ c_0, c_1, \theta_1, \theta_2, \dots, \theta_n \}$$

c.f., **K - M**推定量:  $\tilde{h}_i = r_i / N_i$





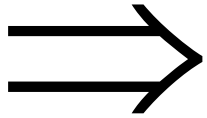
$$\ln\left(\frac{h}{1-h}\right) = \begin{cases} c_0 + c_1 x : \text{ロジスティック回帰 (GLM)} \\ c_0 + c_1 x + \frac{1}{12} \sum_{i=1}^n \theta_i (x - x_i)_+^3 : \text{平滑化スプライン (GAM)} \end{cases}$$

$$h = \sum_{j=0}^J \left\{ \frac{1}{1 + \exp\left[-\sum_{j=0}^J \left\{ \frac{\beta_j}{1 + \exp(-\alpha_j x)} \right\} \right]} \right\} : \text{ニューラルネット}$$

$h$  = 離散ハザード率

$x$  = 為替

## 平滑化スプライン



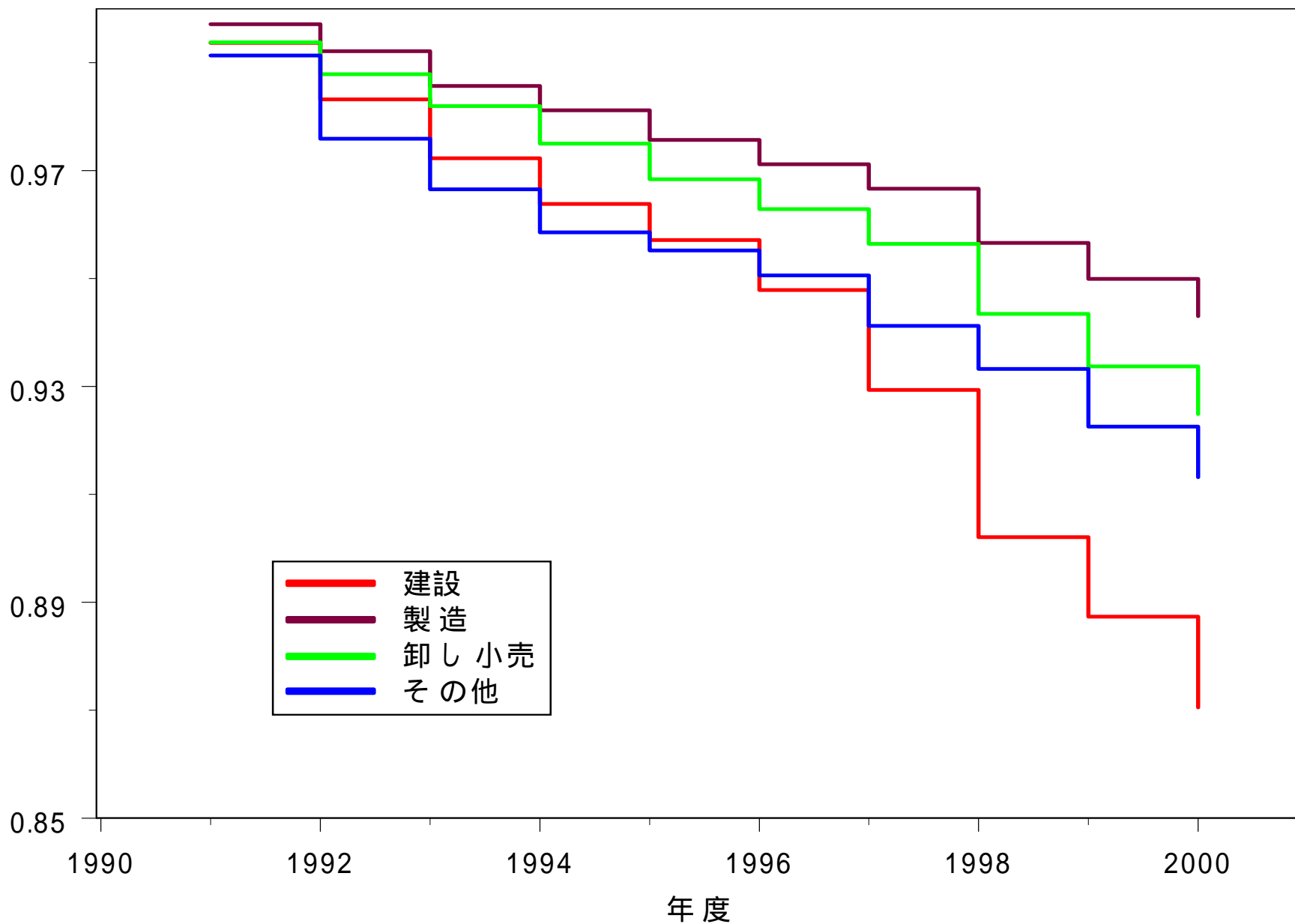
局所評点化法は、 $\lambda$  カルティ付き対数尤度

$$\sum_{i=1}^n \{r_i \ln h_i + (N_i - r_i) \ln(1 - h_i)\} + \frac{1}{2} \sum_{j=1}^p \lambda_j \int \{f_j''(t)\}^2 dt$$

を最大化

H-T,p.149;H-T-F,5.6節

生存率



# 生存関数の比較

# モデル適合度(逸脱度)

業種	GAM( <i>d.f.</i> )	GLM( <i>d.f.</i> =7)
全体	1.0016(1.0056)	45.91
建設	0.9311(1.3023)	23.09
製造	1.5991(1.0211)	15.62
卸し・小売	1.6641(1.0080)	16.10
その他	1.5596(1.0042)	16.18

# 非線形効果の有無

業種	時点	為替
全体	*	**
建設		**
製造		**
卸し・小売		*
その他	*	

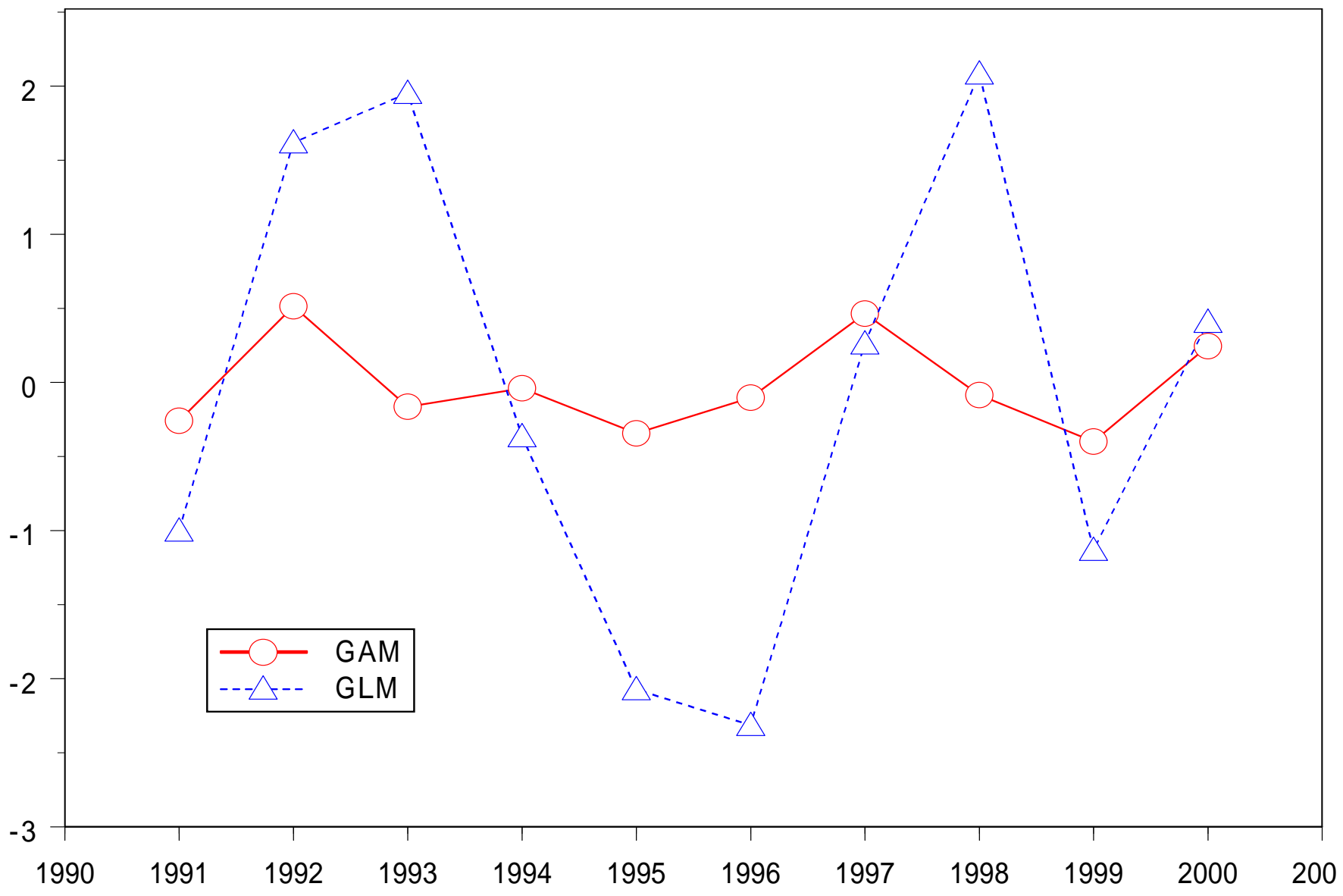
\*\* : 1% 有意、 \* : 5% 有意、    : 10% 有意

# 残差分析

符号付き逸脱度残差

$$Dev^{<i>} = \text{sgn}\left(r_i - N_i \hat{h}_i\right) \sqrt{2r_i \ln\left(\frac{r_i}{N_i \hat{h}_i}\right) + 2(N_i - r_i) \ln\left(\frac{N_i - r_i}{N_i - N_i \hat{h}_i}\right)}$$

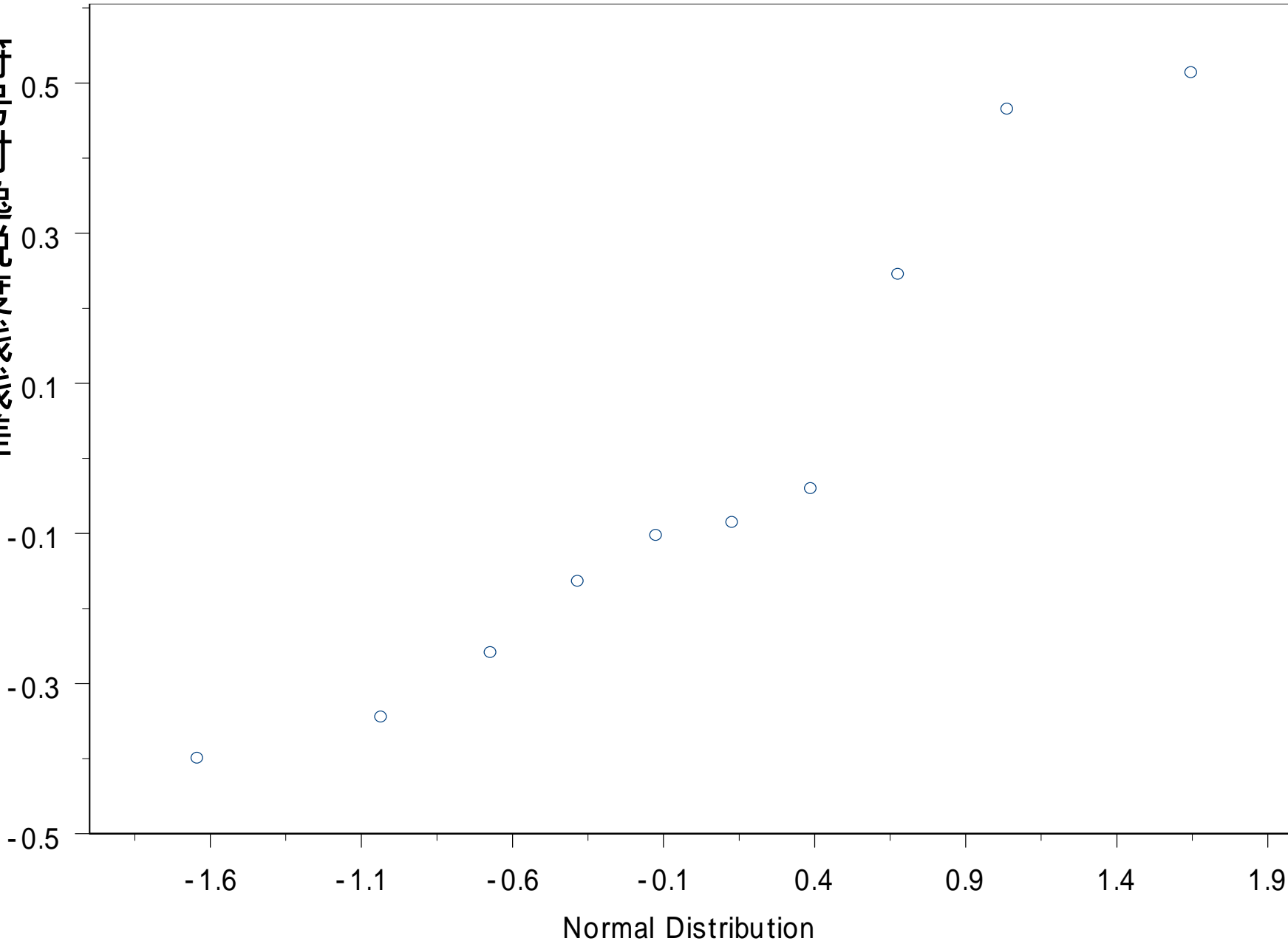
$\sim N(0,1)$



逸脱度残差(建設)

年度

符号付逸脱度残差



$Q-Q$ プロット(平滑化スプライン)



Formula



Variable

Choose Variables:

- 倒産率
- 時点
- 名目為替レート
- 重み

Transformation

Special Term

Term Category:

Option:

Format:

Add

Add

Response

Linear: (+)

Interaction: (:)

Spline: (s)

Loess: (lo)

Offset: (offset)

Remove

Remove Intercept

Term:

Remove

Formula:

倒産率~s(時点)+s(名目為替レート)

OK

Cancel

Apply



current

Help

# 限界(条件付き)倒産確率の推定

生存率の推定

共変量として入力

$$S(x, a_l) = \prod_{1 \leq i \leq l} \{1 - h_i(x, a_i)\}$$

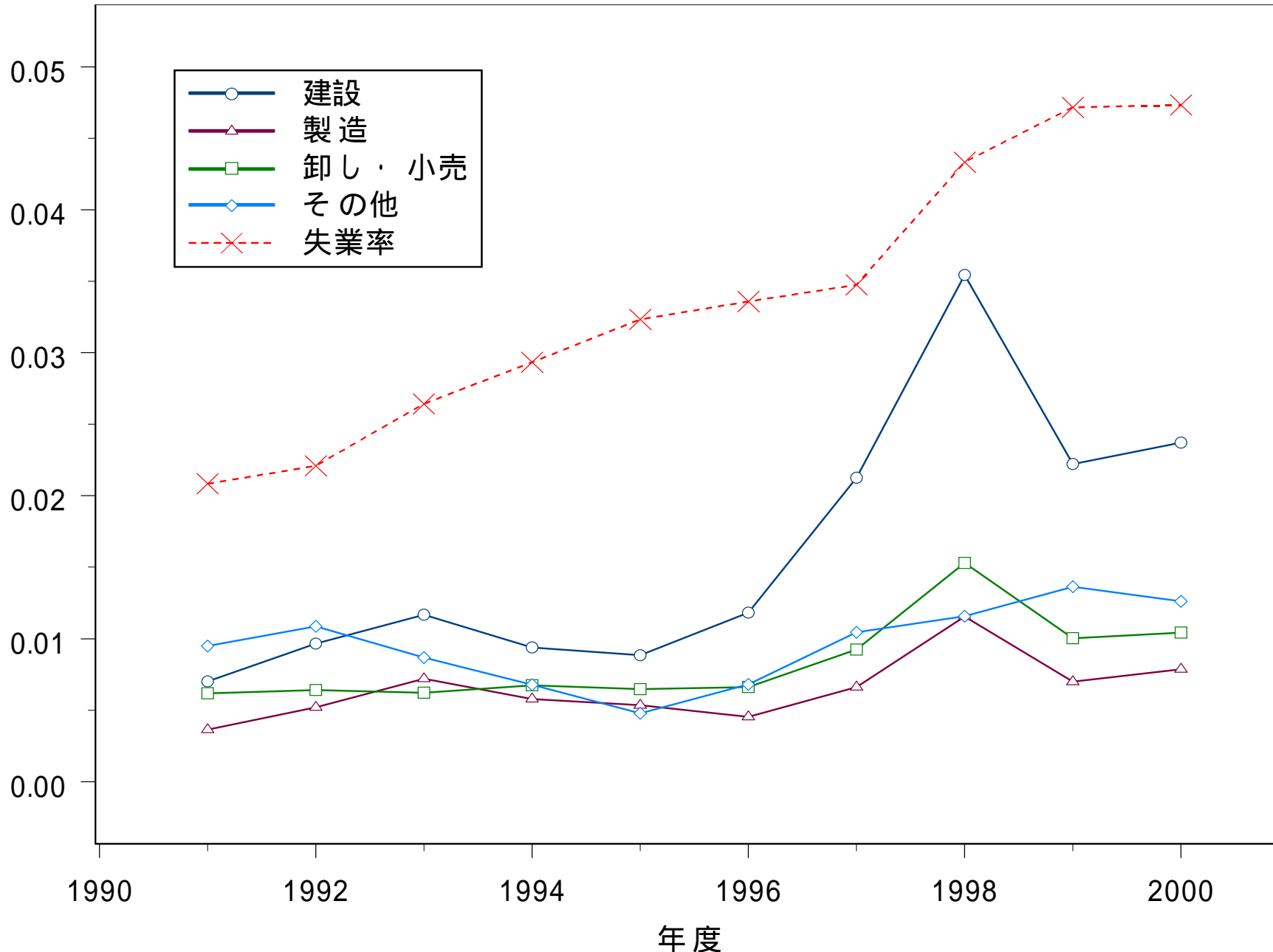
$a_l$ まで生存した企業の次の $\Delta t$ (1年)後の限界倒産確率

$$\Pr(a_l, a_l + \Delta t) = \frac{S(x, a_l) - S(x, a_l + \Delta t)}{S(x, a_l)}$$

c.f., ハザード(瞬間倒産率):

$$h(a_l) = \lim_{\Delta t \rightarrow 0} \left\{ \frac{1}{\Delta t} \left( \frac{S(x, a_l) - S(x, a_l + \Delta t)}{S(x, a_l)} \right) \right\}$$

# 倒産確率



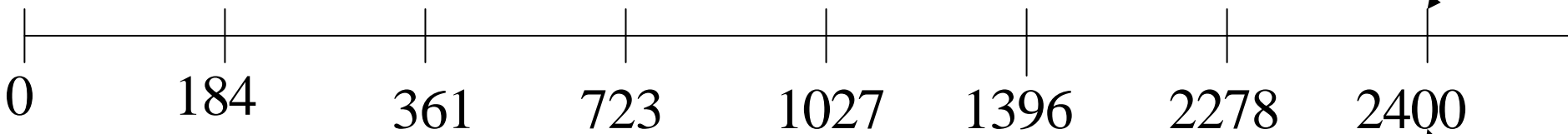
## 限界倒産確率の比較

# PBC(原発性胆汁性肝硬変)データ(辻谷ら,2005)

時間依存性変数

No.	生存時間	打ち切り (=0)	共 変 量 (初診時の値)		
			年齢	プロトロンビン時間	ビリルビン値
1	400	0	58.8	12.2	14.5
.	.	.	.	.	.
9	2400	1	42.5	11.0	3.2
.	.	.	.	.	.
312	788	0	33.2	10.8	6.4

# 患者#9の入力データ(死亡例)

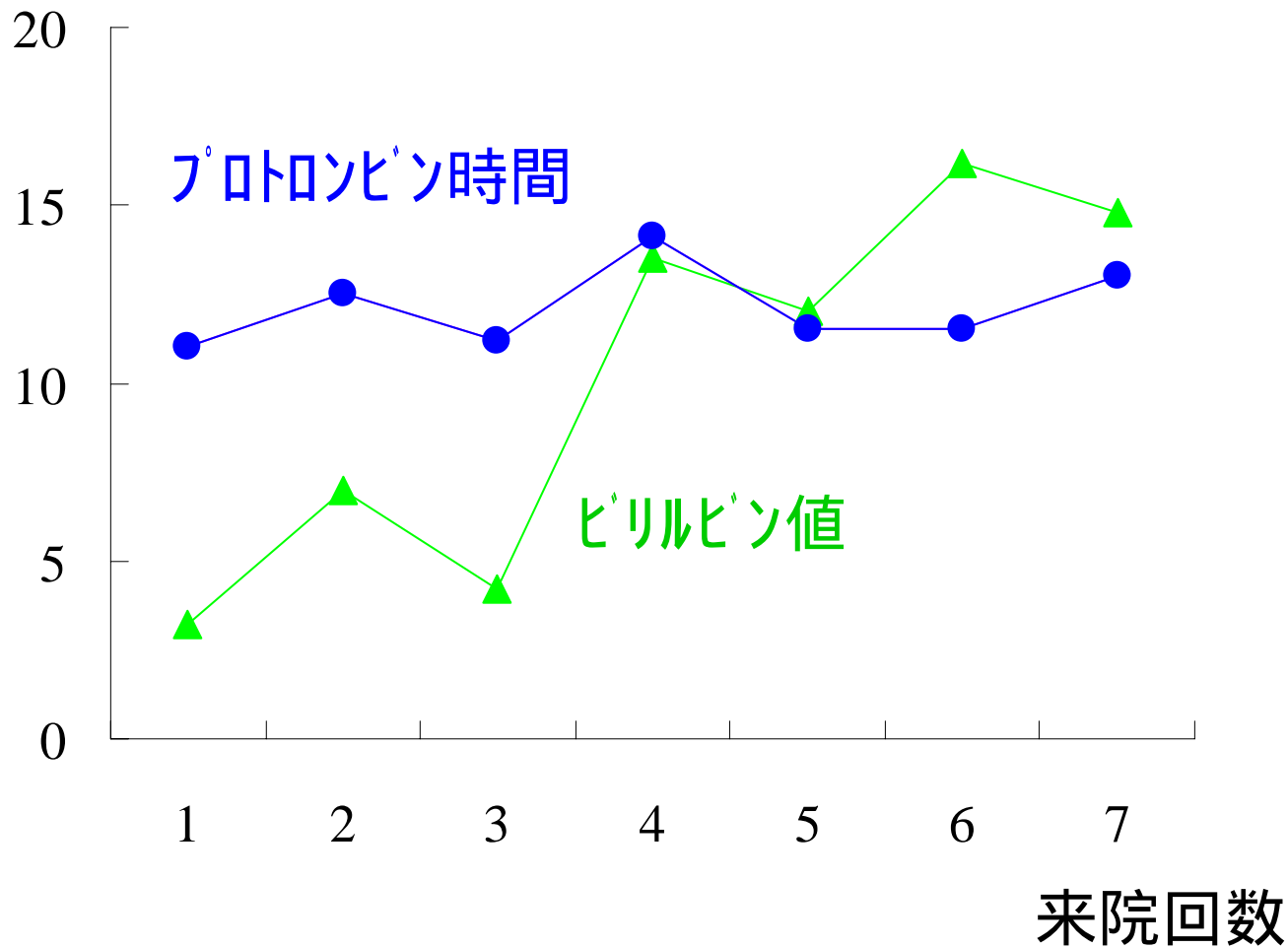


時間 区間 $l_9$	生存時間 $a_l^{<9>}$	年齢 $x_{l1}^{<9>}$	プロトン ピン時間 $x_{l2}^{<9>}$	ビリルビン 値 $x_{l3}^{<9>}$	応答 $\delta^{<9>}$
1	92.0	42.5	11.0	3.2	0
2	272.5	43.0	12.5	7.0	0
3	542.0	43.5	11.2	4.2	0
4	875.0	44.5	14.1	13.5	0
5	1211.5	45.4	11.5	12.0	0
6	1837.0	46.4	11.5	16.2	0
7	2339.0	48.8	13.0	14.8	1

来院日

死亡

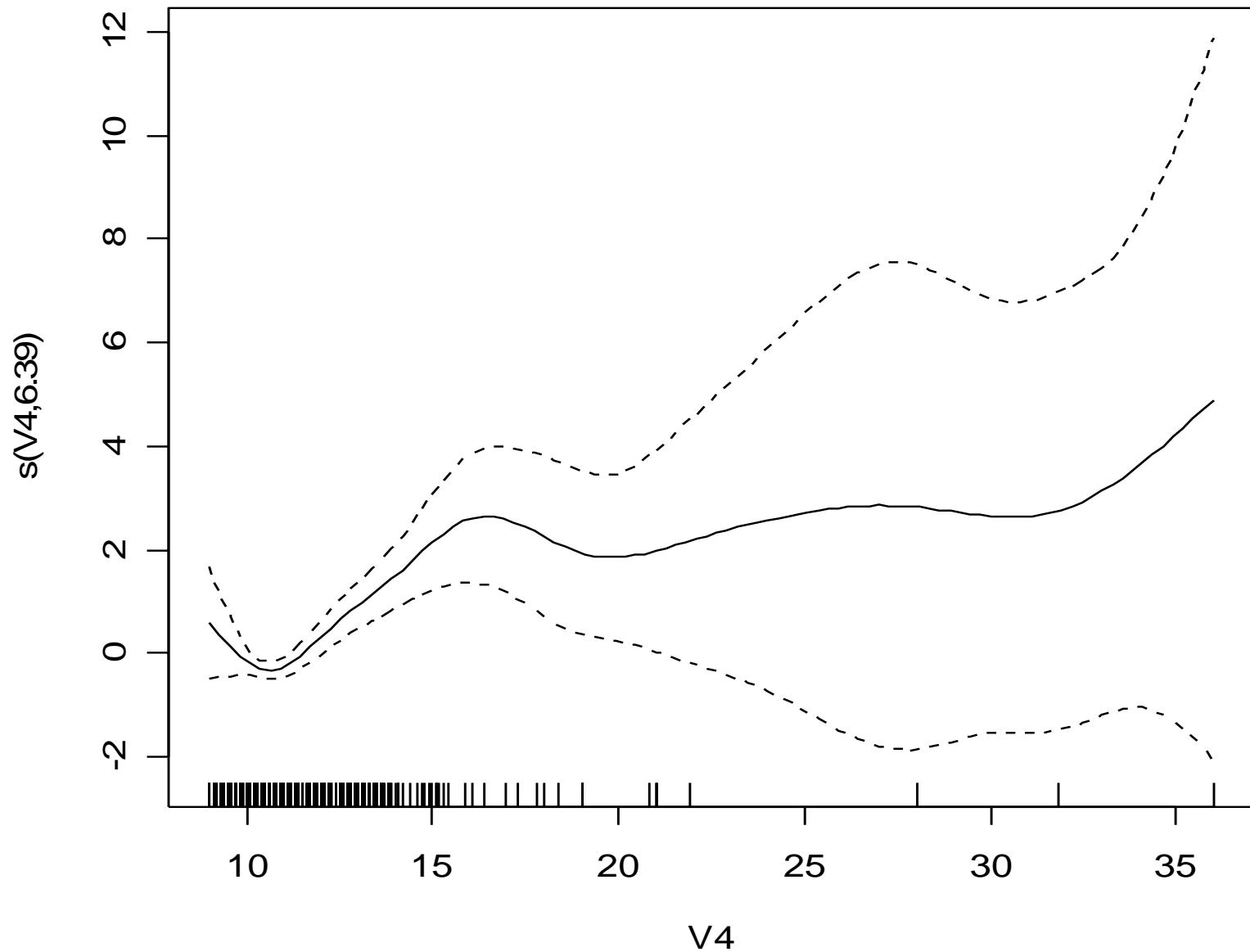
プロトンピン時間  
ビリルビン値



患者#9の時間依存型データ

# 共変量の平滑関数に対する有意性検定

共変量	edf	Chi.sq
年齢	1.000	43.679**
プロトロンビオン時間	3.254	138.66**
ビリubin値	6.386	44.128**



平滑化スプラインを当てはめた結果(ビリビン値)



# 限界(条件付き)倒産確率の推定

## 生存率の推定

$$S(\mathbf{x}, a_l) = \prod_{1 \leq i \leq l} \{1 - h_i^{<d>}(\mathbf{x}, a_i)\}$$

## $\Delta t$ (6ヶ月)後の条件付き生存率

$$\Pr(a, a + \Delta t) = \frac{S(\mathbf{x}, a + \Delta t)}{S(\mathbf{x}, a)}, a_l = \frac{t_{l-1} + t_l}{2}, t_0 = 0$$

# 短期予後予測(患者#9の6ヶ月後の生存率)

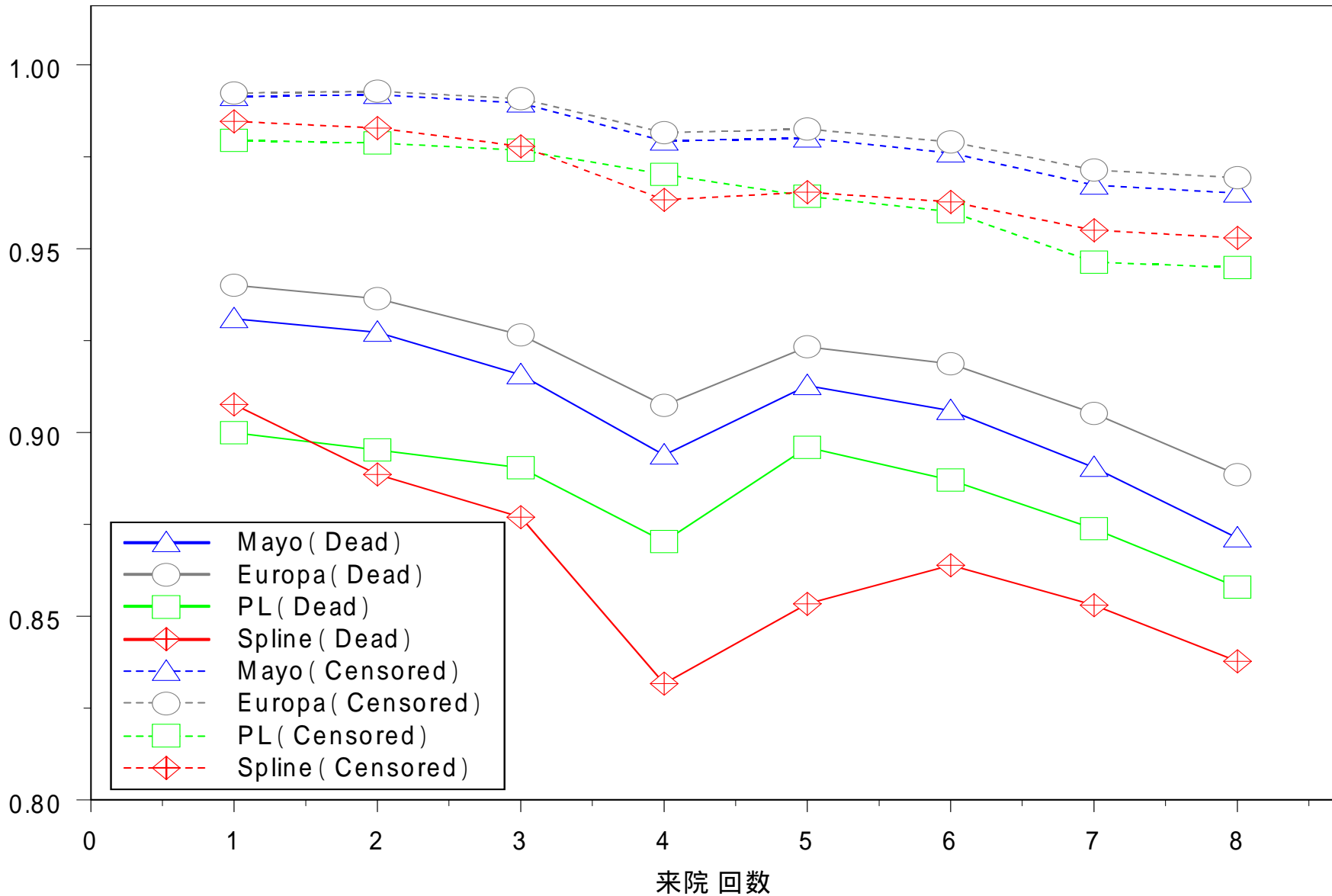
共変量		来院回数						
		1	2	3	4	5	6	7
年齢		42.5	43.0	43.5	44.5	45.3	46.3	48.7
Bili		3.2	7.0	4.2	13.5	12.0	16.2	14.8
ProTime		11.0	12.5	11.2	14.1	11.5	11.5	13.0
生存率	Mayo	.987	.939	.980	.767	.903	.851	.763
	スプライン	.987	.912	.972	.517	.823	.794	.627

## 群 $g$ における $\Delta t$ 後の条件付き生存率の平均値

$$S_g(l) = \frac{1}{n_l^{[g]}} \sum_{d=1}^{n_l^{[g]}} \Pr_d^{[g]}(l, l + \Delta t), \quad g = 1, 2; l = 1, 2, \dots$$

$\Pr_d^{[g]}(l, l + \Delta t)$ : 来院回数に $l$ における $\Delta t$ 後の条件付き生存率

条件付き生存率



4手法による全312例中の死亡例、打ち切り例に対する条件付き生存率の比較

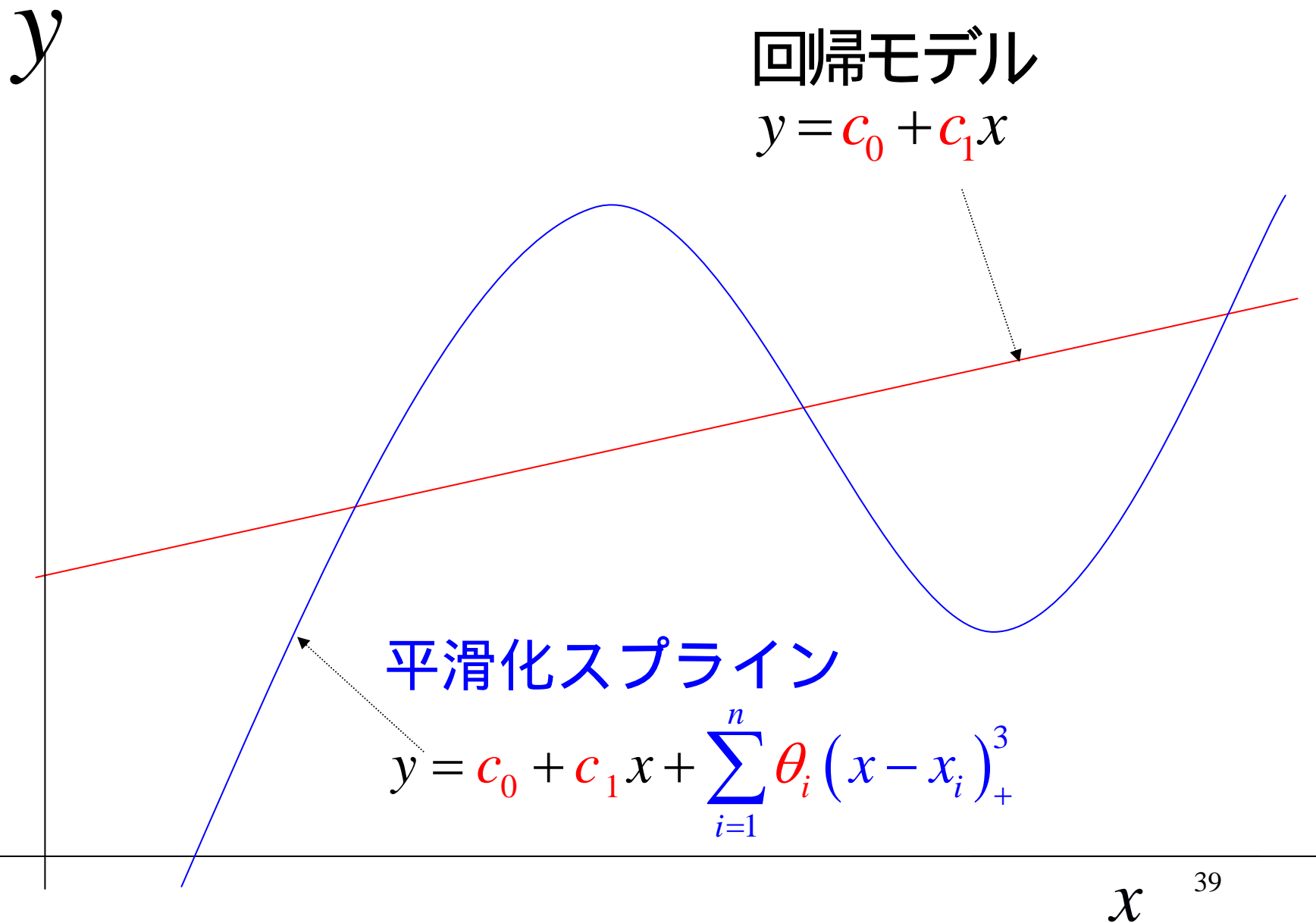
## 平滑化スプラインの特徴

1. 時間依存型変数の取扱いが容易
2. 観測期間全体を考慮に入れて解析
3. 比例ハザード性の制約を受けない
4. 共変量の非線形性を考慮
5. 死亡例が少ない場合も適用可
6.  $ML$ 解が一意
7. 計算時間が短い

# 参考文献

- [1] Biganzoli, E., Boracchi, P., Marriani, L. and Marubini, E. (1988): Feed forward neural networks for the analysis of censored survival data: a partial logistic approach, *Statist. Medicine*, **17**, pp.1169-1186.
- [2] Cox, D.R. (1975): Partial likelihood, *Biometrika*, **62**, pp.269-276.
- [3] Efron, B. (1988): Logistic regression, survival analysis, and Kaplan-Meier curv, *J. Amer. Statist. Assoc.*, **83**, pp.414-425.
- [4] 森平総一郎ら(2002):倒産確率の期間構造推定と信用リスクのある債権の評価、第17回JAFEE大会
- [5] Murtaugh, P.A. et al., (1994): Primary biliary cirrhosis: Prediction of short-term survival based on repeated patient visits, *Hepatology*, **20**, pp.126-134.
- [6] 辻谷将明, 左近賢人(2005): 時間依存型共変量を伴う生存データの解析, *応用統計学*, **34**, 15-29.

# 付録A. 平滑化スプライン(応答が連続値)



データ :  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

$$y = c_0 + c_1 x + \frac{1}{12} \sum_{i=1}^n \theta_i |x - x_i|^3$$

平滑化パラメータ

$$\begin{pmatrix} \hat{\theta} \\ \hat{c} \end{pmatrix} = \begin{pmatrix} R + \lambda I_n & Q^t \\ Q & \mathbf{0} \end{pmatrix}^{-1} \begin{pmatrix} y \\ \mathbf{0} \end{pmatrix}$$



$$\mathbf{y} = (y_1, y_2, \dots, y_n)^t, \mathbf{Q} = \begin{pmatrix} 1 & 1 & \cdot & \cdot & 1 \\ x_1 & x_2 & \cdot & \cdot & x_n \end{pmatrix}$$

$$\mathbf{R} = \begin{pmatrix} 0 & \frac{|x_1 - x_2|^3}{12} & \cdot & \frac{|x_1 - x_n|^3}{12} \\ \frac{|x_2 - x_1|^3}{12} & 0 & \cdot & \frac{|x_2 - x_n|^3}{12} \\ \cdot & \cdot & \cdot & \cdot \\ \frac{|x_n - x_1|^3}{12} & \frac{|x_n - x_2|^3}{12} & \cdot & 0 \end{pmatrix} \leftarrow N\text{-ポイント}$$

$$\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n)^t, \mathbf{c} = (c_0, c_1)^t$$

$\lambda (\geq 0)$ を大きくする  $\Leftrightarrow$  滑らかな曲線を求めることに重点をおく

$f(x)$ の曲率

$$\sum_{i=1}^n \left[ \ln \left( \frac{h_i}{1-h_i} \right) - f(x_i) \right]^2$$

+  $\lambda$

$$\int \{ f''(x) \}^2 dt$$

$\Rightarrow$  最小になる  $f(x)$

小さいほどモデルの当てはまりは良い

曲げ弾性エネルギー (小さいほど滑らかな曲線さ)

$\Leftrightarrow$  3次の自然スプライン

竹澤(上):p.198

$\lambda (\geq 0)$  を大きくする  $\Leftrightarrow$  滑らかな曲線を求めることに重点をおく

$$\sum_{i=1}^n [y - f(x_i)]^2 + \lambda \int \{f''(x)\}^2 dt \Rightarrow \text{最小になる平滑化曲線}$$

$\Leftrightarrow$  3次の自然スプライン

$$y = c_0 + c_1 x + \sum_{i=1}^n \theta_i (x - x_i)_+^3$$
$$(x - x_i)_+^3 = \begin{cases} (x - x_i)^3 & : x \geq x_i \\ 0 & : x < x_i \end{cases}$$

## 付録B: **二値**応答の場合(共変量2個)へ拡張

$$\ln\left(\frac{h_i}{1-h_i}\right) = \alpha + s_1(x_{1i}) + s_2(x_{2i}), i = 1, 2, \dots, n$$

$$= \alpha + \beta_1 x_{1i} + f(x_{1i}) + \beta_2 x_{2i} + f(x_{2i})$$
$$= \underbrace{\alpha + \beta_1 x_{1i} + \beta_2 x_{2i}} + \underbrace{f(x_{1i}) + f(x_{2i})}$$

GLMの重付き  
最小二乗法

後退当てはめによ  
る局所評点化法

# 局所評点化法

データ:  $(\mathbf{x}_i, \delta_i)$ ,  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})$ ,  $\delta_i = \begin{cases} 1: \text{個体 } i \text{ が死亡} \\ 0: \text{その他} \end{cases}$

$$\text{モデル: } \ln \left\{ h_i / (1 - h_i) \right\} = \alpha + \sum_{j=1}^p s_j(x_{ij}), i = 1, 2, \dots, n$$

## ステップ 1 (初期化)

$$\hat{s}_j(x_{ij}) \equiv 0, i = 1, 2, \dots, n; j = 1, 2, \dots, p$$

$$\hat{\alpha} = \ln \left( \frac{\bar{\delta}}{1 - \bar{\delta}} \right), \quad \bar{\delta} = \sum_{i=1}^n \delta_i / n$$

## ステップ 2

$$\hat{\eta}_i = \hat{\alpha} + \sum_{j=1}^p \hat{s}_j(x_{ij}), \quad \hat{h}_i = \frac{1}{1 + \exp(-\hat{\eta}_i)}$$

### ステップ 3(反復)

$$(a) z_i = \hat{\eta}_i + \frac{\delta_i - \hat{h}_i}{\hat{h}_i (1 - \hat{h}_i)} : \text{調整従属変数}$$

$$(b) w_i = \hat{h}_i (1 - \hat{h}_i) : \text{重み}$$

(c)  $(x_1, z_1), (x_2, z_2), \dots, (x_n, z_n)$  に後退当てはめ法を適用し、 $j = 1, 2, \dots, p$  について  $\hat{f}_j$  を算出する。そして、 $\hat{h}_i$  を推定する。

### ステップ 4(収束判定)

$$\text{逸脱度 } D(\delta, \hat{h}) = -2 \sum_{i=1}^n \left[ \delta_i \ln \hat{h}_i + (1 - \delta_i) \ln (1 - \hat{h}_i) \right]$$

が収束するまでステップ 3 を繰返す

## 後退当てはめ法(ガウス・ザイデル法)

$$z_i = \alpha + f_1(x_{1i}) + f_2(x_{2i})$$

ステップ 1(初期化)  $\hat{f}_2^{(0)} \equiv 0$

ステップ 2  $x_{1i}$  に対して  $z_i - \hat{f}_2^{(0)}(x_{2i}) - \hat{\alpha}$  を  
平滑化  $\xrightarrow{\text{平滑化スプライン}}$   $\hat{f}_1^{(1)}$

ステップ 3  $x_{2i}$  に対して  $z_i - \hat{f}_1^{(1)}(x_{1i}) - \hat{\alpha}$  を  
平滑化  $\xrightarrow{\text{平滑化スプライン}}$   $\hat{f}_2^{(1)}$

ステップ 4  $x_{1i}$  に対して  $z_i - \hat{f}_2^{(1)}(x_{2i}) - \hat{\alpha}$  を  
平滑化  $\xrightarrow{\text{平滑化スプライン}}$   $\hat{f}_1^{(2)}$

ステップ 5 以下

$\hat{f}_1^{(m)}(x_{1i}), \hat{f}_2^{(m)}(x_{2i})$  が収束するまで繰り返す 47

# 短期予後予測モデル

## (1) 日本肝移植適応研究会モデル (1979 ~ 1990; 141例)

死亡6カ月前のデータと、生存例のat randomなデータを入力して、ロジスティック回帰でsimple analysis(「患者対照研究」:8.1節)

$$\lambda = -4.333 + 1.2739 \ln(T - Bili) + 4.4880 \ln(GOT / GTP)$$

$\Rightarrow DR = \frac{1}{1 + \exp(-\lambda)}$ : 任意観測時より6ヶ月後の予想死亡率  $\Rightarrow$  50%を超えると移植考慮

## (2) Mayo updated model (1974 ~ 1984; 312例; D- $\alpha$ ニシラミン trial)

任意観測時より、3~24(6)ヶ月の予測生存率

$$PI(t) = 0.051 \times (age) + 1.209 \times \ln(T - Bil) + 2.754 \times \ln(PT) - 3.304 \times \ln(Alb) + 0.675 \times edema\ index$$

$\Rightarrow S(t, x) = \{S_0(t)\}^{\exp\{PI(t) - 6.119\}}$ : 20%以下になると移植考慮

## (3) Europe new version model (1971 ~ 1983; 248例; アサチオプリン trial)

任意観測時より、1, 3, 6ヶ月の予測生存率

$$PI(t) = 2.53 \times \{\ln(T - Bil) - 1.53\} + 1.39(\text{if } asites\ \text{present}) - 0.085 \times (Alb - 34.3) \\ + 0.040 \times (age - 55) + 0.65(\text{if } GI\ \text{bleeding}\ \text{present})$$

$\Rightarrow P(t, t+h) = \exp(-\lambda_0 \cdot h \cdot \exp[PI(t)])$ ,  $\lambda_0: \hat{\Lambda}_0(t)$  のグラフより推定