

公開講座 2000年7月6日
統計数理研究所 下平英寿

- (I) モデル信頼集合
- (II) 情報量規準の拡張

参考: 「モデル選択理論の新展開」
(統計数理 47, 1999)

C1

1 モデル信頼集合

- モデル選択 \hat{k} — 点推定
AICの差が小さければ, どちらのモデルもおなじくらい良いはず
- モデル信頼集合 \mathcal{I} — 区間推定
モデル選択のバラツキを考慮したうえで, 良いモデルを複数選ぶ

とくにノンネストなモデルの比較に有効

C3

モデル選択, 情報量規準の研究の紹介

1. モデル信頼集合

変数選択問題

モデル選択のバラツキと一貫性

モデル選択の信頼性

2. 情報量規準の拡張

さまざまな推測方式に応じて規準を構成する

罰金付き最尤法, ベイズ予測分布, 不完全データ, 実験計画

C2

1.1 回帰モデルの変数選択問題 (例: 説明変数が3個)

$$Y \sim N(\mu(X_a, X_b, X_c), \sigma^2)$$
$$\mu(x_a, x_b, x_c) = \beta_0 + \beta_a x_a + \beta_b x_b + \beta_c x_c$$

変数選択問題

$$x_a, x_b, x_c$$

の中から適切な変数だけを選ぶ.

モデル α は選択する変数の集合を表す

$$\alpha \subset \{a, b, c\}$$

パラメタのとり得る値とパラメタ数は

$$\Theta_\alpha = \{(\sigma, \beta) | \sigma > 0, \beta_i = 0, i \in \alpha^c\}, \quad |\alpha| + 2$$

C4

- 候補となるモデル α の集合を \mathcal{M} で表す .

$$\alpha \in \mathcal{M}$$

[例] フルモデル $\{a, b, c\}$

すべての組み合わせ $\mathcal{M} = \{\phi, \{a\}, \{b\}, \{c\}, \{a, b\}, \{a, c\}, \{b, c\}, \{a, b, c\}\}$

2個の変数を使う $\mathcal{M} = \{\{a, b\}, \{a, c\}, \{b, c\}\}$

- 候補モデルに番号をつける

$$\alpha(k), \quad k \in \{1, 2, \dots, K\}$$

各モデルは集合 α または番号 k で示される

- パラメトリック・モデル

$$M_\alpha = f(\cdot | \Theta_\alpha), \quad \alpha \in \mathcal{M}$$

$$M_1, M_2, \dots, M_K$$

C5

モデルの包含関係

二つのモデル M_α と M_β を考える

- $\alpha \subset \beta$ のとき , M_α は M_β の部分モデル

このとき M_α は M_β に「ネスト」していると言う

[例] $\phi \subset \{a\} \subset \{a, b\} \subset \{a, b, c\}$

- $\alpha \subset \beta$ でも $\beta \subset \alpha$ でも無いとき , M_α と M_β は「ノンネスト」と言う

[例] $\{a, b\}$ と $\{b, c\}$ は互いにネストの関係に無い (ノンネスト)

C6

モデル選択とモデル信頼集合

- 期待平均対数尤度

$$l_n^*(\alpha) \quad \text{または} \quad l_n^*(k)$$

を最小にする α^* または k^* が「良いモデル」と考える .

- バイアス補正した対数尤度

$$l(\hat{\theta}_\alpha) - \dim M_\alpha \quad \text{または} \quad l(\hat{\theta}_k) - \dim M_k$$

をおおきくするモデル $\hat{\alpha}$ または \hat{k} を選択する .

モデル信頼集合 \mathcal{T} は , \mathcal{M} の部分集合 ($\mathcal{T} \subset \mathcal{M}$) で ,

$$P\{\alpha^* \in \mathcal{T}\} \quad \text{または} \quad P\{k^* \in \mathcal{T}\} \geq 1 - P^*$$

を満たすもの . ただし P^* は有意水準

C7

数値例

HALD のセメントデータ

目的変数: セメント発熱量

説明変数: 4種の主成分の混合率

データ数: $n = 13$

モデル数: $K = 16$ (すべての説明変数の組み合わせ $2^4 = 16$ 個)

BOSTON の住宅価格データ

目的変数: 住宅価格 (各地区の中央値の対数)

説明変数: 各地区の13個の指標

a=犯罪率, f=平均部屋数, k=教師数:生徒数の比, m=社会階層の構成比

データ数: $n = 506$

モデル数: $K = 286$ (説明変数を3個だけ使う)

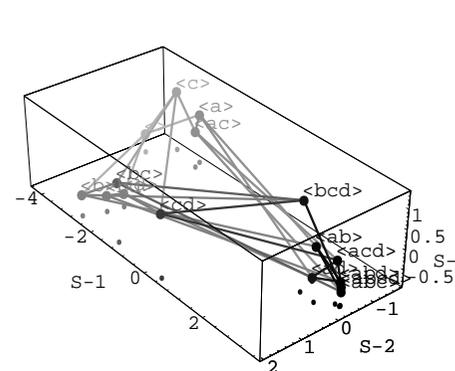
C8

HALDデータのAICとp-値

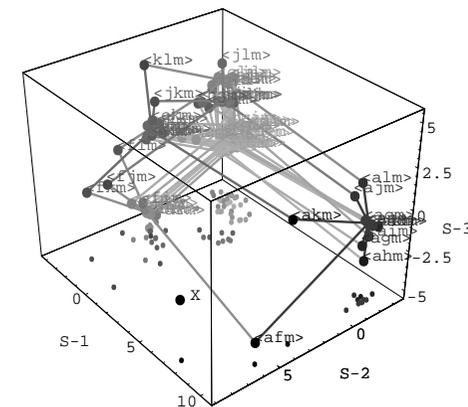
順位	α	ΔAIC	正規検定	多重比較	選択確率	LR検定
1	<abd>	0	—	0.972	0.180	0.863
2	<abc>	0.04	0.493	0.938	0.256	0.796
3	<ab>	0.45	0.438	0.911	0.255	0.290
4	<acd>	0.75	0.382	0.924	0.158	0.376
5	<abcd>	1.97	0.088	0.451	0.002	—
6	<ad>	3.77	0.169	0.540	0.094	0.055
7	<bcd>	5.60	0.136	0.353	0.052	0.018
8	<cd>	14.88	0.019	0.074	0.004	0.000
9	<bc>	26.06	0.001	0.000	0.000	0.000
10	<d>	33.88	0.000	0.000	0.000	0.000
11		34.20	0.000	0.000	0.000	0.000
12	<bd>	35.66	0.000	0.000	0.000	0.000
13	<a>	38.55	0.000	0.000	0.000	0.000
14	<ac>	40.14	0.000	0.000	0.000	0.000
15	<c>	44.09	0.000	0.000	0.000	0.000
16	<>	46.47	0.000	0.000	0.000	0.000

C9

モデル地図



HALDデータ



BOSTONデータ

C11

BOSTONデータにおける上位20個のモデルのAICとp-値

順位	α	LK	p-value × 100					LR
			BA	BP	KH	MC	MS	
1	<afm>	0	70.3	48.8	—	99.3	99.4	0.0
2	<akm>	0.9	29.7	39.5	46.3	98.8	98.7	0.0
3	<ahm>	9.3	0.0	7.3	18.0	86.2	77.1	0.0
4	<agm>	12.4	0.0	0.5	9.0	79.8	53.9	0.0
5	<adm>	12.7	0.0	1.5	11.1	80.7	55.1	0.0
6	<alm>	16.1	0.0	0.0	5.0	74.0	22.1	0.0
7	<ajm>	17.7	0.0	0.0	2.4	71.2	4.6	0.0
8	<abm>	19.2	0.0	0.0	1.9	66.6	10.0	0.0
9	<aim>	19.8	0.0	0.0	1.4	65.8	8.3	0.0
10	<acm>	20.4	0.0	0.0	1.0	64.3	3.0	0.0
11	<aem>	20.5	0.0	0.0	1.1	63.8	3.5	0.0
12	<klm>	23.4	0.0	1.4	8.3	55.9	41.5	0.0
13	<fjm>	24.4	0.0	0.2	2.9	55.7	26.0	0.0
14	<jkm>	27.0	0.0	0.0	3.2	47.9	12.5	0.0
15	<fkm>	27.1	0.0	0.1	3.4	47.9	31.4	0.0
16	<flm>	29.9	0.0	0.4	3.7	40.6	27.7	0.0
17	<hjm>	30.3	0.0	0.1	2.6	40.1	20.5	0.0
18	<dkm>	31.0	0.0	0.0	2.8	38.1	15.7	0.0
19	<bkm>	31.3	0.0	0.0	2.3	37.2	13.0	0.0
20	<ekm>	32.4	0.0	0.0	1.7	34.7	6.6	0.0

LK: $\Delta \log L$
 BA: ベイズ事後確率
 BP: 選択確率
 KH: 正規検定
 MC: 多重比較 (おもみ無し)
 MS: 多重比較
 LR: 尤度比検定

C10

分子系統樹

ミトコンドリアのDNAシーケンス

```

1                               h=20                               50
human  ANLLLLIVPILIAMAFMLMIFERKILGYMQLRKGPNVVGPYGLLQPFADAMKLFYKPEPLKP
seal   INIISLIIPILLAVAFLLTVERKVLGYMQLRKGPNI VGPYGLLQPIADAVKLFYKPEPLRP
cow    INILMLIIPILLAVAFLLTVERKVLGYMQLRKGPNVVGPYGLLQPIADAIKLFYKPEPLRP
rabbit INTLLLILPVLLAMAFLLTVERKILGYMQLRKGPNI VGPYGLLQPIADAIKLFYKPEPLRP
mouse  INILTLVLPILIAMAFLLTVERKILGYMQLRKGPNI VGPYGLLQPFADAMKLFYKPEPLRP
opossum INLLMYIIPILLAVAFLLTVERKVLGYMQFRKGPNI VGPYGLLQPFADALKLFYKPEPLRP
    
```

シーケンス数: 6種の哺乳類

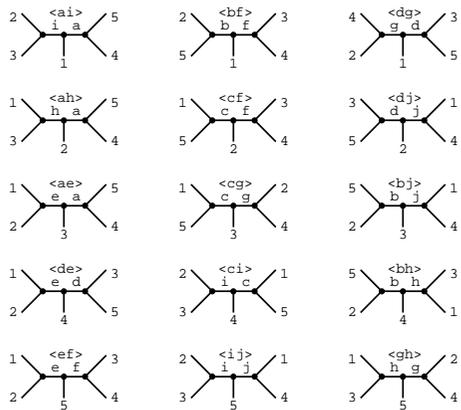
human, (harbor seal, cow), rabbit, mouse, opossum

説明変数: 10個の "split"

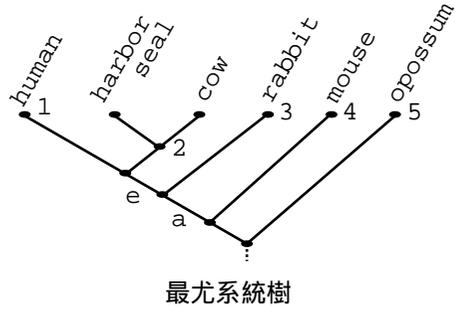
データ数: $n = 3416$ のアミノ酸

モデル数: $K = 15$ の系統樹

C12



可能な系統樹



a = {123|45}, b = {134|25}, c = {234|15}, d = {124|35}, e = {12|345},
 f = {34|125}, g = {24|135}, h = {13|245}, i = {23|145}, j = {14|235}.

C13

α	$\Delta \log L$	$\hat{\theta}_\alpha$
full	95.2	(6.14, 2.32, 2.78, 2.28, 2.24, 2.94, 1.21, 1.90, 2.13, 2.23)
{aefi}	61.3	(5.27, , , , 1.70, 3.07, , , 2.43,)
{ae}	38.5	(5.86, , , , 2.23, , , , ,)
{bcef}	34.6	(, 2.17, 3.05, , 2.51, 1.43, , , ,)
{a}	28.4	(6.18, , , , , , , , ,)
{ef}	21.0	(, , , , , 1.88, 2.31, , , ,)
{e}	12.7	(, , , , , 2.52, , , , ,)
{}	0	(, , , , , , , , ,)

C15

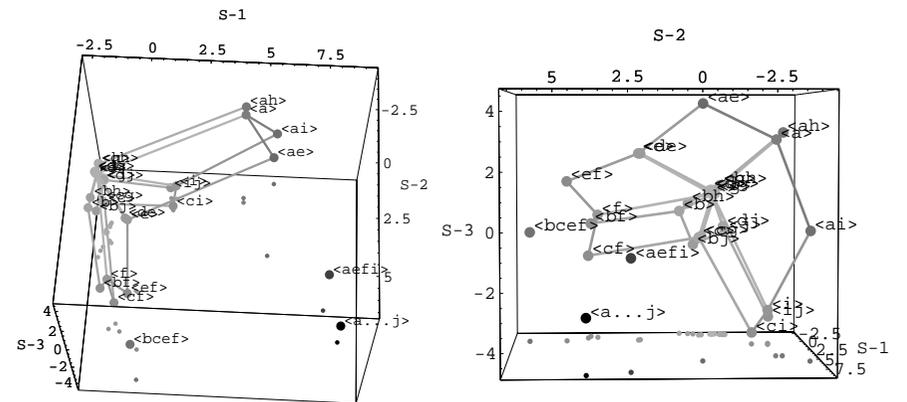
Mammal Phylogeny: *p*-values

α	$\Delta \log L$	<i>p</i> -values								tree
		LR	BA	BP	KH	MC	MS	MB		
1	{a, e}	0.0	.000	.934	.583	-	.941	.944	.85	(((12)3)4)5
2	{a, i}	2.7	.000	.065	.317	.360	.811	.805	.57	((1(23))4)5
3	{a, h}	7.4	.000	.001	.038	.121	.577	.422	.27	(((13)2)4)5
4	{e, f}	17.6	.000	.000	.012	.040	.169	.203	.15	((12)(34)5)
5	{c, f}	18.9	.000	.000	.030	.066	.139	.296	.18	(1(2(34))5)
6	{c, i}	20.1	.000	.000	.006	.050	.109	.100	.09	(1((23)4)5)
7	{b, f}	20.6	.000	.000	.011	.048	.107	.248	.19	((1(34))2)5
8	{i, j}	22.2	.000	.000	.001	.032	.070	.048	.08	((14)(23)5)
9	{d, e}	25.4	.000	.000	.000	.001	.029	.013	.01	(((12)4)3)5
10	{b, j}	26.3	.000	.000	.002	.018	.032	.124	.09	(((14)3)2)5
11	{b, h}	28.9	.000	.000	.000	.008	.017	.069	.08	(((13)4)2)5
12	{d, j}	31.6	.000	.000	.000	.003	.006	.032	.04	(((14)2)3)5
13	{c, g}	31.7	.000	.000	.000	.003	.006	.035	.04	(1((24)3)5)
14	{g, h}	34.7	.000	.000	.000	.001	.002	.012	.03	((13)(24)5)
15	{d, g}	36.2	.000	.000	.000	.000	.001	.007	.02	((1(24))3)5

LR:尤度比検定, BA:ベイズ事後確率, BP:選択確率, KH:正規検定,
 MC:多重比較(おもみ無し), MS:多重比較, MB:New Test

C14

モデル地図



C16

1.2 モデル選択のバラツキと一致性

- 選択されたモデル \hat{k} はデータ x の関数

$$\hat{k} = \hat{k}(x)$$

したがって $\hat{k}(X)$ は確率変数

- $\hat{k}(X)$ は k^* の周辺に分布している

- モデル選択の一致性

$$P(\hat{k}(X) = k^*) \rightarrow 1 \quad (n \rightarrow \infty)$$

AICによるモデル選択は、一致性がない場合があることが知られている

これは必ずしも AIC の欠点ではない

C17

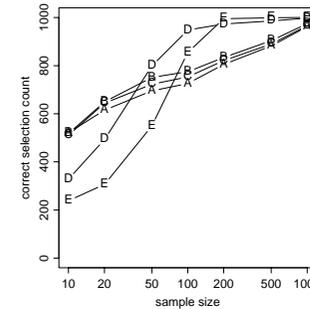
モデル選択の一致性 (変数選択の例)

シミュレーション: 一番良いモデルが選ばれた回数 (1000回中)

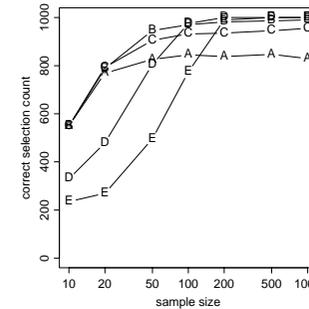
データサイズ $n = 10, 20, 50, 100, 200, 500, 1000$

$$IC_\alpha = -2 \times l(\hat{\theta}_\alpha) + c_n \dim M_\alpha$$

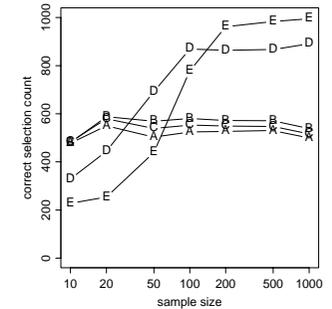
係数 c_n : $A = 2$ (AIC), $B = \log n$ (BIC), $C = 2n^{0.1}$, $D = 2n^{0.5}$, $E = 2n^{0.6}$



ケース1



ケース2



ケース3

C19

モデル選択のバラツキ (教科書のフーリエ級数の例)

表 4.6 AIC 最小となる k ($=\hat{k}$) の分布 ($n = 500$)

k_A	2	4	6	8	10	12	14	16	18	20	22	計
度数	0	0	0	326	342	173	91	27	17	14	10	1000

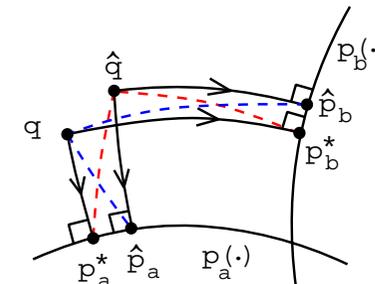
表 4.7 $n = 2000$ のときの \hat{k} の分布

k_A	2	4	6	8	10	12	14	16	18	20	22	計
度数	0	0	0	6	210	359	264	68	58	21	14	1000

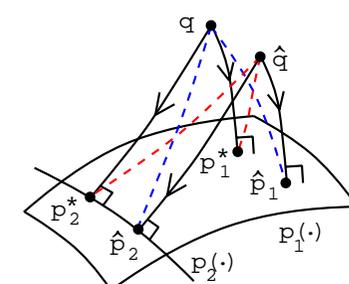
なお, $l_{500}^*(k)$ 最小は $k^* = 10$, $l_{2000}^*(k)$ 最小は $k^* = 12$ のとき

C18

モデルと真の分布の関係 (一般の場合)



ケース1 (ノンネスト)



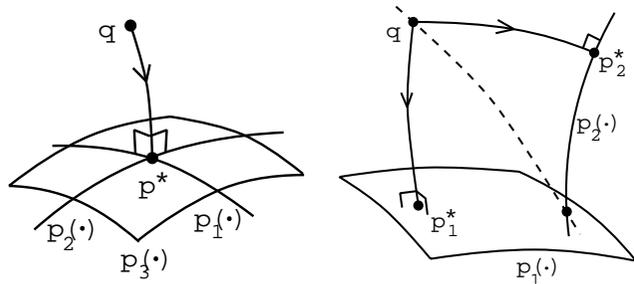
ケース1' (ネスト)

数値例のパラメタ値 $\beta = (\beta_a, \beta_b, \beta_c)$

ケース1: $\beta = (1.0, 0.9, 0.1)$, $\mathcal{M} = \{\phi, \{a\}, \{b, c\}\}$

C20

モデルと真の分布の関係 (特殊な場合)



ケース2

ケース3

ケース2: $\beta = (1, 1, 0)$, $\mathcal{M} = \{\phi, \{a\}, \{a, c\}\}$

ケース3: $\beta = (1, 1, 0)$, $\mathcal{M} = \{\phi, \{a\}, \{b, c\}\}$

C21

ICの差の第1項

$$l(\theta_\alpha^*) - l(\theta_\beta^*) = \sum_{i=1}^n (\log f(x_i | \theta_\alpha^*) - \log f(x_i | \theta_\beta^*))$$

中心極限定理より, 十分大きな n で

$$l(\theta_\alpha^*) - l(\theta_\beta^*) \sim N(l^*(\theta_\alpha^*) - l^*(\theta_\beta^*), nJ_{\alpha,\beta})$$

$$l^*(\theta_\alpha^*) - l^*(\theta_\beta^*) = n \times (I(g; f(\theta_\beta^*)) - I(g; f(\theta_\alpha^*)))$$

$$\begin{aligned} J_{\alpha,\beta} &= V[\log f(X | \theta_\alpha^*) - \log f(X | \theta_\beta^*)] \\ &\approx I(f(\theta_\alpha^*), f(\theta_\beta^*)) + I(f(\theta_\beta^*), f(\theta_\alpha^*)) \end{aligned}$$

分散の推定は

$$\hat{J}_{\alpha,\beta} = \frac{1}{n} \sum_{i=1}^n (\log f(x_i | \hat{\theta}_\alpha) - \log f(x_i | \hat{\theta}_\beta) - (l(\hat{\theta}_\alpha) - l(\hat{\theta}_\beta)) / n)^2$$

C23

ICの差

二つのモデル M_α と M_β

$$\Delta IC = 2 \times (l(\hat{\theta}_\alpha) - l(\hat{\theta}_\beta)) - c_n \times (\dim M_\alpha - \dim M_\beta)$$

$\Delta IC > 0$ なら M_α を選択, $\Delta IC < 0$ なら M_β を選択

$$\begin{aligned} \Delta IC &\approx 2 \times (l(\theta_\alpha^*) - l(\theta_\beta^*)) \\ &\quad + (\Delta(\theta_\alpha^*, \hat{\theta}_\alpha) - \Delta(\theta_\beta^*, \hat{\theta}_\beta)) - c_n \times (\dim M_\alpha - \dim M_\beta) \end{aligned}$$

ただし対数尤度の $\hat{\theta}_\alpha$ の周りでのテーラ展開は (θ を M_α の自由パラメタに制限して)

$$\begin{aligned} l(\theta) &\approx l(\hat{\theta}_\alpha) + (\theta - \hat{\theta}_\alpha) \left[\frac{\partial l}{\partial \theta'} \right]_{\hat{\theta}_\alpha} \\ &\quad + \frac{1}{2} (\theta - \hat{\theta}_\alpha) \left[\frac{\partial^2 l}{\partial \theta' \partial \theta} \right]_{\hat{\theta}_\alpha} (\theta - \hat{\theta}_\alpha)' \end{aligned}$$

これより

$$l(\hat{\theta}_\alpha) \approx l(\theta_\alpha^*) + \frac{1}{2} \Delta(\theta_\alpha^*, \hat{\theta}_\alpha)$$

C22

ICの差 ケース1

十分大きな n では第1項が支配的

$$\frac{1}{2} \Delta IC \approx (l(\theta_\alpha^*) - l(\theta_\beta^*)) \sim N(l^*(\theta_\alpha^*) - l^*(\theta_\beta^*), nJ_{\alpha,\beta})$$

$I(g; f(\theta_\beta^*)) - I(g; f(\theta_\alpha^*)) > 0$ なら

$$P(\Delta IC > 0) \rightarrow 1$$

$I(g; f(\theta_\beta^*)) - I(g; f(\theta_\alpha^*)) < 0$ なら

$$P(\Delta IC < 0) \rightarrow 1$$

もし $I(g; f(\theta_\beta^*)) - I(g; f(\theta_\alpha^*)) = 0$ ならば, ケース2かケース3に相当

C24

ICの差 ケース2

$f(\cdot|\theta_\alpha^*) = f(\cdot|\theta_\beta^*)$ のとき ,

$$J_{\alpha,\beta} = V \left[\log f(X|\theta_\alpha^*) - \log f(X|\theta_\beta^*) \right] = 0$$

$$I(g; f(\theta_\beta^*)) - I(g; f(\theta_\alpha^*)) = 0$$

xなので ΔIC の第1項は0になる

$$\Delta IC \approx \left(\Delta(\theta_\alpha^*, \hat{\theta}_\alpha) - \Delta(\theta_\beta^*, \hat{\theta}_\beta) \right) - c_n \times (\dim M_\alpha - \dim M_\beta)$$

$$\Delta(\theta_\alpha^*, \hat{\theta}_\alpha) \sim \chi_{\dim M_\alpha}^2, \quad \Delta(\theta_\beta^*, \hat{\theta}_\beta) \sim \chi_{\dim M_\beta}^2$$

もし

$$d = \dim M_\alpha - \dim M_\beta > 0$$

ならば M_β の方がよいモデル . もし $c_n \rightarrow \infty$ ならば M_β を選択する確率は

$$P(\Delta IC < 0) \rightarrow 1$$

C25

ICの差 ケース3

$I(g; f(\theta_\beta^*)) - I(g; f(\theta_\alpha^*)) = 0$ かつ $f(\cdot|\theta_\alpha^*) \neq f(\cdot|\theta_\beta^*)$ のとき

$$\Delta IC \approx 2 \times \left(l(\theta_\alpha^*) - l(\theta_\beta^*) \right) - c_n \times (\dim M_\alpha - \dim M_\beta)$$

でケース1とほぼ同じなのだが ,

$$l^*(\theta_\alpha^*) - l^*(\theta_\beta^*) = 0$$

なので

$$l(\theta_\alpha^*) - l(\theta_\beta^*) \sim N \left(0, nJ_{\alpha,\beta} \right)$$

もし

$$d = \dim M_\alpha - \dim M_\beta > 0$$

ならば M_β の方がよいモデル . もし $c_n/n^{\frac{1}{2}} \rightarrow \infty$ ならば M_β を選択する確率は

$$P(\Delta IC < 0) \rightarrow 1$$

C27

ICの差 ケース2 (ネスト)

もしネスト ($M_\beta \subset M_\alpha$) ならば

$$\Delta IC \sim \chi_d^2 - c_n d, \quad d = \dim M_\alpha - \dim M_\beta$$

M_β を選ぶ確率は

$$P(\Delta IC < 0) = P(\chi_d^2 < c_n d)$$

$c_n = 2, d = 1$ なら $P(\Delta IC < 0) \approx 0.84$.

$c_n = 2, d = 10$ なら $P(\Delta IC < 0) \approx 0.97$.

固定した d について $c_n \rightarrow \infty$ のとき $P(\Delta IC < 0) \rightarrow 1$.

C26

モデル選択が一致性を持つための条件

モデル選択規準として

$$IC_\alpha = -2 \times l(\hat{\theta}_\alpha) + c_n \dim M_\alpha$$

を使うことにする .

ケース1

$$c_n/n \rightarrow 0$$

ケース2

$$c_n \rightarrow \infty$$

ケース3

$$c_n/n^{\frac{1}{2}} \rightarrow \infty$$

AIC ($c_n = 2$) はケース1で一致性がある

BIC ($c_n = \log n$) はケース1と2で一致性がある

C28

1.3 モデル選択の信頼性

モデル選択の一致性の議論は実用上あまり意味が無い

- 一般的な状況(ケース1)ではどのモデル選択規準も一致性をもつ
- 有限の n ではど規準を使ってもモデル選択のバラツキがある

モデル選択のバラツキを定量的に評価して、選択の信頼性を表す数値で示す

P_1, P_2, \dots, P_K ; 各モデルの p -value, p -値

モデルをひとつだけ選ぶのではなく、同程度に良いモデルをすべて選び出す

モデル信頼集合

C29

確率値 (p -value)

各モデル M_k , $k = 1, \dots, K$ に $[0, 1]$ の範囲の数値

$$P_k$$

を与え、それが 1 に近いほど良いモデルである可能性を表す

いろいろな考え方があ

- モデルが選択される確率の推定 (ブートストラップ選択確率)
- モデル信頼集合 (それに対応する検定の確率値 p -value)

C30

リサンプリング

データ

$$\mathbf{x} = \{x_1, \dots, x_n\}$$

選択されたモデル

$$\hat{k}(\mathbf{x})$$

データからのリサンプリング (各 \tilde{x}_i は \mathbf{x} からランダムに選ぶ)

$$\tilde{\mathbf{x}} = \{\tilde{x}_1, \dots, \tilde{x}_n\}$$

選択されるモデル

$$\hat{k}(\tilde{\mathbf{x}})$$

C31

ブートストラップ選択確率

リサンプリングを多数繰り返す (例えば $B = 1000$)

$$\tilde{\mathbf{x}}[1], \tilde{\mathbf{x}}[2], \dots, \tilde{\mathbf{x}}[B]$$

選択されるモデル

$$\hat{k}(\tilde{\mathbf{x}}[1]), \hat{k}(\tilde{\mathbf{x}}[2]), \dots, \hat{k}(\tilde{\mathbf{x}}[B])$$

モデル M_k のブートストラップ選択確率

$$P_k = \frac{\#\{\hat{k}(\tilde{\mathbf{x}}[b]) = k, b = 1, \dots, B\}}{B}, \quad k = 1, \dots, K$$

$$\sum_{k=1}^K P_k = 1, \quad 0 \leq P_k \leq 1$$

すなわち M_k が選択された頻度 P_k が大きいモデルほど良いモデルと考えられる。

C32

モデル選択と仮説検定

- パラメトリック・モデル

$$M_\alpha = f(\cdot | \Theta_\alpha), \quad \alpha \in \mathcal{M}$$

$$M_1, M_2, \dots, M_K$$

- モデル選択 \hat{k} では、どの M_k が真の分布に近いかを推定する。

相対的にどのモデルが良いか？

- 仮説検定では、どの M_k が真の分布を含むかを調べる。

どのモデルが正しいか？

C33

モデル選択，仮説検定，モデル選択の検定

これらは互いに関連しているが，同じものではない。

- モデル選択 — どのモデルが（相対的に）良いかの推定
- 仮説検定 — どのモデルが正しいかの検定
- モデル選択の検定 — どのモデルが（相対的に）良いかの検定

モデルを比較する目的はなにか

良いモデルを選択して予測精度をあげる？ — モデル選択よりモデル混合

モデル選択を通して対象の理解を深める？ — 解釈可能なモデルを利用

C35

モデル選択の検定

- モデル M_k が M_1, \dots, M_K の中で一番良いという仮説を

$$H_k : k^* = k$$

であらわす． ($l_n^*(k) = l_n^*(k^*)$ となるタイの k すべてで H_k がいえるとする．)

- モデル選択の検定

$$H_1, H_2, \dots, H_K$$

のうち、どの仮説が正しいか？

- モデル信頼集合

$$\mathcal{T} = \{k \in \{1, \dots, K\} \mid H_k \text{ が有意水準 } P^* \text{ で棄却されない}\}$$

C34

AICの差の有意性

二つのモデル M_α と M_β

$$\Delta \text{AIC} = 2 \times (l(\hat{\theta}_\alpha) - l(\hat{\theta}_\beta)) - 2 \times (\dim M_\alpha - \dim M_\beta)$$

$\Delta \text{AIC} > 0$ なら M_α を選択， $\Delta \text{AIC} < 0$ なら M_β を選択

AICの差の有意性を検定する — 統計 (Linhart 1988)，計量経済 (Vuong 1989)，分子生物 (Kishino-Hasegawa 1989) などによって提案

$l_n^*(\alpha) = l_n^*(\beta)$ かつ $I(f(\theta_\alpha^*); f(\theta_\beta^*)) > 0$ のとき，十分大きな n に対して

$$\frac{\Delta \text{AIC}}{\sqrt{\hat{V}\{\Delta \text{AIC}\}}} \sim N(0, 1)$$

ネストの場合はあまり良い近似ではないが，ノンネストの場合に有効。

分散の推定は

$$\hat{V}\{\Delta \text{AIC}\} = 4n\hat{J}_{\alpha,\beta} + v_{\alpha,\beta}$$

C36

モデル選択の検定 (二つのモデル)

$$AIC_\beta - AIC_\alpha > 0$$

のとき M_α が選択されるが,

$$\frac{AIC_\beta - AIC_\alpha}{\sqrt{\hat{V}\{AIC_\beta - AIC_\alpha\}}} < c$$

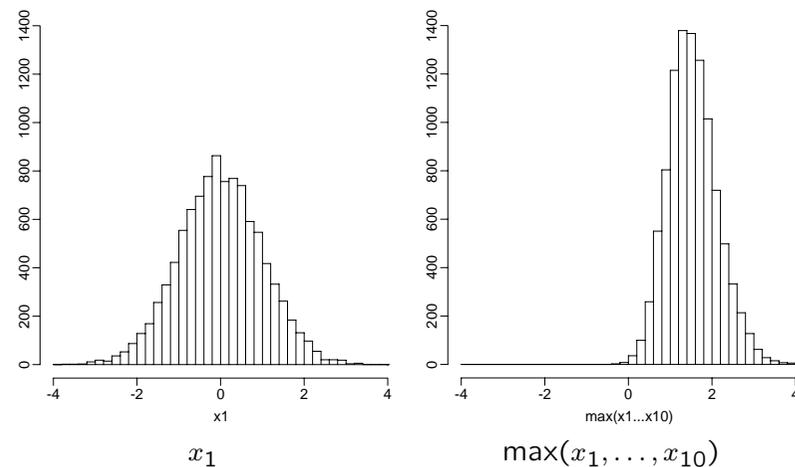
ならば M_α と M_β の良さの差は有意ではないと判断される。ただし c は有意水準 P^* より

$$\Phi(c) = 1 - P^*$$

によって定める。

C37

選択バイアス



$$(x_1, \dots, x_{10}) \sim N(0, I) \quad 10000 \text{ samples}$$

C39

モデル選択の検定 ($K \geq 3$) — 選択バイアスを無視した場合

K 個のモデルを **AIC** の小さい順に並べて,

$$M_1, M_2, \dots, M_K$$

$$AIC_1 \leq AIC_2 \leq \dots \leq AIC_K$$

このうち,

$$\frac{AIC_k - AIC_1}{\sqrt{\hat{V}\{AIC_k - AIC_1\}}} < c$$

を満たす M_k は, M_1 とのモデルの良さの差が有意でないと考えられる。

実際には, $AIC_k - AIC_1$ は選択バイアスによって大きくなる傾向がある。

$$AIC_k - AIC_1 = \max_{k'=1, \dots, K} (AIC_k - AIC_{k'})$$

C38

多重比較による p -値の計算

$$aic = (AIC_1, AIC_2, \dots, AIC_K) \sim N(E\{aic\}, V\{aic\})$$

と近似的にみなし, $E\{AIC_k\}$ を最小にする k^* の信頼集合を「多重比較」で構成する。

$$\delta_k(aic) = \max_{k'=1, \dots, k-1, k+1, \dots, K} \frac{AIC_k - AIC_{k'}}{\sqrt{\hat{V}\{AIC_k - AIC_{k'}\}}}$$

(1) リサンプリングを B 回行い, 各モデルの **AIC** の複製を計算する。

$$aic(\tilde{x}[b]), \quad b = 1, \dots, B$$

(2) **AIC** の平均 $\hat{E}\{aic\}$ を求め, それを引いたうえで δ_k の複製を計算する。

$$\tilde{\delta}_k[b] = \delta_k(aic(\tilde{x}[b]) - \hat{E}\{aic\}), \quad b = 1, \dots, B$$

(3) 各モデル M_k 毎に, H_k の検定に関する p -値を計算する。

$$P_k = \frac{\#\{\tilde{\delta}_k[b] > \delta_k(aic), \quad b = 1, \dots, B\}}{B}, \quad k = 1, 2, \dots, K$$

C40

H_k の検定

あらかじめ与えた有意水準 P^* に対して,

$P_k < P^*$ なら H_k を棄却する

$P_k \geq P^*$ なら H_k を棄却しない

ただし, $H_k: \delta_k(E\{aic\}) \leq 0$, もしくは

$$H_k: E(AIC_1) \geq E(AIC_k), \dots, E(AIC_K) \geq E(AIC_k)$$

である. もし H_k が真ならば, H_k が誤って棄却される確率は

$$P(P_k < P^*) \leq P^*$$

等号は

$$E(AIC_1) = E(AIC_2) = \dots = E(AIC_K)$$

のとき.

C41

モデル信頼集合の構成

モデル信頼集合は

$$\mathcal{T} = \{k \in \{1, 2, \dots, K\} \mid P_k \geq P^*\}$$

これが k^* を含む確率 (被覆確率) は

$$P(k^* \in \mathcal{T}) \geq 1 - P^*$$

(証明)

$$k^* \in \mathcal{T} \Leftrightarrow P_{k^*} \geq P^* \Leftrightarrow H_{k^*} \text{が棄却されない}$$

\mathcal{T} に含まれるモデルは, \hat{k} より有意に悪いとは言えない
 \mathcal{T} に含まれないモデルは, \hat{k} より有意に悪いとは言える

下平 (1993), Shimodaira (1998) など

C42

確率分布のベクトル表現

モデル $f(\cdot|\theta)$ を n 次元ベクトル

$$\xi(\theta) = (\log f(x_1|\theta), \dots, \log f(x_n|\theta))$$

で表す.

各モデル M_k は $\dim M_k$ 次元の集合

$$\{\xi(\theta) \mid \theta \in \Theta_k\}$$

最尤モデルは

$$\xi(\hat{\theta}_k)$$

で表される.

実際の描画では主成分分析 (PCA) 等により, 低次元に射影する

C43

2 情報量規準の拡張

真の分布: $X_1, X_2, \dots, X_n \sim g(\cdot)$

モデル: $X_1, X_2, \dots, X_n \sim f(\cdot|\theta)$

データ: $\mathbf{x} = (x_1, x_2, \dots, x_n)$

予測分布: $f(\cdot|\hat{\theta}(\mathbf{x}))$

$f(\cdot|\hat{\theta}(\mathbf{x}))$ が平均的にどれだけ $g(\cdot)$ に近いか?

$$E \{I(g(\cdot); f(\cdot|\hat{\theta}(\mathbf{X})))\}$$

の推定量としての情報量規準

データから予測分布を与える方式 (手続き) には, さまざまな形式が考えられる

$$f(\cdot|\mathbf{x})$$

の平均的な良さ

$$E \{I(g(\cdot); f(\cdot|\mathbf{X}))\}$$

の推定量は, 情報量規準の拡張になる.

C44

色々な情報量規準

- パラメタ推定による予測分布の場合
AIC, TIC, GIC, EIC, RIC (罰金付き最尤法) など
- ベイズ予測分布の場合
 曲率との関係
- 不完全データの場合
EM アルゴリズムとの関係, **PDIO**
- 説明変数の分布が変化する場合
 実験計画, 重み付き最尤法
- ベイズ的方法

C45

M-推定量

$$\sum_{i=1}^n \psi(x_i | \hat{\theta}) = 0$$

$$T^{(1)}(x|g) = M(\psi|g)^{-1} \psi(x|\theta^*); \quad M(\psi|g) = -E \left\{ \frac{\partial \psi(X|\theta)}{\partial \theta'} \Big|_{\theta^*} \right\}$$

$$GIC = -l(\hat{\theta}) + \text{tr} \left(E \left\{ \psi(X|\theta^*) \frac{\partial \log p(X|\theta)}{\partial \theta'} \Big|_{\theta^*} \right\} \cdot M(\psi|g)^{-1} \right)$$

罰金付き最尤法

$$l(\theta) - \lambda k(\theta)$$

を最大にする推定量で, λ をどうやって選ぶか?

$$\psi(x|\theta) = \frac{\partial \log f(x|\theta)}{\partial \theta} - \lambda \frac{\partial k(\theta)}{\partial \theta}$$

C47

2.1 パラメタ推定による予測分布の場合

分布 $g(\cdot)$ の空間から Θ への汎関数

$$T(g(\cdot))$$

を考える. 特に, モデルが真の分布を含む場合は,

$$T(f(\cdot|\theta^*)) = \theta^*$$

とする. パラメタの推定量として

$$T(\hat{g}(\cdot))$$

を用いる.

$$GIC = -l(T(\hat{g})) + E \left\{ \frac{\partial \log f(X|\theta)}{\partial \theta'} \Big|_{\theta^*} T^{(1)}(X|g) \right\}$$

解析的に解く代わりに, ブートストラップを使って数値的に計算する方法がある (**EIC**).

C46

2.2 ベイズ予測分布の場合

パラメタの事前分布: $f(\theta)$

パラメタの事後分布: $f(\theta|x) \propto f(x|\theta)f(\theta)$

$$f(\theta|x) = \frac{f(x|\theta)f(\theta)}{\int f(x|\theta)f(\theta) d\theta}$$

ベイズ予測分布 $f(\cdot|x)$ は

$$f(z|x) = \int f(z|\theta)f(\theta|x) d\theta$$

$$IC = - \sum_{i=1}^n \log f(x_i|x) + \text{tr}(GH^{-1})$$

C48

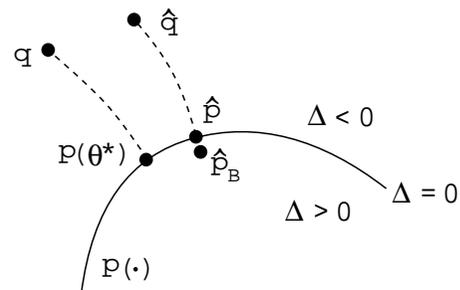
[参考] ベイズ予測分布と最尤推定の関係

ベイズ予測分布のICは書き直せて、

$$IC \approx - \sum_{i=1}^n \log f(x_i|\hat{\theta}) + \frac{1}{2} (\text{tr}(GH^{-1}) + \dim \theta)$$

一方、最尤推定のICは

$$IC = - \sum_{i=1}^n \log f(x_i|\hat{\theta}) + \text{tr}(GH^{-1})$$



したがってこれらの差は $\Delta/2$, ただし
 $\Delta = \text{tr}(GH^{-1}) - \dim \theta$

C49

2.3 不完全データの場合

完全データ $x = (y, z)$ のうち y しか観測出来ず、 z が観測できない確率変数

完全データのモデル: $f(x|\theta)$
 観測データのモデル: $f(y|\theta)$

$$f(y|\theta) = \int f(y, z|\theta) dz$$

最尤推定 $\hat{\theta}$ は次の対数尤度から求める (例えばEMアルゴリズム等を使って)

$$l_Y(\theta) = \sum_{i=1}^n \log f(y_i|\theta)$$

AIC は $f(y|\hat{\theta})$ の $g(y)$ からの平均的な隔たりを測る

$$AIC = -l_Y(\hat{\theta}) + \dim \theta$$

ただし

$$g(y) = \int g(y, z) dz$$

C50

完全データに関するモデルの良さを測る規準

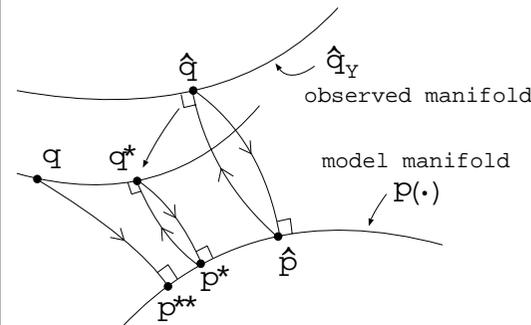
モデル $f(x|\hat{\theta})$ の真の分布 $g(x)$ からの平均的な隔たりを測るには

$$IC = -l_Y(\hat{\theta}) + \text{tr}(I_X I_Y^{-1})$$

が使える.

$$I_X(\theta) = - \int f(x|\theta) \frac{\partial^2 \log f(x|\theta)}{\partial \theta \partial \theta'} dx, \quad I_Y(\theta) = - \int f(y|\theta) \frac{\partial^2 \log f(y|\theta)}{\partial \theta \partial \theta'} dy$$

はそれぞれ x と y に関する θ の Fisher 情報行列.



情報量の差

$$I_{Z|Y} = I_X - I_Y$$

とおくと、

$$\text{tr}(I_X I_Y^{-1}) = \dim \theta + \text{tr}(I_{Z|Y} I_Y^{-1})$$

C51

2.4 説明変数の分布が変化する場合

説明変数 z の分布がデータを観測した時と予測分布の良さを評価する時で異なる場合

回帰モデル: $f(y|z, \theta)$
 目的変数: y
 説明変数: z

$$z \sim g_0(z) \text{ データ}, \quad z \sim g_1(z) \text{ 評価}$$

重み付き対数尤度

$$l_w(\theta) = \sum_{i=1}^n w(z_i) \log f(y_t|z_t, \theta)$$

これを最大にする重み付き最尤推定を $\hat{\theta}_w$ とかく. 予測分布

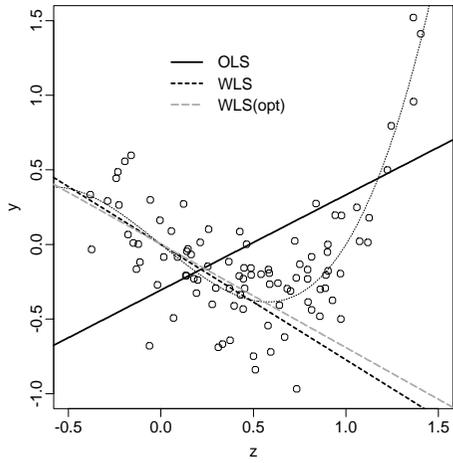
$$f(y|z, \hat{\theta}_w) g_1(z)$$

がどれだけ $g(y|z) g_1(z)$ に近いかを評価する.

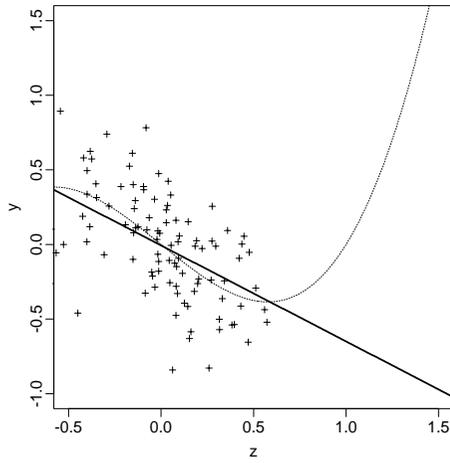
C52

数値例：多項式回帰モデル

$$y = -z + z^3 + \epsilon, \quad \epsilon \sim N(0, 0.3^2)$$



$z \sim N(0.5, 0.5^2)$
 $\lambda = 0, 0.77, 1$



$z \sim N(0, 0.3^2)$
 $\lambda = 0$

C53

2.5 ベイズ的方法

$$x_1, x_2, \dots, x_n, \dots \sim f(x|\theta)$$

事前分布 $f(\theta)$ を考えると

$$(x_1, \dots, x_n) \sim f(x_1, \dots, x_n) = \int f(x_1|\theta) \cdots f(x_n|\theta) f(\theta) d\theta$$

十分大きな n に対して

$$\text{BIC} = - \sum_{i=1}^n \log f(x_i|\hat{\theta}) + \frac{\log n}{2} \dim \theta$$

は $-\log f(x_1, \dots, x_n)$ の近似になる。(あまり良い近似ではないが.)

BIC を小さくするモデルは, 近似的に事後確率を大きくする

C55

重み付き最尤推定の良さを測る規準

データが $(y, z) \sim g(y|z)g_0(z)$ のとき,

$f(y|z, \hat{\theta}_w)g_1(z)$ が $g(y|z)g_1(z)$ にどれだけ近いかを測るには

$$\text{IC}_w = - \sum_{i=1}^n \frac{g_1(z_i)}{g_0(z_i)} \log f(y_i|z_i, \hat{\theta}_w) + \text{tr}(J_w H_w^{-1})$$

が使える. ただし

$$\hat{J}_w = \frac{1}{n} \sum_{i=1}^n \frac{g_1(z_i)}{g_0(z_i)} w(z_i) \frac{\partial \log f(y_i|z_i, \theta)}{\partial \theta} \Big|_{\hat{\theta}_w} \frac{\partial \log f(y_i|z_i, \theta)}{\partial \theta'} \Big|_{\hat{\theta}_w}$$

$$\hat{H}_w = - \frac{1}{n} \sum_{i=1}^n w(z_i) \frac{\partial^2 \log f(y_i|z_i, \theta)}{\partial \theta \partial \theta'} \Big|_{\hat{\theta}_w}$$

実際には重み関数として

$$w(z) = \left(\frac{g_1(z)}{g_0(z)} \right)^\lambda, \quad \lambda \in [0, 1]$$

C54

AIC と **BIC** の関係

$$f(x_1, \dots, x_n) = f(x_n|x_1, \dots, x_{n-1})f(x_{n-1}|x_1, \dots, x_{n-2}) \cdots f(x_3|x_1, x_2)f(x_2|x_1)f(x_1)$$

より

$$-\log f(x_1, \dots, x_n) = - \sum_{i=1}^n \log f(x_i|x_1, \dots, x_{i-1})$$

各時刻でのベイズ予測分布の良さは **IC** で測れて, 各項に

$$\frac{1}{2i} \dim \theta$$

の補正項が出て来る. これを $i = 1, \dots, n$ まで足し合わせて

$$\sum_{i=1}^n \frac{1}{2i} \dim \theta \approx \frac{\log n}{2} \dim \theta$$

という **BIC** の第2項が出て来る. これに対して, **AIC** は

$$x_1, \dots, x_n$$

から与える予測分布 $f(x_{n+1}|\hat{\theta}(x_1, \dots, x_n))$ の良さを測っている.

C56