

線形回帰モデルの変数選択法と
それによる解析例

1. ティーラの説明：三浦

2. 線形回帰モデルの説明：大嶋

3. Mallowsの C_p 規準：黒木，岡田

4. 変数選択法と解析例：浅野，辻

アメリカ合衆国の主要60都市における①
大気汚染の健康に及ぼす影響の一部

City RA E D ...log₁₀S O₂ MORT

OH 36 11.4 ... 1.77 921.9

NY 35 11.0 ... 1.59 997.9

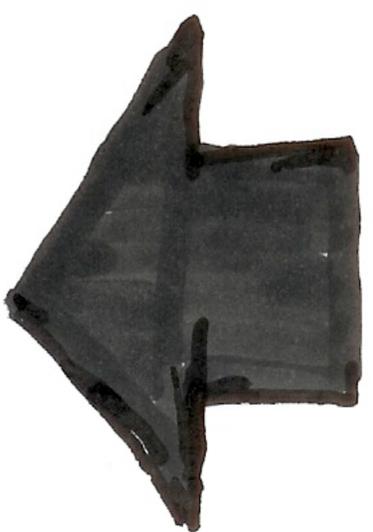
PA 44 9.8 ... 1.52 962.4

線形回帰モデル

②

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i$$

$i = 1, 2, \dots, n$



$$y = X\beta + \varepsilon$$

$n \times 1$

$n \times (p+1)$

$(p+1) \times 1$

$n \times 1$

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \beta_0 + x_{11}\beta_1 + \dots + x_{1p}\beta_p \\ \beta_0 + x_{21}\beta_1 + \dots + x_{2p}\beta_p \\ \vdots \\ \beta_0 + x_{n1}\beta_1 + \dots + x_{np}\beta_p \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix} \quad (3)$$

$$= \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

代表的なモデル規準

④

AkaikeのAIC規準

→ "小モデル" 良いモデル

$$AIC = -2Q(\beta) + 2(p+1)$$

= -2(最大尤度) + 2(モデルに含まれる独立パラメータ)

MallowsのCp規準

→ "小モデル" 良いモデル

$$C_p = \frac{RSS_p}{\hat{\tau}_2^2} + 2(p+1) - n$$

= $\frac{RSS_p}{\hat{\tau}_2^2} + 2 \times$ (モデルに含まれる独立パラメータ) - n

真のモデル

$$E(y) = \theta$$

$$VAR(y) = \sigma^2 I_n$$

$y = X\beta + \varepsilon$ の回帰モデル

$$E(y) = X\beta$$

$$VAR(y) = \sigma^2 I_n$$

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

⑥

$$\hat{\gamma} = \hat{\theta} = X ({}^t X X)^{-1} {}^t X y = H y$$

H: ハット行列

直の世界で言平価すると $E(\hat{\theta}) = H \theta$ かつ
推定値の良さを言平価する尺度として

$$\hat{\theta} - \theta = (\hat{\theta} - E(\hat{\theta})) + (E(\hat{\theta}) - \theta)$$

$$\Delta_p = E {}^t (\hat{\theta} - \theta) (\hat{\theta} - \theta)$$

$$= (p+1) \tau^2 + {}^t \theta (I-H) \theta$$

ノジビタ

ノビタカ

残差平方和 $RSS_p = \sum (y_i - \hat{\theta})^2$ (7)
真の世界では

$$E(RSS_p) = (n-p-1) \tau^2 + \theta^t \theta (I-H) \theta$$

Δ_p の不偏推定量は

$$RSS_p + 12(p+1) - n \tau^2$$

Mallows の C_p 規準

$$C_p = \frac{RSS_p}{\tau^2} + 2(p+1) - n = \left(\begin{array}{l} \text{F値に含めた} \\ \text{独立パラメータ数} \end{array} \right) - n$$

変数選択法

⑧

フルモテールのサイズ... p , サブモテールを合わせた数... 2^{p-1}

● 総当たり法

2^{p-1} 通り ($p=6$ のとき 63 通り)

すべての組み合わせの中から検討する方法

...とても大変

● 逐次選択法

逐次的に変数を選択する方法

・変数増加法 ・変数減少法

変数増加法のアルゴリズム

⑨

Step 1: $i \leftarrow 1$ (i : 選択された変数の個数)

・MORTとの相関係数が

最大となる説明変数を選ぶ。

・ C_i を求める。

Step 2: $i \leftarrow i + 1$

・最大の寄与率 ($\frac{RSS_i}{\sum (y_i - \hat{y}_i)^2}$) となる説明変数を選ぶ。
・ C_i を求める。

Step 3: if $C_{i-1} > C_i$: Step 2へ戻る
else $C_{i-1} \leq C_i$: 選択を終了する

変数減少法のアルゴリズム

⑩

Step 1: $i \leftarrow p$ (i : 選択された変数の個数)

寄与率と C_p を計算する。

Step 2: $i \leftarrow i - 1$

フルモデルの寄与率からの減少量が
最小となる変数を選ぶ。

Step 3: $C_i \leq C_{i+1}$ YES

NO \rightarrow 終了。