

# 具体的模型选择方法

## 四. 模型选择方法

Some popular model selection approaches

合格的模型选择方法应满足一些性质，如 选择一致性 (selection consistency)

假设  $\exists A_0 \subset P = \{2, 3, \dots, p\}$  (去除第一个截距项)，使对  $\forall j \in A_0$  有  $\beta_{[j]} \neq 0$   
而  $\forall j \notin A_0$  有  $\beta_{[j]} = 0$

合格的模型选择方法应当基于数据产生一个  $\hat{A}_n \subset \{2, 3, \dots, p\}$  满足

$$\lim_{n \rightarrow \infty} P(\hat{A}_n = A_0) = 1$$

### 一. 全子集回归 All-subset Regression

考虑  $P = \{2, 3, \dots, p\}$  的所有非空子集，使用标准  $M(\cdot)$  选择：

$$\hat{A}_n = \arg \min_{A \in \mathcal{P}} M(A)$$

下面描述一些记号：

(1) 记  $x_{iA}$  为  $x_i = (x_{i1}, \dots, x_{ip})^\top$  的子向量，只含  $x_{ij}$  :  $j \in \{1\} \cup A$  的分量

(2) 记  $\hat{\beta}_A$  表示数据集  $\{(y_i, x_{iA}) : i \in \mathcal{N}\}$  的 OLS 估计  $\mathcal{N} = \{1, 2, \dots, n\}$

之后，便要决定利用如何的标准  $M(\cdot)$

## 1.1 预测误差 prediction error

对于  $K \geq 2$ , 定义  $N = \{1, 2, \dots, n\}$  为  $K$  个互斥子集  $N_1, \dots, N_K$

对于  $k=1, 2, \dots, K$ , 记  $\hat{\beta}_A^{(-k)}$  为使用  $\{(y_i, x_{iA}) : i \notin N_k\}$  估计出的 OLS, 则有:

$$M(A) = \sum_{k=1}^K \sum_{i \in N_k} (y_i - x_{iA}^\top \cdot \hat{\beta}_A^{(-k)})^2$$

这便是交叉验证 (cross validation) 策略, 同样它不限于线性模型:

$$M(A) = \sum_{k=1}^K \sum_{i \in N_k} (y_i - g_A^{(-k)}(x_{iA}))^2$$

## 1.2 调整样本决定系数 Adjusted $R^2$

若以  $R^2$  作为  $M(\cdot)$ , 则对于  $A_1 \subset A_2$ , 必然有  $M(A_1) \leq M(A_2)$   
故要对自变量个数施加惩罚

令  $R^2(A)$  为使用  $\{(y_i, x_{iA}) : i \in N\}$  计算出的  $R^2$ , 调整为:

$$R^2_A(A) = 1 - \frac{(n-1) \cdot (1 - R^2(A))}{n - |A| - 1} = 1 - \frac{(n-1) \cdot SSE_A}{(n - |A| - 1) \cdot SST}$$

$$1 - \frac{(n-1)(1 - \frac{SST - SSE_A}{SST})}{n - |A| - 1} = 1 - \frac{(n-1) SSE_A}{(n - |A| - 1) SST}$$

其中:

$$SSE_A = \sum_{i=1}^n (y_i - x_{iA}^\top \cdot \hat{\beta}_A)^2$$

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

Tips:  $\begin{aligned} \Rightarrow SSE_A \downarrow &\Rightarrow R_A^2(A) \uparrow \\ A \uparrow & \\ \Rightarrow |A| \uparrow &\Rightarrow R_A^2(A) \downarrow \end{aligned}$

故可将  $M(A) = -R_A^2(A)$  作为标准

### 1.3 赤池信息量 Akaike information criterion, AIC

AIC 是基于 Likelihood 函数值的选择标准。

Definition AIC 定义：假设希望基于数据  $Z$  估计参数  $\theta \in \mathbb{R}^d$  而  $Z$  的似然函数为  $L(Z; \theta)$ ，由此得到 MLE 为  $\hat{\theta}$ ，则 AIC 为：

$$AIC = -2 \log L(Z; \hat{\theta}) + 2d$$

以上 AIC 由 2 部分组成：第一部分表示模型拟合的好坏。 $L(Z; \hat{\theta})$  越大，拟合越好；第二部分是对模型复杂的惩罚。

在正态假设下，基于  $\{(y_i, x_{iA}) : i \in N\}$  的选模型  $A$  的参数 MLE 分别为 OLS:  $\hat{\beta}_A$   $\hat{O}_A^2 = \frac{SSE_A}{n}$  代入  $L(Z; \hat{\theta})$  并去除无关项，得

$$AIC(A) = n \log SSE_A + 2 \cdot |A|$$

特别地： $d = |A| + 2(p+截距 + \hat{\sigma}_e^2)$

故最小化  $M(A) = AIC(A)$

## 1.4 $C_p$ 统计量 $C_p$ -Statistic

若  $A$  正确，则  $x_{iA}^\top \hat{\beta}_A$  可视为  $E(y_i) = x_i^\top \beta$  的估计 ( $i \in \mathcal{N}$ )

而估计偏差平方和与随机误差方差的比值为：

$$T_p(A) = \frac{1}{\hat{\sigma}^2} \sum_{i=1}^n (x_{iA}^\top \hat{\beta}_A - x_i^\top \beta)^2$$

可以证明：

$$E[T_p(A)] = \frac{E(SSE_A)}{\hat{\sigma}^2} - n + 2(|A| + 1)$$

忽略无关的项得  $C_p$  统计量

$$C_p(A) = \frac{SSE_A}{\hat{\sigma}^2} + 2|A|$$

其中  $\hat{\sigma}^2$  为全模型  $\sigma^2$  的无偏估计。

$$\hat{\sigma}^2 = \frac{1}{n-p} \sum_{i=1}^n (y_i - x_i^\top \hat{\beta})^2$$

故选择最小化  $M(A) = C_p(A)$

Tips. 由定理 3，若  $A$  正确，有  $\beta_{ij} = 0 \quad \forall j \notin A$ ，则

$$E(SSE_A) = (n - |A| - 1)\sigma^2$$

结合  $E(T_p(A))$  知

$$E[T_p(A)] = |A| + 1$$

放

$$\begin{aligned}M(A) &= |SSE_A/\hat{\sigma}^2 - n + 2(|A|+1) - (|A|+1)| \\&= |SSE_A/\hat{\sigma}^2 - n + |A|+1|\end{aligned}$$

也是合适的

## 二、逐步回归 Stepwise Regression