

Fitting linear models by Least squares

Part 3

Given: TS $(x_1, y_1) \dots (x_n, y_n)$

Goal: Generate prediction rule $f(x)$

Simple approach:

- Try linear rule $\ell(x) = b_0 + b_1 x_1 + \dots + b_p x_p$
- Find coefficient vector \hat{b} that minimizes resubstitution error:

$$\rightarrow \hat{b} = \underset{b}{\operatorname{argmin}} \|y - Xb\|^2 \quad \text{where}$$

$$\text{design matrix } X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}_{n \times p}$$

We always denote # of columns by p and assume $X' = (1 \dots 1)^T$

Facts

- $\hat{y} = \text{projection of } y \text{ on } [X' \dots X^p]$
- If $\operatorname{rank}(X) = p$ then
$$\hat{b} = (X^T X)^{-1} X^T y$$
$$\hat{y} = X \hat{b} = X (X^T X)^{-1} X^T y \quad \leftarrow \text{not operational! if } \operatorname{rank}(X) < p$$
- If $\operatorname{rank}(X) < p$ then \hat{b} not unique; there are ∞ many ways to write \hat{y} as a lincom of $X' \dots X^p$

There is a way to compute \hat{y} even for rank deficient case. May discuss later

For the moment let's assume $\text{rank}(X) = p$

Then
$$\hat{y} = \underbrace{X(X^T X)^{-1} X^T}_{n \times n} y$$

H is called "hat matrix"

H is a projection matrix

$$H = H^T \quad \text{symmetric}$$

$$H^2 = H \quad \text{idempotent}$$

Assume $y_i = f(x_i) + \varepsilon_i$
 ε_i iid $E(\varepsilon_i) = 0$
 $V(\varepsilon_i) = \sigma^2$

Then the expected squared estimation error
(expectation is taken over the ε 's; x_1, \dots, x_n
are held fixed).

$$\begin{aligned} ESE &= \frac{1}{n} E(\|f - \hat{y}\|^2) \\ &= \frac{1}{n} E(\|f - H y\|^2) \\ &= \frac{1}{n} E(\|f - H f - H \varepsilon\|^2) \\ &= \frac{1}{n} \left[E(\|f - H f\|^2) + E(\|H \varepsilon\|^2) \right] \\ &\quad + \text{cross-term} = 0 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n} \left[\|(I-H)\hat{f}\|^2 + E(\boldsymbol{\varepsilon}^T H^T H \boldsymbol{\varepsilon}) \right] \\
&= \frac{1}{n} \left[\|(I-H)\hat{f}\|^2 + E \sum_{i,j} H_{ij} \varepsilon_i \varepsilon_j \right] \\
&= \frac{1}{n} \left[\underbrace{\|(I-H)\hat{f}\|^2}_{\text{squared bias}} + \underbrace{\varepsilon^2 \sum_{i,j} H_{ij}}_{\text{variance}} \right]
\end{aligned}$$

trace $H = p$

Suppose $\hat{f} \in [X^1 \dots X^p]$ \Rightarrow squared bias = 0
 Estimate is unbiased

Suppose X^p unnecessary because
 $\hat{f} \in [X^1 \dots X^{p-1}]$. Then performance is better
 because variance term will decrease
 and the bias² won't increase.