

①

22 April 04

Cross-validation, The practice of partitioning a sample of data into subsamples such that analysis is initially performed on a single subsample, while further subsamples are retained "blind" in order for subsequent use in confirming and validating the initial analysis.

Tied in with the idea of prediction error.

Loss function $L(y, \hat{y})$

e.g. $(y - \hat{y})^2$

y : future response, \hat{y} : prediction from model

prediction error, Δ

$$E(y - \hat{y})^2$$

in classification problem

$$\text{Pr}\{\hat{y} \neq y\}$$

$$E L(y, \hat{y})$$

error rate

(2)

22 April 04

e.g. multiple regression

Data y_i, x_i^T

Future observation y_0, x_0^T

$$y_0 = x_0^T \beta + \epsilon_0$$

Aggregate prediction error Δ

$$\frac{1}{n} \sum E(y_0 - \hat{y}_0)^2$$

(*)

$$\hat{y}_0 = x_0^T \hat{\beta}$$

$$y_0 - \hat{y}_0 = x_0^T (\beta - \hat{\beta}) + \epsilon_0$$

$$E(y_0 - \hat{y}_0)^2 = x_0^T (X^T X)^{-1} x_0 \sigma^2 + \sigma^2$$

\sum is over x_1, \dots, x_n

$$(*) = \frac{1}{n} \sum_i \text{tr} \{ (X^T X)^{-1} x_i x_i^T \} \sigma^2 + \sigma^2$$

$$= \left(\frac{p+1}{n} \right) \sigma^2$$

Estimate by $\left(\frac{p+1}{n} \right) s^2$

(3)

22 April 04

Estimate by deleting cases in turn

$$\hat{\Delta}_{cv} = \frac{1}{n} \sum_i (y_i - x_i^T \hat{\beta}_{-i})^2$$

$$\hat{\beta} - \hat{\beta}_{-i} = (X^T X)^{-1} x_i (y_i - x_i^T \hat{\beta}) / (1 - H_{ii})$$

$$H_{ii}: \text{leverage} \quad H = X(X^T X)^{-1} X^T$$

$$\hat{\Delta}_{cv} = \frac{1}{n} \sum_i (y_i - x_i^T \hat{\beta})^2 / (1 - H_{ii})^2$$

$$\text{Here } m(x_i) = x_i^T \hat{\beta}$$

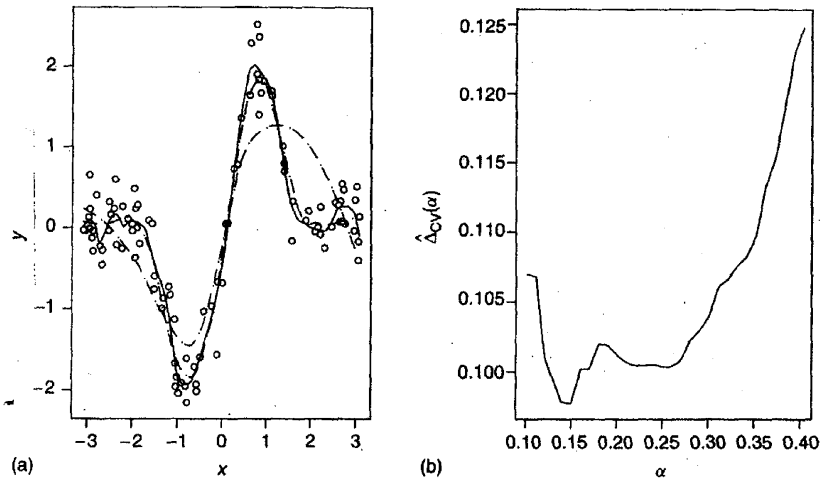
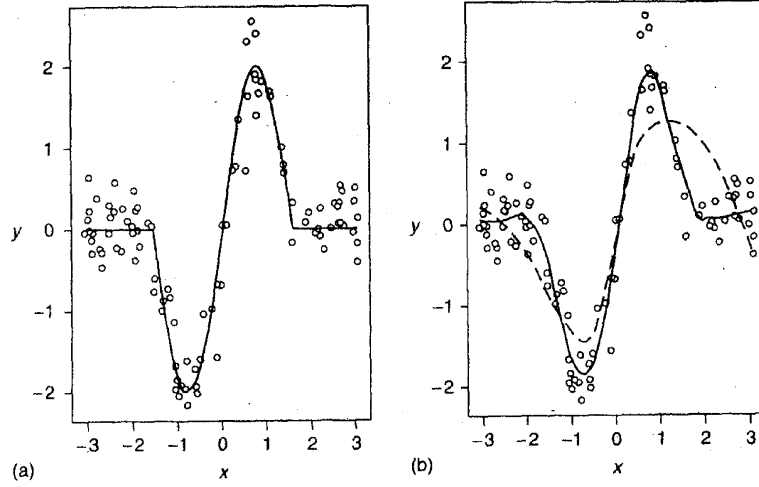
In case of $EY = m(x, \alpha)$ and

$$\hat{y} = \hat{m}(x, \alpha) = H(\alpha)y$$

have

$$\hat{\Delta}_{cv} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{m}(x_i, \alpha))^2 / (1 - H_{ii}(\alpha))^2$$

Can be used to choose α .



(a) LOESS fits based on $\alpha \in \left\{ \frac{1}{2}, \frac{2}{3} \right\}$ (dashed) and on cross-validated $\alpha_{CV}(S) = 0.15$ (solid); (b) $\hat{\Delta}_{CV}(\alpha)$ for $.11, \dots, 0.4$

22 April 04

Cross-validation in

Splus

cv.tree()

supsmu(, span = "cv")

ucv MASS

bcv MASS

smooth.spline()

R

validate.tree(Design)

bcv (MASS)

ucv (MASS)

gv.glm(boot)

knncv (class)

(fun Fcts)

xpred, rpart (rpart)

crossval (supclust)

cv.tree (tree)

mgcv ?