

Problem 1 (Ch. 2.4 p. 18-)

a) ~~loss~~, ~~loss~~ = loss of

a) $L(Y, f(X)) =$ loss if we predict by $f(X)$
and Y is correct.

Squared error loss: $L(Y, f(X)) = (Y - f(X))^2$

Optimal?

Choose f by minimizing

$$EPE(f) = E(Y - f(X))^2$$

expected
prediction
error

$$= E\{Y\} = E_X\{E_{Y|X}\}$$

$$= E\{E[(Y - f(X))^2 | X]\}$$

given

So ~~for~~ ~~at~~ $X=x$ we can ~~minimize~~ let ~~be~~ $f(X) = c$
where

$$c = \underset{c}{\operatorname{argmin}} \underbrace{E((Y - c)^2 | X=x)} \equiv \underset{c}{\operatorname{argmin}} g(c)$$

(see: $g(c)$ is a  parabola

$$g'(c) = E[-2(Y - c) | X=x] = 0$$

$$\Rightarrow c = E(Y | X=x). \quad \text{Q.E.D.}$$

b) $f(x) = x^T \beta$

Then can consider

$$EPE(\beta) = E(Y - X^T \beta)^2$$

Now set $\frac{\partial}{\partial \beta} EPE(\beta) = E \left[\frac{\partial}{\partial \beta} (Y - X^T \beta)^2 \right] = 0$

~~$= E \left[\frac{\partial}{\partial \beta} (Y - X^T \beta)^2 \right]$~~

On component form:

$$= E \left[\frac{\partial}{\partial \beta_j} \left(Y - \sum_{i=1}^p X_i \beta_i \right)^2 \right]$$

$$= E \left[-2 X_j \left(Y - \sum_{i=1}^p X_i \beta_i \right) \right] ; j=1, \dots, p$$

so equation is

$$E[X_j Y] = E \left[X_j \sum_{i=1}^p X_i \beta_i \right] ; j=1, \dots, p$$

or $E[X Y] = E[X X^T \beta] = E[X X^T] \beta$

so $\beta = (E[X X^T])^{-1} E(X Y)$ (see (2.16) p. 19)

c) Let $Y \in \{1, 2, \dots, k\}$.

$$L(Y, f(X)) = \begin{cases} 0 & \text{if } f(X) = Y \\ 1 & \text{if } f(X) \neq Y \end{cases}$$

$$EPE(f) = P[f(X) \neq Y]$$

$$= \sum_x \sum_{y=1}^k$$

$$= \sum_x \sum_{y=1}^k L(y, f(x)) P(X=x, Y=y)$$

$$= \sum_x \sum_{y=1}^k L(y, f(x)) P(X=x) P(Y=y | X=x)$$

$$= \sum_x P(X=x) \sum_{y=1}^k \underbrace{L(y, f(x)) P(Y=y | X=x)}_{\substack{\text{if } y=f(x) \\ \text{then } 1 \\ \text{else } 0}}$$

$$= \sum_x P(X=x) (1 - P(Y=f(x) | X=x))$$

$$= 1 - \sum_x P(Y=f(x) | X=x) P(X=x)$$

Thus - for each x let $f(x) = y$ where

$$\underline{\underline{y = \operatorname{argmax}_y P(Y=y | X=x)}}$$

d) Measures: Let \mathcal{T} = training set

$$Err_{\mathcal{T}} = E[L(Y, \hat{f}(X)) | \mathcal{T}]$$

So $\hat{f}(\cdot)$ is a function of training data, while (X, Y) is drawn from joint distr.

If we take expectation over \mathcal{T} also, we get

$$Err = E[Err_{\mathcal{T}}] = E[L(Y, \hat{f}(X))]$$

where randomness is (for fixed N)

- ① Draw training sample and compute $\hat{f}(\cdot)$
- ② Draw a new pair (X, Y) .

Training error:

$$\overline{err} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}(x_i))$$

is too optimistic because same data is used to ~~est~~ compute $\hat{f}(\cdot)$ as is used to measure fit.

Improvement: Estimate optimism :

$$\text{Define } Err_{in} = \frac{1}{N} \sum_{i=1}^N E_{Y_i^0} [L(Y_i^0, \hat{f}(x_i)) | \mathcal{T}]$$

where the same x_i as in training set \mathcal{T} are used, but where Y_i^0 are drawn "at each x_i ", i.e. from cond. distr. ~~given~~ of Y given $X = x_i$.

Define

$$op \equiv Err_{in} - \bar{err}$$

and $\omega \equiv E_y(op)$

where E_y is expected value over the training set with a fixed x_1, x_2, \dots, x_N

For many loss functions we have

$$E_y(Err_{in}) = E_y(\bar{err}) + \underbrace{\frac{2}{N} \sum_{i=1}^N Cov(\hat{y}_i, y_i)}_{\omega}$$

$$\hat{\omega} \approx E_y(\bar{err}) + 2 \cdot \frac{d}{N} \sigma_\varepsilon^2$$

often when model is

$$Y = f(X) + \varepsilon$$

so we can estimate

$$Err_{in} \approx \bar{err} + 2 \cdot \frac{d}{N} \sigma_\varepsilon^2$$

e) K-fold Cross-Validation

Divide data in K approx. equal sets ($K=5$ or 10 usual)

Let $\kappa: \{1, \dots, N\} \rightarrow \{1, \dots, K\}$ be the indexing function

$$CV(\hat{f}) = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{-\kappa(i)}(x_i))$$

is the CV-estimate of prediction error.

Bootstrapping:

From training data $(x_1, y_1), \dots, (x_N, y_N)$
draw with replacement N times:

$$(x_1^*, y_1^*), \dots, (x_N^*, y_N^*)$$

Do this B times!

Use these to ~~fit~~ ^{compute} $\hat{f}^{*b}(\cdot)$

Then measure of prediction error is e.g.

$$\hat{Err}^{(.632)} = 0.368 \bar{err} + 0.632 \cdot \hat{Err}^{(1)}$$

where

$$\hat{Err}^{(1)} = \frac{1}{N} \sum_{i=1}^N \frac{1}{C^{-i}} \sum_{b \in C^{-i}} L(y_i, \hat{f}^{*b}(x_i))$$

These are appropriate when there are too few data to have separate validation and test sets: Use all data as training data!

Problem 2

a) Cubic spline is given by:

- K knots $\xi_1, \xi_2, \dots, \xi_K$
- Between knots given by 3rd order polynomials
- Continuous at knots
Continuous first and second derivatives at knots.

Natural Cubic spline: linear outside boundary knots.

Dimension of cubic splines:

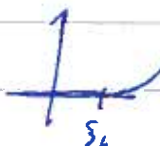
$$\begin{array}{r} K+1 \text{ regions with } 4 \text{ coefficients} \quad 4K+4 \\ - K \text{ knots} \cdot 3 \text{ restrictions} \quad -3K \end{array}$$

$\Rightarrow K+4$ basis functions needed.

These ~~are~~ ^{can be} given as in exercise

(Must check that they are

- linearly independent
- requirements at knots:

Consider $g(x) = (x - \xi_k)_+^3$  So $g(\xi_k+) = g(\xi_k-) = 0$

$$g'(x) = 3(x - \xi_k)_+^2$$

$$g'(\xi_k+) = g'(\xi_k-) = 0$$

$$g''(x) = 6(x - \xi_k)_+$$

$$g''(\xi_k+) = g''(\xi_k-) = 0.$$

b) Natural requires

for $x < \xi_1$ is $f(x)$ linear
now the given $f(x)$ is for $x < \xi_1$

$$f(x) = \sum_{j=0}^3 \beta_j x^j$$

so linear $\Rightarrow \beta_2 = \beta_3 = 0.$

For $x > \xi_k$ is

$$f(x) = \beta_0 + \beta_1 x + \sum_{k=1}^K \theta_k (x - \xi_k)^3 \quad (\text{without } +\text{'s})$$

So

$$f'(x) = \beta_1 + \sum_{k=1}^K \theta_k 3(x - \xi_k)^2$$

$$f''(x) = \sum_{k=1}^K 6\theta_k (x - \xi_k)$$

$$f'''(x) = \sum_{k=1}^K 6\theta_k$$

Hence: $f'''(x) = 0 \Rightarrow \sum_{k=1}^K \theta_k = 0$

$$f''(x) = 0 \Rightarrow \underbrace{\sum_{k=1}^K \theta_k x}_{=0} = \sum_{k=1}^K \theta_k \xi_k$$

$$\Rightarrow \sum_{k=1}^K \theta_k \xi_k = 0$$

c) For natural splines we need

$$\underbrace{N+4 - 4}_{\text{cubic spline}} = N \text{ basis functions}$$

coeff. for x^2 and x^3 terms are 0 outside both ξ_1 and ξ_K .

Now note that

$$\sum_{k=1}^K \theta_k = 0 \Rightarrow \theta_K = -\sum_{k=1}^{K-1} \theta_k \quad (*)$$

$$\sum_{k=1}^K \xi_k \theta_k = 0 \Rightarrow \sum_{k=1}^{K-1} \xi_k \theta_k + \xi_K \theta_K = 0$$

\Rightarrow using $(*)$

$$\sum_{k=1}^{K-1} (\xi_k - \xi_K) \theta_k = 0 \quad (**)$$

or $\sum_{k=1}^{K-1} (\xi_k - \xi_K) \theta_k = -(\xi_{K-1} - \xi_K) \theta_{K-1}$

Now under constraints in b) is

$$f(x) = \beta_0 + \beta_1 x + \sum_{k=1}^K \theta_k (x - \xi_k)_+^3 = \beta_0 + \beta_1 x + \sum_{k=1}^{K-1} \theta_k (x - \xi_k)_+^3 + \theta_K (x - \xi_K)_+^3$$

$$= \beta_0 + \beta_1 x + \sum_{k=1}^{K-1} \theta_k (x - \xi_k)_+^3 - \sum_{k=1}^{K-1} \theta_k (x - \xi_K)_+^3$$

$$= \beta_0 + \beta_1 x + \sum_{k=1}^{K-1} \theta_k \left[(x - \xi_k)_+^3 - (x - \xi_K)_+^3 \right]$$

$$= \beta_0 + \beta_1 x + \sum_{k=1}^{K-1} \theta_k (\xi_K - \xi_k) d_k(x)$$

~~But using (*) this is~~

$$= \beta_0 + \beta_1 x + \sum_{k=1}^{K-2} \theta_k (\xi_K - \xi_k) d_k(x)$$

$$+ \theta_{K-1} (\xi_K - \xi_{K-1}) d_{K-1}(x)$$

~~$= \sum_{k=1}^{K-2} (\xi_k - \xi_{k+1}) \theta_{k+1} d_{k+1}(x)$~~

by (*) is this = $\sum_{k=1}^{K-2} (\xi_k - \xi_{k+1}) \theta_{k+1} d_{k+1}(x)$

$$= \beta_0 + \beta_1 x + \sum_{k=1}^{K-2} \theta_k (\xi_K - \xi_k) (d_k(x) - d_{k-1}(x))$$

$$= \beta_0 N_1(x) + \beta_1 N_2(x) + \sum_{k=1}^{K-2} \theta_k (\xi_K - \xi_k) N_{k+2}(x)$$

which gives a linear combination } Q.E.D.
of the $N_1(\cdot), \dots, N_K(\cdot)$

Alternatively:

$$\beta_0 + \beta_1 x + \sum_{k=1}^{K-2} \alpha_k N_{k+2}(x)$$

$$= \beta_0 + \beta_1 x + \sum_{k=1}^{K-2} \alpha_k [d_k(x) - d_{k-1}(x)]$$

$$= \beta_0 + \beta_1 x + \sum_{k=1}^{K-2} \frac{\alpha_k}{s_K - s_k} d_k(x) - d_{K-1}(x) \sum_{k=1}^{K-2} \alpha_k$$

$$= \beta_0 + \beta_1 x + \sum_{k=1}^{K-2} \frac{\alpha_k}{s_K - s_k} (x - s_k)^3 + \sum_{k=1}^{K-2} \frac{\alpha_k}{s_K - s_k} (x - s_k)^3$$

$$- \frac{(x - s_{K-1})^3}{s_K - s_{K-1}} \sum_{k=1}^{K-2} \alpha_k + \frac{(x - s_K)^3}{s_K - s_{K-1}} \sum_{k=1}^{K-2} \alpha_k$$

~~Coeff to~~ θ_k

For $k \leq K-2$: $\theta_k = \frac{\alpha_k}{s_K - s_k}$

$$\theta_{K-1} = - \left(\sum_{k=1}^{K-2} \alpha_k \right) \cdot \frac{1}{s_K - s_{K-1}}$$

$$\theta_K = - \sum_{k=1}^{K-2} \frac{\alpha_k}{s_K - s_k} + \frac{\sum_{k=1}^{K-2} \alpha_k}{s_K - s_{K-1}}$$

Now

$$\sum_{k=1}^K \theta_k = \sum_{k=1}^{K-2} \frac{\alpha_k}{\xi_K - \xi_k} - \frac{\sum_{k=1}^{K-2} \alpha_k}{\xi_K - \xi_{K-1}} \quad \text{--- det same = } \underline{\underline{0}}$$

$i \theta_k$

$$\sum_{k=1}^K \xi_k \theta_k = \sum_{k=1}^{K-2} \frac{\xi_k \alpha_k}{\xi_K - \xi_k} - \frac{\left(\sum_{k=1}^{K-2} \alpha_k \right) \xi_{K-1}}{\xi_K - \xi_{K-1}} + \frac{\left(\sum_{k=1}^{K-2} \alpha_k \right) \xi_K}{\xi_K - \xi_{K-1}}$$

$$\left(\sum_{k=1}^{K-2} \alpha_k \right)$$

$$= \sum_{k=1}^{K-2} \frac{\alpha_k}{\xi_K - \xi_k} (\xi_k - \xi_K)$$

$$= - \sum_{k=1}^{K-2} \alpha_k$$

$$\underline{\underline{0}}$$