

Ex 2.3 Derive  $d(p, N) = \left(1 - \frac{1}{2}\right)^{1/p}$

Treat the distance between data points and the origin as a random variable  $X$ , since the data points are equally distributed, the cdf of  $X$  is

$$F(X < x) = x^p, \quad x \in [0, 1]$$

Then the closest distance from data points to the origin is order statistics, whose cdf is

$$F_1(x) = 1 - (1 - F(x))^N$$

and whose median satisfies  $F(d(p, N)) = \frac{1}{2}$ , then it is derived.

Ex. 2.5

- (a) Derive equation (2.27). The last line makes use of (3.8) through a conditioning argument.
- (b) Derive equation (2.28), making use of the *cyclic* property of the trace operator [ $\text{trace}(AB) = \text{trace}(BA)$ ], and its linearity (which allows us to interchange the order of trace and expectation).

$$\begin{aligned} \text{(a) } \text{EPE}(x_0) &= E_{y_0|x_0} E_{\mathcal{T}} (\hat{y}_0 - y_0)^2 \\ &= \text{Var}(y_0|x_0) + E_{\mathcal{T}} [\hat{y}_0 - E_{\mathcal{T}} \hat{y}_0]^2 + [E_{\mathcal{T}} \hat{y}_0 - x_0^T \beta]^2 \\ &= \text{Var}(y_0|x_0) + \text{Var}_{\mathcal{T}}(\hat{y}_0) + \text{Bias}^2(\hat{y}_0) \\ &= \sigma^2 + E_{\mathcal{T}} x_0^T (\mathbf{X}^T \mathbf{X})^{-1} x_0 \sigma^2 + 0^2. \end{aligned} \quad (2.27)$$

$$\begin{aligned} &= E_{y_0|x_0} E_{\mathcal{T}} (x_0^T \beta + \varepsilon - \hat{y}_0)^2 \\ &= E_{y_0|x_0} E_{\mathcal{T}} [(x_0^T \beta - \hat{y}_0)^2 + 2\varepsilon(x_0^T \beta - \hat{y}_0) + \varepsilon^2] \\ &= E_{\mathcal{T}} [x_0^T \beta - \hat{y}_0]^2 + 2 E_{y_0|x_0} [\varepsilon(x_0^T \beta - \hat{y}_0)] \\ &\quad + E_{y_0|x_0} (\varepsilon^2) \\ &= E_{y_0|x_0} (\varepsilon^2) + E_{\mathcal{T}} [x_0^T \beta - E(\hat{y}_0) + E(\hat{y}_0) - \hat{y}_0]^2 \end{aligned}$$

$E_{y_0|x_0}(\varepsilon) = 0$

$$\begin{aligned} \text{(b) } E_{x_0} \text{EPE}(x_0) &\stackrel{\substack{X_0(p \times 1) \\ \text{tr} \times \text{tr}}}{\sim} E_{x_0} x_0^T \text{Cov}(X)^{-1} x_0 \sigma^2 / N + \sigma^2 \\ &= \text{trace}[\text{Cov}(X)^{-1} \text{Cov}(x_0)] \sigma^2 / N + \sigma^2 \\ &= \sigma^2 (p/N) + \sigma^2. \end{aligned} \quad (2.28)$$

Since  $\text{trace}(AB) = \text{trace}(BA)$

$$\begin{aligned} &= \text{trace}(E_{x_0} (x_0^T \text{Cov}(X)^{-1} x_0) \sigma^2 / N + \sigma^2) \\ &= \text{trace}(E_{x_0} (x_0 x_0^T \text{Cov}(X)^{-1}) \sigma^2 / N + \sigma^2) \\ &= \text{trace}(E_{x_0} [x_0 x_0^T (\mathbf{X}^T \mathbf{X})^{-1}] \sigma^2 / N + \sigma^2) \\ &\quad \sigma^2 \mathbf{I}_p \sigma^2 \mathbf{I}_p \\ &= \sigma^2 \sigma^{-2} \text{trace}(\mathbf{I}_p) \cdot \sigma^2 / N + \sigma^2 \\ &= \sigma^2 (p/N) + \sigma^2 \end{aligned}$$

Ex 2.6

Ex. 2.6 Consider a regression problem with inputs  $x_i$  and outputs  $y_i$ , and a parameterized model  $f_{\theta}(x)$  to be fit by least squares. Show that if there are observations with *tied* or *identical* values of  $x$ , then the fit can be obtained from a reduced weighted least squares problem.

Ex. 2.6 OLS:  $\hat{f}_{\theta}(x) = X(X^T X)^{-1} X^T Y$  (1)

for example,  $X_M = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$  while  $X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow x_3 = x_1 + x_2$

where we derived the weight  $W = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$

$$\text{Then } X_M = W X_M = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_1 + x_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

plug it in (1):  $\hat{f}_{\theta}(x) = W X_M (X_M^T W^T W X_M)^{-1} X_M^T W^T Y$   
which is a weighted LS fit.

Obviously, the size of  $X_M$  is smaller than that of  $X$ ,  
 $\hat{f}_{\theta}(x)$  is a reduced weighted LS.