

$$1. (a) \frac{\partial}{\partial \beta} E[(Y - X^T \beta)^2]$$

$$= \frac{\partial}{\partial \beta} E[Y^T Y - 2Y^T X^T \beta + (X^T \beta)^T \cdot X^T \beta] = E[-2XY + 2XX^T \beta] = 0$$

$$\therefore 2E[XX^T]\beta = 2E[XY]$$

$$\therefore E[XX^T]\beta = E[XY]$$

$$\therefore \beta = E^{-1}[XX^T]E[XY]$$

(b) The nearest neighbors calculate the average number of the labels.

of some points that are close to aim point x .

The least squares use the function $\beta = (X^T X)^{-1} X^T y$ to get the parameters β to minimize the squared error loss and then approximate the regression function.

Difference: ① The nearest neighbors is local because for a new point x , we only consider the points that is close to x . It has no

consumption for model. When dimension grows, we need exponential increase of sample points to approximate the regression function. It's

Bias will not change too much but Var will be higher when dimension grows.

② The least square is global because it considers all the sample points in the training set. It consumes the model linearly before training. When dimension grows, we need linear increase of sample points to approximate the regression function.

$$(c) EPE(f) = E(|Y - f(x)|)$$

$$= E_X E_{Y|X}(|Y - f(x)| | X)$$

$$\therefore f(x) = \arg\min_f E_{Y|X}(|Y - f(x)| | X=x)$$

$$= \arg\min_f \int_y |y - f(x)| Pr(y|x) dy$$

$$= \arg\min_f \left(\int_{y \geq f(x)} (y - f(x)) Pr(y|x) dy + \int_{y < f(x)} (f(x) - y) Pr(y|x) dy \right)$$

$$= \arg\min_f \left(\underbrace{\int_{y \geq f(x)} y Pr(y|x) dy - \int_{y < f(x)} y Pr(y|x) dy + f(x) Pr(y < f(x) | X) - f(x) Pr(y \geq f(x) | X)}_{=0} \right)$$

$$\therefore \frac{\partial A}{\partial f(x)} = Pr(y < f(x) | X) - Pr(y \geq f(x) | X)$$

$$= 2 Pr(y < f(x) | X) - 1 = 0$$

$$\therefore Pr(y < f(x) | X) = \frac{1}{2}$$

$$\therefore f(x) = \text{median}(Y | X=x)$$

$$\begin{aligned}
 2. (a) \quad \|\hat{y} - \hat{X}w\|_2^2 &= \sum_{i=1}^n (y_i - x_i w_i)^2 + \sum_{i=n+1}^{d+n} (y_i - x_i w_i)^2 \\
 &= \sum_{i=1}^n (y_i - x_i w_i)^2 + \sum_{i=n+1}^{n+d} (\sqrt{\lambda} w_i)^2 \\
 &= \|y - Xw\|_2^2 + \lambda \|w\|_2^2
 \end{aligned}$$

$$\therefore \frac{\partial \|y - \hat{X}w\|_2^2}{\partial w} = \frac{\partial (\|y - Xw\|_2^2 + \lambda \|w\|_2^2)}{\partial w} = \frac{\partial (y^T y - 2y^T Xw + w^T X^T X w + \lambda w^T w)}{\partial w}$$

$$= -2X^T y + 2X^T X w + 2\lambda w = 0$$

$$\therefore (X^T X + \lambda I_d)w = X^T y$$

$$\therefore w = (X^T X + \lambda I_d)^{-1} X^T y$$

$$(b) \quad \hat{X}\beta = y \quad \therefore \beta \in \mathbb{R}^{n+d}$$

$$\text{Suppose } \beta = \begin{bmatrix} \beta_d \\ \beta_n \end{bmatrix}$$

$$\therefore \hat{X}\beta = [X \quad \alpha I_n] \begin{bmatrix} \beta_d \\ \beta_n \end{bmatrix}$$

$$= X\beta_d + \alpha I_n \beta_n$$

$$= X\beta_d + \alpha \beta_n = y$$

$$\therefore \beta_n = \frac{1}{\alpha} (y - X\beta_d)$$

$$\therefore \|\beta\|^2 = \|\beta_d\|^2 + \|\beta_n\|^2$$

$$= \|\beta_d\|^2 + \frac{1}{\alpha^2} \|y - X\beta_d\|^2$$

$$\therefore \frac{\partial \|\beta\|^2}{\partial \beta_d} = 2\beta_d + \frac{1}{\alpha^2} \cdot 2(-X^T y + X^T X \beta_d)$$

$$= 2\beta_d + \frac{2}{\alpha^2} (-X^T y + X^T X \beta_d) = 0$$

$$\therefore \beta_d = (X^T X + \alpha^2 I_d)^{-1} X^T y$$

$$\therefore \alpha = \sqrt{\lambda}$$

$$\therefore \beta_n = \frac{y - X(X^T X + \lambda I_d)^{-1} X^T y}{\sqrt{\lambda}}$$

$$(c) \text{ As } X' = XV, w' = V^T w$$

$$\therefore \|y - X'w'\|_2^2 + \lambda \|w'\|_2^2 = \|y - XVV^T w\|_2^2 + \lambda \|V^T w\|_2^2$$

$$\text{As } X = U\Sigma V^T \text{ and } V \text{ is identity}$$

$$\therefore \|y - X'w'\|_2^2 + \lambda \|w'\|_2^2 = \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

$$= \|y - U\Sigma V^T w\|_2^2 + \lambda \|w\|_2^2$$

$$= y^T y - 2y^T U\Sigma w + (U\Sigma w)^T \cdot U\Sigma w + \lambda w^T w$$

$$\therefore \frac{\partial (\|y - X'w\|_2^2 + \lambda \|w\|_2^2)}{\partial w} = -2(y^T \cdot U \Sigma)^T + 2 \Sigma^T \Sigma w + 2 \lambda w = 0$$

$$\therefore w = (\Sigma^T \Sigma + \lambda I)^{-1} \cdot (y^T U \Sigma V^T)^T$$

$$= (\Sigma^T \Sigma + \lambda I)^{-1} \cdot X^T y$$

$$\text{As } \Sigma^T \Sigma + \lambda \cdot I = \begin{bmatrix} \sigma_1^2 + \lambda & & \\ & \sigma_2^2 + \lambda & \\ & & \ddots \\ & & & \sigma_d^2 + \lambda \end{bmatrix}$$

$$\therefore (\Sigma^T \Sigma + \lambda I)^{-1} = \begin{bmatrix} \frac{1}{\sigma_1^2 + \lambda} & & \\ & \frac{1}{\sigma_2^2 + \lambda} & \\ & & \ddots \\ & & & \frac{1}{\sigma_d^2 + \lambda} \end{bmatrix}$$

$$\therefore w_j = w_j = \sum_{i=1}^n \frac{1}{\sigma_i^2 + \lambda} \cdot x_{ij} \cdot y_i$$

$$3. (a) ① E[\hat{X} - \mu] = E\left[\frac{x_1 + x_2 + \dots + x_n}{n} - \mu\right] = \frac{1}{n}E[x_1] + \frac{1}{n}E[x_2] + \dots + \frac{1}{n}E[x_n] - \mu$$

$$\because E[x_i] = \mu \text{ for every } i \text{ in } 1, 2, \dots, n$$

$$\therefore \text{Bias}[\hat{X}] = \frac{1}{n} \cdot n \cdot \mu - \mu = 0$$

$$\text{Var}[\hat{X}] = \text{Var}\left[\frac{x_1 + x_2 + \dots + x_n}{n}\right] = \frac{1}{n^2} \text{Var}[x_1] + \frac{1}{n^2} \text{Var}[x_2] + \dots + \frac{1}{n^2} \text{Var}[x_n]$$

$$\because \text{Var}[x_i] = \sigma^2 \text{ for every } i \text{ in } 1, 2, \dots, n$$

$$\therefore \text{Var}[\hat{X}] = \frac{1}{n^2} \cdot n \cdot \sigma^2 = \frac{\sigma^2}{n}$$

$$② E[\hat{X} - \mu] = E\left[\frac{x_1 + x_2 + \dots + x_n}{n + n_0} - \mu\right]$$

$$\text{The same as above, } \text{Bias}[\hat{X}] = \frac{1}{n + n_0} \cdot n \cdot \mu - \mu = -\frac{n_0}{n + n_0} \mu$$

$$\text{Var}[\hat{X}] = \text{Var}\left[\frac{x_1 + x_2 + \dots + x_n}{n + n_0}\right]$$

$$\text{The same as above, } \text{Var}[\hat{X}] = \frac{1}{(n + n_0)^2} \cdot n \cdot \sigma^2 = \frac{n \cdot \sigma^2}{(n + n_0)^2}$$

$$(b) E[(\hat{X} - \mu)^2] = E[(\hat{X} - E(\hat{X}) + E(\hat{X}) - \mu)^2]$$

$$= E[(\hat{X} - E(\hat{X}))^2] + 2E[(\hat{X} - E(\hat{X})) (E(\hat{X}) - \mu)] + E[(E(\hat{X}) - \mu)^2]$$

$$= \text{Var}[\hat{X}] + 2(E[\hat{X}] - E[E(\hat{X})]) \{E(\hat{X}) - \mu\} + E[(\text{Bias}[\hat{X}])^2]$$

$$= \text{Var}[\hat{X}] + (\text{Bias}[\hat{X}])^2$$

$$E[(\hat{X} - X')^2] = E[(\hat{X} - E(\hat{X}) + E(\hat{X}) - X')^2]$$

$$= \text{Var}[\hat{X}] + 2E[(\hat{X} - E(\hat{X})) (E(\hat{X}) - X')] + E[(E(\hat{X}) - X')^2]$$

$$= \text{Var}[\hat{X}] + 2E[\hat{X}E(\hat{X}) - \hat{X}X' - E(\hat{X})^2 + E(\hat{X})X'] + E[E(\hat{X})^2 + X'^2 - 2E(\hat{X})X']$$

$$= \text{Var}[\hat{X}] + 2[E(\hat{X})^2 - E(\hat{X}X') - E(\hat{X})^2 + E(\hat{X}) \cdot E(X')] + E(\hat{X})^2 + E(X'^2)$$

$$= \text{Var}[\hat{X}] - 2E(\hat{X}X') + E(\hat{X})^2 + E(X'^2) - 2E(\hat{X})E(X')$$

$$= \text{Var}[\hat{X}] + \text{Bias}[\hat{X}]^2 + 2E[\hat{X}]\mu - \mu^2 - 2E[\hat{X}X'] + E[X'^2]$$

$$= \text{Var}[\hat{X}] + \text{Bias}[\hat{X}]^2 - 2\text{Cov}(\hat{X}, X') + \text{Var}[X']$$

$$\text{As } \hat{X} \text{ and } X' \text{ are i.i.d, so } \text{Cov}(\hat{X}, X') = 0$$

$$\therefore E[(\hat{X} - X')^2] = \text{Var}[\hat{X}] + \text{Bias}[\hat{X}]^2 + \text{Var}[X']$$

$E[(\hat{X} - X')^2]$ is bigger than $E[(\hat{X} - \mu)^2]$ and the difference between them is the variance of X'

$$(c). E[(\hat{x} - \mu)^2] = \text{Var}[\hat{x}] + (\text{Bias}[\hat{x}])^2$$

$$\text{way 1: } E[(\hat{x} - \mu)^2] = \frac{\sigma^2}{n} + 0 = \frac{\sigma^2}{n}$$

$$\text{way 2: } E[(\bar{x} - \mu)^2] = \frac{n\sigma^2}{(n+n_0)^2} + \frac{n_0^2}{(n+n_0)^2} \mu^2$$

$$= \frac{n\sigma^2 + n_0^2 \mu^2}{(n+n_0)^2}$$