## Statistical decision Theory

$\underbrace{X \in \mathbb{R}^p}_{\text{input}}, \underbrace{Y \in \mathbb{R}}_{\substack{\text{output} \\ \text{(regression)}}}$ $\underbrace{Pr(x,y)}$ ← Assumption - Datas are drawn from a distribution.

→ Joint distribution → $\{(x_1, y_1)$

Goal: $f(x): \mathbb{R}^p \to \mathbb{R}$

$\hat{y} = f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 \cdots + \beta_p x_p$

$(x_n, y_n)\}$

$x = (x_1, x_2, x_3, \cdots x_p)^T$

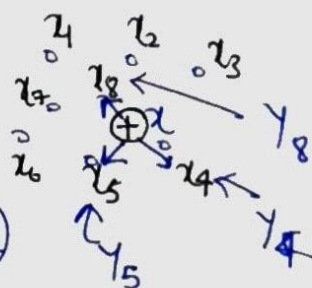$= \beta_0 + \sum_{j=1}^{p} x_j \beta_j$

set $\Rightarrow x_0 = 1$

$= \sum_{j=0}^{p} x_j \beta_j$ $\quad$ [Linear Regression]

### K- Nearest Neighbour:

$\hat{y}(x) = \frac{1}{K} \sum_{x_i \in N(x)} y_i$

$K = 3$



$\hat{y}(1) = \frac{1}{3} \sum (y_5 + y_8 + y_4)$

Loss $f^n$: $L(Y, f(x))$

Squared error: $(Y - f(x))^2$

Expected prediction Error (EPE (f)) = $E\{(Y - f(x))^2\}$

$= \int (Y - f(x))^2 \cdot pr(dx, dy)$

$Pr(x,y) = P_r(Y|x) P_r(x)$ �� $\rightarrow$

$= E_x E_{Y/x} ([Y - f(x)]^2 | x)$

minimise over diff values of c. ~~Getrido min~~ single point)

Take that value of c to calculate $f(x)$

$f(x) = \text{arg min}_{c} E_X E_{Y/x} ([Y - c]^2 | X = x_0)$

(Conditioning on a point)

[1]

The function $f(x)$ that minimises mean-squared error is given as

$$f(x) = E[Y|X=x] \quad \text{(Conditional Expectation or regression function)}$$

$$\boxed{f(x) = \text{Avg}(Y_i \mid x_i = x)} \rightarrow \text{we suppose for } x_i = 3,$$

$$f(x) = \frac{Y_i + Y_i' + Y_i'' + Y_i'''}{4}$$

$$= \frac{(2+3+7+9)}{4}$$

$Y_i = 2$
$Y_i' = 3$
$Y_i'' = 7$
$Y_i''' = 9$

1. Problem: If we have only one sample of $x$. Its difficult to ~~get~~ avg.

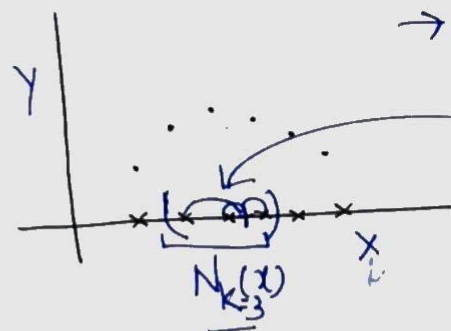2. unable to predict on unseen data.

$\rightarrow \boxed{\underline{\text{NOT WORKING}}}$

Conclusion:

Instead of Conditioning on single data point, we will Condition on a region (many data points)

$\rightarrow$ Neighborhood of $x$

$$\hat{f}(x) = \text{Avg}(Y_i \mid x_i \in N_k(x))$$

Assumption: output of the function will be constant on that region.

$\rightarrow$ output will be same for

this region



$N_{k=3}(x)$   $x_i$

As $K \uparrow$, estimate **more** stable

As $N, K \to \infty$, $K/N \to 0$, $\hat{f}(x) \Longrightarrow E[Y|x=x]$

[2]

In Case of linear regression,

$$f(x) = x^T \beta \qquad X = \begin{bmatrix} x_{11}, x_{12}, \cdots x_{1p} \\ x_{21}, x_{22}, \cdots x_{2p} \\ \\ x_{n1}, x_{n2}, \cdots x_{np} \end{bmatrix}_{m \times p}$$

$$(Y - X\beta)^2 = EPE(\hat{f})$$

$$\boxed{\hat{\beta} = (X^T X)^{-1} X^T y}$$

$$(Y - X\beta)^T (Y - X\beta) = Y^T(Y - X\beta) - (X\beta)^T(Y - X\beta)$$

$$= Y^T Y - Y^T X\beta - \beta^T X^T Y + \beta^T X^T X \beta$$

$$E(\hat{f}) = Y^T Y - 2Y^T X\beta + \beta^T X^T X \beta$$

$$\frac{\partial E}{\partial \beta} = 0 - 2Y^T X + 2X^T X\beta = 0$$

$$\Rightarrow X^T X\beta = Y X^T y$$

$$\boxed{\hat{\beta} = (X^T X)^{-1} X^T y}$$

## Classification

$$x \in \mathbb{R}^p, \quad G \in \mathcal{G} \qquad pr(x, G) \qquad \begin{cases} (x_1, G_1), \\ (x_2, G_2) \\ (x_3, G_3) \\ \vdots \\ (x_n, G_{12}) \end{cases} y$$

$$f(x): \mathbb{R}^p \to \mathcal{G}$$

LOSS = is a $(K \times K)$ matrix where

$\boxed{\text{(zero on diagonal)}}$  $K = Card(\mathcal{G})$

$\boxed{L(k, \ell) \text{ Cost of classifying } k \text{ as } \ell.}$

$0 - 1$ loss function, $k = 3$

$$\begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}_{3 \times 3}$$

3

$$EPE(\hat{f}) = E\left[L\left(G,\hat{f}\right)\right] \rightarrow \text{Discrete distribution)}$$
$$= E_x E_{G|x}\left\{\left[G,\hat{f}\right]\mid x\right\}$$
$$= E_x \sum_{k=1}^{K} L\left[k,\hat{f}(x)\right] pr(k|x)$$

original class

$$\hat{f}(x) = \arg\min_{g} \sum_{k=1}^{K} L(k,g)\cdot pr(k|x=x)$$

$\hookrightarrow$ predicted class

### 0-1 class   3 classes

$$pr(1|x) = 0.6$$
$$pr(2|x) = 0.2$$
$$pr(3|x) = 0.2$$

let's $\underline{\underline{g=2}}$

$$\begin{array}{c} \downarrow \\ \rightarrow \\ \rightarrow \end{array}\begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & \phi \\ \phi & \phi & 0 \end{bmatrix}\hspace{-1em}\begin{array}{c} L \end{array}$$

$$\hat{f}(x) = \sum_{k=1}^{K} L(k,g)\, pr(k|x=x)$$

for $g=2,$  $L(1,2)pr(1|x) + L(2,2)pr(2|x) + L(3,2)pr(3|x)$
$$= 0.8$$

$g=1,$  $L(1,1)pr(1|x) + L(2,1)pr(2|x) + L(3,1)pr(3|x)$

### Baiyes optimal classifier

$$\hat{f}(x) = \arg\max_{g} pr(g|x)$$

K-NN: (pick "k" nearest neighbour and take majority)

$\hat{f}(x) \approx \hat{y}$

Binary classifier

$(x_1, 0)$

$(x_2, 1)$

$(x_3, 0)$

$f(x) \geq 0.5$ {class-1}

$(x_1, 1)$

$< 0.5$ {class 0}

$(x_5, 0)$

Bias- variance    week-2 lecture7