

Home Work - 1

1)

$$\text{Cost function } L = (\hat{y}(\vec{x}) - y)^2$$

Average cost error squared function

$$E[L] = \iint (y - \hat{y}(\vec{x}))^2 \cdot p(\vec{x}, y) \cdot d\vec{x} dy$$

For finding optimal estimator, keep $\frac{\partial E(L)}{\partial \hat{y}(\vec{x})} = 0$

$$\Rightarrow \frac{\partial E(L)}{\partial \hat{y}(\vec{x})} = 2 \int (y - \hat{y}(\vec{x})) \cdot p(\vec{x}, y) \cdot dy = 0$$

Since $\hat{y}(\vec{x})$ is not dependent on y

$$\begin{aligned} \Rightarrow \int \hat{y}(\vec{x}) \cdot p(\vec{x}, y) \cdot dy &= \hat{y}(\vec{x}) \cdot \int p(\vec{x}, y) dy \\ &= \hat{y}(\vec{x}) \cdot p(\vec{x}) \end{aligned}$$

$$\Rightarrow \int y \cdot p(\vec{x}, y) dy = \hat{y}(\vec{x}) \cdot p(\vec{x}) = 0$$

$$\hat{y}(\vec{x}) = \frac{\int y \cdot p(\vec{x}, y) dy}{p(\vec{x})}$$

$$= \int y \cdot p(y|\vec{x}) dy$$

$$\therefore \hat{y}(\vec{x}) = E[y|\vec{x}]$$

2)

$(\hat{y}(\vec{x}) - y(\vec{x}))^2 \rightarrow$ squared error of arbitrary model.

* $\hat{y}(\vec{x})$ Defined estimator optimal $y^*(\vec{x}) = E(y|\vec{x})$

$$\rightarrow (\hat{y}(\vec{x}) - y(\vec{x}))^2 = (\hat{y}(\vec{x}) - y^*(\vec{x}) + y^*(\vec{x}) - y)^2$$

$$E[(\hat{y}(\vec{x}) - y)^2] = E[(\hat{y}(\vec{x}) - y^*(\vec{x}) + y^*(\vec{x}) - y)^2]$$

$$= E[(\hat{y}(\vec{x}) - y^*(\vec{x}))^2 + (y^*(\vec{x}) - y)^2 + 2(\hat{y}(\vec{x}) - y^*(\vec{x}))(y^*(\vec{x}) - y)]$$

$$\Rightarrow E[(\hat{y}(\vec{x}) - y)^2] = E[(\hat{y}(\vec{x}) - y^*(\vec{x}))^2] + E[(y^*(\vec{x}) - y)^2] + 2E[(\hat{y}(\vec{x}) - y^*(\vec{x}))(y^*(\vec{x}) - y)] \quad (1)$$

consider third term in (1)

$$2E[(\hat{y}(\vec{x}) - y^*(\vec{x}))(y^*(\vec{x}) - y)] = 2 \iint (\hat{y}(\vec{x}) - y^*(\vec{x}))(y^*(\vec{x}) - y) \cdot p(\vec{x}, y) d\vec{x} dy$$

$$= 2 \int (\hat{y}(\vec{x}) - y^*(\vec{x})) \cdot \int (y^*(\vec{x}) - y) \cdot p(\vec{x}, y) dy d\vec{x}$$

$$= 2 \int (\hat{y}(\vec{x}) - y^*(\vec{x})) \cdot \int (y^*(\vec{x}) - y) \cdot p(\vec{x}) \cdot p(y|\vec{x}) dy d\vec{x}$$

$$= 2 \int (\hat{y}(\vec{x}) - y^*(\vec{x})) \cdot \int (y^*(\vec{x}) \cdot p(y|\vec{x}) dy - y \cdot p(y|\vec{x}) dy) p(\vec{x}) d\vec{x}$$

$$= 2 \int (\hat{y}(\vec{x}) - y^*(\vec{x})) \left[(y^*(\vec{x}) - E(y|\vec{x})) \cdot p(\vec{x}) \right] d\vec{x}$$

$$= 0$$

$$E[(\hat{y}(\vec{x}) - y)^2] = E[(\hat{y}(\vec{x}) - y^*(\vec{x}))^2] + \underbrace{E[(y^*(\vec{x}) - y)^2]}_{\text{Noise}} \quad (2)$$

Let us consider, a case where model depends on dataset D .

$\hat{y}_D(\vec{x})$ be model from D .

From 2nd term in ②:

$$E_D[(\hat{y}_D(\vec{x}) - y^*(\vec{x}))^2] = ? \quad \text{--- ③}$$

$$\begin{aligned} (\hat{y}_D(\vec{x}) - y^*(\vec{x}))^2 &= (\hat{y}_D(\vec{x}) - E_D[\hat{y}_D(\vec{x})] + E_D[\hat{y}_D(\vec{x})] - y^*(\vec{x}))^2 \\ &= (\hat{y}_D(\vec{x}) - E_D[\hat{y}_D(\vec{x})])^2 + E_D[(E_D[\hat{y}_D(\vec{x})] - y^*(\vec{x}))^2] \\ &\quad + 2 \cdot (\hat{y}_D(\vec{x}) - E_D[\hat{y}_D(\vec{x})]) \cdot (E_D[\hat{y}_D(\vec{x})] - y^*(\vec{x})) \quad \text{④} \end{aligned}$$

Putting ④ in ③

$$\begin{aligned} E_D[(\hat{y}_D(\vec{x}) - y^*(\vec{x}))^2] &= E_D[(E_D[\hat{y}_D(\vec{x})] - y^*(\vec{x}))^2] + \\ &\quad E_D[(\hat{y}_D(\vec{x}) - E_D[\hat{y}_D(\vec{x})])^2] + \\ &\quad 2 \cdot E_D[(\hat{y}_D(\vec{x}) - E_D[\hat{y}_D(\vec{x})]) \cdot (E_D[\hat{y}_D(\vec{x})] - y^*(\vec{x}))] \end{aligned}$$

$$\Rightarrow E_D[(\hat{y}_D(\vec{x}) - y^*(\vec{x}))^2] = \underbrace{(E_D[\hat{y}_D(\vec{x})] - y^*(\vec{x}))^2}_{\text{(bias)}^2} + \underbrace{E_D[(\hat{y}_D(\vec{x}) - E_D[\hat{y}_D(\vec{x})])^2]}_{\text{Variance}}$$

$$\therefore E[(\hat{y}(\vec{x}) - y)^2] = (\text{bias})^2 + \text{variance} + \text{noise}.$$

3)

K-class linear discriminant classifier.

All the classes are distinguished by hyper plane.

$$\hat{y}_k(\vec{x}) = \vec{w}_k^T \cdot \vec{x} + w_{k0}$$

If $\hat{y}_k(\vec{x}) > \hat{y}_j(\vec{x})$ for all $j \neq k$, then output class is k

↳ decision rule.

$$X = \begin{bmatrix} x^{(1)T} \\ x^{(2)T} \\ \vdots \\ x^{(N)T} \end{bmatrix}_{N \times (d+1)}; \quad \vec{W} = \begin{bmatrix} w_0^{(1)} & w_0^{(2)} & \dots & w_0^{(k)} \\ \vdots & \vdots & \ddots & \vdots \\ w_d^{(1)} & w_d^{(2)} & \dots & w_d^{(k)} \end{bmatrix}_{(d+1) \times k}$$

$$\vec{y} = \begin{bmatrix} 0 & 0 & \dots & 1 & 0 & 0 \end{bmatrix}_{1 \times (d+1)}$$

↳ k^{th} class.

Sum of squared error $E(\vec{w}) = \frac{1}{2} (\vec{y} - \vec{W}^T X)^T (\vec{y} - \vec{W}^T X)$

$$E(\vec{w}) = (\vec{y} - \vec{W}^T X)^T \cdot (\vec{y} - \vec{W}^T X)$$

$$\frac{\partial E}{\partial \vec{w}} = 0 = \frac{\partial}{\partial \vec{w}} \left(\vec{y}^T \cdot \vec{y} - \vec{y}^T \cdot \vec{W}^T X - \vec{W}^T X^T \cdot \vec{y} - \vec{W}^T X^T X \vec{W} \right)$$

$$0 = 2 (\vec{X}^T \vec{X} \vec{W} - \vec{X}^T \vec{y})$$

$$\vec{W} = (\vec{X}^T \vec{X})^{-1} \cdot \vec{X}^T \vec{y}$$

↳ optimal solⁿ.

4) Fisher's Linear Discriminant of 2-class classifier can.

$C_0, C_1 \rightarrow$ two classes ; $x^{(i)}, y^{(i)} \rightarrow$ training data.

N - training samples

Projecting $x^{(i)}$ onto 1-dimension axis

$$z^{(i)} = w^T \cdot x^{(i)} \quad (w^T \rightarrow \text{projection matrix})$$

We would want the projected data to be as separable as possible when compared to variance.

\Rightarrow same as making the difference of respective classes maximum relative to variance.

$m_0 = w^T \cdot M_0$, $M_0, M_1 \rightarrow$ means in 2 dimension.

$m_1 = w^T \cdot M_1$, $m_0, m_1 \rightarrow$ means in 1 dimension.

$$\text{Variance } S_0^2 = \sum_{x_i \in C_0} (w^T \cdot x^{(i)} - m_0)^2$$

$$S_1^2 = \sum_{x_i \in C_1} (w^T \cdot x^{(i)} - m_1)^2$$

$$\text{Maximum } J(w) = \frac{(m_1 - m_0)^2}{S_0^2 + S_1^2}$$

$$J(w) = \frac{(w^T M_1 - w^T M_0) \cdot (w^T M_1 - w^T M_0)^T}{\sum_{x_i \in C_0} [w^T (x^{(i)} - M_0)]^2 + \sum_{x_i \in C_1} [w^T (x^{(i)} - M_1)]^2}$$

$$= J(w) = \frac{w^T \cdot (M_1 - M_0) \cdot (M_1 - M_0)^T \cdot w}{\sum_{x_i \in C_0} w^T \cdot (x^{(i)} - M_0) (x^{(i)} - M_0)^T \cdot w + \sum_{x_i \in C_1} w^T \cdot (x^{(i)} - M_1) (x^{(i)} - M_1)^T \cdot w}$$

$$\therefore J(w) = \frac{w^T \cdot S_B \cdot w}{w^T \cdot S_W \cdot w} \quad ; S_B, S_W \rightarrow \text{symmetric.}$$

$$\frac{\partial J(\vec{w})}{\partial w} = \frac{\frac{\partial}{\partial w} (w^T \cdot S_B w)}{w^T \cdot S_w \cdot w} - \frac{w^T \cdot S_B \cdot w}{(w^T \cdot S_w \cdot w)^2} \cdot \frac{\partial (w^T \cdot S_w \cdot w)}{\partial w}$$

$$\Rightarrow \frac{\partial J(\vec{w})}{\partial w} = \frac{(S_B + S_B^T) \cdot w}{w^T \cdot S_w \cdot w} - \frac{w^T \cdot S_B \cdot w}{(w^T \cdot S_w \cdot w)^2} (S_w + S_w^T) \cdot w$$

$$\frac{\partial J(\vec{w})}{\partial w} = 0$$

$$\Rightarrow (2 S_B w) \cdot (w^T \cdot S_w \cdot w) - (w^T \cdot S_B \cdot w) (2 S_w \cdot w) = 0$$

$$\Rightarrow (w^T \cdot S_w \cdot w) (S_B w) = (w^T \cdot S_B \cdot w) (S_w w)$$

$$\Rightarrow S_B w = \left(\frac{w^T \cdot S_B \cdot w}{w^T \cdot S_w \cdot w} \right) \cdot S_w \cdot w$$

$$S_B = (M_1 - M_0) \cdot (M_1 - M_0)^T \cdot w$$

$S_B \cdot w$ is in direction of $(M_1 - M_0)$

$$S_w \cdot w = \lambda \cdot (S_B w)$$

$S_w \cdot w$ is also in direction of $(M_1 - M_0)$

$$w = S_w^{-1} \cdot (M_1 - M_0)$$

Optimal solution.

$$[(M_1 - M_0)^T (M_1 - M_0)] w = [(M_1 - M_0)^T (M_1 - M_0)] w$$

$$w^T (M_1 - M_0) (M_1 - M_0)^T w = (w)^T I \cdot w$$

5) Recall $y^*(\vec{x}) = E[y|\vec{x}]$ is the optimal estimator

Average squared error $E[(y - \hat{y}(\vec{x}))^2]$

$$y^*(\vec{x}) = \arg \min_{\hat{y}(\vec{x})} E[L(y, \hat{y}(\vec{x}))]$$

loss function

zero-one loss function

$$L(y, \hat{y}(\vec{x})) = \begin{cases} 0 & , y = \hat{y}(\vec{x}) \\ 1 & , y \neq \hat{y}(\vec{x}) \end{cases}$$

$$E_{xy}[L(y, \hat{y}(\vec{x}))] = E_x \left[\sum_{y \in C_K} L(y, \hat{y}(\vec{x})) \cdot p(y = K|\vec{x}) \right]$$

$$y^*(\vec{x}) = \arg \min_{\hat{y}(\vec{x})} E_x \left[\sum_{y \in C_K} L(y, \hat{y}(\vec{x})) \cdot p(y = K|\vec{x}) \right]$$

$$= \arg \min_{\hat{y}(\vec{x})} E_x \left[\begin{matrix} L(y=1, \hat{y}(\vec{x})) \cdot p(y=1|\vec{x}) + \\ \vdots \\ L(y=K, \hat{y}(\vec{x})) \cdot p(y=K|\vec{x}) \end{matrix} \right]$$

Let $\hat{y}(\vec{x}) = K'$

$$y^*(\vec{x}) = \arg \min_{\hat{y}(\vec{x})} E_x [1 - p(y = K'|\vec{x})]$$

$$= \arg \min_{\hat{y}(\vec{x})} E_x [1] - E_x [p(y = K'|\vec{x})]$$

$$\therefore y^*(\vec{x}) = \arg \max_K p(y = K|\vec{x})$$