1.)

Loss $= E\left(\left(y(x) - \hat{y}(x)\right)^2\right)$ where

$\hat{y}(x)$ is the estimator

$\Rightarrow L = \iint \left(y(x) - \hat{y}(x)\right)^2 p(x,y) \, dy \, dx$

$\dfrac{\partial L}{\partial \hat{y}} = 0$

$\Rightarrow \displaystyle\int_x \left(\int_y 2(y-\hat{y}) \, p(x,y) \, dy\right) dx = 0$

$\Rightarrow \displaystyle\int_y 2\left(y-\hat{y}\right) p(x,y) \, dy = 0$

$\Rightarrow \displaystyle\int y \, p(x,y) \, dy = \hat{y}\left(\int p(x,y) \, dy\right)$

$\Rightarrow \left(p(x)\right) E(y|x) = \hat{y}\left(p(x)\right)$

$\Rightarrow \boxed{\hat{y}(x) = E(y|x)}$

**2.)** <u>Bias-Variance</u> :-

Say $\hat{y}_D(\vec{x})$ is estimator for the data $D$ & $y^* = E(y|\vec{x})$

$\Rightarrow$ The error is $E\left(\left(y - \hat{y}_D(\vec{x})\right)^2\right)$

$$E\left(\left(y - \hat{y}_D(\vec{x})\right)^2\right) = E\left(\left(y - y^* + y^* - \hat{y}_D(x)\right)^2\right)$$

$$= E\left(\left(y - y^*\right)^2 + \left(y^* - \hat{y}_D(x)\right)^2\right.$$
$$\left. + 2\left(y - y^*\right)\left(y^* - \hat{y}_D(x)\right)\right)$$

$$= E\left(\left(y - y^*\right)^2\right) + E\left(\left(y^* - \hat{y}_D\right)^2\right)$$

$$+ 2\iint \left(y - y^*\right)\left(y^* - \hat{y}_D\right)p(x,y)\,dy\,dx$$

$$\iint (y - y^*)(y^* - \hat{y}_D)\, p(x,y) = \int(y^* - \hat{y}_D)\left(\int(y - y^*)p(y|x)\,dy\right)dx$$

$$= \int(y^* - \hat{y}_D)\left(E(y|x) - y^*\right) dx \overset{0}{}$$

$$= 0$$

$\Rightarrow E\left(y - \hat{y}_D\right) =$ Noise $+ E\left(\left(y^* - \hat{y}_D\right)^2\right)$

Say $E_D(\hat{y}_D)$ is the mean of the estimators $\hat{y}_D$

Scanned by CamScanner

$$\Rightarrow E_D\left(\left(y^* - \hat{y}_D\right)^2\right)$$

$$= E_D\left(\left(y^* - E_D(\hat{y}_D) + E_D(\hat{y}_D) - y_D\right)^2\right)$$

$$= E_D\left(\left(y^* - E_D(\hat{y}_D)\right)^2\right)$$

$$+ E_D\left(\left(y_D - E_D(\hat{y}_D)\right)^2\right) \xrightarrow{\quad\quad} 0$$

$$+ 2E_D\left[\left(y^* - E_D(\hat{y}_D)\right)\left(y_D - E_D(\hat{y}_D)\right)\right]$$

$$\therefore E_D\left(y_D - E_D(\hat{y}_D)\right) = 0$$

$$= \underbrace{\left(y^* - E_D(\hat{y}_D)\right)^2}_{\text{Bias}}$$

$$+ \text{variance}$$

$$\therefore E\left(\left(y - \hat{y}_D\right)^2\right)$$

$$= (\text{Bias})^2 + \text{Variance} + \text{Noise}$$

3.)

$$L = \text{tr}\left[(\tilde{X}\tilde{W} - Y)^T(\tilde{X}\tilde{W} - Y)\right]$$

where

$\tilde{X}$'s $k^{th}$ column is $(1, x^T)^T$

& $\tilde{W}$'s " " " $(W_{ko}, w_k^T)^T$

$$\Rightarrow L = \sum_{\ell}\sum_{j}(\tilde{X}\tilde{W} - Y)^T_{\ell j} \cdot (\tilde{X}\tilde{W} - Y)_{j\ell}$$

$$= \sum_{\ell}\sum_{j}(\tilde{X}\tilde{W} - Y)^2_{j\ell}$$

$$= \sum_{\ell}\sum_{j}\left((\tilde{X}\tilde{W})_{j\ell} - Y_{j\ell}\right)^2$$

$$= \sum_{\ell}\sum_{j}\left(\sum_{k}\tilde{X}_{jk} \cdot W_{k\ell} - Y_{j\ell}\right)^2$$

$$\frac{\partial L}{\partial W_{k\ell}} = 0 \Rightarrow \sum_{\ell}\sum_{j}(\tilde{X}_{jk})\left(\sum_{k}\tilde{X}_{jk}\cdot W_{k\ell} - Y_{j\ell}\right) = 0$$

$$\Rightarrow \sum_{\ell}\sum_{j}(\tilde{X}^T)_{kj}\left((\tilde{X}W)_{j\ell} - Y_{j\ell}\right) = 0$$

$$\Rightarrow \tilde{X}^T(\tilde{X}W) - (\tilde{X})^T Y = 0$$

$$\Rightarrow \boxed{W = (\tilde{X}^T\tilde{X})^{-1}(\tilde{X})^T y}$$

$$\therefore \; w = \left( x^T x \right)^{-1} x^T y$$

The above is same as that multiple output case in Linear regression.

④. Fischer's Linear Discriminant :-

> $y = (\vec{w})^T \vec{x}$ is basically projecting a $(D+1)$ dimension vector to one dimension

$\vec{w}$ has many possibilities.

> we choose a $\vec{w}$ such that the intra-class variance is minimized & inter-class is maximized.

> Consider 2 classes $C_1$, $C_2$ whose means are given by

$$\vec{m_1} = \frac{1}{N_1} \sum_{n \in C_1} \vec{x_n}$$

$$\vec{m_2} = \frac{1}{N_2} \sum_{n \in C_2} \vec{x_n}$$

we would like to choose a vector $\vec{w}$ such that

$$m_2 - m_1 = \vec{w}^T \left( \vec{m_2} - \vec{m_1} \right) \text{ is}$$

maximized.

> This can be done by having arbitrarily large $w$ (which is not preferred, cause it might lead to overfit.)

So we constrain $w$ to have $\leq$ unit length i.e., $\sum_i w_i^2 \leq 1$ $\Rightarrow \vec{w} \propto \left( \vec{m_2} - \vec{m_1} \right)$

> the within class variance is given by

$$s_k^2 = \sum_{y_n \in C_k} (y_n - m_k)^2$$

• where
$$y_n = \vec{w}^T \vec{x}$$
$$m_k = \vec{w}^T \vec{m}$$

> The fisher criterion is defined as

$$J(\vec{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

But

$$m_2 - m_1 = w^T(\overline{m}_2 - \overline{m}_1) = (\overline{m}_2 - \overline{m}_1)^T w$$

$$\Rightarrow (m_2 - m_1)^2 = \left(w^T(\overline{m}_2 - \overline{m}_1)\right)\left(w^T(\overline{m}_2 - \overline{m}_1)\right)$$

$$\Rightarrow (m_2 - m_1)^2 = w^T \boxed{(\overline{m}_2 - \overline{m}_1)(\overline{m}_2 - \overline{m}_1)^T} w$$

$$\downarrow$$

$$S_B$$

$$S_1^2 + S_2^2 = w^T S_w w$$

where

$$\boxed{S_w = \sum_{n \in C_1}(x_n - m_1)(x_n - m_1)^T + \sum_{n \in C_2}(x_n - m_2)(x_n - m_2)^T}$$

$$\downarrow$$

within class variance of full data

$$\therefore J(\vec{w}) = \frac{w^T S_B w}{w^T S_w w}$$

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

$$\Rightarrow \left(w^T S_w w\right) S_B w = \left(w^T S_B w\right) S_w w$$

we fixed $|w| \leq 1$

So dropping the scalars

$$\Rightarrow \left(\overline{m}_2 - \overline{m}_1\right)\underbrace{\left(\overline{m}_2 - \overline{m}_1\right)^T w}_{\text{constant}} = S_w w$$

$$\Rightarrow \boxed{w \propto S_w^{-1}\left(\overline{m}_2 - \overline{m}_1\right)}$$

# KNN :-

› What was observed is that data is distributed almost similarly about origin.

› So if the test sample is close to origin, we need to have a bigger $k$ value to make a concrete prediction.

5.)

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

$$E_{xy}\left(L(y, \hat{y}(\vec{x}))\right)$$

$$= E_x\left[\sum_{y \in C_k} L(y, \hat{y}(\vec{x})) \cdot P(y=k|\vec{x})\right]$$

$$\because E_{xy}(f(x,y)) = E_x\left(E_{y|x}(f(x, y))\right)$$

$\Rightarrow$ we need to find $\hat{y}(\vec{x}) = y^*$ such that the above expectation is minimized.

$$\sum_{y \in C_k} L(y, \hat{y}(\vec{x})) \cdot P(y=k|\vec{x})$$

$\qquad$ If $\quad y = \hat{y}(\vec{x})$, $\quad$ then $L(y, \hat{y}) = 0$

$\Rightarrow$ Summation turns out to be

$$\sum_{y \in C_k \; ; \; y \neq \hat{y}} P(y=k|\vec{x}) = 1 - P(\hat{y}=k|\vec{x})$$

$$\Rightarrow y^* = \underset{\hat{y}(\vec{x})}{\arg\min} \; E_X\left(1 - P\left(y = \hat{y}(\vec{x}) \mid \vec{x}\right)\right)$$

Say $\hat{y} = k$

This is same as maximizing

$$P(y = k \mid \vec{x}) \quad \text{for} \quad y \in C_k$$

i.e.,

$$\boxed{y^* = \underset{y \in C_k}{\arg\max} \; P(y = k \mid \vec{x})}$$