Nakka Chaluadri
EE16BTECH 11022

① $$p(x,y) = p(y|x) \, p(x) = p(x|y) \, p(y)$$

by Bayes Theorem

And, in the statistical setting, we aim to get an estimate for the absolute best model $\hat{y}(\vec{x})$ for SSE cost func

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

$$E(\vec{x}) = \iint \left(y(\vec{x}) - \hat{y}(\vec{x})\right)^2 p(x,y) \, d\vec{x} \, dy$$

for the best $\hat{y}(\vec{x})$, minimize $\vec{E}(\vec{x})$ wrt $\hat{y}$

$$\frac{d}{d\hat{y}(\vec{x})} E(\vec{x}) = 0$$

$\Rightarrow$
$$\frac{d}{d\hat{y}(\vec{x})} \iint (y - \hat{y})^2 \, p(x,y) \, d\vec{x} \, dy = 0$$

$$\Rightarrow \iint -2(y - \hat{y}) \, p(x,y) \, d\vec{x} \, dy = 0$$

$$\rightarrow \iint y\, p(x,y)\, dx\, dy = \iint \hat{y}\, p(x,y)\, dx\, dy$$

$$p(x,y) = p(y|x)\, p(x)$$

$$\iint y\, p(y|\vec{x})\, p(\vec{x})\, (y - \hat{y})\, d\vec{x}\, dy = 0$$

$$\iint \left( y\, p(y|\vec{x})\, dy \right) p(x)\, dx = \iint \hat{y}\, p(x,y)\, dx\, dy$$

$$= \iint E(y|x)\, p(x)\, dx = \hat{y} \iint p(x,y)\, dx\, dy$$

$E(y|x)$ is independent of $\vec{x}$

$\iint p(x,y)\, dx\, dy = 1$  &  $\int p(x)\, dx = 1$

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

② Let $y$ be the labels for the given data

i.e., $y(x)$

$\hat{y}(x)$ is the arbitrary model's label estimate

$y^*(x)$ is the absolute best model for the given data

then :-

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

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

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

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

Consider $E\left((y-y^*)(y^*-\hat{y})\right)$

$$= \iint (y-y^*)(y^*-\hat{y}) \, p(x,y) \, dx \, dy$$

$$= \int_x \left[\int_y (y-y^*)(y^*-\hat{y}) \, p(y|x) \, dy\right] p(x) \, dx$$

$$= \int_x (y^*-\hat{y}) \left[\int_y y - y^* \, p(y|x) \, dy\right] p(x) \, dx$$

$\int p(y|a)\, dy = 1$ & $y^*$ is independent of $y$

$$= \int_x y^* - \hat{y}\left[\int (y\, p(y|a)\, dy) - y^*\right] p(a)\, da$$

& $\int y\, p(y|a)\, dy = y^*$

$\rightarrow E\left((y-y^*)(y^* - \hat{y})\right) = 0.$

Consider $E\left((y^* - \hat{y})^2\right)$

We actually have access to say $D$ datasets of the same experiment. Then

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

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

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

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

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

Again consider $E_D\left[\left(\hat{y}_D - E_D(\hat{y}(a))\right)\left(E_D(\hat{y}(a)) - y^*\right)\right]$

$$= \int\left(\hat{y}_D - E_D(\hat{y}(a))\right)\left(E_D(\hat{y}(a)) - y^*\right) p(a, \hat{y})\, da\, d\hat{y}$$

$$= \int\left(E_D(\hat{y}(a)) - y^*\right)\int\underbrace{\left(\hat{y}_D - E_D(\hat{y}(a))\right)}_{\nearrow 0}\, p(y)\, dy \cdot p(a|y)\, da$$

$$= 0$$

$$\Rightarrow \quad E_D\big((y-\hat{y})^2\big) = E\big((y-y^*)^2\big) + E_D\big((\hat{y}-E_D(\hat{y}))^2\big)$$

$$+ E_D\big((E_D(\hat{y})-y^*)^2\big)$$

$$\left\{ \begin{array}{l} \text{here} \quad E\big((y-y^*)^2\big) = \text{noise} \\[2ex] E\big[(\hat{y}-E_D(\hat{y}(x)))^2\big] = \text{Variance} \\[2ex] E\big[(E_D(\hat{y})-y^*)^2\big] = \text{Bias.} \end{array} \right\}.$$

⑤ $\Delta f = \begin{cases} 0 & \hat{y} = y \\ 1 & else \end{cases}$   $\hat{y} \rightarrow$ estimate
   $y \rightarrow$ label

for the statistical setting

$\begin{cases} \rightarrow p(y|x) \text{ must be maximized} \\ \quad p(y|x) \propto p(x|y)\, p(y) \end{cases}$   $\begin{cases} \text{We know the} \\ \text{samples follow} \\ \text{some distribution.} \end{cases}$

and for optimal solution

$y^*$ for which $\underset{\hat{y}}{\text{argmin}} \left( \mathbf{E}\, \Delta f(y,\hat{y}) \right)$

$$= E_{XY}\left( \Delta f(y,\hat{y}) \right) = E_X\left[ E_{Y|X}\, \Delta f(y,\hat{y}) \right]$$

$$= E_X\left[ \int \Delta f(y,\hat{y})\, p(y|x)\, dy \right]$$

now since discrete, $\int \rightarrow \Sigma$

$$= E_X\left[ \sum_{y \in C_x} \Delta f(y,\hat{y})\, p(y|x) \right]$$

Let $\hat{y} = k'$  $k' \in C_x$

$$= E_X\left[ \Delta f(y,k')\, pr(y=k'|x) + \Delta f(y,k')\, pr(y=2|x) \right. $$
$$\left. + \cdots \right]$$

$$= E_X\left[ \sum_{y \in C_x} pr(y=k) - pr(y=k'|x) \right]$$

$$= E_x \left( 1 - P_r(y=k'|x) \right)$$

for optimal solutions

$$E_x \left[ 1 - P_r(y=k'|x) \right] \text{ is minimized}$$

$$\Rightarrow \quad P_r(y=k'|x) \text{ is maximized}$$

$$\therefore \quad y^*(\vec{x}) = \underset{y=c_k}{\arg\max} \; p\left(y=k|\vec{x}\right)$$

which is intuitive and what we try to do in each and every classification problem → maximize the prob. of a label given $\vec{x}$

⑧ K-Class Linear Discriminant Classifier

Each class is seperated by a hyperplane

$$y_k(x) = w_k^T x + w_{k_0}$$

Note: Till now we've been using $\hat{y} = xW$

$$x \rightarrow \begin{bmatrix} x_0^1 & x_1^1 & \cdots & x_d^1 \\ \vdots & & & \end{bmatrix}$$

here $x = \begin{bmatrix} x_0^1 \\ x_1^1 \\ \vdots \\ x_d^1 \end{bmatrix}$ &

$$X = \begin{bmatrix} x_0^1 & x_1^1 & \cdots & x_d^1 \\ x_0^2 & x_1^2 & \cdots & x_d^2 \\ \vdots & & & \end{bmatrix}$$

$x \rightarrow$ Column vector of components

$X \rightarrow$ whole data in standard format

To merge bias term, consider

$$\tilde{W} = \begin{bmatrix} W_{10} & W_{20} & - & W_{N0} \\ W_{11} & W_{21} & & W_{N1} \\ \vdots & \vdots & \ddots & \vdots \\ W_{1d} & W_{2d} & - & W_{Nd} \end{bmatrix}$$

$\Rightarrow$ Cost function $= \mathbb{E}\left((y - \hat{y})^2\right)$

$= $ ~~[scribbled out]~~ Consider $t_n$ as the

estimate vector for $x_n$ which will be of a

column of K elements, each element being either 0 or 1

(9) Consider the SSE cost function

$$CF = E(\omega) = \sum_{i=1}^{N} \left[ y^{(i)} - \hat{y}^{(i)} \right]^2$$

for $i^{th}$ row of $y$ and $k^{th}$ term (dimension)

$$CF_{ik} = \left( y_k^{(i)} - \sum_{j=0}^{d} x_j^{(i)} \omega_{jk} \right)^2$$

for optimal $\omega^*$

$$\frac{\partial CF_{ik}}{\partial \omega_{jk}} = 0 \Rightarrow -2x_j^{(i)} \left( y_k^{(i)} - \sum_{j=0}^{d} x_j^{(i)} \omega_{jk} \right) = 0$$

Ref for
multidimensional
labels proof

→ In matrix format

$$X^T \left( Y - X\omega^* \right) = 0$$

which yields $\boxed{\omega^* = (X^T X)^{-1} X^T Y}$

In case of basis functions;

$$\boxed{\omega^* = (\phi^T \phi)^{-1} \phi^T Y}$$

their orders being $d+1 \times K$ and

$(M+1) \times K$ respectively

# Fischer's Linear Discriminant

Consider two class classification problem.

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

$\vec{m_2} = \dfrac{1}{N_1} \sum_{n \in C_2} \vec{x_n}$

The main idea of Fischer's LDA is to prevent overlap of classes during dimensionality reduction

To get to a single dimension from multiple dimension vector, consider its component on a vector W of the same dimensionality

$$m_2 - m_1 = W^T \left( \vec{m_2} - \vec{m_1} \right)$$

$$m_k = W^T \vec{m_k}$$

To restrict the ~~vector~~ itself, consider $\|w\| = 1$

To reduce ~~~~ ~~~~ intra-class covariance and increase inter-class variance, ~~Fische~~ Fischer's criterion is used.

$$J(w) = \frac{(m_2 - m_1)^2}{S_1^2 + S_2^2}$$

$$(m_2 - m_1)^2 = \left(w^T \vec{m_k}\right)^2$$

$$= w^T \vec{m_k} \, \vec{m_k}^T w .$$

$$\vec{m_k} \, \vec{m_k}^T = S_B = \text{between class covariance matrix}$$

$||^y \qquad S_1^2 + S_2^2$

$$= \sum_{n \in G} (y_1 - m_1)^2 + \sum_{n \in G_2} (y_2 - m_2)^2$$

$$= \sum_{n \in G} (w^T x_n - m_1)^2 + \sum_{n \in G_2} (w^T x_n - m_2)^2$$

$$= \sum_{n \in G} (w^T x_n - m_1)(w^T x_n - m_1)^T + \sum_{n \in G_2} (w^T x_n - m_2)(w^T x_n - m_2)^T$$

$$= w^T S_w w$$

$$\boxed{S_w = \sum_{n \in G} (\vec{x}_n - \vec{m}_1)(\vec{x}_n - \vec{m}_1)^T + \sum_{n \in G_2} (\vec{x}_n - \vec{m}_2)(\vec{x}_n - \vec{m}_2)^T}$$

(within class variance)

Note that $S_B \vec{w} = (\vec{m}_2 - \vec{m}_1)(m_2 - m_1)^T w$

$$= (\vec{m}_2 - \vec{m}_1) \underbrace{\left( m_k^T \right)}_{\text{scalar}}$$

$$\Rightarrow \boxed{S_B(\vec{w}) \text{ is in direction of } (\vec{m}_2 - \vec{m}_1)}$$

for best discrimination

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

$$\frac{d}{dw}\left(\frac{\vec{w}^{T}S_{B}\vec{w}}{\vec{w}^{T}S_{w}\vec{w}}\right)=0$$

$$\Rightarrow \left(w^{T}S_{w}w\right)\left(S_{B}w\right)=\left(w^{T}S_{B}w\right)\left(S_{w}w\right)$$

$$S_{B}\vec{w}\propto \vec{m_{2}}-\vec{m_{1}}$$

$$w^{T}S_{w}w \ \& \ w^{T}S_{B}w \text{ are scalars}$$

$$\Rightarrow \boxed{w=S_{w}^{-1}\left(\vec{m_{2}}-\vec{m_{1}}\right)}$$

which is known as fischer's linear discri-
minant