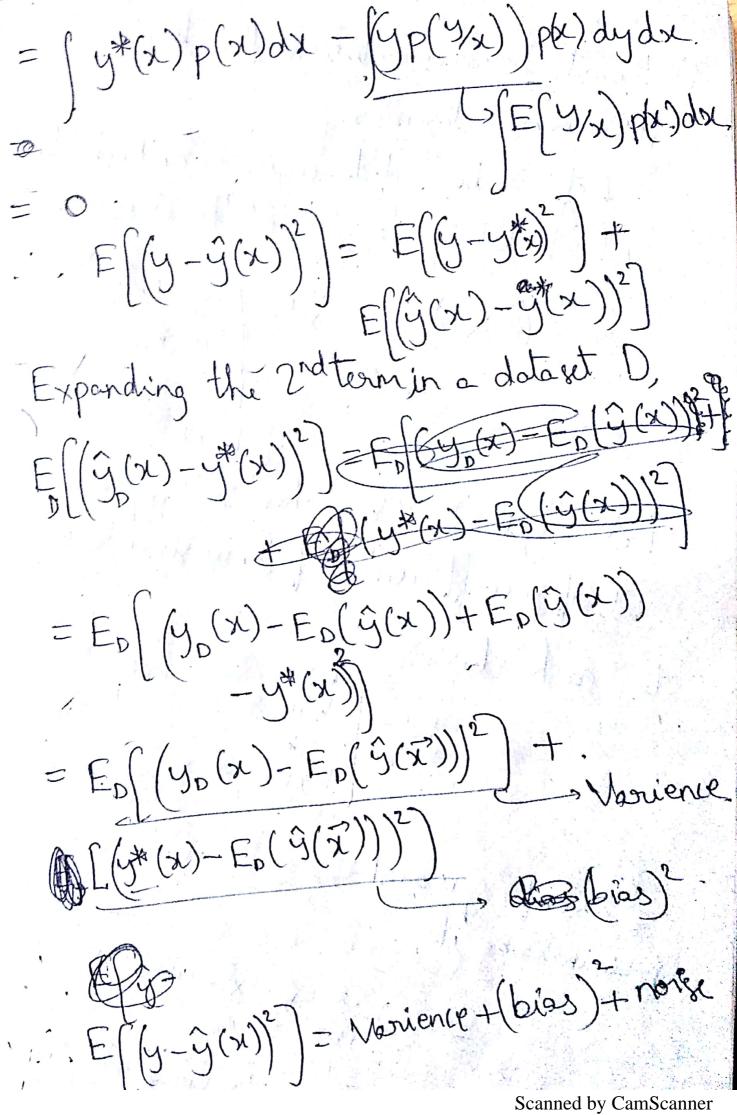
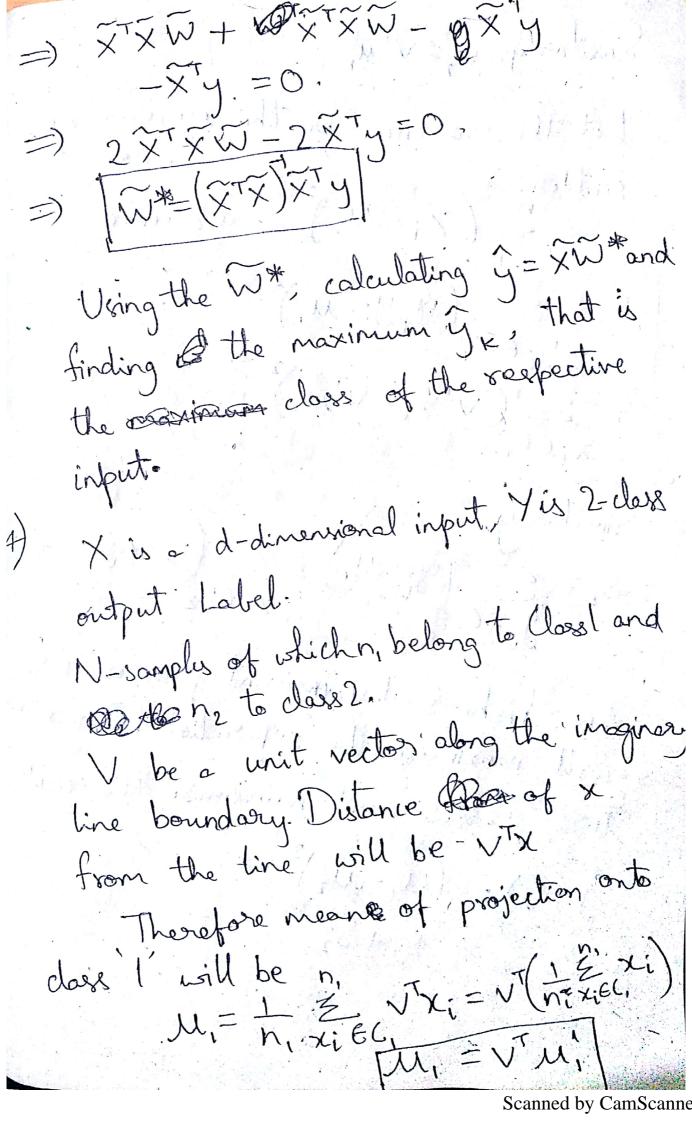
We know that 1= 19-912  $E(L) = \iint E(y-\hat{y})^2 p(x,y) dxdy.$ =  $\int \int |y-\hat{y}(x)|^2 p(x,y) dxdy.$ For the Averaged loss to be minimum. differentiale E(L) us. st.  $\hat{y}(x)$ .  $\frac{\partial E(L)}{\partial \hat{g}(x)} = 2 \int (y - \hat{g}(x)) p(x,y) dx dy$  $=) \int y p(x,y) dy - \int \hat{y}(x) p(x,y) dy = 0$  $=\int \int y p(x,y) dy - \int \hat{y}(x) p(x,y) dy = 0.$ 

=) 
$$\hat{g}(x)p(x) = [y p(x,y)dy]$$
  
=)  $\hat{g}(x)$  =  $[y p(x,y)]$  dy  
=)  $[\hat{g}(x)] = E[(y-y^* + y^* - \hat{g})^2]$   
= $[(y-y^*)^2] + E[(y^* - \hat{g})^2] + E[(y^* - \hat{g})^2]$   
Simplifying the third term,  
 $E[(y-y^*)(y^* - \hat{g})] = [[\hat{g}(x) - y^*(x)], (y^*(y)y)]$   
=  $[(\hat{g}(x) - y^*(x))](y^*(x) - y)$  p(x,y) dix dy  
= $[(\hat{g}(x) - y^*(x))](y^*(x) - y)$  p(x,y) dix dy  
= $[(y^*(x) - y), (y^*(x) - y), (y^*(x)) + y^*(x)]$ 



3) Least Squares Solution for k-class discriminant classifier. Let X be d-dimensional and y be labels which belong to one of the classes.  $\widetilde{X} = (L \times T)^{T}$   $\widetilde{X}^{T} = \begin{pmatrix} 1 & 1 & \dots & 1 \\ \chi_{11} & \chi_{21} & \dots & \chi_{NI} \end{pmatrix}$ (Xid - - - - Xnd) Let  $\hat{y}(x) = XW$  when  $W^T$  represents the set of coefficients of dimensions k wirt to each class. Sun of squares. F-tru (XW-y) T(XW-y)
Minimising it,  $\frac{\partial E}{\partial w} = 0. \Rightarrow \frac{\partial}{\partial w} (\widetilde{w} \cdot \widetilde{x} - y) (\widetilde{x} \cdot \widetilde{w} - y)^{=0}$ 

Scanned by CamScanner



Scanned by CamScanner

Similarly, u, = VTu Let the vorience be the measure of S= & (Zi-Mz) Vorience S.= ( ( 9 yi-u, )  $= \underbrace{\sum_{x \in C} (v^{T}x_{i} - v^{T}u_{i}^{T})^{2}}_{x \in C}$ Similarly,  $S_2 = \underbrace{\xi}_{\text{titel}} \left( y_i - \lambda l_z \right) = \underbrace{\xi}_{\text{titel}} \left( \sqrt{\chi_i} - \sqrt{\lambda l_z} \right)$ In Fisher's discreninant process, we will maximise the separation between the tuso classes and minimize the scattering  $J(v) = \frac{(u_1 - u_1)^2}{2}$ 5,+52

Scanned by CamScanner

Let 
$$8_1 = \underbrace{\sum_{x_i \in C_1} (x_i - u_i^*)^* (x_i - u_i^*)^*}_{x_i \in C_1} (x_i - u_i^*)^* (x_i - u_i^*)^*$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

Similarly,  $\underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* (x_i - u_i^*)^* }_{x_i \in C_1} (x_i - u_i^*)^* V.$ 

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

Similarly,  $\underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$ 

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

Similarly,  $\underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$ 

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

Similarly,  $\underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$ 

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

Similarly,  $\underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$ 

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

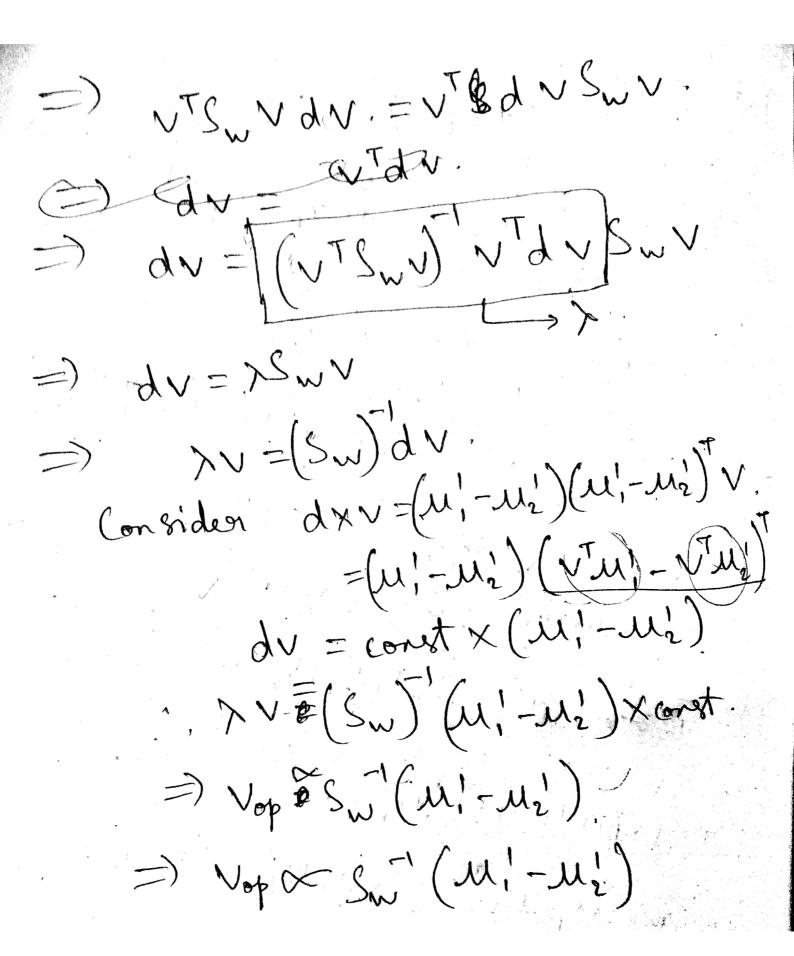
$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^* V.$$

$$= \underbrace{\sum_{x_i \in C_2} (x_i - u_i^*)^* }_{y \in C_1} (x_i - u_i^*)^*$$



$$L(\hat{y}, y) = \begin{cases} 0 & y = \hat{y} \\ \text{else} \end{cases}$$

$$E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ E \times \left[ L(\hat{y}, \hat{y}) \right] \right]$$

$$= E \times \left[ E \times \left[ L(\hat{y}, \hat{y}) \right], p(\hat{y} = \hat{y} / x) \right]$$

$$= E \times \left[ E \times \left[ L(\hat{y}, \hat{y}) \right], p(\hat{y} = \hat{y} / x) \right]$$

$$= E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ (1 - p(\hat{y} = \hat{y} / x)) \right]$$

$$= E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ (1 - p(\hat{y} = \hat{y} / x)) \right]$$

$$= E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ (1 - p(\hat{y} = \hat{y} / x)) \right]$$

$$= E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ (1 - p(\hat{y} = \hat{y} / x)) \right]$$

$$= E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ (1 - p(\hat{y} = \hat{y} / x)) \right]$$

$$= E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ (1 - p(\hat{y} = \hat{y} / x)) \right]$$

$$= E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ L(\hat{y}, \hat{y}) \right]$$

$$= E \times \left[ E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ L(\hat{y}, \hat{y}) \right]$$

$$= E \times \left[ E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ L(\hat{y}, \hat{y}) \right]$$

$$= E \times \left[ E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ L(\hat{y}, \hat{y}) \right]$$

$$= E \times \left[ E \times \left[ L(\hat{y}, \hat{y}) \right] = L \times \left[ L(\hat{y}, \hat{y}) \right]$$

$$= E \times \left[ E \times \left[ L(\hat{y}, \hat{y}) \right] = E \times \left[ L(\hat{y}, \hat{y}) \right] = L \times \left[ L($$