$$\frac{H\omega^{-2}}{2}$$

$$\frac{5.5a}{600} \text{ Political}$$

$$\frac{1}{600} = E\left[\left(y - \hat{y} - \hat{y} - \hat{y}\right)^{2}\right]$$
where $\hat{y} = (y - \hat{y} - \hat{y})^{2}$ Proposition

$$L = \left[\left(y - \hat{y} - \hat{y} - \hat{y}\right)^{2}\right] + \left[y - \hat{y} - \hat{y} - \hat{y}\right] + \left[y -$$

Assuming,
Adata set D and your in the estimator associated with the data set D

The erotor is
$$\left[\mathbb{E} \left[\left(\mathbf{j} - \hat{\mathbf{y}}_{0} \hat{\mathbf{n}} \right) \right]^{2} \right]$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbf{j}_{0} \hat{\mathbf{n}} - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right] \right]^{2}$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) \right]^{2} + \mathbb{E} \left[\left(\mathbf{j}_{0} \hat{\mathbf{n}} \hat{\mathbf{j}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right] \right]^{2}$$

$$+ 2 \int \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) \left(\mathbf{j}_{0} \hat{\mathbf{n}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right) \right] P(\hat{\mathbf{n}}, \mathbf{j}) d\mathbf{j} d\mathbf{n}$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) \right] \left(\mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j}_{0} \hat{\mathbf{n}} \hat{\mathbf{j}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right] \right]$$

$$+ 2 \int \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) \left(\mathbf{j}_{0} \hat{\mathbf{n}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right] P(\hat{\mathbf{n}}, \mathbf{j}) d\mathbf{j} d\mathbf{n}$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) \right] \left(\mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j}_{0} \hat{\mathbf{n}} \hat{\mathbf{j}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right] \right]$$

$$+ 2 \int \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) \left(\mathbf{j}_{0} \hat{\mathbf{n}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right) \right] P(\hat{\mathbf{n}}, \mathbf{j}) d\mathbf{j} d\mathbf{n}$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right) \right] P(\hat{\mathbf{n}}, \mathbf{j}) d\mathbf{j} d\mathbf{n}$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right] \right]$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right] \right]$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right] \right]$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right] \right]$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right] \right]$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) - \hat{\mathbf{j}}_{0} \hat{\mathbf{n}} \right] \right]$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) - \hat{\mathbf{j}_{0}} \hat{\mathbf{n}} \right] \right]$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) - \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) \right]$$

$$= \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) + \mathbb{E} \left[\left(\mathbf{j} - \mathbf{j}_{0} \hat{\mathbf{n}} \right) - \mathbb$$

There is a trade of when we reduce bias, variance shoots and vice venser. This is known as bias variance tradeoff.

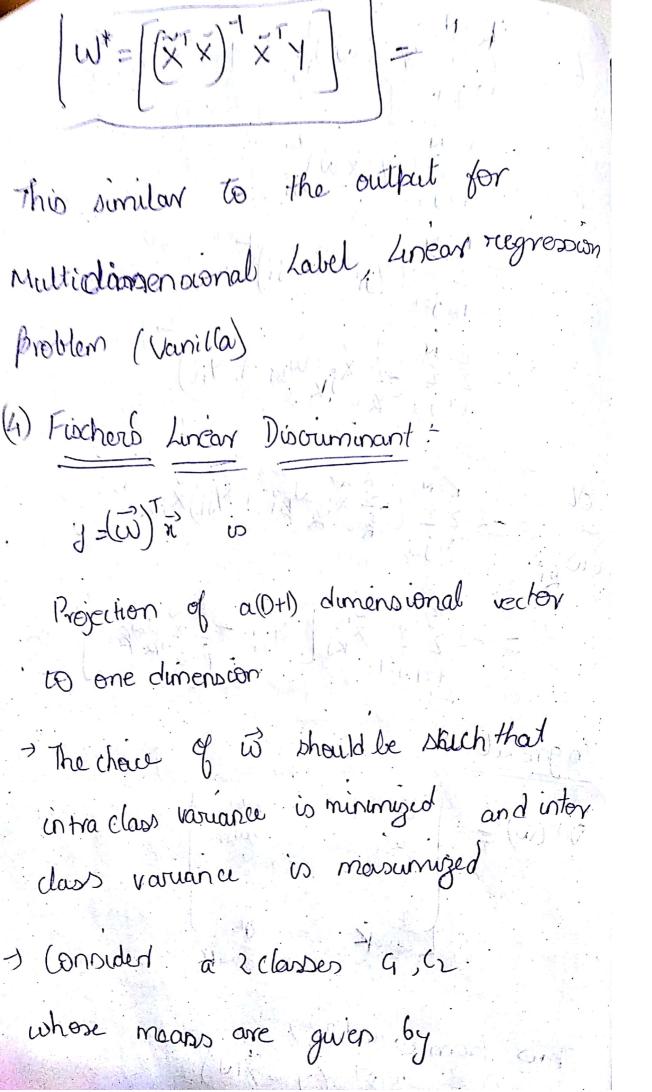
(3)
$$L = Tr \left[Y - \widetilde{X}\widetilde{W} \right] \cdot \left(\widetilde{M} - \widetilde{X}\widetilde{W} \right) \right]$$

where $\widetilde{X} = \left[\begin{array}{c} 1 & X^{(1)} & X^{(2)} & X^{$

Scanned by CamScanner

$$\begin{array}{lll}
\mathcal{L} &= & \left[\sum_{i=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{i}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{i}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{x} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N} \left(\mathbf{y} - \mathbf{y} \, \widetilde{\omega} \right) \right]_{j}^{T} \\
\mathcal{L} &= & \left[\sum_{i=1}^{N} \sum_{j=1}^{N$$

Scanned by CamScanner



m2 = 1. 2 xn we would like to choose a vector w such that $m_2-m_1=wT(m_2-m_1)$ (is maximized. This can be done by having arbitrorily large w which is not preferred because it might kad to overfitting So we constrain to to have curit lingth ie, [[[[]] The within class variance is given by $5_{R}^{2} = \sum_{y_{n} \in C_{k}} (y_{n} - m_{k})^{2}$

where
$$y_n = \vec{w}^T \vec{n}$$
 $m_{1c} = \vec{w}^T \vec{m}$

fiver outerion

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

Now,

 $(m_1 - m_1)^2 = \vec{w}^T (\vec{m}_1 - \vec{m}_1) = (\vec{m}_2 - \vec{m}_1)^T \vec{w}$

Since its a scalar $(\vec{a} - \vec{a})$ where a isocalar.

 $(n_1 - m_1)^2 = \vec{w}^T (m_2 - m_1) (m_1 - m_1)^T \vec{w}$

where

 $s_3 = (m_1 - m_1) (m_1 - m_1)^T = (ehiverantlass)$
 $constructive Mahrlip$
 $s_1^2 + s_2^2 = \sum_{n \in G} (y_1 - m_1)^T + \sum_{n \in G} (\vec{w}_{1n} - m_1)^T + \sum_{n \in G} (\vec{w}_{2n} - m_1)^T + \sum_{n \in G} (\vec{w}_{2n$

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

Eng $L(y,\hat{y}) = E_x \left[E_{Y|x}(L(y,\hat{y},\hat{y})) \right]$

Power of the first state o