

Extension to M-ary Bayes Classifier

Lec 3

$$H_y : \underline{x} \sim p(\underline{x}|y), y=1, \dots, M$$

$$\text{priors } \pi_y = P\{Y=y\}, y=1, \dots, M$$

$$\begin{aligned}\hat{y}_{\text{Bayes}} &= f_{\text{Bayes}}(\underline{x}) = \arg \max_y \pi_y p(\underline{x}|y) \\ &= \arg \max_y (\ln \pi_y + \ln p(\underline{x}|y))\end{aligned}$$

f_{Bayes} minimizes P_e (\equiv Err_{pred})

Multivariate Gaussian Distribution

Recall : Gaussian random variable X with df

$$f_x(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad \begin{matrix} \text{mean} \\ \text{variance} \end{matrix} \quad X \sim N(\mu, \sigma^2)$$

Generalization to d-dimensional vector \underline{x} :

$$\phi(\underline{x}, \underline{\mu}, C) = \frac{1}{(2\pi)^{\frac{d}{2}} \sqrt{|C|}} e^{-\frac{1}{2} (\underline{x}-\underline{\mu})^T C^{-1} (\underline{x}-\underline{\mu})}$$

mean of \underline{x} Covariance matrix of \underline{x}

$$\underline{\mu} = E[\underline{x}], \quad C = E[(\underline{x}-\underline{\mu})(\underline{x}-\underline{\mu})^T]$$

$$\mu_e = E[x_e], \quad C_{e,m} = E[(x_e - \mu_e)(x_m - \mu_m)]$$

Compact Notation : $\underline{x} \sim N(\underline{\mu}, C)$

M-ary Classification with Multivariate Gaussians

$$H_y : p(\underline{x} | y) = \phi(\underline{x}, \underline{\mu}_y, C_y)$$

$$\hat{y}_{\text{Bayes}} = f_{\text{Bayes}}(\underline{x}) = \arg \max_y (\ln \pi_y + \ln p(\underline{x} | y))$$

$$p(\underline{x} | y) = \frac{1}{(2\pi)^{\frac{d}{2}} \sqrt{|C_y|}} e^{-\frac{1}{2} (\underline{x} - \underline{\mu}_y)^T C_y^{-1} (\underline{x} - \underline{\mu}_y)}$$

$$\ln p(\underline{x} | y) = -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |C_y| - \frac{1}{2} (\underline{x} - \underline{\mu}_y)^T C_y^{-1} (\underline{x} - \underline{\mu}_y)$$

$$\hat{y}_{\text{Bayes}} = \arg \max_y \underbrace{\ln \pi_y - \frac{1}{2} \ln |C_y| - \frac{1}{2} (\underline{x} - \underline{\mu}_y)^T C_y^{-1} (\underline{x} - \underline{\mu}_y)}$$

$\delta_y(\underline{x})$: objective function

If C_y 's are different, $\delta_y(\underline{x})$ is quadratic in \underline{x}

If $C_y = C$, for all y , then

$$\frac{1}{2} (\underline{x} - \underline{\mu}_y)^T C_y^{-1} (\underline{x} - \underline{\mu}_y) = \frac{1}{2} \left(\underline{x}^T C^{-1} \underline{x} - \stackrel{\uparrow}{\underline{\mu}_y^T C^{-1} \underline{x}} - \stackrel{\leftarrow}{\underline{x}^T C^{-1} \underline{\mu}_y} + \stackrel{\rightarrow}{\underline{\mu}_y^T C^{-1} \underline{\mu}_y} \right)$$

not fn. of y same

$$\Rightarrow \hat{y}_{\text{Bayes}} = \arg \max_y \underbrace{[\ln \pi_y + \underline{\mu}_y^T C^{-1} \underline{x} - \frac{1}{2} \underline{\mu}_y^T C^{-1} \underline{\mu}_y]}_{\text{linear objective } \delta_y(\underline{x})}$$

Linear Discriminant Analysis (LDA) Classifier

- Assuming $C_y = C$, for all y ,

$$f_{\text{Bayes}}(\underline{x}) = \arg \max_y \left[\ln \pi_y + \underline{\mu}_y^T C^{-1} \underline{x} - \frac{1}{2} \underline{\mu}_y^T C^{-1} \underline{\mu}_y \right]$$

- LDA approach replaces $\{\pi_e, \underline{\mu}_e\}_{e=1}^M$ and C by estimates $\{\hat{\pi}_e, \hat{\mu}_e\}_{e=1}^M$ and \hat{C} that are obtained from training data $\mathcal{T} = \{\underline{x}_i, y_i\}_{i=1}^N$

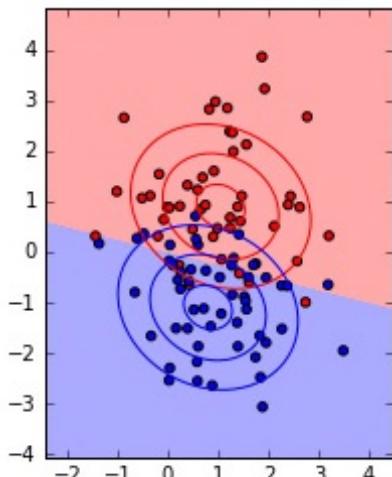
Estimation Let $N_e = \# \text{ training data with class label } e$

$$\hat{\pi}_e = \frac{N_e}{N}, \quad \hat{\mu}_e = \frac{\sum_{i:y_i=e} \underline{x}_i}{N}$$

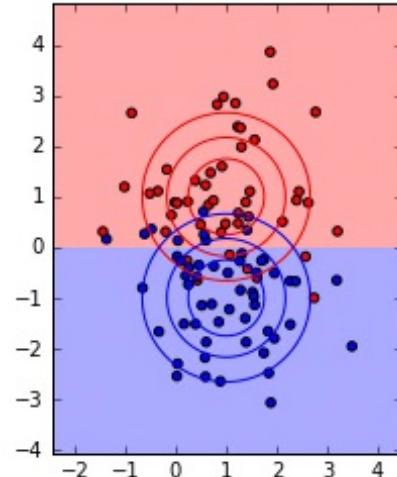
$$\hat{C} = \frac{\sum_{e=1}^M \sum_{i:y_i=e} (\underline{x}_i - \hat{\mu}_e)(\underline{x}_i - \hat{\mu}_e)^T}{(N-M) \leftarrow \text{or } N}$$

$$\text{Let } \delta_e(\underline{x}) = \ln \hat{\pi}_e + \hat{\mu}_e^T \hat{C}^{-1} \underline{x} - \frac{1}{2} \hat{\mu}_e^T \hat{C}^{-1} \hat{\mu}_e$$

$$\text{Then } \hat{y}_{\text{LDA}} = f_{\text{LDA}}(\underline{x}) = \arg \max_e \delta_e(\underline{x})$$



LDA



Bayes

Linear Classifiers

- Linear classifiers (e.g. LDA) are classifiers for which the objective $\delta_\ell(\underline{x})$ is linear

$$\delta_\ell(\underline{x}) = \underline{a}_\ell^\top (\underline{x}) + b_\ell, \quad \ell=1, 2, \dots, M$$

For LDA, $\underline{a}_\ell^\top = \hat{\underline{\mu}}_\ell^\top \hat{\underline{C}}^{-1}$, $b_\ell = \ln \hat{\pi}_\ell - \frac{1}{2} \hat{\underline{\mu}}_\ell^\top \hat{\underline{C}}^{-1} \hat{\underline{\mu}}_\ell$

$$\hat{y}_{\text{lin}} = f_{\text{lin}}(\underline{x}) = \arg \max_\ell (\underline{a}_\ell^\top (\underline{x}) + b_\ell)$$

- Decision between classes ℓ and m :

$$\underline{a}_\ell^\top \underline{x} + b_\ell \stackrel{\ell}{\geq_m} \underline{a}_m^\top \underline{x} + b_m$$

$$\equiv \underline{\omega}_{\ell m}^\top \underline{x} \stackrel{\ell}{\geq_m} \beta_{\ell m} \quad \begin{aligned} \underline{\omega}_{\ell m} &= \underline{a}_\ell - \underline{a}_m \\ \beta_{\ell m} &= b_m - b_\ell \end{aligned}$$

- Define $g_{\ell m}(\underline{x}) = \underline{\omega}_{\ell m}^\top \underline{x} - \beta_{\ell m} = -g_{m \ell}(\underline{x})$

\swarrow linear discriminant function

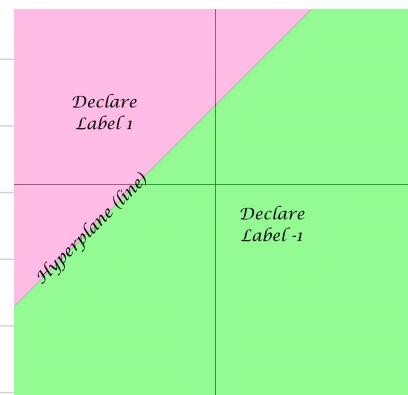
Then

$$g_{\ell m}(\underline{x}) \stackrel{\ell}{\geq_m} 0. \quad \text{Note: } g_{\ell m}(\underline{x}) = \delta_\ell(\underline{x}) - \delta_m(\underline{x})$$

- $g_{\ell m}(\underline{x}) = 0$ defines a hyperplane in \mathbb{R}^d .

Binary Linear Classifier

$$\hat{y}_{\text{lin}} = f_{\text{lin}}(\underline{x}) = \begin{cases} 1 & \text{if } \underline{\omega}^\top \underline{x} \geq \beta \\ -1 & \text{if } \underline{\omega}^\top \underline{x} < \beta \end{cases}$$



Example $M = 3, d = 2, \underline{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$$\delta_1(\underline{x}) = x_1 + 2x_2 + 1$$

$$\delta_2(\underline{x}) = 2x_1 - x_2 + 2$$

$$\delta_3(\underline{x}) = -x_1 + 2x_2 + 1$$

Linear discriminant functions :

$$g_{12}(\underline{x}) = -x_1 + 3x_2 - 1$$

$$g_{23}(\underline{x}) = 3x_1 - 3x_2 + 1$$

$$g_{31}(\underline{x}) = -2x_1$$

Suppose we want to classify $\underline{x} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$.

1. Using $\delta_e(\underline{x})$:

$$\delta_1(\underline{x}) = 4, \delta_2(\underline{x}) = 3, \delta_3(\underline{x}) = 0$$

$$\Rightarrow \hat{y} = \arg \max_l \delta_l(\underline{x}) = 1$$

2. Using $g_{\text{elm}}(\underline{x})$:

$$g_{12}(\underline{x}) = 2 > 0$$

$$g_{23}(\underline{x}) = 1 > 0$$

$$g_{31}(\underline{x}) = -2 < 0 \Rightarrow g_{13}(\underline{x}) > 0$$

$$\Rightarrow \hat{y} = 1$$