

# Linear Discriminant Analysis

Multi-class classification

$$\mathbf{x} = (x_1, \dots, x_k) \rightsquigarrow P(t|\mathbf{x}), \quad t = 1, 2, \dots, s$$

- We studied logistic regression.
- The **Linear Discriminant Analysis (LDA)** is based on Bayesian inference.

- $\pi_t$  prior probability that an observation belongs to class  $t$

$f_t(\mathbf{x}) := P(\mathbf{x}|t)$  likelihood

Bayes' Theorem:

$$P(t|\mathbf{x}) \propto \pi_t f_t(\mathbf{x})$$

- LDA assumes

(1)  $f_t(\mathbf{x})$  is normal, i.e.,  $f_t(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_t, \Sigma_t)$

(2) all the covariances are the same,

$$\text{i.e., } \Sigma := \Sigma_1 = \dots = \Sigma_S.$$

- Recall  $D$ -dimensional Gaussian distribution:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}$$

where  $\boldsymbol{\mu}$  is the mean and  $\Sigma$  is the covariance.

- We have

$$\begin{aligned} \ln P(t|\mathbf{x}) &= \ln \pi_t + \ln f_t(\mathbf{x}) + (\text{constant}) \\ &= \ln \pi_t - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_t)^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_t) + (\text{constant}) \\ &= \ln \pi_t - \frac{1}{2}\boldsymbol{\mu}_t^\top \Sigma^{-1} \boldsymbol{\mu}_t + \mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}_t + (\text{constant}). \end{aligned}$$

- Define the **discriminant function** by

$$\delta_t(\mathbf{x}) := \ln \pi_t - \frac{1}{2} \boldsymbol{\mu}_t^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_t + \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_t$$

for  $t = 1, \dots, s$ .

- Given  $\mathbf{x}$ , if  $\delta_{t_*}(\mathbf{x})$  is the largest, observation  $\mathbf{x}$  belongs to class  $t_*$  with largest probability.

~~~~~> Classification

- In LDA, we use the training data to approximate  $\delta_t(\mathbf{x})$ .

Given  $\mathcal{D} = \mathcal{D}_1 \sqcup \mathcal{D}_2 \sqcup \dots \sqcup \mathcal{D}_s$  (disjoint union),

set

$$N_t := \#(\mathcal{D}_t), \quad t = 1, 2, \dots, s, \quad N := N_1 + \dots + N_s.$$

We make the following estimates:

$$\hat{\pi}_t = N_t/N,$$

$$\hat{\boldsymbol{\mu}}_t = \frac{1}{N_t} \sum_{\mathbf{x} \in \mathcal{D}_t} \mathbf{x},$$

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{N-s} \sum_{t=1}^s \sum_{\mathbf{x} \in \mathcal{D}_t} (\mathbf{x} - \boldsymbol{\mu}_t)(\mathbf{x} - \boldsymbol{\mu}_t)^\top.$$

An approximation of the discriminant function is given by

$$\hat{\delta}_t(\mathbf{x}) := \ln \hat{\pi}_t - \frac{1}{2} \hat{\boldsymbol{\mu}}_t^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_t + \mathbf{x}^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_t.$$

- When  $\pi_1 = \pi_2 = \dots = \pi_s$ , the **decision boundaries** are given by

$$-\frac{1}{2}\hat{\boldsymbol{\mu}}_i^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_i + \mathbf{x}^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_i = -\frac{1}{2}\hat{\boldsymbol{\mu}}_j^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_j + \mathbf{x}^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}}_j$$

for  $i \neq j$ .