

## Exercise Sheet 3

### Exercise 1: Fisher Discriminant (10 + 10 + 10 P)

The objective function to find the Fisher Discriminant has the form

$$\max_{\mathbf{w}} \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}}$$

where  $\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^\top$  is the between-class scatter matrix and  $\mathbf{S}_W$  is within-class scatter matrix, assumed to be positive definite. Because there are infinitely many solutions (multiplying  $\mathbf{w}$  by a scalar doesn't change the objective), we can extend the objective with a constraint, e.g. that enforces  $\mathbf{w}^\top \mathbf{S}_W \mathbf{w} = 1$ .

- (a) *Reformulate* the problem above as an optimization problem with a quadratic objective and a quadratic constraint.
- (b) *Show* using the method of Lagrange multipliers that the solution of the reformulated problem is also a solution of the generalized eigenvalue problem:

$$\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w}$$

- (c) Show that the solution of this optimization problem is equivalent (up to a scaling factor) to

$$\mathbf{w}^* = \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$$

### Exercise 2: Bounding the Error (10 + 10 P)

The direction learned by the Fisher discriminant is equivalent to that of an optimal classifier when the class-conditioned data densities are Gaussian with same covariance. In this particular setting, we can derive a bound on the classification error which gives us insight into the effect of the mean and covariance parameters on the error.

Consider two data generating distributions  $P(\mathbf{x} \mid \omega_1) = \mathcal{N}(\boldsymbol{\mu}, \Sigma)$  and  $P(\mathbf{x} \mid \omega_2) = \mathcal{N}(-\boldsymbol{\mu}, \Sigma)$  with  $\mathbf{x} \in \mathbb{R}^d$ . Recall that the Bayes error rate is given by:

$$P(\text{error}) = \int_{\mathbf{x}} P(\text{error} \mid \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

- (a) Show that the conditional error can be upper-bounded as:

$$P(\text{error} \mid \mathbf{x}) \leq \sqrt{P(\omega_1 \mid \mathbf{x}) P(\omega_2 \mid \mathbf{x})}$$

- (b) Show that the Bayes error rate can then be upper-bounded by:

$$P(\text{error}) \leq \sqrt{P(\omega_1) P(\omega_2)} \cdot \exp\left(-\frac{1}{2} \boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{\mu}\right)$$

### Exercise 3: Fisher Discriminant (10 + 10 P)

Consider the case of two classes  $\omega_1$  and  $\omega_2$  with associated data generating probabilities

$$p(\mathbf{x} \mid \omega_1) = \mathcal{N}\left(\begin{bmatrix} -1 \\ -1 \end{bmatrix}, \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}\right) \quad \text{and} \quad p(\mathbf{x} \mid \omega_2) = \mathcal{N}\left(\begin{bmatrix} +1 \\ +1 \end{bmatrix}, \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

- (a) Find for this dataset the Fisher discriminant  $\mathbf{w}$  (i.e. the projection  $y = \mathbf{w}^\top \mathbf{x}$  under which the ratio between inter-class and intra-class variability is maximized).
- (b) Find a projection for which the ratio is minimized.

### Exercise 4: Programming (30 P)

Download the programming files on ISIS and follow the instructions.

**Exercise 1: Fisher Discriminant (10 + 10 + 10 P)**

The objective function to find the Fisher Discriminant has the form

$$\max_{\mathbf{w}} \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}}$$

where  $\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^\top$  is the between-class scatter matrix and  $\mathbf{S}_W$  is within-class scatter matrix, assumed to be positive definite. Because there are infinitely many solutions (multiplying  $\mathbf{w}$  by a scalar doesn't change the objective), we can extend the objective with a constraint, e.g. that enforces  $\mathbf{w}^\top \mathbf{S}_W \mathbf{w} = 1$ .

- (a) Reformulate the problem above as an optimization problem with a quadratic objective and a quadratic constraint.

**Solution:**

The problem can be formulated as

$$\begin{array}{ll} \max_{\mathbf{w}} & \mathbf{w}^\top \mathbf{S}_B \mathbf{w} \\ \text{s.t.} & \mathbf{w}^\top \mathbf{S}_W \mathbf{w} = 1 \end{array} \quad \Rightarrow \quad \begin{array}{ll} \max_{\mathbf{w}} & \mathbf{w}^\top \mathbf{S}_B \mathbf{w} \\ \text{s.t.} & 1 - \mathbf{w}^\top \mathbf{S}_W \mathbf{w} = 0 \end{array}$$

- (b) Show using the method of Lagrange multipliers that the solution of the reformulated problem is also a solution of the generalized eigenvalue problem:

$$\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w}$$

**Solution.**

We can formulate the Lagrange function:

$$\mathcal{L}(\lambda) = \mathbf{w}^\top \mathbf{S}_B \mathbf{w} + \lambda (1 - \mathbf{w}^\top \mathbf{S}_W \mathbf{w})$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{S}_B \mathbf{w} - \lambda \mathbf{S}_W \mathbf{w}$$

$$\Downarrow \text{ let } \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0$$

$$\therefore \mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w} \rightarrow \text{generalized Eigenvalue problem.}$$

- (c) Show that the solution of this optimization problem is equivalent (up to a scaling factor) to

$$\mathbf{w}^* = \mathbf{S}_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$$

**Solution:**

日期: /

From above we have:

$$S_B \omega = \lambda S_\omega \omega$$

$\Downarrow$  left multiplied by  $S_\omega^{-1}$

$$\therefore S_\omega^{-1} S_B \omega = \lambda S_\omega^{-1} S_\omega \omega$$

$\Downarrow$

$$S_\omega^{-1} S_B \omega = \lambda I \omega$$

$$\Downarrow S_B = (m_2 - m_1) (m_2 - m_1)^T$$

$$S_\omega^{-1} (m_2 - m_1) \underbrace{(m_2 - m_1)^T \omega}_{\text{scaler}} = \lambda \omega$$

$\Downarrow$   
 $\beta$

$\Downarrow$

$$S_\omega^{-1} (m_2 - m_1) \cdot \beta = \lambda \omega$$

$\therefore$  suppose the scaling factor is  $\lambda^* = \frac{\lambda}{\beta}$

$\Downarrow$

$$\omega^* = \frac{1}{\lambda^*} S_\omega^{-1} (m_2 - m_1)$$

$\Downarrow$

$\alpha$  (scaling factor)

$$\therefore \omega^* = \alpha S_\omega^{-1} (m_2 - m_1)$$

$\Downarrow$

$$\omega^* = S_\omega^{-1} (m_2 - m_1)$$

### Exercise 2: Bounding the Error (10 + 10 P)

The direction learned by the Fisher discriminant is equivalent to that of an optimal classifier when the class-conditioned data densities are Gaussian with same covariance. In this particular setting, we can derive a bound on the classification error which gives us insight into the effect of the mean and covariance parameters on the error.

Consider two data generating distributions  $P(\mathbf{x} | \omega_1) = \mathcal{N}(\boldsymbol{\mu}, \Sigma)$  and  $P(\mathbf{x} | \omega_2) = \mathcal{N}(-\boldsymbol{\mu}, \Sigma)$  with  $\mathbf{x} \in \mathbb{R}^d$ . Recall that the Bayes error rate is given by:

$$P(\text{error}) = \int_{\mathbf{x}} P(\text{error} | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

(a) Show that the conditional error can be upper-bounded as:

$$P(\text{error} | \mathbf{x}) \leq \sqrt{P(\omega_1 | \mathbf{x}) P(\omega_2 | \mathbf{x})}$$

Solution:

From sheet 1 we have the definition:

$$P(\text{error} | \mathbf{x}) = \min(P(\omega_1 | \mathbf{x}), P(\omega_2 | \mathbf{x}))$$

And we know that there is an inequality called generalized mean inequality:

$$\min(x_1, \dots, x_n) \leq M_p(x_1, \dots, x_n) \leq \max(x_1, \dots, x_n)$$

$$\text{where } M_p(x_1, \dots, x_n) = \left( \frac{1}{n} \sum_{i=1}^n x_i^p \right)^{\frac{1}{p}} \quad p > 0$$

$$M_0(x_1, \dots, x_n) = \left( \prod_{i=1}^n x_i \right)^{\frac{1}{n}} \Rightarrow \text{equal to geometric mean}$$

$\therefore$  we can derive:

$$\min(P(\omega_1 | \mathbf{x}), P(\omega_2 | \mathbf{x})) \leq M_0(P(\omega_1 | \mathbf{x}), P(\omega_2 | \mathbf{x}))$$

$\Downarrow$

$$\min(P(\omega_1 | \mathbf{x}), P(\omega_2 | \mathbf{x})) \leq (P(\omega_1 | \mathbf{x}) P(\omega_2 | \mathbf{x}))^{\frac{1}{2}}$$

$\Downarrow$

$$\min(P(\omega_1 | \mathbf{x}), P(\omega_2 | \mathbf{x})) \leq \sqrt{P(\omega_1 | \mathbf{x}) P(\omega_2 | \mathbf{x})}$$

$\Downarrow$

$$P(\text{error} | \mathbf{x}) \leq \sqrt{P(\omega_1 | \mathbf{x}) P(\omega_2 | \mathbf{x})}$$

(b) Show that the Bayes error rate can then be upper-bounded by:

$$P(\text{error}) \leq \sqrt{P(\omega_1)P(\omega_2)} \cdot \exp\left(-\frac{1}{2}\mu^T \Sigma^{-1} \mu\right)$$

Solution:

$$P(\text{error}) \leq \int_{\mathcal{X}} P(\text{error}|x) P(x) dx$$

$$P(\text{error}) \leq \int_{\mathcal{X}} \sqrt{P(\omega_1|x)P(\omega_2|x)} P(x) dx$$

$$= \int_{\mathcal{X}} \sqrt{\frac{P(x|\omega_1)P(\omega_1)}{P(x)} \cdot \frac{P(x|\omega_2)P(\omega_2)}{P(x)}} P(x) dx$$

$$= \int_{\mathcal{X}} \sqrt{P(x|\omega_1)P(\omega_1)P(x|\omega_2)P(\omega_2)} dx$$

$$= \sqrt{P(\omega_1)P(\omega_2)} \int_{\mathcal{X}} \sqrt{P(x|\omega_1)P(x|\omega_2)} dx$$

Since  $P(x|\omega_1)$  and  $P(x|\omega_2)$  have opposite mean and same variance Gaussian.

$$\therefore \int_{\mathcal{X}} \sqrt{P(x|\omega_1)P(x|\omega_2)} dx$$

$$= \int_{\mathcal{X}} \sqrt{\frac{1}{2^d \sqrt{2\pi} (\det(\Sigma))}} \exp\left(-\frac{1}{2}((x-\mu)^T \Sigma^{-1}(x-\mu) + (x+\mu)^T \Sigma^{-1}(x+\mu))\right) dx$$

$$= \int_{\mathcal{X}} \frac{1}{\sqrt{2\pi} \det(\Sigma)} \sqrt{\exp\left(-\frac{1}{2}(2x^T \Sigma^{-1}x + 2\mu^T \Sigma^{-1}\mu)\right)} dx$$

$$= \int_{\mathcal{X}} \frac{1}{\sqrt{2\pi} \det(\Sigma)} \sqrt{\exp(-x^T \Sigma^{-1}x) \cdot \exp(-\mu^T \Sigma^{-1}\mu)} dx$$

$$= \exp\left(-\frac{1}{2}\mu^T \Sigma^{-1}\mu\right) \cdot \underbrace{\int_{\mathcal{X}} \frac{1}{\sqrt{2\pi} \det(\Sigma)} \exp\left(-\frac{1}{2}x^T \Sigma^{-1}x\right) dx}_{N(0, \Sigma)}$$

日期:

$$\dots = \exp(-\frac{1}{2}\mu^T \Sigma^{-1} \mu) \cdot 1 \quad (\text{integrate the pdf on the whole interval})$$

$$\therefore p(\text{error}) \leq \sqrt{P(w_1)P(w_2)} \exp(-\frac{1}{2}\mu^T \Sigma^{-1} \mu)$$

$\therefore$  proved!

**Exercise 3: Fisher Discriminant (10 + 10 P)**

Consider the case of two classes  $\omega_1$  and  $\omega_2$  with associated data generating probabilities

$$p(\mathbf{x} | \omega_1) = \mathcal{N}\left(\begin{bmatrix} -1 \\ -1 \end{bmatrix}, \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}\right) \quad \text{and} \quad p(\mathbf{x} | \omega_2) = \mathcal{N}\left(\begin{bmatrix} +1 \\ +1 \end{bmatrix}, \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

- (a) Find for this dataset the Fisher discriminant  $\mathbf{w}$  (i.e. the projection  $y = \mathbf{w}^\top \mathbf{x}$  under which the ratio between inter-class and intra-class variability is maximized).

Solution:

$$\begin{aligned} \mathbf{w}^* &= S_w^{-1} (\mathbf{m}_2 - \mathbf{m}_1) \\ &= (\Sigma_1 + \Sigma_2)^{-1} (\mathbf{m}_2 - \mathbf{m}_1) \\ &= \left( \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \right)^{-1} \cdot \begin{bmatrix} 2 \\ 2 \end{bmatrix} \\ &= \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix}^{-1} \cdot \begin{bmatrix} 2 \\ 2 \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix} \end{aligned}$$

$$\therefore \mathbf{w}^* \text{ could be } \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix}$$

- (b) Find a projection for which the ratio is minimized.

Solution:

$$S_w^{-1} S_B \mathbf{w} = \lambda \mathbf{w}$$

$$\therefore S_w^{-1} = \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & \frac{1}{2} \end{bmatrix}, \quad S_B = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^\top = \begin{bmatrix} 4 & 4 \\ 4 & 4 \end{bmatrix}$$

$$\therefore S_w^{-1} S_B = \begin{bmatrix} 1 & 1 \\ 2 & 2 \end{bmatrix} = A$$

$$\therefore A \mathbf{w} = \lambda \mathbf{w}$$

$$\therefore |\lambda I - A| = (\lambda - 1)(\lambda - 2) - 2 = \lambda^2 - 3\lambda = 0$$

$$\therefore \textcircled{1} \lambda_1 = 0 \rightarrow \mathbf{w} \text{ could be } \begin{bmatrix} 2 \\ -2 \end{bmatrix}$$

$$\textcircled{2} \lambda_2 = 3 \rightarrow \mathbf{w} \text{ could be } \begin{bmatrix} 1 \\ 2 \end{bmatrix} \text{ or } \begin{bmatrix} 2 \\ 4 \end{bmatrix} \dots$$