

## 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

$$\max_w \frac{w^T S_B w}{w^T S_W w}$$

where  $S_B = (m_2 - m_1)(m_2 - m_1)^T$  is the between-class scatter matrix and  $S_W$  is within-class scatter matrix, assumed to be positive definite. Because there are infinitely many solutions (multiplying  $w$  by a scalar doesn't change the objective), we can extend the objective with a constraint, e.g. that enforces  $w^T S_W w = 1$ .

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

$$\operatorname{argmax}_w w^T S_B w \quad \text{with} \quad w^T S_W w = 1$$

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

$$S_B w = \lambda S_W w$$

$$\mathcal{L}(w, \lambda) = w^T S_B w + \lambda (1 - w^T S_W w)$$

$$\frac{\partial}{\partial \lambda} \mathcal{L}(w, \lambda) = 1 - w^T S_W w = 0$$

$$\nabla_w \mathcal{L}(w, \lambda) = 2S_B w - 2\lambda S_W w = 0 \Rightarrow S_B w = \lambda S_W w$$

Why  $\nabla_w (w^T S_B w) = 2S_B w$  and  $\nabla_w (w^T S_W w) = 2S_W w$ :

$$\begin{aligned} w &= \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \quad S_B = \begin{bmatrix} S_{11} & \dots & S_{1d} \\ \vdots & \ddots & \vdots \\ S_{d1} & \dots & S_{dd} \end{bmatrix} \\ w^T S_B w &= [w_1 \ w_2 \ \dots \ w_d] \begin{bmatrix} S_{11} & S_{12} & \dots & S_{1d} \\ S_{21} & \ddots & & \\ \vdots & & \ddots & \\ S_{d1} & & & S_{dd} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} \\ &= \underbrace{[w_1 S_{11} + w_2 S_{12} + \dots + w_d S_{1d} \quad \dots \quad w_1 S_{d1} + w_2 S_{d2} + \dots + w_d S_{dd}]}_d \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} \\ &= [w_1 \cdot (w_1 S_{11} + w_2 S_{12} + \dots + w_d S_{1d}) + w_2 (\dots) + \dots + w_d (w_1 S_{d1} + w_2 S_{d2} + \dots + w_d S_{dd})] \\ \nabla \frac{f}{f} (w^T S_B w) &= \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix} = \begin{bmatrix} 2S_{11} \cdot w_1 + 2w_2 S_{12} + \dots + 2w_d S_{1d} \\ \vdots \\ 2w_1 S_{d1} + \dots + 2w_d S_{dd} \end{bmatrix} \\ &= 2S_B w \end{aligned}$$

A more formal proof is provided below:

Why  $\nabla_W(W^T S_B W) = 2 S_B W$  and  $\nabla_W(W^T S_W W) = 2 S_W W$ :

formal proof:

We show that  $\nabla_W(W^T S W) = 2 S W$  if  $S = S^T$ .  $W^T S W$  can be rewritten as

$$W^T S W = \sum_i w_i \sum_j S_{ij} w_j = \sum_i \sum_j w_i S_{ij} w_j$$

This is derived with respect to some  $w_k$ . Each  $w_i S_{ij} w_j$  is looked at separately.

If  $i=j$ , then  $w_i S_{ij} w_j = w_i \cdot S_{ii} w_i = S_{ii} w_i^2$ . This square is only non-zero when derived,

if  $i=j=k$ ; in this case we get  $\frac{\partial}{\partial w_k} (S_{kk} w_k^2) = 2 \cdot S_{kk} w_k$ . If  $i \neq j$  all terms where neither  $i$  nor  $j$

are  $k$  will be zero if derived with respect to  $w_k$ . Only pairs  $\frac{\partial}{\partial w_k} (w_i \cdot S_{ik} \cdot w_k) = S_{ik} \cdot w_i$  and

$\frac{\partial}{\partial w_k} (w_k \cdot S_{kj} \cdot w_j) = S_{kj} w_j$  are non-zero, because the others are constants. Because  $S$  is

symmetric, we get  $S_{ik} \cdot w_i = S_{ki} \cdot w_i$ . Therefore the following holds

$$\frac{\partial}{\partial w_k} \left( \sum_i \sum_j w_i S_{ij} w_j \right) = \sum_i 2 \cdot S_{ki} w_i = 2 \cdot \sum_i S_{ki} w_i = 2 \cdot S W$$

Therefore  $\nabla_W(W^T S W) = 2 S W$  if  $S = S^T$ .

Both  $S_B = (m_2 - m_1)(m_2 - m_1)^T$  and  $S_W = S_1 + S_2$ .  $S_W$  are symmetric, because the scatter matrix is symmetric:

$$\begin{aligned} S_k &= \sum_{i \in \mathcal{C}_k} (x_i - \mu_k)(x_i - \mu_k)^T \\ &= \sum_{i \in \mathcal{C}_k} ((x_i - \mu_k)^T (x_i - \mu_k))^T \\ &= \left( \sum_{i \in \mathcal{C}_k} (x_i - \mu_k)^T (x_i - \mu_k) \right)^T \\ &= \left( \sum_{i \in \mathcal{C}_k} (x_i - \mu_k)(x_i - \mu_k)^T \right)^T = S_k^T \end{aligned}$$

Using the lemma that was proven above results in  $\nabla_W(W^T S_B W) = 2 S_B W$  and

$$\nabla_W(W^T S_W W) = 2 S_W W.$$

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

$$w^* = S_W^{-1}(m_2 - m_1)$$

$$S_B \cdot w = \lambda S_W w$$

$$(m_2 - m_1)(m_2 - m_1)^T \cdot w = \lambda S_W \cdot w$$

$$\Leftrightarrow S_W^{-1} (m_2 - m_1)(m_2 - m_1)^T \cdot w = \lambda w$$

( $S_W^{-1}$  exists because  $S_W$  is assumed to be positive definite)

$$\Leftrightarrow S_W^{-1} (m_2 - m_1) \cdot \underbrace{\lambda (m_2 - m_1)^T \cdot w}_{\text{scalar}} = w$$

$$\Leftrightarrow w^* = S_W^{-1} (m_2 - m_1)$$

## 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:

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

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

Case 1. if  $P(\omega_1 | \mathbf{x}) \leq P(\omega_2 | \mathbf{x})$

$$\begin{aligned} P(\text{error} | \mathbf{x}) &= P(\omega_1 | \mathbf{x}) = \sqrt{P(\omega_1 | \mathbf{x}) P(\omega_1 | \mathbf{x})} \\ &\leq \sqrt{P(\omega_1 | \mathbf{x}) P(\omega_2 | \mathbf{x})} \end{aligned}$$

Case 2. if  $P(\omega_1 | \mathbf{x}) \geq P(\omega_2 | \mathbf{x})$

$$\begin{aligned} P(\text{error} | \mathbf{x}) &= P(\omega_2 | \mathbf{x}) = \sqrt{P(\omega_2 | \mathbf{x}) P(\omega_2 | \mathbf{x})} \\ &\leq \sqrt{P(\omega_1 | \mathbf{x}) P(\omega_2 | \mathbf{x})} \end{aligned}$$

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

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

$$\begin{aligned} P(\text{error}) &= \int_{\mathbf{x}} P(\text{error} | \mathbf{x}) P(\mathbf{x}) d\mathbf{x} \\ &\leq \int_{\mathbf{x}} \sqrt{P(\omega_1 | \mathbf{x}) P(\omega_2 | \mathbf{x})} P(\mathbf{x}) d\mathbf{x} \\ &= \int_{\mathbf{x}} \sqrt{\frac{P(\mathbf{x} | \omega_1) P(\omega_1)}{P(\mathbf{x})} \frac{P(\mathbf{x} | \omega_2) P(\omega_2)}{P(\mathbf{x})}} P(\mathbf{x}) d\mathbf{x} \\ &= \sqrt{P(\omega_1) P(\omega_2)} \cdot \int_{\mathbf{x}} \sqrt{P(\mathbf{x} | \omega_1) P(\mathbf{x} | \omega_2)} d\mathbf{x} \\ &= \sqrt{P(\omega_1) P(\omega_2)} \int_{\mathbf{x}} \sqrt{\left(\frac{1}{(2\pi)^d \det(\Sigma)}\right)^2 e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})} \cdot e^{-\frac{1}{2}(\mathbf{x}+\boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}+\boldsymbol{\mu})}} d\mathbf{x} \\ &= \sqrt{P(\omega_1) P(\omega_2)} \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \int_{\mathbf{x}} \sqrt{e^{-\frac{1}{2}[(\mathbf{x}-\boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}) + (\mathbf{x}+\boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}+\boldsymbol{\mu})]}} d\mathbf{x} \end{aligned}$$

specifically we have:  $(\mathbf{x}-\boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}) + (\mathbf{x}+\boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}+\boldsymbol{\mu})$

$$\begin{aligned} &= \mathbf{x}^T \Sigma^{-1} \mathbf{x} + \boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu} - \cancel{\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}} - \cancel{\boldsymbol{\mu}^T \Sigma^{-1} \mathbf{x}} + \mathbf{x}^T \Sigma^{-1} \mathbf{x} + \boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu} + \cancel{\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}} + \cancel{\boldsymbol{\mu}^T \Sigma^{-1} \mathbf{x}} \\ &= 2(\mathbf{x}^T \Sigma^{-1} \mathbf{x} + \boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu}) \end{aligned}$$

$$\begin{aligned}
P(\text{error}) &= \sqrt{P(W_1)P(W_2)} \cdot \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \int_{\mathbf{x}} \sqrt{e^{-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x} + \mu^T \Sigma^{-1} \mu}} d\mathbf{x} \\
&= \sqrt{P(W_1)P(W_2)} \cdot e^{-\frac{1}{2} (\mu^T \Sigma^{-1} \mu)} \cdot \underbrace{\int_{\mathbf{x}} \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} e^{-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}} d\mathbf{x}}_{=1 \quad \text{the density of normal distribution } \mathcal{N}(0, \Sigma) \text{ integrate over the entire space.}} \\
&= \sqrt{P(W_1)P(W_2)} \cdot e^{-\frac{1}{2} (\mu^T \Sigma^{-1} \mu)}
\end{aligned}$$

### 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  $y = \mathbf{w}^T \mathbf{x}$  under which the ratio between inter-class and intra-class variability is maximized).

$$\begin{aligned}
S_W &= S_1 + S_2 = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix} & A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \\
S_W^{-1} &= \frac{1}{8} \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix} = \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} & A^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} \\
\mathbf{w} &= S_W^{-1} (\mu_2 - \mu_1) = \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 1 - (-1) \\ 1 - (-1) \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}
\end{aligned}$$

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

$$S_B = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \begin{bmatrix} 2 & 2 \end{bmatrix} = \begin{bmatrix} 4 & 4 \\ 4 & 4 \end{bmatrix}$$

$$\begin{aligned}
S_B \mathbf{w} &= \lambda S_W \mathbf{w} \\
S_W^{-1} S_B \mathbf{w} &= \lambda \mathbf{w} \\
\begin{bmatrix} \frac{1}{4} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 4 & 4 \\ 4 & 4 \end{bmatrix} \mathbf{w} &= \lambda \mathbf{w}
\end{aligned}$$

$$\begin{aligned}
\begin{bmatrix} 1 & 1 \\ 2 & 2 \end{bmatrix} \mathbf{w} &= \lambda \mathbf{w} \\
\det \begin{bmatrix} 1-\lambda & 1 \\ 2 & 2-\lambda \end{bmatrix} &= (1-\lambda)(2-\lambda) - 2 = \lambda^2 - 3\lambda + 2 - 2 = 0 \\
\lambda_1 &= 0 \quad \lambda_2 = 3
\end{aligned}$$

① eigenvector for  $\lambda_1 = 0$ :

$$\begin{aligned}
\begin{bmatrix} 1-0 & 1 \\ 2 & 2-0 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\
v_1 + v_2 &= 0 \\
2v_1 + 2v_2 &= 0 \Rightarrow \vec{v}_1 = \begin{pmatrix} v_1 \\ -v_1 \end{pmatrix} \xrightarrow{v_1=1} \begin{pmatrix} 1 \\ -1 \end{pmatrix}
\end{aligned}$$

② eigenvector for  $\lambda_2 = 3$ :

$$\begin{aligned}
\begin{bmatrix} 1-3 & 1 \\ 2 & 2-3 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\
\begin{bmatrix} -2 & 1 \\ 2 & -1 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow \vec{v}_2 = \begin{pmatrix} v_1 \\ 2v_1 \end{pmatrix} \xrightarrow{v_1=\frac{1}{2}} \begin{pmatrix} \frac{1}{2} \\ 1 \end{pmatrix}
\end{aligned}$$

$$\textcircled{1} w = \vec{v}_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

$$J_1 = \frac{w^T S_B w}{w^T S_W w} = w^T S_B w = (1 \ -1) \begin{bmatrix} 8 & 8 \\ 8 & 8 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} \\ = [0 \ 0] \begin{bmatrix} 1 \\ -1 \end{bmatrix} = 0$$

$$\textcircled{2} w = \vec{v}_2 = \begin{pmatrix} \frac{1}{2} \\ 1 \end{pmatrix}$$

$$J_2 = \begin{bmatrix} \frac{1}{2} & 1 \end{bmatrix} \begin{bmatrix} 8 & 8 \\ 8 & 8 \end{bmatrix} \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix} \\ = [12 \ 12] \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix} = 18$$

$$J_1 < J_2$$

→ Now we choose the corresponding eigenvector  $\vec{v}_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$  as  $w$ ,

the ratio is minimized