

Fisher Linear Discriminant

In this exercise, we apply Fisher Linear Discriminant as described in Chapter 3.8.2 of Duda et al. on the UCI Abalone dataset. A description of the dataset is given at the page <https://archive.ics.uci.edu/ml/datasets/Abalone>. The following two methods are provided for your convenience:

- `utils.Abalone.__init__(self)` reads the Abalone data and instantiates two data matrices corresponding to: *infant (I)*, *non-infant (N)*.
- `utils.Abalone.plot(self,w)` produces a histogram of the data when projected onto a vector `w`, and where each class is shown in a different color.

Sample code that makes use of these two methods is given below. It loads the data, looks at the shape of instantiated matrices, and plots the projection on the first dimension of the data representing the length of the abalone.

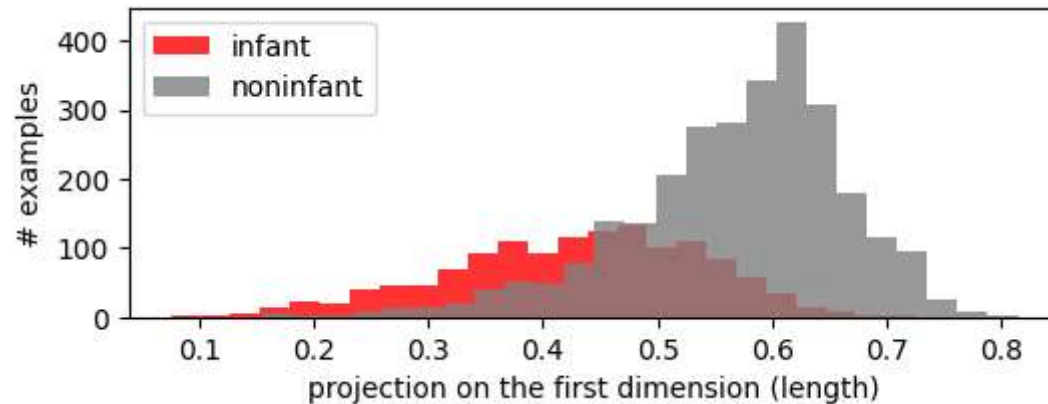
```
In [1]: %matplotlib inline
import utils, numpy
from matplotlib import pyplot as plt

# Load the data
abalone = utils.Abalone()

# Print dataset size for each class
print(abalone.I.shape, abalone.N.shape)

# Project data on the first dimension
w1 = numpy.array([1,0,0,0,0,0,0])
abalone.plot(w1, 'projection on the first dimension (length)')

(1342, 7) (2835, 7)
```



Implementation (10 + 5 + 5 = 20 P)

- Create a function `w = fisher(X1,X2)` that takes as input the data for two classes and returns the Fisher linear discriminant.
- Create a function `objective(X1,X2,w)` that evaluates the objective defined in Equation 96 of Duda et al. for an arbitrary projection vector `w`.
- Create a function `z = phi(X)` that returns a quadratic expansion for each data point `x` in the dataset. Such expansion consists of the vector `x` itself, to which we concatenate the vector of all pairwise products between elements of `x`. In other words, letting $x = (x_1, \dots, x_d)$ denote the d -dimensional data point, the quadratic expansion for this data point is a $d \cdot (d + 3)/2$ dimensional vector given by $\phi(x) = (x_i)_{1 \leq i \leq d} \cup (x_i x_j)_{1 \leq i \leq j \leq d}$. For example, the quadratic expansion for $d = 2$ is $(x_1, x_2, x_1^2, x_2^2, x_1 x_2)$.

```
In [4]: def fisher(X1,X2):
        ##### Replace by your code
        mean1 = numpy.mean(X1, axis=0)
        mean2 = numpy.mean(X2, axis=0)
        mean_diff = mean1 - mean2

        # Here I get the term " * (X1.shape[0] - 1)" from Nico, but I don't get that.
        # Since in the slides I can find any clue about this scaling vector.
        # But after adding this term the results seem to be identical to the given ones.
        cov1 = numpy.cov(X1, rowvar=False, bias = True) * (X1.shape[0] - 1)
        cov2 = numpy.cov(X2, rowvar=False, bias = True) * (X2.shape[0] - 1)
        S_w = cov1 + cov2
```

```

S_w_inv = numpy.linalg.inv(S_w)

return S_w_inv@mean_diff
#####

def objective(X1,X2,w):
    ##### Replace by your code
    mean1 = numpy.mean(X1, axis=0)
    mean2 = numpy.mean(X2, axis=0)
    mean_diff = mean1 - mean2

    S_B = numpy.outer(mean_diff, mean_diff)

    cov1 = numpy.cov(X1, rowvar=False, bias = True) * (X1.shape[0] - 1)
    cov2 = numpy.cov(X2, rowvar=False, bias = True) * (X2.shape[0] - 1)
    S_w = cov1 + cov2

    max_mean_class_diff = w.T @ S_B @ w
    min_inclass_var = w.T @ S_w @ w

    return (max_mean_class_diff / min_inclass_var)
#####

def expand(X):
    # Get the number of samples (n) and dimensionality (d) of the input data
    n, d = X.shape

    # Calculate the number of quadratic terms: (d*(d+3))/2
    expanded_dim = (d * (d + 3)) // 2

    # Initialize an empty array to hold the expanded dataset
    Z = numpy.zeros((n, expanded_dim))

    # Iterate through each data point
    for i, x in enumerate(X):
        # Copy the original vector
        z = list(x)

        # Append all pairwise products (including squares) to the expansion
        for j in range(d):
            for k in range(j, d):
                z.append(x[j] * x[k])

    # Assign the expanded vector to the result array

```

```

        Z[i] = z

    return Z

def my_plot(Z1, Z2, w, name):
    plt.figure(figsize=(6,2))
    plt.xlabel(name)
    plt.ylabel('# examples')
    plt.hist(numpy.dot(Z1,w), bins=25, alpha=0.8, color='red', label='infant')
    plt.hist(numpy.dot(Z2,w), bins=25, alpha=0.8, color='gray', label='noninfant')
    plt.legend()
    plt.show()

```

Analysis (5 + 5 = 10 P)

- **Print value of the objective function and the histogram for several values of w :**

- w is a canonical coordinate vector for the first feature (length).
- w is the difference between the mean vectors of the two classes.
- w is the Fisher linear discriminant.
- w is the Fisher linear discriminant (after quadratic expansion of the data).

```

In [6]: ##### REPLACE BY YOUR CODE
%matplotlib inline

w1 = numpy.array([1,0,0,0,0,0,0])
print('First dimension:')
print(objective(abalone.I, abalone.N, w1))
abalone.plot(w1, 'projection on the first dimension (length)')

mean1 = numpy.mean(abalone.I, axis=0)
mean2 = numpy.mean(abalone.N, axis=0)
mean_diff = mean1 - mean2
w2 = mean_diff
print('Difference:')
print(objective(abalone.I, abalone.N, w2))
abalone.plot(w2, 'projection on the difference between the mean vectors')

w3 = fisher(abalone.I, abalone.N)

```

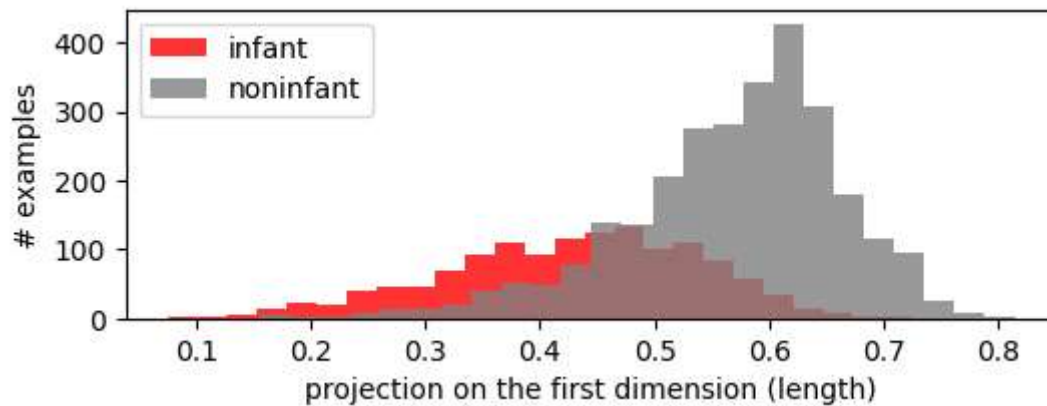
```

print('Fisher linear discriminant:')
print(objective(abalone.I, abalone.N, w3))
abalone.plot(w3, 'projection on the Fisher linear discriminant')

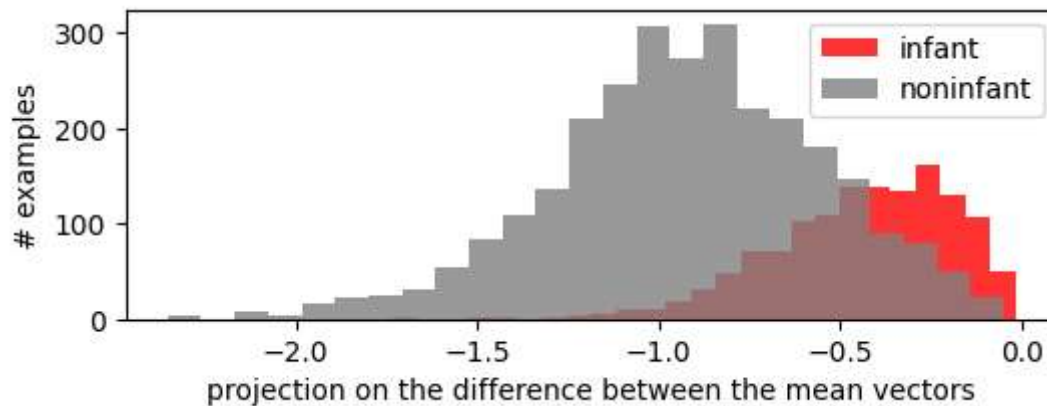
Z1 = expand(abalone.I)
Z2 = expand(abalone.N)
w4 = fisher(Z1, Z2)
print('Fisher linear discriminant(after expand):')
print(objective(Z1, Z2, w4))
my_plot(Z1, Z2, w4, 'projection on the first dimension (length)')
#####

```

First dimension:
0.00048003649875237174

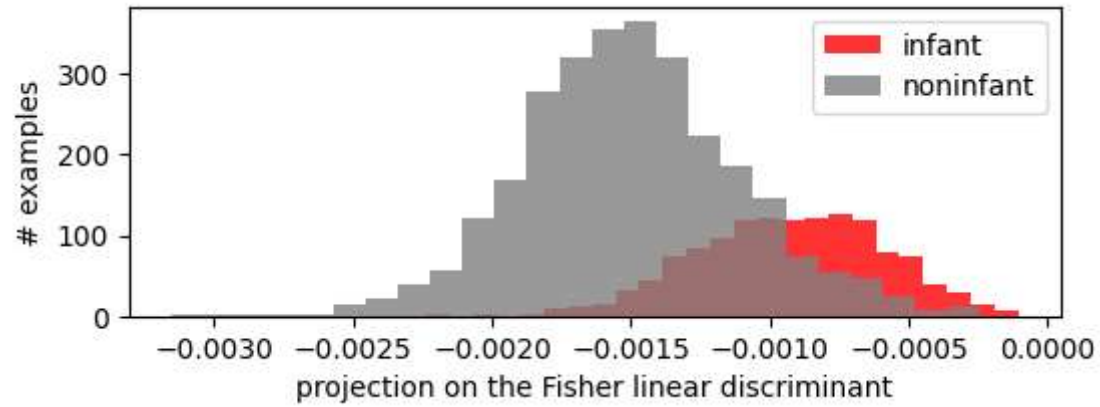


Difference:
0.0004997853617622609



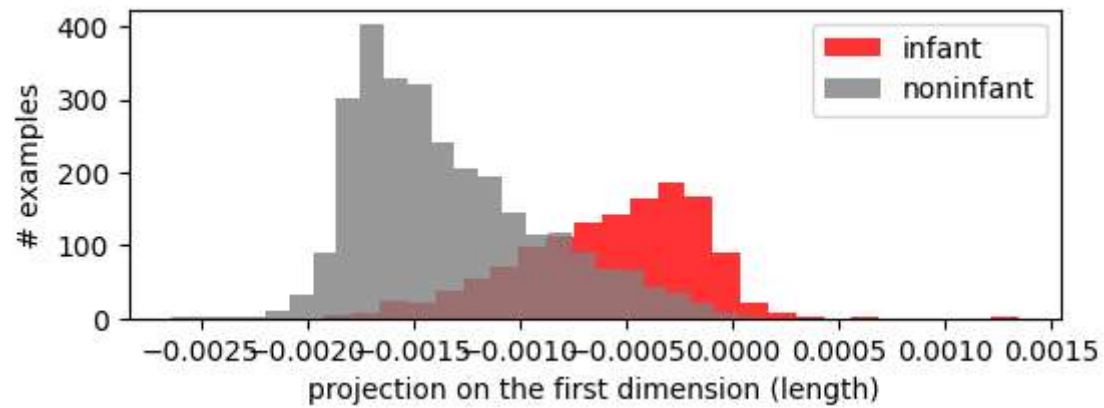
Fisher linear discriminant:

0.0005694412070089378



Fisher linear discriminant(after expand):

0.0007663654723881523



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 ω_1 and ω_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:

爽筒の 大

From above we have:

$$S_B \omega = \lambda S_w \omega$$

\Downarrow left multiplied by S_w^{-1}

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

\Downarrow

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

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

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

\Downarrow
 β

\Downarrow

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

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

\Downarrow

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

\Downarrow

α (scaling factor)

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

\Downarrow

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

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:

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

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

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

(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:

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

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

$$= \int_{\mathbf{x}} \sqrt{\frac{P(\mathbf{x}|\omega_1)P(\omega_1)}{P(\mathbf{x})} \cdot \frac{P(\mathbf{x}|\omega_2)P(\omega_2)}{P(\mathbf{x})}} P(\mathbf{x}) d\mathbf{x}$$

$$= \int_{\mathbf{x}} \sqrt{P(\mathbf{x}|\omega_1)P(\omega_1)P(\mathbf{x}|\omega_2)P(\omega_2)} d\mathbf{x}$$

$$= \sqrt{P(\omega_1)P(\omega_2)} \int_{\mathbf{x}} \sqrt{P(\mathbf{x}|\omega_1)P(\mathbf{x}|\omega_2)} d\mathbf{x}$$

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

$$\therefore \int_{\mathbf{x}} \sqrt{P(\mathbf{x}|\omega_1)P(\mathbf{x}|\omega_2)} d\mathbf{x}$$

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

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

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

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

爽簡明 天 $\frac{1}{\| \cdot \|}$
 $\therefore = \exp(-\frac{1}{2}\mu^T \Sigma^{-1} \mu) \cdot 1$ (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 ω_1 and ω_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:

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

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

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

Solution:

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

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

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

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

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

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

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