

Exercise Sheet 4

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

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

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:

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

- (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

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

- (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.

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> (<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]:

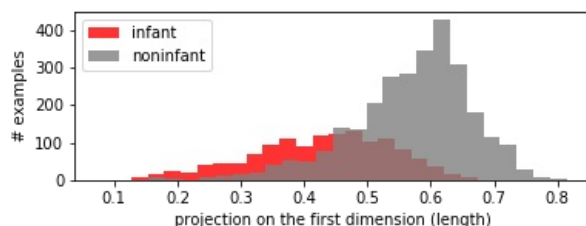
```
%matplotlib inline
import utils,numpy

# 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 [2]:

```
def fisher(X1,X2):
    ##### Replace by your code
    import solutions
    return solutions.fisher(X1,X2)
    #####

def objective(X1,X2,w):
    ##### Replace by your code
    import solutions
    return solutions.objective(X1,X2,w)
    #####

def expand(X):
    ##### Replace by your code
    import solutions
    return solutions.expand(X)
    #####
```

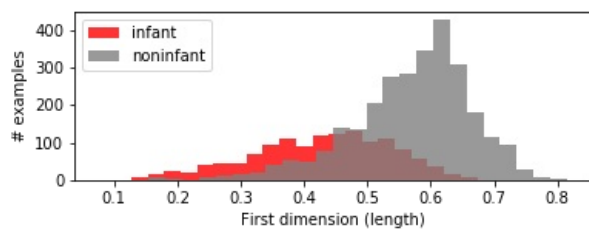
Analysis (5 + 5 = 10 P)

- **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).

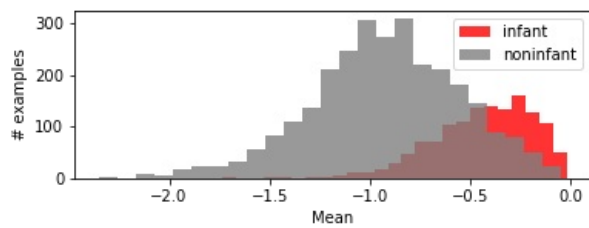
In [3]:

```
##### REPLACE BY YOUR CODE
%matplotlib inline
import solutions
solutions.analysis()
#####
```

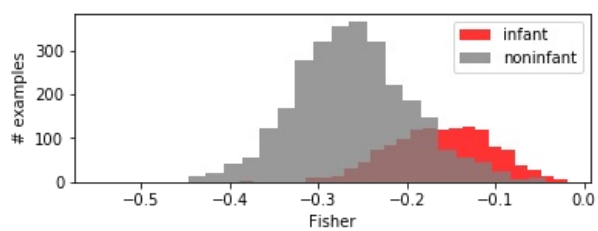
First dimension (length): 0.00048



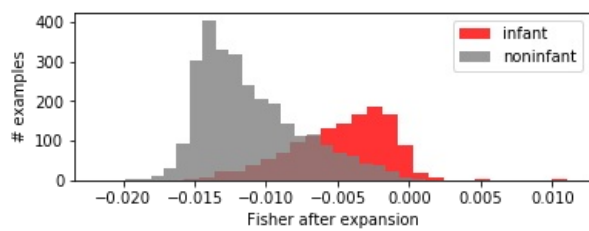
Means Linear: 0.00050



Fisher: 0.00057



Fisher after expansion: 0.00077



Sheet 4

\bar{x} : vector

X : matrice

1

a) $f = \max_w \frac{\bar{w}^T \underline{S}_B \bar{w}}{\bar{w}^T \underline{S}_W \bar{w}}$ with the constraint 1) $\bar{w}^T \underline{S}_W \bar{w} = 1$

1) $= \max_w \frac{\bar{w}^T \underline{S}_B \bar{w}}{1} = \max_w \bar{w}^T \underline{S}_B \bar{w}$

b) $f(\bar{w}, \lambda) = \bar{w}^T \underline{S}_B \bar{w} - \lambda (\bar{w}^T \underline{S}_W \bar{w} - 1)$

$$\frac{df}{d\bar{w}} = 2 \cdot \underline{S}_B \cdot \bar{w} - 2\lambda \underline{S}_W \cdot \bar{w} \stackrel{!}{=} 0$$

$$\Leftrightarrow \underline{S}_B \cdot \bar{w} = \lambda \underline{S}_W \bar{w}$$

c) $\Leftrightarrow \frac{1}{\lambda} \underline{S}_B \bar{w} = \underline{S}_W \bar{w}$

$$\Leftrightarrow \frac{1}{\lambda} \underline{S}_W^{-1} \underline{S}_B \bar{w} = \bar{w}$$

$$\Leftrightarrow \frac{1}{\lambda} \underline{S}_W^{-1} (\bar{m}_1 - \bar{m}_2) (\bar{m}_1 - \bar{m}_2)^T \bar{w} = \bar{w}$$

$$\Leftrightarrow \bar{w} = \text{const.} \cdot \underline{S}_W^{-1} (\bar{m}_1 - \bar{m}_2)$$

$$\Leftrightarrow \bar{w} \propto \underline{S}_W^{-1} (\bar{m}_1 - \bar{m}_2)$$

2

a) According to Sheet 1 exercise 1, the Bayes error $P(\text{error}|\bar{x})$ is $\min[P(w_1|\bar{x}), P(w_2|\bar{x})]$ for two classes w_1 and w_2 :

$$P(\text{error}|\bar{x}) = \min[P(w_1|\bar{x}), P(w_2|\bar{x})]$$

$$= M_{-\infty}[P(w_1|\bar{x}), P(w_2|\bar{x})]$$

with $M_{-\infty}$ being the generalized mean (see relation 1a) of sheet 1)

$$\leq M_0[P(w_1|\bar{x}), P(w_2|\bar{x})]$$

because $M_p \leq M_q$ if $p < q$

$$= \sqrt{P(w_1|\bar{x}) P(w_2|\bar{x})}$$

b) According to Eq. 1 in sheet 1:

$$P(\text{error}) = \int P(\text{error}|\bar{x}) p(\bar{x}) d\bar{x}$$

according to 2a):

$$\leq \int \sqrt{P(w_1|\bar{x}) P(w_2|\bar{x})} p(\bar{x}) d\bar{x}$$

according to Bayes' rule:

$$= \int \frac{P(\bar{x}|w_1) P(w_1)}{P(\bar{x})} \frac{P(\bar{x}|w_2) P(w_2)}{P(\bar{x})} p(\bar{x}) d\bar{x}$$

$$= \sqrt{P(w_1) P(w_2)} \int \sqrt{P(\bar{x}|w_1) P(\bar{x}|w_2)} d\bar{x}$$

$$= \sqrt{P(w_1) P(w_2)} \int \sqrt{\mathcal{N}(\bar{x}|\bar{\mu}_1, \Sigma_1) \cdot \mathcal{N}(\bar{x}|\bar{\mu}_2, \Sigma_2)} d\bar{x}$$

$$= \sqrt{P(w_1) P(w_2)} \int \frac{e^{-\frac{1}{2}(\bar{x}-\bar{\mu}_1)^T \Sigma_1^{-1} (\bar{x}-\bar{\mu}_1)}}{\sqrt{(2\pi)^d |\Sigma_1|}} \cdot \frac{e^{-\frac{1}{2}(\bar{x}-\bar{\mu}_2)^T \Sigma_2^{-1} (\bar{x}-\bar{\mu}_2)}}{\sqrt{(2\pi)^d |\Sigma_2|}} d\bar{x}$$

$$= \sqrt{P(w_1) P(w_2)} \int \frac{e^{-\frac{1}{2}\bar{x}^T \Sigma^{-1} \bar{x}} e^{-\frac{1}{2}\bar{\mu}_1^T \Sigma^{-1} \bar{\mu}_2}}{\sqrt{(2\pi)^d |\Sigma|}} d\bar{x} = \sqrt{P(w_1) P(w_2)} e^{-\frac{1}{2}\bar{\mu}_1^T \Sigma^{-1} \bar{\mu}_2} \int \frac{e^{-\frac{1}{2}\bar{x}^T \Sigma^{-1} \bar{x}}}{\sqrt{(2\pi)^d |\Sigma|}} d\bar{x}$$

$$= \sqrt{P(w_1) P(w_2)} e^{-\frac{1}{2}\bar{\mu}_1^T \Sigma^{-1} \bar{\mu}_2}$$

3

a) The within-class scatter matrix \underline{S}_w is:

$$\underline{S}_w = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}$$

The difference between the means is:

$$\mu_2 - \mu_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix} - \begin{pmatrix} -1 \\ -1 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

According to 1.c):

$$w = \underline{S}_w^{-1} (\mu_2 - \mu_1) = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 2 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

b)

The projection $w = \begin{pmatrix} -1 \\ 1 \end{pmatrix}$ minimizes the ratio

$$\text{because } w^T \mu_1 = (-1 \ 1) \begin{pmatrix} -1 \\ -1 \end{pmatrix} = 0$$

$$\text{and } w^T \mu_2 = (-1 \ 1) \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 0$$