

## T-61.231 Principles of Pattern Recognition

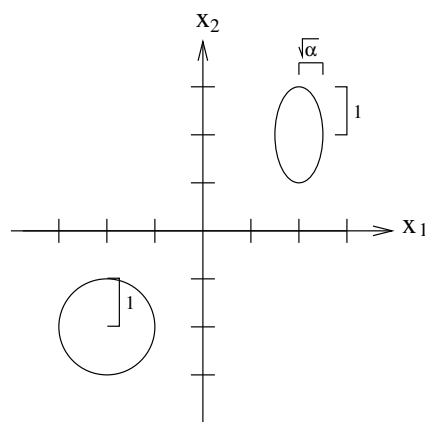
Answers to exercise 5: 21.10.2002

1. A Gaussian distribution has two parameters: a mean vector  $\bar{\mu}$  and covariance matrix  $\Sigma$ .

$$\hat{\mu}_i = \bar{m}_i \text{ and } \hat{\Sigma} = \frac{1}{n_i - 1} \sum_{k=1}^{n_i} (\bar{x}_k - \hat{\mu}_i)(\bar{x}_k - \hat{\mu}_i)^T$$

- a) Lets assume that there are approximately as many samples from both classes ( $n_1 \approx n_2$ ). Now we can draw both distributions in the same picture and use the same scale for both of them.

The density of class 1 is symmetric as the diagonal elements of  $S_1$  are equal. As for class 2, the density is expanded in the direction of the “width” of the distribution on the  $x_1$ -axis depends on  $\alpha$ .



b)  $\hat{w} = S_W^{-1}(\bar{m}_1 - \bar{m}_2)$

$$S_W = S_1 + S_2 = \begin{bmatrix} 1 + \alpha & 0 \\ 0 & 2 \end{bmatrix} \Rightarrow S_W^{-1} = \frac{1}{2(1 + \alpha)} \begin{bmatrix} 2 & 0 \\ 0 & 1 + \alpha \end{bmatrix} = \begin{bmatrix} \frac{1}{1 + \alpha} & 0 \\ 0 & \frac{1}{2} \end{bmatrix}$$

$$\hat{w} = \begin{bmatrix} \frac{1}{1 + \alpha} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} (-2 - 2) \\ (-2 - 2) \end{bmatrix} = -2 \begin{bmatrix} \frac{2}{1 + \alpha} \\ 1 \end{bmatrix}$$

- c) To determine the eigenvectors of  $S_W^{-1}S_B$  we first calculate the between-class scatter matrix  $S_B$  from  $S_B = (\bar{m}_1 - \bar{m}_2)(\bar{m}_1 - \bar{m}_2)^T$ .

$$S_B = \begin{bmatrix} (-2 - 2) \\ (-2 - 2) \end{bmatrix} \begin{bmatrix} (-2 - 2) & (-2 - 2) \end{bmatrix} = \begin{bmatrix} 16 & 16 \\ 16 & 16 \end{bmatrix}$$

Now

$$S_W^{-1}S_B = \begin{bmatrix} \frac{1}{1 + \alpha} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 16 & 16 \\ 16 & 16 \end{bmatrix} = \begin{bmatrix} \frac{16}{1 + \alpha} & \frac{16}{1 + \alpha} \\ 8 & 8 \end{bmatrix}$$

The eigenvectors for matrix  $A$  are defined as  $A\bar{x} = \lambda\bar{x} \Rightarrow (A - \lambda I)\bar{x} = \bar{0}$ . A nontrivial solution exists, if  $\det(A - \lambda I) = 0$

$$\begin{vmatrix} \frac{16}{1 + \alpha} - \lambda & \frac{16}{1 + \alpha} \\ 8 & 8 - \lambda \end{vmatrix} = 0 \Leftrightarrow \left(\frac{16}{1 + \alpha} - \lambda\right)(8 - \lambda) - 8 \frac{16}{1 + \alpha} = 0 \Leftrightarrow \lambda^2 - \left(8 + \frac{16}{1 + \alpha}\right)\lambda = 0$$

Thus the eigenvalues are  $\lambda_1 = 0$  and  $\lambda_2 = 8 + \frac{16}{1 + \alpha}$  (Note that  $\alpha = \sigma_1^2 > 0$ ).

The eigenvector corresponding to the larger eigenvalue:  $A\bar{e} = \lambda_2\bar{e}$ , where  $\bar{e} = [e_1 \ e_2]^T$ . Thus

$$\begin{bmatrix} \frac{16}{1+\alpha} & \frac{16}{1+\alpha} \\ 8 & 8 \end{bmatrix} \begin{bmatrix} e_1 \\ e_2 \end{bmatrix} = \left(8 + \frac{16}{1+\alpha}\right) \begin{bmatrix} e_1 \\ e_2 \end{bmatrix}$$

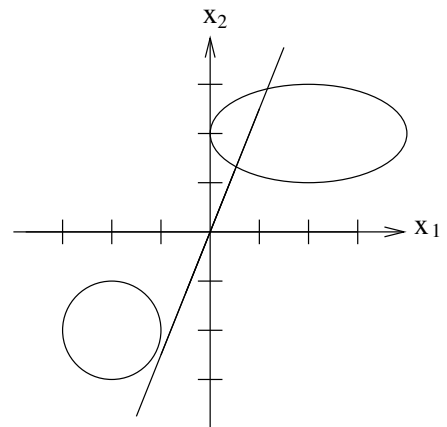
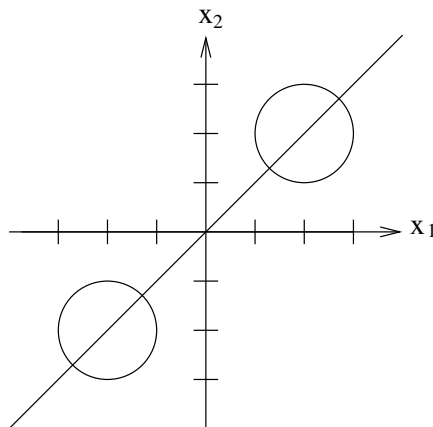
Calculating from the lower row (the same result can be obtained from the upper row, too)

$$8e_1 + 8e_2 = \left(8 + \frac{16}{1+\alpha}\right)e_2 \Leftrightarrow e_1 = \frac{2}{1+\alpha}e_2 \Rightarrow e = \begin{bmatrix} \frac{2}{1+\alpha} \\ 1 \end{bmatrix}$$

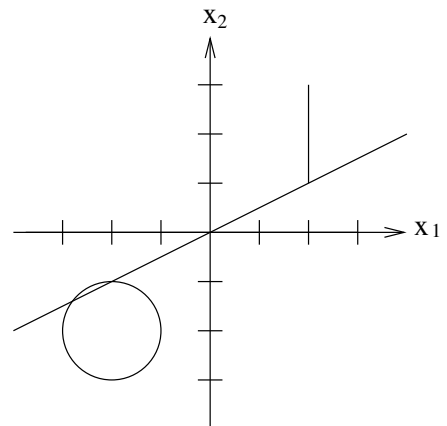
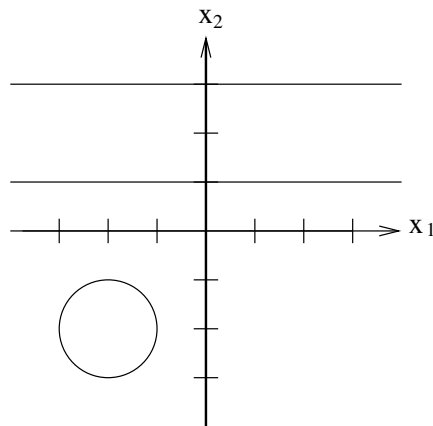
As can be seen, both methods produced a vector of the same orientation, as was to be expected.

- d) The Fisher linear discriminant is thus  $\hat{w} = \begin{bmatrix} \frac{2}{1+\alpha} \\ 1 \end{bmatrix}$ ,  $\underline{m}_1 = (-2 \ -2)^T$ ,  $\underline{m}_2 = (2 \ 2)^T$ ,  
 $\mathbf{S}_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  and  $\mathbf{S}_2 = \begin{bmatrix} \alpha & 0 \\ 0 & 1 \end{bmatrix}$ .

$$\alpha = 1 \Rightarrow \hat{w} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \text{ and } S_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} : \quad \alpha = 4 \Rightarrow \hat{w} = \begin{bmatrix} 2/5 \\ 1 \end{bmatrix} \text{ and } S_2 = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix} :$$



$$\alpha \rightarrow \infty \Rightarrow \hat{w} \rightarrow \begin{bmatrix} 0 \\ 1 \end{bmatrix} \text{ and } S_2 \rightarrow \begin{bmatrix} \infty & 0 \\ 0 & 1 \end{bmatrix} : \quad \alpha \rightarrow 0 \Rightarrow \hat{w} \rightarrow \begin{bmatrix} 2 \\ 1 \end{bmatrix} \text{ and } S_2 \rightarrow \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} :$$



The discriminants seem quite valid, as it can be seen that by projecting the distributions onto the discriminant in all cases the distributions become well separable.

2. Fisher's linear discriminant is based on projecting the original data onto a line in the direction of  $\bar{w}$  ( $\|\bar{w}\| = 1$ ) so that the criterion function

$$J(\bar{w}) = \frac{(m_1 - m_2)^2}{\hat{\sigma}_1^2 + \hat{\sigma}_2^2},$$

where  $m_i$  and  $\hat{\sigma}_i^2$  are the projected mean and within-class scatter of the projected data for class  $i$ , is maximized. In this exercise, we consider an alternative criterion function  $J'(\bar{w}) = (m_1 - m_2)^2$  consisting of only the nominator of  $J(\bar{w})$ .

$$J'(\bar{w}) = (m_1 - m_2)^2 = (\bar{w}^T \bar{m}_1 - \bar{w}^T \bar{m}_2)^2 = [\bar{w}^T (\bar{m}_1 - \bar{m}_2)]^2$$

This is clearly maximized as  $\bar{w} \parallel (\bar{m}_1 - \bar{m}_2)$  i.e. when Fisher's discriminant is parallel to the line fixed by the sample means  $\bar{m}_1$  and  $\bar{m}_2$ .

3. Let us denote the function slightly differently for the proof; let  $l = d - 1$  and  $O(N, l) = C(N, d)$ , where

$$C(N, d) = 2 \sum_{k=0}^{d-1} \binom{N-1}{k}.$$

In my opinion this is more illustrative, as we are actually using the dimension of the space  $d$  instead of constantly using  $l$  to denote  $l - 1$  dimensional space.  $C(N, d)$  tells us the number of groupings that can be formed by  $d$ -dimensional hyperplanes to separate the  $N$  points into two classes.

First we need to prove that  $C(N + 1, d) = C(N, d) + C(N, d - 1)$ .

Let  $C(N, d)$  be a separable set of dichotomies  $X$ . Let's take a new point  $x_{N+1}$  so that  $X \cup \{x_{N+1}\}$  is in the general position (well distributed). Let there be a vector  $w$  that divides  $X$  into two sets  $X = \{X^+, X^-\}$  so that  $w \cdot x > t \Rightarrow x \in X^+$  and  $w \cdot x < t \Rightarrow x \in X^-$ , where  $t$  is a scalar.

If  $\{X^+, X^-\}$  is separable, must also either  $\{X^+ \cup \{x_{N+1}\}, X^-\}$  or  $\{X^+, X^- \cup \{x_{N+1}\}\}$  be separable. However, they both are separable if and only if  $\exists w$  that is a vector that separates  $\{X^+, X^-\}$  in a  $(d - 1)$  dimensional space and is orthogonal to  $x_{N+1}$

To prove the prior statement regarding  $w$ , let the set of separating vectors  $W = \{w : w \cdot x > t, x \in X^+; w \cdot x < t, x \in X^-\}$ . The set  $\{X^+ \cup \{x_{N+1}\}, X^-\}$  is homogeneously separable if and only if  $\exists w \in W$  so that  $w \cdot x_{N+1} > t$ , and equivalently  $\{X^+, X^- \cup \{x_{N+1}\}\}$  is homogeneously separable if and only if  $\exists w \in W$  so that  $w \cdot x_{N+1} < t$ . Let the sets be linearly separable with  $w_1$  and  $w_2$ , respectively. Then  $w^* = (-w_2 \cdot x_{N+1})w_1 + (w_1 \cdot x_{N+1})w_2$  separates  $\{X^+, X^-\}$  by the hyperplane  $\{x : w^* \cdot x = t\}$  passing through  $x_{N+1}$ . Conversely, if the sets  $\{X^+, X^-\}$  are homogeneously linearly separable by a hyperplane containing  $x_{N+1}$ , then  $\exists w^* \in W$  so that  $w^* \cdot x = t$ . Since  $W$  is an open set,  $\exists \epsilon > 0$  so that  $w^* + \epsilon x_{N+1}$  and  $w^* - \epsilon x_{N+1}$  are in  $W$ . Hence  $\{X^+ \cup \{x_{N+1}\}, X^-\}$  and  $\{X^+, X^- \cup \{x_{N+1}\}\}$  are homogeneously linearly separable by  $w^* + \epsilon x_{N+1}$  and  $w^* - \epsilon x_{N+1}$ , respectively.

So the set can be separated if and only if  $\exists w$  so that the projection onto a  $(d - 1)$  dimensional subspace is separable. By the induction hypothesis there are  $C(N, d - 1)$  such separable dichotomies. Hence,

$$C(N + 1, d) = C(N, d) + C(N, d - 1)$$

By repeatedly applying of this to the terms on the right we obtain

$$C(N, d) = \sum_{k=0}^{N-1} \binom{N-1}{k} C(1, d - k)$$

Now, as it is obvious that one point can be separated in two ways if the dimension is greater or equal to 1 and no separation can be made when the dimension is below one, or

$$C(1, m) = \begin{cases} 2, m \geq 1 \\ 0, m < 1 \end{cases}$$

The original theorem follows by separating the part of the sum where  $d-k < 1 \Leftrightarrow k > d-1$ :

$$C(N, d) = 2 \sum_{k=0}^{d-1} \binom{N-1}{k} + 0 \cdot \sum_{k=d}^{N-1} \binom{N-1}{k} = 2 \sum_{k=0}^{d-1} \binom{N-1}{k} \Leftrightarrow O(n, l) = 2 \sum_{k=0}^l \binom{N-1}{k}$$