

2.6 Discriminant Functions for the Normal Density

2.6 Discriminant Functions for the Normal Density

In Sect. 2.4.1 we saw that the minimum-error-rate classification can be achieved by use of the discriminant functions

$$g_i(\mathbf{x}) = \ln p(\mathbf{x}|\omega_i) + \ln P(\omega_i). \quad (46)$$

This expression can be readily evaluated if the densities $p(\mathbf{x}|\omega_i)$ are multivariate normal, i.e., if $p(\mathbf{x}|\omega_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$. In this case, then, from Eq. 37 we have

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^t \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i). \quad (47)$$

Let us examine the discriminant function and resulting classification for a number of special cases.

2.6.1 Case 1: $\Sigma_i = \sigma^2 \mathbf{I}$

The simplest case occurs when the features are statistically independent, and when each feature has the same variance σ^2 .

$$g_i(\mathbf{x}) = -\frac{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{2\sigma^2} + \ln P(\omega_i), \quad (48)$$

Simplificando

we obtain the equivalent **linear discriminant functions**

$$g_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + w_{i0}, \quad (51)$$

where

$$\mathbf{w}_i = \frac{1}{\sigma^2} \boldsymbol{\mu}_i \quad (52)$$

and

$$w_{i0} = \frac{-1}{2\sigma^2} \boldsymbol{\mu}_i^t \boldsymbol{\mu}_i + \ln P(\omega_i). \quad (53)$$

We call w_{i0} the *threshold* or *bias* in the i th direction.

2.6.1 Case 1: $\Sigma_i = \sigma^2 \mathbf{I}$

Simplificando más

$$\mathbf{w}^t(\mathbf{x} - \mathbf{x}_0) = 0, \quad (54)$$

where

$$\mathbf{w} = \boldsymbol{\mu}_i - \boldsymbol{\mu}_j \quad (55)$$

and

$$\mathbf{x}_0 = \frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \frac{\sigma^2}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2} \ln \frac{P(\omega_i)}{P(\omega_j)} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j). \quad (56)$$

This equation defines a hyperplane through the point \mathbf{x}_0 and orthogonal to the vector \mathbf{w} . Since $\mathbf{w} = \boldsymbol{\mu}_i - \boldsymbol{\mu}_j$, the hyperplane separating \mathcal{R}_i and \mathcal{R}_j is orthogonal to the line linking the means. If $P(\omega_i) = P(\omega_j)$, the second term on the right of Eq. 56 vanishes, and thus the point \mathbf{x}_0 is halfway between the means, and the hyperplane is the perpendicular bisector of the line between the means (Fig. 2.11). If $P(\omega_i) \neq P(\omega_j)$, the point \mathbf{x}_0 shifts away from the more likely mean

2.6.1 Case 1: $\Sigma_i = \sigma^2 \mathbf{I}$

If the prior probabilities $P(\omega_i)$ are the same for all c classes, then the $\ln P(\omega_i)$ term becomes another unimportant additive constant that can be ignored.

When this happens, the optimum decision rule can be stated very simply:

to classify a feature vector \mathbf{x} , measure the Euclidean distance $\|\mathbf{x} - \mu_i\|$ from each \mathbf{x} to each of the c mean vectors, and assign \mathbf{x} to the category of the nearest mean.

Such a classifier is called a **minimum distance classifier**.

If each mean vector is thought of as being an ideal prototype or template for patterns in its class, then this is essentially a **template-matching** procedure.

2.6.1 Case 1: $\Sigma_i = \sigma^2 \mathbf{I}$

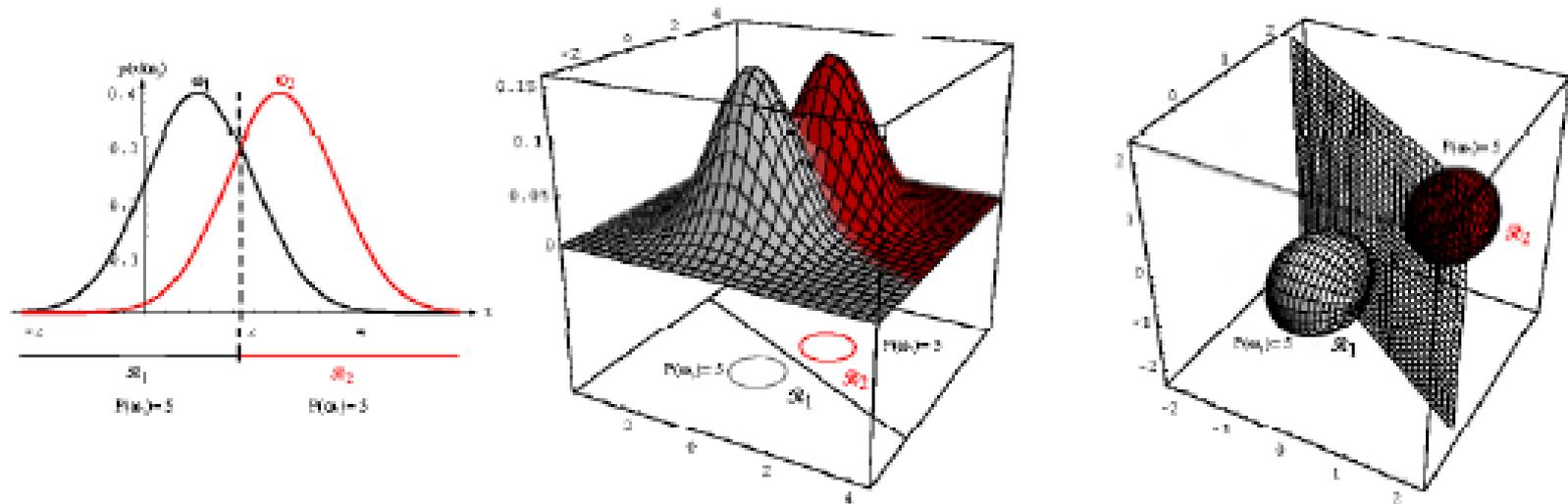


Figure 2.10: If the covariances of two distributions are equal and proportional to the identity matrix, then the distributions are spherical in d dimensions, and the boundary is a generalized hyperplane of $d - 1$ dimensions, perpendicular to the line separating the means. In these 1-, 2-, and 3-dimensional examples, we indicate $p(\mathbf{x}|\omega_i)$ and the boundaries for the case $P(\omega_1) = P(\omega_2)$. In the 3-dimensional case, the grid plane separates \mathcal{R}_1 from \mathcal{R}_2 .

2.6.1 Case 1: $\Sigma_i = \sigma^2 \mathbf{I}$

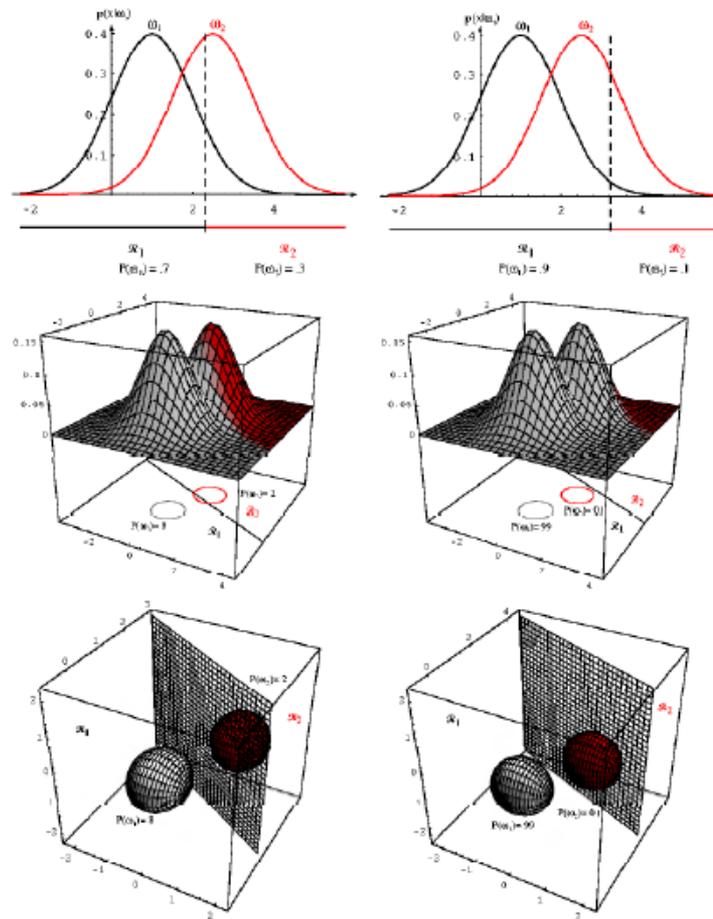


Figure 2.11: As the priors are changed, the decision boundary shifts; for sufficiently disparate priors the boundary will not lie between the means of these 1-, 2- and 3-dimensional spherical Gaussian distributions.

2.6.2 Case 2: $\Sigma_i = \Sigma$

Another simple case arises when the covariance matrices for all of the classes are identical but otherwise arbitrary. Geometrically, this corresponds to the situation in which the samples fall in hyperellipsoidal clusters of equal size and shape, the cluster for the i th class being centered about the mean vector μ_i . Since both $|\Sigma_i|$ and the $(d/2) \ln 2\pi$ term in Eq. 47 are independent of i , they can be ignored as superfluous additive constants. This simplification leads to the discriminant functions

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mu_i)^t \Sigma^{-1}(\mathbf{x} - \mu_i) + \ln P(\omega_i). \quad (57)$$

Simplificando

$$g_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + w_{i0}, \quad (58)$$

where

$$\mathbf{w}_i = \Sigma^{-1} \mu_i \quad (59)$$

and

$$w_{i0} = -\frac{1}{2} \mu_i^t \Sigma^{-1} \mu_i + \ln P(\omega_i). \quad (60)$$

2.6.2 Case 2: $\Sigma_i = \Sigma$

Since the discriminants are linear, the resulting decision boundaries are again **hyperplanes**.

The hyperplane separating R_i and R_j is **generally not orthogonal** to the line between the means.

However, it does intersect that line at the point x_0 which is halfway between the means if the prior probabilities are equal. If the prior probabilities are not equal, the optimal boundary hyperplane is shifted away from the more likely mean (Fig. 2.12).

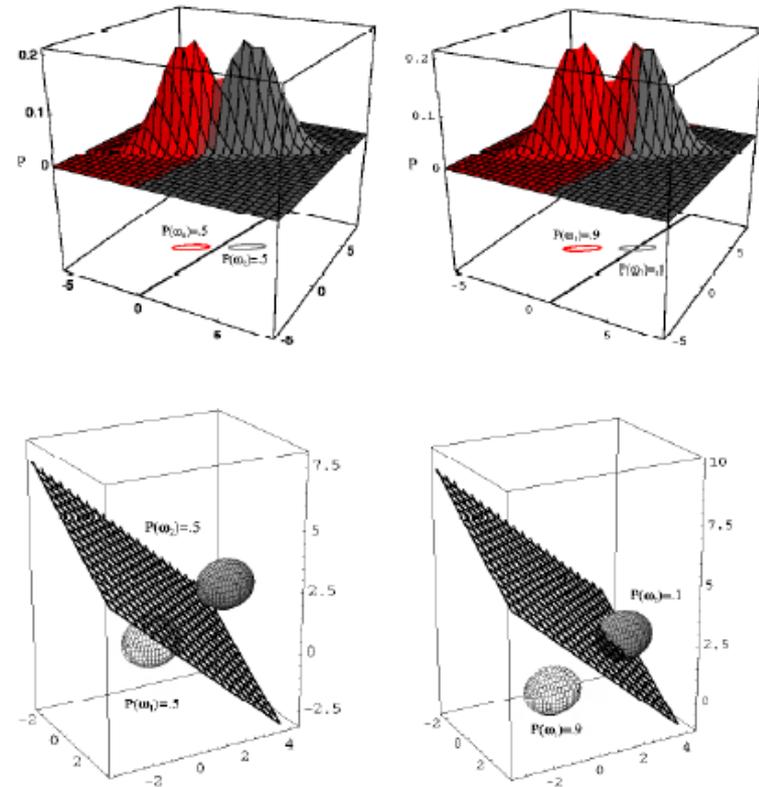


Figure 2.12: Probability densities (indicated by the surfaces in two dimensions and ellipsoidal surfaces in three dimensions) and decision regions for equal but asymmetric Gaussian distributions. The decision hyperplanes need not be perpendicular to the line connecting the means.

2.6.3 Case 3: $\Sigma_i = \text{arbitrary}$

In the general multivariate normal case, the covariance matrices are different for each category.

The resulting discriminant functions are inherently quadratic.

$$g_i(\mathbf{x}) = \mathbf{x}^t \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^t \mathbf{x} + w_{i0}, \quad (64)$$

where

$$\mathbf{W}_i = -\frac{1}{2} \Sigma_i^{-1}, \quad (65)$$

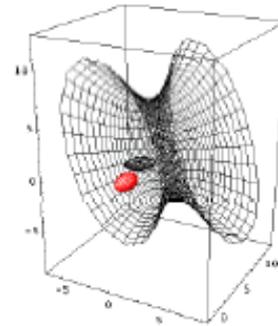
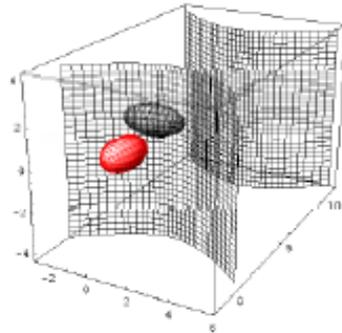
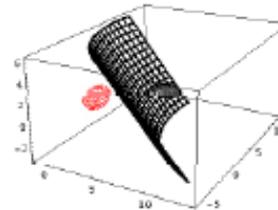
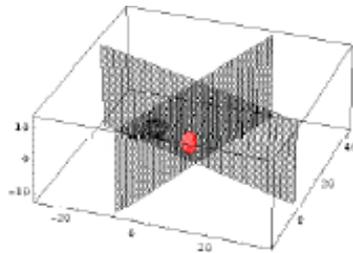
$$\mathbf{w}_i = \Sigma_i^{-1} \boldsymbol{\mu}_i \quad (66)$$

and

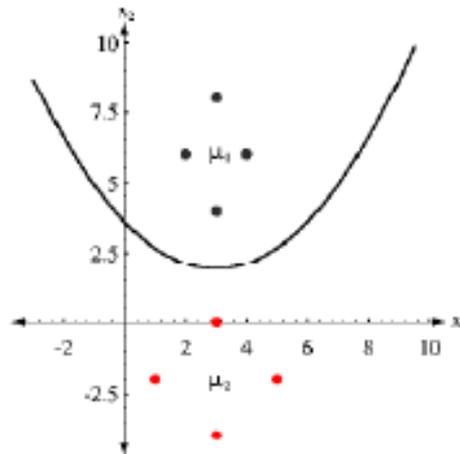
$$w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^t \Sigma_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i). \quad (67)$$

2.6.3 Case 3: $\Sigma_i = \text{arbitrary}$

The decision surfaces are hyperquadrics, and can assume any of the general forms: hyperplanes, pairs of hyperplanes, hyperspheres, hyperellipsoids, hyperparaboloids...



To clarify these ideas, we explicitly calculate the decision boundary for the two category two-dimensional data.



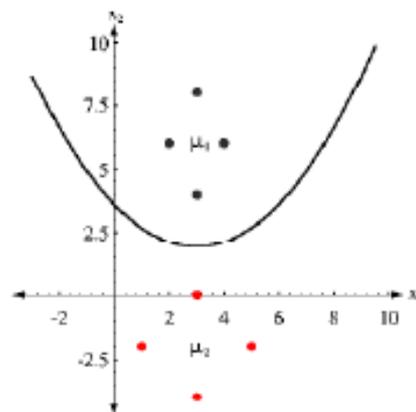
The computed Bayes decision boundary for two Gaussian distributions, each based on four data points.

Let w_1 be the set of the four black points, and w_2 the red points. For now we simply assume that we need merely calculate the means and covariances

$$\mu_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}; \quad \Sigma_1 = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix} \text{ and } \mu_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix}; \quad \Sigma_2 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}.$$

We assume equal prior probabilities, $P(\omega_1) = P(\omega_2) = 0.5$, and substitute these into the general form for a discriminant, Eqs. 64 – 67, setting $g_1(\mathbf{x}) = g_2(\mathbf{x})$ to obtain the decision boundary:

$$x_2 = 3.514 - 1.125x_1 + 0.1875x_1^2.$$



The computed Bayes decision boundary for two Gaussian distributions, each based on four data points.