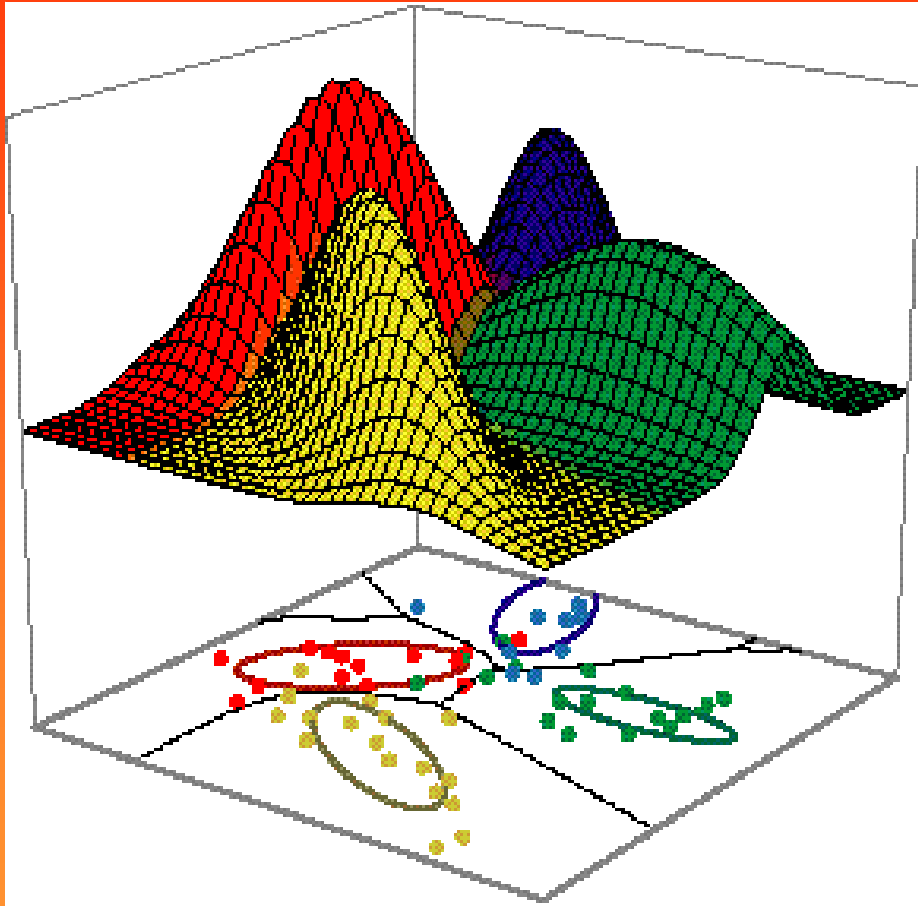


Pattern Classification



All materials in these slides were taken from
Pattern Classification (2nd ed) by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000
with the permission of the authors and the publisher

Chapter 2

Bayesian Decision Theory

(Sections 2-6,2-9)

- Discriminant Functions for the Normal Density
- Bayes Decision Theory – Discrete Features

2.6 Discriminant Functions for the Normal Density

- We saw that the minimum error-rate classification can be achieved by the discriminant function

$$g_i(x) = \ln p(x | \omega_i) + \ln P(\omega_i)$$

- Case of multivariate normal

$$p(x | \omega_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[-\frac{1}{2} (x - \mu_i)^t \Sigma_i^{-1} (x - \mu_i) \right]$$



$$g_i(x) = -\frac{1}{2} (x - \mu_i)^t (\Sigma_i)^{-1} (x - \mu_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

Case $\Sigma_i = \sigma^2.I$ (I stands for the identity matrix)

$$g_i(x) = -\frac{\|X - \mu_i\|^2}{2\sigma^2} + \ln P(\omega_i)$$



- Case $\Sigma_i = \sigma^2.I$ (I stands for the identity matrix)

$g_i(x) = w_i^t x + w_{i0}$ (*linear discriminant function*)

where :

$$w_i = \frac{\mu_i}{\sigma^2}; \quad w_{i0} = -\frac{1}{2\sigma^2} \mu_i^t \mu_i + \ln P(\omega_i)$$

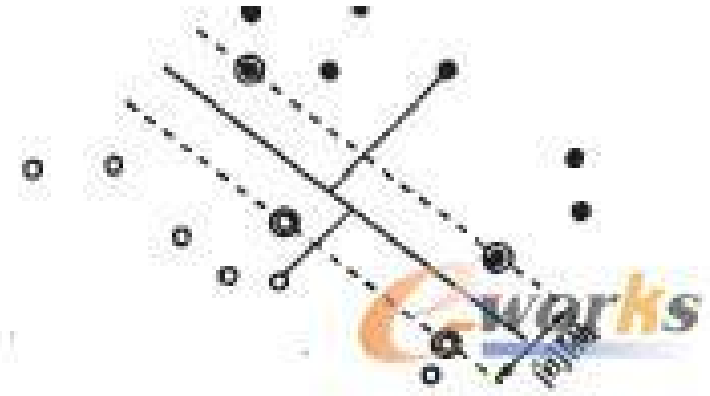
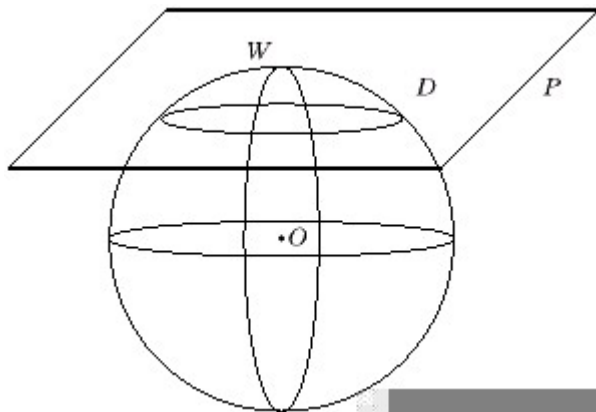
(ω_{i0} is called the threshold for the i th category!)



□ A classifier that uses linear discriminant functions is called “a linear machine”

□ The decision surfaces for a linear machine are pieces of hyperplanes defined by:

$$g_i(x) = g_j(x)$$



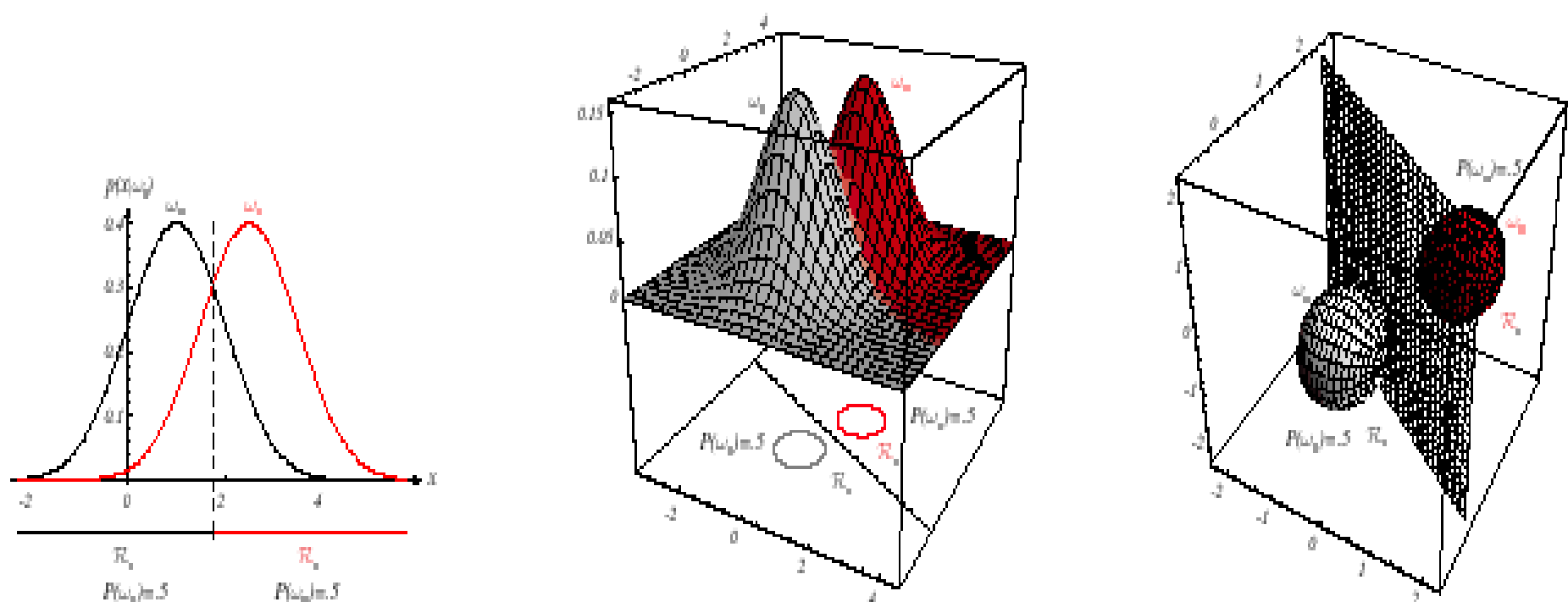


FIGURE 2.10. If the covariance matrices for two distributions are equal and proportional to the identity matrix, then the distributions are spherical in d dimensions, and the boundary is a generalized hyperplane of $d - 1$ dimensions, perpendicular to the line separating the means. In these one-, two-, and three-dimensional examples, we indicate $p(\mathbf{x}|\omega_i)$ and the boundaries for the case $P(\omega_1) = P(\omega_2)$. In the three-dimensional case, the grid plane separates \mathcal{R}_1 from \mathcal{R}_2 . From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



- The hyperplane separating \mathcal{R}_i and \mathcal{R}_j is always orthogonal to the line linking the means!

$$g_i(x) = g_j(x)$$

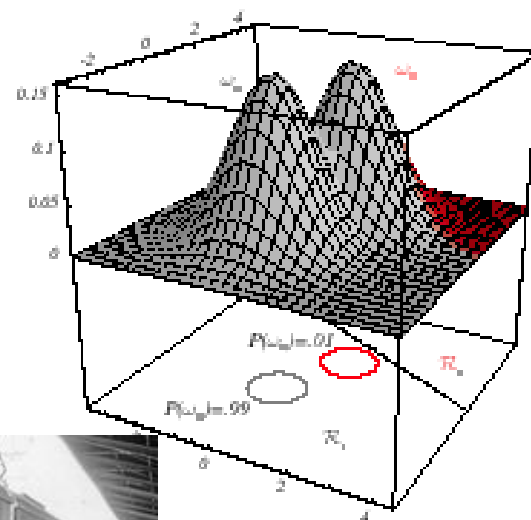
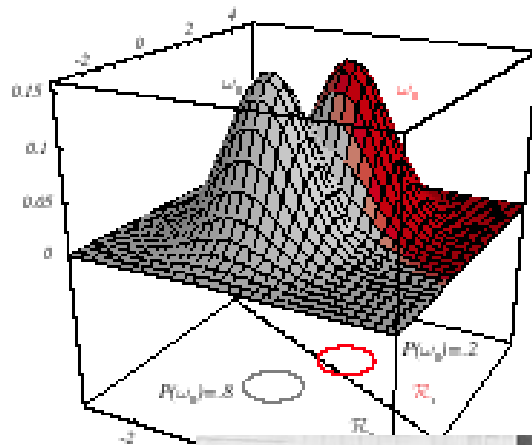
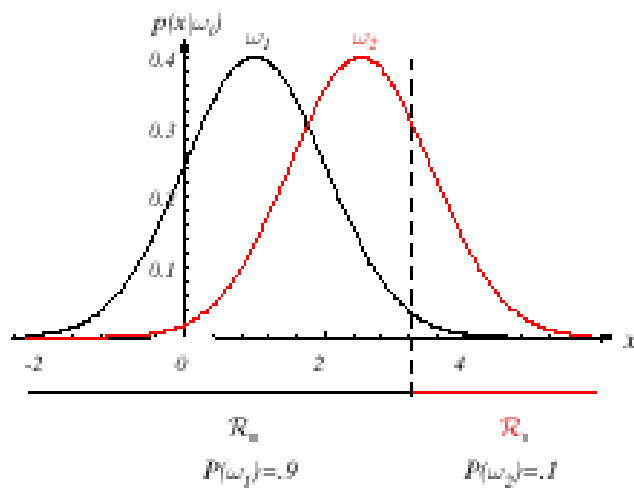
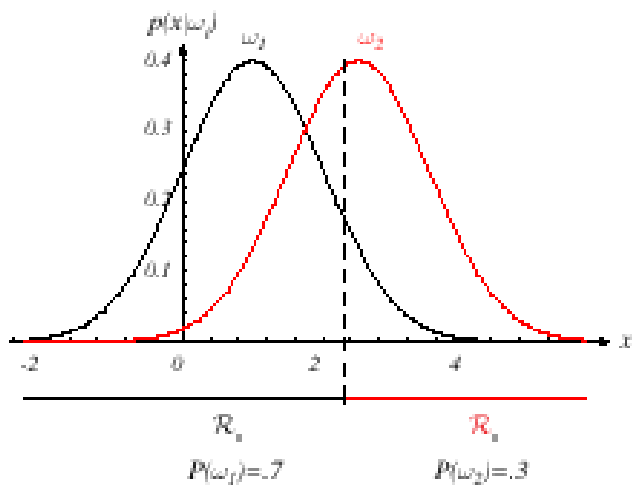


$$W^t (x - x_0) = 0$$

$$W = \mu_i - \mu_j$$

$$x_0 = \frac{1}{2}(\mu_i + \mu_j) - \frac{\sigma^2}{\|\mu_i - \mu_j\|^2} \ln \frac{P(\omega_i)}{P(\omega_j)} (\mu_i - \mu_j)$$

$$\text{if } P(\omega_i) = P(\omega_j) \text{ then } x_0 = \frac{1}{2}(\mu_i + \mu_j)$$



part 3)

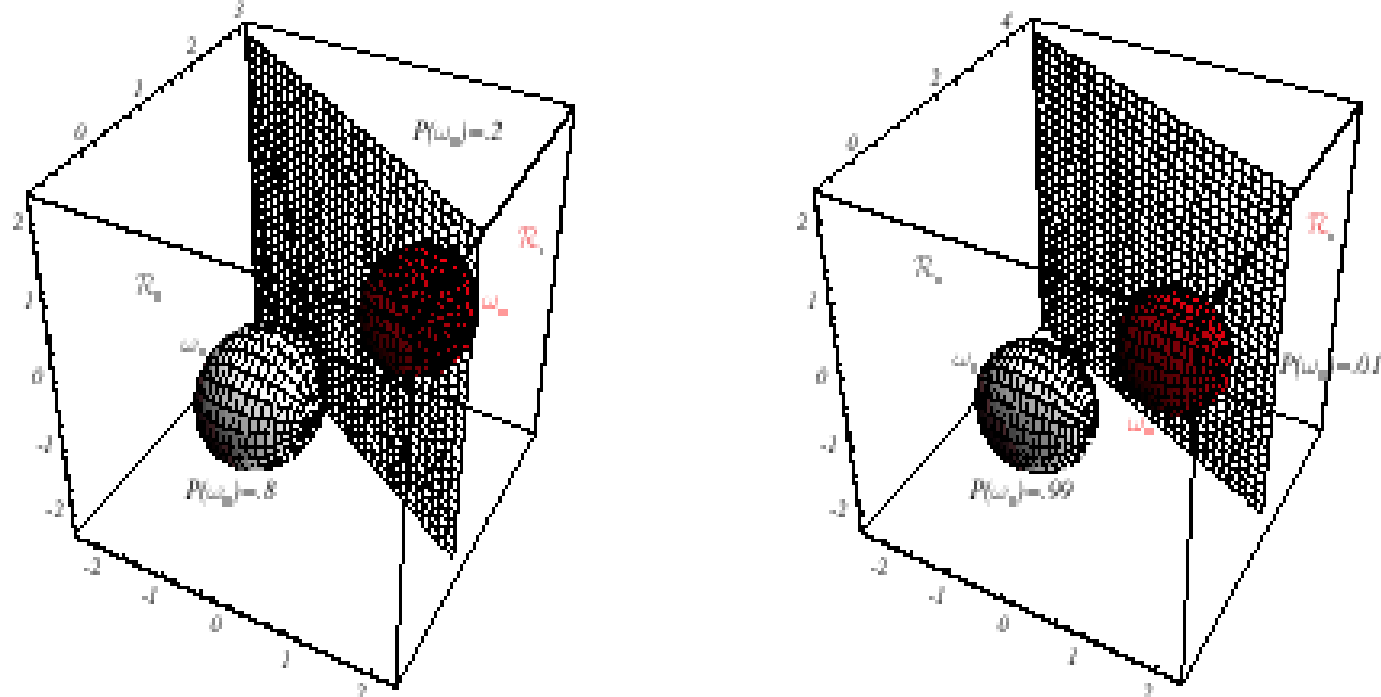


FIGURE 2.11. As the priors are changed, the decision boundary shifts; for sufficiently disparate priors the boundary will not lie between the means of these one-, two- and three-dimensional spherical Gaussian distributions. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



- Case $\Sigma_i = \Sigma$ (covariance of all classes are identical but arbitrary!)



$$g_i(x) = -\frac{1}{2}(X - \mu_i)^T \Sigma^{-1}(X - \mu_i) + \ln P(\omega_i)$$

- Case $\Sigma_i = \Sigma$ (covariance of all classes are identical but arbitrary!)

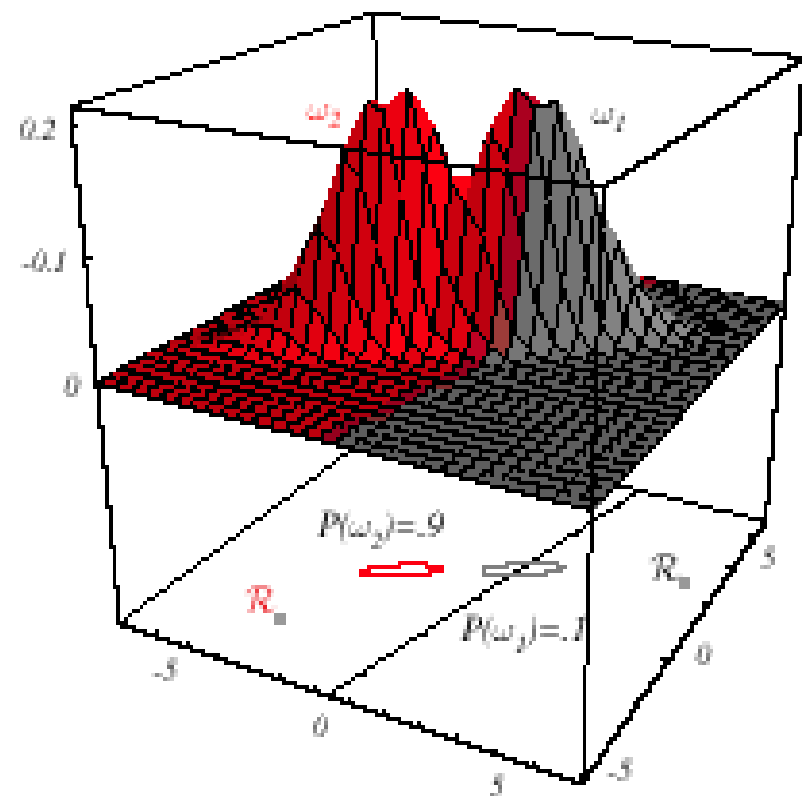
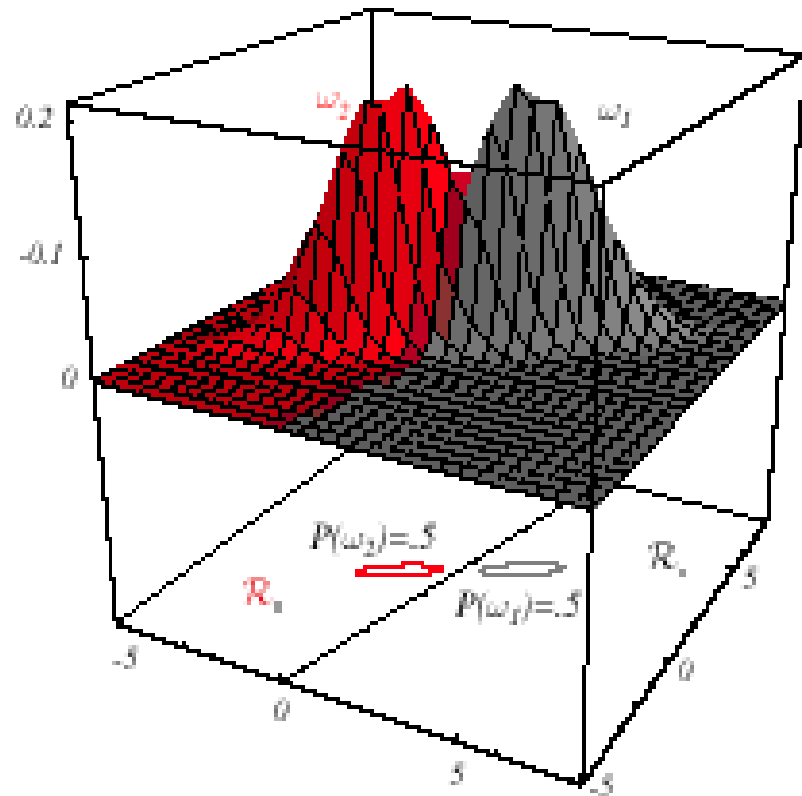
- Hyperplane separating \mathcal{R}_i and \mathcal{R}_j

$$w^t (x - x_0) = 0$$

$$w = \Sigma^{-1}(\mu_i - \mu_j)$$

$$x_0 = \frac{1}{2}(\mu_i + \mu_j) - \frac{\ln[P(\omega_i) / P(\omega_j)]}{(\mu_i - \mu_j)^t \Sigma^{-1}(\mu_i - \mu_j)} \cdot (\mu_i - \mu_j)$$

(the hyperplane separating \mathcal{R}_i and \mathcal{R}_j is generally not orthogonal to the line between the means!)



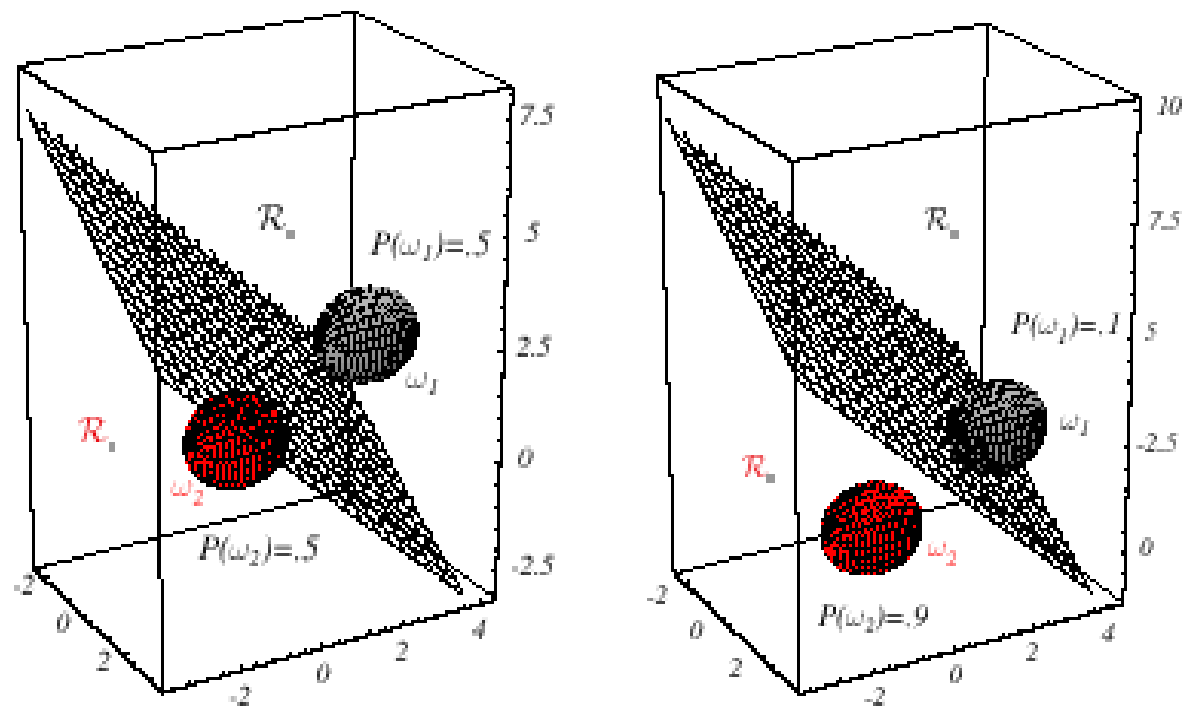


FIGURE 2.12. Probability densities (indicated by the surfaces in two dimensions and ellipsoidal surfaces in three dimensions) and decision regions for equal but asymmetric Gaussian distributions. The decision hyperplanes need not be perpendicular to the line connecting the means. From: Richard O. Duda, Peter E. Hart, and David C. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



■ Case $\Sigma_i = \text{arbitrary}$

- The covariance matrices are different for each category

$$g_i(x) = x^t W_i x + w_i^t x + w_{i0}$$

where :

$$W_i = -\frac{1}{2} \Sigma_i^{-1}$$

$$w_i = \Sigma_i^{-1} \mu_i$$

$$w_{i0} = -\frac{1}{2} \mu_i^t \Sigma_i^{-1} \mu_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$



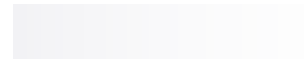
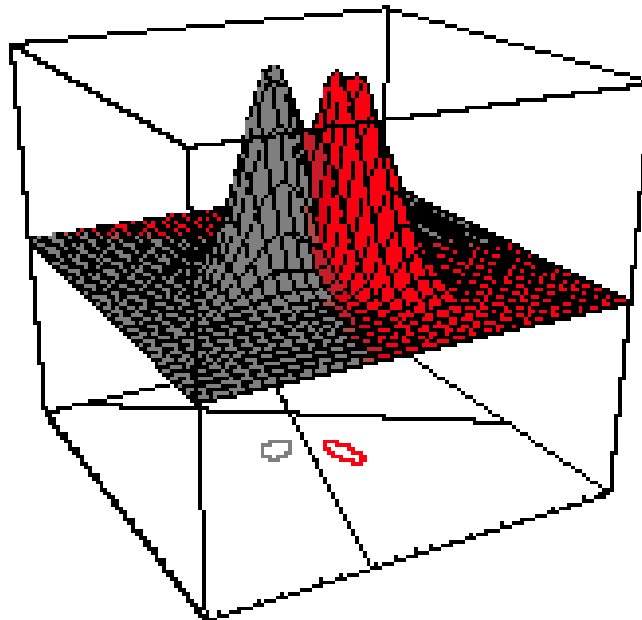
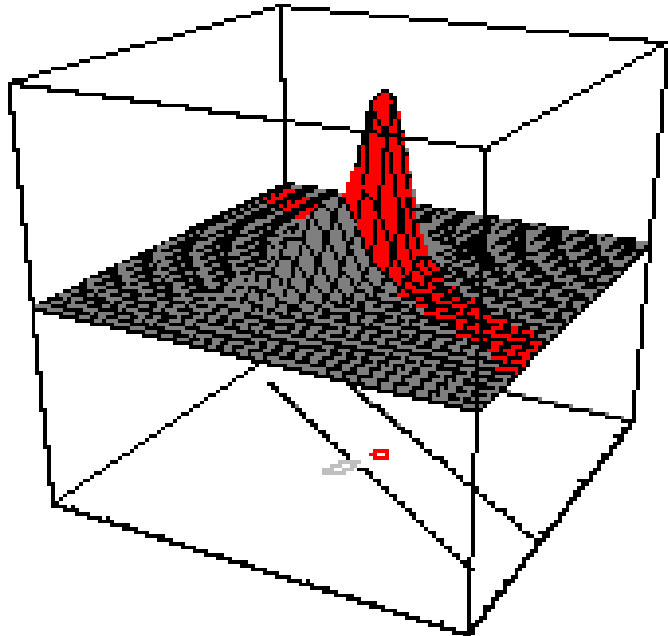
(Hyperquadrics which are: hyperplanes, pairs of hyperplanes, hyperspheres, hyperellipsoids, hyperparaboloids, hyperboloids)

Case $\Sigma_i = \text{arbitrary}$

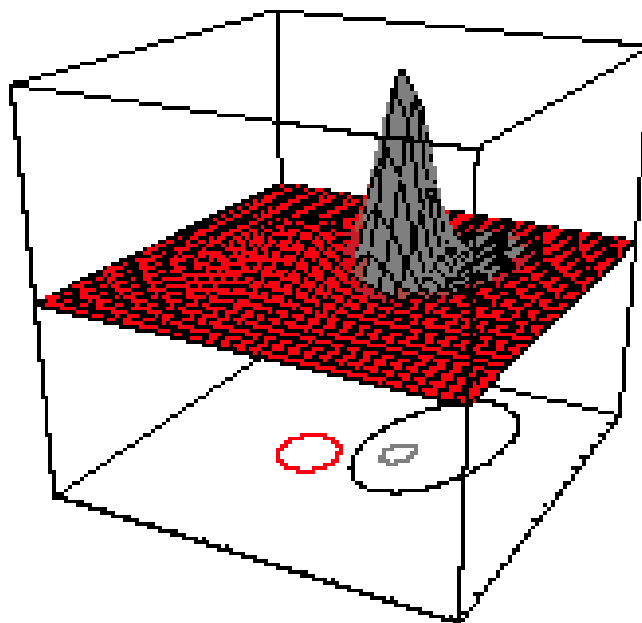
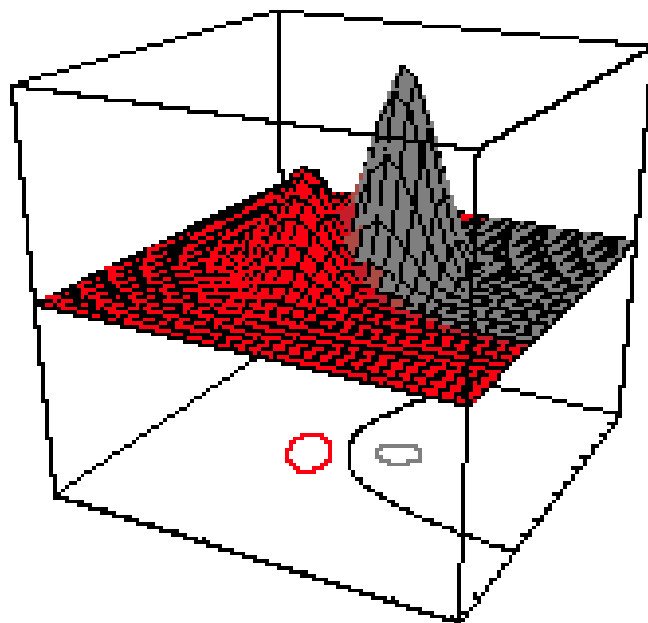
The covariance matrices are
different for each category

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^t (\Sigma_i)^{-1} (x - \mu_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$





吃饭



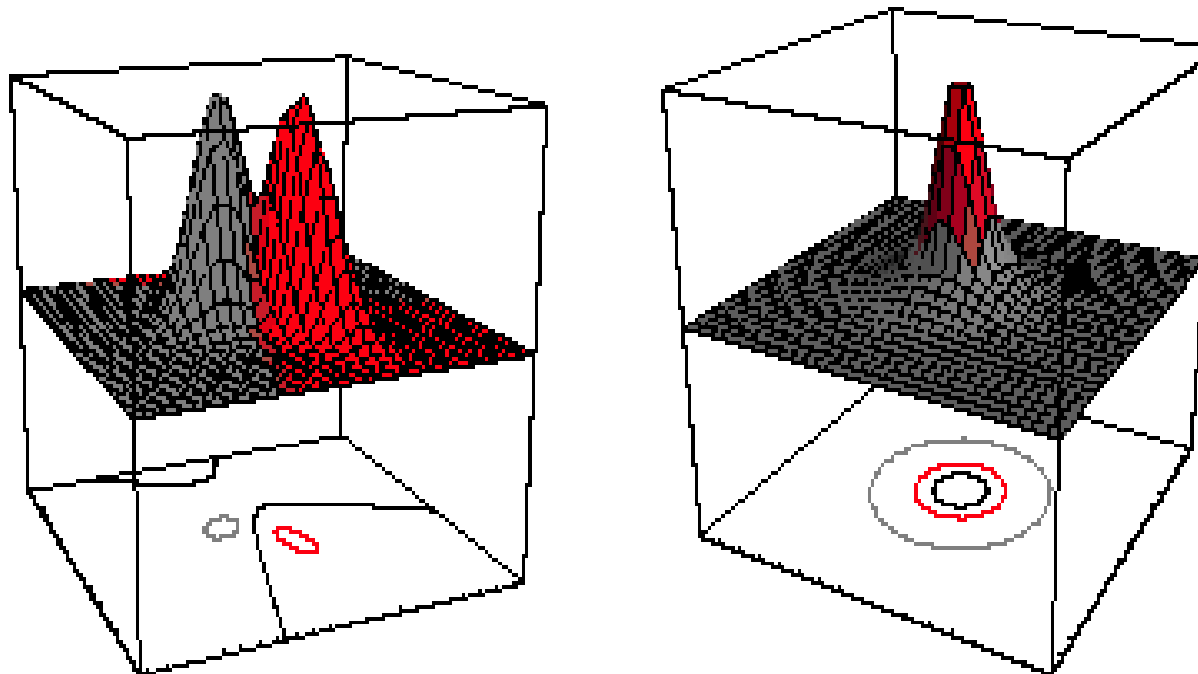


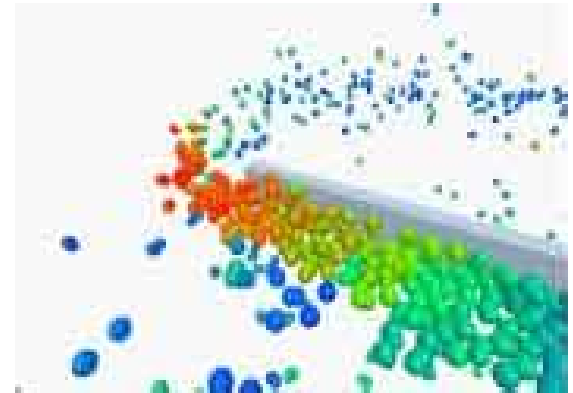
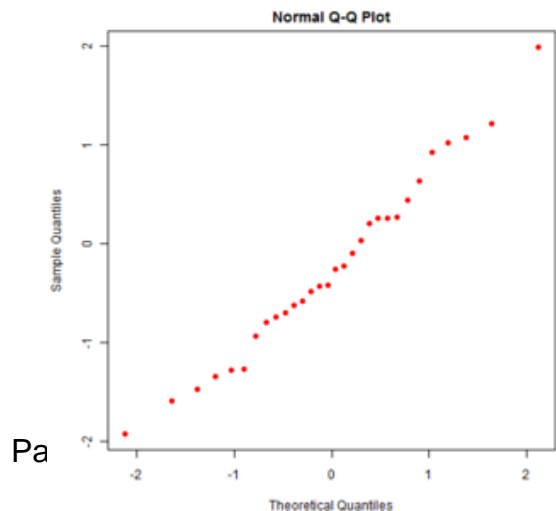
FIGURE 2.14. Arbitrary Gaussian distributions lead to Bayes decision boundaries that are general hyperquadrics. Conversely, given any hyperquadric, one can find two Gaussian distributions whose Bayes decision boundary is that hyperquadric. These variances are indicated by the contours of constant probability density. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

2.9 Bayes Decision Theory – Discrete Features

- Components of x are binary or integer valued, x can take only one of m discrete values

$$V_1, V_2, \dots, V_m$$

$$P(\omega_j | x) = \frac{P(x | \omega_j)P(\omega_j)}{P(x)}$$



It's hard to evaluate the probability of discrete high-dimensional sample

- $x=[\text{he and she first went shopping.....}]$
- How to evaluate the value of $p(x)$?



■ Case of independent binary features in 2 category problem

Let $x = [x_1, x_2, \dots, x_d]^t$ where each x_i is either 0 or 1, with probabilities:

$$p_i = P(x_i = 1 \mid \omega_1) \quad (P(x_i = 0 \mid \omega_1) = 1 - p_i)$$

$$q_i = P(x_i = 1 \mid \omega_2) \quad (P(x_i = 0 \mid \omega_2) = 1 - q_i)$$

$$P(X \mid \omega_1) = \prod_{i=1}^d p_i^{x_i} (1 - p_i)^{1-x_i}$$
$$P(X \mid \omega_2) = \prod_{i=1}^d q_i^{x_i} (1 - q_i)^{1-x_i}$$
$$\frac{P(X \mid \omega_1)}{P(X \mid \omega_2)} = \prod_{i=1}^d \left(\frac{p_i}{q_i}\right)^{x_i} \left(\frac{1-p_i}{1-q_i}\right)^{1-x_i}$$



- The discriminant function in this case is:

$$g(x) = \ln \frac{p(x | \omega_1)}{p(x | \omega_2)} + \ln \frac{p(\omega_1)}{p(\omega_2)}$$

$$g(x) = \sum_{i=1}^d w_i x_i + w_0$$

where :

$$w_i = \ln \frac{p_i(1 - q_i)}{q_i(1 - p_i)} \quad i = 1, \dots, d$$

and :

$$w_0 = \sum_{i=1}^d \ln \frac{1 - p_i}{1 - q_i} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

decide ω_1 if $g(x) > 0$ and ω_2 if $g(x) \leq 0$



Assumption

- $V=[v_1, v_2, \dots, v_m]$

- $P(v/w_i)=P(v_1/w_i)\dots P(v_m/w_i)$

- *For natural language processing,*
 v_1, v_2, \dots, v_m *denote different words*



- Bayes Decision Theory on discrete features have been widely used in natural language processing, speech recognition, text classification, OCR et al.



整体的NLP技术体系

