

# Chapter 2 (part 3)

## Bayesian Decision Theory

### (Sections 2-6,2-9)

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

# Discriminant Functions for the Normal Density

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

$$g_i(x) = \ln P(x | w_i) + \ln P(w_i)$$

- Case of multivariate normal

$$P(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (x - m)^T \Sigma^{-1} (x - m) \right]$$

$$g_i(x) = -\frac{1}{2} (x - m_i)^T \Sigma_i^{-1} (x - m_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(w_i)$$

$$P(x) = \frac{1}{(2p)^{d/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (x - \mathbf{m})^T \Sigma^{-1} (x - \mathbf{m}) \right]$$

- Mahalanobis distance

The points of constant density are hyperellipsoids for which the quadratic form is constant

$$r^2 = \left[ (x - \mathbf{m})^T \Sigma^{-1} (x - \mathbf{m}) \right]$$

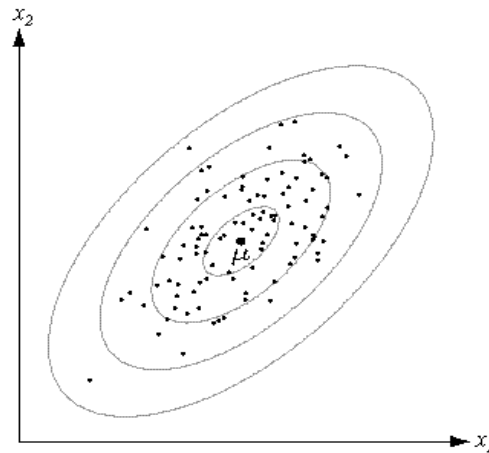


FIGURE 2.9. Samples drawn from a two-dimensional Gaussian lie in a cloud centered on the mean  $\mu$ . The ellipses show lines of equal probability density of the Gaussian. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- Euclidean distance

$$r^2 = \left[ (x - \mathbf{m})^T (x - \mathbf{m}) \right]$$

## First (independent) case

- Case  $S_i = s^2.I$  (I stands for the identity matrix)

$$g_i(x) = -\frac{1}{2}(x - \mathbf{m}_i)^T \Sigma_i^{-1}(x - \mathbf{m}_i) - \frac{d}{2} \ln 2p - \frac{1}{2} \ln |\Sigma_i| + \ln P(w_i)$$

Ignoring quadratic term  $X^T X$ , the same for all classes, we get a linear discriminant function:

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

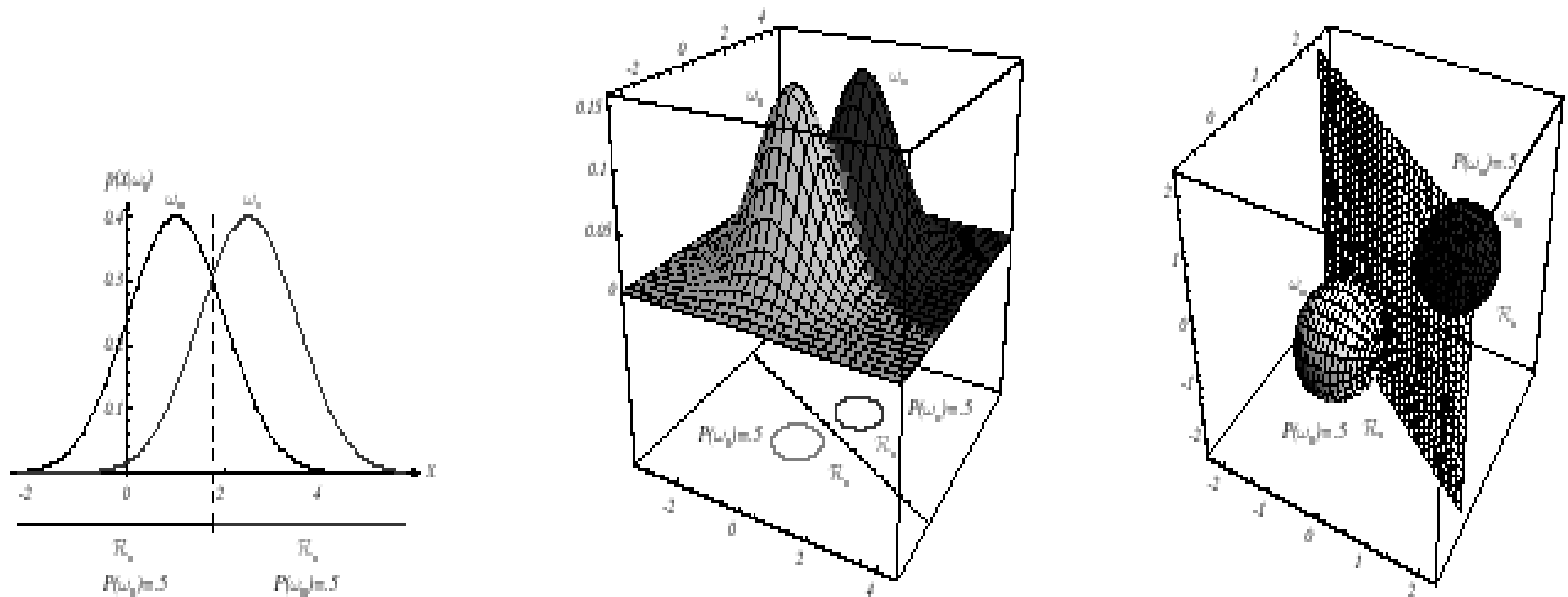
where :

$$w_i = \frac{\mathbf{m}_i}{s^2}; \quad w_{i0} = -\frac{1}{2s^2} \mathbf{m}_i^t \mathbf{m}_i + \ln P(w_i)$$

( $w_{i0}$  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  $R_i$  and  $R_j$

$$x_0 = \frac{1}{2}(\mathbf{m}_i + \mathbf{m}_j) - \frac{S^2}{\|\mathbf{m}_i - \mathbf{m}_j\|^2} \ln \frac{P(w_i)}{P(w_j)} (\mathbf{m}_i - \mathbf{m}_j)$$

always orthogonal to the line linking the means!

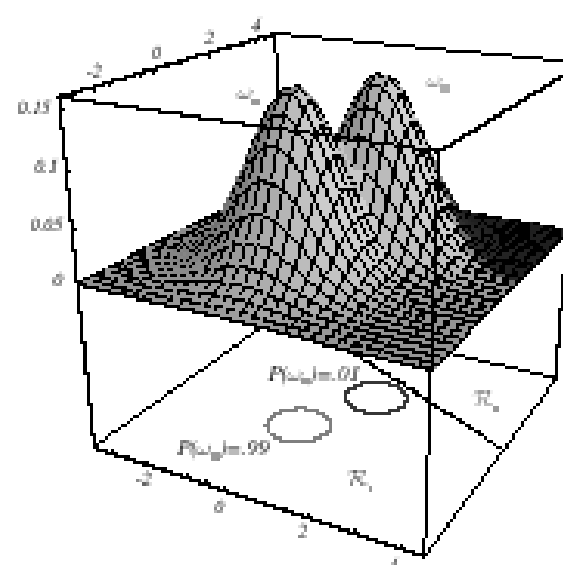
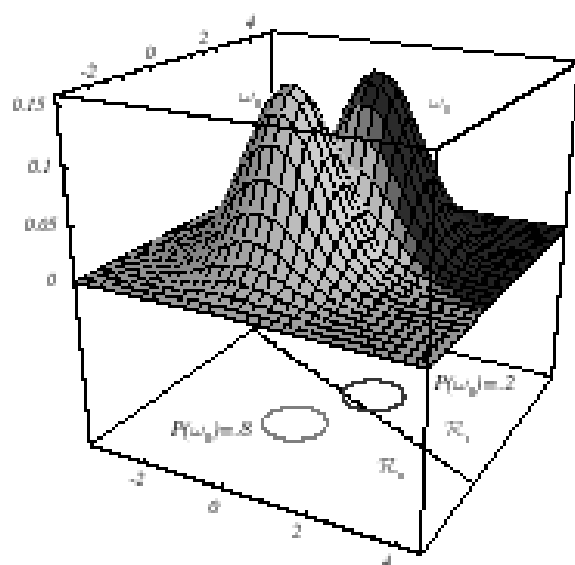
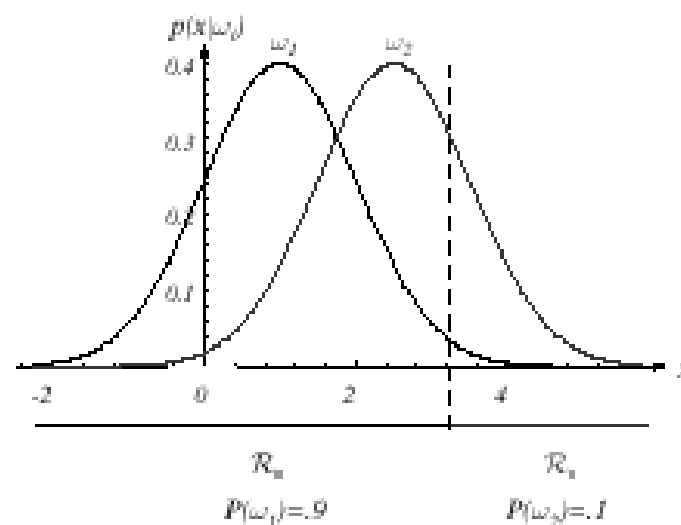
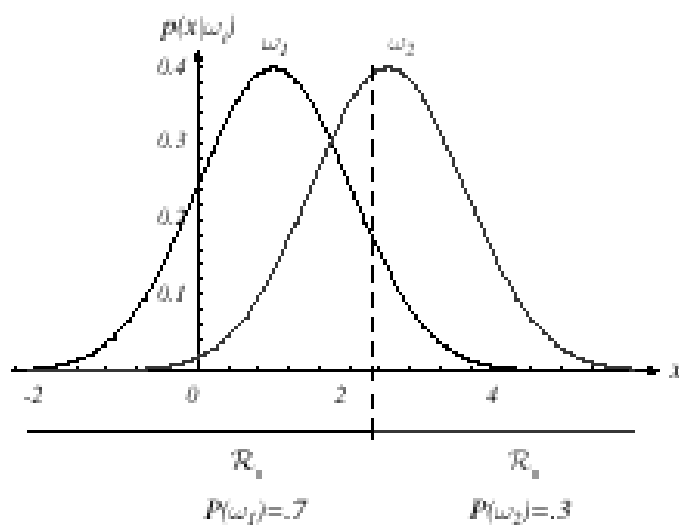
$$\text{if } P(w_i) = P(w_j) \text{ then } g_i(x) = (x - \mathbf{m}_i)^t (x - \mathbf{m}_i)$$

$$\text{and } x_0 = \frac{1}{2}(\mathbf{m}_i + \mathbf{m}_j)$$

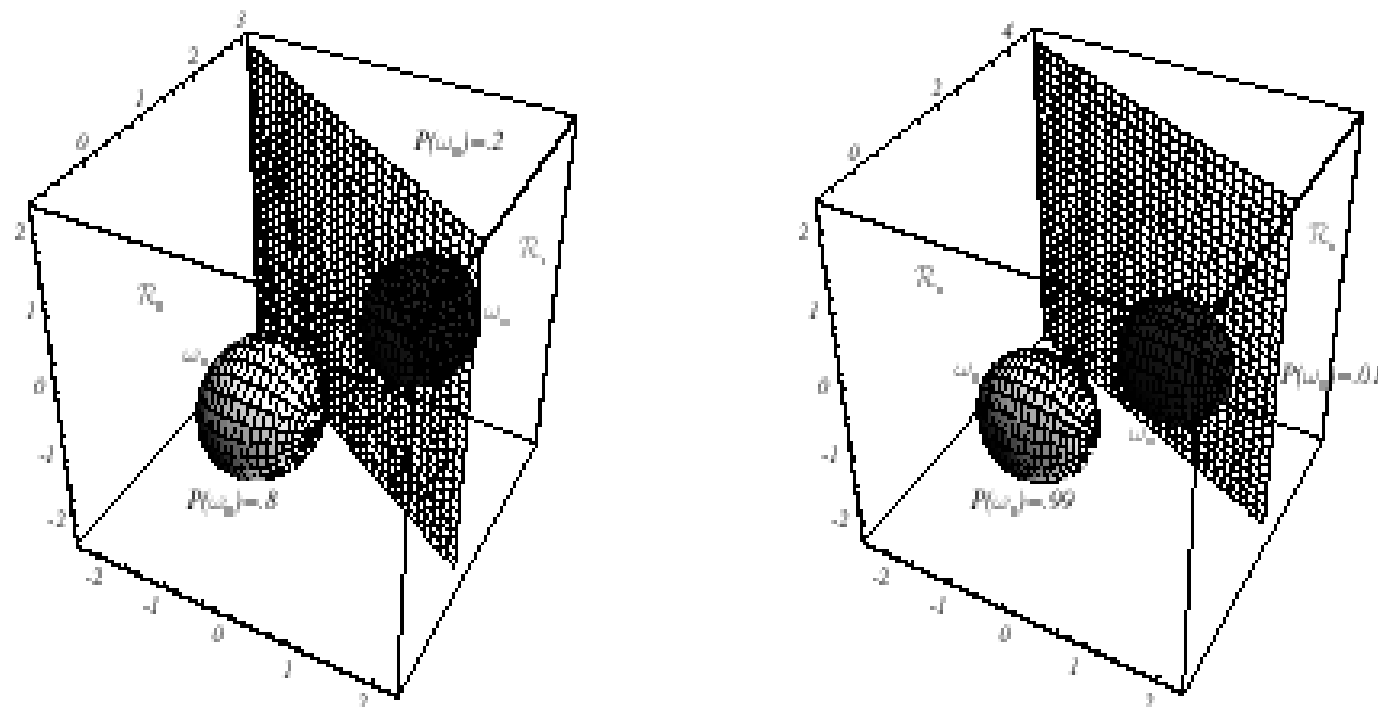
If prior probabilities are the same for all classes, then the input pattern sample should be classified into the class minimizing the Euclidean distance to its mean.

### Minimum Distance Classifier :

- adopts a template-matching approach. The mean  $\mu_k$  for each class  $C_k$  is assigned during training.
- For each new pattern sample, extract the feature vector and compute its Euclidean distance to class mean. Then, classify sample into the class minimizing this distance.







**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.

## Second Case

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

$$g_i(x) = w_i^t x + w_{i0} \text{ (linear)}$$

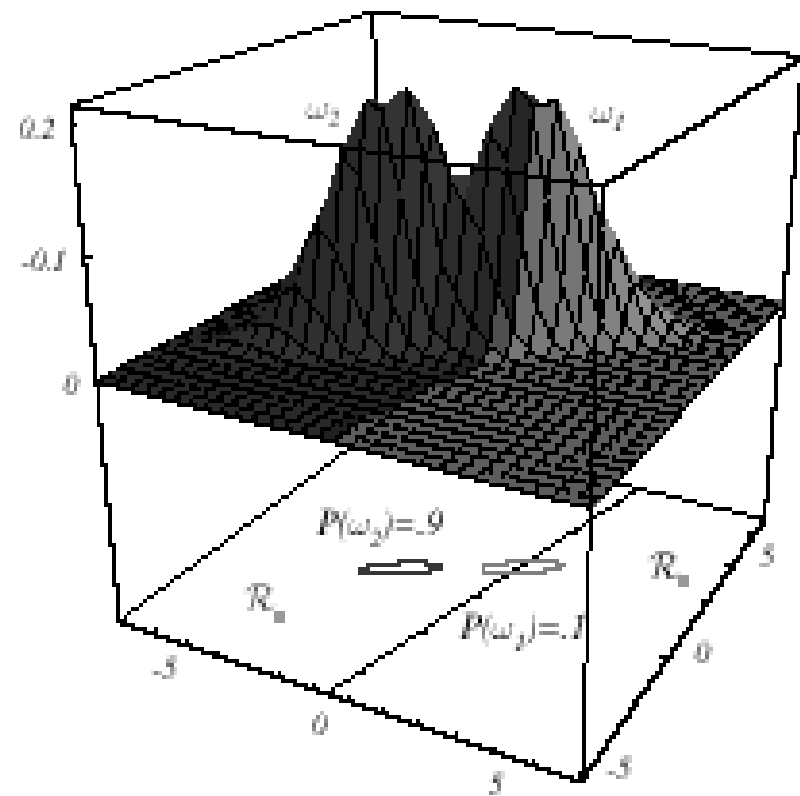
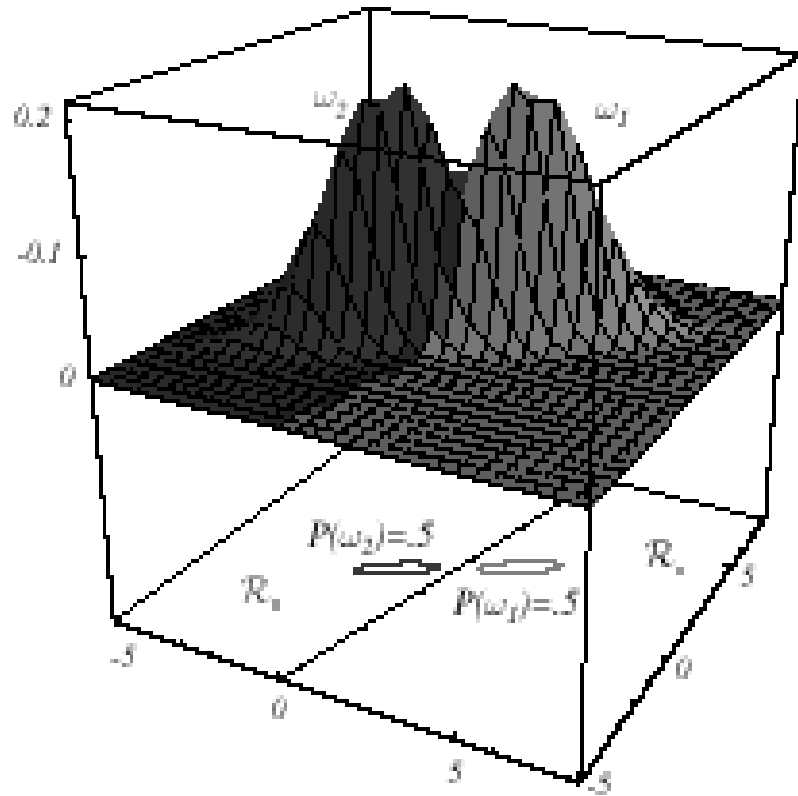
where :

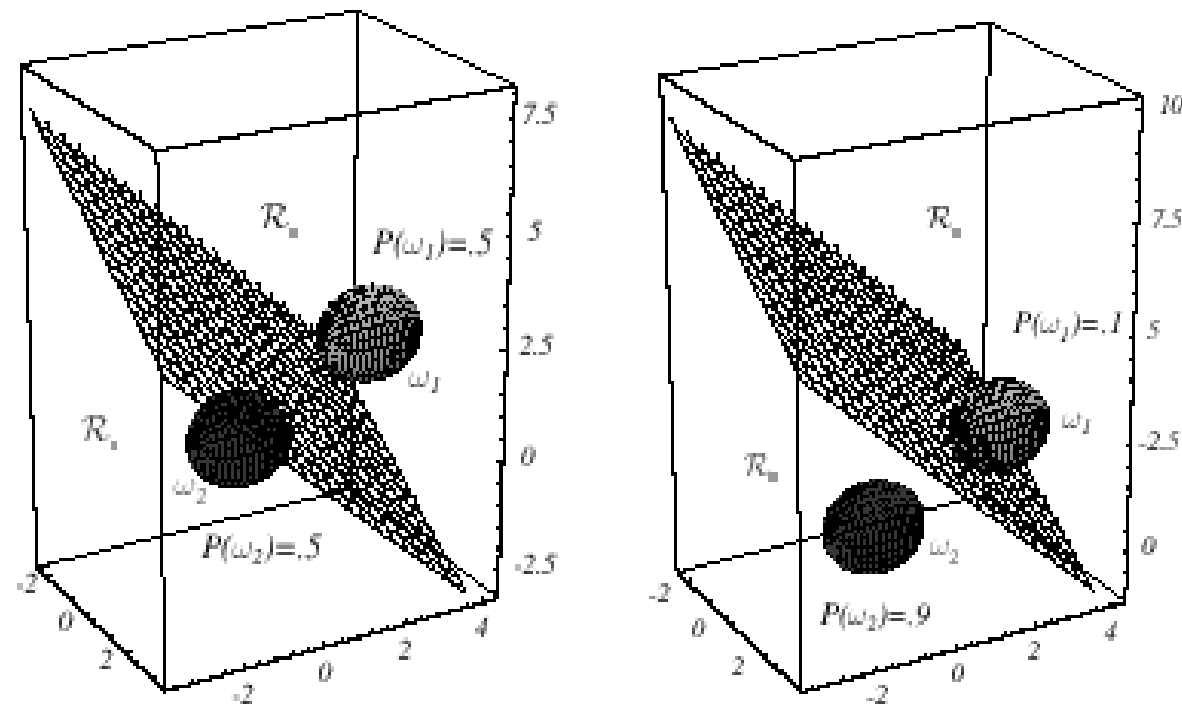
$$w_i = \Sigma^{-1} m_i; \quad w_{i0} = -\frac{1}{2} m_i^t \Sigma^{-1} m_i + \ln P(w_i)$$

- Hyperplane separating  $R_i$  and  $R_j$

$$x_0 = \frac{1}{2}(m_i + m_j) - \frac{\ln[P(w_i)/P(w_j)]}{(m_i - m_j)^t \Sigma^{-1} (m_i - m_j)} \cdot (m_i - m_j)$$

(the hyperplane separating  $R_i$  and  $R_j$  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 G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- Case  $\Sigma_i = \text{arbitrary}$ 
  - The covariance matrices are different for each category

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

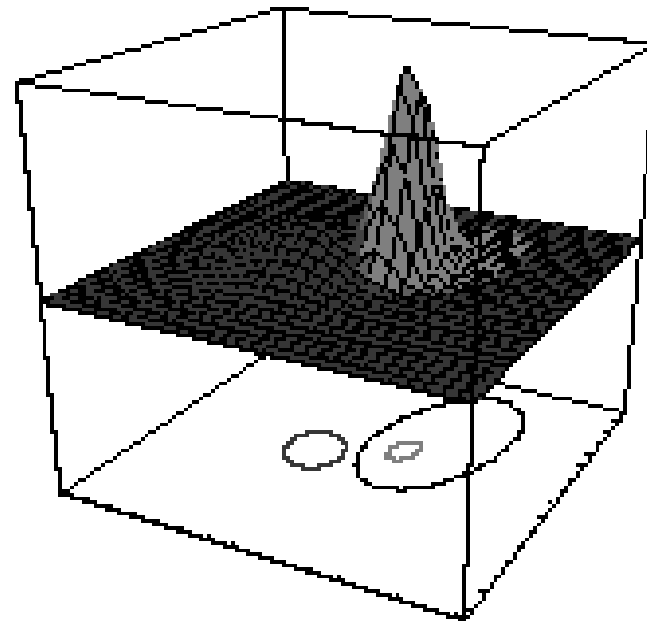
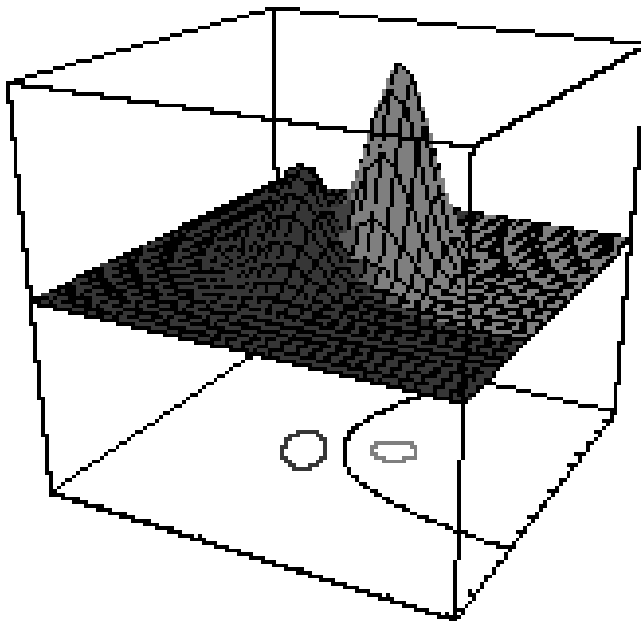
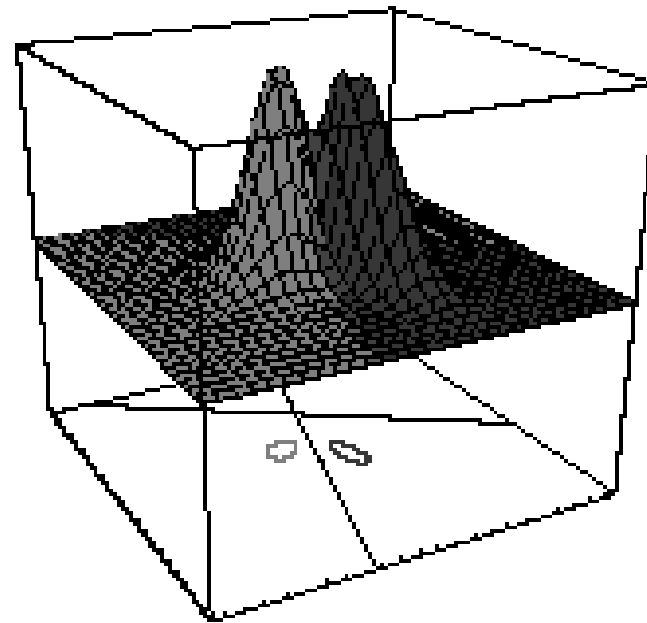
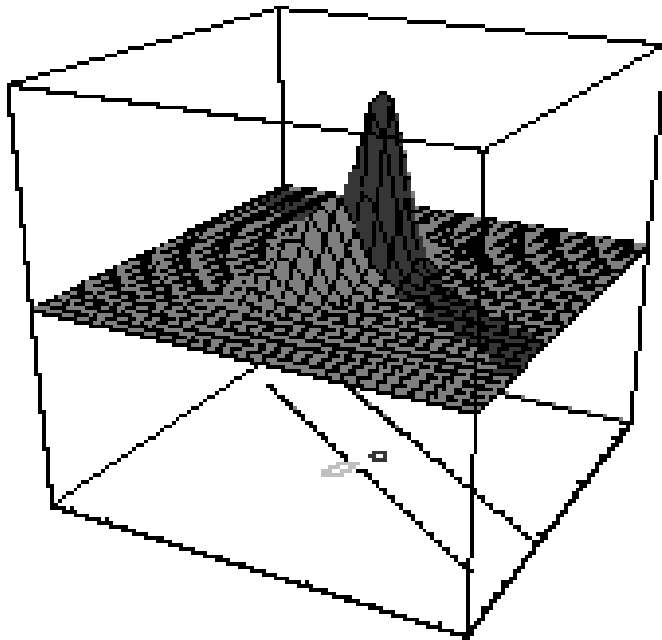
where :

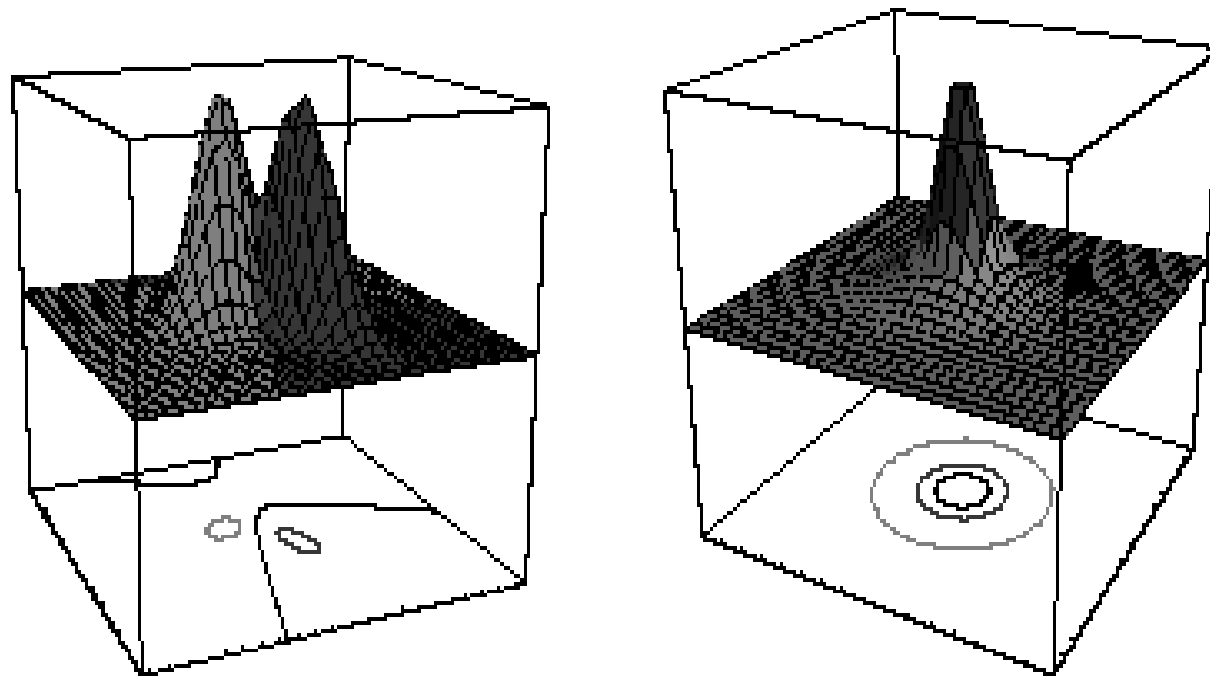
$$W_i = -\frac{1}{2} \Sigma_i^{-1}$$

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

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

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





**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.

## Decision Regions Case $\Sigma_i = \text{arbitrary}$

$$g(x) = g_i(x) - g_j(x) = 0$$

$$g(x) = x^t A x + b^t x + c = 0$$

$$A = \Sigma_j^{-1} - \Sigma_i^{-1}$$

$$b = -2 \Sigma_j^{-1} m_j + 2 \Sigma_i^{-1} m_i$$

$$c = 2 \ln \frac{P(w_i)}{P(w_j)} + \ln \frac{|\Sigma_j|}{|\Sigma_i|} + m_j^t \Sigma_j^{-1} m_j - m_i^t \Sigma_i^{-1} m_i$$



## Example: 2 features, 2 categories

$$m_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}, \quad \Sigma_1 = \begin{bmatrix} 1/2 & 0 \\ 0 & 2 \end{bmatrix}, \quad m_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix} \quad \text{and} \quad \Sigma_2 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\Sigma_1^{-1} = \begin{bmatrix} 2 & 0 \\ 0 & 1/2 \end{bmatrix}, \quad \Sigma_2^{-1} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix}, \quad P(w_1) = P(w_2)$$

$$x^t A x + b^t x + c = 0$$

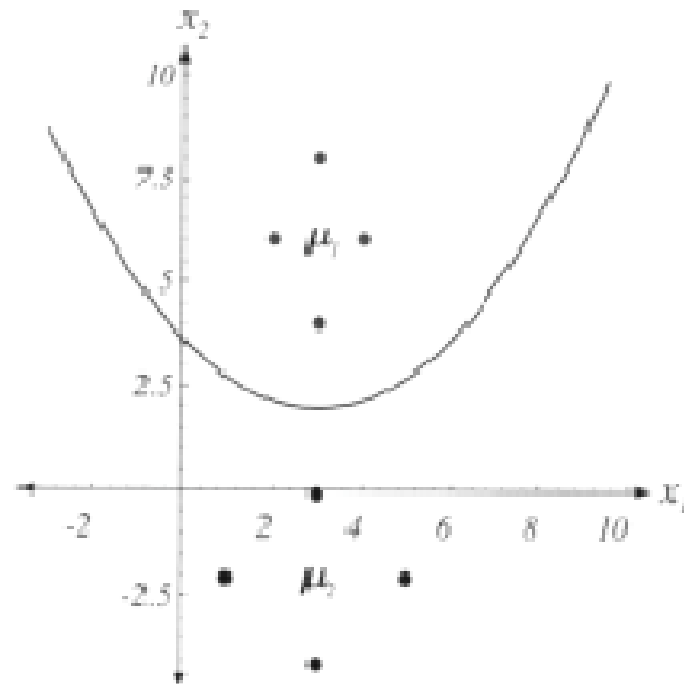
$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} -3/2 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 2 \\ 0 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} 3 \\ 6 \end{bmatrix} - 2 \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} 3 \\ -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$+ \ln(4) + \begin{bmatrix} 3 & -2 \end{bmatrix} \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} 3 \\ -2 \end{bmatrix} - \begin{bmatrix} 3 & 6 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} 3 \\ 6 \end{bmatrix} = 0$$

$$-\frac{3}{2}x_1^2 + 9x_1 + 8x_2 - 28.1137 = 0$$

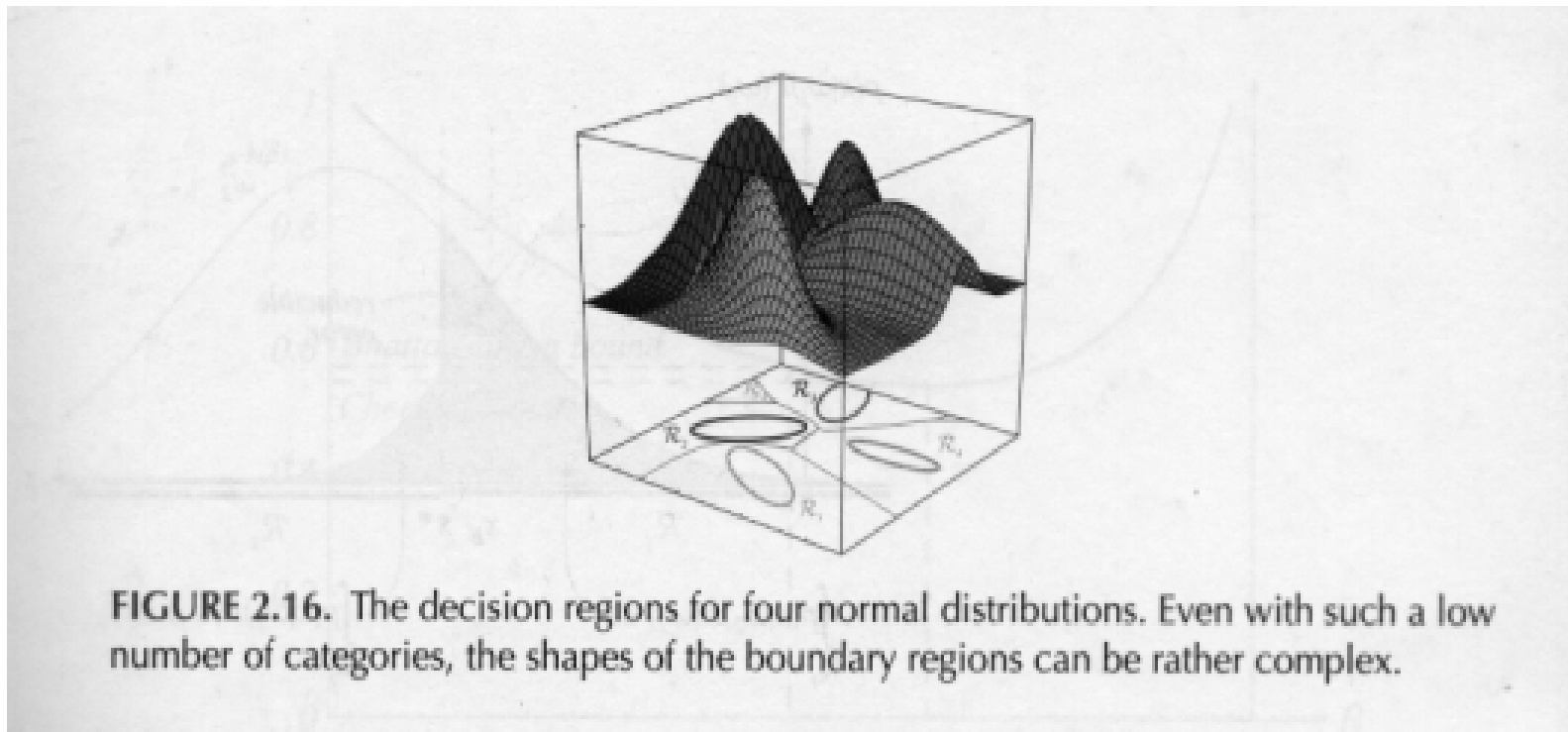
$$x_2 = 3.514 - 1.125x_1 + 0.1875x_1^2$$

# Example (cont.)



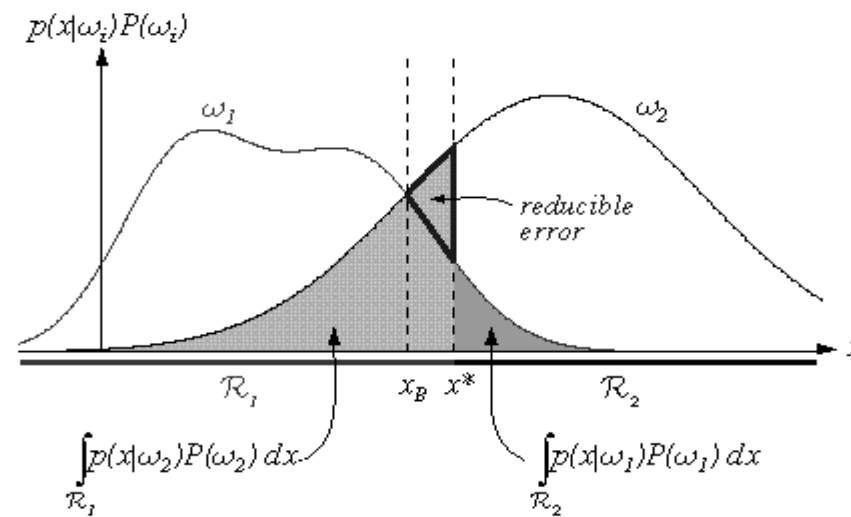
# Discriminant functions for a multiple category system

Between two categories the boundary decision is a hyperquadric



# Error Probabilities

$P(\text{error}) = \text{Error total} = \text{Sum of all error probabilities}$



**FIGURE 2.17.** Components of the probability of error for equal priors and (nonoptimal) decision point  $x^*$ . The pink area corresponds to the probability of errors for deciding  $\omega_1$  when the state of nature is in fact  $\omega_2$ ; the gray area represents the converse, as given in Eq. 70. If the decision boundary is instead at the point of equal posterior probabilities,  $x_B$ , then this reducible error is eliminated and the total shaded area is the minimum possible; this is the Bayes decision and gives the Bayes error rate. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

# Error Bounds

In most cases, it is difficult to compute the integrals for complex forms

What is the error bound?

Chernoff bound:  $\min[a,b] \leq a^\beta b^{1-\beta}$ ,  $0 \leq \beta \leq 1$ ,  $a, b \geq 0$

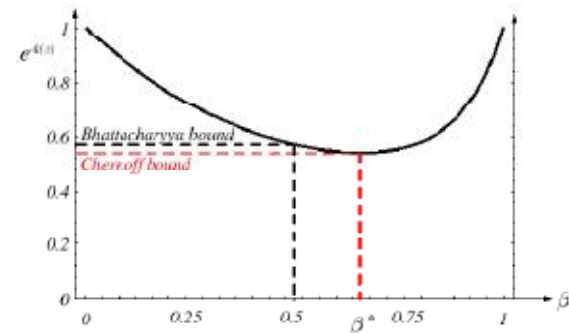
$$\begin{aligned} P(\text{error}) &= \min(P(x|\omega_1)P(\omega_1), P(x|\omega_2)P(\omega_2)) \\ &\leq \int P^\beta(\omega_1)P^{1-\beta}(\omega_2) P^\beta(x|\omega_1)P^{1-\beta}(x|\omega_2)dx \\ &= P^\beta(\omega_1)P^{1-\beta}(\omega_2)e^{-K(\beta)} \text{ for normal distributions} \end{aligned}$$

$K(\beta)$ : Equation 75

Bhattacharyya Bound

$$\beta = 1/2$$

for the Gaussian case:



$$K\left(\frac{1}{2}\right) = \frac{1}{8}(\mathbf{m}_2 - \mathbf{m}_1)^t \left[ \frac{\Sigma_1 + \Sigma_2}{2} \right]^{-1} (\mathbf{m}_2 - \mathbf{m}_1) + \frac{1}{2} \ln \frac{\left| \frac{\Sigma_1 + \Sigma_2}{2} \right|}{\sqrt{|\Sigma_1| |\Sigma_2|}}$$

# Bayes Decision Theory – Discrete Features

- Components of  $\mathbf{x}$  are binary or integer valued,  $\mathbf{x}$  can take only one of  $m$  discrete values

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

- Case of independent binary features in 2 category problem

$$P(\mathbf{x} | \omega_j) = \prod_{i=1}^d P(x_i | \omega_j)$$

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

$$p_i = P(x_i = 1 | w_1)$$

$$q_i = P(x_i = 1 | w_2)$$

# Case of independent binary features in 2 categories

$$P(\mathbf{x} | \omega_1) = \prod_{i=1}^d p_i^{x_i} (1 - p_i)^{1-x_i}$$

$$P(\mathbf{x} | \omega_2) = \prod_{i=1}^d q_i^{x_i} (1 - q_i)^{1-x_i}$$

$$p_i = P(x_i = 1 | \omega_1) \text{ et } q_i = P(x_i = 1 | \omega_2)$$

$$g(\mathbf{x}) = P(\omega_1 | \mathbf{x}) - P(\omega_2 | \mathbf{x}) \Rightarrow g(\mathbf{x}) = \ln \frac{P(\mathbf{x} | \omega_1)}{P(\mathbf{x} | \omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

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

$$g(\mathbf{x}) = \sum_{i=1}^d \left[ x_i \ln \frac{p_i}{q_i} + (1 - x_i) \ln \frac{1 - p_i}{1 - q_i} \right] + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

- The discriminant function in this case is:

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

where :

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

and :

$$w_0 = \sum_{i=1}^d \ln \frac{1 - p_i}{1 - q_i} + \ln \frac{P(w_1)}{P(w_2)}$$

decide  $w_1$  if  $g(x) > 0$  and  $w_2$  if  $g(x) \leq 0$

§ If  $p_i = q_i$ ,  $g(x)$  depends only on  $p(\omega_1)$  and  $p(\omega_2)$

§ If  $p_i > q_i$ ,  $1 - p_i < 1 - q_i$  so  $w_i > 0$ .  $x_i = 1$  results in favor of  $\omega_1$ .