

# Pattern Recognition (Αναγνώριση Προτύπων)

## Bayesian Decision Theory (Μπεϋζιανή Θεωρία Αποφάσεων)

Panos Trahanias

UNIVERSITY OF CRETE  
DEPARTMENT of COMPUTER  
SCIENCE



# *Bayes Decision theory*

- Statistically optimal classification.
- Based on the probabilistic description of the classification problem/task.
- It assumes:
  - Classification problem can be probabilistically stated
  - Relevant values and probability functions are known (not valid in practice).

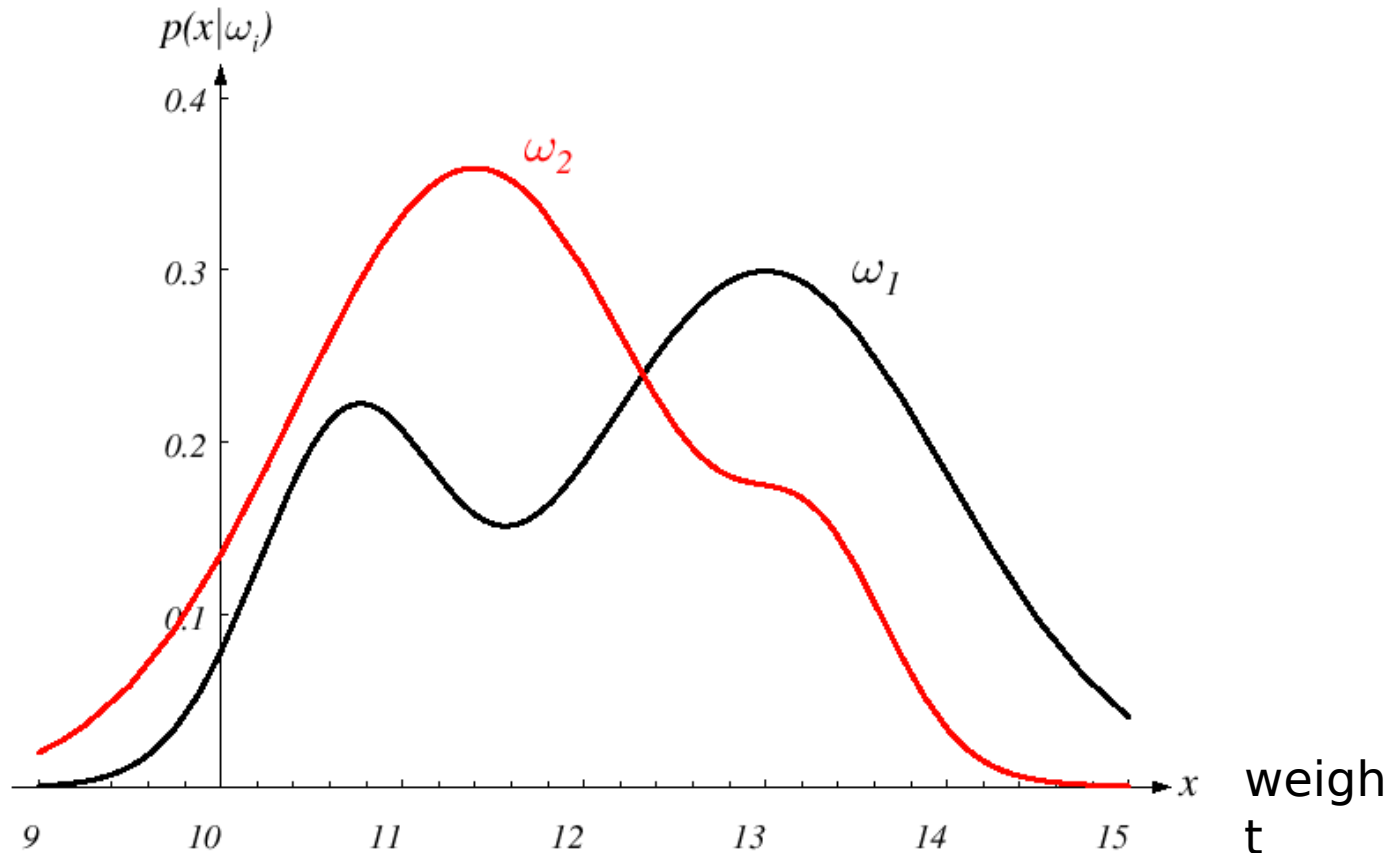


# Class Conditional Probabilities

$\omega_1$ : Sea bass

$\omega_2$ : Salmon

- $p(x | \omega_2)$ : Conditional Probability Density Function (PDF) of the variable  $x$  given the state of nature.
- Likelihood: Given that salmon is observed, what is the probability that its weight is between 11 and 12?





# Definitions and Bayes Decision Rule

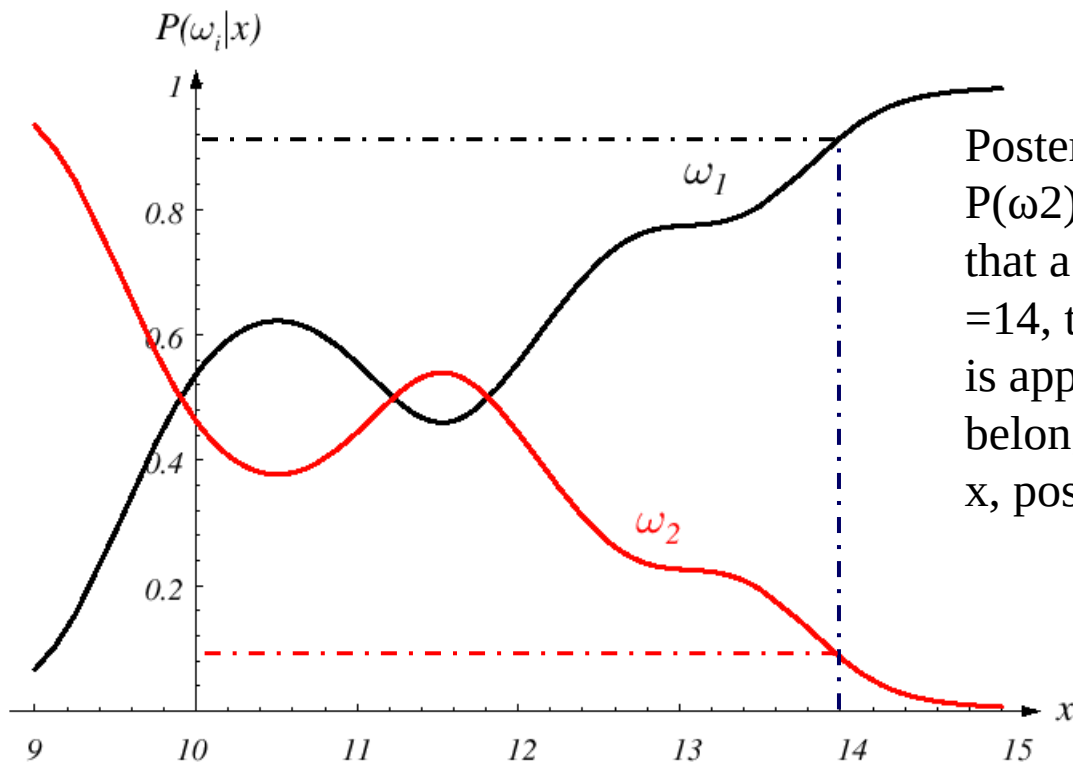
- State of nature (Κατάσταση της φύσης)
- Prior (Εκ των προτέρων πιθανότητα)
- Posterior (Εκ των υστέρων πιθανότητα)
- Likelihood (Πιθανοφάνεια)
- Evidence

$$\begin{aligned} P(\omega_j | x) &= \frac{p(x | \omega_j) \cdot P(\omega_j)}{\sum_{x \in X} p(x | \omega_j) \cdot P(\omega_j)} \\ &= \frac{p(x | \omega_j) \cdot P(\omega_j)}{p(x)} \end{aligned}$$



# Posterior Probabilities

Bayes rule facilitates estimation/computation of posterior probabilities (otherwise hard to compute), given prior probability, likelihood and evidence.



Posterior probabilities when  $P(\omega_1) = 2/3$  and  $P(\omega_2) = 1/3$ . Given e.g.

that a pattern is measured with feature value  $x = 14$ , the probability that it belongs to class  $\omega_2$  is approx. 0.08, whereas the probability that it belongs to  $\omega_1$  is 0.92.

For each  $x$ , posterior probabilities sum up to 1.0.



## Selection of the class that has the Highest posterior probability!!!

Choose  $\omega_i$  if  $P(\omega_i | x) > P(\omega_j | x)$  for all  $i = 1, 2, \dots, c$

$$P(\text{error}) = \min [ P(\omega_1 | x), P(\omega_2 | x), \dots, P(\omega_c | x) ]$$

In case of multiple features,  $\mathbf{x} = \{x_1, x_2, \dots, x_d\}$  then

Choose  $\omega_i$  if  $P(\omega_i | \mathbf{x}) > P(\omega_j | \mathbf{x})$  for all  $i = 1, 2, \dots, c$

$$P(\text{error}) = \min [ P(\omega_1 | \mathbf{x}), P(\omega_2 | \mathbf{x}), \dots, P(\omega_c | \mathbf{x}) ]$$



# The Loss Function

- Mathematical description of the cost of each choice.
  - Are some choices more “expensive” than others?
- ✓  $\{\omega_1, \omega_2, \dots, \omega_c\}$ : Set of physical states (classes)
  - ✓  $\mathbf{x} = [x_1, \dots, x_d]^T$ : Feature vector
  - ✓  $\{\alpha_1, \alpha_2, \dots, \alpha_a\}$ : Set of actions. Note that ‘a’ does not have to be the same as ‘c’, as we can perform more or less actions than the number of classes. For example, rejection is also a possible action.
  - ✓  $\lambda(\alpha_i | \omega_j)$ : Cost (κόστος) of the action  $\alpha_i$  when the real class is  $\omega_j$ .
  - ✓  $R(\alpha_i | \mathbf{x})$ : conditional risk – Expected loss for action  $\alpha_i$ .

$$R(\alpha_i | \mathbf{x}) = \sum_{j=1}^c \lambda(\alpha_i | \omega_j) \cdot P(\omega_j | \mathbf{x})$$

Bayes decision chooses the action that  
minimizes the conditional risk!



# *Bayes decision based on the conditional risk*

1. Calculation of conditional risk  $R(\alpha_i | \mathbf{x})$  for each action.
2. Selection of the action with the lowest conditional risk.  
Suppose that it is action  $k$

3. The overall risk is:

$$R = \int_{\mathbf{x} \in X} R(\alpha_k | \mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x}$$

4. This is the Bayesian risk, the lowest possible risk that any classifier can have!
5. E.g. classification into one of two classes:

$$\frac{p(x / \omega_1)}{p(x / \omega_2)} \stackrel{\omega_1}{\geq} \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \frac{P(\omega_2)}{P(\omega_1)}$$





# Minimum Error Rate Classification

- If as action  $\alpha_i$  we choose the classification in the class  $\omega_j$ , and if all the costs of incorrect classification are equal to one, we have the so-called symmetric or 0-1 selection:

$$\lambda(\alpha_i | \omega_j) = \begin{cases} 0, & \text{if } i = j \\ 1, & \text{if } i \neq j \end{cases}$$

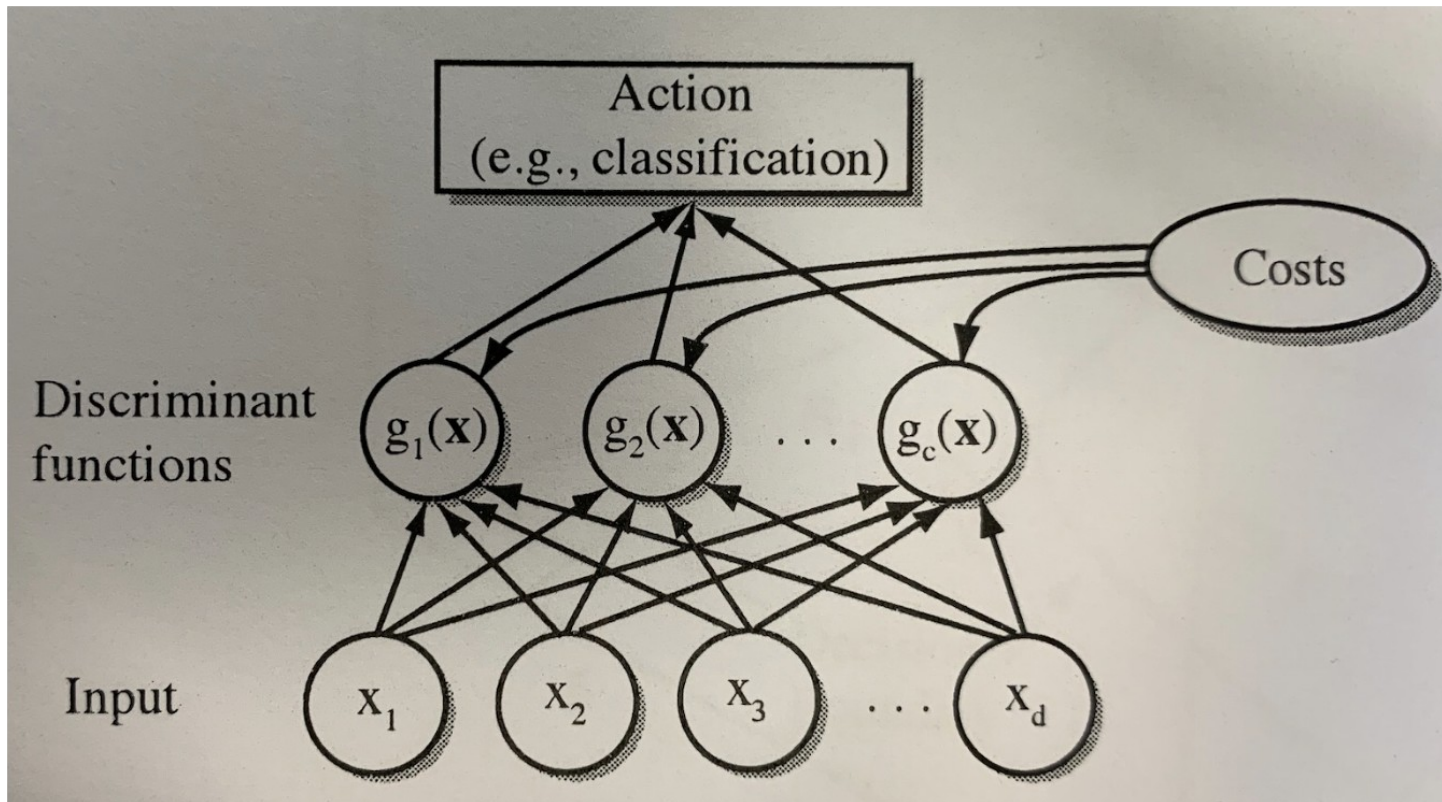
- This cost function stipulates zero loss for correct classification, and unit loss for incorrect classification. The corresponding conditional risk corresponding to this cost function is:

$$R(\alpha_i | \mathbf{x}) = \sum_{\substack{j \neq i \\ j=1, \dots, c}} P(\omega_j | \mathbf{x}) = 1 - P(\omega_i | \mathbf{x})$$

which is exactly the probability of error. Obviously, to minimize the risk, we should select the class that maximizes the posterior probability!!!



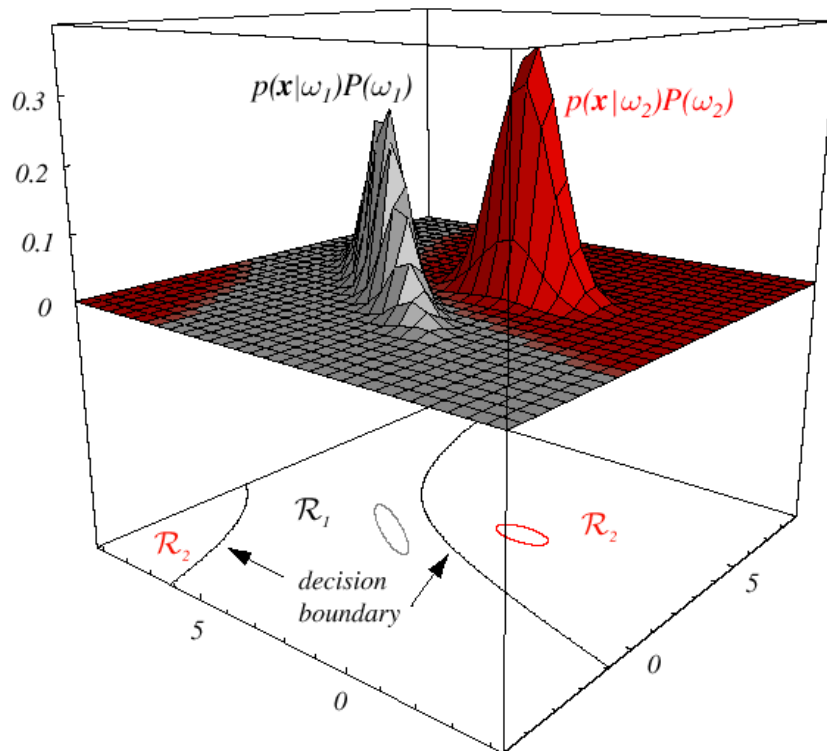
# *Discriminant Based Classification*





# Discriminant Based Classification

- The discrimination function  $g(\mathbf{x})$ , separates the classes from each other. This function corresponds the input vector to a class according to the definition: Select class  $i$  if :  
$$g_i(\mathbf{x}) > g_j(\mathbf{x}) \quad \forall i \neq j, \quad i, j = 1, 2, \dots, c$$
- The Bayes rule can be implemented in the form of discrimination functions  
$$g_i(\mathbf{x}) = P(\omega_i | \mathbf{x})$$



Any discrimination function creates  $c$  decision areas,  $R_1, \dots, R_c$ , which are separated by **decision surfaces** (επιφάνειες απόφασης).

Decision areas are not required to be continuous.

Decision surfaces satisfy the:

$$g_i(\mathbf{x}) = g_j(\mathbf{x})$$



## *Discriminant Based Classification*

$$g_i(\mathbf{x}) = P(\omega_i | \mathbf{x})$$

$$g_i(\mathbf{x}) = P(\omega_i | \mathbf{x}) = \frac{p(\mathbf{x} | \omega_i) P(\omega_i)}{\sum_{j=1}^c p(\mathbf{x} | \omega_j) P(\omega_j)}$$

$$g_i(\mathbf{x}) = p(\mathbf{x} | \omega_i) P(\omega_i)$$

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



# Gaussian PDF

- If the likelihood functions follow the multidimensional Gaussian, then the discrimination function takes the form

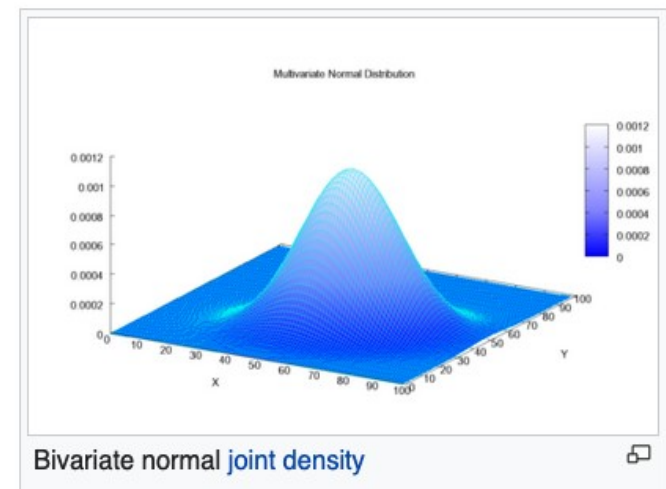
$$g_i(\mathbf{x}) = -\frac{1}{2} \left[ (\mathbf{x} - \mu_i)^T \cdot \Sigma_i^{-1} \cdot (\mathbf{x} - \mu_i) \right] - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

- There are 3 cases depending on the covariance matrix

$$p(\mathbf{x}) \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}):$$

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

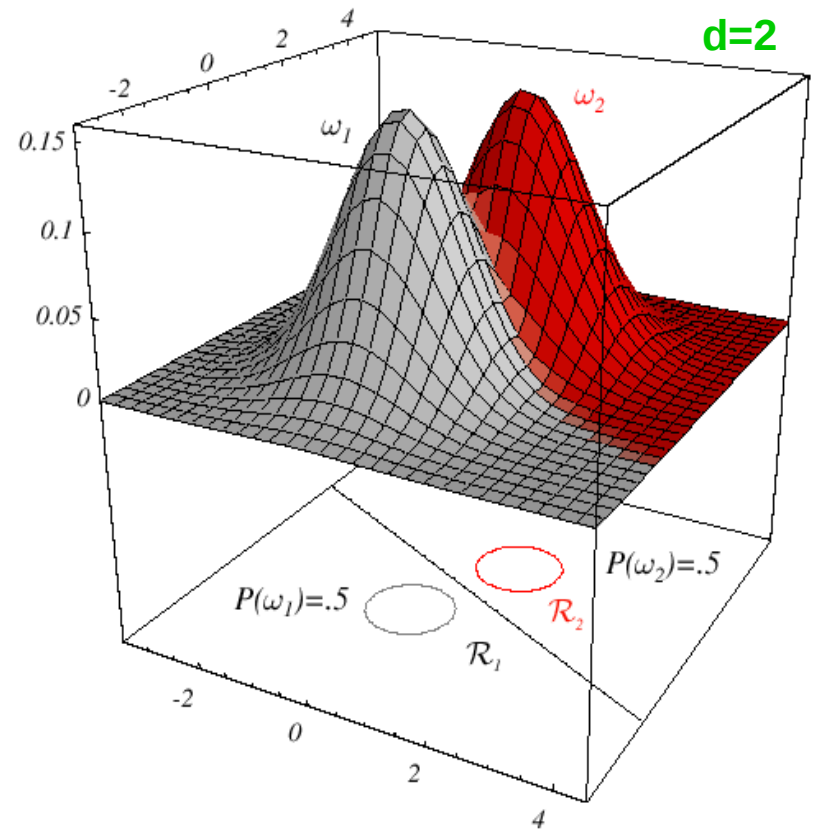
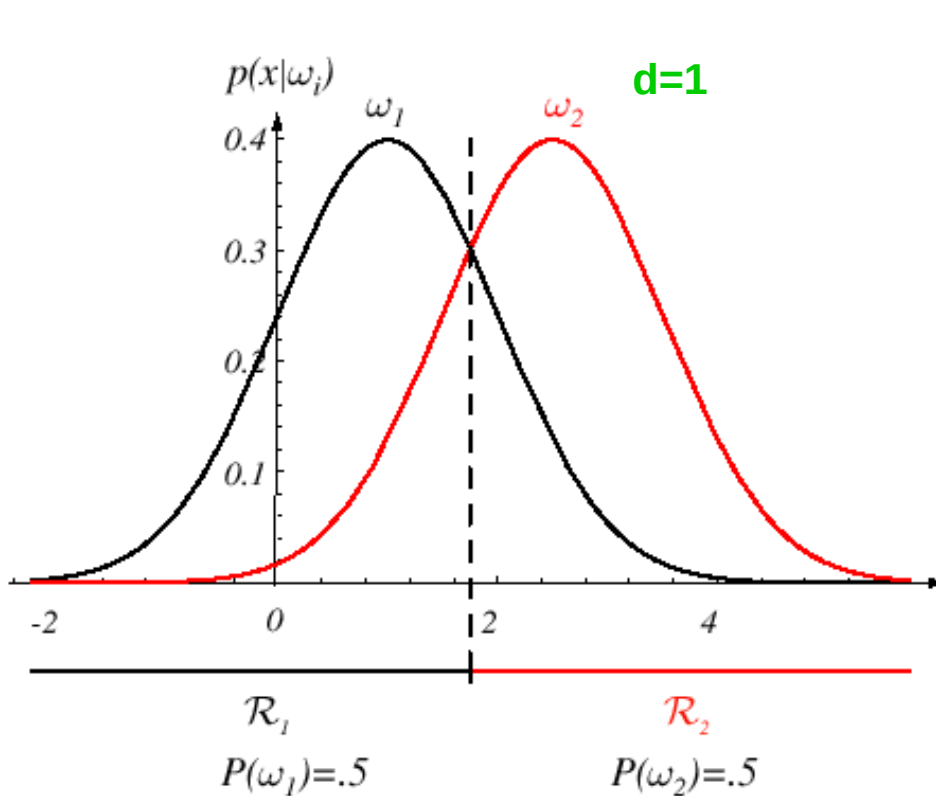
Multivariate





# Case 1: $\Sigma_i = \sigma^2 I$

- The features are statistically independent, and all have the same covariance: The samples are in **super-spheres** of equal size, and the decision surfaces are **super-planes** of dimension  $d-1$ .







# Case 1: $\Sigma_i = \sigma^2 I$

## Covariance Matrix

$$S = \begin{bmatrix} s_1^2 & s_{12} & s_{13} & \dots & s_{1p} \\ s_{21} & s_2^2 & s_{23} & \dots & s_{2p} \\ s_{31} & s_{32} & s_3^2 & \dots & s_{3p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ s_{p1} & s_{p2} & s_{p3} & \dots & s_p^2 \end{bmatrix}$$

where  $s_j^2 = \left(\frac{1}{n}\right) \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$  : variance of  $j^{\text{th}}$  variable

$$s_{jk} = \left(\frac{1}{n}\right) \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k) :$$

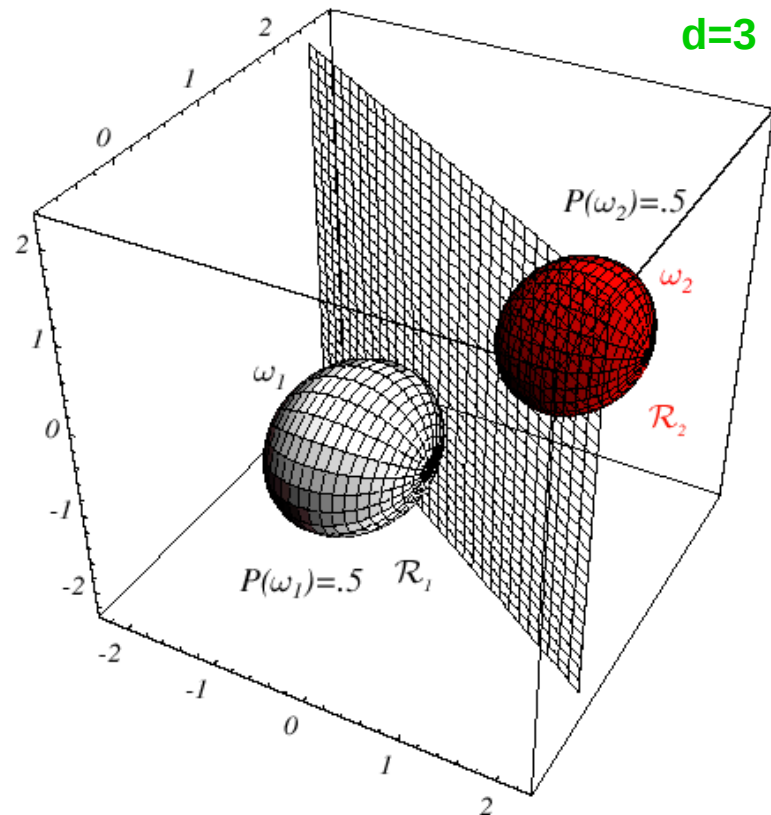
covariance b/w  
 $j^{\text{th}}$  &  $k^{\text{th}}$  variables

$$\bar{x}_j = \left(\frac{1}{n}\right) \sum_{i=1}^n x_{ij} : \text{mean of } j^{\text{th}} \text{ variable}$$



# Case 1: $\Sigma_i = \sigma^2 I$

- When  $d = 3$ , the samples are in **spheres** of equal size, and the decision surfaces are **flat**.





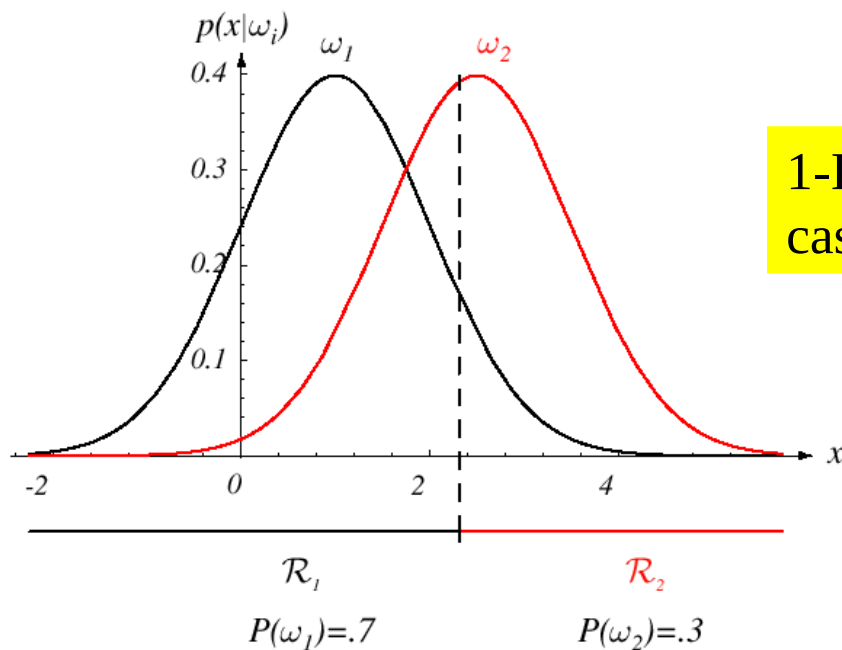


# Case 1: $\Sigma_i = \sigma^2 I$

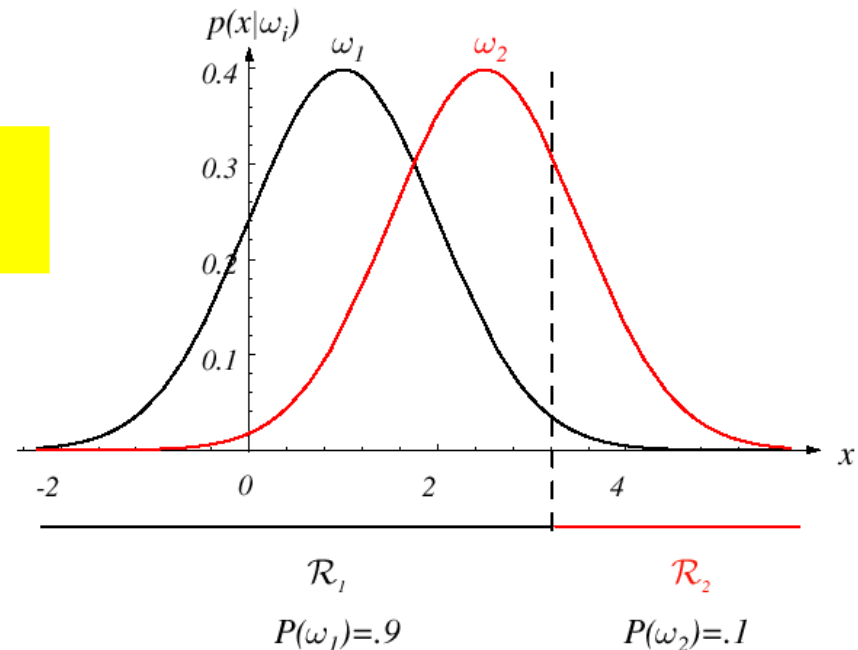
➤ This case creates linear discriminant functions:

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

$$\mathbf{w}_i = \frac{1}{\sigma^2} \boldsymbol{\mu}_i, \quad w_{i0} = \frac{-1}{2\sigma^2} \boldsymbol{\mu}_i^T \cdot \boldsymbol{\mu}_i + \ln P(\omega_i)$$



1-D  
case



**Notice that the prior probabilities move the threshold away from the most probable mean.**



## Case 1: $\Sigma_i = \sigma^2 I$

$$g_i(\mathbf{x}) = -\frac{1}{2} \left[ (\mathbf{x} - \mu_i)^T \cdot \Sigma_i^{-1} \cdot (\mathbf{x} - \mu_i) \right] - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

$$g_i(\mathbf{x}) = -\frac{1}{2\sigma^2} [\mathbf{x}^t \mathbf{x} - 2\mu_i^t \mathbf{x} + \mu_i^t \mu_i] + \ln P(\omega_i),$$

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

$$\mathbf{w}_i = \frac{1}{\sigma^2} \mu_i$$

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

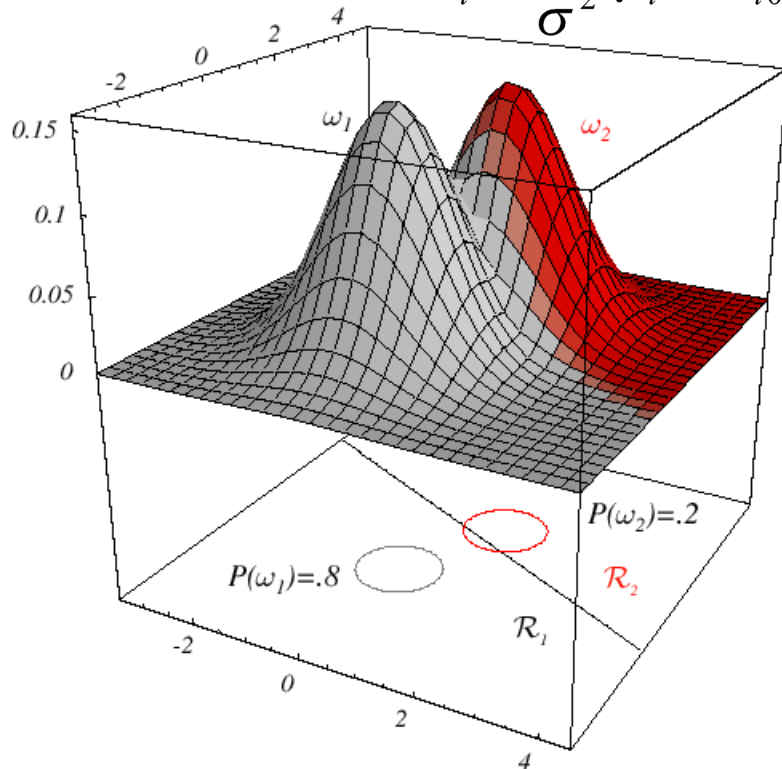


# Case 1: $\Sigma_i = \sigma^2 I$

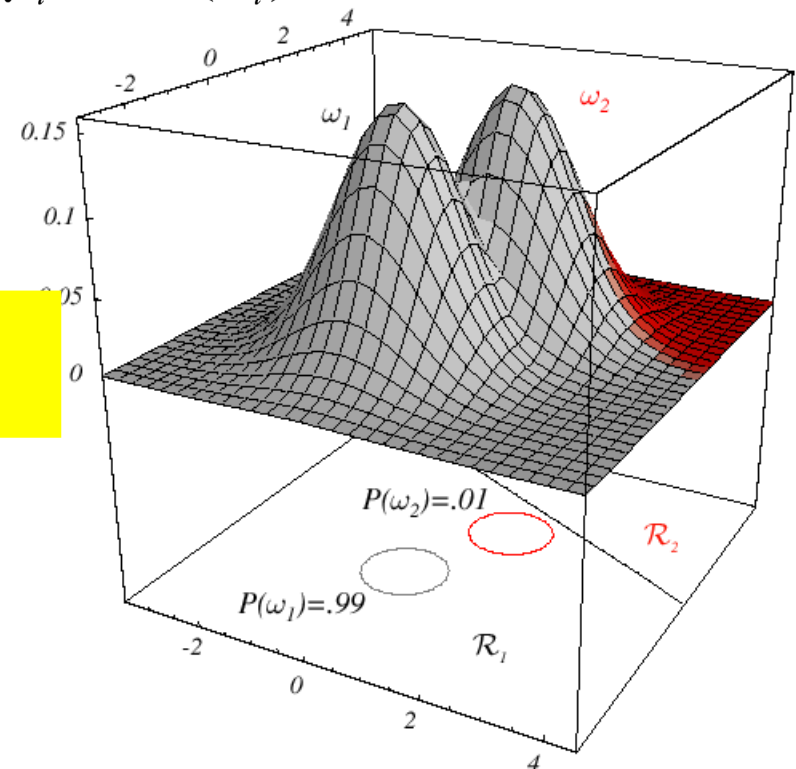
➤ This case creates linear discrimination functions:

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

$$\mathbf{w}_i = \frac{1}{\sigma^2} \boldsymbol{\mu}_i, \quad w_{i0} = \frac{-1}{2\sigma^2} \boldsymbol{\mu}_i^T \cdot \boldsymbol{\mu}_i + \ln P(\omega_i)$$



2-D  
case



Notice that the prior probabilities move the decision line away from the most probable mean.

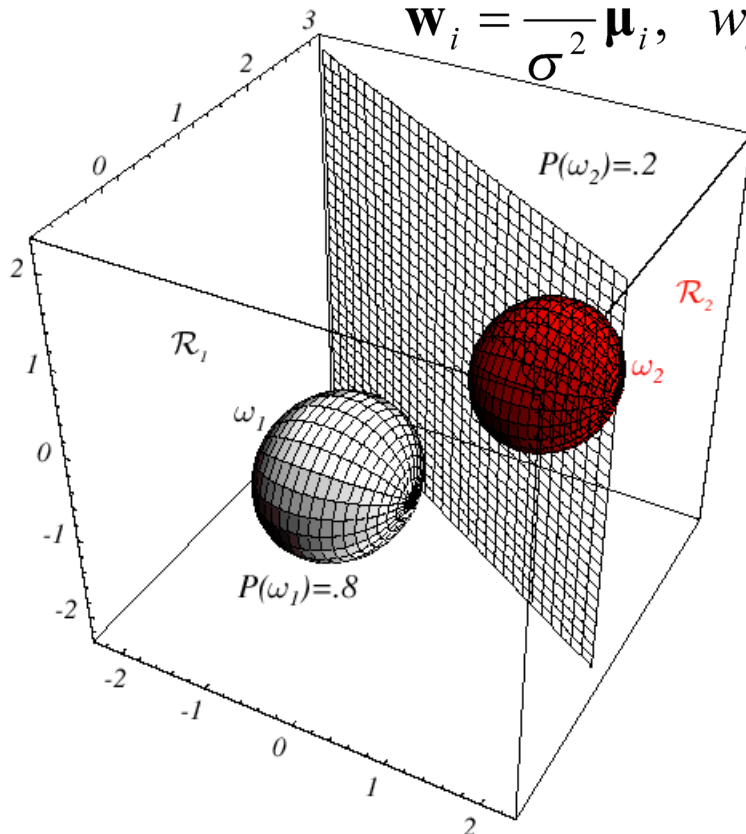


# Case 1: $\Sigma_i = \sigma^2 I$

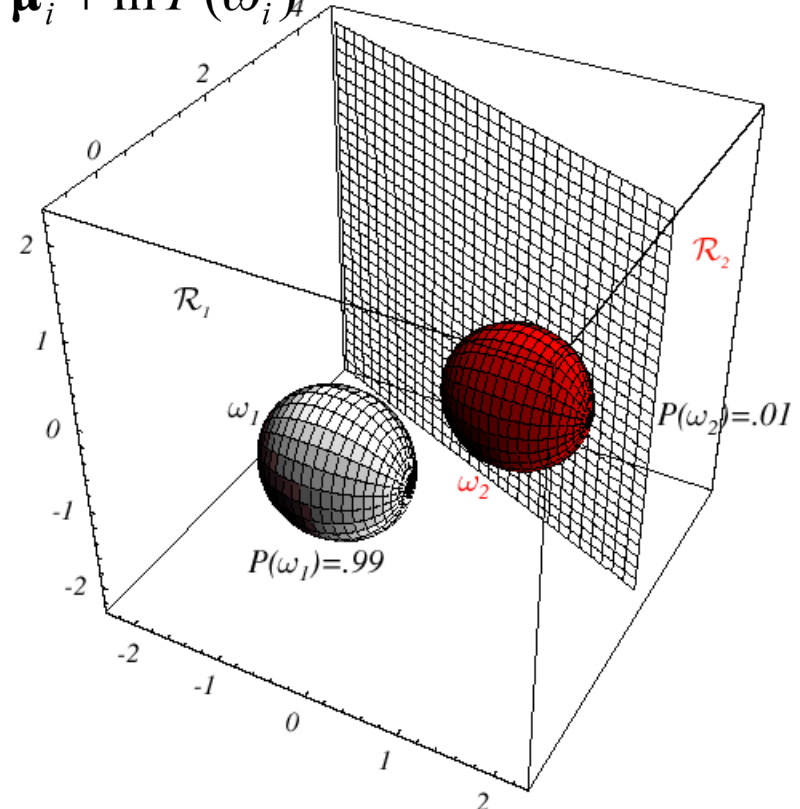
➤ This case creates linear discrimination functions:

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

$$\mathbf{w}_i = \frac{1}{\sigma^2} \boldsymbol{\mu}_i, \quad w_{i0} = \frac{-1}{2\sigma^2} \boldsymbol{\mu}_i^T \cdot \boldsymbol{\mu}_i + \ln P(\omega_i)$$



3-D  
case



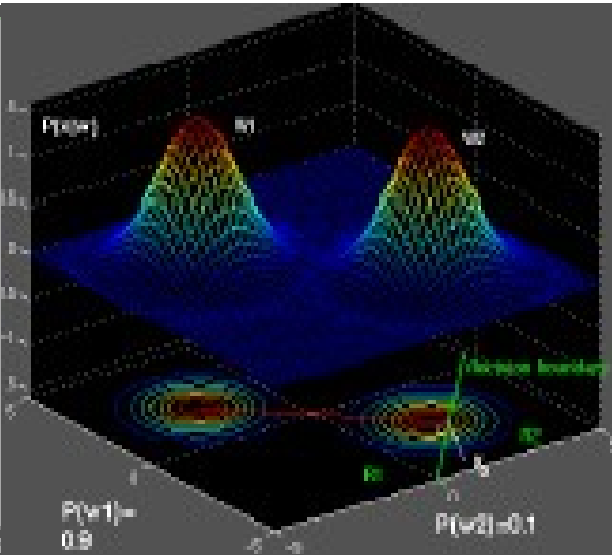
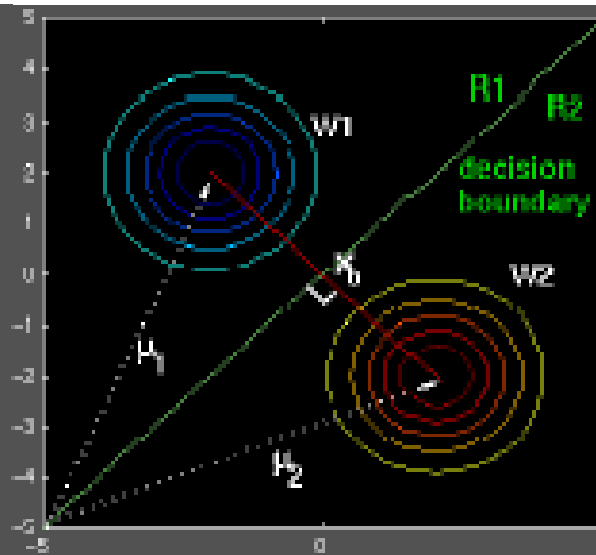
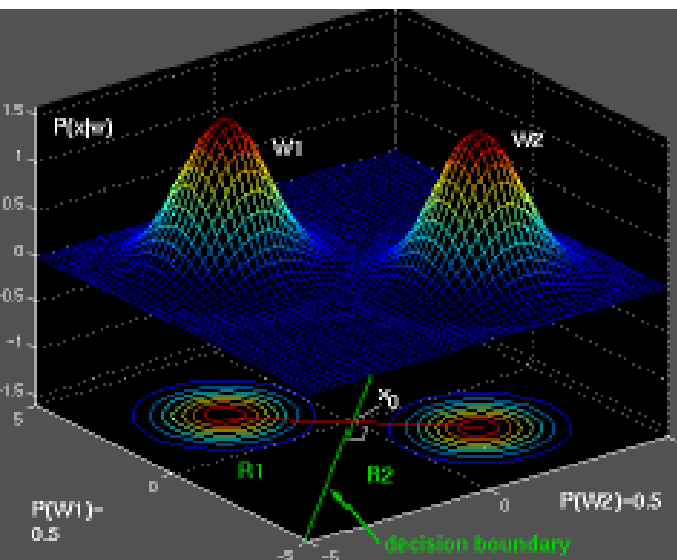
Notice that the prior probabilities move the decision plane away from the most probable mean.



# Case 1: $\Sigma_i = \sigma^2 I$

The decision surfaces are hyper-planes defined by the linear equations  $g_i(\mathbf{x}) = g_j(\mathbf{x})$ , written as  $\mathbf{w}^T(\mathbf{x} - \mathbf{x}_0) = 0$  where:

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



Decision Surface: Hyper-plane that passes through the point  $\mathbf{x}_0$  and is perpendicular to the vector  $\mathbf{w}$  that connects the mean values  $\boldsymbol{\mu}_i$  και  $\boldsymbol{\mu}_j$ .



Case 1:  $\Sigma_i = \sigma^2 I$

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

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

$$\mathbf{w} = \mu_i - \mu_j$$

$$\mathbf{x}_0 = \frac{1}{2}(\mu_i + \mu_j) - \frac{\sigma^2}{\|\mu_i - \mu_j\|^2} \ln \frac{P(\omega_i)}{P(\omega_j)} (\mu_i - \mu_j).$$



## Case 2: $\Sigma_i = \Sigma$

➤ The covariance matrices are arbitrary, but the same for all classes.

The features create hyper-elliptical groups of the same size and shape with centers  $\mu_i$ .

➤ Linear decision functions → Hyper-planes as decision surfaces

$$g_i(\mathbf{x}) = -\frac{1}{2} \left[ (\mathbf{x} - \mu_i)^T \cdot \Sigma_i^{-1} \cdot (\mathbf{x} - \mu_i) \right] - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$



$$g_i(\mathbf{x}) = -\frac{1}{2} \left[ (\mathbf{x} - \mu_i)^T \cdot \Sigma_i^{-1} \cdot (\mathbf{x} - \mu_i) \right] + \ln P(\omega_i)$$



Squared Mahalanobis distance

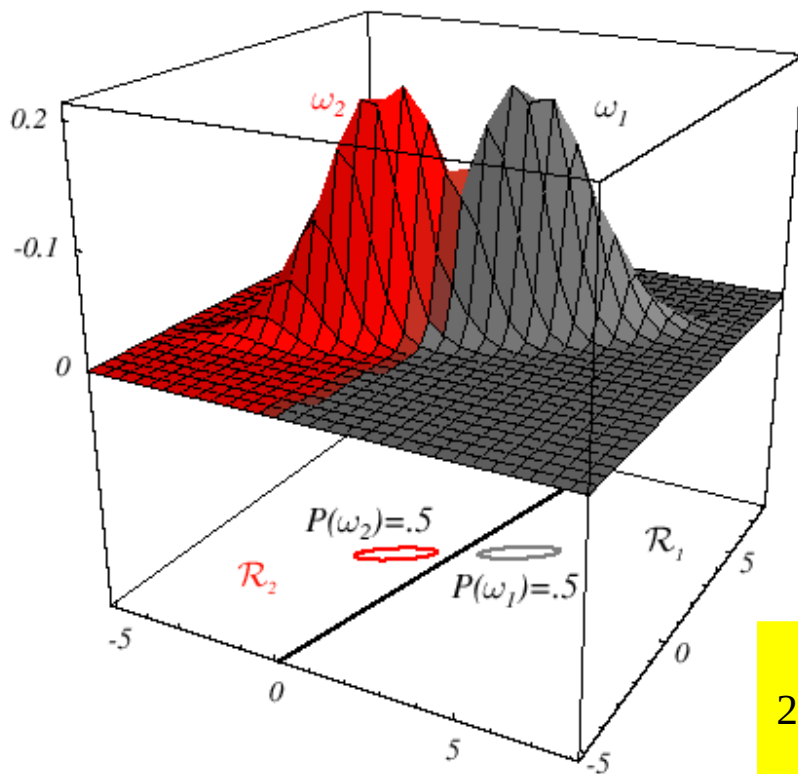




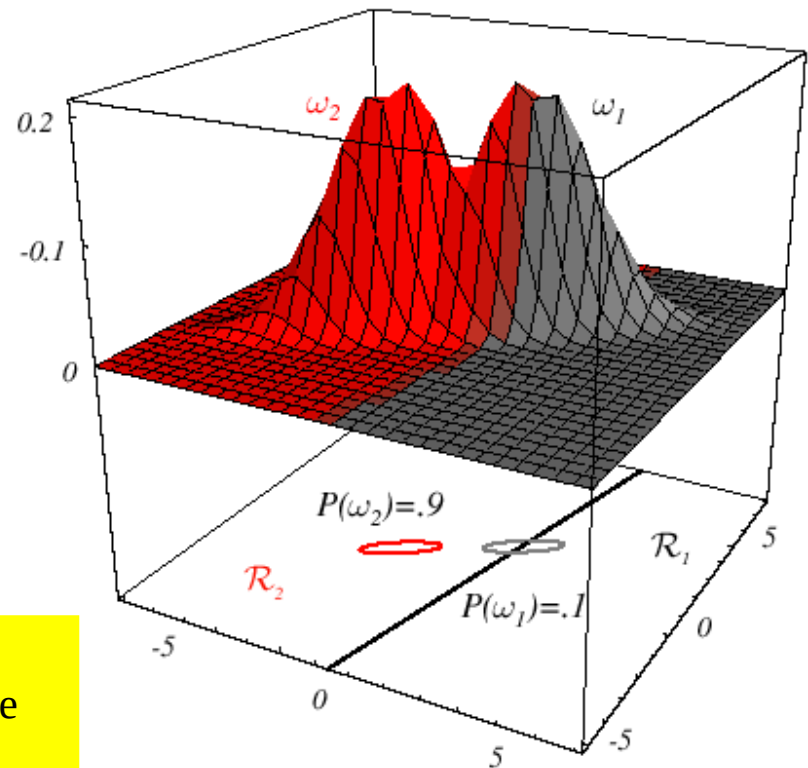
## Case 2: $\Sigma_i = \Sigma$

- The covariance matrices are arbitrary, but the same for all classes. The features create hyper-elliptical groups of the same size and shape with centers  $\mu_i$ .
- Linear decision functions → Hyper-planes as decision surfaces

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0} \quad \mathbf{w}_i = \Sigma^{-1} \mu_i, \quad w_{i0} = -\frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i + \ln P(\omega_i)$$



2D Case

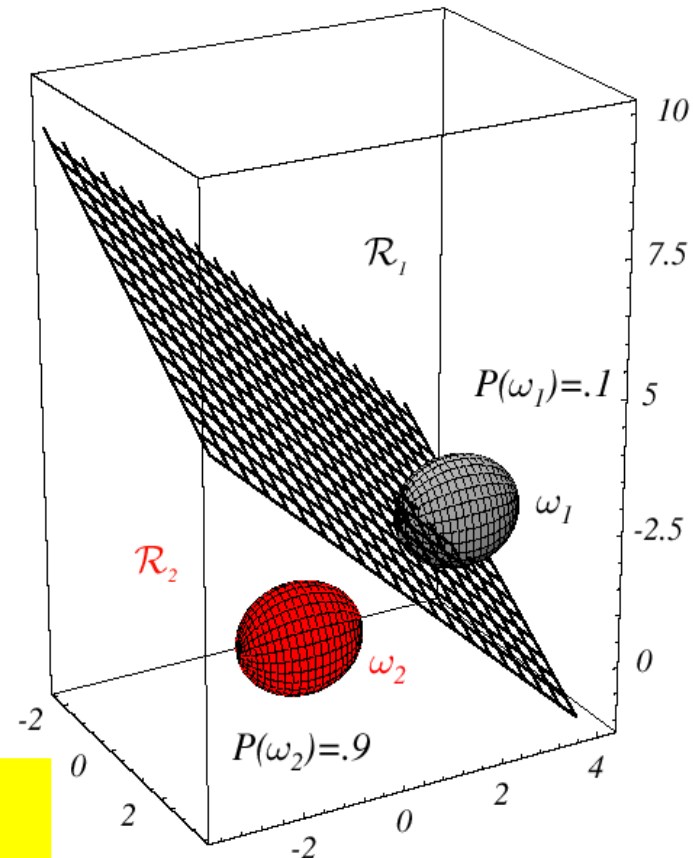
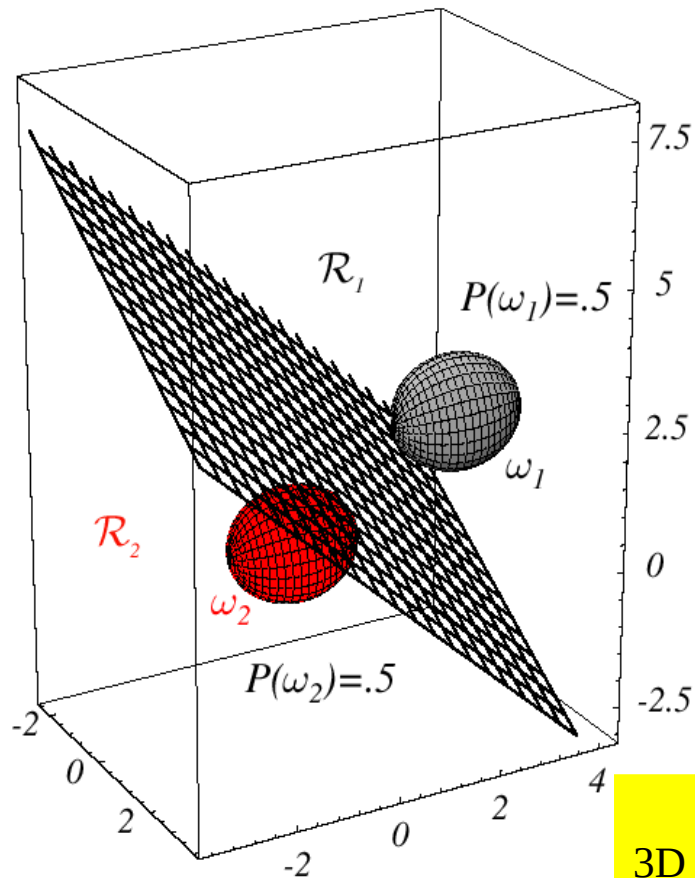






## Case 2: $\Sigma_i = \Sigma$

- The covariance matrices are arbitrary, but the same for all classes. The features create hyper-elliptical groups of the same size and shape with centers  $\mu_i$ .
- Linear decision functions → Hyper-planes as decision surfaces



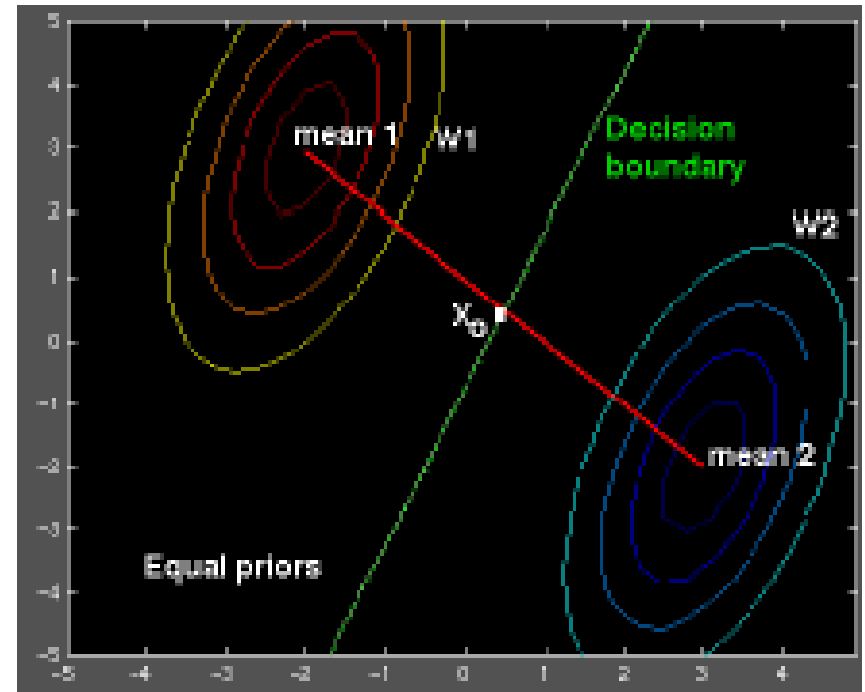
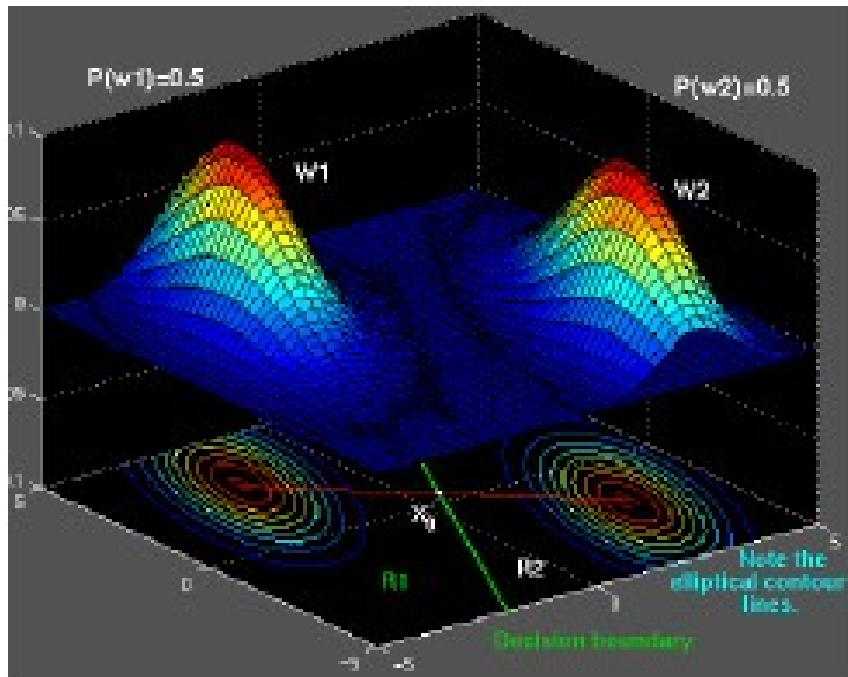
3D Case



## Case 2: $\Sigma_i = \Sigma$

➤ Decision surfaces  $\mathbf{w}^T(\mathbf{x} - \mathbf{x}_0) = 0$  where:

$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j), \quad \mathbf{x}_0 = \frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \frac{\ln(P(\omega_i)/P(\omega_j))}{(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^T \Sigma^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)$$



Since  $\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)$ , the decision hyper-plane is **not** perpendicular to the vector  $\mathbf{w}$  that connects the mean values  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\mu}_j$ .



## Case 2: $\Sigma_i = \Sigma$

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

...  $P(\omega_i)$  are the same for all  $c$  classes then the  $\ln P(\omega_i)$

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

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

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

the discriminants are linear, the resulting decision boundaries are again linear (Fig. 2.10). If  $\mathcal{R}_i$  and  $\mathcal{R}_j$  are contiguous, the boundary between them is given by the equation

$$\mathbf{w}^t (\mathbf{x} - \mathbf{x}_0) = 0, \quad (61)$$

$$\mathbf{w} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) \quad (62)$$

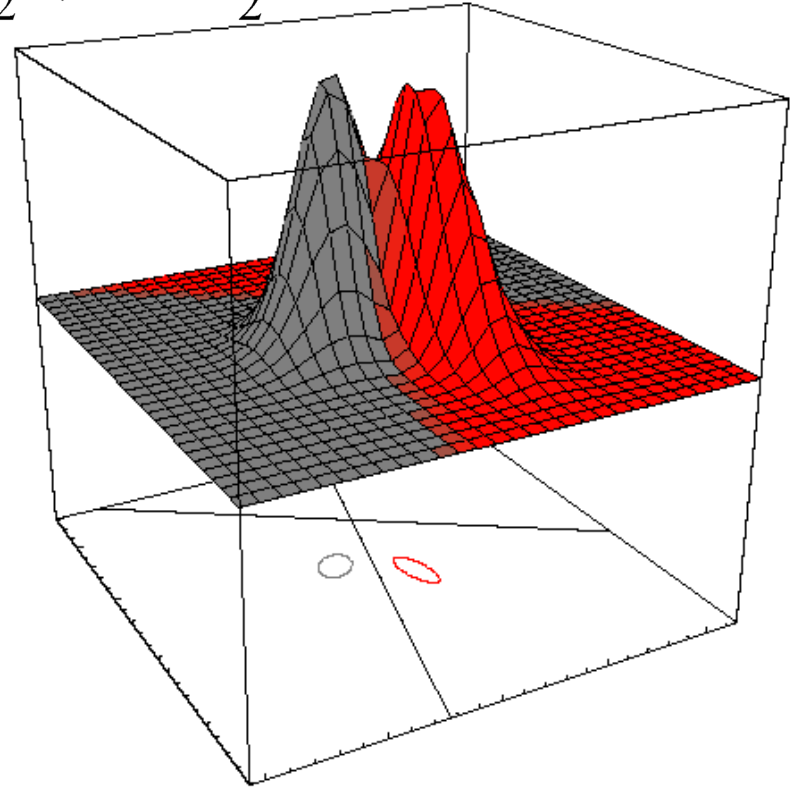
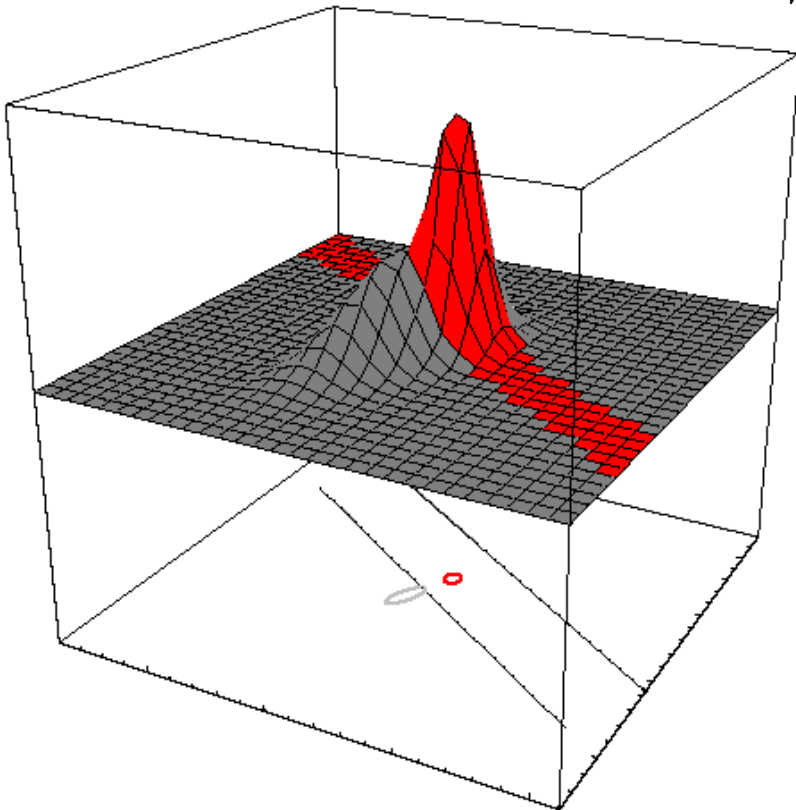
$$\mathbf{x}_0 = \frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \frac{\ln [P(\omega_i)/P(\omega_j)]}{(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^t \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j). \quad (63)$$



## Case 3: $\Sigma_i = \text{any}$

- Non-linear but squared decision functions.
- Decision surfaces *hyperquadrics* (*hyper-elliptical, hyper-paraboloid, etc.*).

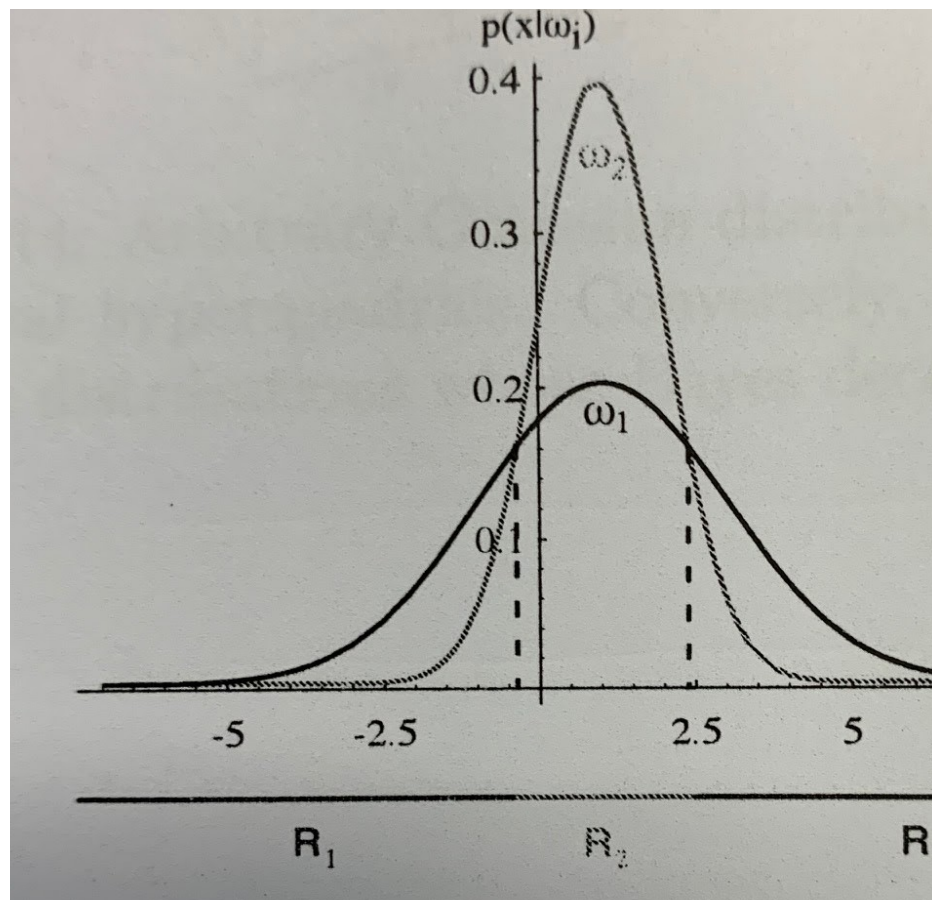
$$g_i(\mathbf{x}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0} \quad \mathbf{W}_i = -\frac{1}{2} \Sigma_i^{-1}, \quad \mathbf{w}_i = \Sigma_i^{-1} \boldsymbol{\mu}_i$$
$$w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^T \Sigma_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$





## Case 3: $\Sigma_i = \text{any}$

- Non-simply connected decision regions can arise in one dimension for Gaussians having unequal variance







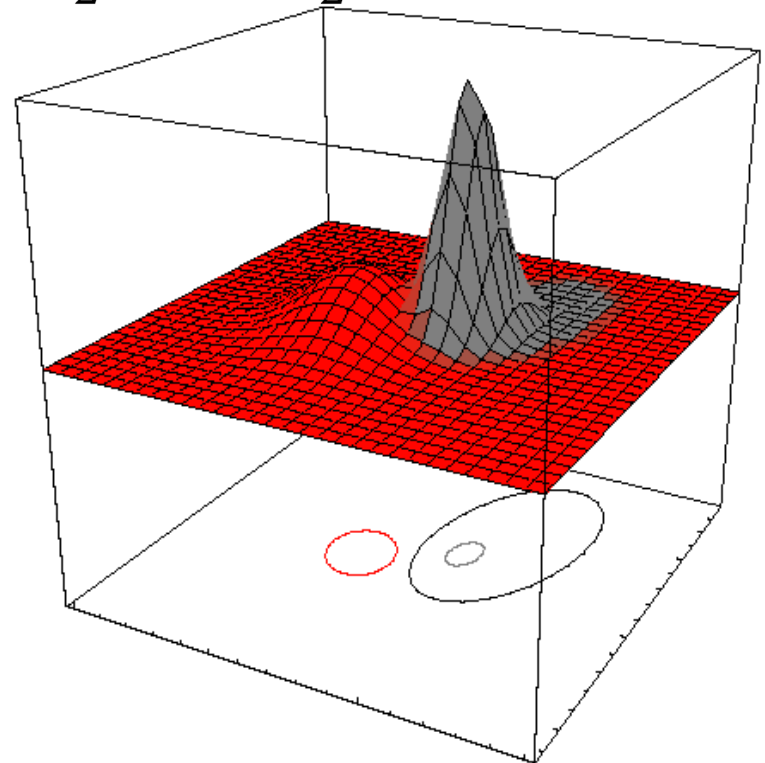
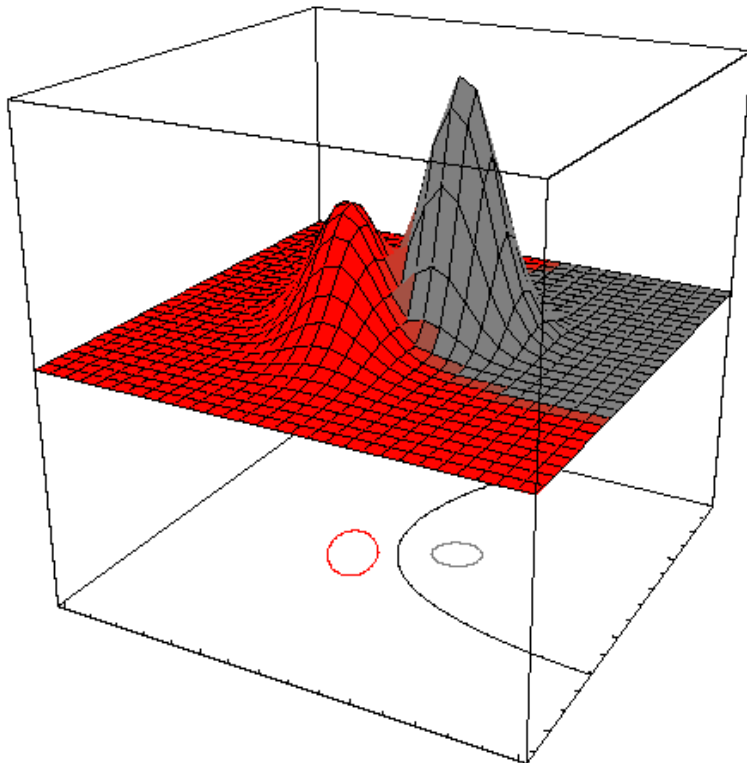
## Case 3: $\Sigma_i = \text{any}$

- Non-linear but squared decision functions.
- Decision surfaces *hyperquadratics* (*hyper-elliptical, hyper-paraboloid, etc.*).

$$g_i(\mathbf{x}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

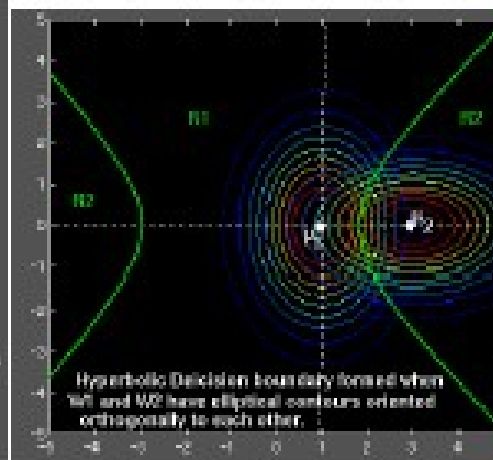
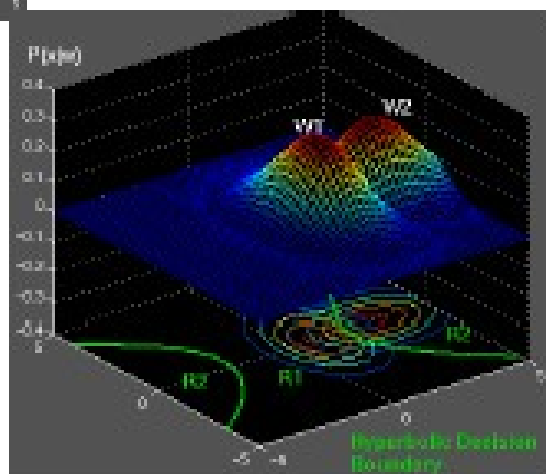
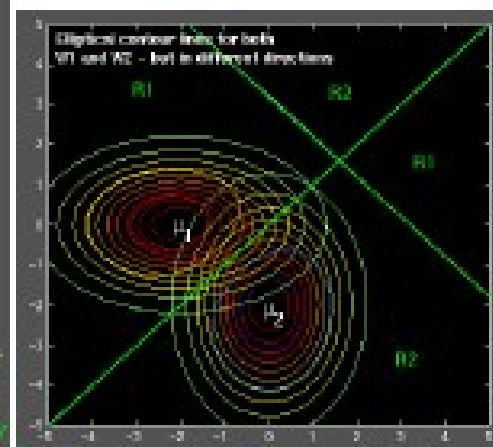
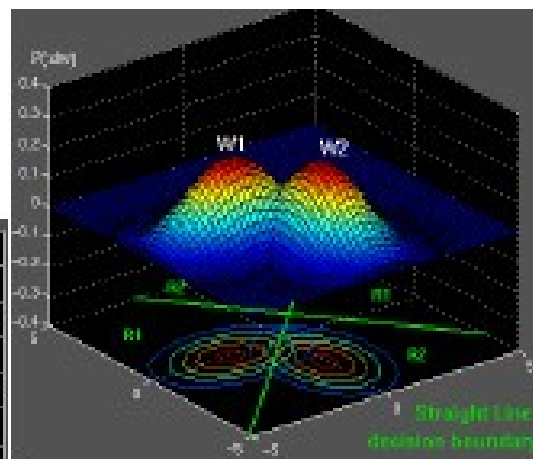
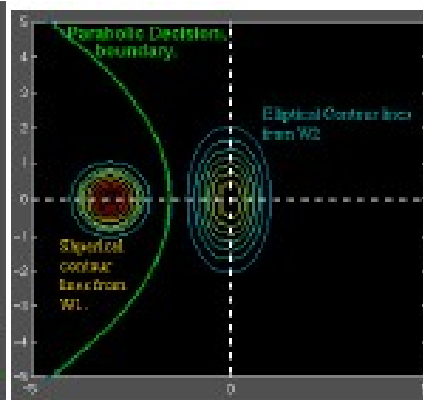
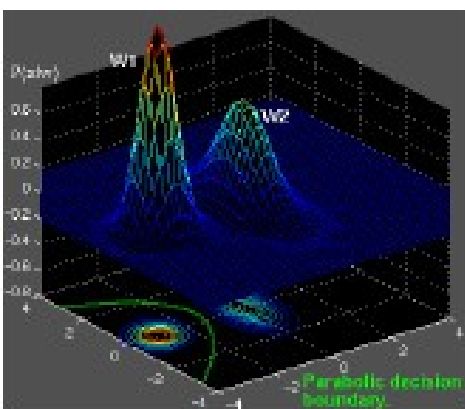
$$\mathbf{W}_i = -\frac{1}{2} \Sigma_i^{-1}, \quad \mathbf{w}_i = \Sigma_i^{-1} \boldsymbol{\mu}_i$$

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





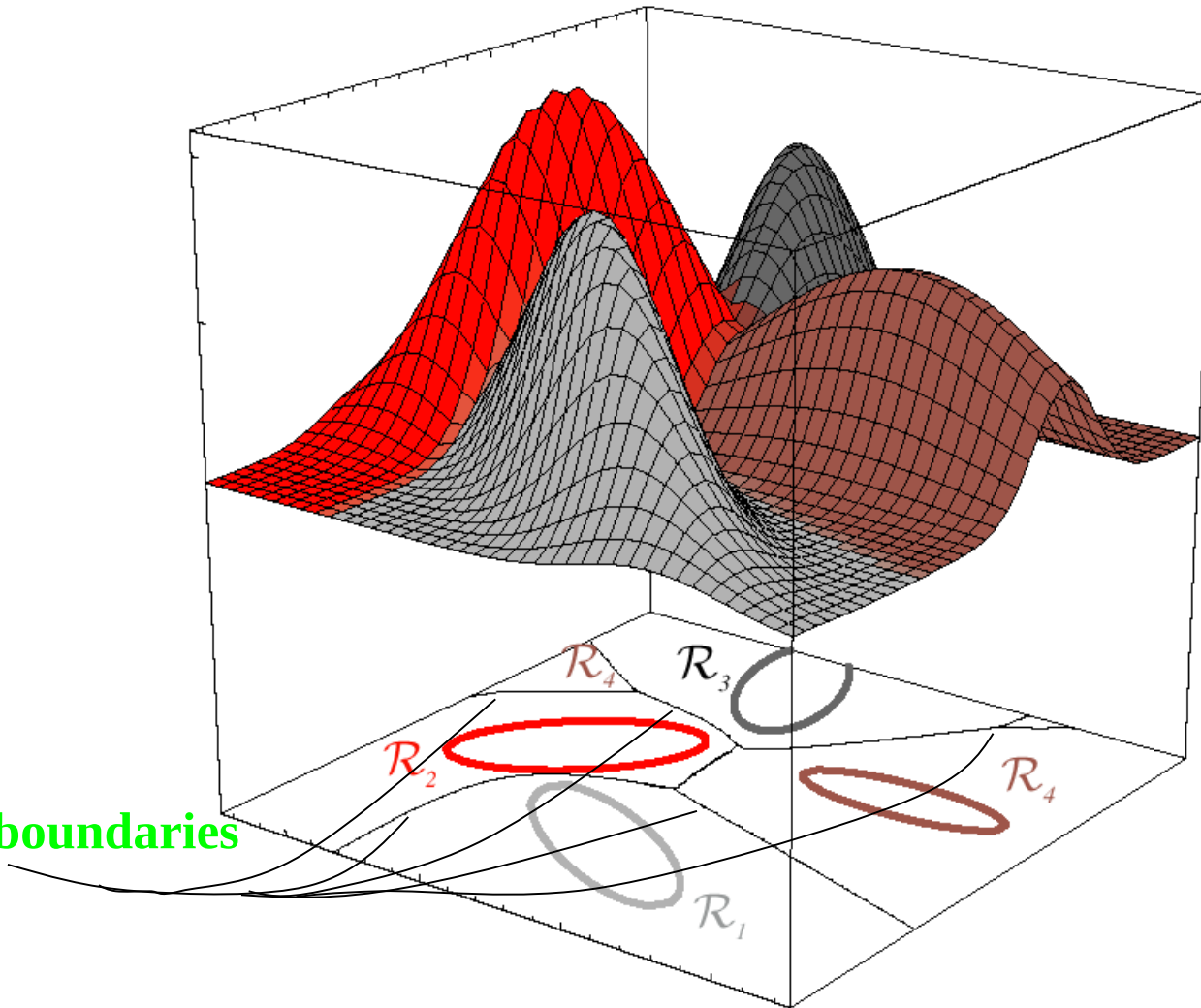
# Case 3: $\Sigma_i = \text{any}$





## Case 3: $\Sigma_i = \text{any}$

- In the case of multiple classes, the boundaries are even more complex:

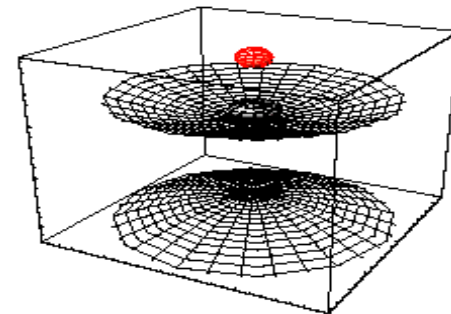
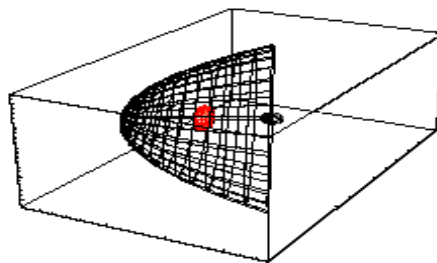
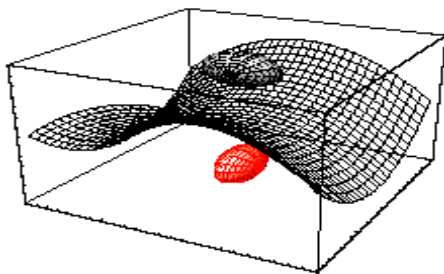
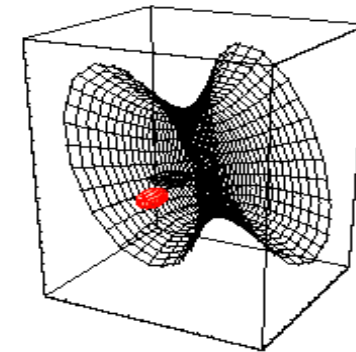
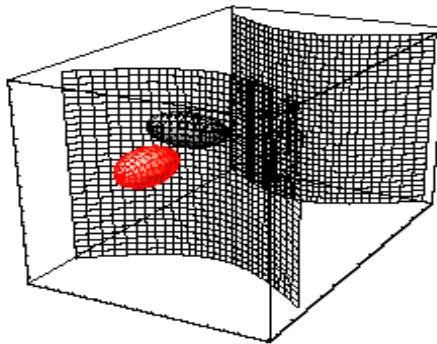
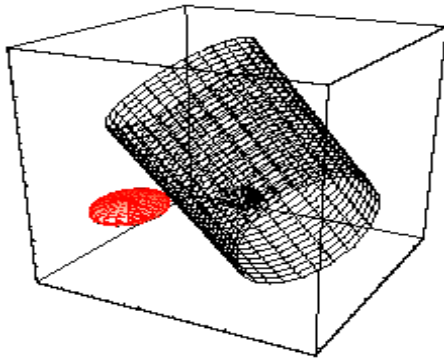
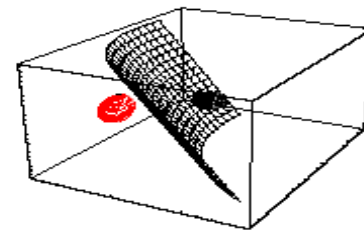
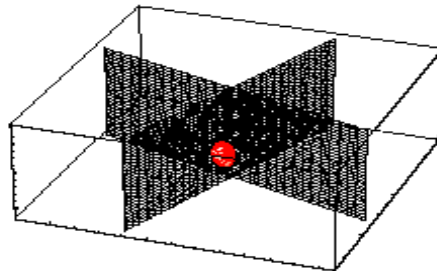
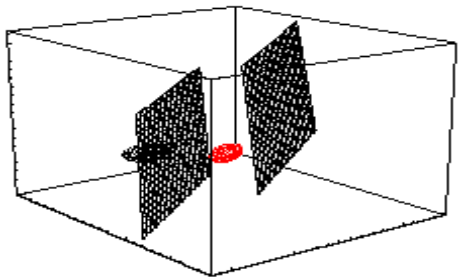






# Case 3: $\Sigma_i = \text{any}$

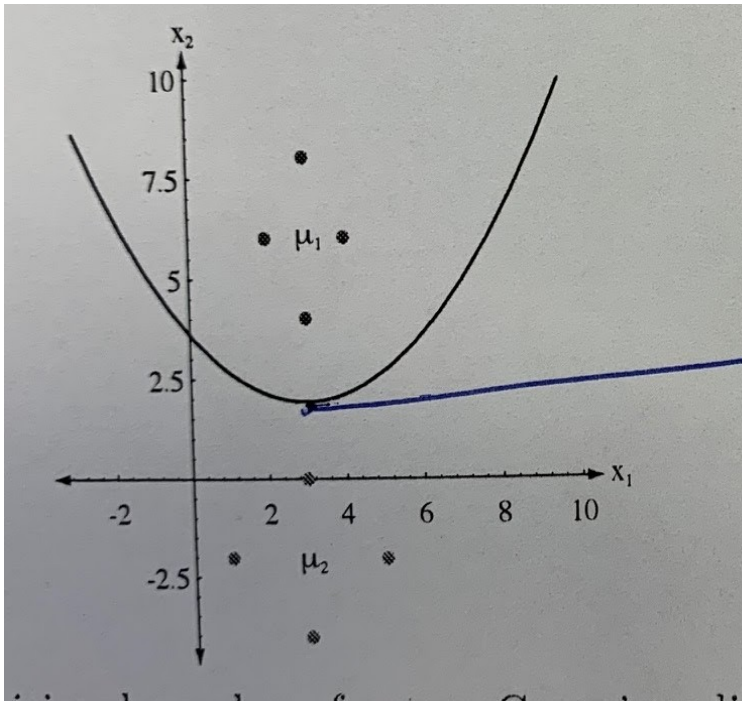
3-D





### Case 3: $\Sigma_i = \text{any}$

*Example:  
Decision Regions  
for 2-D data*

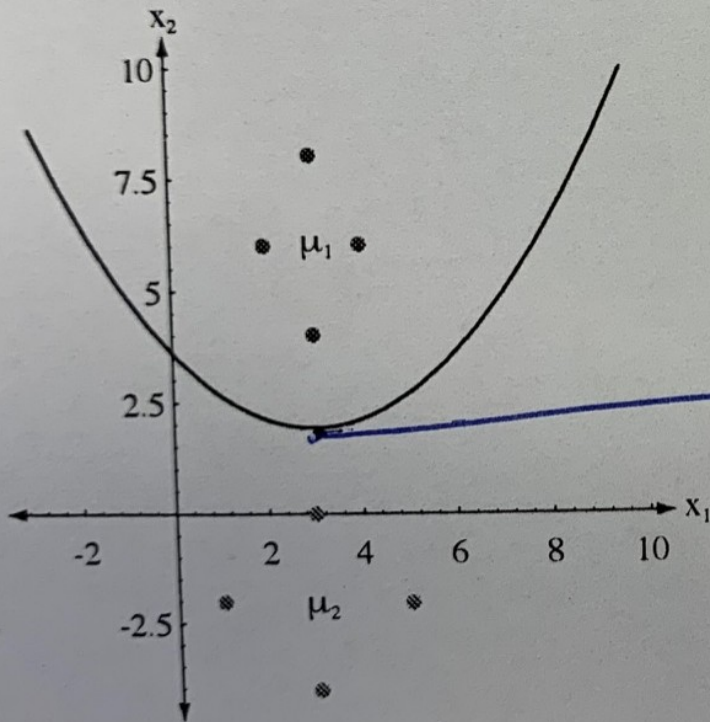


$$\mu_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}; \quad \Sigma_1 = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix} \text{ and } \mu_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix}; \quad \Sigma_2 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}.$$

$$\Sigma_1^{-1} = \begin{pmatrix} 2 & 0 \\ 0 & 1/2 \end{pmatrix} \text{ and } \Sigma_2^{-1} = \begin{pmatrix} 1/2 & 0 \\ 0 & 1/2 \end{pmatrix}.$$



### Case 3: $\Sigma_i = \text{any}$



(3  
1,83)

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



# Error Probabilities

➤ In the case of two classes, there are two cases of error:

○  $\mathbf{x}$  is at  $R_1$  while  $K\tau\Phi$  is  $w_2$ , and vice versa.

