

---

# Bayesian Classifier

---

Eun Yi Kim

---



Artificial Intelligence  
& Computer Vision  
Laboratory

---

# I N D E X

---

Bayesian Decision Theory

Simple Classification

A More General Classification

Discriminant Function

Evaluation

---



Artificial Intelligence  
& Computer Vision  
Laboratory

# Bayesian Decision Theory



- Design classifiers to make **decisions** subject to minimizing an expected **"risk"**.
  - The simplest **risk** is the **classification error**.
  - When misclassification errors are **not** equally important, the **risk** can include the **cost** associated with different misclassification errors.

# Terminology



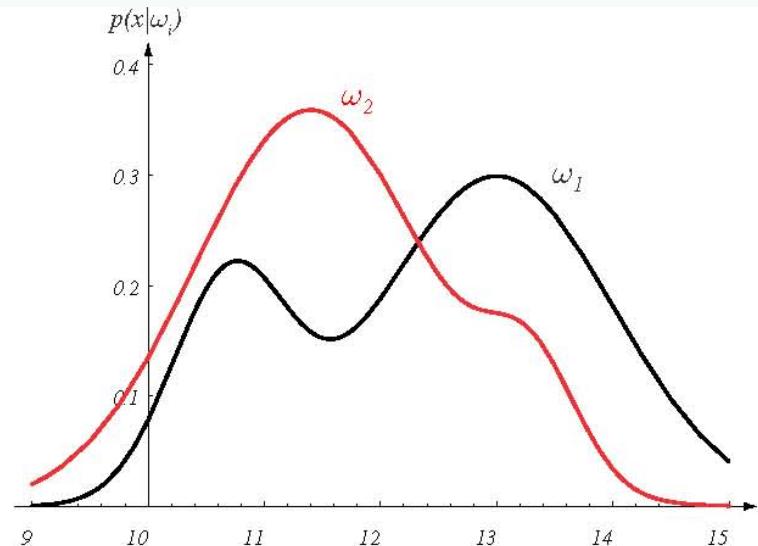
- State of nature  $\omega$  (*class label*):
  - e.g.,  $\omega_1$  for sea bass,  $\omega_2$  for salmon
- Probabilities  $P(\omega_1)$  and  $P(\omega_2)$  (*priors*):
  - e.g., prior knowledge of how likely is to get a sea bass or a salmon
- Probability density function  $p(x)$  (*evidence*):
  - e.g., how frequently we will measure a pattern with **feature value  $x$**  (e.g.,  $x$  corresponds to lightness)



# Terminology

- Conditional probability density  $p(x|\omega_j)$  (*likelihood*) :
  - e.g., how frequently we will measure a pattern with **feature value  $x$**  given that the pattern belongs to **class  $\omega_j$**

e.g., lightness distributions between salmon/sea-bass populations



**FIGURE 2.1.** Hypothetical class-conditional probability density functions show the probability density of measuring a particular feature value  $x$  given the pattern is in category  $\omega_i$ . If  $x$  represents the lightness of a fish, the two curves might describe the difference in lightness of populations of two types of fish. Density functions are normalized, and thus the area under each curve is 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons,

# Terminology



- Conditional probability  $P(\omega_j/x)$  (*posterior*) :
  - e.g., the probability that the fish belongs to **class**  $\omega_j$  given **feature**  $x$ .

# Decision Rule using Prior Probability Only



Decide  $\omega_1$  if  $P(\omega_1) > P(\omega_2)$ ; otherwise decide  $\omega_2$

$$P(error) = \begin{cases} P(\omega_1) & \text{if we decide } \omega_2 \\ P(\omega_2) & \text{if we decide } \omega_1 \end{cases}$$

or  $P(error) = \min[P(\omega_1), P(\omega_2)]$

- Favours the most likely class.
- This rule will be making the same decision all times.
  - i.e., optimum if no other information is available

# Decision Rule using Conditional Probability



- Using Bayes' rule:

$$P(\omega_j / x) = \frac{p(x / \omega_j)P(\omega_j)}{p(x)} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

where  $p(x) = \sum_{j=1}^2 p(x / \omega_j)P(\omega_j)$  (i.e., scale factor – sum of probs = 1)

Decide  $\omega_1$  if  $P(\omega_1 / x) > P(\omega_2 / x)$ ; otherwise decide  $\omega_2$   
or

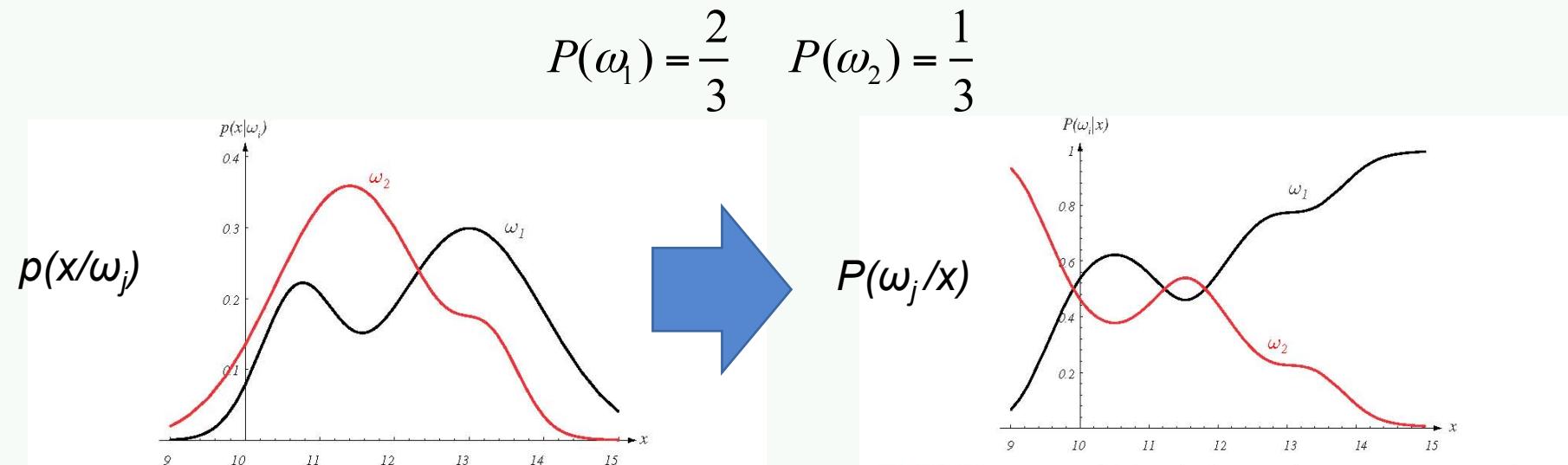
Decide  $\omega_1$  if  $p(x/\omega_1)P(\omega_1) > p(x/\omega_2)P(\omega_2)$ ; otherwise decide  $\omega_2$   
or

Decide  $\omega_1$  if  $p(x/\omega_1)/p(x/\omega_2) > P(\omega_2)/P(\omega_1)$ ; otherwise decide  $\omega_2$   
likelihood ratio threshold

# Decision Rule using Conditional Probability



Artificial Intelligence  
& Computer Vision  
Laboratory



**FIGURE 2.1.** Hypothetical class-conditional probability density functions show the probability density of measuring a particular feature value  $x$  given the pattern is in category  $\omega_i$ . If  $x$  represents the lightness of a fish, the two curves might describe the difference in lightness of populations of two types of fish. Density functions are normalized, and thus the area under each curve is 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons,

**FIGURE 2.2.** Posterior probabilities for the particular priors  $P(\omega_1) = 2/3$  and  $P(\omega_2) = 1/3$  for the class-conditional probability densities shown in Fig. 2.1. Thus in this case, given that a pattern is measured to have feature value  $x = 14$ , the probability it is in category  $\omega_2$  is roughly 0.08, and that it is in  $\omega_1$  is 0.92. At every  $x$ , the posteriors sum to 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

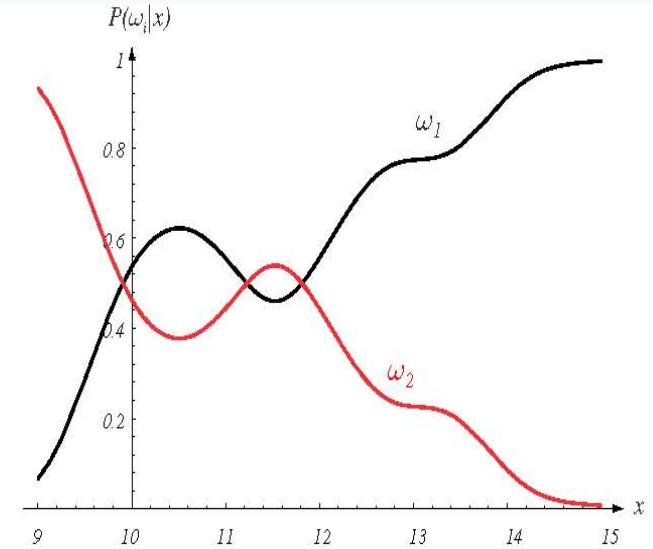


# Probability of Error

- The probability of error is defined as:

$$P(\text{error} / x) = \begin{cases} P(\omega_1 / x) \text{ if we decide } \omega_2 \\ P(\omega_2 / x) \text{ if we decide } \omega_1 \end{cases}$$

or  $P(\text{error}/x) = \min[P(\omega_1/x), P(\omega_2/x)]$



- What is the average probability error?

$$P(\text{error}) = \int_{-\infty}^{\infty} P(\text{error}, x) dx = \int_{-\infty}^{\infty} P(\text{error} / x) p(x) dx$$

- The Bayes rule is optimum, that is, it minimizes the average probability error!

# Where do Probabilities come from?



Artificial Intelligence  
& Computer Vision  
Laboratory

- There are two competitive answers:
  - (1) Relative frequency (**objective**) approach.
    - Compute probabilities from experiments.
  - (2) Bayesian (**subjective**) approach.
    - Compute probabilities from models.



# Example: Objective approach

- Classify cars whether they are more or less than \$50K:
  - Classes:  $C_1$  if price > \$50K,  $C_2$  if price <= \$50K
  - Feature:  $x$ , the **height** of a car
- Use the Bayes' rule to compute the posterior probabilities:

$$P(C_i/x) = \frac{p(x/C_i)P(C_i)}{p(x)}$$

- We need to estimate  $p(x/C_1)$ ,  $p(x/C_2)$ ,  $P(C_1)$ ,  $P(C_2)$

# Example: Objective approach

- Collect data
  - Ask drivers how much their car was and measure height.
- Determine **prior** probabilities  $P(C_1)$ ,  $P(C_2)$ 
  - e.g., 1209 samples:  $\#C_1=221$   $\#C_2=988$

$$P(C_1) = \frac{221}{1209} = 0.183$$

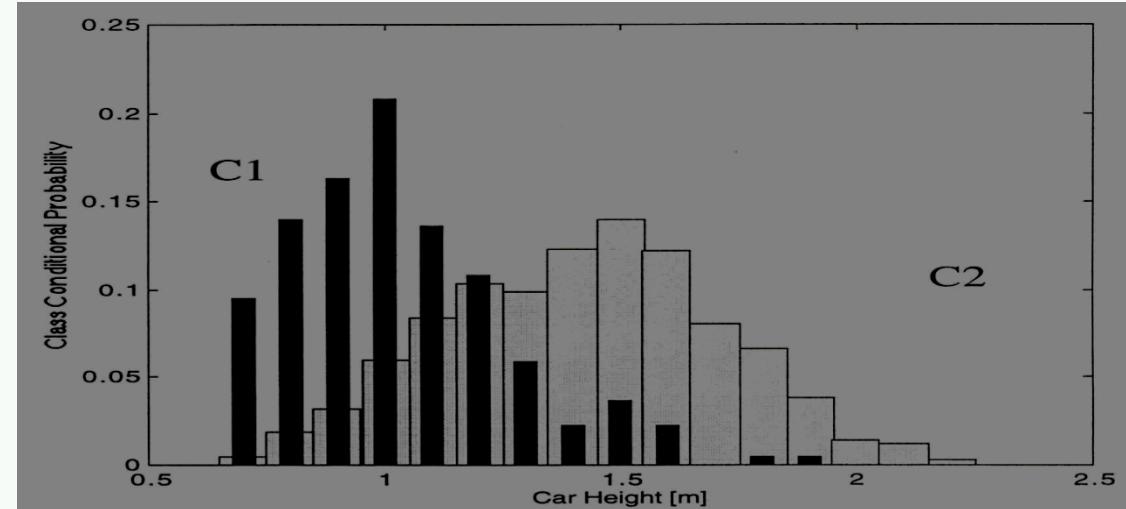
$$P(C_2) = \frac{988}{1209} = 0.817$$



# Example: Objective approach

- Determine class conditional probabilities (likelihood)
  - Discretize car height into bins and compute **normalized histogram**.

$$p(x / C_i)$$



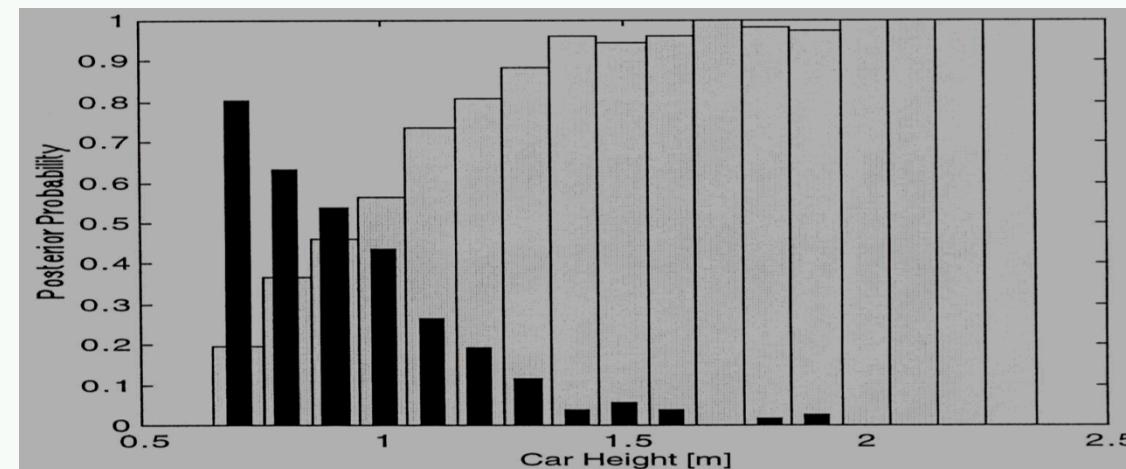


# Example: Objective approach

- Calculate the **posterior** probability for each bin, e.g.:

$$\begin{aligned} P(C_1 / x = 1.0) &= \frac{p(x = 1.0 / C_1)P(C_1)}{p(x = 1.0 / C_1)P(C_1) + p(x = 1.0 / C_2)P(C_2)} = \\ &= \frac{0.2081 * 0.183}{0.2081 * 0.183 + 0.0597 * 0.817} = 0.438 \end{aligned}$$

$P(C_i / x)$



# A More General Theory



Artificial Intelligence  
& Computer Vision  
Laboratory

- Use **more** than one features.
- Allow **more** than two categories.
- Allow **actions** other than classifying the input to one of the possible categories (e.g., **rejection**).
- Employ a more general error function (i.e., conditional “**risk**”) by associating a “**cost**” (based on a “**loss**” function) with different errors.

# Terminology

- Features form a vector
- A set of  $c$  categories  $\omega_1, \omega_2, \dots, \omega_c$
- A finite set of  $l$  actions  $\alpha_1, \alpha_2, \dots, \alpha_l$  (typically  $l \geq c$ )
- A *loss* function  $\lambda(\alpha_i / \omega_j)$ 
  - the *cost* associated with taking action  $\alpha_i$  when the correct classification category is  $\omega_j$
- Conditional risk  $R(\alpha_i / \mathbf{x})$  – expected loss of taking action  $\alpha_i$  given  $\mathbf{x}$
- Classification is now performed using the  $R(\alpha_i / \mathbf{x})$  instead of  $P(\omega_i / \mathbf{x})$

# Conditional Risk

- Suppose we take **action**  $\alpha_i$ , when  $\mathbf{x}$  is observed.
- The **conditional risk** (or **expected loss**) with taking **action**  $\alpha_i$ , is defined as:

$$R(a_i / \mathbf{x}) = \sum_{j=1}^c \lambda(a_i / \omega_j) P(\omega_j / \mathbf{x})$$

where  $P(\omega_j / \mathbf{x}) = \frac{p(\mathbf{x} / \omega_j) P(\omega_j)}{p(\mathbf{x})}$

# Overall Risk

- The **overall risk** is defined as:

$$R = \int R(a(\mathbf{x}) / \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

, where  $a(\mathbf{x})$  is a general **decision rule** that determines which action  $\alpha_1, \alpha_2, \dots, \alpha_l$  to take for every  $\mathbf{x}$ .

- Use the Bayes rule to minimize  $R$ .

# Overall Risk

- The *Bayes rule* minimizes  $R$  by:
  - (i) Computing  $R(\alpha_i/x)$  for every  $\alpha_i$  given an  $x$
  - (ii) Choosing the action  $\alpha_i$  with the minimum conditional risk  $R(\alpha_i/x)$
- The resulting minimum  $R^*$  is called *Bayes risk* and is the **best** performance that can be achieved:

$$R^* = \min R$$

# Example: Two-category classification



- Define
  - $\alpha_1$ : decide  $\omega_1$
  - $\alpha_2$ : decide  $\omega_2$
  - $\lambda_{ij} = \lambda(\alpha_i / \omega_j)$
- The conditional risks are:

$$R(a_i / \mathbf{x}) = \sum_{j=1}^c \lambda(a_i / \omega_j) P(\omega_j / \mathbf{x})$$

$$R(a_1 / \mathbf{x}) = \lambda_{11} P(\omega_1 / \mathbf{x}) + \lambda_{12} P(\omega_2 / \mathbf{x})$$
$$R(a_2 / \mathbf{x}) = \lambda_{21} P(\omega_1 / \mathbf{x}) + \lambda_{22} P(\omega_2 / \mathbf{x})$$

# Example: Two-category classification



- Minimum risk decision rule:

**Decide**  $\omega_1$  if  $R(a_1/\mathbf{x}) < R(a_2/\mathbf{x})$ ; otherwise decide  $\omega_2$

or, using 
$$R(a_1/\mathbf{x}) = \lambda_{11}P(\omega_1/\mathbf{x}) + \lambda_{12}P(\omega_2/\mathbf{x})$$
$$R(a_2/\mathbf{x}) = \lambda_{21}P(\omega_1/\mathbf{x}) + \lambda_{22}P(\omega_2/\mathbf{x})$$

**Decide**  $\omega_1$  if  $(\lambda_{21} - \lambda_{11})P(\omega_1/\mathbf{x}) > (\lambda_{12} - \lambda_{22})P(\omega_2/\mathbf{x})$ ; otherwise decide  $\omega_2$

or, using 
$$P(\omega_j/\mathbf{x}) = \frac{p(\mathbf{x}/\omega_j)P(\omega_j)}{p(\mathbf{x})}$$

**Decide**  $\omega_1$  if  $\frac{p(\mathbf{x}/\omega_1)}{p(\mathbf{x}/\omega_2)} > \frac{(\lambda_{12} - \lambda_{22})}{(\lambda_{21} - \lambda_{11})} \frac{P(\omega_2)}{P(\omega_1)}$ ; otherwise decide  $\omega_2$

likelihood ratio      threshold

# Example: Two-category classification

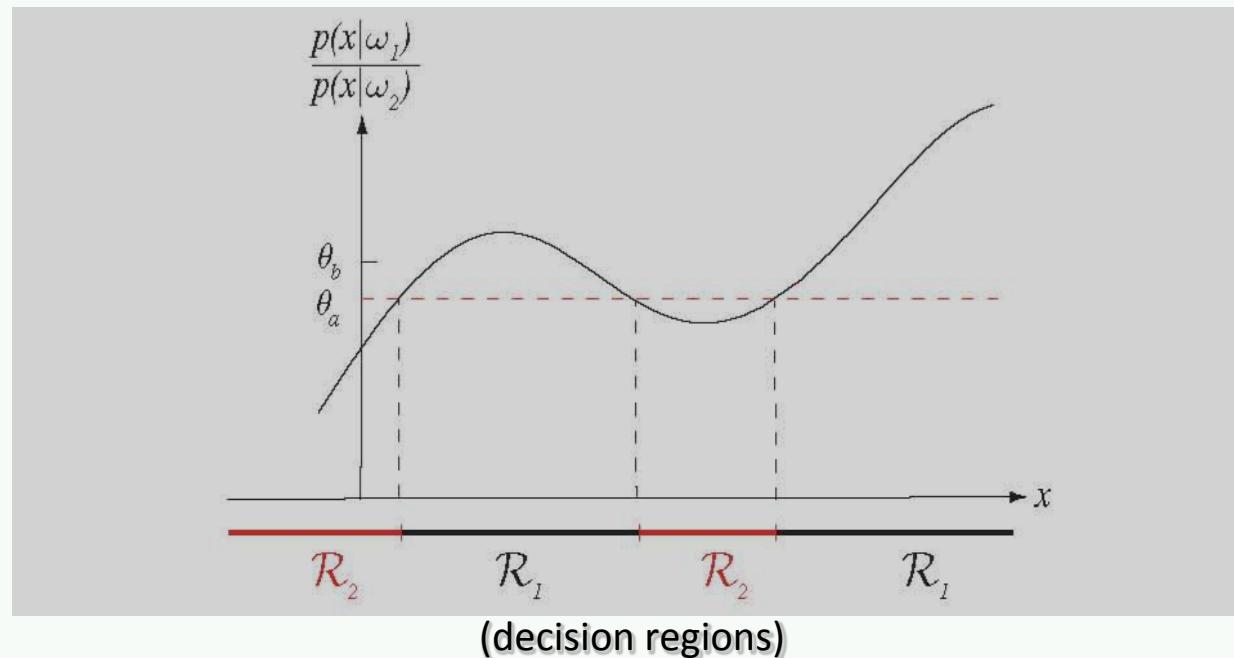


- Assuming **general loss**:

**Decide**  $\omega_1$  if  $\frac{p(\mathbf{x}/\omega_1)}{p(\mathbf{x}/\omega_2)} > \frac{(\lambda_{12} - \lambda_{22})}{(\lambda_{21} - \lambda_{11})} \frac{P(\omega_2)}{P(\omega_1)}$ ; otherwise decide  $\omega_2$

- Assuming **zero-one loss**:

**Decide**  $\omega_1$  if  $p(x/\omega_1)/p(x/\omega_2) > P(\omega_2)/P(\omega_1)$  otherwise **decide**  $\omega_2$



$$\theta_a = P(\omega_2)/P(\omega_1)$$

$$\theta_b = \frac{P(\omega_2)(\lambda_{12} - \lambda_{22})}{P(\omega_1)(\lambda_{21} - \lambda_{11})}$$

# Special case: Zero-One Loss Function



Artificial Intelligence  
& Computer Vision  
Laboratory

- Assign the **same loss** to all errors:

$$\lambda(a_i/\omega_j) = \begin{cases} 0 & i = j \\ 1 & i \neq j \end{cases}$$

- What is the conditional risk in this case?

$$R(a_i/\mathbf{x}) = \sum_{j=1}^c \lambda(a_i/\omega_j) P(\omega_j/\mathbf{x}) = \sum_{i \neq j} P(\omega_j/\mathbf{x}) = 1 - P(\omega_i/\mathbf{x})$$

# Special case: Zero-One Loss Function



- The **decision rule** becomes:

**Decide  $\omega_1$**  if  $R(a_1/\mathbf{x}) < R(a_2/\mathbf{x})$ ; otherwise decide  $\omega_2$

**or**    **Decide  $\omega_1$**  if  $1 - P(\omega_1/\mathbf{x}) < 1 - P(\omega_2/\mathbf{x})$ ; otherwise decide  $\omega_2$

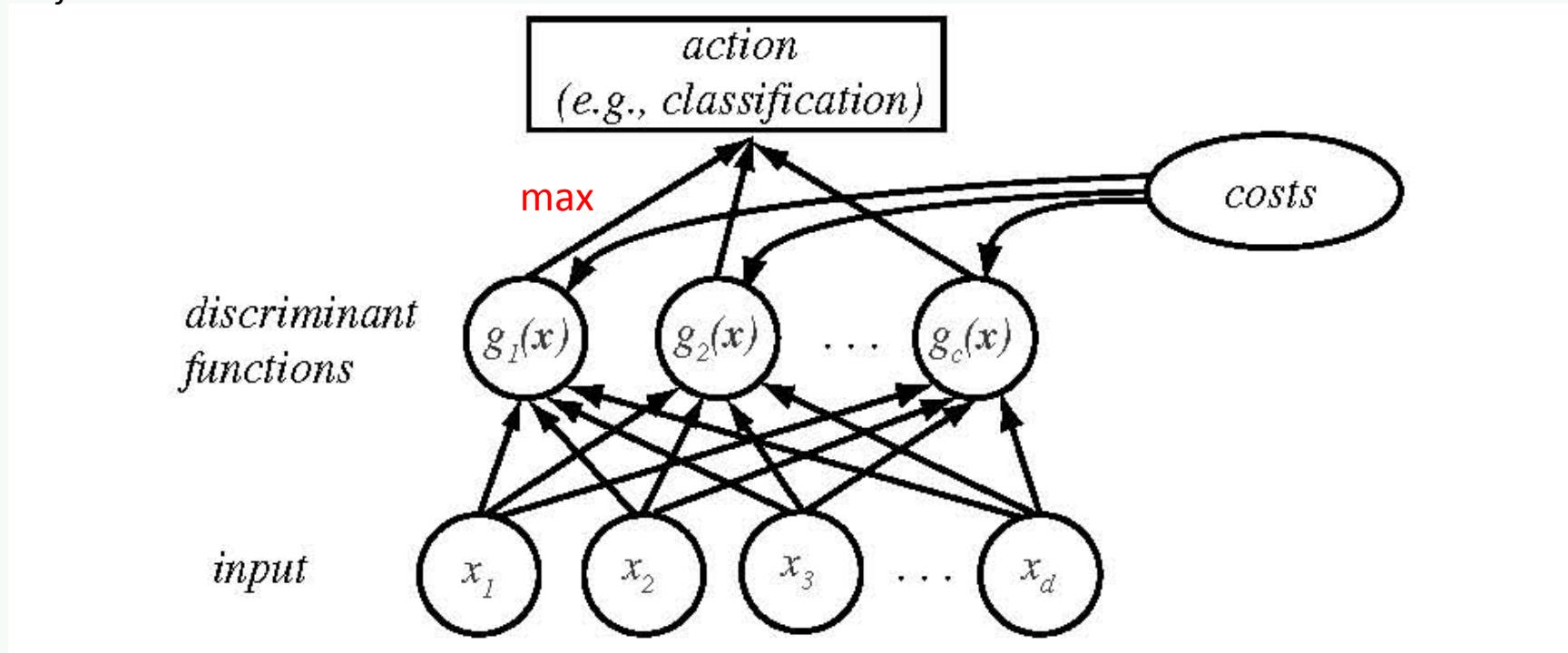
**or**    **Decide  $\omega_1$**  if  $P(\omega_1/\mathbf{x}) > P(\omega_2/\mathbf{x})$ ; otherwise decide  $\omega_2$

- The **overall risk** becomes the **average probability error!**



# Discriminant Functions

- Represent a classifier through **discriminant functions**  
 $g_i(x), i = 1, \dots, c$
- A feature vector  $\mathbf{x}$  is assigned to class  $\omega_i$  if:  
$$g_i(\mathbf{x}) > g_j(\mathbf{x}) \text{ for all } j \neq i$$



# Discriminants for Bayes Classifier



Artificial Intelligence  
& Computer Vision  
Laboratory

- Assuming a **general loss** function:

$$g_i(\mathbf{x}) = -R(\alpha_i / \mathbf{x})$$

- Assuming the **zero-one loss** function:

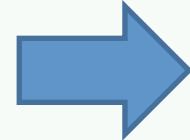
$$g_i(\mathbf{x}) = P(\omega_i / \mathbf{x})$$

# Discriminants for Bayes Classifier



- Replacing  $g_i(\mathbf{x})$  with  $f(g_i(\mathbf{x}))$ , where  $f()$  is monotonically increasing, does not change the classification results.

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



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

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

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

we'll use this extensively!

# Case of two categories

- More common to use a single discriminant function (*dichotomizer*) instead of two:

$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$$

**Decide**  $\omega_1$  if  $g(\mathbf{x}) > 0$ ; otherwise decide  $\omega_2$

Examples:

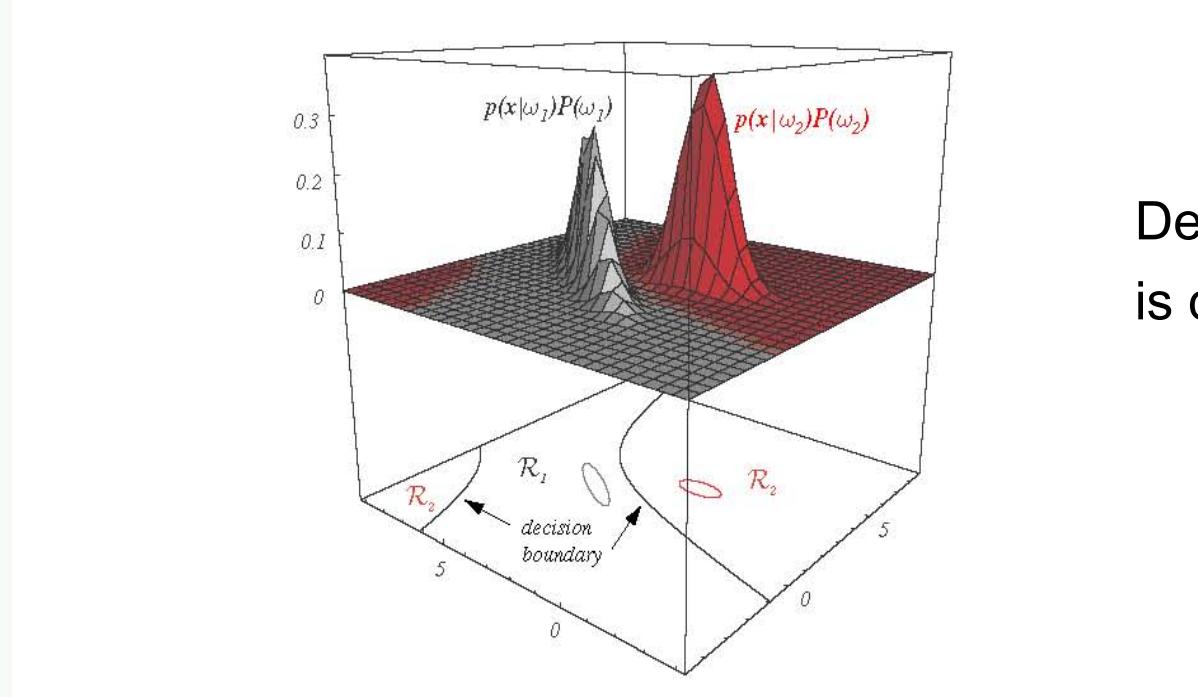
$$g(\mathbf{x}) = P(\omega_1 / \mathbf{x}) - P(\omega_2 / \mathbf{x})$$

$$g(\mathbf{x}) = \ln \frac{p(\mathbf{x} / \omega_1)}{p(\mathbf{x} / \omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

# Decision Regions and Boundaries



- Discriminants divide the feature space in *decision regions*  $R_1, R_2, \dots, R_c$ , separated by *decision boundaries*.



Decision boundary  
is defined by:

$$g_1(\mathbf{x}) = g_2(\mathbf{x})$$

# Discriminant Function for Multivariate Gaussian Density



Artificial Intelligence  
& Computer Vision  
Laboratory

$$N(\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \mu)^t \Sigma^{-1} (\mathbf{x} - \mu)\right]$$

- Consider the following discriminant function:

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

- If  $p(\mathbf{x} / \omega_i) \sim N(\mu_i, \Sigma_i)$ , then

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

# Multivariate Gaussian Density: Case I



Artificial Intelligence  
& Computer Vision  
Laboratory

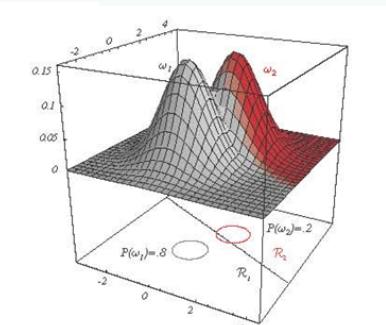
$$g_i(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^t \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i)$$

- $\boldsymbol{\Sigma}_i = \sigma^2 \mathbf{I}$  (diagonal matrix)
  - This is true when features are **uncorrelated** (or **statistically independent**) with **same variance**
  - The clusters have **spherical shape** and same size (centered at  $\boldsymbol{\mu}_i$ )

- If we disregard  $\frac{d}{2} \ln 2\pi$  and  $\frac{1}{2} \ln |\boldsymbol{\Sigma}_i|$  (constants):

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

where  $\|\mathbf{x} - \boldsymbol{\mu}_i\|^2 = (\mathbf{x} - \boldsymbol{\mu}_i)^t (\mathbf{x} - \boldsymbol{\mu}_i)$



# Multivariate Gaussian Density: Case I



Artificial Intelligence  
& Computer Vision  
Laboratory

- Expanding the above expression:

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

- Disregarding  $\mathbf{x}^t \mathbf{x}$  (constant), we get a linear discriminant:

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

where  $\mathbf{w}_i = \frac{1}{\sigma^2} \mu_i$ , and  $w_{i0} = -\frac{1}{2\sigma^2} \mu_i^t \mu_i + \ln P(\omega_i)$

# Multivariate Gaussian Density: Case I



Artificial Intelligence  
& Computer Vision  
Laboratory

- Decision boundary is determined by hyperplanes; setting  $g_i(\mathbf{x}) = g_j(\mathbf{x})$ :

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

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

# Multivariate Gaussian Density: Case I



Artificial Intelligence  
& Computer Vision  
Laboratory

- Properties of decision boundary:
  - It passes through  $\mathbf{x}_0$
  - It is orthogonal to the line linking the means.
  - What happens when  $P(\omega_i) = P(\omega_j)$  ?
  - If  $P(\omega_i) \neq P(\omega_j)$ , then  $\mathbf{x}_0$  shifts away from the most likely category.
  - If  $\sigma$  is very small, the position of the boundary is insensitive to  $P(\omega_i)$  and  $P(\omega_j)$

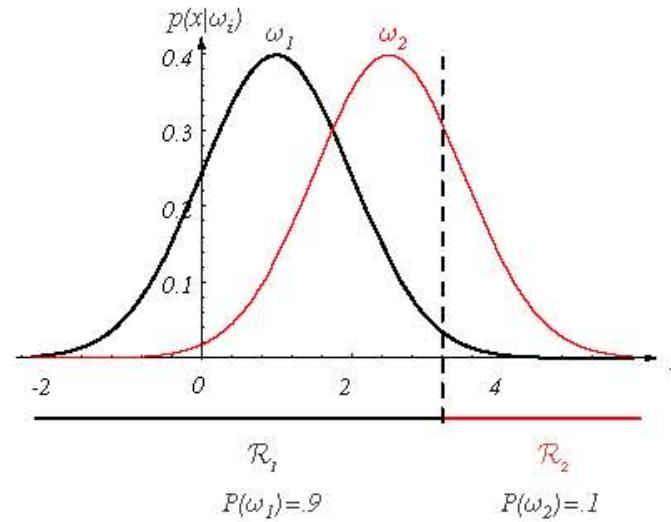
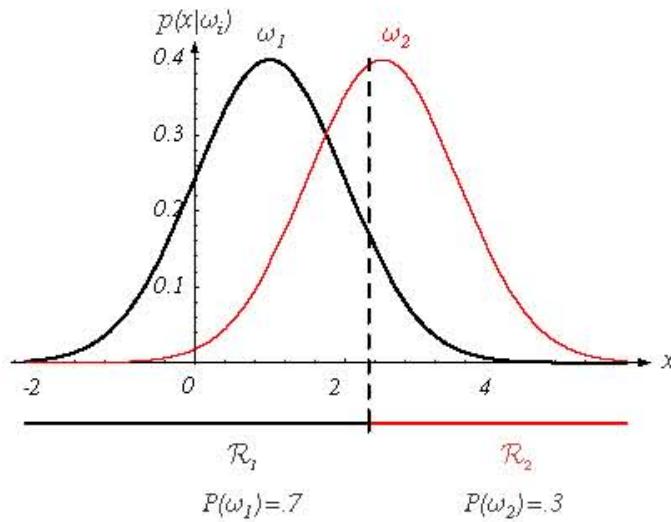
$$\mathbf{w}^t(\mathbf{x} - \mathbf{x}_0) = 0$$

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

# Multivariate Gaussian Density: Case I



Artificial Intelligence  
& Computer Vision  
Laboratory

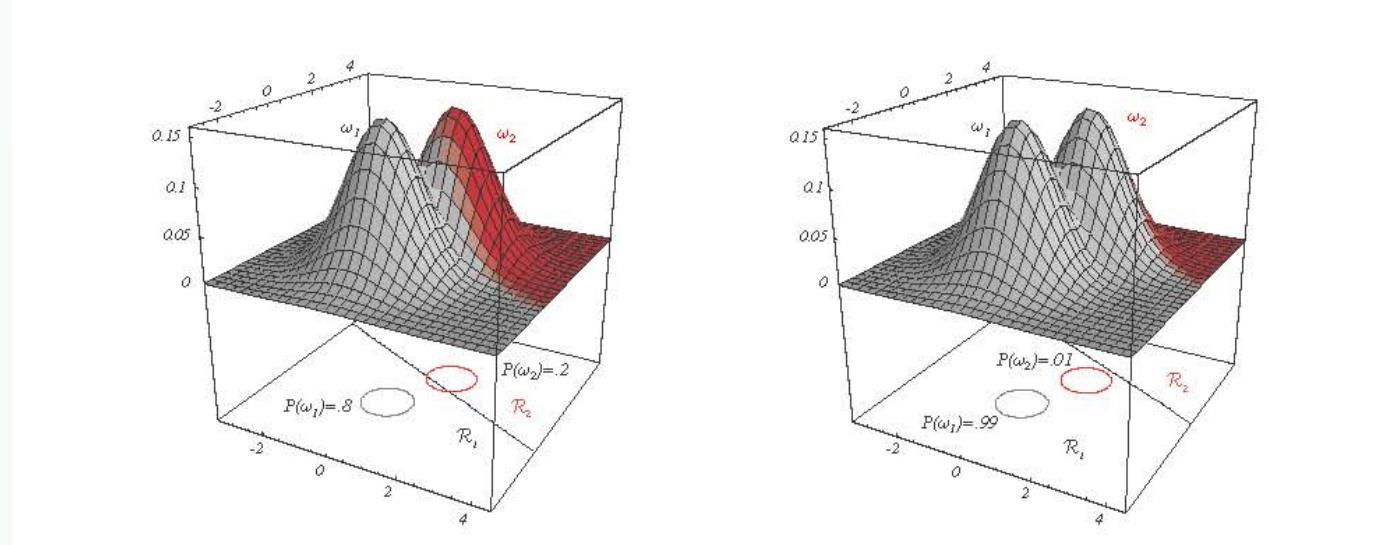


If  $P(\omega_i) \neq P(\omega_j)$ , then  $\mathbf{x}_0$  shifts away from the most likely category.

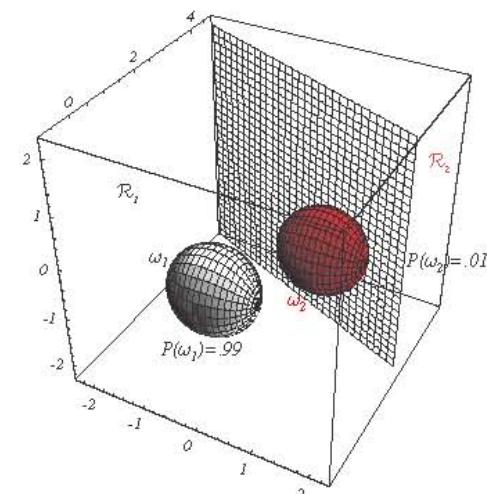
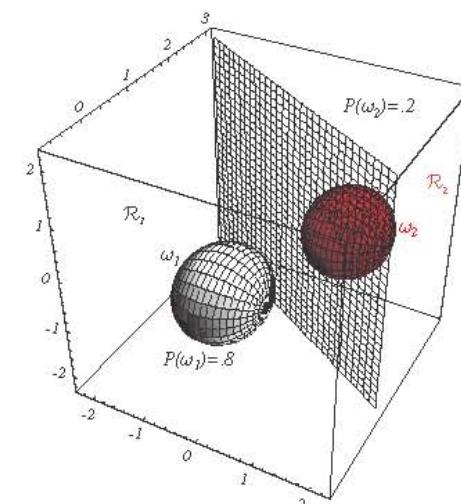
# Multivariate Gaussian Density: Case I



Artificial Intelligence  
& Computer Vision  
Laboratory



If  $P(\omega_i) \neq P(\omega_j)$ , then  $x_0$  shifts away  
from the most likely category.



# Multivariate Gaussian Density: Case I



Artificial Intelligence  
& Computer Vision  
Laboratory

- When  $P(\omega_i)$  are equal, then the discriminant becomes:

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

- This is the **Euclidean distance** (or **minimum distance classifier**)
- Recall the assumptions for **Case I**:
  - Diagonal covariance matrix (true for uncorrelated or statistically independent features).
  - All features have the same variance.

# Multivariate Gaussian Density: Case II



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

- $\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}$  (not necessarily diagonal)

- The clusters have hyperellipsoidal shape and same size (centered at  $\boldsymbol{\mu}$ ).

- If we disregard  $\frac{d}{2} \ln 2\pi$  and  $\frac{1}{2} \ln |\boldsymbol{\Sigma}_i|$  (constants):

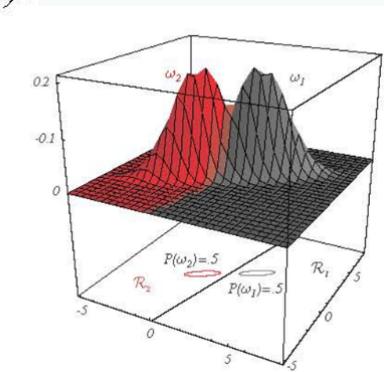
$$g_i(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \ln P(\omega_i)$$

- Expanding the above expression and disregarding the quadratic term:

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

(linear discriminant)

where  $\mathbf{w}_i = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i$ , and  $w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^t \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \ln P(\omega_i)$



# Multivariate Gaussian Density: Case II



Artificial Intelligence  
& Computer Vision  
Laboratory

- Decision boundary is determined by hyperplanes; setting  $g_i(\mathbf{x}) = g_j(\mathbf{x})$ :

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

where  $\mathbf{w} = \Sigma^{-1}(\mu_i - \mu_j)$  and  $\mathbf{x}_0 = \frac{1}{2}(\mu_i + \mu_j) - \frac{\ln[P(\omega_i)/P(\omega_j)]}{(\mu_i - \mu_j)^t \Sigma^{-1}(\mu_i - \mu_j)}(\mu_i - \mu_j)$

# Multivariate Gaussian Density: Case II



Artificial Intelligence  
& Computer Vision  
Laboratory

- Properties of hyperplane (decision boundary):
  - It passes through  $\mathbf{x}_0$
  - It is **not** orthogonal to the line linking the means.
  - What happens when  $P(\omega_i) = P(\omega_j)$  ?
  - If  $P(\omega_i) \neq P(\omega_j)$ , then  $\mathbf{x}_0$  shifts away from the most likely category.

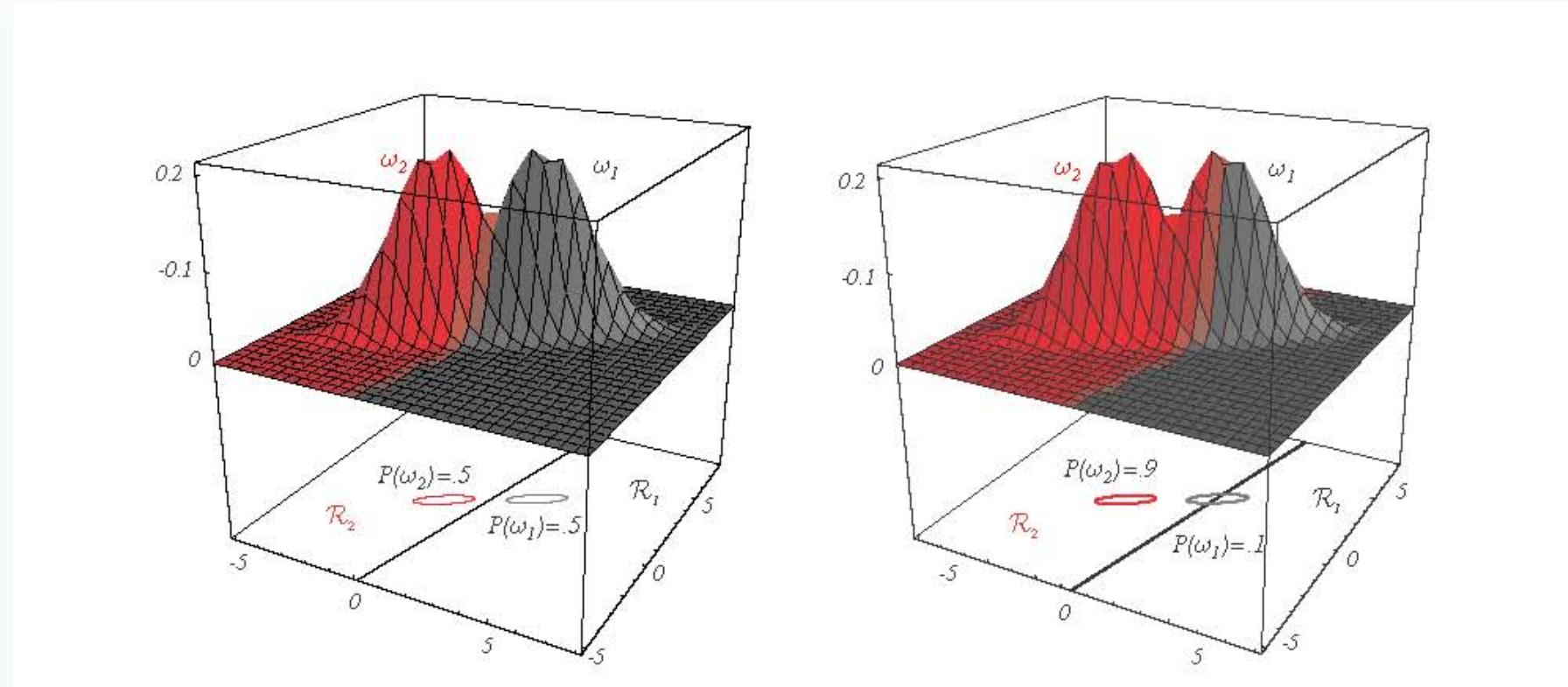
$$\mathbf{w}^t(\mathbf{x} - \mathbf{x}_0) = 0$$

where  $\mathbf{w} = \Sigma^{-1}(\mu_i - \mu_j)$  and  $\mathbf{x}_0 = \frac{1}{2}(\mu_i + \mu_j) - \frac{\ln[P(\omega_i)/P(\omega_j)]}{(\mu_i - \mu_j)^t \Sigma^{-1}(\mu_i - \mu_j)}(\mu_i - \mu_j)$

# Multivariate Gaussian Density: Case II



Artificial Intelligence  
& Computer Vision  
Laboratory

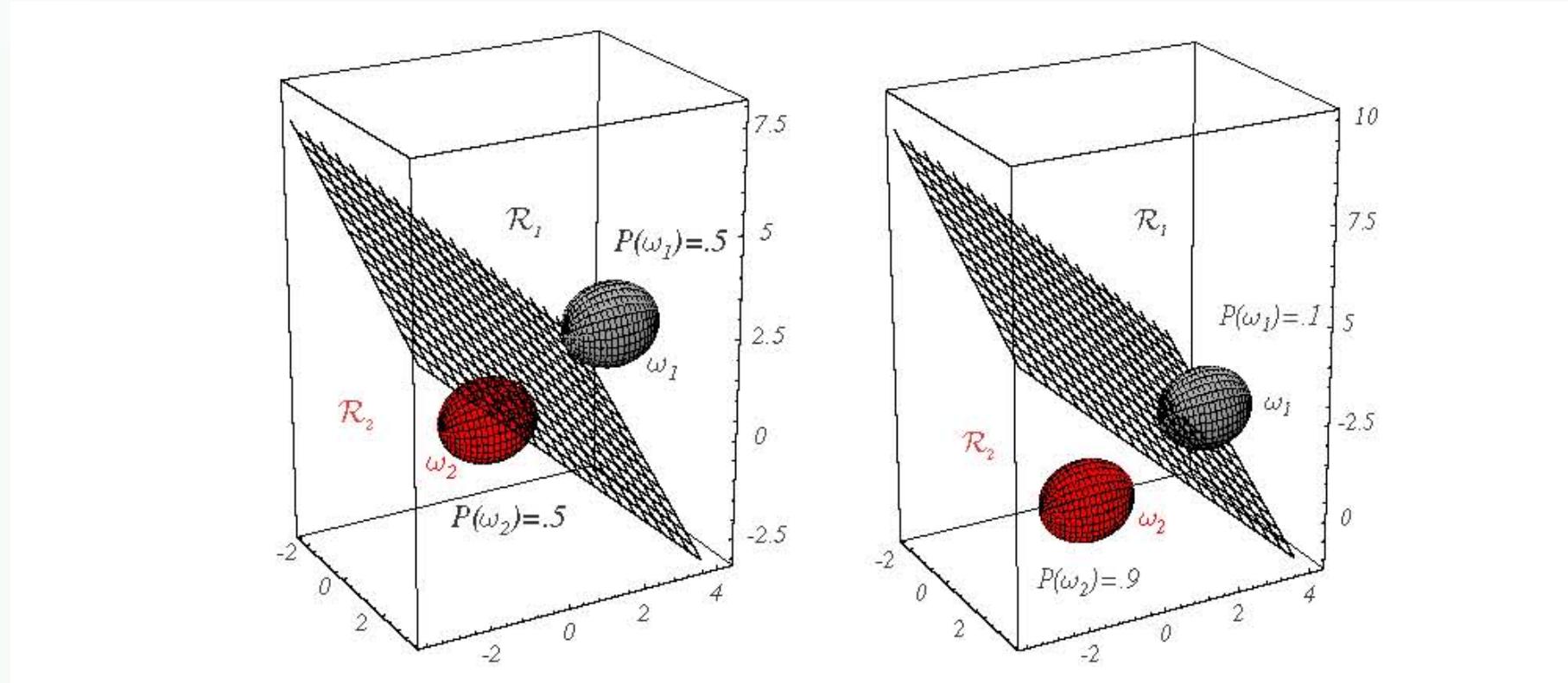


If  $P(\omega_i) \neq P(\omega_j)$ , then  $\mathbf{x}_0$  shifts away from the most likely category.

# Multivariate Gaussian Density: Case II



Artificial Intelligence  
& Computer Vision  
Laboratory



If  $P(\omega_i) \neq P(\omega_j)$ , then  $\mathbf{x}_0$  shifts away from the most likely category.

# Multivariate Gaussian Density: Case II



Artificial Intelligence  
& Computer Vision  
Laboratory

- When  $P(\omega_i)$  are equal, then the discriminant becomes:

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



$$g_i(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)$$

- This is known as the **Mahalanobis distance** (i.e., **Mahalanobis distance classifier**)

# Multivariate Gaussian Density: Case III



Artificial Intelligence  
& Computer Vision  
Laboratory

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

- $\boldsymbol{\Sigma}_i$  = arbitrary (each class has its own covariance matrix)

- The clusters have different shapes and sizes (centered at  $\boldsymbol{\mu}$ ).

- If we disregard  $\frac{d}{2} \ln 2\pi$  (constant):

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

(quadratic discriminant)

where  $\mathbf{W}_i = -\frac{1}{2} \boldsymbol{\Sigma}_i^{-1}$ ,  $\mathbf{w}_i = \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i$ , and  $w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^t \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i)$

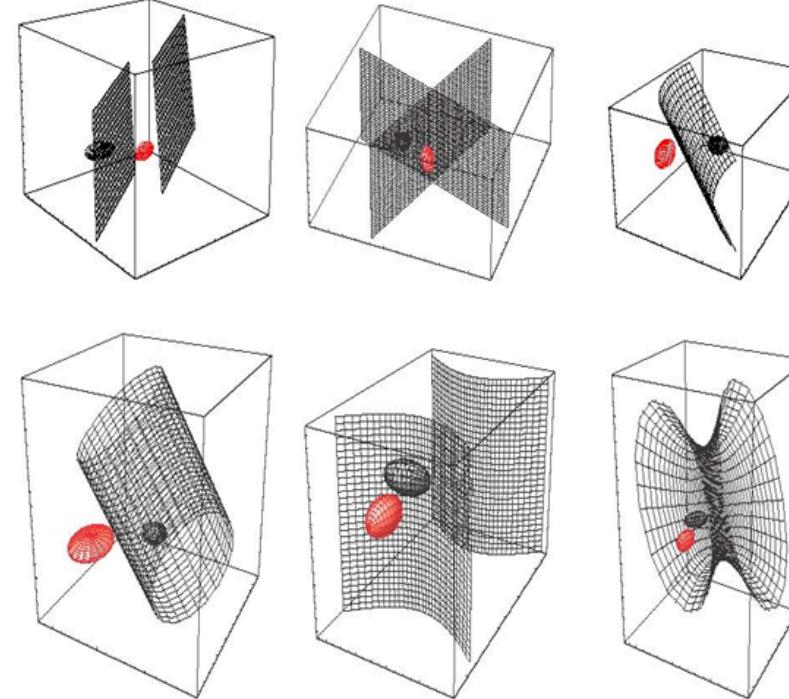
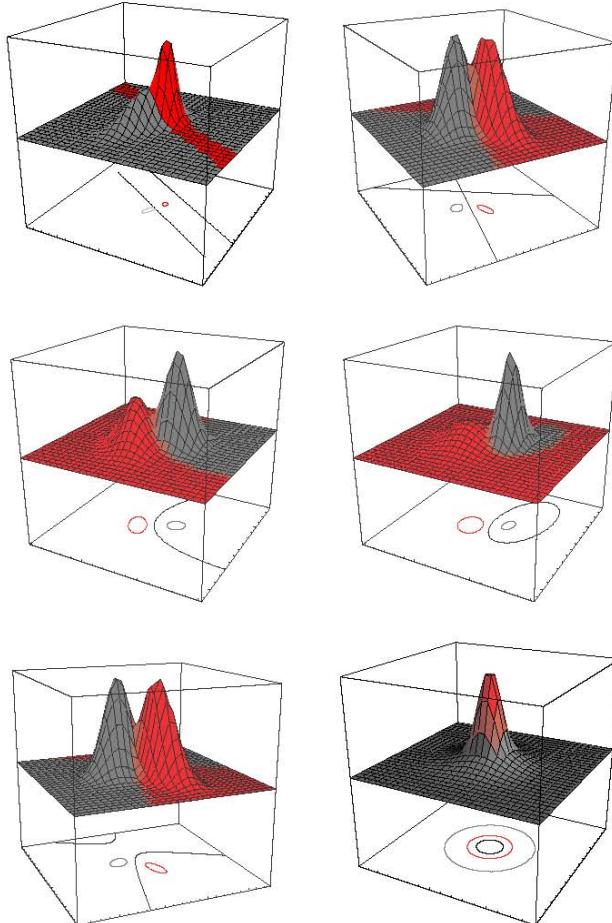
- Decision boundary is determined by hyperquadrics; setting  $g_i(\mathbf{x}) = g_j(\mathbf{x})$

e.g., hyperplanes, pairs of hyperplanes, hyperspheres, hyperellipsoids, hyperparaboloids etc.

# Multivariate Gaussian Density: Case III



Artificial Intelligence  
& Computer Vision  
Laboratory



**non-linear decision  
boundaries**

**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.

# Multivariate Gaussian Density: Case III



Artificial Intelligence  
& Computer Vision  
Laboratory

## Example

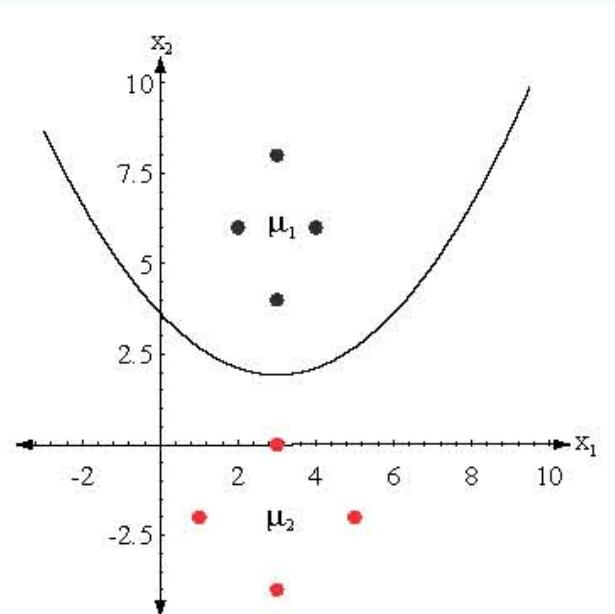
$$\mu_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}; \quad \Sigma_1 = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix} \text{ and } \mu_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix}; \quad \Sigma_2 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}.$$

Decision boundary:

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

$$P(\omega_1) = P(\omega_2)$$

boundary does  
**not** pass through  
midpoint of  $\mu_1, \mu_2$



# Error Bounds

- Exact error calculations could be difficult – easier to estimate **error bounds**.

$$P(\text{error}) = \int P(\text{error}, \mathbf{x}) d\mathbf{x} = \int P(\text{error}/\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

$$P(\text{error}/\mathbf{x}) = \begin{cases} P(\omega_1/\mathbf{x}) & \text{if we decide } \omega_2 \\ P(\omega_2/\mathbf{x}) & \text{if we decide } \omega_1 \end{cases} \quad \text{or } \min[P(\omega_1/\mathbf{x}), P(\omega_2/\mathbf{x})]$$

$$P(\omega_j / \mathbf{x}) = \frac{p(\mathbf{x}/\omega_j)P(\omega_j)}{p(\mathbf{x})}$$

- Using the inequality:

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

$$P(\text{error}) = \int \min[p(\mathbf{x}/\omega_1)P(\omega_1), p(\mathbf{x}/\omega_2)P(\omega_2)] d\mathbf{x} \leq$$

$$P^\beta(\omega_1)P^{1-\beta}(\omega_2) \int p^\beta(\mathbf{x}/\omega_1) p^{1-\beta}(\mathbf{x}/\omega_2) d\mathbf{x} = e^{-\kappa(\beta)}$$



# Error Bounds

- If the class conditional distributions are **Gaussian**, then

$$\int p^\beta(\mathbf{x}/\omega_1) p^{1-\beta}(\mathbf{x}/\omega_2) d\mathbf{x} = e^{-k(\beta)}$$

where:

$$k(\beta) = \frac{\beta(1-\beta)}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t [(1-\beta)\boldsymbol{\Sigma}_1 + \beta\boldsymbol{\Sigma}_2]^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ + \frac{1}{2} \ln \frac{|(1-\beta)\boldsymbol{\Sigma}_1 + \beta\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|^{1-\beta} |\boldsymbol{\Sigma}_2|^\beta}$$

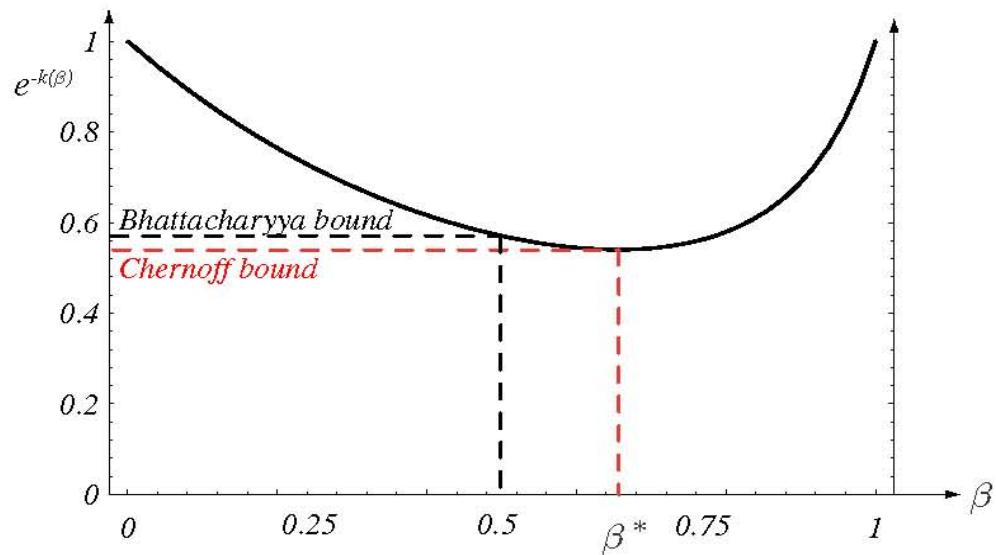
**Warning:** certain prints of our textbook have a typo!

determinant



# Error Bounds

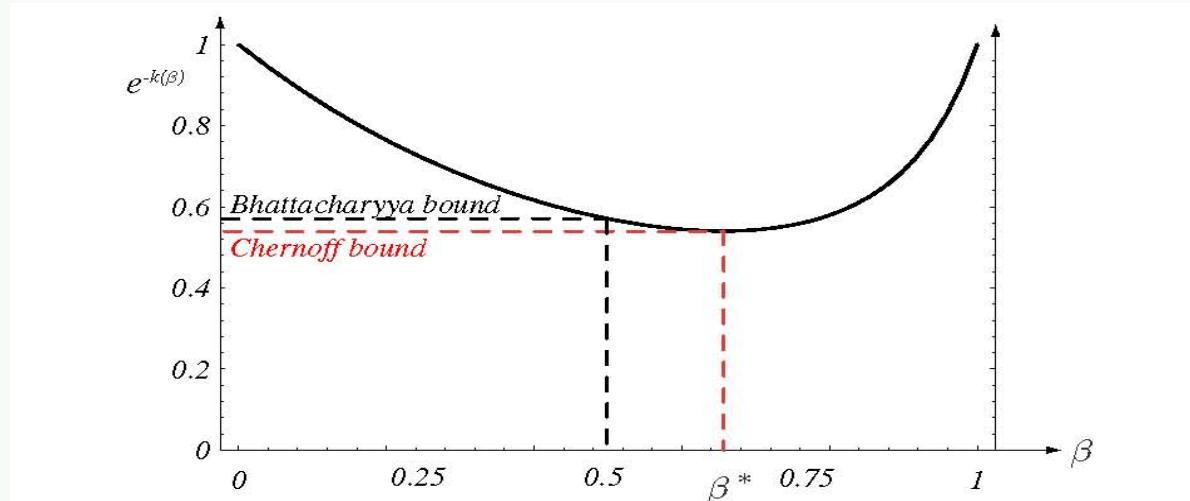
- The *Chernoff* bound is obtained by minimizing  $e^{-k(\beta)}$ 
  - This is a 1-D optimization problem, regardless to the dimensionality of the class conditional densities.



**FIGURE 2.18.** The Chernoff error bound is never looser than the Bhattacharyya bound. For this example, the Chernoff bound happens to be at  $\beta^* = 0.66$ , and is slightly tighter than the Bhattacharyya bound ( $\beta = 0.5$ ). From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

# Error Bounds

- The *Bhattacharyya* bound is obtained by setting  $\beta=0.5$ 
  - Easier to compute than Chernoff error but looser.



**FIGURE 2.18.** The Chernoff error bound is never looser than the Bhattacharyya bound. For this example, the Chernoff bound happens to be at  $\beta^* = 0.66$ , and is slightly tighter than the Bhattacharyya bound ( $\beta = 0.5$ ). From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- Note:** the Chernoff and Bhattacharyya bounds will not be good error bounds if the densities are **not** Gaussian.

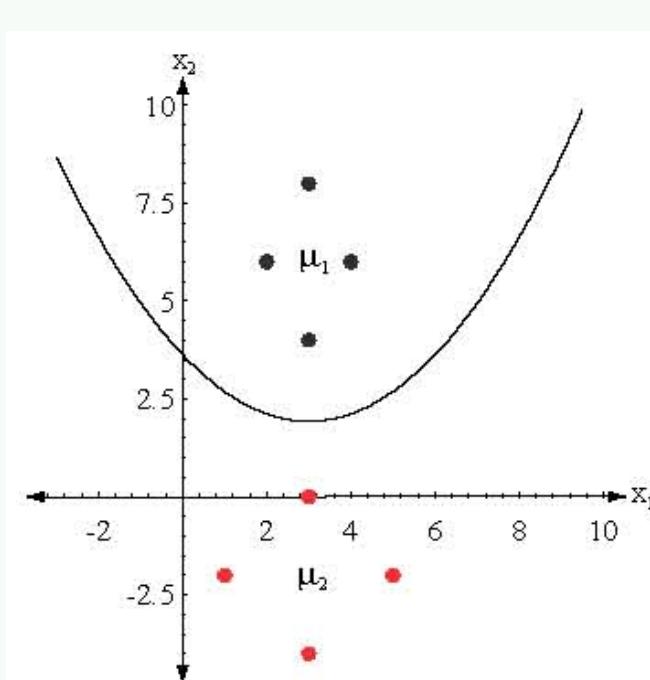
# Multivariate Gaussian Density: Case III



Artificial Intelligence  
& Computer Vision  
Laboratory

## Example

$$k(\beta) = \frac{\beta(1-\beta)}{2} (\mu_1 - \mu_2)^t [(1-\beta)\Sigma_1 + \beta\Sigma_2]^{-1} (\mu_1 - \mu_2) \\ + \frac{1}{2} \ln \frac{|(1-\beta)\Sigma_1 + \beta\Sigma_2|}{|\Sigma_1|^{1-\beta} |\Sigma_2|^\beta}.$$



$$\mu_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}; \quad \Sigma_1 = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix}$$

$$\mu_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix}; \quad \Sigma_2 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}.$$

Bhattacharyya error:

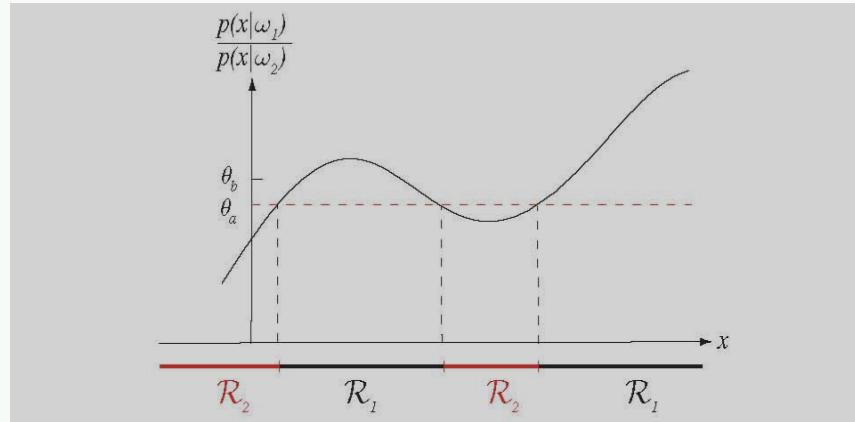
$$k(0.5)=4.06$$

$$P(error) \leq 0.0087$$

# Receiver Operating Characteristic (ROC) Curve



- Every classifier typically employs some kind of a **threshold**.



$$\theta_a = P(\omega_2) / P(\omega_1)$$

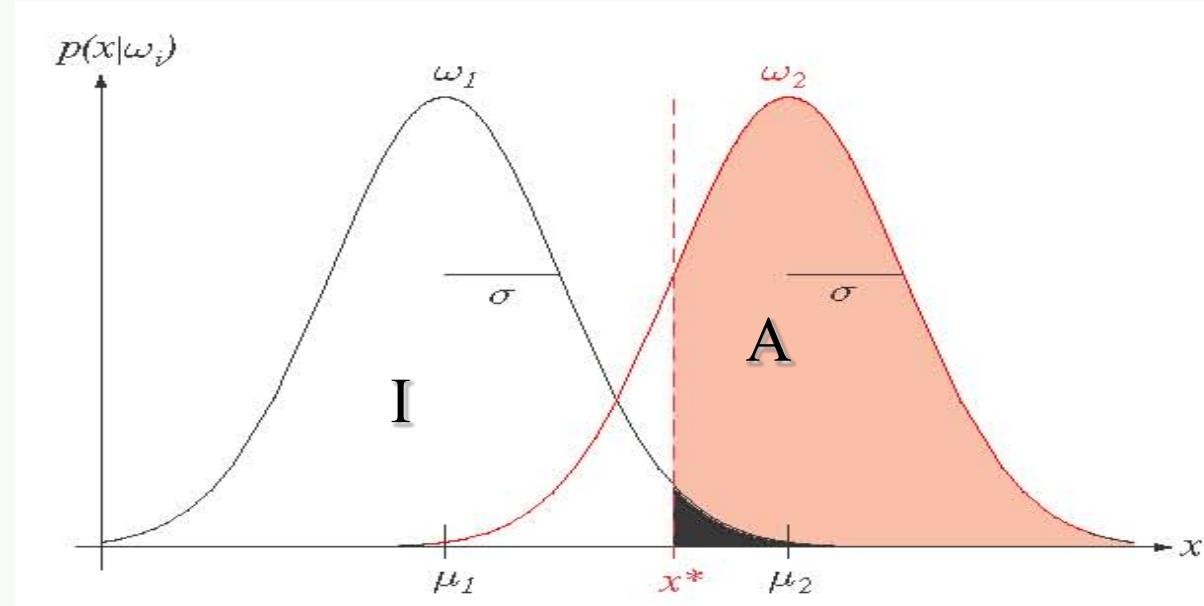
$$\theta_b = \frac{P(\omega_2)(\lambda_{12} - \lambda_{22})}{P(\omega_1)(\lambda_{21} - \lambda_{11})}$$

- Changing the threshold can affect the performance of the classifier.
- ROC curves allow us to evaluate/compare the performance of a classifier using **different** thresholds.



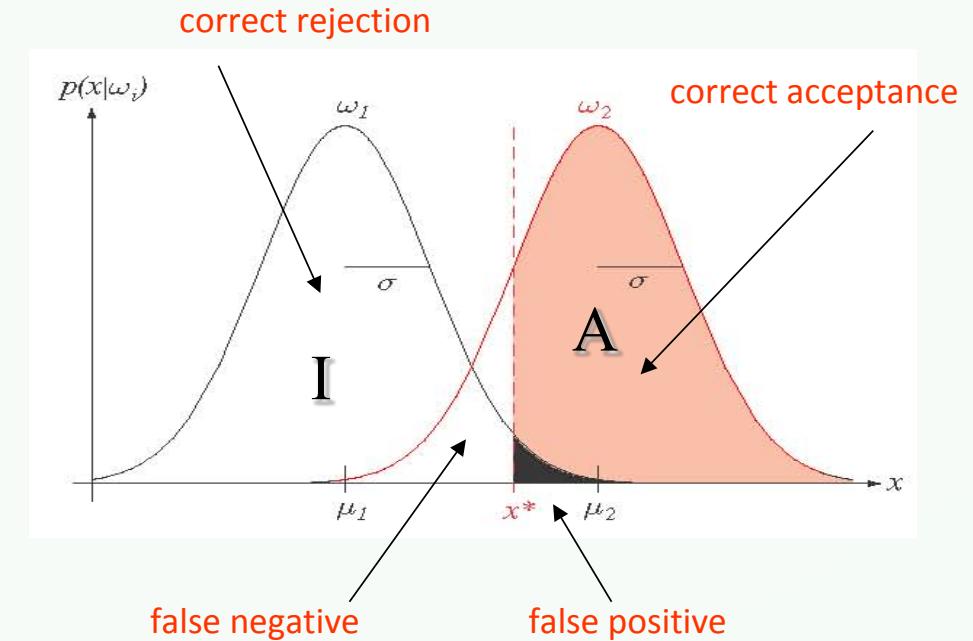
# Example: Person Authentication

- Authenticate a person using biometrics (e.g., fingerprints).
- There are two possible distributions (i.e., classes):
  - Authentic* (A) and *Impostor* (I)



# Example: Person Authentication

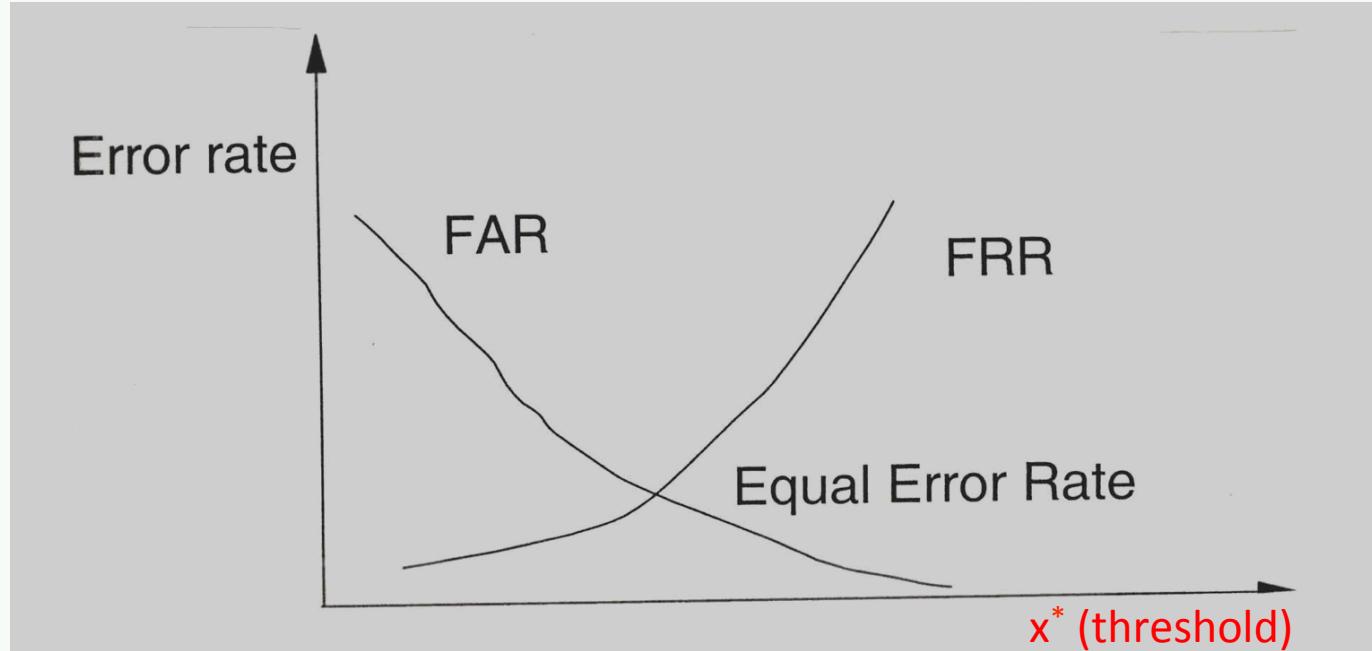
- Possible decisions:
  - (1) **correct acceptance (true positive)**:
    - X belongs to A, and we decide A
  - (2) **incorrect acceptance (false positive)**:
    - X belongs to I, and we decide A
  - (3) **correct rejection (true negative)**:
    - X belongs to I, and we decide I
  - (4) **incorrect rejection (false negative)**:
    - X belongs to A, and we decide I





# Error vs Threshold

ROC Curve



**FAR:** False Accept Rate (False Positive)

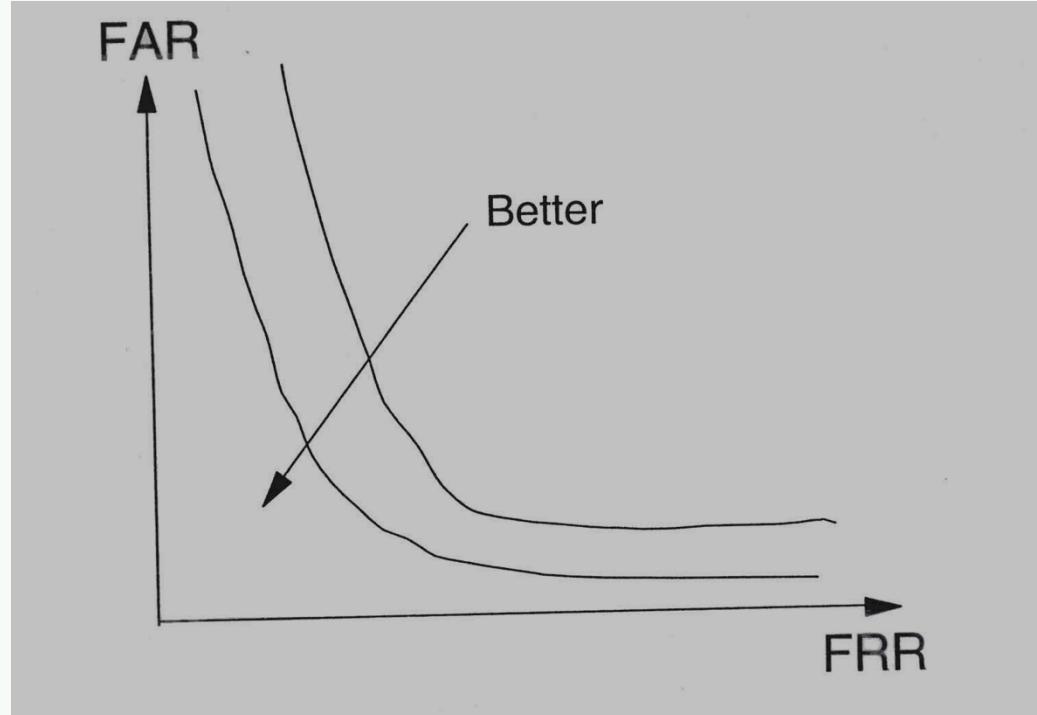
**FRR:** False Reject Rate (False Negative)

# False Negatives vs False Positives



Artificial Intelligence  
& Computer Vision  
Laboratory

## ROC Curve



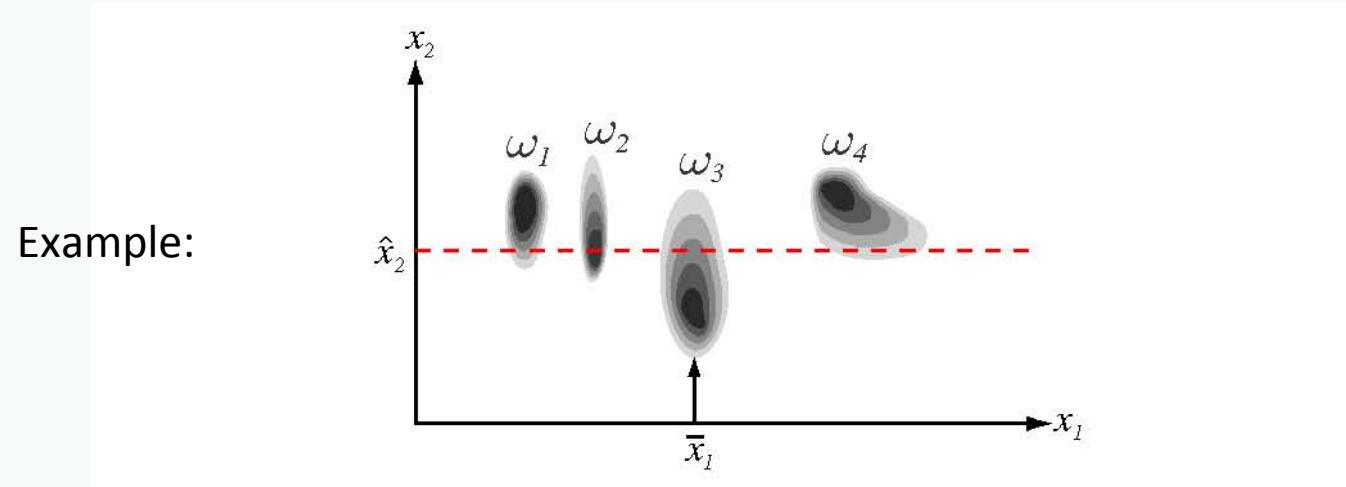
More common to plot  
**FRR** vs **FAR**

**FAR:** False Accept Rate (False Positive)  
**FRR:** False Reject Rate (False Negative)

# Missing Features



- Suppose  $\mathbf{x}=(x_1, x_2)$  is a test vector where  $x_1$  is missing and  $x_2 = \hat{x}_2$ 
  - how would we classify it?



- If we set  $x_1$  equal to the average value, we will classify  $\mathbf{x}$  as  $\omega_3$
- But  $p(\hat{x}_2 / \omega_2)$  is larger; should we classify  $\mathbf{x}$  as  $\omega_2$  ?

# Marginalize Posterior Probability



- Suppose  $\mathbf{x} = [\mathbf{x}_g, \mathbf{x}_b]$  ( $\mathbf{x}_g$ : good features,  $\mathbf{x}_b$ : bad features)
- Derive the Bayes rule using the good features:

$$P(\omega_i/\mathbf{x}_g) \stackrel{p}{=} \frac{P(\omega_i, \mathbf{x}_g)}{p(\mathbf{x}_g)} = \frac{\int P(\omega_i, \mathbf{x}_g, \mathbf{x}_b) d\mathbf{x}_b}{p(\mathbf{x}_g)} =$$
$$\frac{\int P(\omega_i/\mathbf{x}_g, \mathbf{x}_b) p(\mathbf{x}_g, \mathbf{x}_b) d\mathbf{x}_b}{p(\mathbf{x}_g)} = \frac{\int P(\omega_i/\mathbf{x}_g, \mathbf{x}_b) p(\mathbf{x}) d\mathbf{x}_b}{\int p(\mathbf{x}) d\mathbf{x}_b}$$

marginalize  
over “bad”  
features.

**Decide**  $\omega_1$  if  $P(\omega_1/\mathbf{x}_g) > P(\omega_2/\mathbf{x}_g)$ ; otherwise decide  $\omega_2$

# Quiz



- **When:**
- **What:** Bayesian Decision Theory (case studies will not be included in the quiz)