



Chapter 2

Bayesian Decision Theory

Bayesian Decision Theory

MIMA

- Bayesian decision theory is a statistical approach to data mining/pattern recognition
- Mathematical foundation for decision making
- Using probabilistic approach to help making decision so as to minimize the risk (cost).

Bayesian Decision Theory

MIMA

- Basic Assumptions
 - The decision problem is posed (formalized) in **probabilistic terms**
 - All the relevant probability values are known
- Key Principle
 - **Bayes Theorem**

Preliminaries and Notations

MIMA

$\omega_i \in \{\omega_1, \omega_2, \dots, \omega_c\}$: a state of nature

$P(\omega_i)$: prior probability

\mathbf{x} : feature vector

$p(\mathbf{x})$: evidence probability

$p(\mathbf{x} | \omega_i)$: class-conditional density / likelihood

$P(\omega_i | \mathbf{x})$: posterior probability

Decision Before Observation

MIMA

- The Problem
 - To make a decision where
 - Prior probability is known
 - No observation is allowed
- Naïve Decision Rule

Decide ω_1 if $P(\omega_1) > P(\omega_2)$, otherwise ω_2

- This is the best we can do without observation
- Fixed prior probabilities -> Same decisions all time

Bayes Theorem

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

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



Thomas Bayes
(1702-1761)

Decision After Observation

MIMA

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

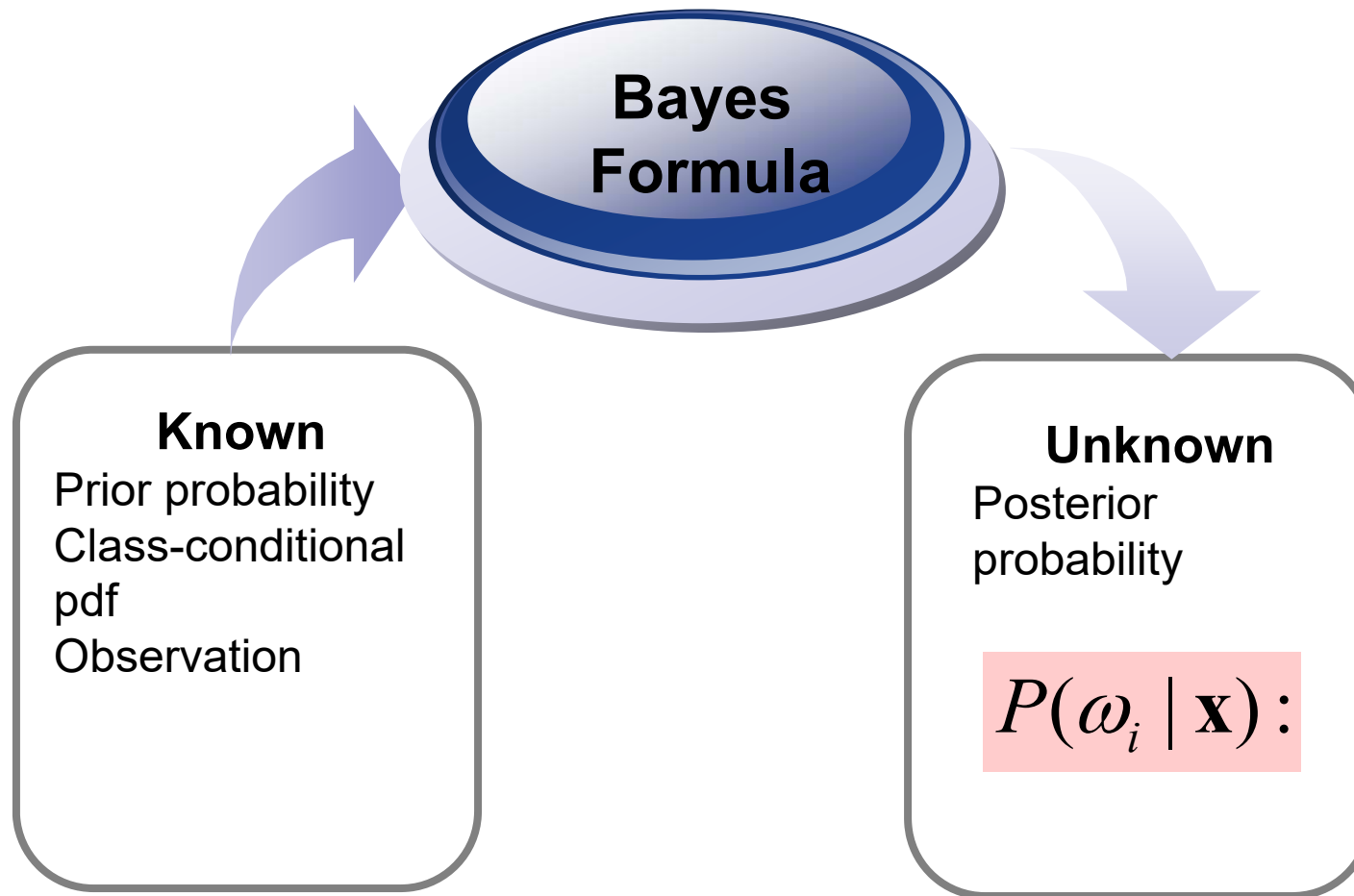
*unimportant in
making decision*

$$\mathcal{D}(\mathbf{x}) = \arg \max_{\omega_i} P(\omega_i | \mathbf{x})$$

Decision After Observation

MIMA

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

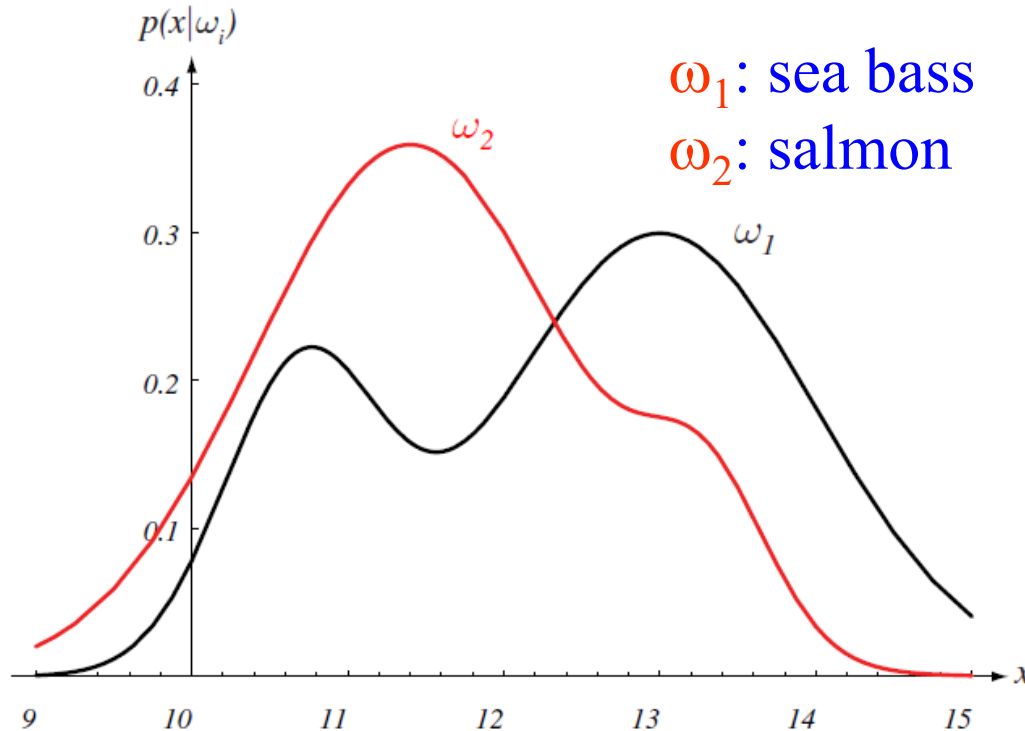


Special Cases

$$P(\omega_i | \mathbf{x}) = \frac{p(\mathbf{x} | \omega_i)P(\omega_i)}{p(\mathbf{x})} \left(\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \right)$$

- Case I: Equal prior probability
 - $P(\omega_1)=P(\omega_2)=\dots=P(\omega_c)=1/c$
 - Depends on the likelihood $p(\mathbf{x}|\omega_j)$
- Case II: Equal likelihood
 - $p(\mathbf{x}|\omega_1)=p(\mathbf{x}|\omega_2)=\dots=p(\mathbf{x}|\omega_c)$
 - Degenerate to naïve decision rule
- Normally, prior probability and likelihood function together in Bayesian decision process

An example



$$P(\omega_1) = 2/3$$

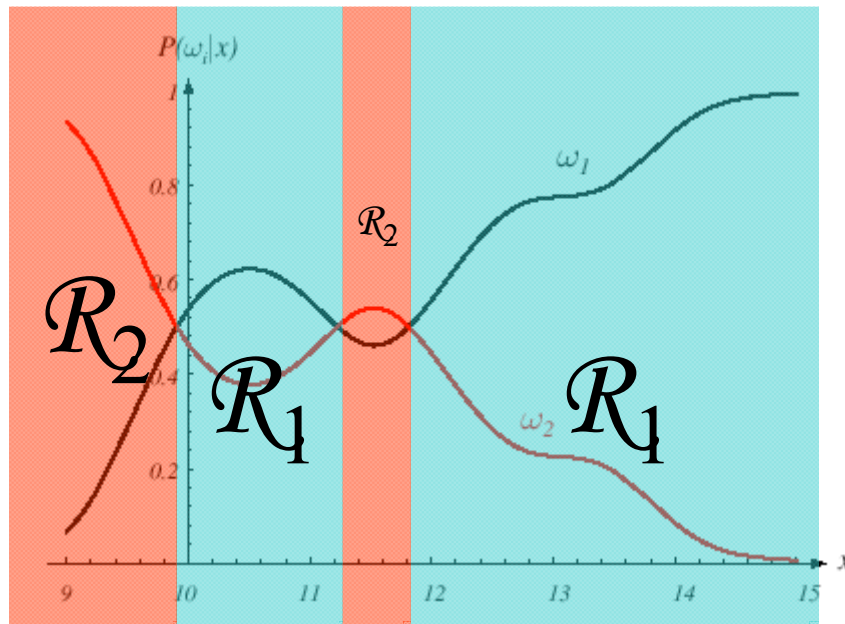
$$P(\omega_2) = 1/3$$

What will the posterior probability for either type of fish look like?

class-conditional pdf for *lightness*

Decide ω_1 if $p(\mathbf{x}|\omega_1)P(\omega_1) > p(\mathbf{x}|\omega_2)P(\omega_2)$; otherwise decide ω_2

An example



posterior probability for either type of fish

h-axis: lightness of fish scales
v-axis: posterior probability for each type of fish

Black curve: sea bass

Red curve: salmon

➤ For each value of x , the higher curve yields the output of Bayesian decision

➤ For each value of x , the posteriors of either curve sum to 1.0

Another Example

■ Problem statement

- A new medical test is used to detect whether a patient has a certain cancer or not, whose test result is either + (positive) or – (negative)
- For patient with this cancer, the probability of returning positive test result is 0.98
- For patient without this cancer, the probability of returning negative test result is 0.97
- The probability for any person to have this cancer is 0.008

■ Question

- If positive test result is returned, does she/he have cancer?

Another Example (Cont.)

ω_1 : cancer ω_2 : no cancer $x \in \{+, -\}$

$$P(\omega_1) = 0.008 \quad P(\omega_2) = 1 - P(\omega_1) = 0.992$$

$$P(+ | \omega_1) = 0.98 \quad P(- | \omega_1) = 1 - P(+ | \omega_1) = 0.02$$

$$P(- | \omega_2) = 0.97 \quad P(+ | \omega_2) = 1 - P(- | \omega_2) = 0.03$$

$$\begin{aligned} P(\omega_1 | +) &= \frac{P(\omega_1)P(+ | \omega_1)}{P(+)} = \frac{P(\omega_1)P(+ | \omega_1)}{P(\omega_1)P(+ | \omega_1) + P(\omega_2)P(+ | \omega_2)} \\ &= \frac{0.008 \times 0.98}{0.008 \times 0.98 + 0.992 \times 0.03} = 0.2085 \end{aligned}$$

$$P(\omega_2 | +) = 1 - P(\omega_1 | +) = 0.7915$$

$P(\omega_2 | +) > P(\omega_1 | +)$
No cancer!

Feasibility of Bayes Formula

MIMA

$$P(\omega_i | \mathbf{x}) = \frac{p(\mathbf{x} | \omega_i)P(\omega_i)}{p(\mathbf{x})} \left(\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \right)$$

- To compute posterior probability, we need to know **prior probability** and **likelihood**

How do we know these probabilities ?

- A simple solution: Counting Relative frequencies
- An advanced solution: Conduct Density estimation

A Further Example

■ Problem

- Based on the height of a car in some campus, decide whether it costs more than \$50,000 or not

ω_1 : price $>$ \$ 50,000

ω_2 : price \leq \$ 50,000

x : height of a car

Decide ω_1 if $P(\omega_1|\mathbf{x}) > P(\omega_2|\mathbf{x})$;
otherwise decide ω_2

Quantities to know:

$P(\omega_1)$ $P(\omega_2)$ $P(x|\omega_1)$ $P(x|\omega_2)$

How to get them?



Counting relative
frequencies via
collected samples

A Further Example (Cont.)

- Collecting samples
 - Suppose we have randomly picked 1209 cars in the campus, got prices from their owners, and measured their heights
- Compute $P(\omega_1)$ and $P(\omega_2)$

cars in ω_1 : 221

cars in ω_2 : 988

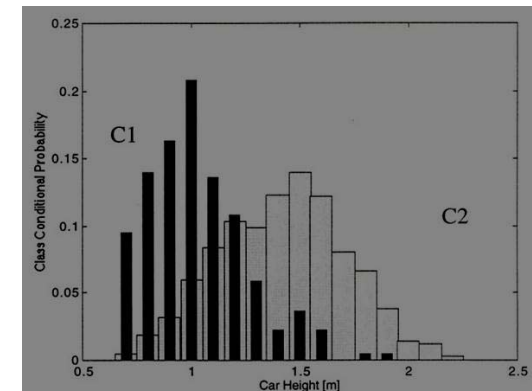
$$P(\omega_1) = \frac{221}{1209} = 0.183$$

$$P(\omega_2) = \frac{988}{1209} = 0.817$$

A Further Example (Cont.)

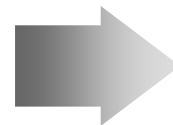
- Compute $P(x|\omega_1)$ $P(x|\omega_2)$
 - Discretize the height spectrum (say [0.5m, 2.5m]) into 20 intervals each with length 0.1m, and then count the number of cars falling into each interval for either class

- Suppose $x = 1.05$, which means that x falls into interval $I_x = [1.0m, 1.1m]$



For ω_1 , # cars in I_x is 46,

For ω_2 , # cars in I_x is 59,



$$P(x = 1.05 | \omega_1) = \frac{46}{221} = 0.2081$$

$$P(x = 1.05 | \omega_2) = \frac{59}{988} = 0.0597$$

A Further Example (Cont.)

■ Question

- For a car with height 1.05m, is its price greater than \$50,000?

$$P(\omega_1) = \frac{221}{1209} = 0.183$$

$$P(\omega_2) = \frac{988}{1209} = 0.817$$

$$P(x = 1.05 | \omega_1) = \frac{46}{221} = 0.2081$$

$$P(x = 1.05 | \omega_2) = \frac{59}{988} = 0.0597$$

$$\frac{P(\omega_2 | x = 1.05)}{P(\omega_1 | x = 1.05)} = \frac{P(\omega_2)P(x = 1.05 | \omega_2)}{P(\omega_1)P(x = 1.05 | \omega_1)}$$

$$= \frac{P(\omega_2)P(x = 1.05 | \omega_2)}{P(\omega_1)P(x = 1.05 | \omega_1)} = \frac{0.817 \times 0.0597}{0.183 \times 0.2081} \Rightarrow P(\omega_1 | \mathbf{x}) < P(\omega_2 | \mathbf{x}),$$

price \leq \$50,000

Is Bayes Decision Rule Optimal

MIMA

- Consider two categories

Decide ω_1 if $P(\omega_1|\mathbf{x}) > P(\omega_2|\mathbf{x})$; otherwise decide ω_2

- When we observe \mathbf{x} , the probability of error is:

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

Thus, under Bayes decision rule, we have

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

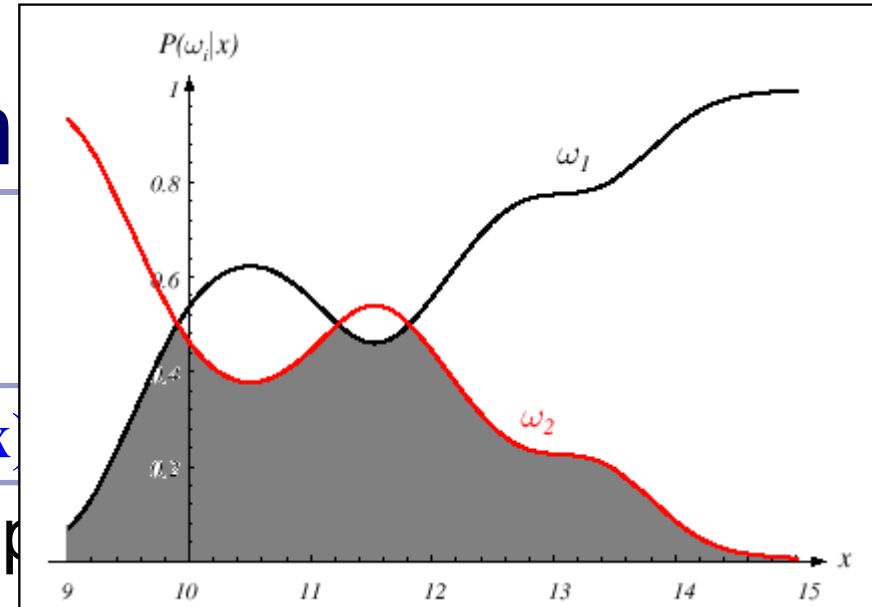
For every \mathbf{x} , we ensure that $P(\text{error}|\mathbf{x})$ is as small as possible

Is Bayes Decision

- Consider two categories

Decide ω_1 if $P(\omega_1|\mathbf{x}) > P(\omega_2|\mathbf{x})$

- When we observe x , the p



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

Thus, under Bayes decision rule, we have

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

For every x , we ensure that $P(\text{error}|x)$ is as small as possible

Generalized Bayes Decision Rule

- Allowing to use more than one feature

$x \in R \Rightarrow x \in R^d$: d-dimensional Euclidean Space

- Allowing more than two states of nature

$\Omega = \{\omega_1, \omega_2, \dots, \omega_c\}$: a set of c states of nature

- Allowing actions other than merely deciding the state of nature

$A = \{\alpha_1, \alpha_2, \dots, \alpha_a\}$: a set of a possible *actions*

Note that $c \neq a$

Generalized Bayes Decision Rule (cont.)

MIMA

- Introducing a loss function more general than the probability of error

$$\lambda : \Omega \times A \rightarrow R \text{ (loss function)}$$

$\lambda_{ij} = \lambda(\omega_j, \alpha_i)$: the loss incurred for taking action α_i
when the state of nature is ω_j

For ease of reference, it
is usually written as:

$$\lambda_{ij} = \lambda(\alpha_i | \omega_j):$$

We want to *minimize the expected loss* in making decision.

Risk

Generalized Bayes Decision Rule (cont.)

- Introduce the probability loss function λ

A simple loss function

Action \ Class	$\alpha_1 =$ "Recipe A"	$\alpha_2 =$ "Recipe B"	$\alpha_3 =$ "No Recipe"
$\omega_1 =$ "cancer"	5	50	10,000
$\omega_2 =$ "no cancer"	60	3	0

than

$\lambda_{ij} = \lambda(\omega_j, \alpha_i)$: the loss incurred for taking action α_i when the state of nature is ω_j

For ease of reference, it is usually written as:

$$\lambda_{ij} = \lambda(\alpha_i | \omega_j):$$

We want to *minimize the expected loss* in making decision.

Risk

Generalized Bayes Decision Rule (cont.)

A simple loss function

MIMA

Action \ Class	$\alpha_1 =$ "Recipe A"	$\alpha_2 =$ "Recipe B"	$\alpha_3 =$ "No Recipe"
$\omega_1 =$ "cancer"	5	50	10,000
$\omega_2 =$ "no cancer"	60	3	0

■ Problem

- Given a particular x , we have to decide which action to take
- To do this, we need to know the loss of taking each action α_i ($1 \leq i \leq a$)

$$\lambda_{ij} = \lambda(\alpha_i | \omega_j):$$

The action
being taken α_i

True state of
nature ω_j

However, the true state
of nature is uncertain

Expected (average)
loss

We want to *minimize the expected loss* in making decision.

Risk

Generalized Bayes Decision Rule (cont.)

MIMA

Given \mathbf{x} , the expected loss (risk) associated with taking action

α_i

- Expected loss

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

The incurred loss of taking action α_i in case of true state of nature being ω_j

The probability of ω_j being the true state of nature

The expected loss is also named as "conditional risk"

Generalized Bayes Decision Rule (cont.)

MIMA

- Suppose we have:

Action \ Class	$\alpha_1 =$ "Recipe A"	$\alpha_2 =$ "Recipe B"	$\alpha_3 =$ "No Recipe"
$\omega_1 =$ "cancer"	5	50	10,000
$\omega_2 =$ "no cancer"	60	3	0

For a particular \mathbf{x} :

$$P(\omega_1/\mathbf{x}) = 0.01$$

$$P(\omega_2/\mathbf{x}) = 0.99$$

$$\begin{aligned} R(\alpha_1 | \mathbf{x}) &= \sum_{j=1}^2 \lambda(\alpha_1 | \omega_j) \cdot P(\omega_j | \mathbf{x}) \\ &= \lambda(\alpha_1 | \omega_1) \cdot P(\omega_1 | \mathbf{x}) + \lambda(\alpha_1 | \omega_2) \cdot P(\omega_2 | \mathbf{x}) \\ &= 5 \times 0.01 + 60 \times 0.99 = 59.45 \end{aligned}$$

Similarly, we can get: $R(\alpha_2 | \mathbf{x}) = 3.47$ $R(\alpha_3 | \mathbf{x}) = 100$

Generalized Bayes Decision Rule (cont.)

MIMA

- 0/1 Loss Function

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

$$\lambda(\alpha_i | \omega_j) = \begin{cases} 0 & \alpha_i \text{ is a correct decision associated with } \omega_j \\ 1 & \text{otherwise} \end{cases}$$

➔ $R(\alpha_i | \mathbf{x}) = P(\text{error} | \mathbf{x})$

Generalized Bayes Decision Rule (cont.)

MIMA

- Bayes decision rule (general case)

$$\alpha(\mathbf{x}) = \arg \min_{\alpha_i \in A} R(\alpha_i | \mathbf{x}) = \arg \min_{\alpha_i \in A} \sum_{j=1}^c \lambda(\alpha_i | \omega_j) P(\omega_j | \mathbf{x})$$

- Overall risk

$$R = \int R(\underbrace{\alpha(\mathbf{x})}_{\text{Decision function}} | x) \cdot p(x) dx$$

Decision function

For every x , we ensure that the conditional risk $R(a(x)|x)$ is as small as possible; Thus, the overall risk over all possible x must be as small as possible.

*The **optimal** one to minimize the overall risk*

*Its resulting overall risk is called the **Bayesian risk***

General Case: Two-Category

$\Omega = \{\omega_1, \omega_2\}$	<i>Loss Function</i>	<i>State of Nature</i>	
		ω_1	ω_2
$A = \{\alpha_1, \alpha_2\}$	<i>Action</i>	α_1	α_2
		λ_{11}	λ_{12}
		λ_{21}	λ_{22}

$$R(\alpha_1 | \mathbf{x}) = \lambda_{11}P(\omega_1 | \mathbf{x}) + \lambda_{12}P(\omega_2 | \mathbf{x})$$

$$R(\alpha_2 | \mathbf{x}) = \lambda_{21}P(\omega_1 | \mathbf{x}) + \lambda_{22}P(\omega_2 | \mathbf{x})$$

General Case: Two-Category

Perform α_1 if $R(\alpha_2|\mathbf{x}) > R(\alpha_1|\mathbf{x})$; otherwise perform α_2

$$\longrightarrow \lambda_{21}P(\omega_1 | \mathbf{x}) + \lambda_{22}P(\omega_2 | \mathbf{x}) > \lambda_{11}P(\omega_1 | \mathbf{x}) + \lambda_{12}P(\omega_2 | \mathbf{x})$$

$$\longrightarrow (\lambda_{21} - \lambda_{11})P(\omega_1 | \mathbf{x}) > (\lambda_{12} - \lambda_{22})P(\omega_2 | \mathbf{x})$$

$$R(\alpha_1 | \mathbf{x}) = \lambda_{11}P(\omega_1 | \mathbf{x}) + \lambda_{12}P(\omega_2 | \mathbf{x})$$

$$R(\alpha_2 | \mathbf{x}) = \lambda_{21}P(\omega_1 | \mathbf{x}) + \lambda_{22}P(\omega_2 | \mathbf{x})$$

General Case: Two-Category

MIMA

Perform α_1 if $R(\alpha_2|\mathbf{x}) > R(\alpha_1|\mathbf{x})$; otherwise perform α_2

$$\longrightarrow \lambda_{21}P(\omega_1 | \mathbf{x}) + \lambda_{22}P(\omega_2 | \mathbf{x}) > \lambda_{11}P(\omega_1 | \mathbf{x}) + \lambda_{12}P(\omega_2 | \mathbf{x})$$

$$\longrightarrow \underbrace{(\lambda_{21} - \lambda_{11})}_{\text{positive}}P(\omega_1 | \mathbf{x}) > \underbrace{(\lambda_{12} - \lambda_{22})}_{\text{positive}}P(\omega_2 | \mathbf{x})$$

*Posterior probabilities are **scaled** before comparison.*

General Case: Two-Category

Perform α_1 if $R(\alpha_2|\mathbf{x}) > R(\alpha_1|\mathbf{x})$; otherwise perform α_2

➔ $\lambda_{21}P(\omega_1 | \mathbf{x}) + \lambda_{22}P(\omega_2 | \mathbf{x}) > \lambda_{11}P(\omega_1 | \mathbf{x}) + \lambda_{12}P(\omega_2 | \mathbf{x})$

➔ $(\lambda_{21} - \lambda_{11})P(\omega_1 | \mathbf{x}) > (\lambda_{12} - \lambda_{22})P(\omega_2 | \mathbf{x})$

➔ $(\lambda_{21} - \lambda_{11})p(\mathbf{x} | \omega_1)P(\omega_1) > (\lambda_{12} - \lambda_{22})p(\mathbf{x} | \omega_2)P(\omega_2)$

➔
$$\frac{p(\mathbf{x} | \omega_1)}{p(\mathbf{x} | \omega_2)} > \frac{(\lambda_{12} - \lambda_{22}) P(\omega_2)}{(\lambda_{21} - \lambda_{11}) P(\omega_1)}$$

General Case: Two-Category

MIMA

→ Perform α_1 if

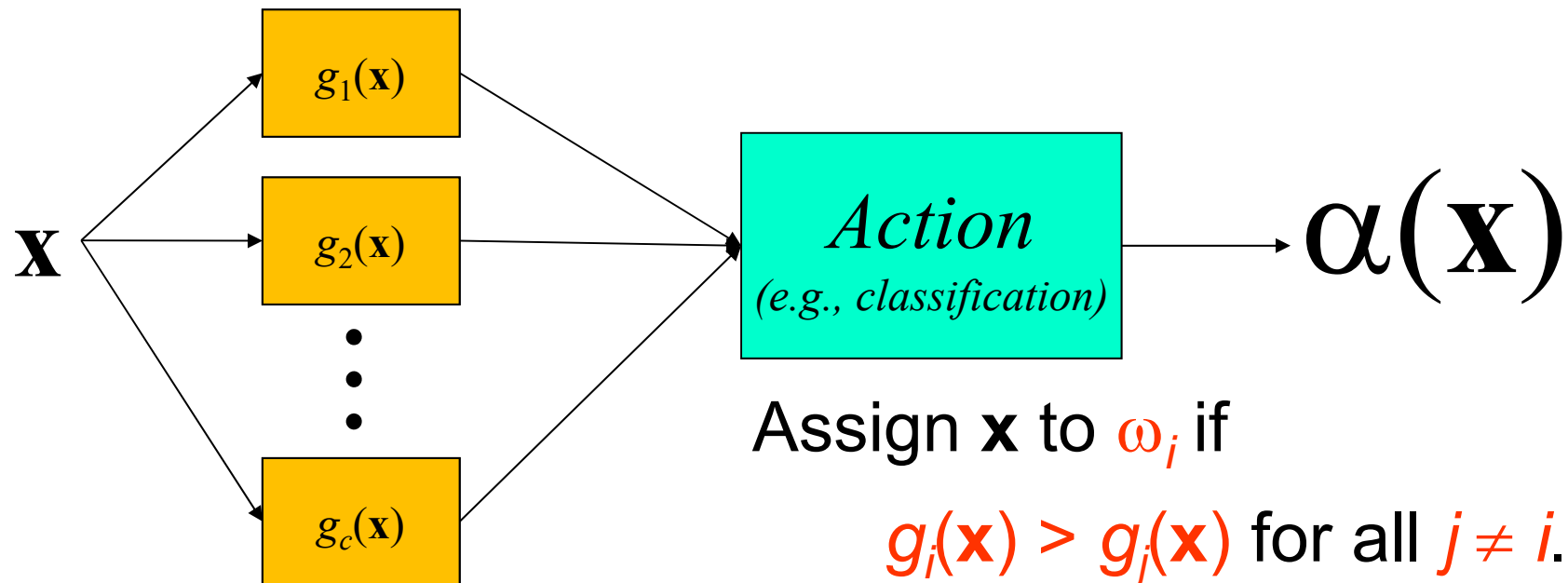
$$\underbrace{\frac{p(\mathbf{x} | \omega_1)}{p(\mathbf{x} | \omega_2)}}_{\text{Likelihood Ratio}} > \underbrace{\frac{(\lambda_{12} - \lambda_{22}) P(\omega_2)}{(\lambda_{21} - \lambda_{11}) P(\omega_1)}}_{\text{Threshold}}$$

Discriminant Function

- Discriminant functions for multiclass

$$g_i(x) : R^d \rightarrow R \quad (1 \leq i \leq c)$$

- One function per category



Discriminant Function

- Minimum Risk Case:

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

- Minimum Error-Rate Case:

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

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

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

Discriminant Function

- Relationship between minimum risk and minimum error rate

$$\begin{aligned} R(\alpha_i | \mathbf{x}) &= \sum_{j=1}^c \lambda(\alpha_i | \omega_j) \cdot P(\omega_j | \mathbf{x}) \\ &= \sum_{j \neq i} \lambda(\alpha_i | \omega_j) \cdot P(\omega_j | \mathbf{x}) + \lambda(\alpha_i | \omega_i) \cdot P(\omega_i | \mathbf{x}) \\ &= \sum_{j \neq i} P(\omega_j | \mathbf{x}) \\ &= 1 - P(\omega_i | \mathbf{x}) \end{aligned}$$

error rate (误差率/错误率)
the probability that action α_i (decide ω_i) is wrong

Discriminant Function

- Various discriminant function
- Identical classification results

If $f(\cdot)$ is a *monotonically increasing* function, then $f(g_i(\cdot))$'s are also be discriminant functions.

- Example

$$f(x) = k \cdot x \quad (k > 0) \quad \Rightarrow \quad f(g_i(x)) = k \cdot g_i(x) \quad (1 \leq i \leq c)$$

$$f(x) = \ln x \quad \Rightarrow \quad f(g_i(x)) = \ln g_i(x) \quad (1 \leq i \leq c)$$

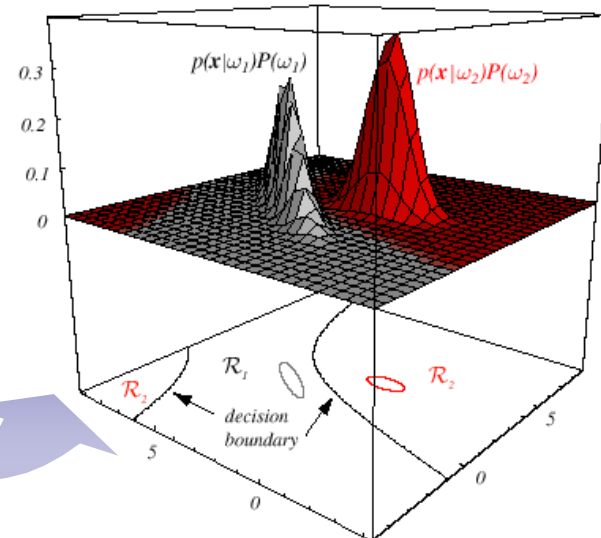
Decision Regions

- c discriminant functions result in c decision regions.

$$\mathcal{R}_i = \{\mathbf{x} \mid g_i(\mathbf{x}) > g_j(\mathbf{x}) \quad \forall j \neq i\}$$

where $\mathcal{R}_i \cap \mathcal{R}_j = \phi$ ($i \neq j$) and $\bigcup_{i=1}^c \mathcal{R}_i = \mathcal{R}^d$

- Decision boundary
 - Decision regions are separated by decision boundaries



Two-category example

The Normal Distribution

Discrete random variable (X) — Assume integer

Probability mass function (pmf): $p(x) = P(X = x)$

Cumulative distribution function (cdf): $F(x) = P(X \leq x) = \sum_{t=-\infty}^x p(t)$

Continuous random variable (X)

Probability density function (pdf): $p(x)$ or $f(x)$ **not a probability**

Cumulative distribution function (cdf): $F(x) = P(X \leq x) = \int_{-\infty}^x p(t)dt$

Expectations

- a.k.a. **expected value, mean or average of a random variable**
- x is a random variable, the expectation of x

$$E[x] = \begin{cases} \sum_{x=-\infty}^{\infty} xp(x) & x \text{ is discrete} \\ \int_{-\infty}^{\infty} xp(x)dx & x \text{ is continuous} \end{cases}$$

The k^{th} moment $E[X^k]$

The 1^{st} moment $\mu_X = E[X]$

The k^{th} central moment $E[(X - \mu_X)^k]$

Important Expectations

■ Mean

$$\mu_X = E[X] = \begin{cases} \sum_{x=-\infty}^{\infty} xp(x) & X \text{ is discrete} \\ \int_{-\infty}^{\infty} xp(x)dx & X \text{ is continuous} \end{cases}$$

■ Variance

$$\sigma_X^2 = Var[X] = E[(X - \mu_X)^2] = \begin{cases} \sum_{x=-\infty}^{\infty} (x - \mu_X)^2 p(x) & X \text{ is discrete} \\ \int_{-\infty}^{\infty} (x - \mu_X)^2 p(x)dx & X \text{ is continuous} \end{cases}$$

Notation: $\sigma^2 = Var[x]$ (σ : standard deviation ?)

Fact: $\sigma^2 = Var[x] = E[x^2] - (E[x])^2$

Entropy

- The *entropy* measures the fundamental *uncertainty* in the value of points selected randomly from a distribution.

$$H[X] = \begin{cases} - \sum_{x=-\infty}^{\infty} p(x) \log p(x) & X \text{ is discrete} \\ - \int_{-\infty}^{\infty} p(x) \log p(x) dx & X \text{ is continuous} \end{cases}$$

Univariate Gaussian Distribution

MIMA

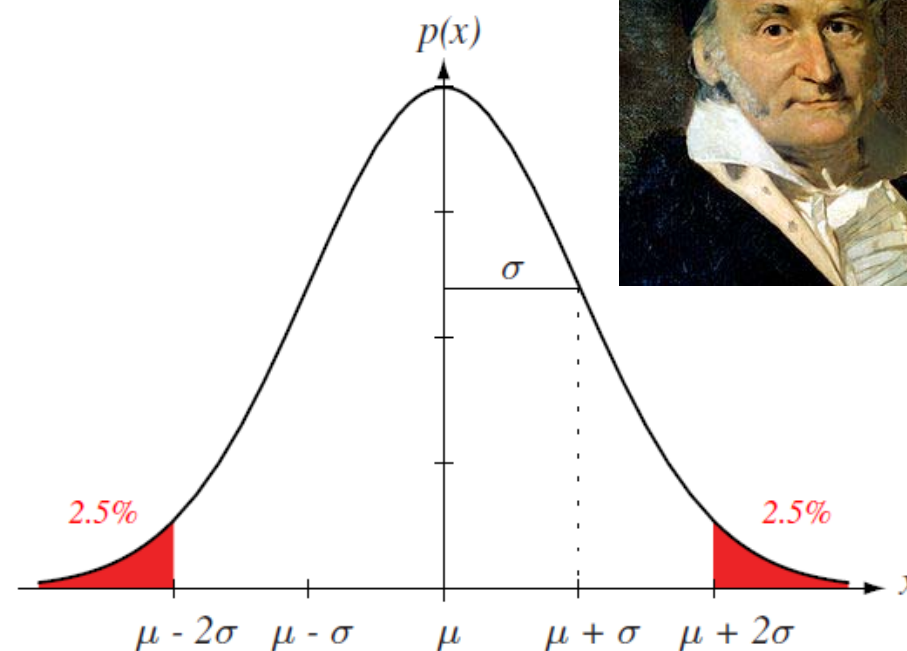
- Gaussian distribution, a.k.a. **Gaussian density**, **normal density**.

$$X \sim N(\mu, \sigma^2)$$

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$E[X] = \mu$$

$$\text{Var}[X] = \sigma^2$$



Univariate Gaussian Distribution

MIMA

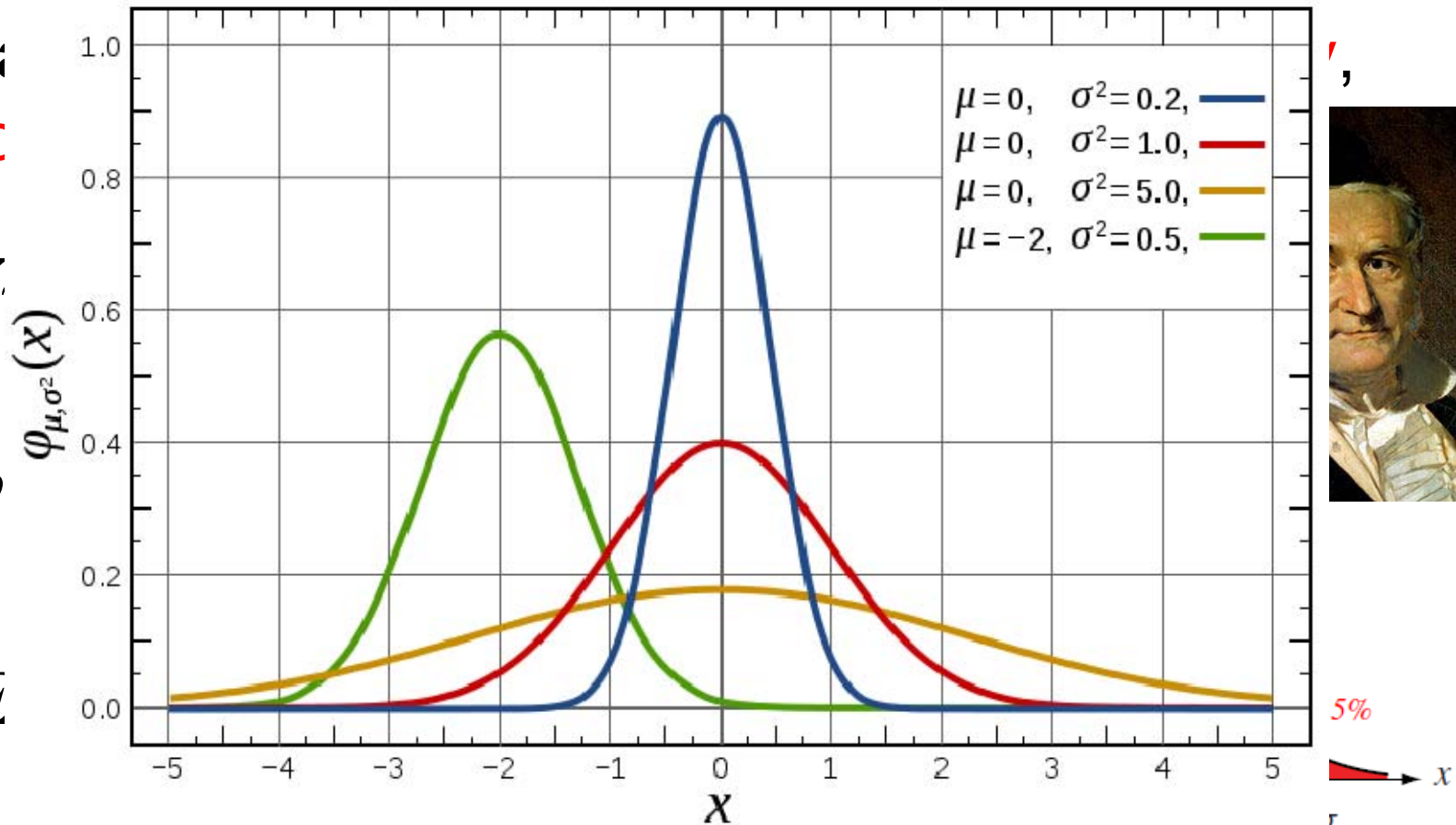
■ Gaussian

nc

X

p

E



$$\text{Var}[X] = \sigma^2$$

Random Vectors

- A d-dimensional random vector is:

$$\mathbf{X} = (x_1, x_2, \dots, x_d)^T \quad \mathbf{X} : \Omega \rightarrow R^d$$

$$X \sim p(X) = p(x_1, x_2, \dots, x_d) \quad (\text{joint pdf})$$

- Expected vector

$$E[\mathbf{X}] = \begin{pmatrix} E[x_1] \\ E[x_2] \\ \vdots \\ E[x_d] \end{pmatrix} \quad E[x_i] = \int_{-\infty}^{+\infty} x_i \underline{p(x_i)} dx_i \quad (1 \leq i \leq d)$$

Marginal pdf on the
ith component.

$$\boldsymbol{\mu} = E[\mathbf{X}] = (\mu_1, \mu_2, \dots, \mu_d)^T$$

Random Vectors

- Covariance matrix

$$\Sigma = E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \Lambda & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \Lambda & \sigma_{2d} \\ \mathbf{M} & \mathbf{M} & \mathbf{O} & \mathbf{M} \\ \sigma_{d1} & \sigma_{d2} & \Lambda & \sigma_d^2 \end{pmatrix}$$

$$\sigma_{ij} = \sigma_{ji} = E[(x_i - \mu_i)(x_j - \mu_j)]$$

$$= \int_{-\infty}^{+\infty} (x_i - \mu_i)(x_j - \mu_j) \underline{p(x_i, x_j)} dx_i dx_j$$

Properties:
Symmetric, Positive semidefinite

Marginal pdf on a pair of random variables (x_i, x_j)

Multivariate Gaussian Distribution

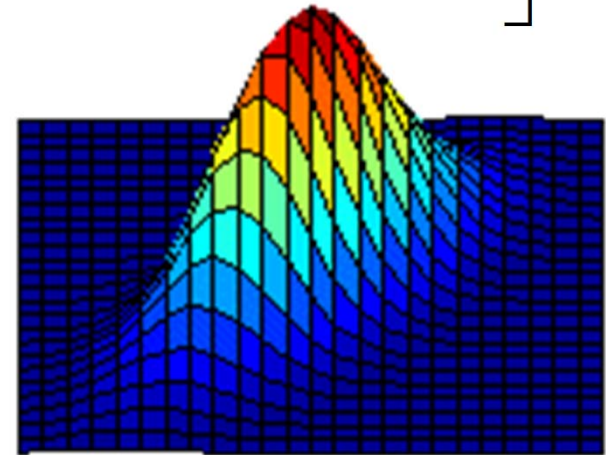
- X is a d -dimensional random vector

$$X \sim N(\mu, \Sigma)$$

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

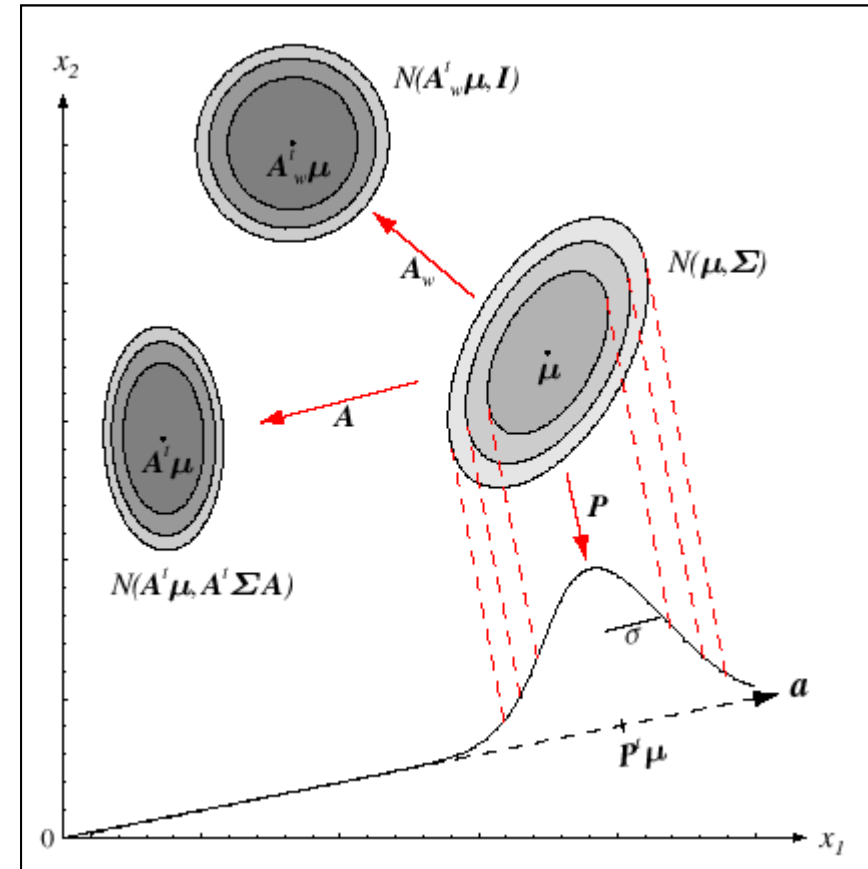
$$E[X] = \mu$$

$$E[(X - \mu)(X - \mu)^T] = \Sigma$$



Properties of $N(\mu, \Sigma)$

- X is a d -dimensional random vector, and
 $X \sim N(\mu, \Sigma)$
- If $Y = AX$, where A is a $d \times k$ matrix, then



$$Y \sim N(A^T \mu, A^T \Sigma A)$$

On Covariance Matrix

- As mentioned before, Σ is *symmetric* and *positive semidefinite*.

$$\Sigma = \Phi \Lambda \Phi^T = \Phi \Lambda^{1/2} \Lambda^{1/2} \Phi^T$$

Φ : *orthonormal* matrix, whose columns are *eigenvectors* of Σ .

Λ : *diagonal* matrix (*eigenvalues*).

- Thus,

$$\Sigma = (\Phi \Lambda^{1/2})(\Phi \Lambda^{1/2})^T$$

Mahalanobis Distance

- Mahalanobis distance

$$r^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

$$X \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$



P.C. Mahalanobis
(1894-1972)

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

*depends on
the value of r^2*

constant

r^2

Discriminant Functions for Gaussian Density

- Minimum-error-rate classification

$$g_i(\mathbf{x}) = P(\omega_i | \mathbf{x}) \quad (1 \leq i \leq c)$$


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


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

$$p(\mathbf{x} | \omega_i) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_i|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)\right]$$

Constant, could be ignored

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

Discriminant Functions for Gaussian Density

■ Three cases

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

- Classes are centered at different mean, and their feature components are pairwise independent have the same variance.

- Case 2 $\Sigma_i = \Sigma$

- Classes are centered at different mean, but have the same variation.

- Case 3 $\Sigma_i \neq \Sigma_j$

- Arbitrary

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

irrelevant

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

$$g_i(\mathbf{x}) = -\frac{1}{2\sigma^2} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 + \ln P(\omega_i) \quad \Sigma_i^{-1} = \frac{1}{\sigma^2} \mathbf{I}$$

$$= -\frac{1}{2\sigma^2} (\underbrace{\mathbf{x}^T \mathbf{x}}_{\text{irrelevant}} - 2\boldsymbol{\mu}_i^T \mathbf{x} + \boldsymbol{\mu}_i^T \boldsymbol{\mu}_i) + \ln P(\omega_i)$$

irrelevant

$$g_i(\mathbf{x}) = \frac{1}{\sigma^2} \boldsymbol{\mu}_i^T \mathbf{x} + \left[-\frac{1}{2\sigma^2} \boldsymbol{\mu}_i^T \boldsymbol{\mu}_i + \ln P(\omega_i) \right]$$

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

$$g_i(\mathbf{x}) = \frac{1}{\sigma^2} \boldsymbol{\mu}_i^T \mathbf{x} + \left[-\frac{1}{2\sigma^2} \boldsymbol{\mu}_i^T \boldsymbol{\mu}_i + \ln P(\omega_i) \right]$$

- It is a linear discriminant function

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

- where

- Weight vector

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

- Threshold/bias

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

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

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

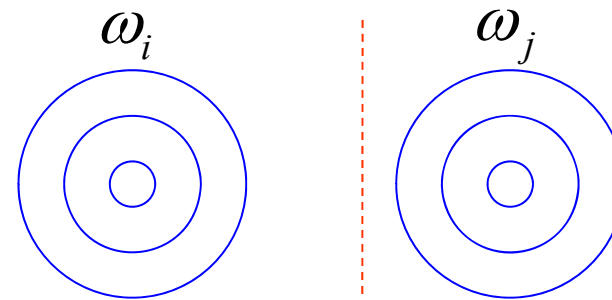
$$\mathbf{w}_i = \frac{1}{\sigma^2} \boldsymbol{\mu}_i$$

$$\mathbf{w}_i^T \mathbf{x} + w_{i0} = \mathbf{w}_j^T \mathbf{x} + w_{j0}$$

$$(\mathbf{w}_i^T - \mathbf{w}_j^T) \mathbf{x} = w_{j0} - w_{i0}$$

$$(\boldsymbol{\mu}_i^T - \boldsymbol{\mu}_j^T) \mathbf{x} = \frac{1}{2} (\boldsymbol{\mu}_i^T \boldsymbol{\mu}_i - \boldsymbol{\mu}_j^T \boldsymbol{\mu}_j) - \sigma^2 \ln \frac{P(\omega_i)}{P(\omega_j)}$$

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



$g_i(\mathbf{x}) = g_j(\mathbf{x})$
Boundary btw.

ω_i and ω_j

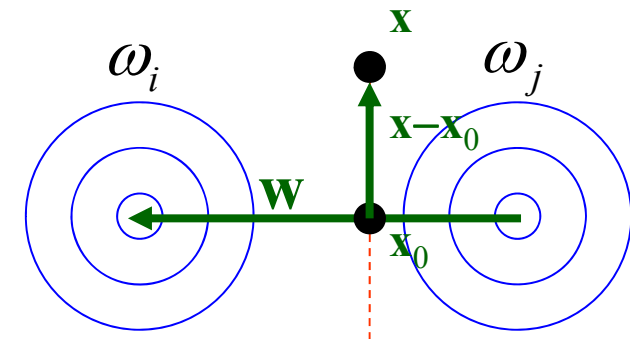
Case 1: $\Sigma_i = \sigma^2 \mathbf{I}$

- The decision boundary will be a **hyperplane perpendicular** to the line btw. the means at somewhere.

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

$$\mathbf{w} = \boldsymbol{\mu}_i - \boldsymbol{\mu}_j$$

$$\mathbf{x}_0 = \underbrace{\frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j)}_{\text{midpoint}} - \underbrace{\frac{\sigma^2}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2} \ln \frac{P(\omega_i)}{P(\omega_j)}}_{0 \text{ if } P(\omega_i)=P(\omega_j)} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)$$

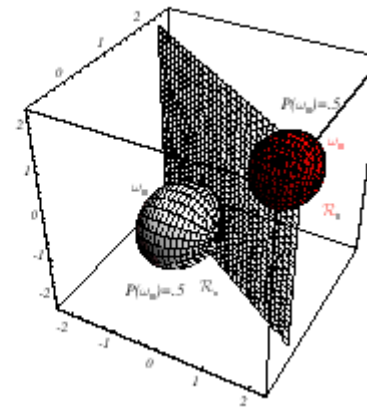
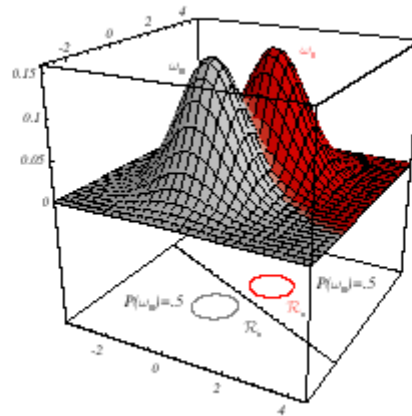
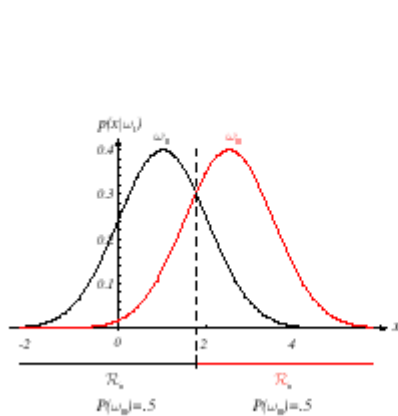


$g_i(\mathbf{x}) = g_j(\mathbf{x})$
Boundary btw.
 ω_i and ω_j

$$\overbrace{(\boldsymbol{\mu}_i^T - \boldsymbol{\mu}_j^T)}^{\mathbf{w}^T} \mathbf{x} = \frac{1}{2} (\boldsymbol{\mu}_i^T - \boldsymbol{\mu}_j^T) (\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \sigma^2 \frac{(\boldsymbol{\mu}_i^T - \boldsymbol{\mu}_j^T) (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2} \ln \frac{P(\omega_i)}{P(\omega_j)}$$

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

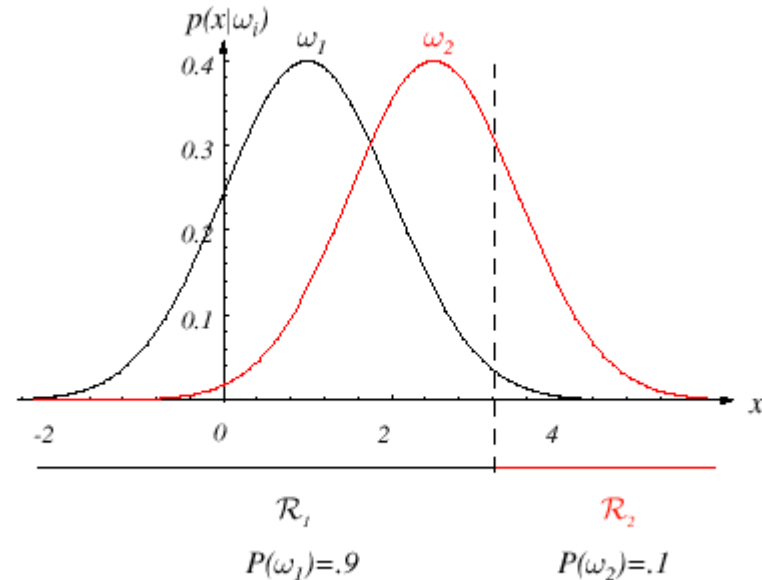
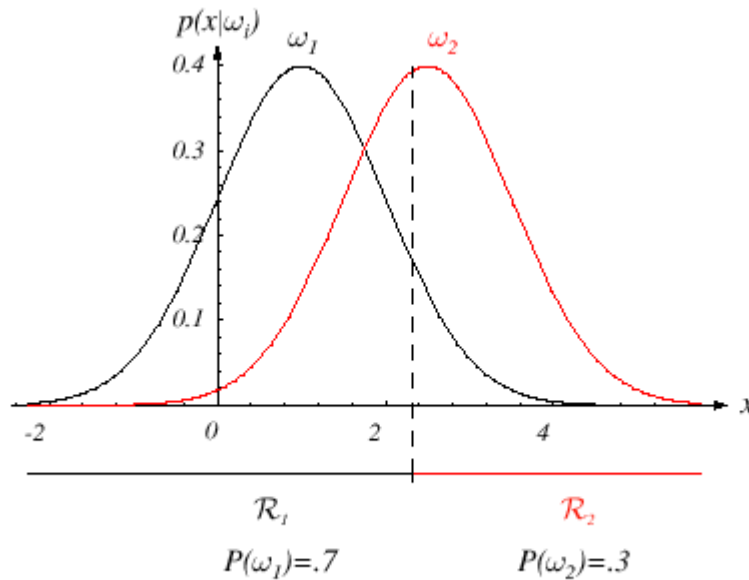
$$P(\omega_1) = P(\omega_2)$$



Minimum distance classifier (template matching)

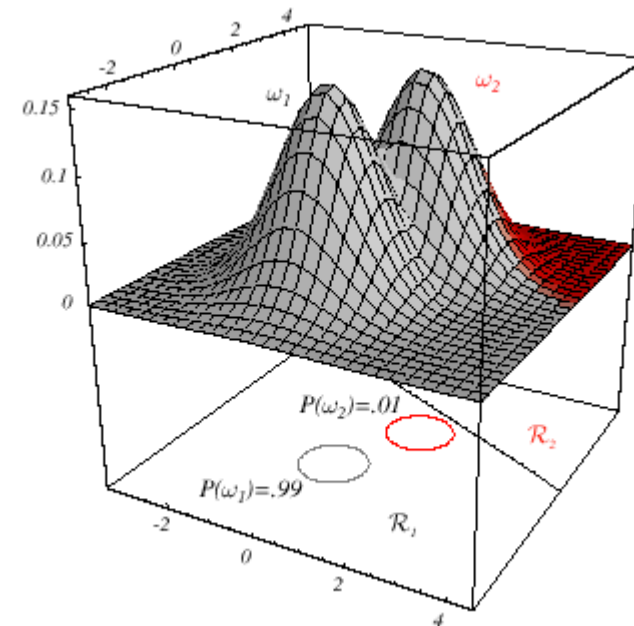
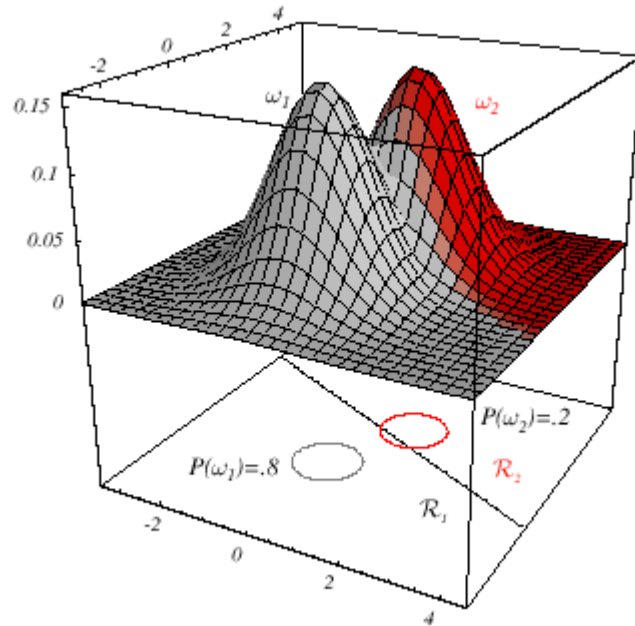
Case 1: $\Sigma_i = \sigma^2 \mathbf{I}$

$$P(\omega_1) > P(\omega_2)$$



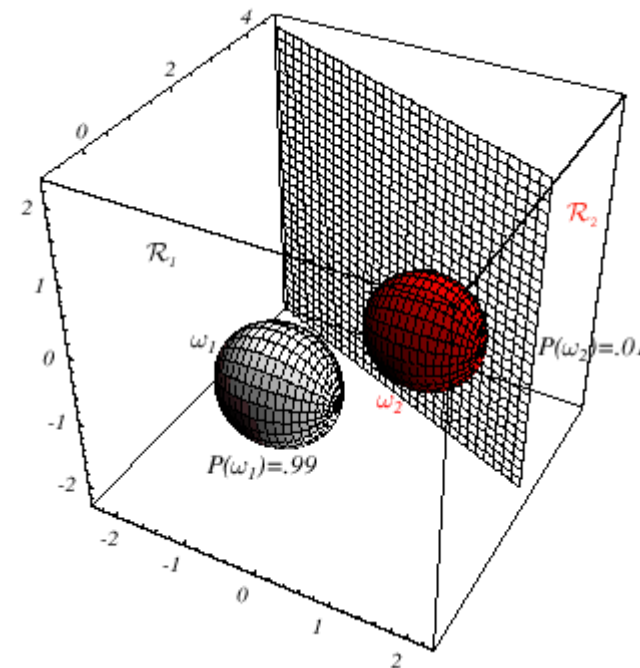
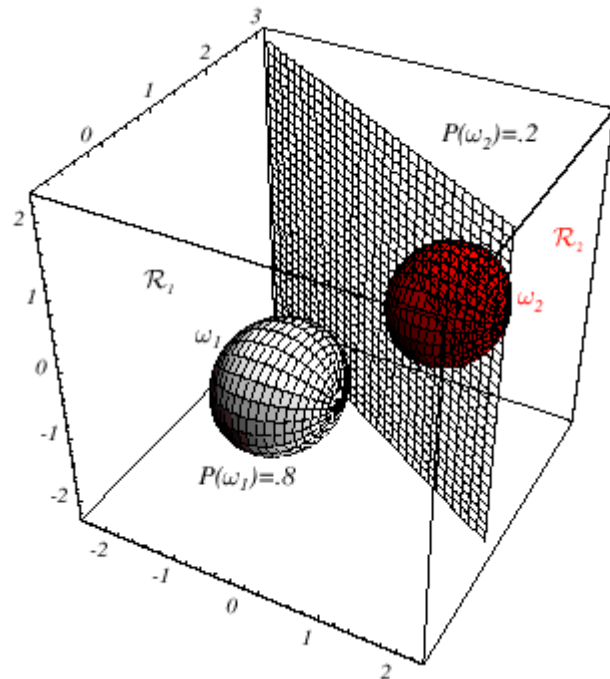
Case 1: $\Sigma_i = \sigma^2 \mathbf{I}$

$$P(\omega_1) > P(\omega_2)$$



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

$$P(\omega_1) > P(\omega_2)$$



Case 2: $\Sigma_i = \Sigma$

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

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

Mahalanobis Distance *Irrelevant if*
 $P(\omega_i) = P(\omega_j) \forall i, j$

$$= -\frac{1}{2}(\underbrace{\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}}_{\text{Irrelevant}} - 2\boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i) + \ln P(\omega_i)$$

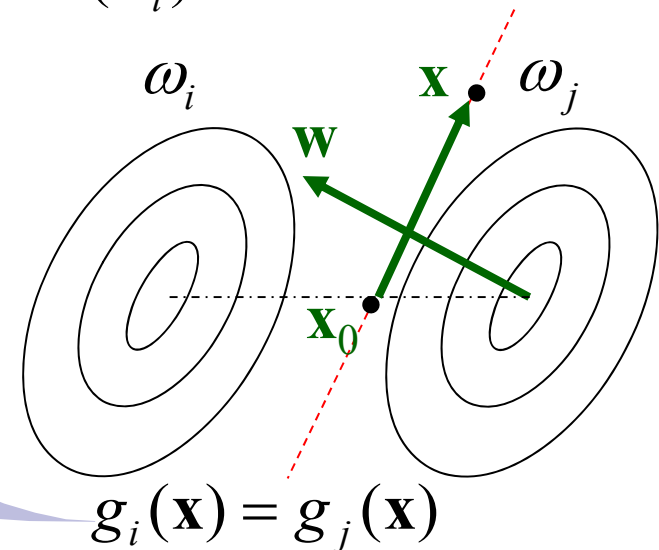
Irrelevant

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0} \quad \left\{ \begin{array}{l} \mathbf{w}_i = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i \\ w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \ln P(\omega_i) \end{array} \right.$$

Case 2: $\Sigma_i = \Sigma$

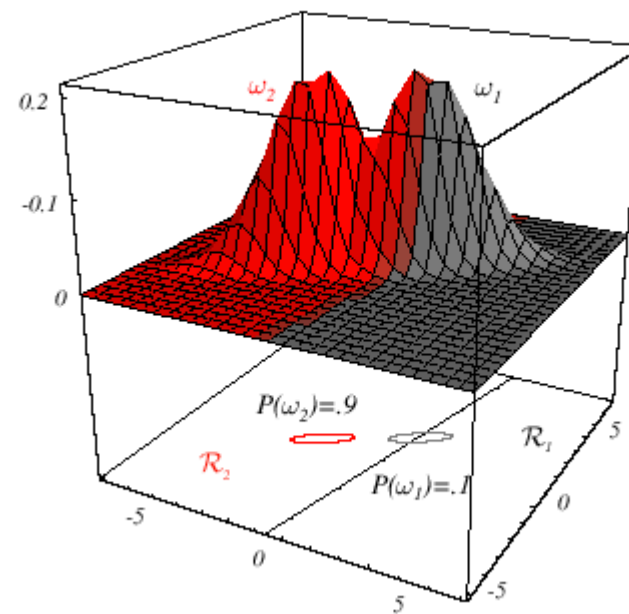
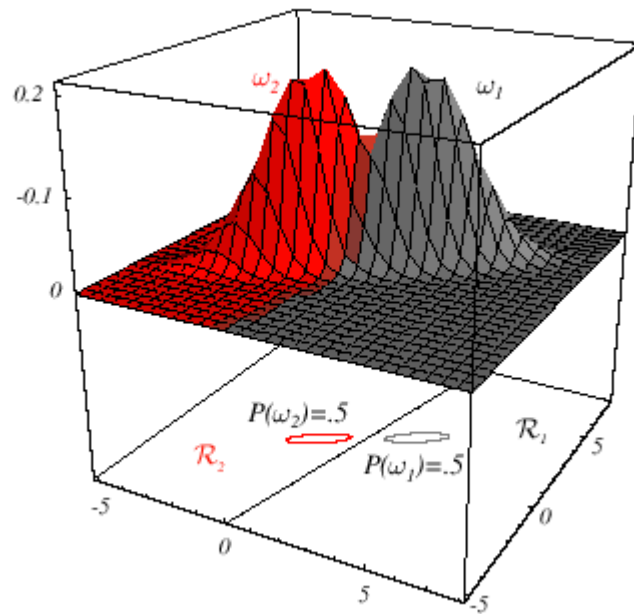
$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0} \begin{cases} \mathbf{w}_i = \Sigma^{-1} \boldsymbol{\mu}_i \\ w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i + \ln P(\omega_i) \end{cases}$$

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

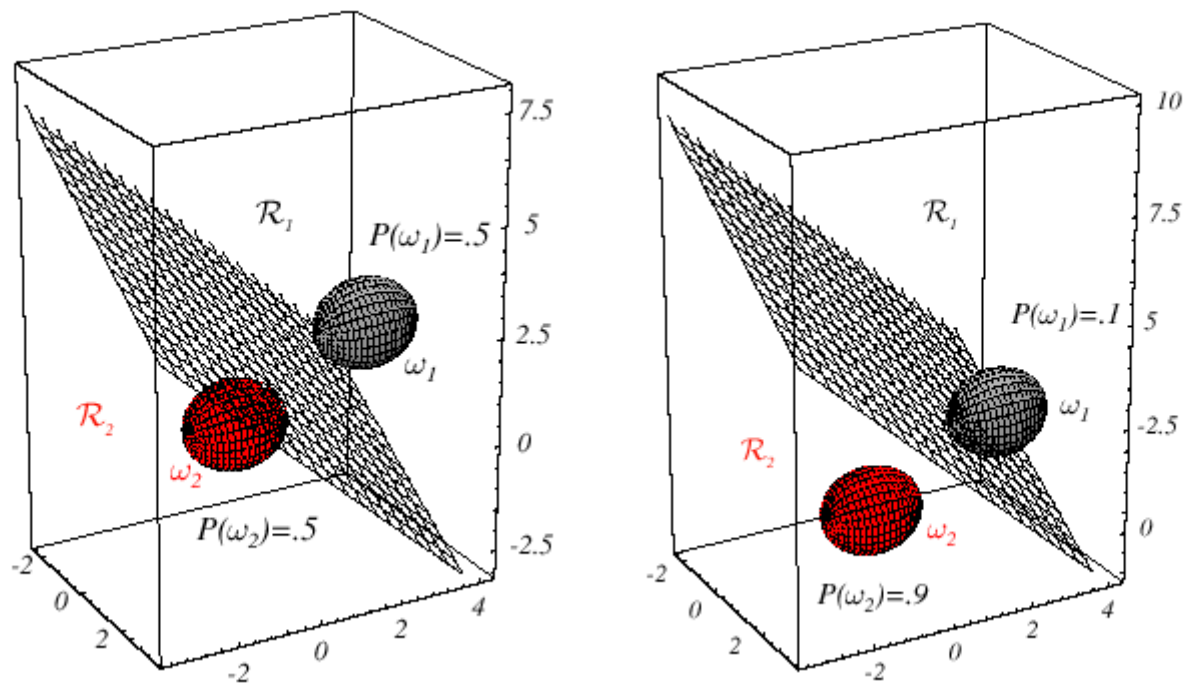


$$\begin{cases} \mathbf{w} = \Sigma^{-1} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) \\ \mathbf{x}_0 = \frac{1}{2} (\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \frac{\ln[P(\omega_i) / P(\omega_j)]}{(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^T \Sigma^{-1} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)} (\boldsymbol{\mu}_i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

Case 2: $\Sigma_i = \Sigma$



Case 2: $\Sigma_i = \Sigma$



Case 3: $\Sigma_i \neq \Sigma_j$

irrelevant

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

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

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

Without this term

In Case 1 and 2

$$\mathbf{W}_i = -\frac{1}{2} \boldsymbol{\Sigma}_i^{-1}$$

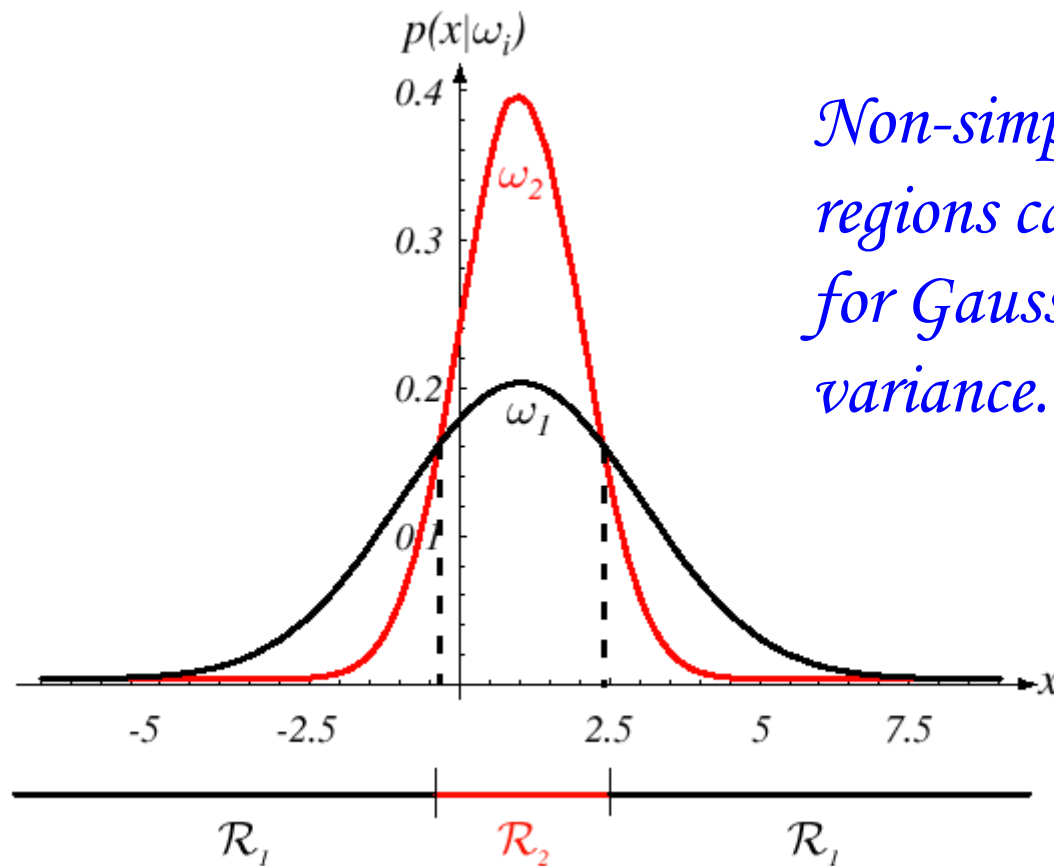
$$\mathbf{w}_i = \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i$$

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

Decision surfaces are hyperquadrics, e.g.,

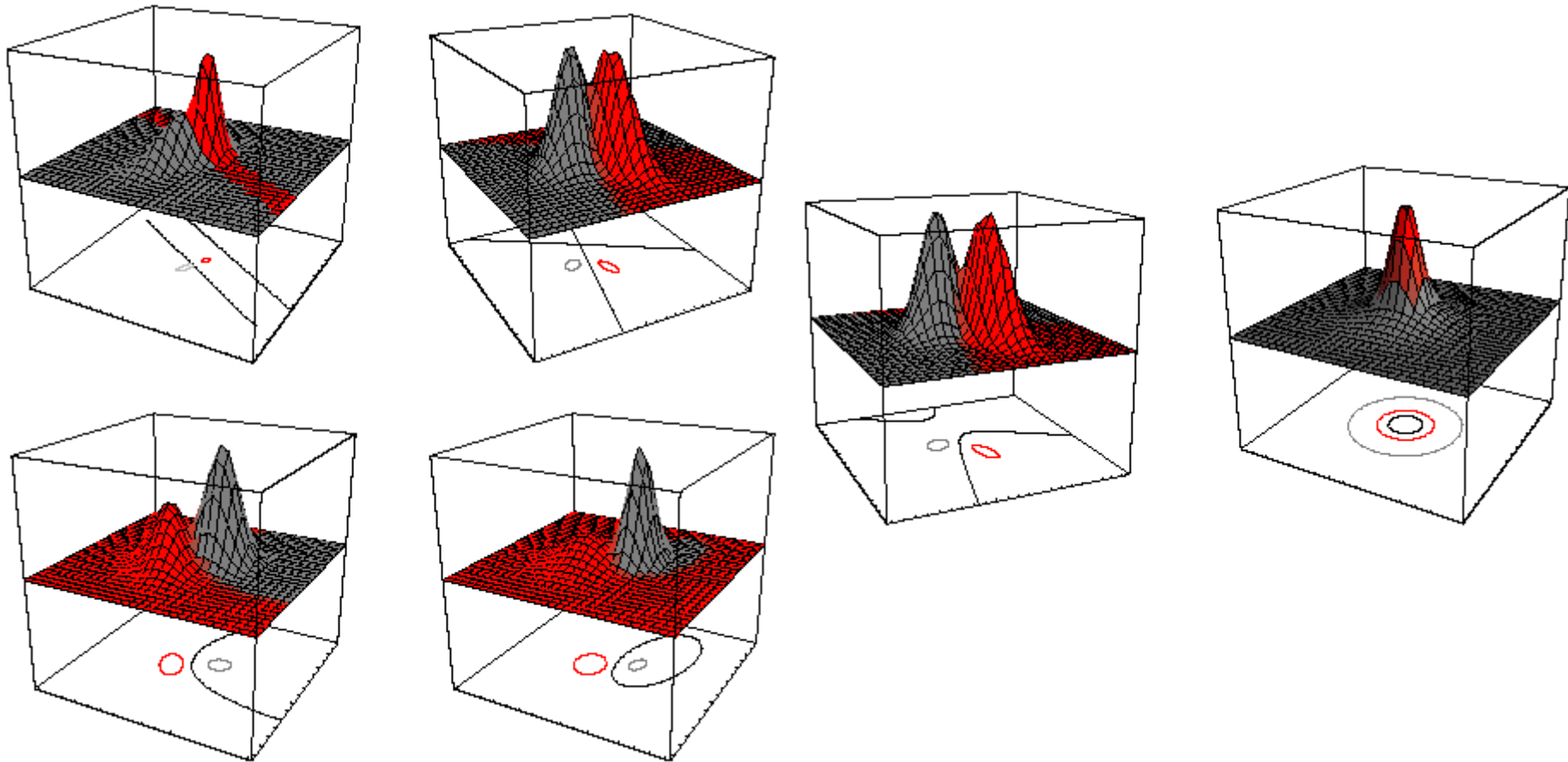
- Hyperplanes
- Hyperspheres
- Hyperellipsoids
- hyperhyperboloids

Case 3: $\Sigma_i \neq \Sigma_j$

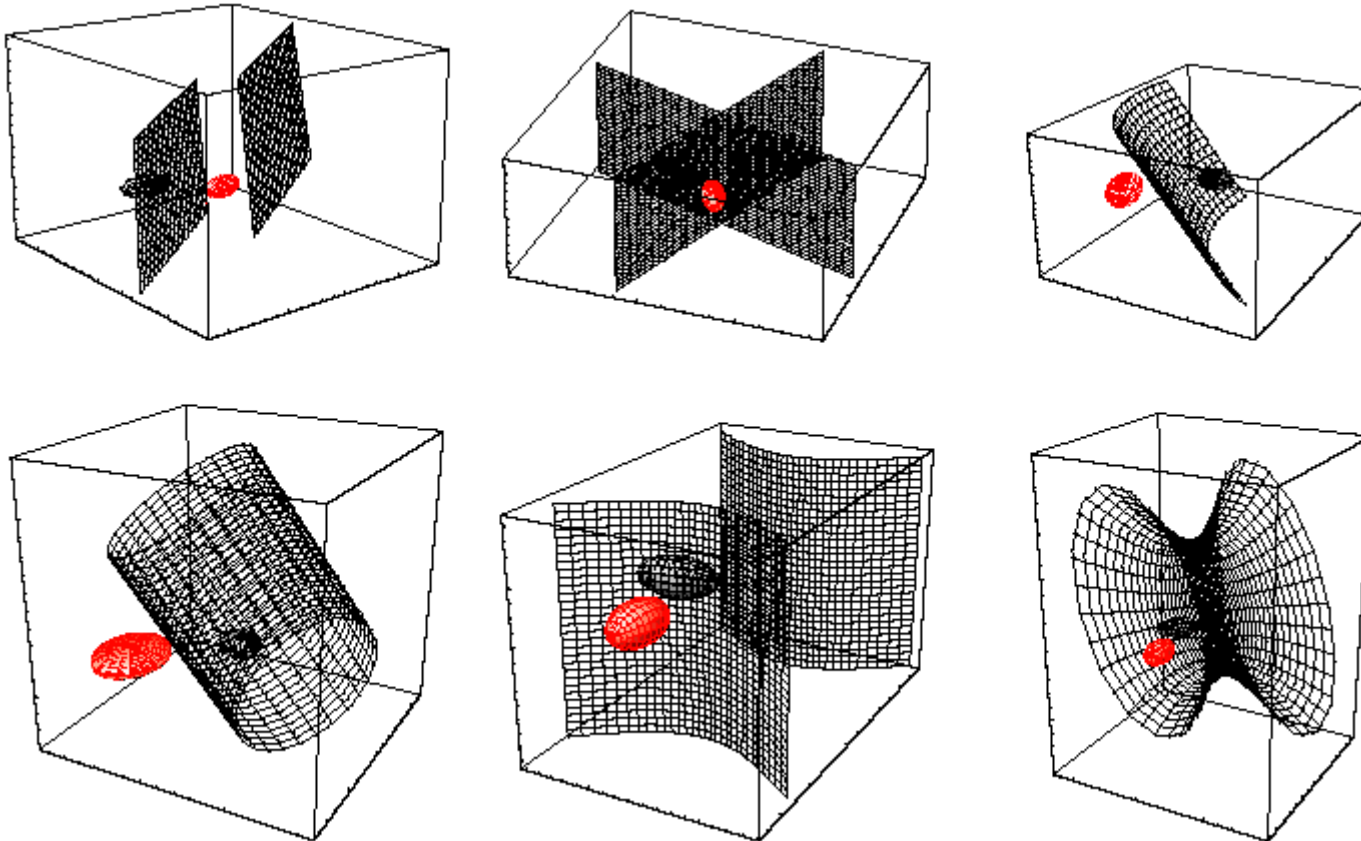


Non-simply connected decision regions can arise in one dimension for Gaussians having unequal variance.

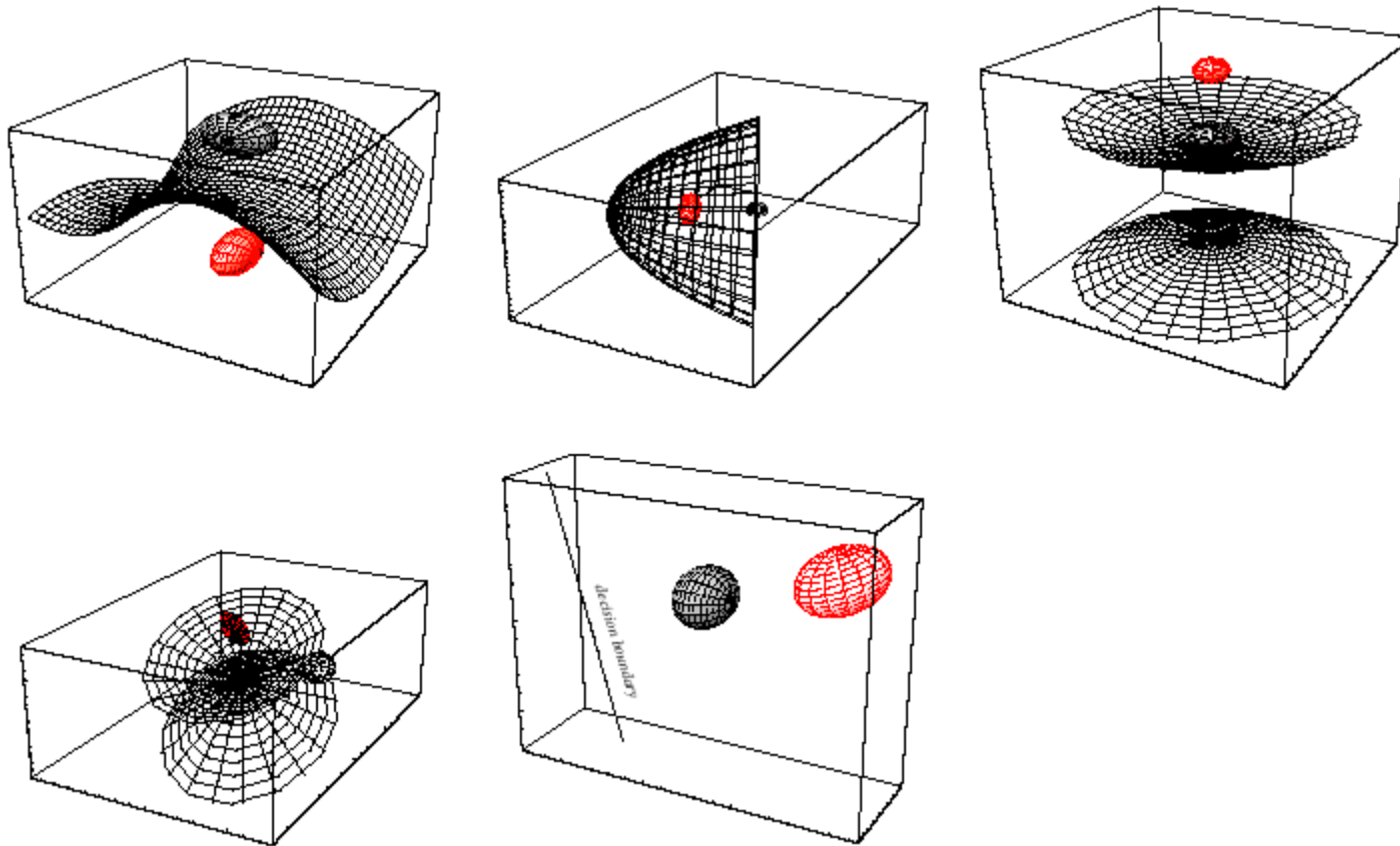
Case 3: $\Sigma_i \neq \Sigma_j$



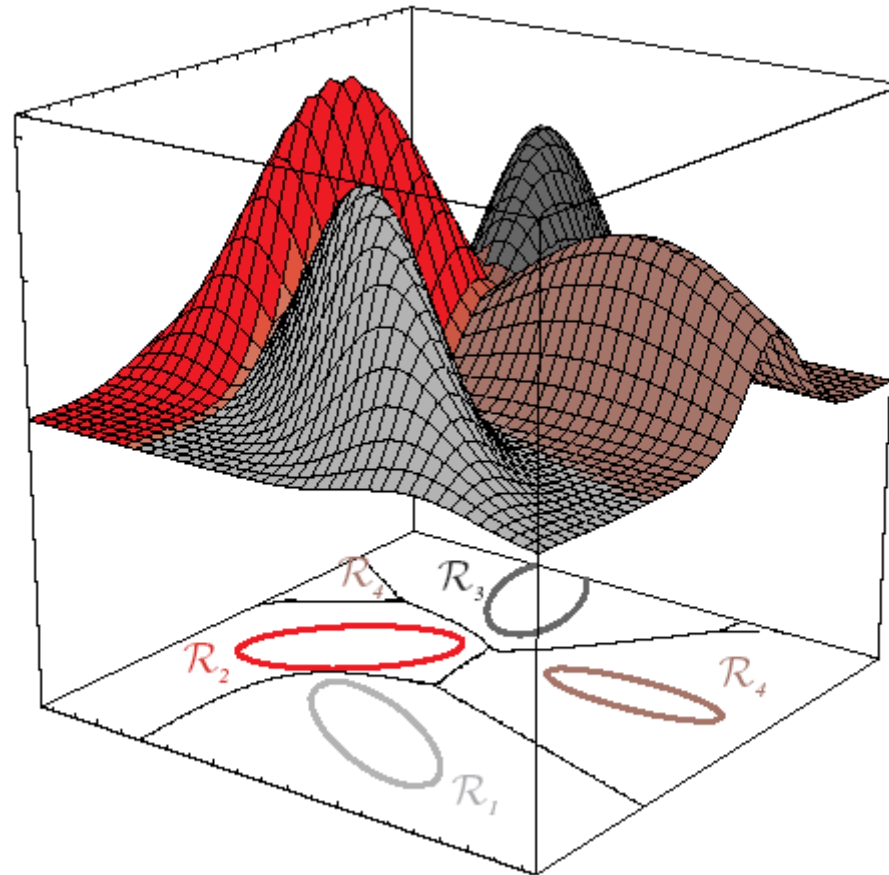
Case 3: $\Sigma_i \neq \Sigma_j$



Case 3: $\Sigma_i \neq \Sigma_j$



Case 3: $\Sigma_i \neq \Sigma_j$



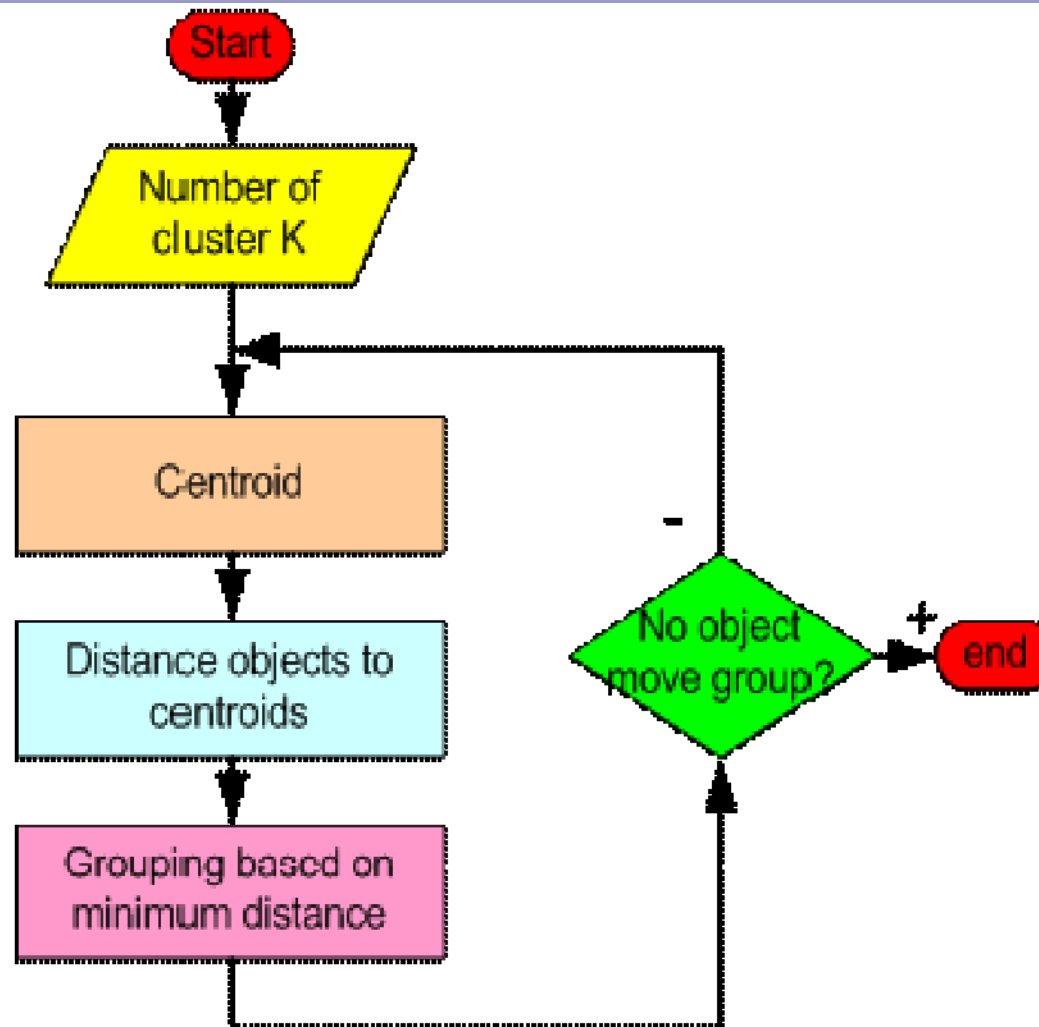
Summary

- Bayesian Decision Theory
 - Basic concepts
 - Bayes theorem
 - Bayes decision rule
- Feasibility of Bayes Decision Rule
 - Prior probability + likelihood
 - Solution I: counting relative frequencies
 - Solution II: conduct density estimation

Summary

- Bayes decision rule: The general scenario
 - Allowing more than one feature
 - Allowing more than two states of nature
 - Allowing actions than merely deciding state of nature
 - Loss function
- Expected loss (conditional risk)
- General Bayes decision rule
- Minimum-error-rate classification
- Discriminant functions
- Gaussian density
- Discriminant functions for Gaussian pdf.

k-means



MIMA Group

[Thank You !]

Any Question?