

# Chapter 2 Bayesian Decision Theory

## Bayesian Decision Theory

Pattern Recognition is a decision process in essence

Bayesian decision theory is a **statistical approach** to pattern recognition

### Basic Assumptions

- The decision problem is posed (formalized) in probabilistic terms
- All the relevant probability values are known

### Key Principle: Bayes Theorem

$$P(H|X) = \frac{P(H)P(X|H)}{P(X)}$$

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

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

- $X$ : the observed sample / **evidence**  
e.g. the length of a fish
- $H$ : the **hypothesis**  
e.g. the fish belongs to the "salmon" category
- $P(H)$ : the **prior probability** (先验概率) that  $H$  holds  
e.g. the probability of catching a salmon
- $P(X|H)$ : the **likelihood** (似然度) of observing  $X$  given that  $H$  holds  
e.g. the probability of observing a 3-inch length fish which is salmon
- $P(X)$ : the **evidence probability** that  $X$  is observed  
e.g. the probability of observing a fish with 3-inch length
- $P(H|X)$ : the **posterior probability** (后验概率) that  $H$  holds given  $X$   
e.g. the probability of  $X$  being salmon given its length is 3-inch

### State of Nature

- Future events that might occur
- State of nature is unpredictable
- Regarded as a random variable

e.g. let  $w$  denote the (discrete) random variable representing the **state of nature (class)** of fish types

$w = w_1$ : sea bass /  $w = w_2$ : salmon

### Prior Probability

Prior probability is the probability distribution which **reflects one's prior knowledge on the random variable**

- the catch produced as much sea bass as salmon  $\rightarrow P(w_1) = P(w_2) = \frac{1}{2}$
- the catch produced more sea bass than salmon  $\rightarrow P(w_1) = \frac{2}{3}; P(w_2) = \frac{1}{3}$

## Class-conditional probability density function

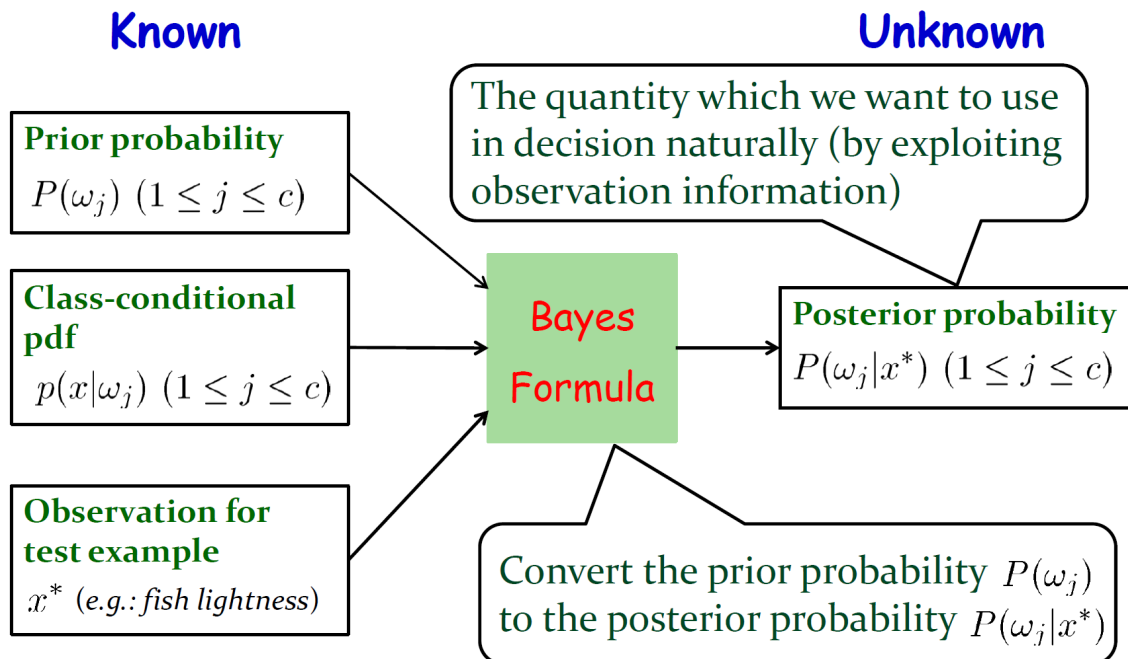
It is a probability density function (pdf) for  $x$  given that the state of nature (class) is  $\omega$

$$p(x|\omega) \geq 0 \quad \int_{-\infty}^{\infty} p(x|\omega) dx = 1$$

The class-conditional pdf describes the difference in the distribution of observations under different classes

$$p(x|\omega_1) \text{ should be different to } p(x|\omega_2)$$

## Decision After Observation



## Bayes Decision Rule

$$\text{if } P(\omega_j|x) > P(\omega_i|x), \forall i \neq j \implies \text{Decide } \omega_j$$

- $P(\omega_j)$  and  $p(x|\omega_j)$  are assumed to be known
- $p(x)$  is **irrelevant** for Bayesian decision but **serving as a normalization factor** and not related to any state of nature)

## Two Special Cases

- Equal prior probability: Depends on the likelihood  $P(x|\omega_j)$

$$P(\omega_1) = P(\omega_2) = \dots = P(\omega_c) = \frac{1}{c}$$

- Equal likelihood: Degenerate to naïve decision rule

$$p(x|\omega_1) = p(x|\omega_2) = \dots = p(x|\omega_c)$$

## Example: Cancer or not?

### 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

Q: If positive test result is returned for some person, does he/she have this kind of cancer or not?

Solution:

Define the notation:

$$\omega_1 : \text{cancer} \quad \omega_2 : \text{no cancer} \quad x \in \{+, -\}$$

Calculate likelihood:

$$\begin{aligned} P(\omega_1) &= 0.008 & P(\omega_2) &= 1 - P(\omega_1) = 0.992 \\ P(+ | \omega_1) &= 0.98 & P(- | \omega_1) &= 1 - P(+ | \omega_1) = 0.02 \\ P(- | \omega_2) &= 0.97 & P(+ | \omega_2) &= 1 - P(- | \omega_2) = 0.03 \end{aligned}$$

Calculate posterior probability:

$$\begin{aligned} P(\omega_1 | +) &= \frac{P(\omega_1) P(+ | \omega_1)}{P(+)} = \frac{P(\omega_1) P(+ | \omega_1)}{P(\omega_1) P(+ | \omega_1) + P(\omega_2) P(+ | \omega_2)} = 0.2085 \\ P(\omega_2 | +) &= 1 - P(\omega_1 | +) = 0.7915 \end{aligned}$$

We find out that

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

So he/she have no cancer!

### Feasibility of Bayes Formula

Get to know the Prior probability  $P(w)$  and likelihood  $p(x | w)$

- Counting relative frequencies (相对频率)  
e.g. Suppose we have randomly picked 1209 cars in the campus, got prices from their owners, and measured their heights

$$\begin{aligned} \text{cars in } w_1 = 221 &\longrightarrow P(\omega_1) = \frac{221}{1209} = 0.183 \\ \text{cars in } w_2 = 988 &\longrightarrow P(\omega_2) = \frac{988}{1209} = 0.817 \end{aligned}$$

- Conduct density estimation (概率密度估计)

**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  $x = 1.05$ , and  $x$  falls into interval  $I_x = [1.0m, 1.1m]$  For  $\omega_1$ ,  $\text{cars in } I_x$  is 46; For  $\omega_2$ ,  $\text{cars in } I_x$  is 59

So

$$\begin{aligned} p(x = 1.05 | w_1) &= \frac{46}{221} = 0.2081 \\ p(x = 1.05 | w_2) &= \frac{59}{988} = 0.0597 \end{aligned}$$

## Is Bayes Decision Rule Optimal?

Bayes Decision Rule (In case of two classes)

if  $P(\omega_1 | x) > P(\omega_2 | x)$  , Decide  $\omega_1$  ; Otherwise  $\omega_2$

Whenever we observe a particular  $x$ , the probability of error is:

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

Under Bayes decision rule, we have

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

The **average probability of error** over all possible  $x$  must be **as small as possible**

## The General Case of Bays Decision Rule

- By allowing to **use more than one feature** (d-dimensional Euclidean space)

$$x \in \mathbf{R} \implies \mathbf{x} \in \mathbf{R}^d$$

- By allowing more than two states of nature (finite set of  $c$  states of nature)

$$\Omega = \{\omega_1, \omega_2, \dots, \omega_c\}$$

- By **allowing actions other than merely deciding the state of nature** (finite set of  $a$  possible actions)

$$\mathcal{A} = \{\alpha_1, \alpha_2, \dots, \alpha_a\}$$

Note that  $c \neq a$

## Loss function

$\lambda(w_j, \alpha_j)$  [also written as  $\lambda(\alpha_i | w_j)$ ] is the loss incurred for taking action  $\alpha_i$  when the state of nature is  $w_j$

$$\lambda : \Omega \times \mathcal{A} \rightarrow R$$

## 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

## Action

Given a particular  $x$ , we have to decide which **action** to take.

With the loss of taking each action  $\alpha_i$ , to minimize the loss  $\lambda(\alpha_i | w_j)$

However, the true state of nature is uncertain: **Expected (average) loss**

## Expected loss (Conditional risk)

Average by **enumerating** over all possible states of nature

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

where  $\lambda(\alpha_i | \omega_j)$  is the incurred loss of taking action  $\alpha_i$  in case of true state of nature being  $\omega_j$ , and  $P(\omega_j | \mathbf{x})$  is the probability of  $\omega_j$  being the true state of nature.

e.g.

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

## Task

find a mapping from patterns to actions

$$\alpha : \mathbf{R}^d \rightarrow \mathcal{A} \text{ (decision function)}$$

In other words, for every  $\mathbf{x}$ , the decision function  $\alpha(\mathbf{x})$  assumes one of the  $a$  actions  $\{\alpha_1, \dots, \alpha_a\}$

## Overall risk R

expected loss with decision function  $\alpha(\cdot)$

$$R = \int R(\alpha(\mathbf{x}) | \mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x}$$

where  $p(\mathbf{x})$  is the pdf for pattern

For every  $\mathbf{x}$ , we ensure that the conditional risk  $R(\alpha(\mathbf{x}) | \mathbf{x})$  is as small as possible, so the risk over all possible  $\mathbf{x}$  must be as small as possible.

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

The resulting overall risk is called the **Bayes risk** (denoted as  $R^*$ ), which is the **best performance** achievable given  $p(\mathbf{x})$  and loss function.

## Example 1: Two-Category Classification

Classification setting

- $\Omega = \{\omega_1, \omega_2\}$ : two states of nature
- $\mathcal{A} = \{\alpha_1, \alpha_2\}$ :  $\alpha_1$  means decide  $\omega_1$ ,  $\alpha_2$  means decide  $\omega_2$
- $\lambda_{ij} = \lambda(\alpha_i | \omega_j)$ : the loss incurred for deciding  $\omega_i$  when the true state of nature is  $\omega_j$

The conditional risk

$$\begin{aligned} R(\alpha_1 | \mathbf{x}) &= \lambda_{11} \cdot P(\omega_1 | \mathbf{x}) + \lambda_{12} \cdot P(\omega_2 | \mathbf{x}) \\ R(\alpha_2 | \mathbf{x}) &= \lambda_{21} \cdot P(\omega_1 | \mathbf{x}) + \lambda_{22} \cdot P(\omega_2 | \mathbf{x}) \end{aligned}$$

If we assume that

$$R(\alpha_1 | \mathbf{x}) < R(\alpha_2 | \mathbf{x})$$

by definition we get

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

by rearrangement we get

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

by Bayes theorem we get

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

because  $\lambda_{21} - \lambda_{11} > 0$ , which means the loss for being error is ordinarily greater than the loss for being correct, we get

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

where  $\frac{p(\mathbf{x} | \omega_1)}{p(\mathbf{x} | \omega_2)}$  is the **likelihood ratio** and  $\frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \cdot \frac{P(\omega_2)}{P(\omega_1)}$  is a **constant** independent of  $\mathbf{x}$

## Example 2: Minimum-Error-Rate Classification

Classification setting

- $\Omega = \{\omega_1, \omega_2, \dots, \omega_c\}$ :  $c$  possible states of nature
- $\mathcal{A} = \{\alpha_1, \alpha_2, \dots, \alpha_c\}$ :  $\alpha_i = \text{decide } \omega_i, 1 \leq i \leq c$

Zero-one (symmetrical) loss function

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

- Assign no loss (e.g. 0) to a correct decision
- Assign a unit loss (e.g. 1) to any incorrect decision (**equal cost**)

Proof

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

$1 - P(\omega_i | \mathbf{x})$  is the **error rate**, the probability that action  $\alpha_i$  is wrong

To **minimum error rate**, decide  $\omega_i$  if  $P(\omega_i | \mathbf{x}) > P(\omega_j | \mathbf{x})$  for all  $j \neq i$

## Minimax Criterion (最小最大化准则)

Generally, we assume that the **prior probabilities** over the states of nature  $\Omega = \{\omega_1, \omega_2, \dots, \omega_c\}$  **are fixed**.

Nonetheless, in some cases we need to design classifiers which can perform well under **varying prior probabilities**.

e.g. the prior probabilities of catching a sea bass or salmon fish might **vary in different regions**

The **minimax criterion** aims to find the classifier which can **minimize the worst overall risk** for **any value of the priors**

## Example of Two-category classification

Suppose the two-category classifier  $\alpha(\cdot)$  decides the feature of  $\omega_1$  in region  $\mathcal{R}_1$  and decides the feature of  $\omega_2$  in region  $\mathcal{R}_2$ . Here,  $\mathcal{R}_1 \cup \mathcal{R}_2 = \mathbf{R}^d$  and  $\mathcal{R}_1 \cap \mathcal{R}_2 = \emptyset$

$$\begin{aligned} R &= \int R(\alpha(\mathbf{x}) | \mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} \\ &= \int_{\mathcal{R}_1} R(\alpha_1 | \mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} + \int_{\mathcal{R}_2} R(\alpha_2 | \mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} \end{aligned}$$

where

$$\begin{aligned} \int_{\mathcal{R}_1} R(\alpha_1 | \mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} &= \int_{\mathcal{R}_1} \sum_{j=1}^2 R(\alpha_1 | \omega_j) \cdot P(\omega_j | \mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} \\ &= \int_{\mathcal{R}_1} \sum_{j=1}^2 \lambda_{1j} \cdot P(\omega_j) \cdot p(\mathbf{x} | \omega_j) d\mathbf{x} \\ &= \int_{\mathcal{R}_1} [\lambda_{11} \cdot P(\omega_1) \cdot p(\mathbf{x} | \omega_1) + \lambda_{12} \cdot P(\omega_2) \cdot p(\mathbf{x} | \omega_2)] d\mathbf{x} \end{aligned}$$

the same

$$\int_{\mathcal{R}_2} R(\alpha_2 | \mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} = \int_{\mathcal{R}_2} [\lambda_{21} \cdot P(\omega_1) \cdot p(\mathbf{x} | \omega_1) + \lambda_{22} \cdot P(\omega_2) \cdot p(\mathbf{x} | \omega_2)] d\mathbf{x}$$

Rewrite the overall risk  $R$  as a function of  $P(\omega_1)$  via

$$\begin{aligned} P(\omega_1) &= 1 - P(\omega_2) \\ \int_{\mathcal{R}_1} p(\mathbf{x} | \omega_1) d\mathbf{x} &= 1 - \int_{\mathcal{R}_2} p(\mathbf{x} | \omega_1) d\mathbf{x} \end{aligned}$$

[In the second equation,  $\int_{\mathcal{R}_1} p(\mathbf{x} | \omega_1) d\mathbf{x}$  stands for the ratio of the sample  $\mathbf{x}$  (its true state of nature is  $\omega_1$ ) have been classified as  $\omega_1$  (view as True Positive); while  $\int_{\mathcal{R}_2} p(\mathbf{x} | \omega_1) d\mathbf{x}$  stands for the ratio of the sample  $\mathbf{x}$  (its true state of nature is  $\omega_1$ ) have been classified as  $\omega_2$  (view as False Positive);]

we get

$$\begin{aligned} R &= \lambda_{22} + (\lambda_{12} - \lambda_{22}) \int_{\mathcal{R}_1} p(\mathbf{x} | \omega_2) d\mathbf{x} \\ &\quad + P(\omega_1) \left[ (\lambda_{11} - \lambda_{22}) + (\lambda_{21} - \lambda_{11}) \int_{\mathcal{R}_2} p(\mathbf{x} | \omega_1) d\mathbf{x} - (\lambda_{12} - \lambda_{22}) \int_{\mathcal{R}_1} p(\mathbf{x} | \omega_2) d\mathbf{x} \right] \end{aligned}$$

where  $R_{mm}$  stands for minimax risk

$$\begin{aligned} R_{mm} &= \lambda_{22} + (\lambda_{12} - \lambda_{22}) \int_{\mathcal{R}_1} p(\mathbf{x} | \omega_2) d\mathbf{x} \\ &= \lambda_{11} + (\lambda_{21} - \lambda_{11}) \int_{\mathcal{R}_2} p(\mathbf{x} | \omega_1) d\mathbf{x} \end{aligned}$$

so for minimax solution we required:

$$\left[ (\lambda_{11} - \lambda_{22}) + (\lambda_{21} - \lambda_{11}) \int_{\mathcal{R}_2} p(\mathbf{x} | \omega_1) d\mathbf{x} - (\lambda_{12} - \lambda_{22}) \int_{\mathcal{R}_1} p(\mathbf{x} | \omega_2) d\mathbf{x} \right] = 0$$

## Discriminant Function (判别函数)

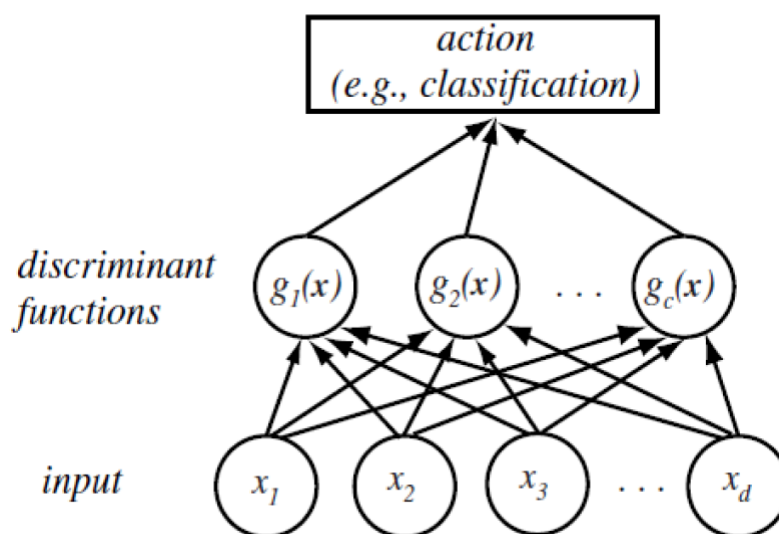
$$g_i : \mathbf{R}^d \rightarrow \mathbf{R} \quad (1 \leq i \leq c)$$

Decide  $\omega_i$  if  $g_i(\mathbf{x}) > g_j(\mathbf{x})$  for all  $j \neq i$

- Useful way to represent classifiers
- One function per category
- Various **discriminant functions** may leads to **Identical classification results**  
 $f(\cdot)$  is a monotonically increasing function (单调递增函数), then

$$f(g_i(\mathbf{x})) \iff g_i(\mathbf{x})$$

which means they are **equivalent in decision**



### Type

- Minimum risk

$$g_i(\mathbf{x}) = -R(\alpha_i | \mathbf{x}) \quad (1 \leq i \leq c)$$

- Minimum-error-rate

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

## Decision region (决策区域)

From  $c$  discriminant functions, get  $c$  decision regions

$$g_i(\cdot) (1 \leq i \leq c) \rightarrow \mathcal{R}_i \subset \mathbf{R}^d (1 \leq i \leq c)$$

### Definition

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

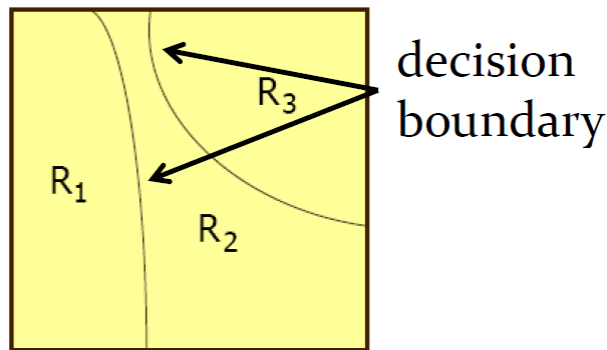
$$\text{where } \mathcal{R}_i \cap \mathcal{R}_j = \emptyset (i \neq j) \text{ and } \bigcup_{i=1}^c \mathcal{R}_i = \mathbf{R}^d$$



## Decision boundary (决策边界)

Surface in feature space where ties occur among several largest discriminant functions

On decision boundaries, more than one  $g_i(\cdot)$  is maximum



## Some Probability

### Expected Value $\mu$

$\sim$  stands for has the distribution

- Discrete

$$x \in \mathcal{X} = \{x_1, x_2, \dots, x_c\}$$

$$x \sim P(\cdot)$$

$$\mathcal{E}[x] = \sum_{x \in \mathcal{X}} x \cdot P(x) = \sum_{i=1}^c x_i \cdot P(x_i)$$

- Continuous

$$x \in \mathbf{R}$$

$$x \sim p(\cdot)$$

$$\mathcal{E}[x] = \int_{-\infty}^{\infty} x \cdot p(x) dx$$

### Variance $\sigma^2$

$$\text{Var}[x] = \mathcal{E}[(x - \mathcal{E}[x])^2]$$

- Discrete

$$\text{Var}[x] = \sum_{i=1}^c (x_i - \mu)^2 \cdot P(x_i)$$

- Continuous

$$\text{Var}[x] = \int_{-\infty}^{\infty} (x - \mu)^2 \cdot p(x) dx$$

## Vector Random Variables

- joint pdf

$$\mathbf{x} \sim p(\mathbf{x}) = p(x_1, x_2, \dots, x_d)$$

- marginal pdf

$$p(\mathbf{x}_1) = \int p(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_2$$

$$(\mathbf{x}_1 \cap \mathbf{x}_2 = \emptyset; \mathbf{x}_1 \cup \mathbf{x}_2 = \mathbf{x})$$

## Expected vector

$$\mathcal{E}[x_i] = \int_{-\infty}^{\infty} x_i \cdot p(x_i) dx_i \quad (1 \leq i \leq d)$$

## Covariance Matrix

### Symmetric & Positive semidefinite

$$\Sigma = [\sigma_{ij}]_{1 \leq i, j \leq d} = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \dots & \sigma_{dd} \end{pmatrix}$$

$$\begin{aligned} \sigma_{ij} &= \sigma_{ji} = \mathcal{E}[(x_i - \mu_i)(x_j - \mu_j)] \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x_i - \mu_i)(x_j - \mu_j) \cdot p(x_i, x_j) dx_i dx_j \end{aligned}$$

$$\sigma_{ii} = \text{Var}[x_i] = \sigma_i^2$$

where  $p(x_i, x_j)$  is **marginal pdf** on a pair of random variables  $(x_i, x_j)$ , exported from **joint pdf**.

## Gaussian Density in Multivariate Case

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

where

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

and

$$\mathbf{x} = (x_1, x_2, \dots, x_d)^t : d\text{-dimensional column vector}$$

$$\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_d)^t : d\text{-dimensional mean vector}$$

$$(\mathbf{x} - \boldsymbol{\mu})^t : 1 \times d \text{ matrix}$$

$$\Sigma^{-1} : d \times d \text{ matrix}$$

$$(\mathbf{x} - \boldsymbol{\mu}) : d \times 1 \text{ matrix}$$

$$(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) : \text{scalar } (1 \times 1 \text{ matrix})$$

because  $\Sigma^{-1}$  is positive definite

$$(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \geq 0$$

## Minimum-error-rate classification

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

which the same as

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

assume that

$$p(\mathbf{x} | \omega_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

so we have

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

where  $\frac{d}{2} \ln 2\pi$  is a constant that could be ignored

**Case 1:  $\boldsymbol{\Sigma}_i = \sigma^2 \mathbf{I}$**

$$\boldsymbol{\Sigma}_i = \sigma^2 \cdot \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix} = \begin{pmatrix} \sigma^2 & & & \\ & \sigma^2 & & \\ & & \ddots & \\ & & & \sigma^2 \end{pmatrix}$$

so we get

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

rearrange

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

where  $\mathbf{x}^t \mathbf{x}$  is the same for all states of nature which could be ignored

Finally we get a **Linear discriminant functions** (线性判别函数)

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

where  $\mathbf{w}_i = \frac{1}{\sigma^2} \boldsymbol{\mu}_i$  is weight vector and  $w_{i0} = -\frac{1}{2\sigma^2} \boldsymbol{\mu}_i^t \boldsymbol{\mu}_i + \ln P(\omega_i)$  is threshold/bias

**Case 2:  $\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}$**

$(\mathbf{x} - \boldsymbol{\mu}_i)^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)$  is call squared **Mahalanobis distance** (马氏距离)

when  $\boldsymbol{\Sigma} = \mathbf{I}$  it reduces to Euclidean distance

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

where  $\mathbf{x}^t \boldsymbol{\Sigma}^{-1} \mathbf{x}$  is the same for all states of nature which could be ignored

Finally we get a **Linear discriminant functions**

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

where  $\mathbf{w}_i = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i$  is weight vector and  $w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^t \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \ln P(\omega_i)$  is threshold/bias

**Case 3:  $\boldsymbol{\Sigma}_i = \text{Arbitrary}$  (任意的)**

Using **quadratic discriminant function** (二次判别函数)

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

where

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

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

$$w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^t \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i) \text{ threshold/bias}$$

