

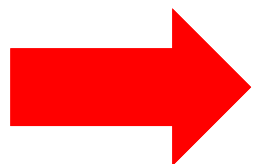


模式识别

中国科学技术大学 汪增福

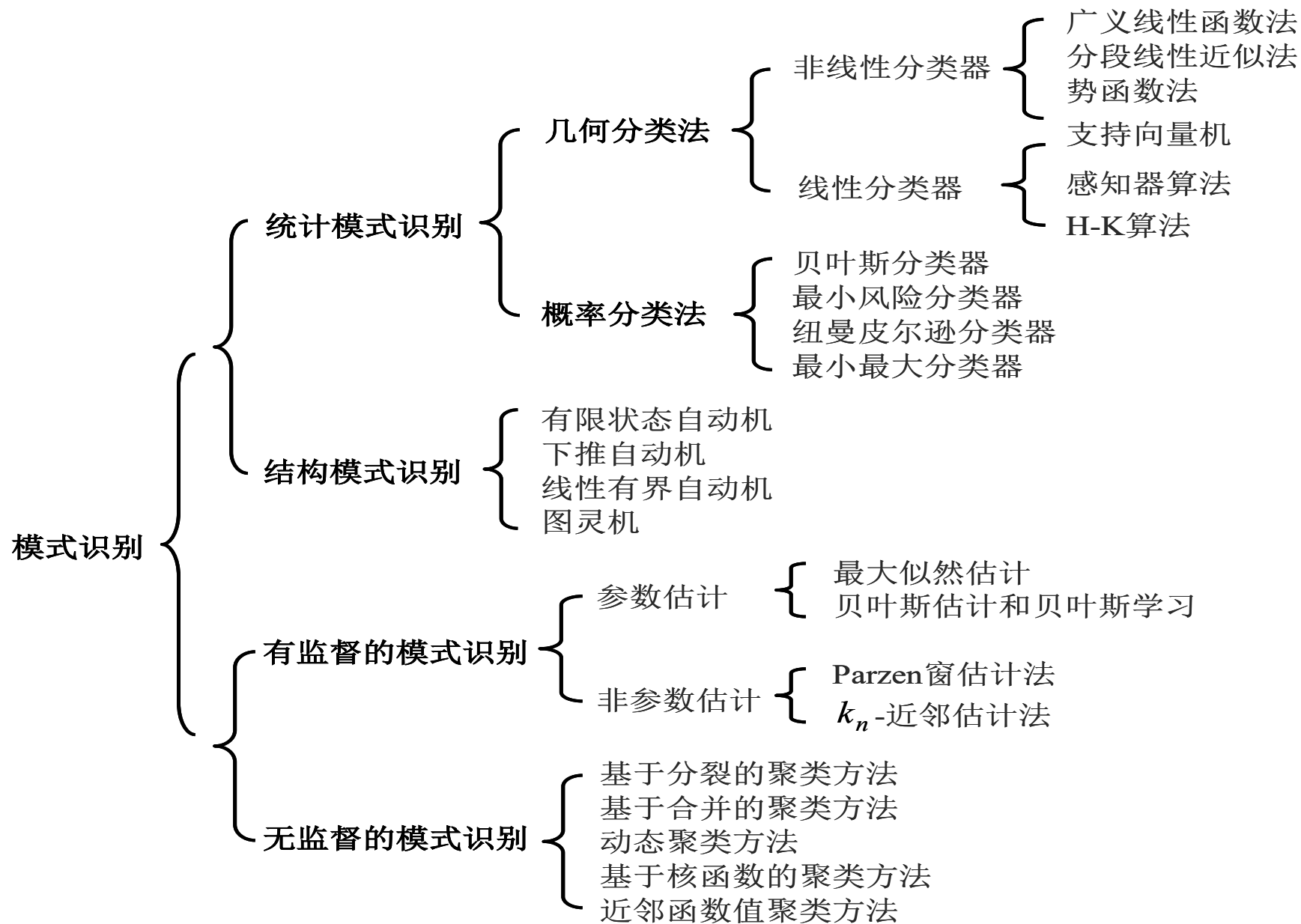
成绩评定方式

- 平时作业，15%。 \longrightarrow x
- 大作业，15%。 \longrightarrow y
- 期末考试（笔试），70%。 \longrightarrow z

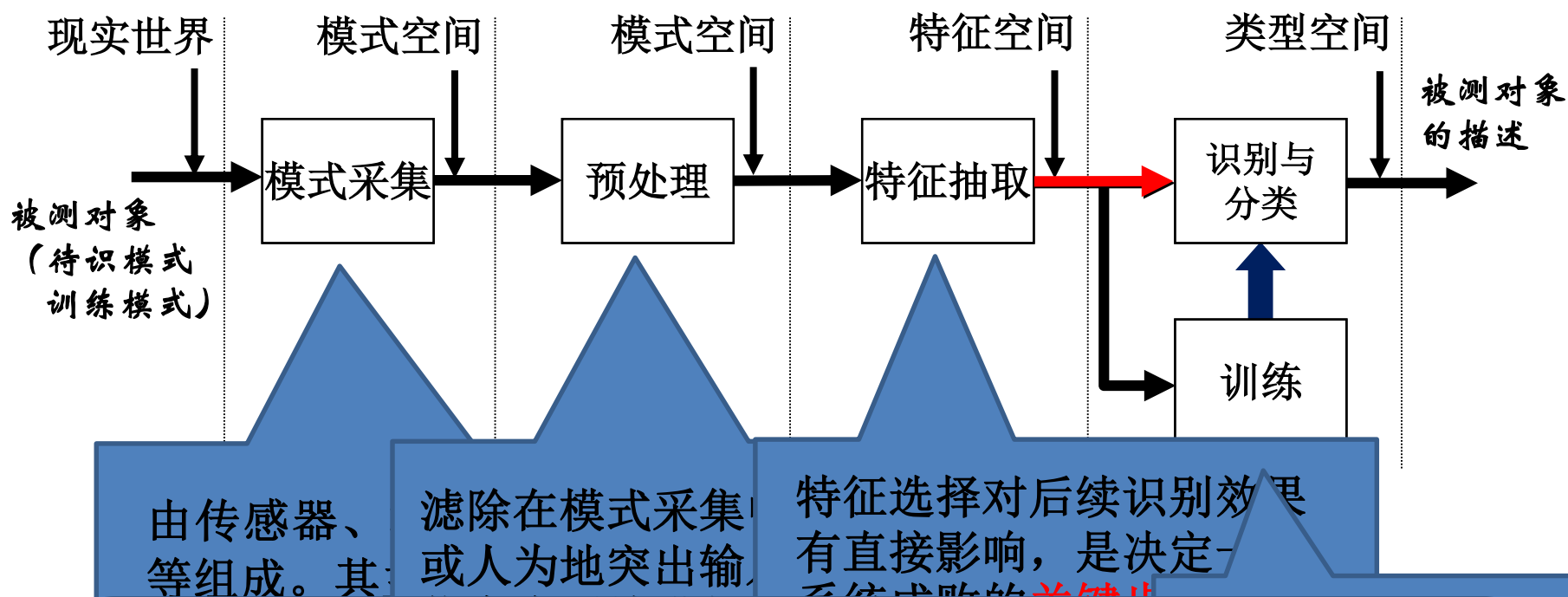


$$\begin{aligned}\text{总成绩} &= \text{三项的加权平均} \\ &= 0.15x + 0.15y + 0.70z\end{aligned}$$

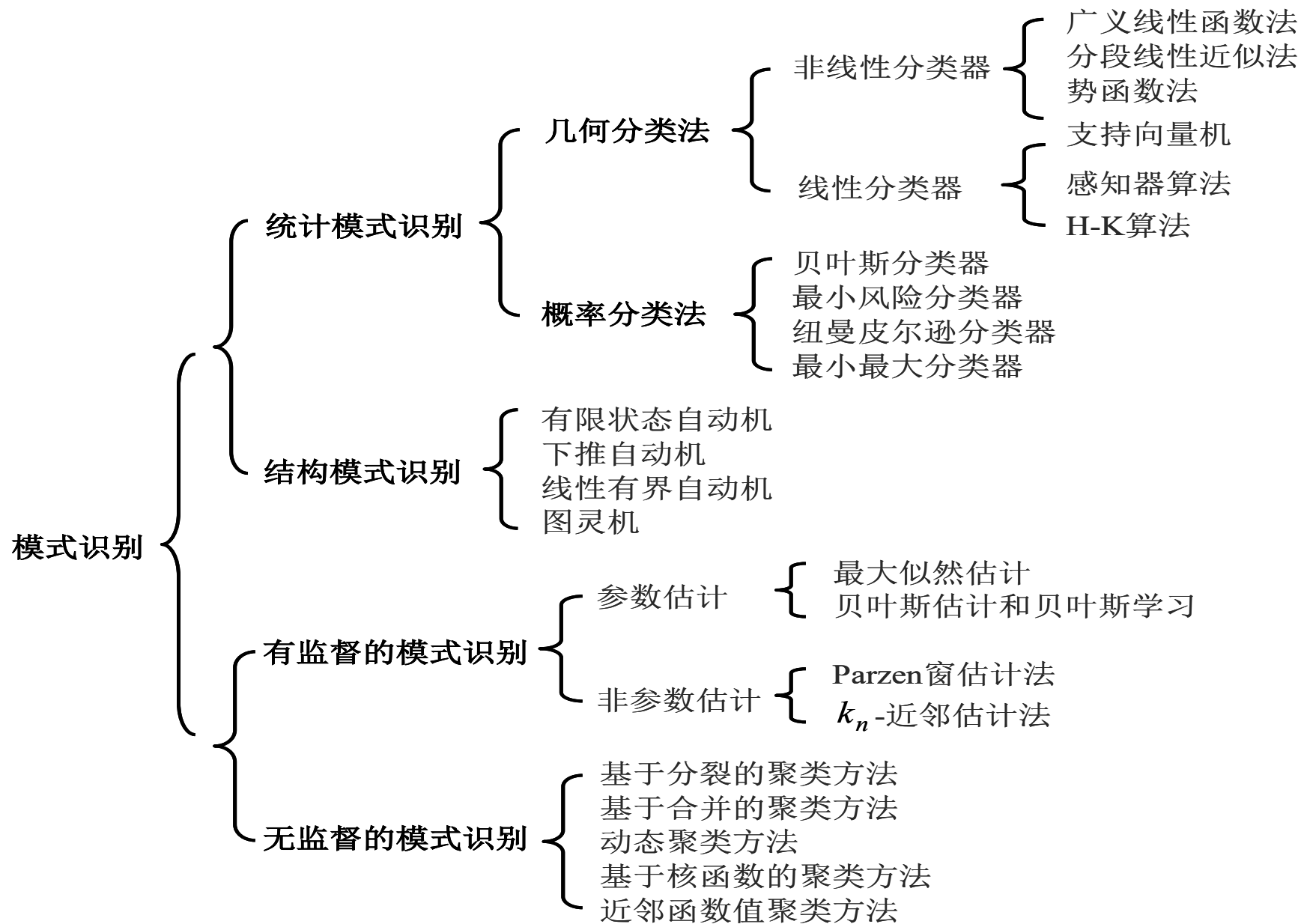
注：期末笔试分两部分：开卷，30分；闭卷，70分。



模式识别系统的基本构成



模式识别是一个过程，它将现实世界中的被测对象通过一系列的变换和处理映射为符号世界中被测对象的分类和描述。



统计模式识别中的几何方法

- 最小距离分类器
- 线性可分情况下的几何分类法
- 非线性可分情况下的几何分类法
- 线性可分问题的非迭代解法
- 最优分类超平面

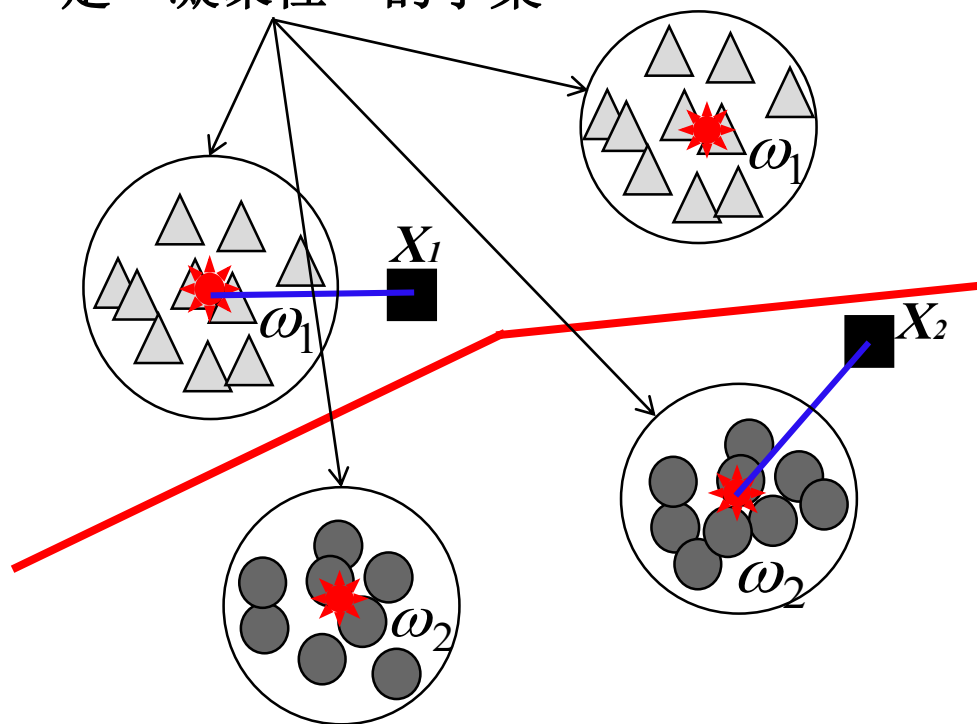
最小距离分类器

● 提高最小距离分类器性能的若干考虑

引入集群操作

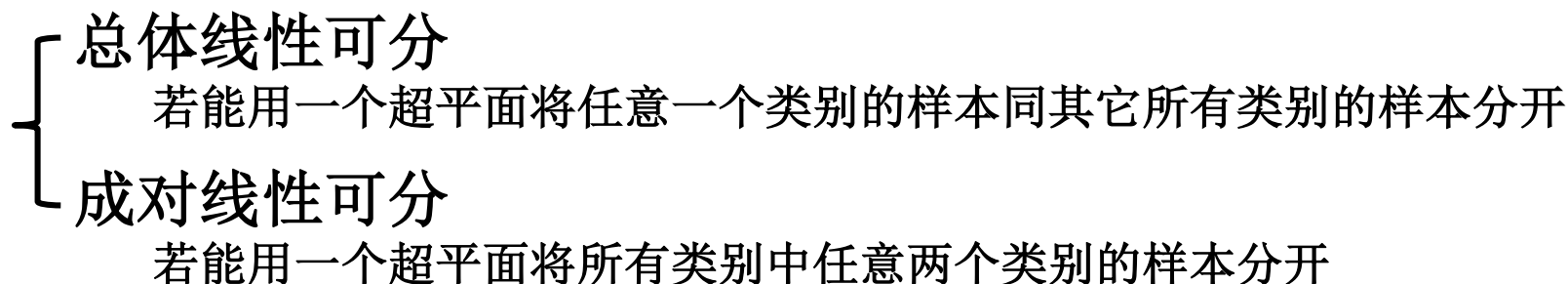
用各个子集的平均样本作为其标准样本
每个子集的样本个数不小于设定的阈值
用所得标准样本集，借助于最近邻法完成分类

具有一定“凝聚性”的子集



线性可分情况下的几何分类法

多类情况下的线性分类器



线性可分

如何得到多类情况下的线性分类器？

➡ 由求解多个两类问题的线性分类器得到。



线性可分情况下的几何分类法

第一种情况：输入样本集合总体线性可分

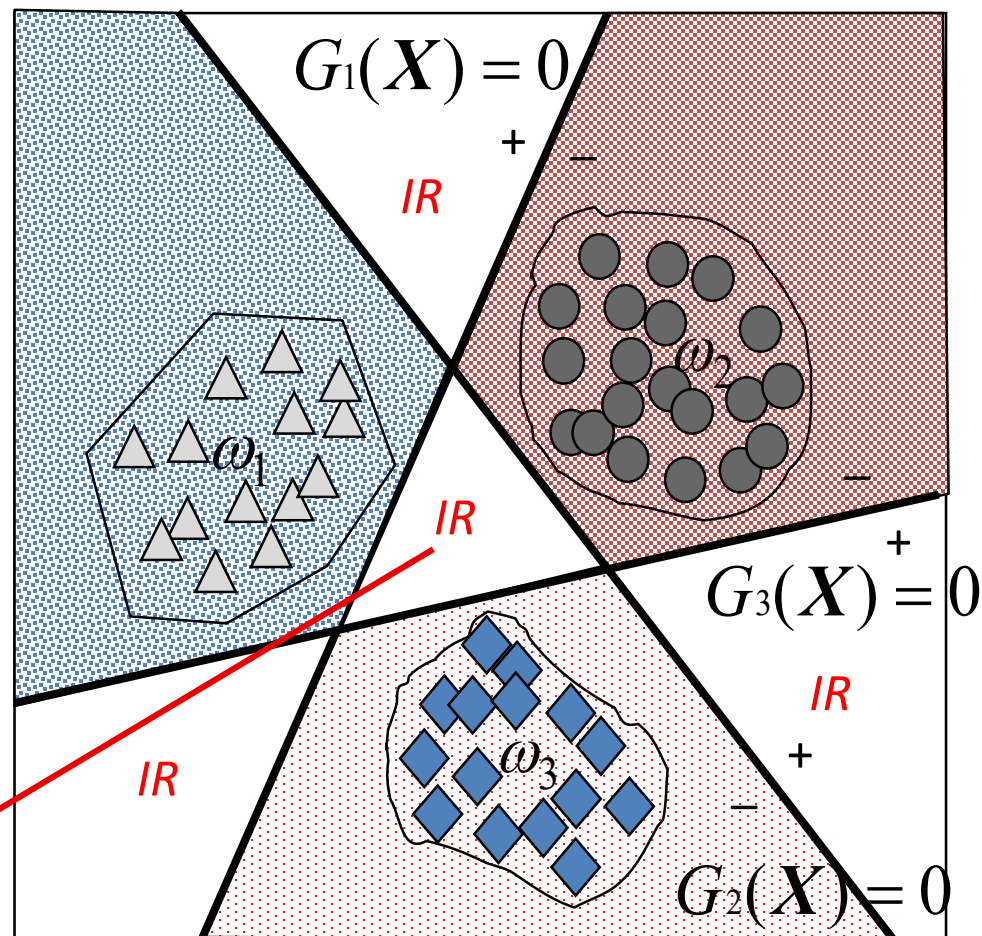
示例：两维三类

$$\left\{ \begin{array}{l} \omega_1 / \bar{\omega}_1 \\ \omega_2 / \bar{\omega}_2 \\ \omega_3 / \bar{\omega}_3 \end{array} \right.$$

决策域 $\omega_i, i=1,2,3$ 的确定

$$G_i(X) > 0, G_j(X) < 0, \forall j \neq i$$

不确定区域



线性可分情况下的几何分类法

第二种情况：输入样本集合成对线性可分

示例：两维三类

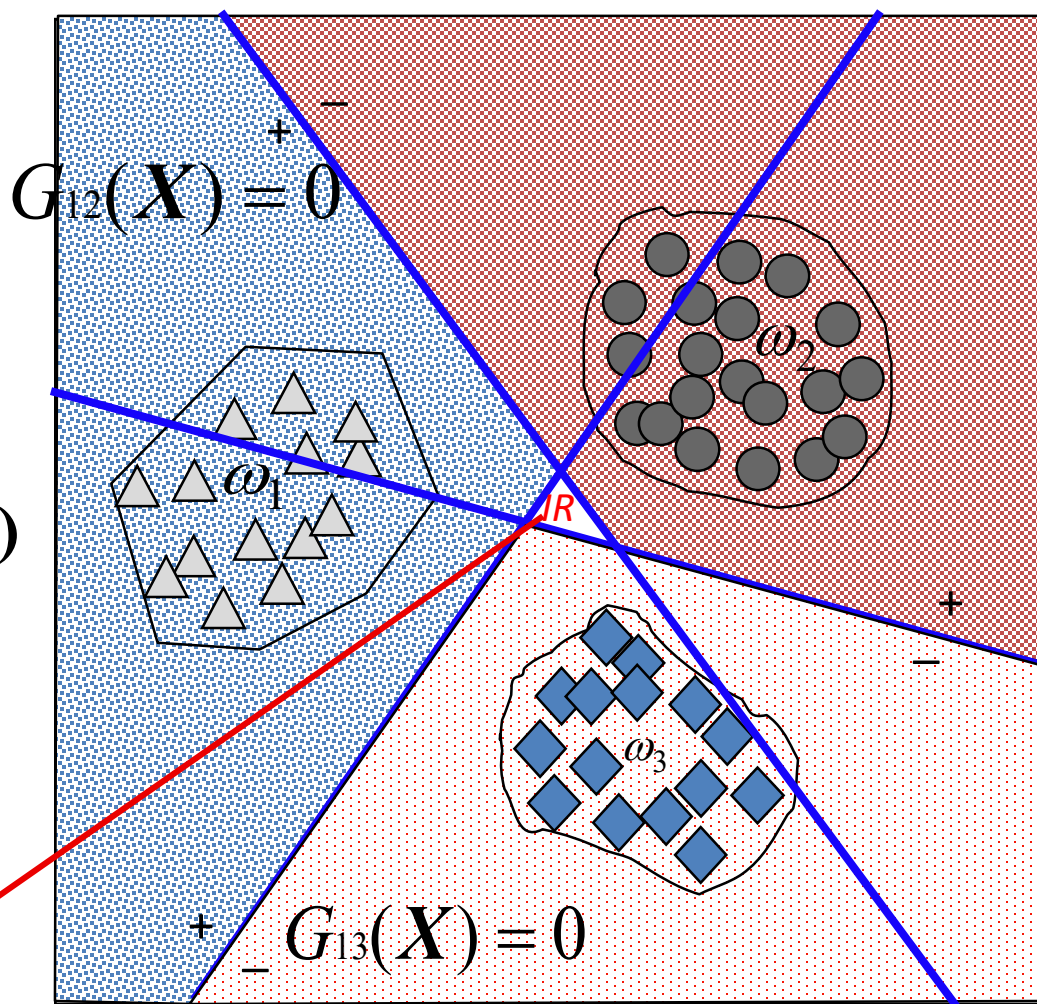
$$\begin{cases} \omega_1 / \omega_2 \\ \omega_1 / \omega_3 \\ \omega_2 / \omega_3 \end{cases}$$

$$G_{ij}(X) = -G_{ji}(X)$$

$\omega_i, i=1,2,3$ 的决策域

$$G_{ij}(X) > 0, \forall j \neq i$$

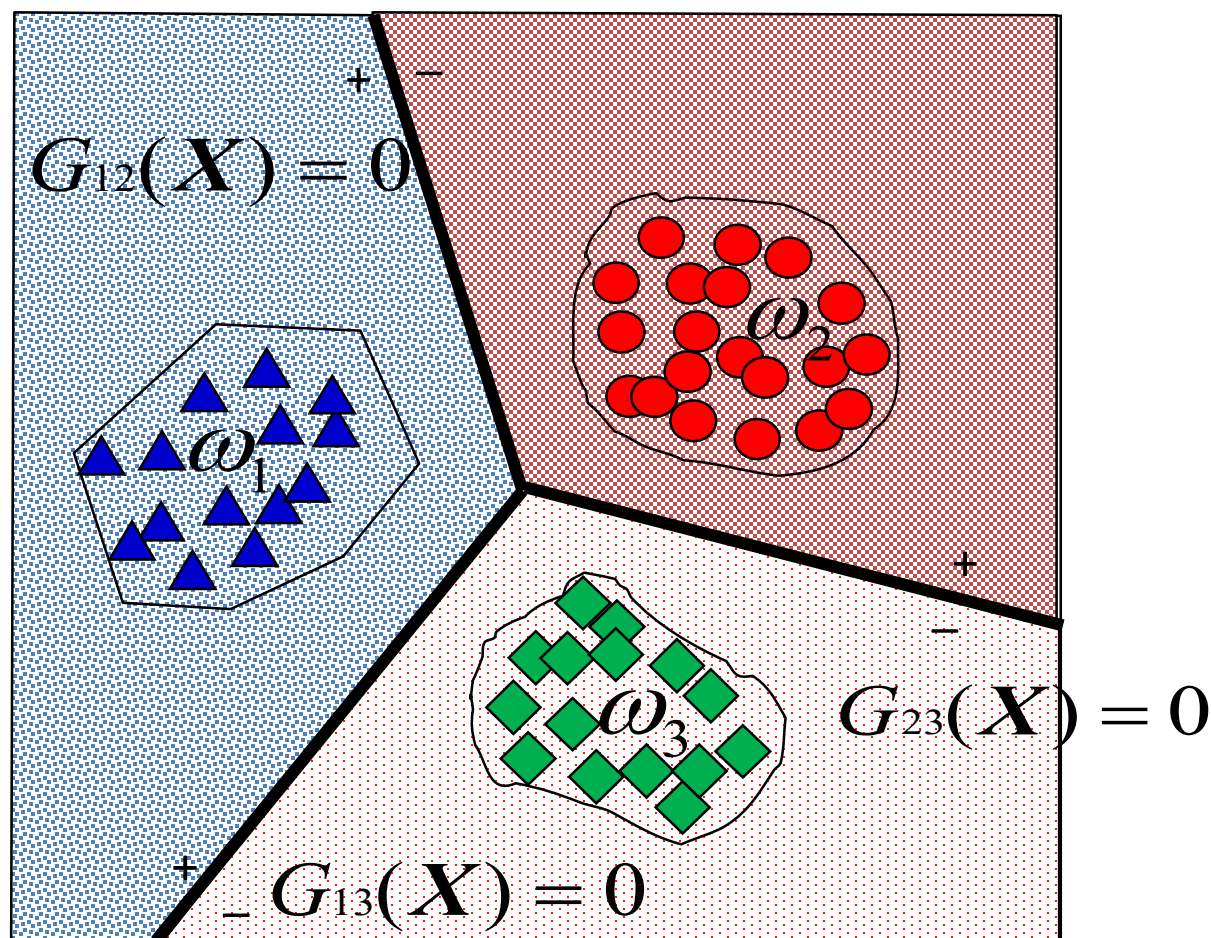
不确定区域



线性可分情况下的几何分类法

第三种情况：最大值判决

- ✓ 只要给定的分类问题是线性可分的，则其决策域不存在不确定区域。



线性可分情况下的几何分类法

感知器算法 ω_i / ω_j 两分问题

Step1. 赋初值：迭代步数 $k=0$ ，固定比例因子 $0 \leq \rho \leq 1$ ；

$W(0)$ ，连续正确分类计数器 $N_c = 0$ ；

Step2. 读入训练样本集合 $\mathcal{X} = \{X_0, X_1, \dots, X_{N-1}\}$ ；

Step3. 取样本 $X = X_{[k]_N}$ $[k]_N = k \bmod(N)$ ；

计算 $G(X) = W(k)^T X$ ；

Step4. 修正权向量：

当 $X \in \omega_i$ 时，

若 $G(X) \leq 0$ ，则 $W(k+1) = W(k) + \rho X$ ， $N_c = 0$ ；

否则 $W(k+1) = W(k)$ ， $N_c + = 1$ ；

当 $X \in \omega_j$ 时，

若 $G(X) \geq 0$ ，则 $W(k+1) = W(k) - \rho X$ ， $N_c = 0$ ；

否则 $W(k+1) = W(k)$ ， $N_c + = 1$ ；

Step5. 若 $N_c \geq N$ ，算法结束；否则 $k = k + 1$ ，返回step3。

线性可分情况下的几何分类法

最小平方误差法：LMSE（Least Mean Square Error）算法

Step1. 由训练样本集构建 $[X]$ 并计算 $[X]^\# = [[X]^T [X]]^{-1} [X]^T$

Step2. 赋初值： $k = 0, 0 \leq \rho \leq 1$
 $b(0) =$ 任一正向量, $W(0) = [X]^\# b(0)$ 。

Step3. 计算 $e(k) = [X]W(k) - b(k)$

Step4. 根据 $e(k)$ 取值情况作进一步处理：

4.1 若各分量取值全部大于或等于0，则判定算法完成；

4.2 若各分量取值有正有负，则继续下一步骤；

4.3 若各分量取值小于或等于0，但不是全部为0，则中止迭代，判定训练样本集不是线性可分的。

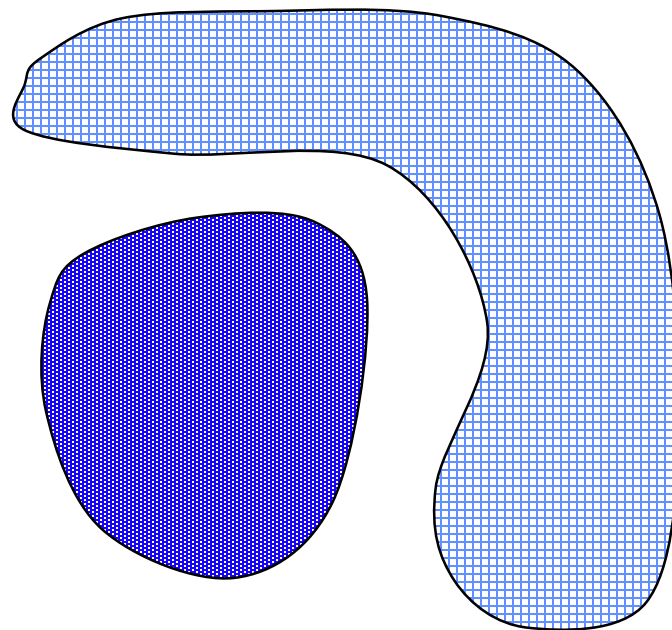
Step5. 完成更新计算

$$\begin{cases} \delta b(k) = \rho(e(k) + |e(k)|) \\ b(k+1) = b(k) + \delta b(k) \\ W(k+1) = W(k) + [X]^\# \delta b(k) \end{cases}$$

Step6. 置 $k=k+1$ ，返回step3。

非线性可分情况下的几何分类法

- 广义线性判别函数法
- 分段线性判别函数近似法
- 非线性判别函数法



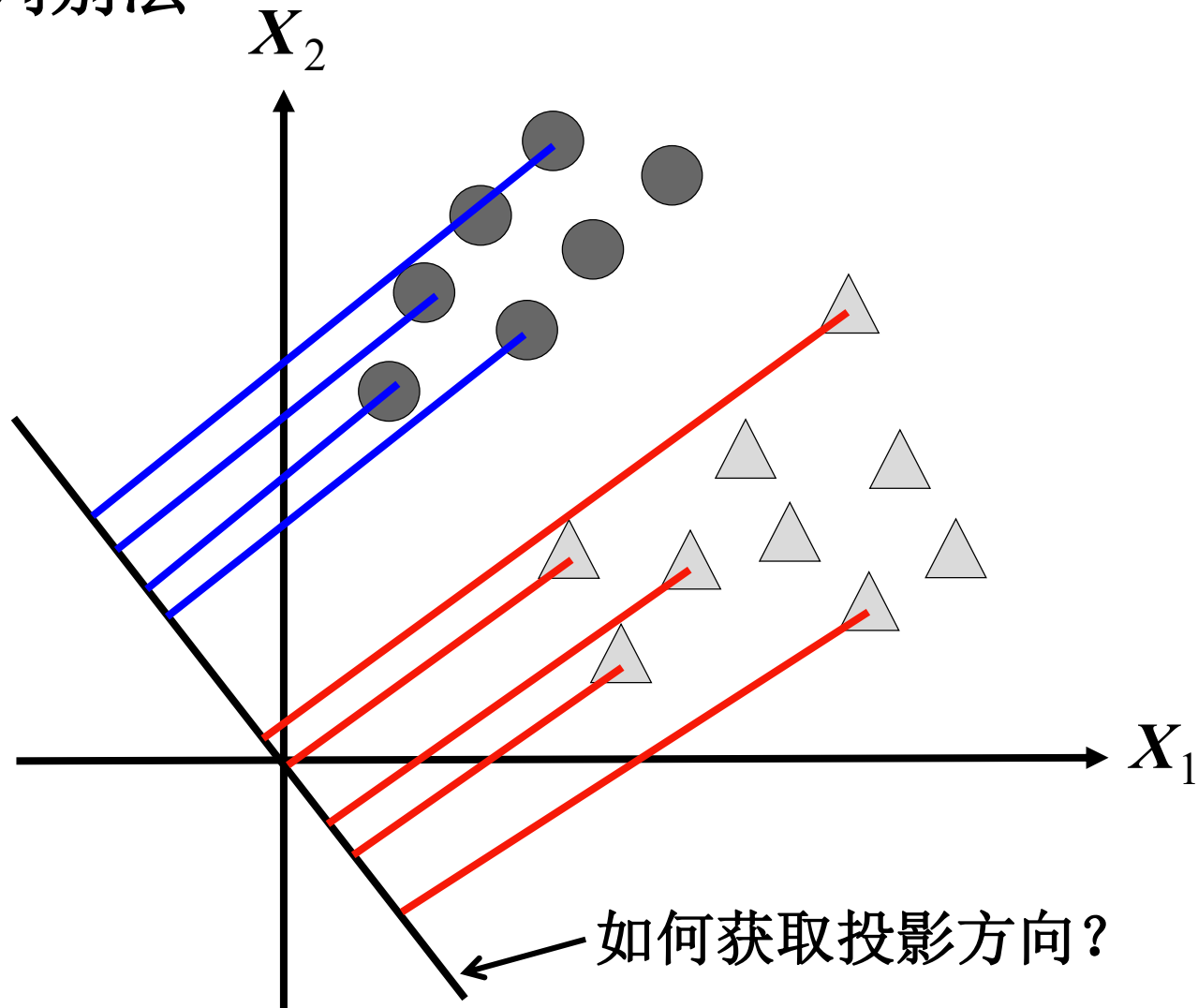
线性可分问题的非迭代解法

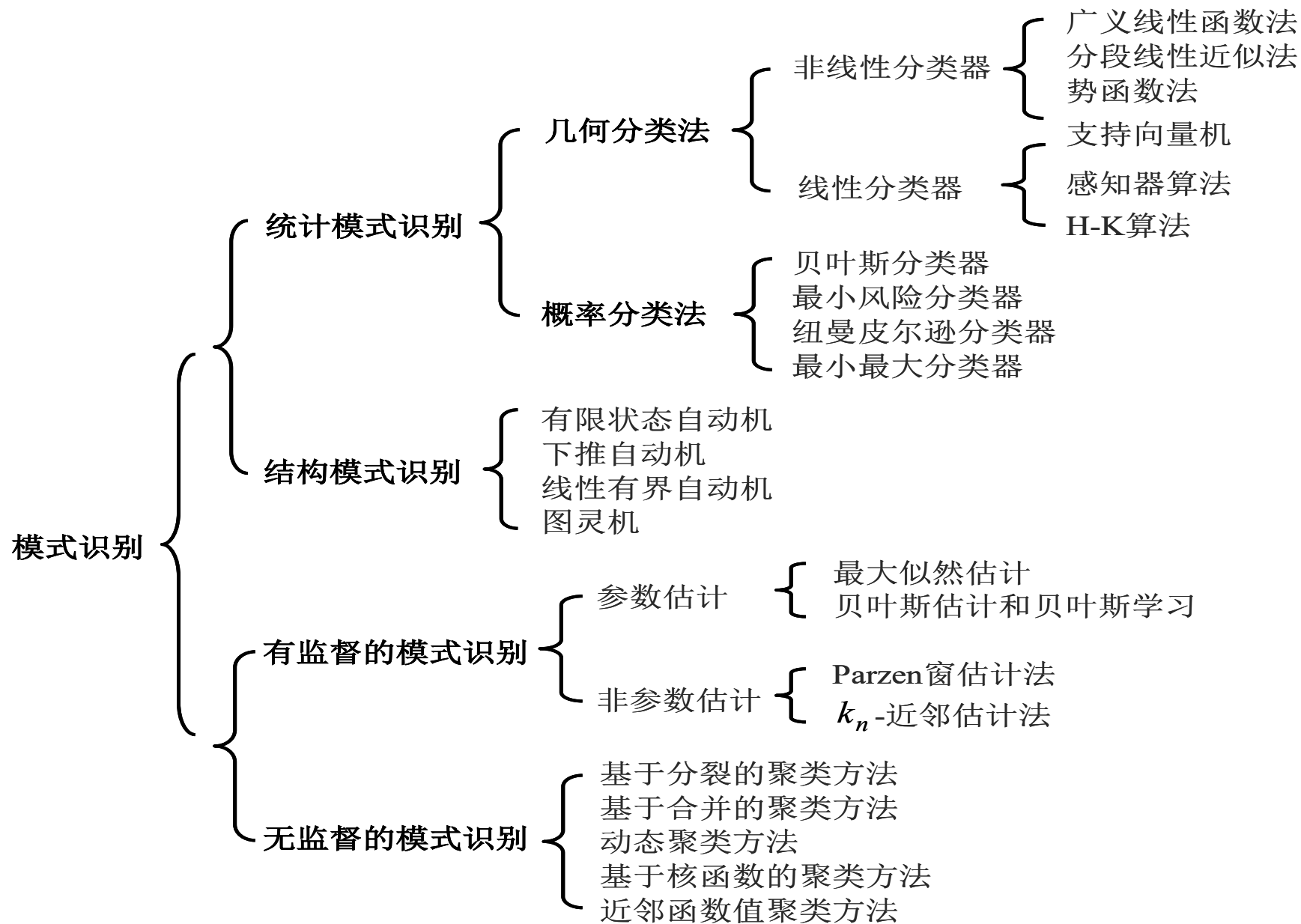
Fisher线性判别法

- 把高维特征空间中的样本投影到一条直线上，实现从高维到一维的数据压缩。
- 改变直线的空间取向，得到不同的投影结果。
- 如果在投影后的直线上训练样本具有很好的分布，则可以通过简单操作实现对输入样本的分类。

线性可分问题的非迭代解法

Fisher线性判别法





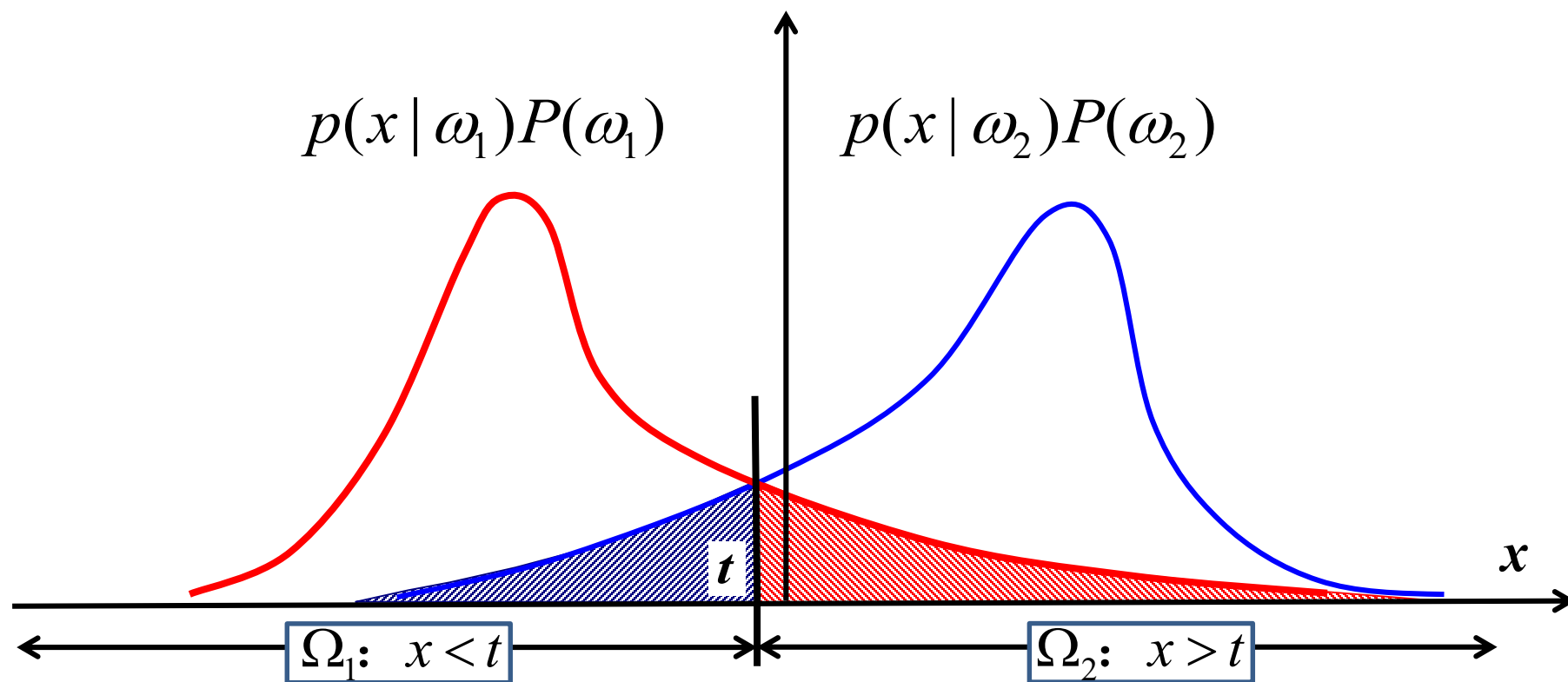
统计模式识别中的概率方法

- 最小错误概率判决准则
- 最小风险判决规则
- **Neyman-Pearson**判决规则
- 最小最大判决规则
- 类条件概率密度的参数估计
- 类条件概率密度的非参数估计

最小错误概率判决准则

一维两类问题示例：分类错误率

$$P(e) = \int_{\Omega_2} P(\omega_1) p(x | \omega_1) dx + \int_{\Omega_1} P(\omega_2) p(x | \omega_2) dx$$



最小风险判决规则

最小风险判决规则

$$R(\alpha_i | X) = \underset{j=1,2,\dots,A}{\text{Minimum}} \{R(\alpha_j | X)\} \Rightarrow X \in \omega_i$$

算法步骤:

Step1. 对观测样本 X , 计算

$$P(\omega_j | X) = \frac{p(X | \omega_j)P(\omega_j)}{\sum_{j=1}^N p(X | \omega_j)P(\omega_j)}, j = 1, 2, \dots, N$$

Step2. 计算各判决的条件平均风险

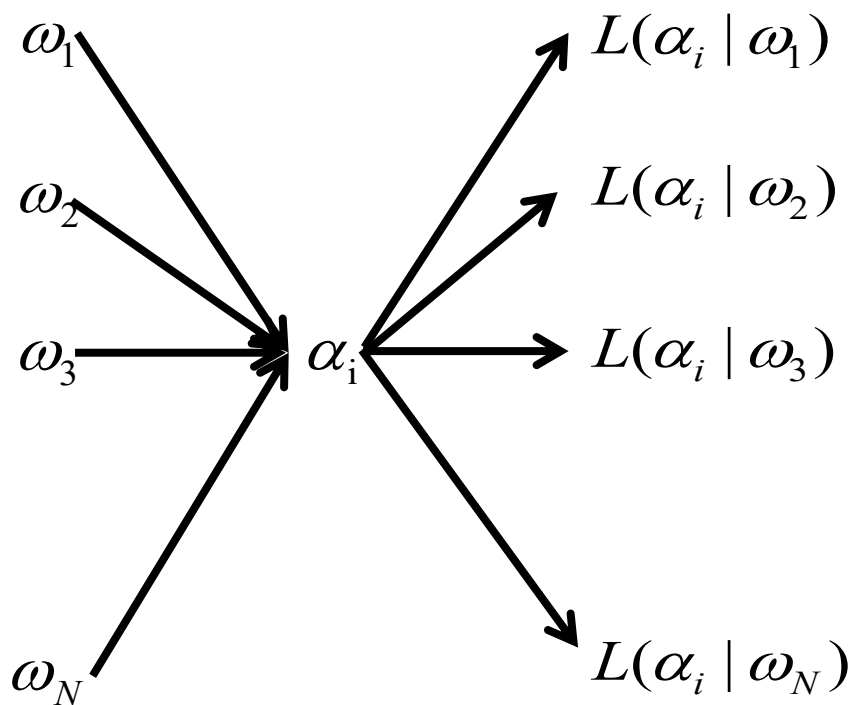
$$R(\alpha_j | X) = \sum_{k=1}^N L(\alpha_j | \omega_k)P(\omega_k | X), j = 1, 2, \dots, N$$

Step3. 将观测样本 X 判属于使条件平均风险最小化的判决所对应的类别。

最小风险判决规则

条件平均风险

类别	判决	风险	判决 α_i 的条件平均风险
----	----	----	-----------------------



观测样本为 X

$P(\omega_j | X), j = 1, 2, \dots, N$

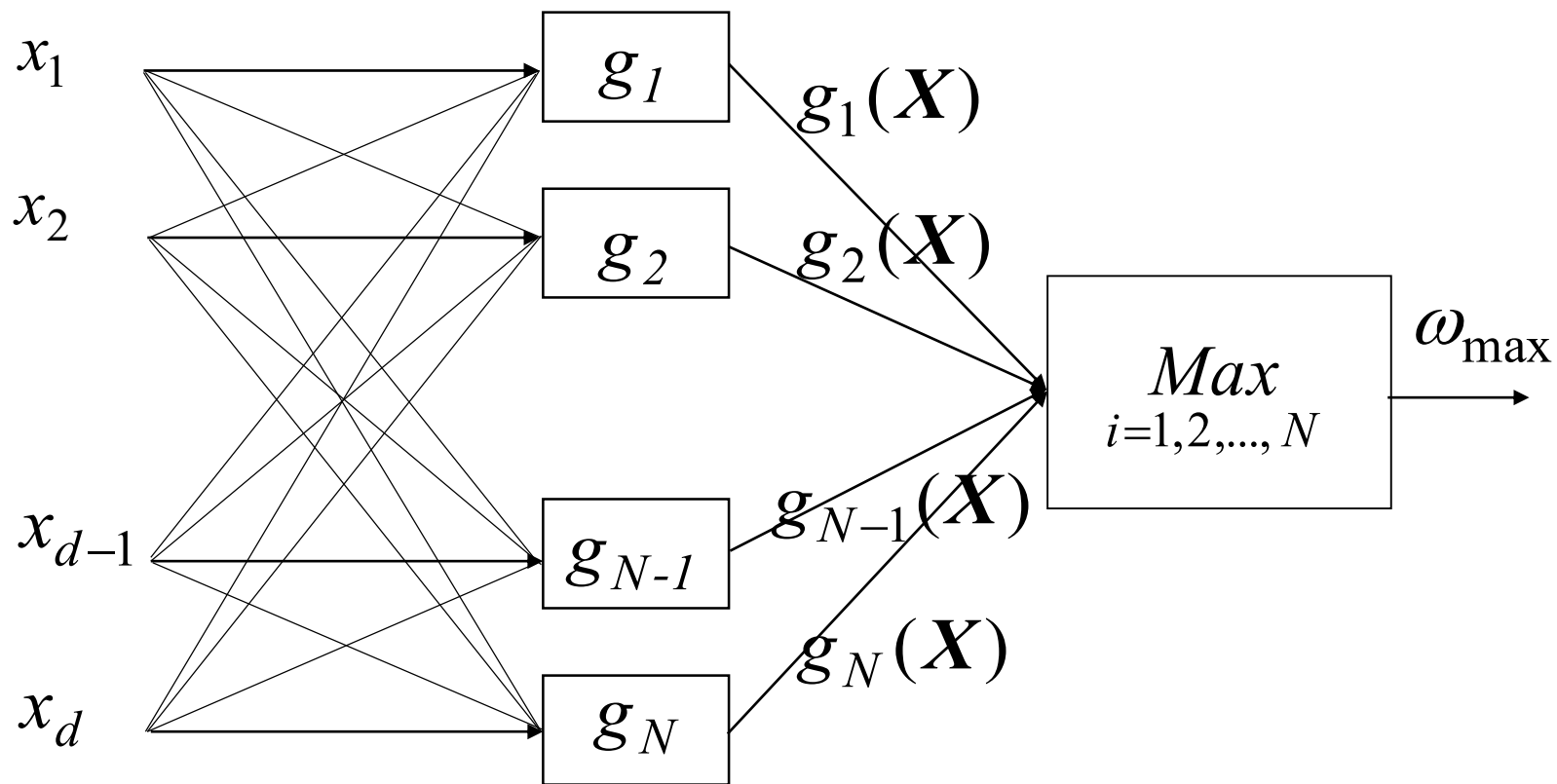
$$R(\alpha_i | X) = \sum_{j=1}^N L(\alpha_i | \omega_j) P(\omega_j | X)$$

$$\sum_{j=1}^N P(\omega_j | X) = 1$$

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

$$\begin{aligned} &= \frac{p(X | \omega_j) P(\omega_j)}{\sum_{j=1}^N p(X | \omega_j) P(\omega_j)} \end{aligned}$$

贝叶斯分类器的一般结构



最大后验概率判决

$$g_i(X) = P(\omega_i | X)$$

最小风险判决

$$g_i(X) = -R(\alpha_i | X)$$

Neyman-Pearson判决规则

➔ 在保持一类分类错误概率不变的条件下，使另一类分类错误概率最小。

用Lagrange乘子法求解

目标函数

约束条件

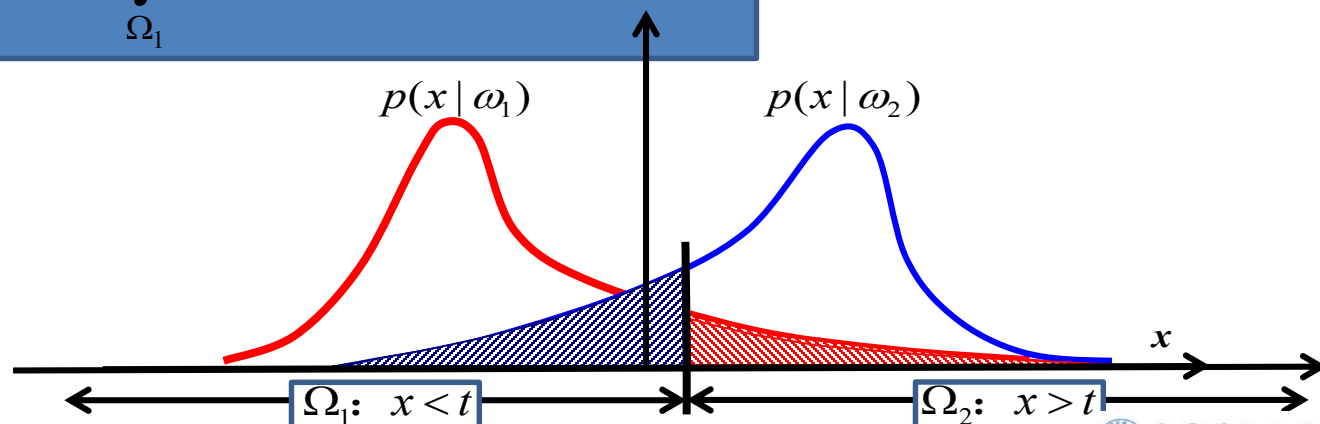
$$J = \int_{\Omega_2} p(x | \omega_1) dx + \lambda \left(\int_{\Omega_1} p(x | \omega_2) dx - \alpha \right)$$

$$= 1 - \int_{\Omega_1} p(x | \omega_1) dx + \lambda \left(\int_{\Omega_1} p(x | \omega_2) dx - \alpha \right)$$

$$= (1 - \lambda \alpha) + \int_{\Omega_1} (\lambda p(x | \omega_2) - p(x | \omega_1)) dx$$

确定 Ω_1

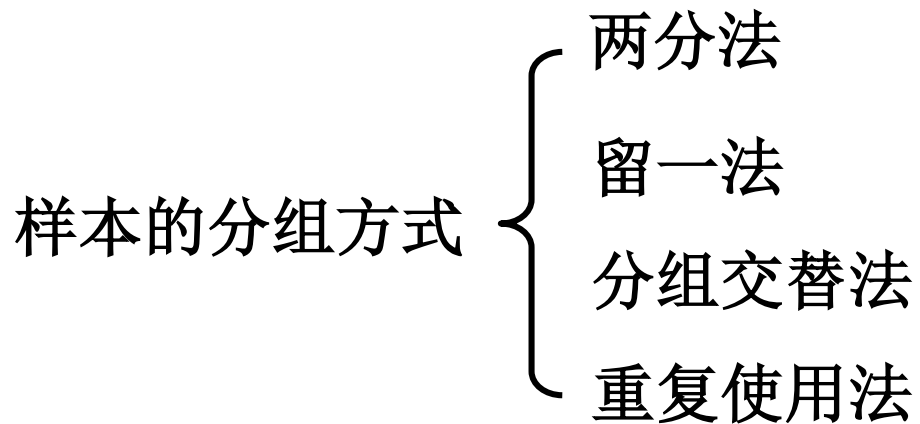
极小化

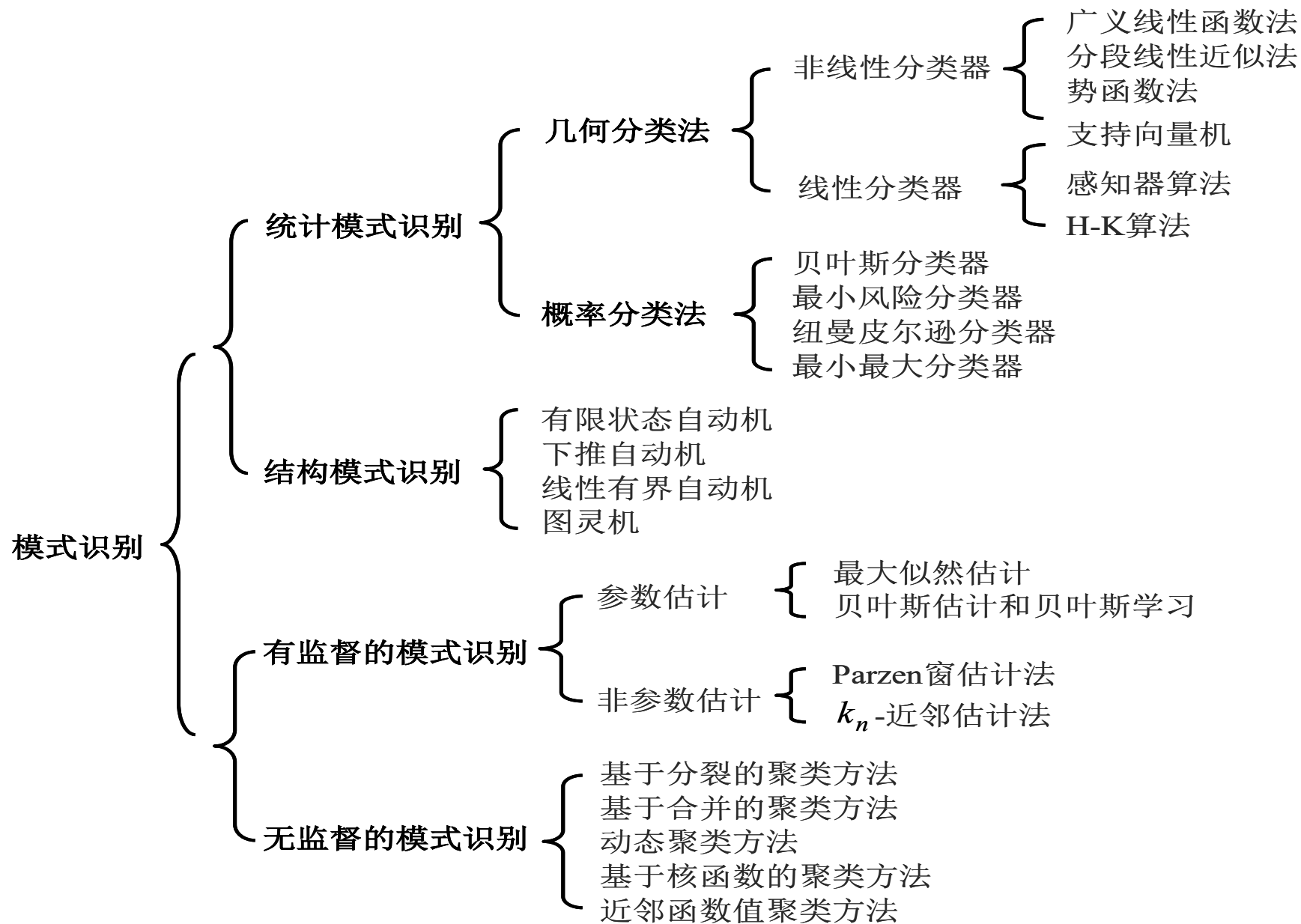


分类器错误率的实验估计

■ 有限样本情况下分类器错误率的实验估计

- ✓ 充分利用已知的有限个样本得到性能好的分类器；
- ✓ 充分利用已知的有限个样本给出可靠的错误率估计。

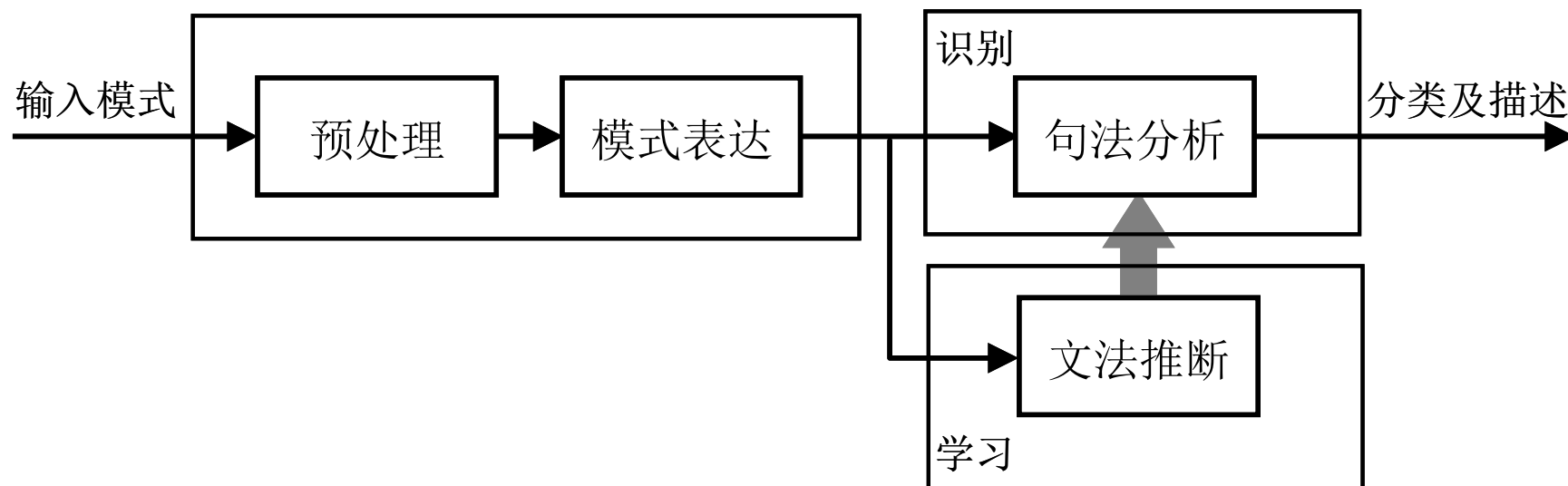




结构模式识别中的句法方法

- 有限状态自动机
- 下推自动机
- 图灵机
- 句法分析
- 文法推断

句法模式识别系统



1. 学习过程（文法推断）

利用已知结构的样本模式来推断产生这些模式的文法规则。

2. 识别过程（句法分析）

用有序字符串表达输入模式，并利用文法规则对其进行句法分析以判断能否由相应的文法所生成。

文法

定义 文法 $G = (N, T, P, S)$ 是一个四元式。其中, N 为 G 的非终结符或变量的有穷集合, T 为 G 的终结符或常量的有穷集合, P 是产生式或再写规则的有穷集合, 而 $S \in N$ 为句子的起始符。

$$N \cap T = \Phi \quad \Sigma = N \cup T \quad \alpha \rightarrow \beta$$

一些约定

大写的拉丁字母 非终结符

小写的拉丁字母 终结符

小写的希腊字母 由非终结符和终结符组成的串

导出=推导=派生

\Rightarrow_G $\Sigma = N \cup T$ 上的一个二元关系

$$\begin{array}{ccc} + & * & n \\ \Rightarrow_G & \Rightarrow_G & \Rightarrow_G \end{array}$$

文法的分类

0型文法 无约束文法/短语结构文法

产生式具有 $\alpha \rightarrow \beta$ 形式的文法称为**0型文法**。

$$\alpha \in \Sigma^+ \quad \beta \in \Sigma^*$$

1型文法 上下文有关文法

产生式具有 $\alpha_1 A \alpha_2 \rightarrow \alpha_1 \beta \alpha_2$ 形式的文法称为**1型文法**。

$$\alpha_1, \alpha_2 \in \Sigma^* \quad A \in N \quad \beta \in \Sigma^+ \quad |A| < |\beta|$$

2型文法 上下文无关文法

产生式具有 $A \rightarrow \beta$ 形式的文法称为**2型文法**。

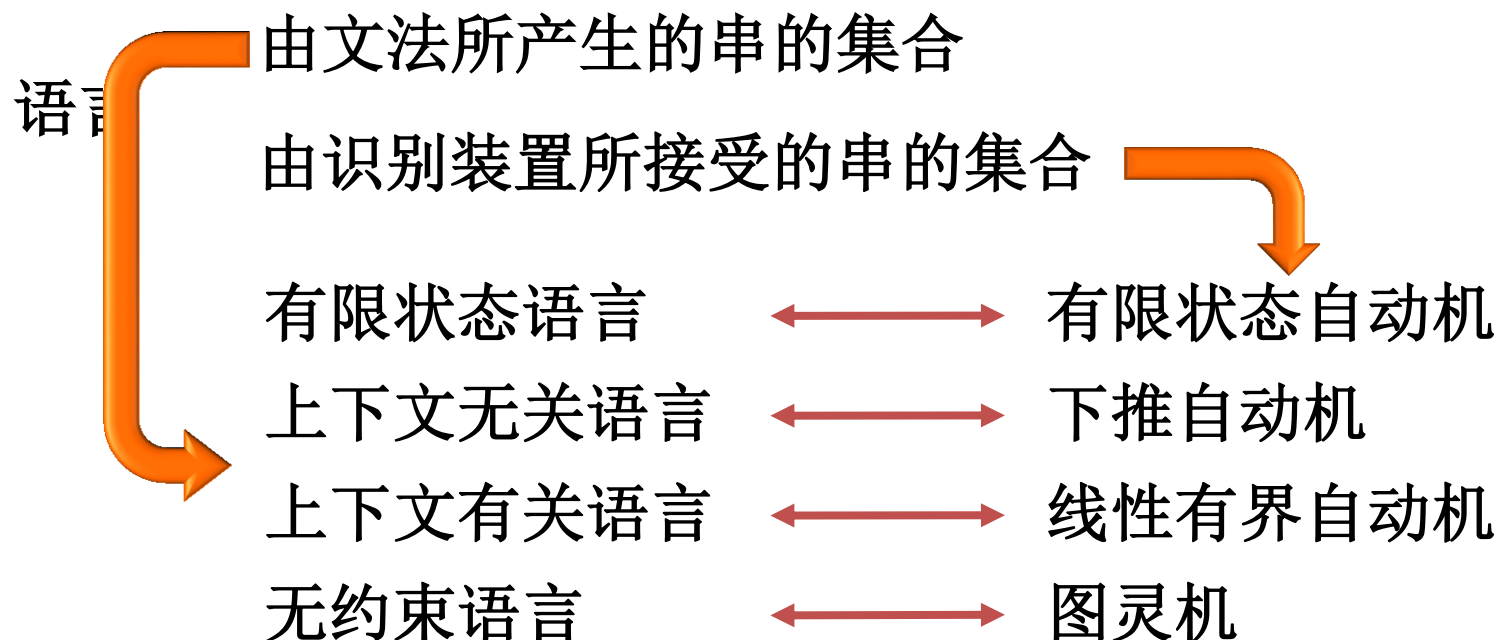
$$A \in N \quad \beta \in \Sigma^+$$

3型文法 有限状态文法/正则文法

产生式具有 $A \rightarrow aB$ 或 $A \rightarrow a$ 形式的文法称为**3型文法**。

$$A, B \in N \quad a \in T$$

自动机



有限状态自动机

确定的有限状态自动机

一个确定的有限状态自动机（**DFA**）是一个五元式：

$$A_f = (Q, \Sigma_I, \delta, q_0, F)$$

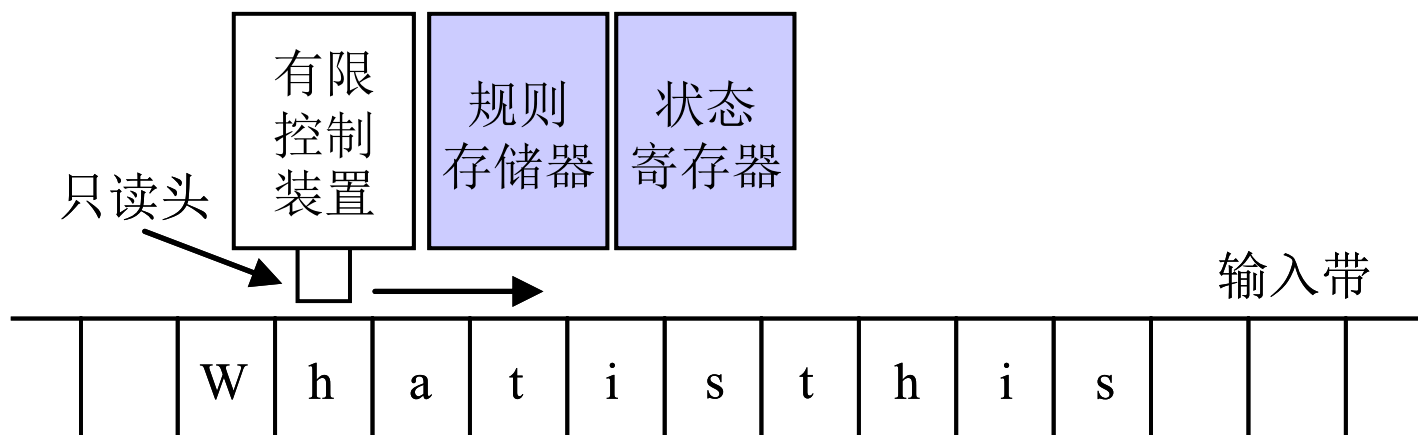
Q 状态的有限集合

Σ_I 输入符号的有限集合，即字母表

δ 映射规则 $\delta : Q \times \Sigma_I \rightarrow Q$

$q_0 \in Q$ 初始状态

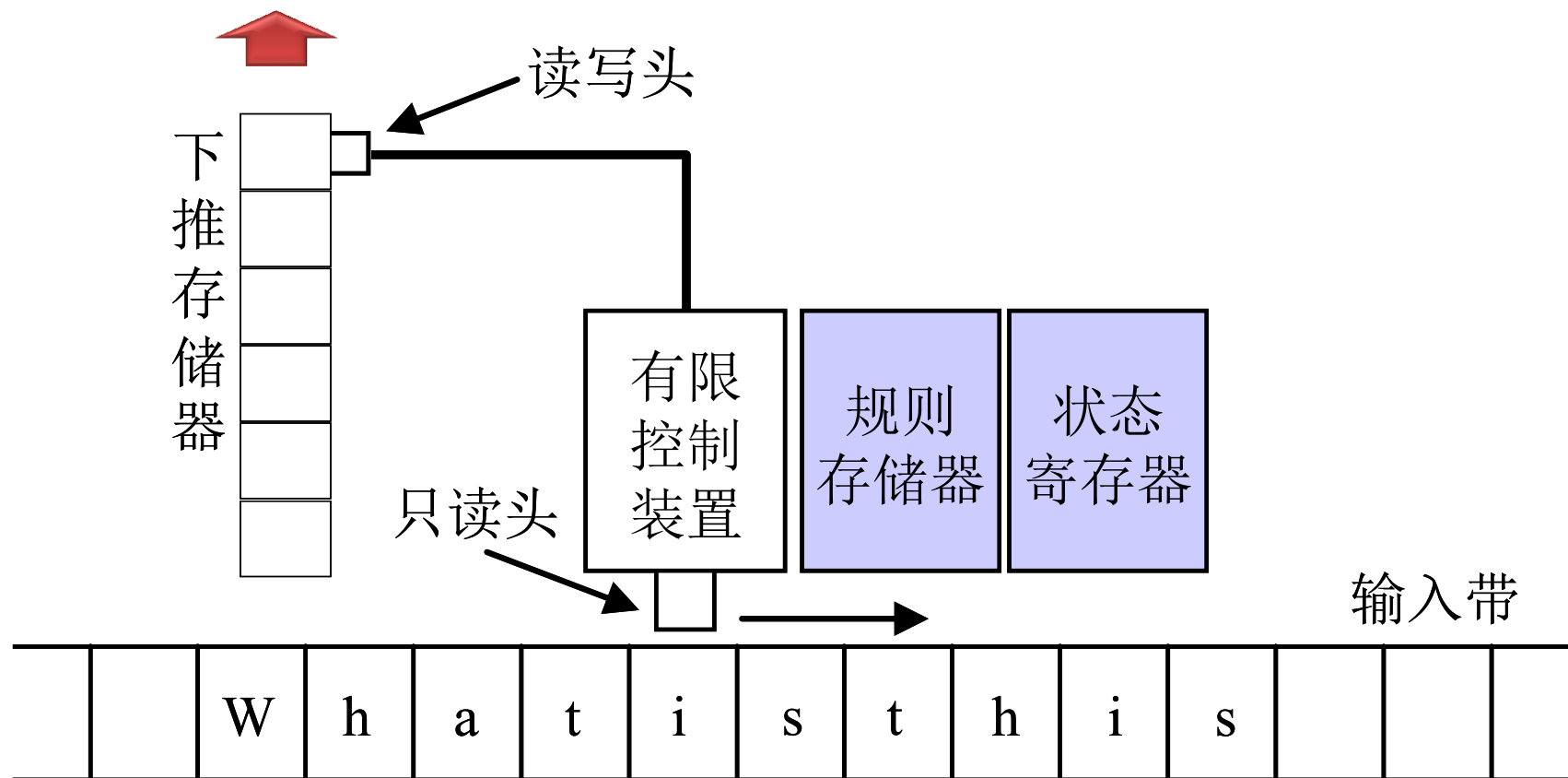
$F \subseteq Q$ 终止状态集合



下推自动机

Push Down Automaton (PDA)

一个长度不受限制的“后入先出”的堆栈



附加有下推存储器的有限状态自动机！

下推自动机

定义 一个非确定的下推自动机（PDA）是一个七元式：

$$A_p = (Q, \Sigma_I, \Gamma, \delta, q_0, Z_0, F)$$

Q 状态的有限集合

Σ_I 输入符号的有限集合，即字母表


Γ 堆栈符号的有限集合

$q_0 \in Q$ 初始状态

$Z_0 \in \Gamma$ 栈底符号

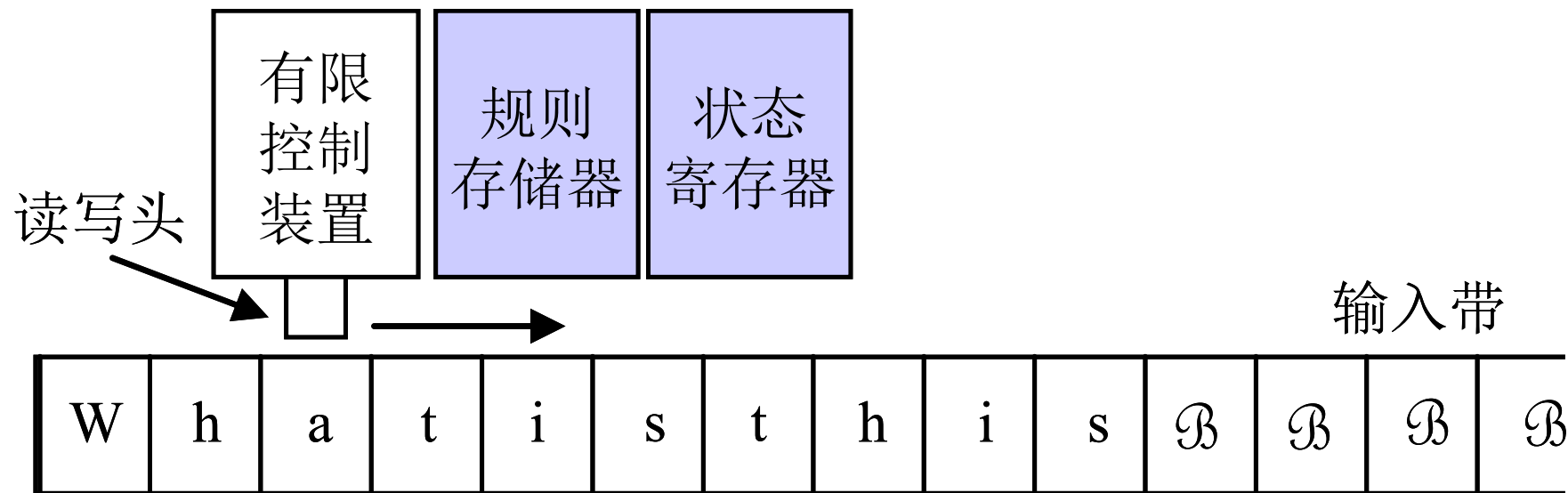
$F \subseteq Q$ 终止状态集合

$$\delta : Q \times (\Sigma_I \cup \{\lambda\}) \times \Gamma \rightarrow 2^{Q \times \Gamma^*}$$


$$\delta(q, a, Z) = \{(p_1, \gamma_1), (p_2, \gamma_2), \dots, (p_m, \gamma_m)\}$$

图灵机

Turing Machine (TM)



- 带读写头并可左右移动的有限状态自动机！
- 有限状态自动机和扩展的下推自动机的合体！

图灵机

定义 一个确定的图灵机 (TM) 是一个七元式:

$$A_T = (Q, \Sigma_I, \mathfrak{B}, \Gamma, \delta, q_0, F)$$

Q 状态的有限集合

Γ 可出现在输入带上的带符号集合

$\Sigma_I \subseteq \Gamma - \{\mathfrak{B}\}$ 输入符号的有限集合

\mathfrak{B} 空白符号

$q_0 \in Q$ 初始状态

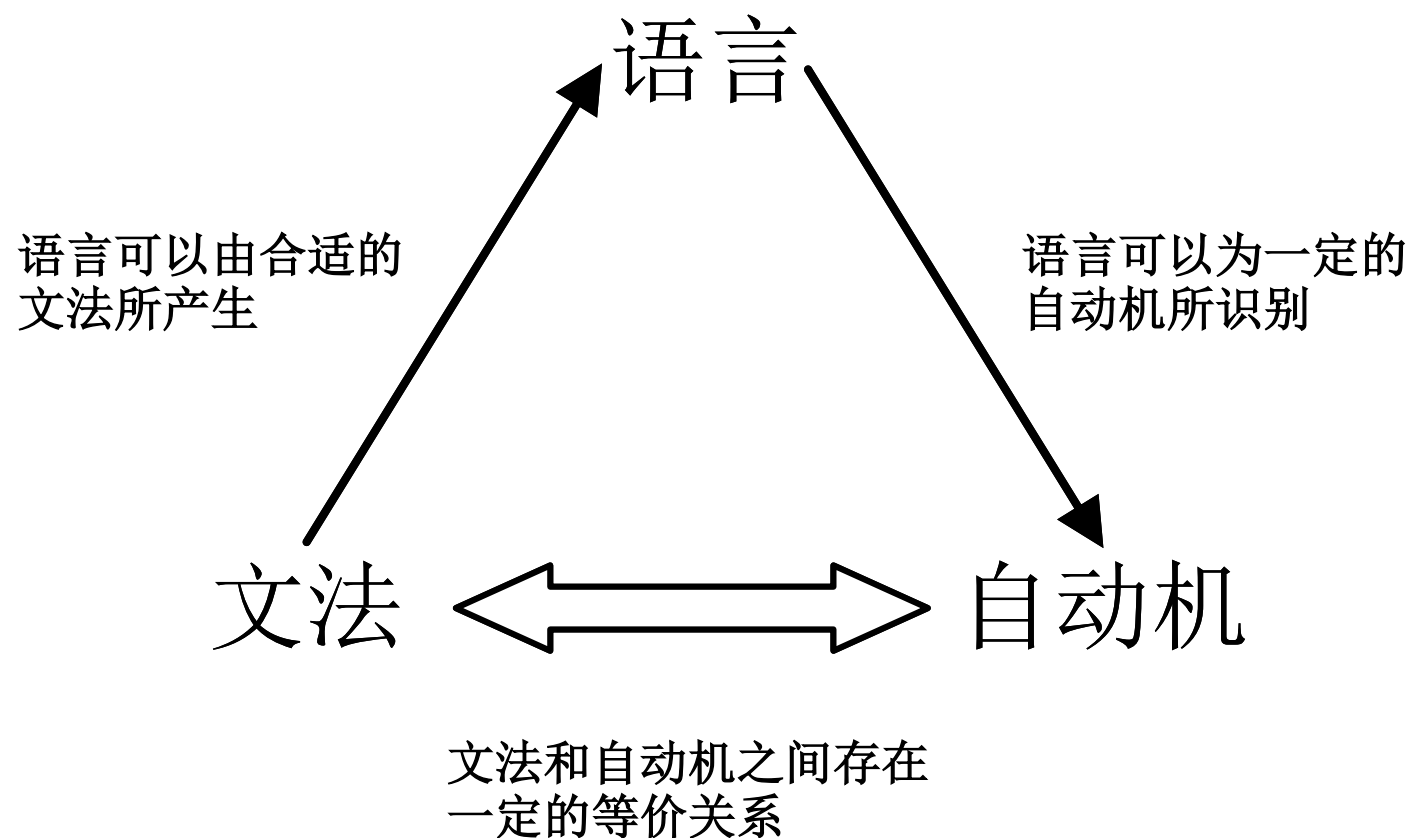
$F \subseteq Q$ 终止状态集合

$\delta: Q \times \Gamma \rightarrow Q \times \Gamma \times \{L, R\}$



$\delta(q, X) = (p, Y, L)$ 或 $\delta(q, X) = (p, Y, R)$

语言、文法和自动机三者之间的关系



句法分析

基于三角表格的反向剖析算法（CYK算法）

适用对象：以**乔姆斯基范式**形式表示的上下文无关文法

设待识别的输入符号串为 $x = a_1a_2 \cdots a_n$

1. 单个符号的派生

$A_{i,1}, 1 \leq i \leq n$ 可派生 a_i 的非终结符

寻找 $A_{i,1} \rightarrow a_i, 1 \leq i \leq n$, 若有, 记录 $A_{i,1}$, 若无, 记 $A_{i,1} = \phi$ 。

2. 两个符号组成的子串的派生

$A_{i,2}, 1 \leq i \leq n-1$ 可派生 $x_{i,2} = a_i a_{i+1}, 1 \leq i \leq n-1$ 的非终结符

寻找 $A_{i,2} \rightarrow A_{i,1} A_{i+1,1}$, 若有, 记录 $A_{i,2}$, 若无, 记 $A_{i,2} = \phi$ 。

3. 三个符号组成的子串的派生

$A_{i,3}, 1 \leq i \leq n-2$ 可派生 $x_{i,3} = a_i a_{i+1} a_{i+2}, 1 \leq i \leq n-2$ 的非终结符

寻找 $A_{i,3} \rightarrow A_{i,1} A_{i+1,2}$, 若有, 记录 $A_{i,3}$, 若无, 记 $A_{i,3} = \phi$ 。
 $A_{i,3} \rightarrow A_{i,2} A_{i+2,1}$

句法分析

基于三角表格的反向剖析算法（CYK算法）

.....

n. 输入符号串 x 的派生

$A_{1,n}$ 可派生 $x = a_1 a_2 \cdots a_n$ 的非终结符

寻找 $\left\{ \begin{array}{l} A_{1,n} \rightarrow A_{1,1} A_{2,n-1} \\ A_{1,n} \rightarrow A_{1,2} A_{3,n-2} \\ \dots\dots\dots \\ A_{1,n} \rightarrow A_{1,n-1} A_{n,1} \end{array} \right. , \text{若有, 记录 } A_{1,n}, \text{若无, 记 } A_{1,n} = \phi.$

最终判决: 若 $A_{1,n}$ 的可能取值中包含 S ,
则表示 $x \in L(G)$, 否则 $x \notin L(G)$ 。

文法推断

正则文法的推断

$$R^+ = \{X_1, X_2, \dots, X_m\} \Rightarrow G = (N, T, P, S)$$

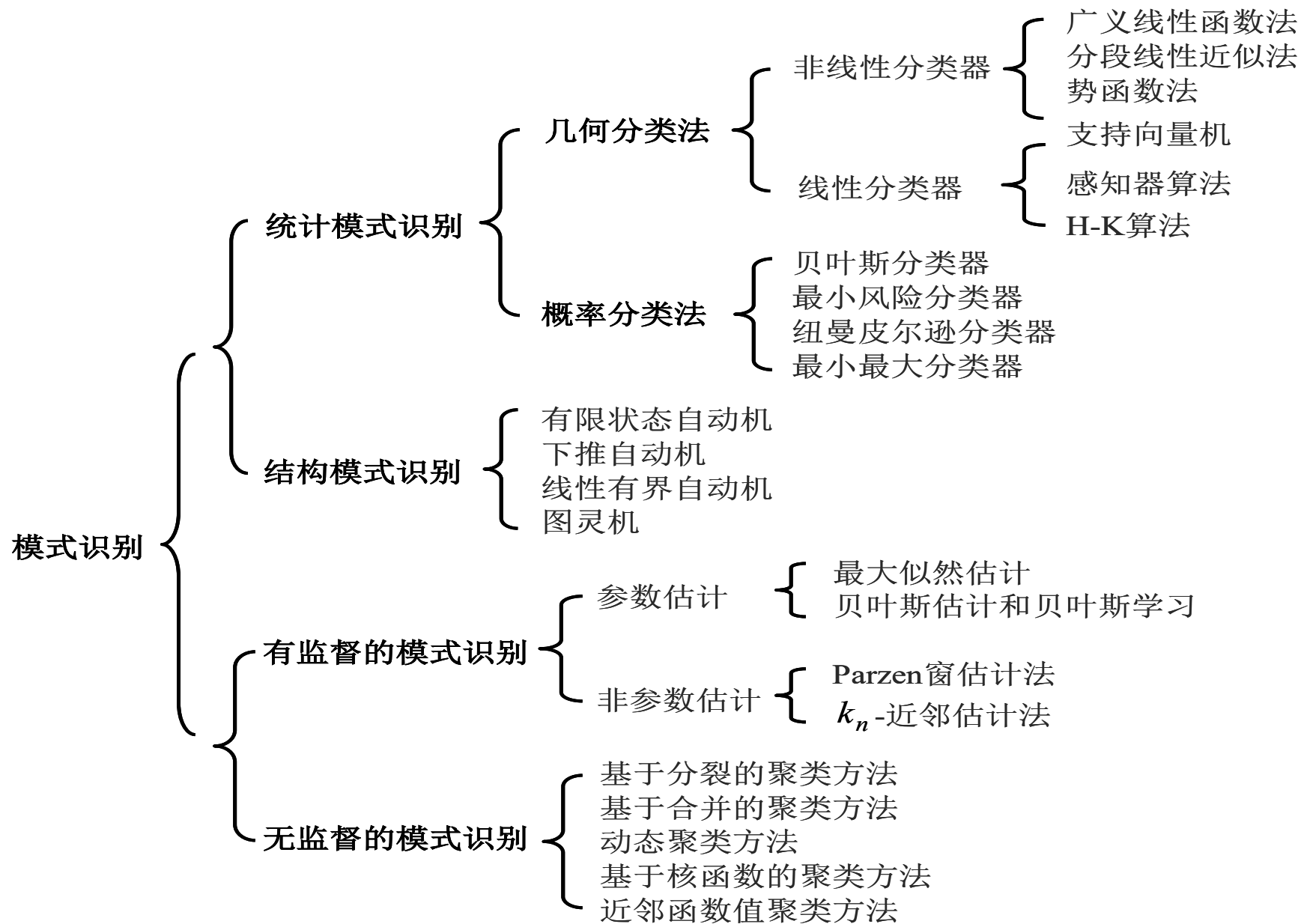
推断步骤

1. 列出 R^+ 中的所有终结符，令其为终结符集合。
2. 对 R^+ 中的符号串 $X_i, i = 1, 2, \dots, m$ ，适当引入非终结符，构建恰好能产生 X_i 的产生式组。
3. 对所得文法中的非终结符和产生式组，进行必要的化简和合并，最终得到所推断的文法。

$$X_i = a_{i1}a_{i2} \cdots a_{in_i}$$

$$S \rightarrow a_{i1}Z_{i1} \quad Z_{i1} \rightarrow a_{i2}Z_{i2}$$

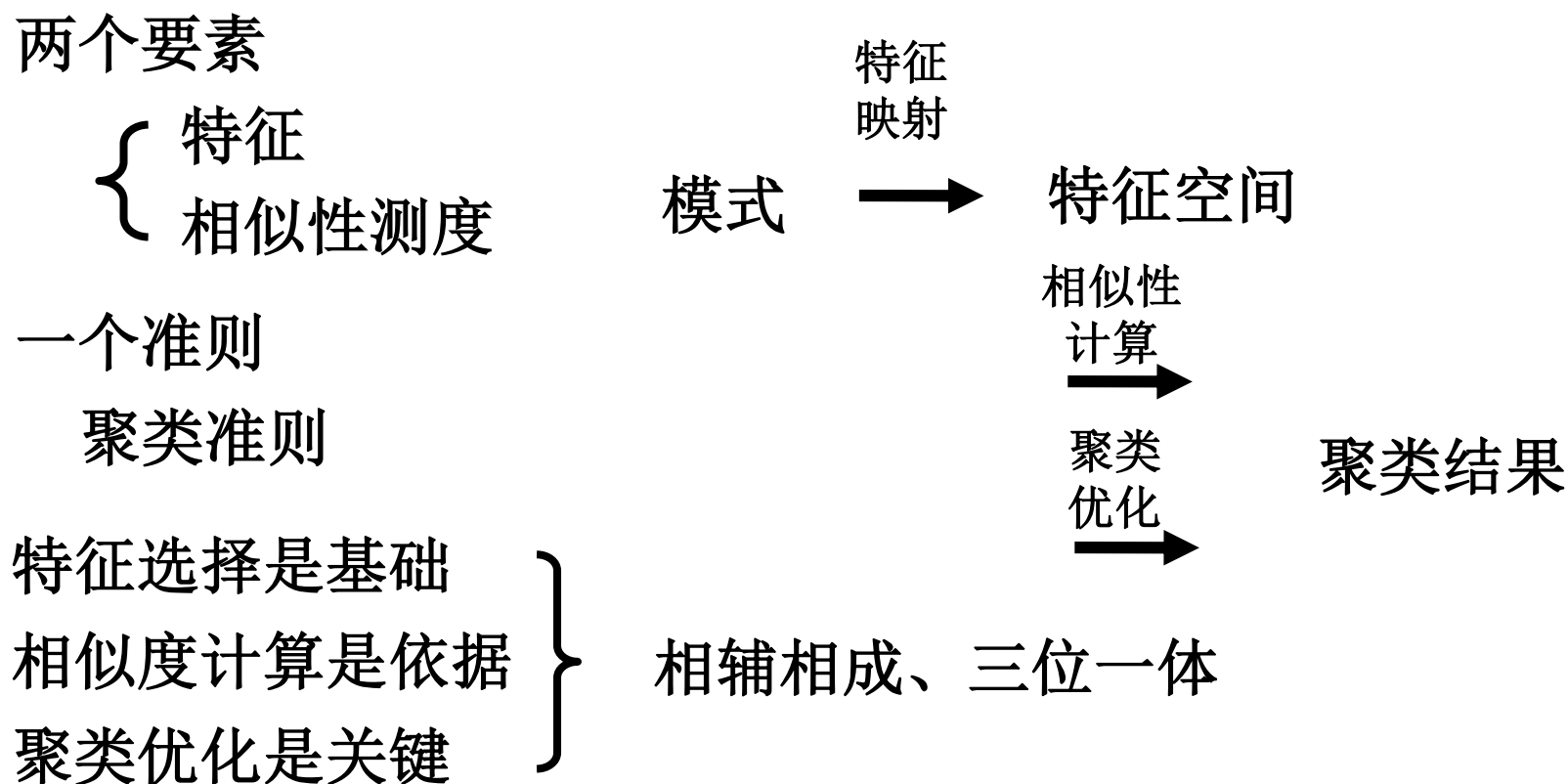
$$\dots\dots\dots Z_{in_i-2} \rightarrow a_{in_i-1}Z_{in_i-1} \quad Z_{in_i-1} \rightarrow a_{in_i}$$



- 聚类准则
- 基于分裂的聚类算法
- 基于合并的聚类算法
- 动态聚类算法
- 近邻函数值准则聚类算法
- 最小张树聚类算法

聚类方法

聚类方法：对于给定特征的两个样本，依照相似性测度计算其相似性，若相似性的度量值大于给定的阈值，则判它们属于同一个类别，否则判它们属于不同的类别。





中国科学技术大学