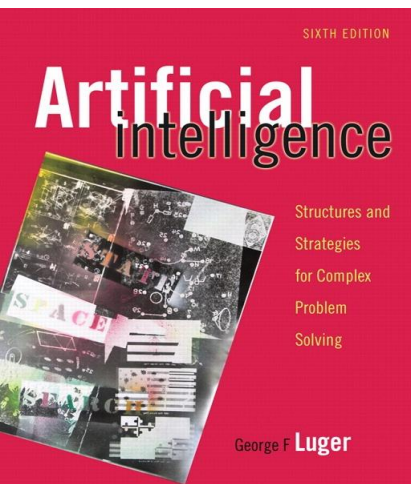


10

Machine Learning: Symbol-Based

10.0	Introduction	10.5	Knowledge and Learning
10.1	A Framework for Symbol-Based Learning	10.6	Unsupervised Learning
		10.7	Reinforcement Learning
10.2	Version Space Search	10.8	Epilogue and References
10.3	The ID3 Decision Tree Induction Algorithm	10.9	Exercises
10.4	Inductive Bias and Learnability		



George F Luger

ARTIFICIAL INTELLIGENCE *6th edition*

Structures and Strategies for Complex Problem Solving

什么是机器学习？

- 机器学习是研究如何使用计算机来模拟人类学习活动的—门学科。——研究计算机获取新知识和新技能、识别现有知识、不断改善性能、实现自我完善的方法。
- 学习是一个有特定目的的知识获取过程，其内在行为是获取知识、积累经验、发现规律，其外在表现是使系统性能得到改进、系统实现自我完善、自适应环境。
- 机器学习是智能系统不断地**积累经验以改善系统性能**的过程。

机器学习的方法

- ❖ 归纳学习
- ❖ 类比学习
- ❖ 基于解释的学习
- ❖ 遗传算法
- ❖ 人工神经网络
- ❖ 随机方法

10 MACHINE LEARNING: SYMBOL-BASED

10.0 Introduction

Section 10.1 outlines a variety of learning tasks.

Section 10.2 and 10.3 examine two algorithms used for concept induction, version space search and ID3.

Section 10.4 considers the role of inductive bias in learning.

The section before 10.6 are supervised learning.

Section 10.6 continues the study of induction by examining unsupervised learning.

In 10.7, we present reinforcement learning.

10.1 A Framework for Symbol-Based Learning

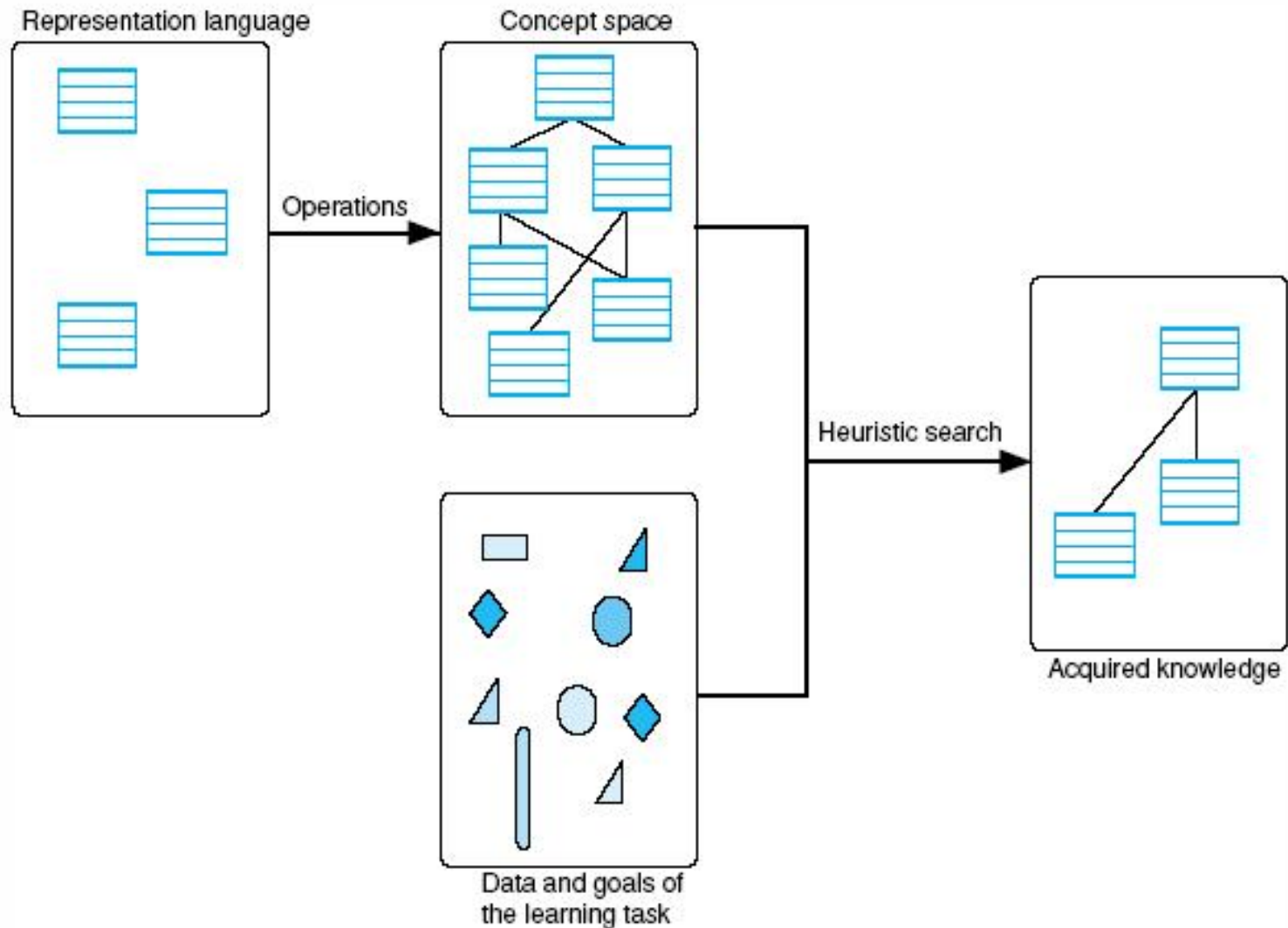


Fig 10.1 A general model of the learning process

1.The data and goals of the learning task.

The concept learning algorithms of section 10.2 and 10.3,begin with a collection of **positive examples** of a **target class**,the goal is to infer a general definition that will allow the learner to recognize future instances of the class.

Explanation-based learning(Section 10.5),attempts to infer a general concept from a single training example and a prior base of domain-specific knowledge.

These algorithms begin with a set of unclassified instances in section 10.6.

2.The representation of learned knowledge

- For examples , two instances of a ball may be represented by:

$\text{Size}(\text{obj1}, \text{small}) \wedge \text{color}(\text{obj1}, \text{red}) \wedge \text{shape}(\text{obj1}, \text{round})$

$\text{Size}(\text{obj2}, \text{large}) \wedge \text{color}(\text{obj2}, \text{red}) \wedge \text{shape}(\text{obj2}, \text{round})$

The general concept of ball could be defined by:

$\text{Size}(X, Y) \wedge \text{color}(X, Z) \wedge \text{shape}(X, \text{round})$

3.A set of operations.

$\text{Size}(\text{obj1}, \text{small}) \wedge \text{color}(\text{obj1}, \text{red}) \wedge \text{shape}(\text{obj1}, \text{round})$

Replacing a single constant with a variable produces the generalizations.

$\text{Size}(\text{obj1}, x) \wedge \text{color}(\text{obj1}, \text{red}) \wedge \text{shape}(\text{obj1}, \text{round})$

$\text{Size}(\text{obj1}, \text{small}) \wedge \text{color}(\text{obj1}, x) \wedge \text{shape}(\text{obj1}, \text{round})$

$\text{Size}(\text{obj1}, \text{small}) \wedge \text{color}(\text{obj1}, \text{red}) \wedge \text{shape}(\text{obj1}, x)$

$\text{Size}(x, \text{small}) \wedge \text{color}(x, \text{red}) \wedge \text{shape}(x, \text{round})$

4.The concept space .

The representation language, together with the operations described above, defines a space of potential concept definitions.

5.Heuristic search.

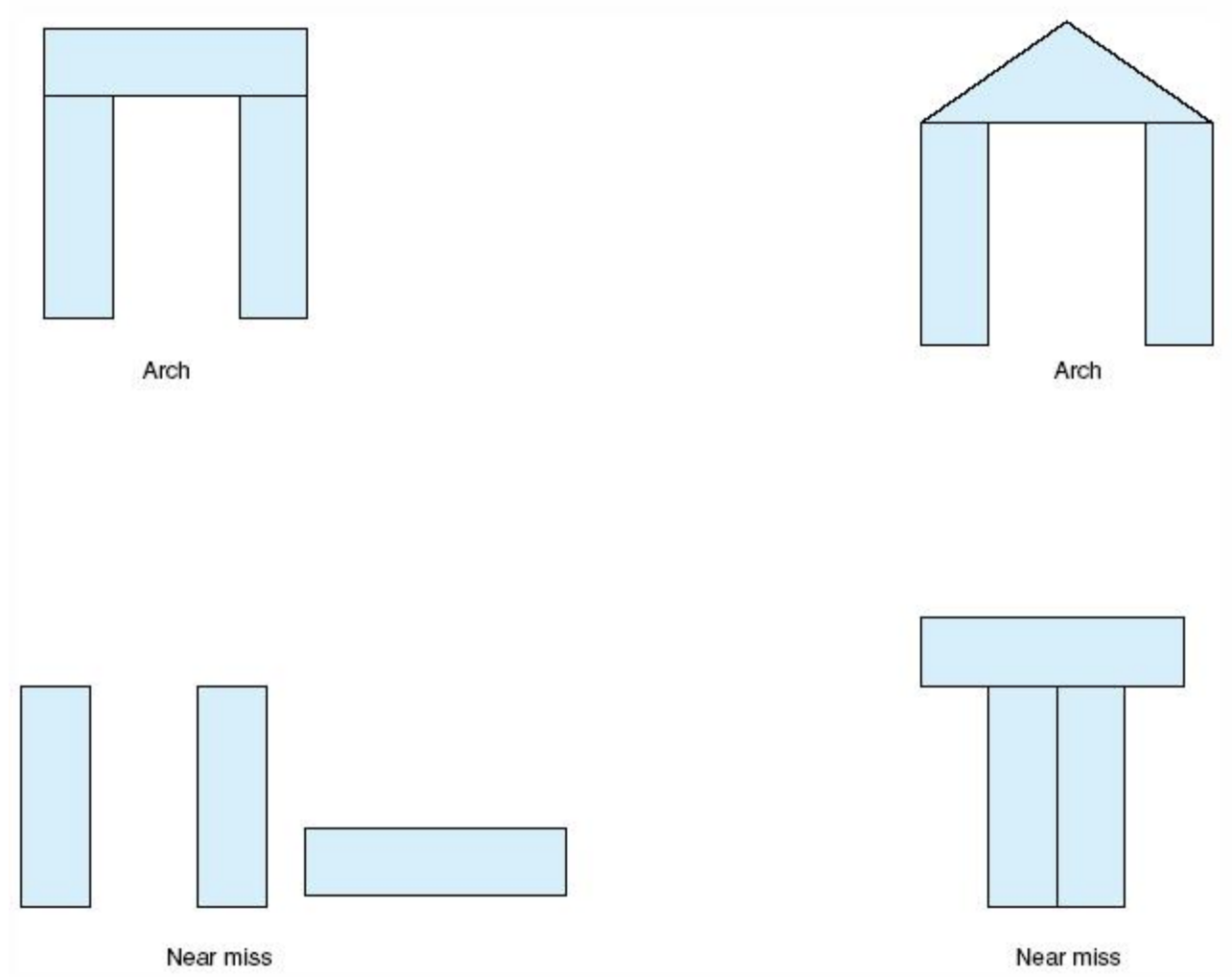
$\text{Size}(\text{obj1}, \text{small}) \wedge \text{color}(\text{obj1}, \text{red}) \wedge \text{shape}(\text{obj1}, \text{round})$

$\text{Size}(\text{obj2}, \text{large}) \wedge \text{color}(\text{obj2}, \text{red}) \wedge \text{shape}(\text{obj2}, \text{round})$

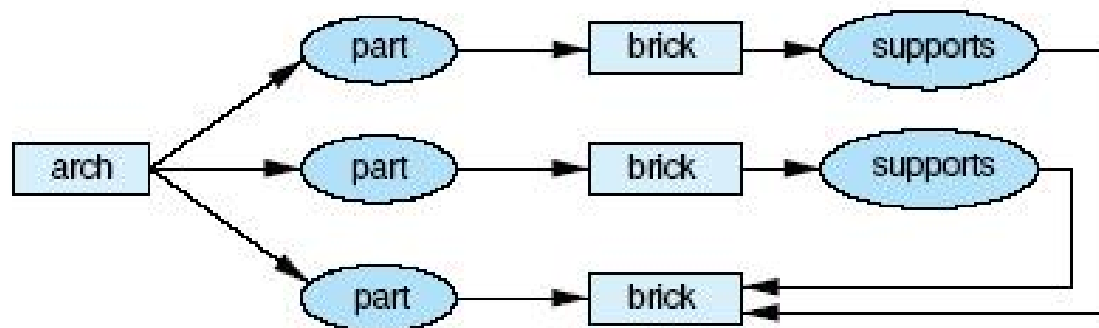
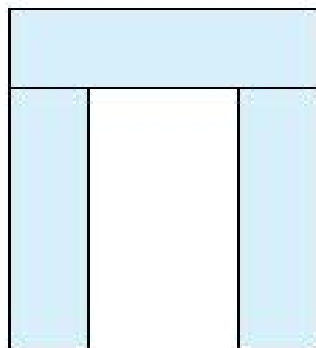
$\text{Size}(X, Y) \wedge \text{color}(X, \text{red}) \wedge \text{shape}(X, \text{round})$

Patrick Winston's program learns general definitions of structural concepts, such as “**arch**” (拱形), in a blocks world. The training data is a series of **positive and negative examples** (正例、反例) of the concept: examples of blocks world structures that fit in the category, along with **near misses** (小差别).

Fig 10.2 Examples and near misses for the concept “arch.”



a. An example of an arch and its network description



b. An example of another arch and its network description

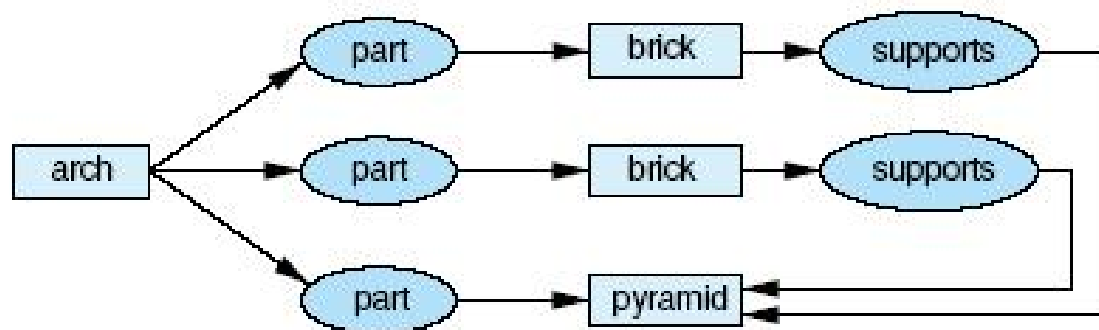
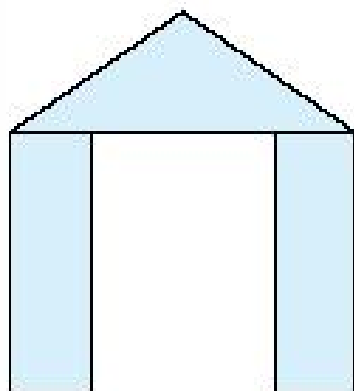
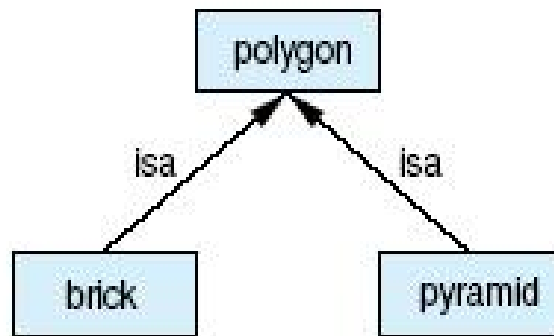


Fig 10.3 generalization of descriptions to include multiple examples (*cont'd*)

c. Given background knowledge that bricks and pyramids are both types of polygons



d. Generalization that includes both examples

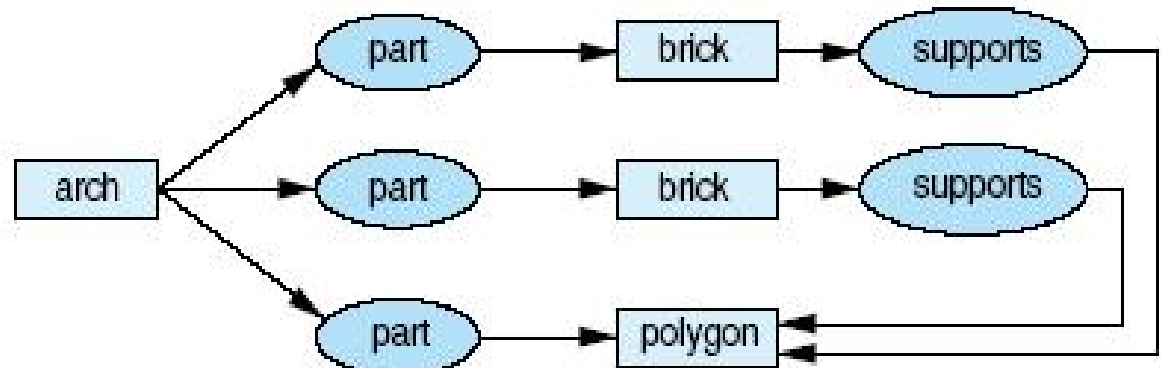
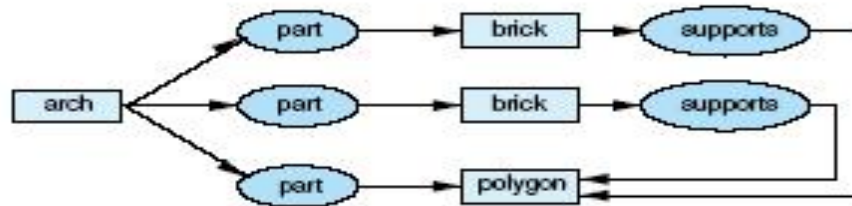
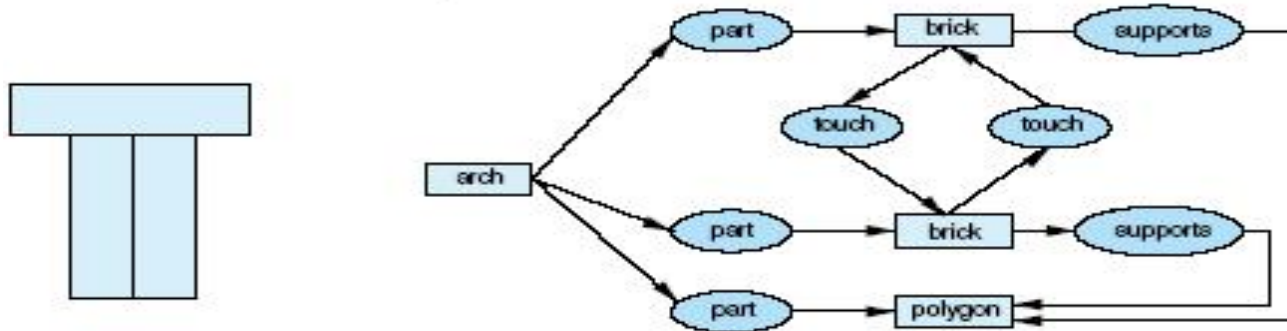


Fig 10.4 Specialization of a description to exclude a near miss. In 10.4c we add constraints to 10.4a so that it can't match with 10.4b.

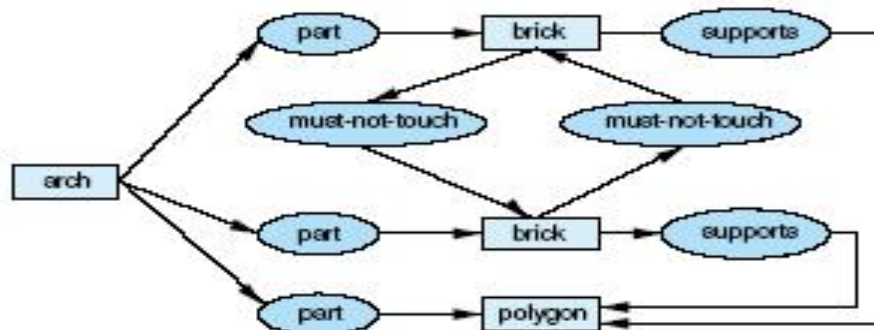
a. Candidate description of an arch



b. A near miss and its description



c. Arch description specialized to exclude the near miss



10.2 Version Space Search

Version space search illustrates the implementation of inductive learning as search through a concept space. Version space search takes advantage of the fact that generalization operations impose an ordering on the concepts in a space, and then uses this ordering to guide the search.

10.2.1 Generalization Operators and the Concept Space

Generalization and specialization are the most common types of operations for defining a concept space. The primary generalization operations used in machine learning are:

1. Replacing constants with variables. For example,

$\text{color}(\text{ball}, \text{red})$ generalizes to $\text{color}(X, \text{red})$

2. Dropping conditions from a conjunctive expression.

$\text{shape}(X, \text{round}) \wedge \text{size}(X, \text{small}) \wedge \text{color}(X, \text{red})$

generalizes to $\text{shape}(X, \text{round}) \wedge \text{color}(X, \text{red})$

3. Adding a disjunct to an expression.

$\text{shape}(X, \text{round}) \wedge \text{size}(X, \text{small}) \wedge \text{color}(X, \text{red})$

generalizes to

$\text{shape}(X, \text{round}) \wedge \text{size}(X, \text{small}) \wedge \text{color}(X, \text{red})$
 $\vee \text{color}(X, \text{blue})$

4. Replacing a property with its parent in a class hierarchy.

$\text{color}(X, \text{red})$ **generalizes to** $\text{color}(X, \text{primary_color})$

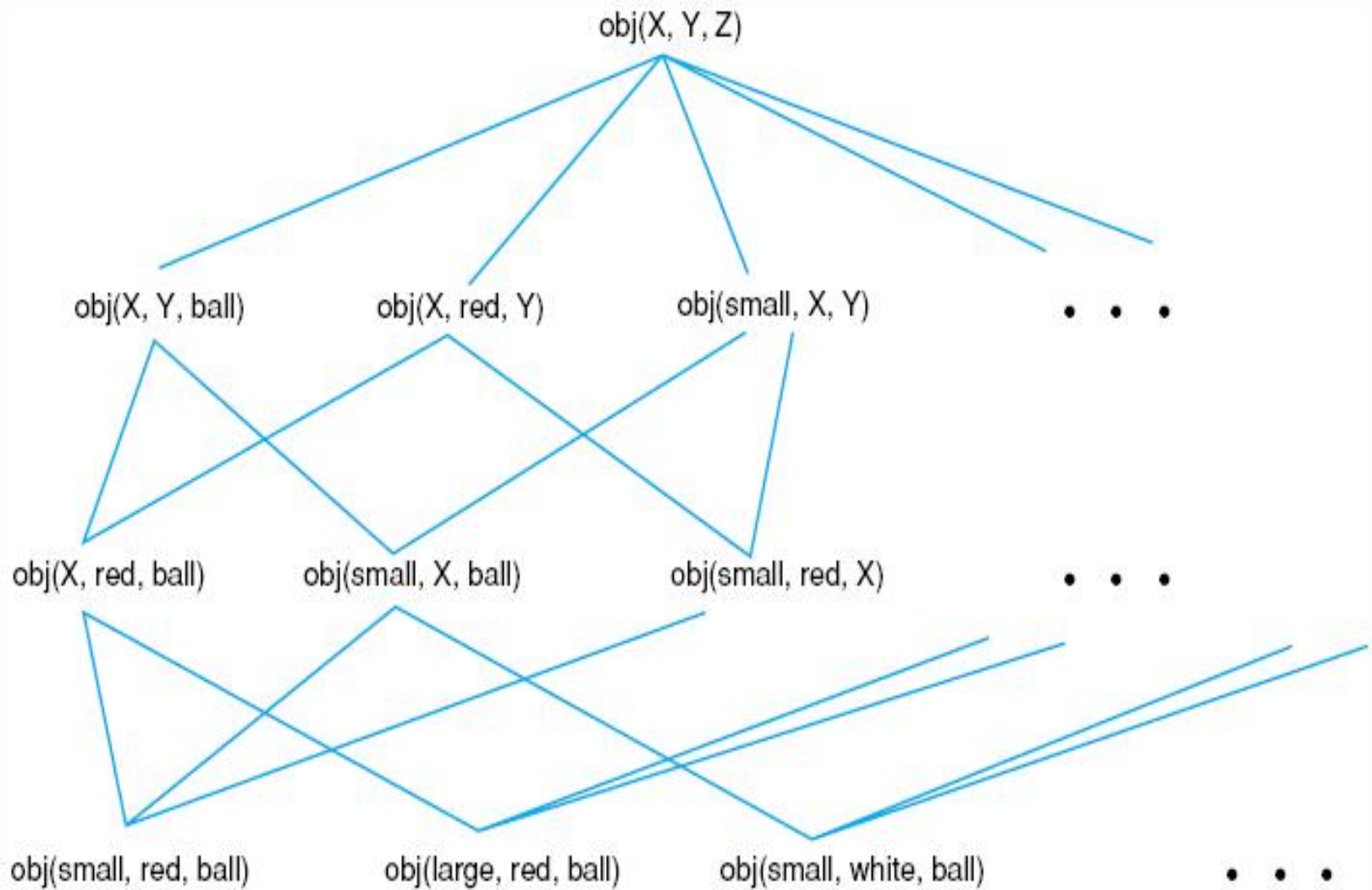
If concept p is more general than concept q , we say that p covers q ($p \supseteq q$). $p \supseteq q$ means that p is **more general** than q .

we formalize this relationship through the notion of covering. For an object X , $p(x) \rightarrow \text{positive}(x)$ and $q(x) \rightarrow \text{positive}(x)$. p covers q if $q(x) \rightarrow \text{positive}(x)$ is a logical consequence of $p(x) \rightarrow \text{positive}(x)$.

color(X,Y) covers color(ball,Z) covers color(ball,red)

size = {large, small}, **color** = {red, white, blue}, **shape** = {ball, brick, cube}, these objects can be represented using the predicate **obj**(size, color, shape).

Fig 10.5 A concept space.



10.2.2 The Candidate Elimination Algorithm

This section presents **three algorithms** for searching the concept space. The first two algorithms reduce the version space in a general to specific direction, respectively. The third algorithm, called candidate elimination, combines these approaches into a bi-directional search.

Defining *specific to general search*, for hypothesis set S as:

Begin

Initialize S to the first positive training instance;

N is the set of all negative instances seen so far;

For each positive instance p

Begin

For every $s \in S$, if s does not match p, replace s with its most specific generalization that matches p;

Delete from S all hypotheses more general than some other hypothesis in S;

Delete from S all hypotheses that match a previously observed negative instance in N;

End;

For every negative instance n

Begin

Delete all members of S that match n;

Add n to N to check future hypotheses for overgeneralization;

End;

End

In this algorithm, negative instances lead to the **specialization** of candidate concepts; the algorithm uses positive instances to eliminate overly specialized concepts.

Begin

Initialize G to contain the most general concept in the space;

P contains all positive examples seen so far;

For each negative instance n

Begin

For each $g \in G$ that matches n, replace g with its most general specializations that do not match n;

Delete from G all hypotheses more specific than some other hypothesis in G;

Delete from G all hypotheses that fail to match some positive example in P;

End;

For each positive instance p

Begin

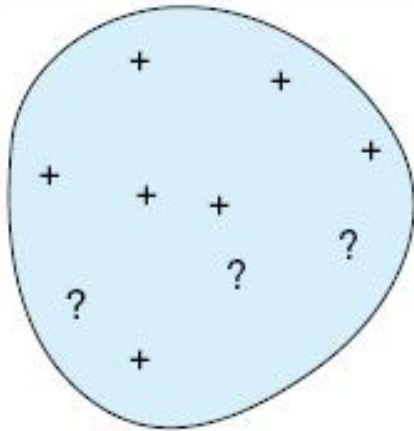
Delete from G all hypotheses that fail to match p;

Add p to P;

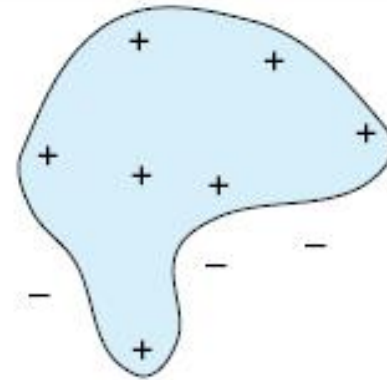
End;

End

Fig 10.6 The role of negative examples in preventing overgeneralization.

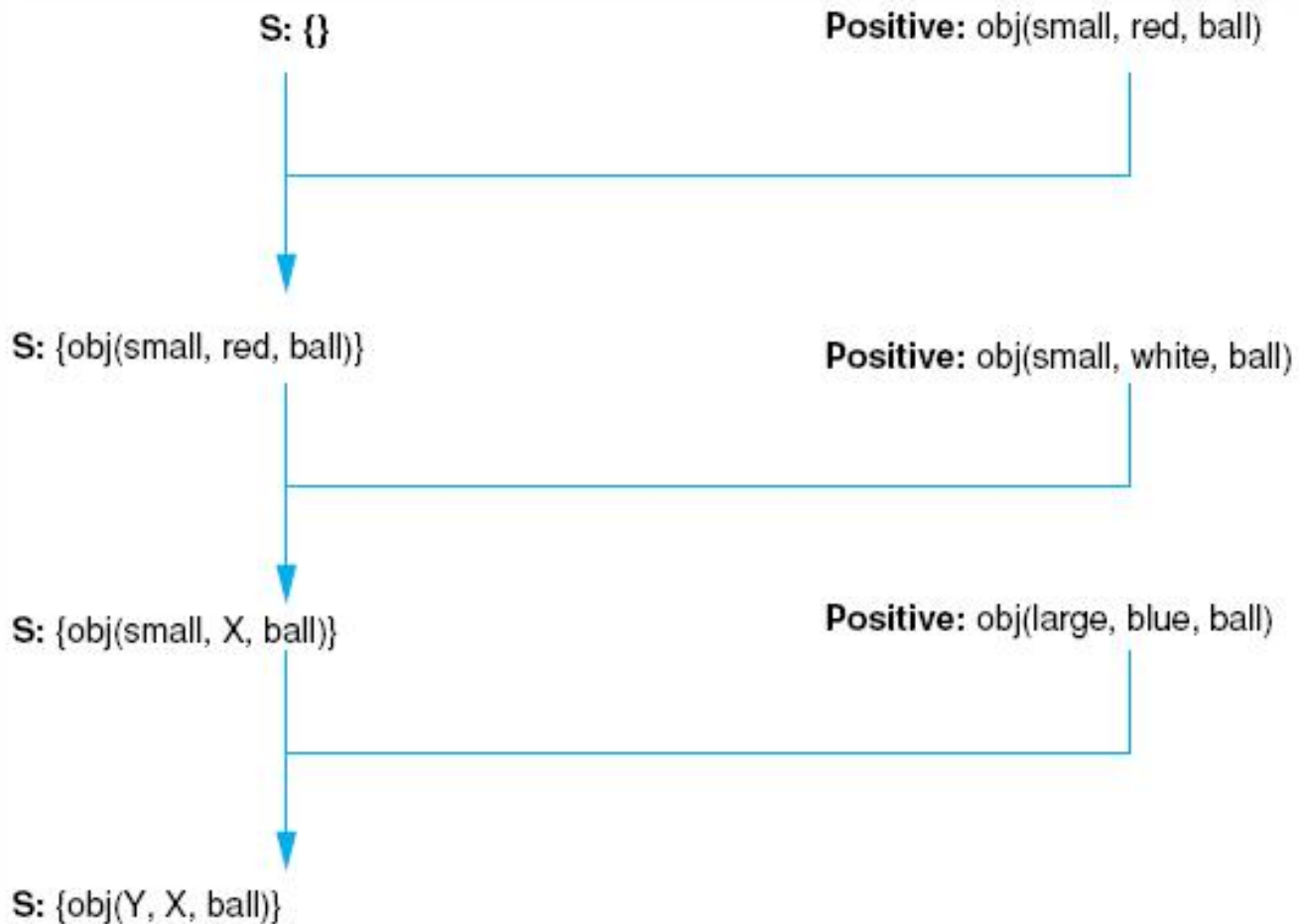


Concept induced from
positive examples only



Concept induced from
positive and negative examples

Fig 10.7 Specific to general search of the version space learning the concept “ball.”



The algorithm specializes **G** and generalizes **S** until they converge on the target concept. The algorithm is defined:

Begin

Initialize G to be the most general concept in the space;

Initialize S to the first positive training instance;

For each new positive instance p

Begin

Delete all members of G that fail to match p;

For every $s \in S$, if s does not match p, replace s with its most specific generalizations that match p;

Delete from S any hypothesis more general than some other hypothesis in S;

Delete from S any hypothesis more general than some hypothesis in G;

End;

For each new negative instance n

Begin

Delete all members of S that match n;

For each $g \in G$ that matches n, replace g with its most general specializations that do not match n;

Delete from G any hypothesis more specific than some other hypothesis in G;

Delete from G any hypothesis more specific than some hypothesis in S;

End;

Fig 10.8 General to specific search of the version space learning the concept “ball.”

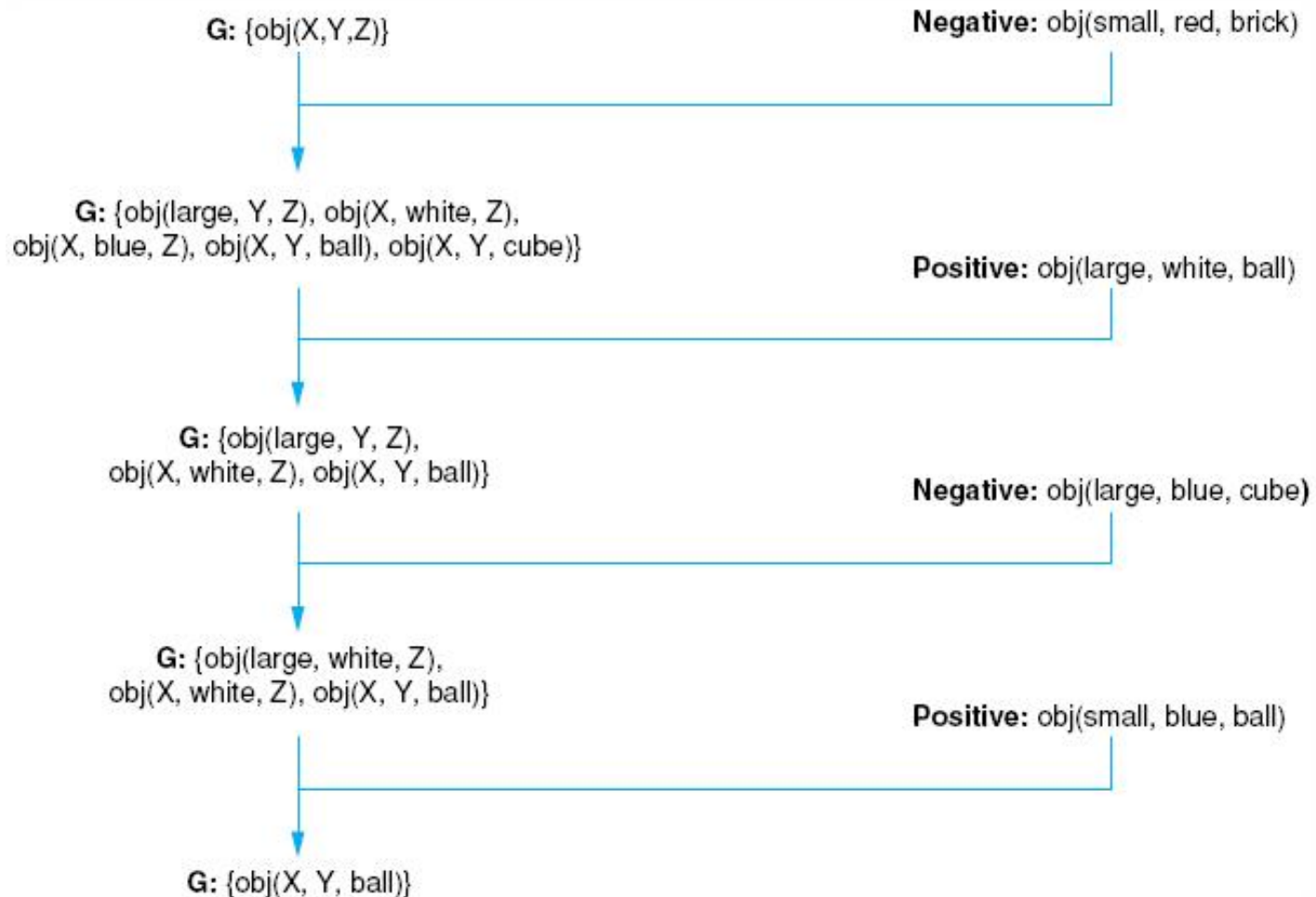


Fig 10.9 The candidate elimination algorithm learning the concept “red ball.”

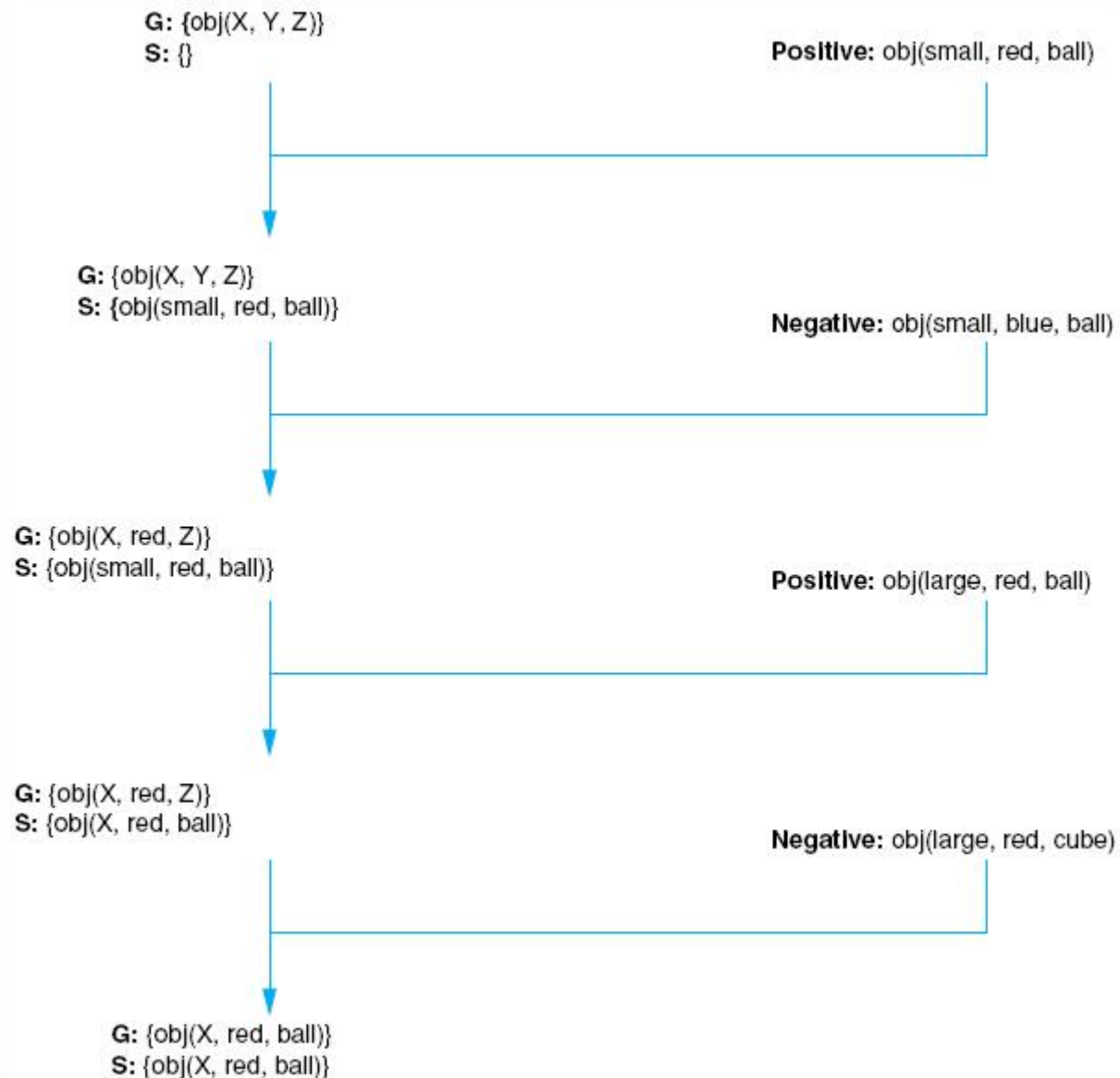
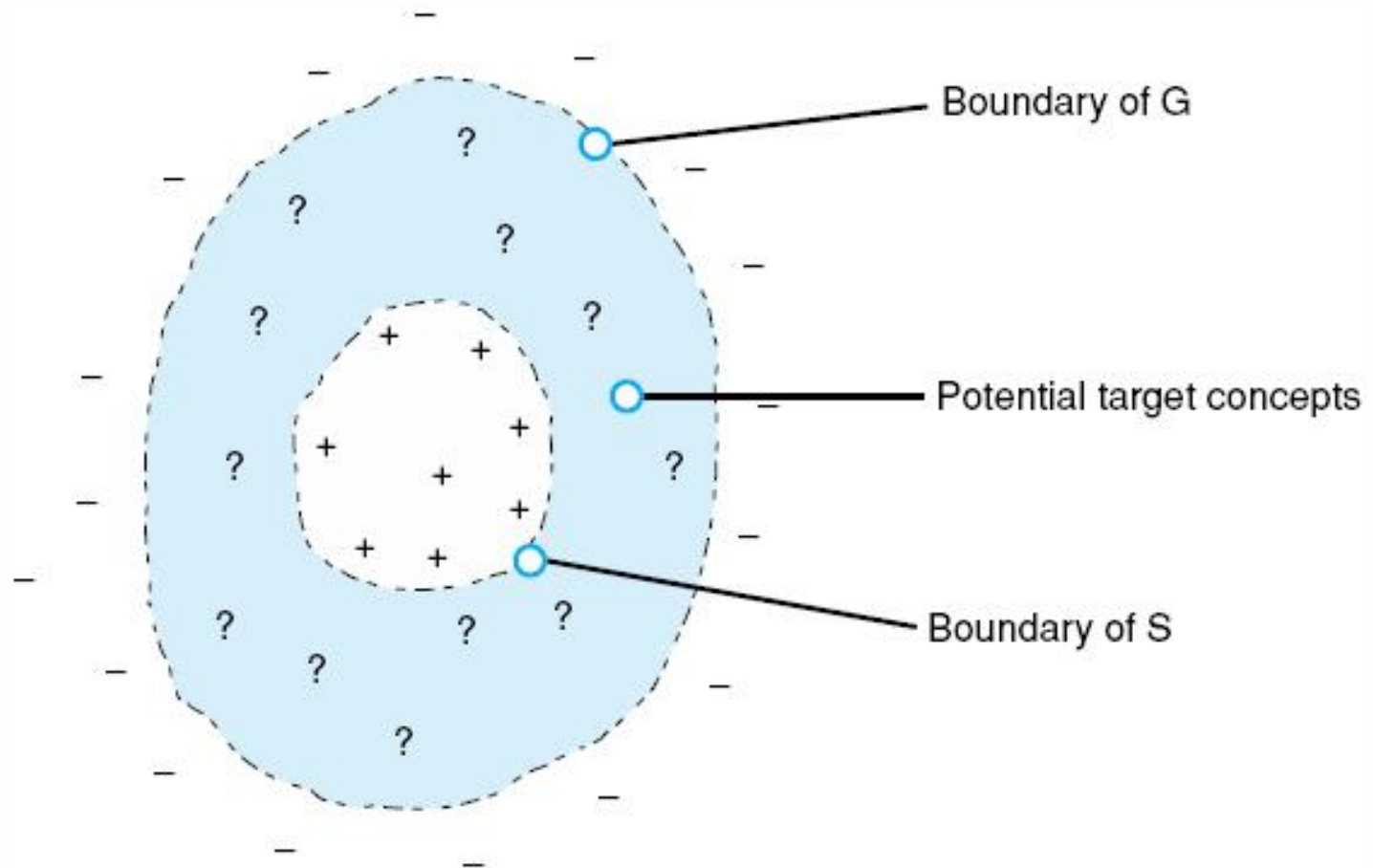


Fig 10.10 Converging boundaries of the G and S sets in the candidate elimination algorithm.



基于描述空间的归纳学习方法

- 在Mitchell提出的描述空间方法中，根据概念的一般性和特殊性把概念空间表示成偏序集的形式。
- 偏序集的最大元素 G 称为零描述。
- 零描述是概念空间中恒真的概念描述，是概念中最一般的概念的集合。
- 偏序集的最小元素 S 是所有训练实例，是概念空间最特殊的概念集合。

例 1

设实例空间是一组实体，每个实体用“大小”和“形状”两个特征描述，若“大小”属性值有：small、large，“形状”有属性值：square、circle。

对这个实例空间，请给出其概念空间的偏序集的最大元素和最小元素。

解: 每个实体有“大小”和“形状”两个属性, 可用 (x,y) 表示实体概念空间的最一般概念, x 描述实体的“大小”, y 描述实体的“形状”, (x,y) 是该实体概念空间中零描述。

偏序集的最大元素 $G=\{(x,y)\}$ 。

偏序集的最小元素 S 是所有的训练实例,

故有:

$S=\{(\text{small square}), (\text{small circle}), (\text{large, square}), (\text{large, small})\}$

- 对于给定的训练实例，利用偏序集可以非常紧凑地表示与训练实例相一致的合理假设集合（一组概念描述） H ， H 由两个子集 G 和 S 所限定， H 是由 G 和 S 以及 G 与 S 之间的元素构成的偏序集。
- 即 $H = G \cup S \cup \{K \mid S < K < G\}$ 其中，“ $<$ ”表示概念空间的偏序关系。
- H 中还有一组与训练实例相一致的概念描述：
 $K = \{(\text{small } y), (\text{large } y), (x \text{ square}), (x \text{ circle})\}$,
 G 、 K 、 S 的偏序关系为 $S < K < G$ 。

概念空间 $H = G \cup K \cup S$ 的偏序集表示如图4。

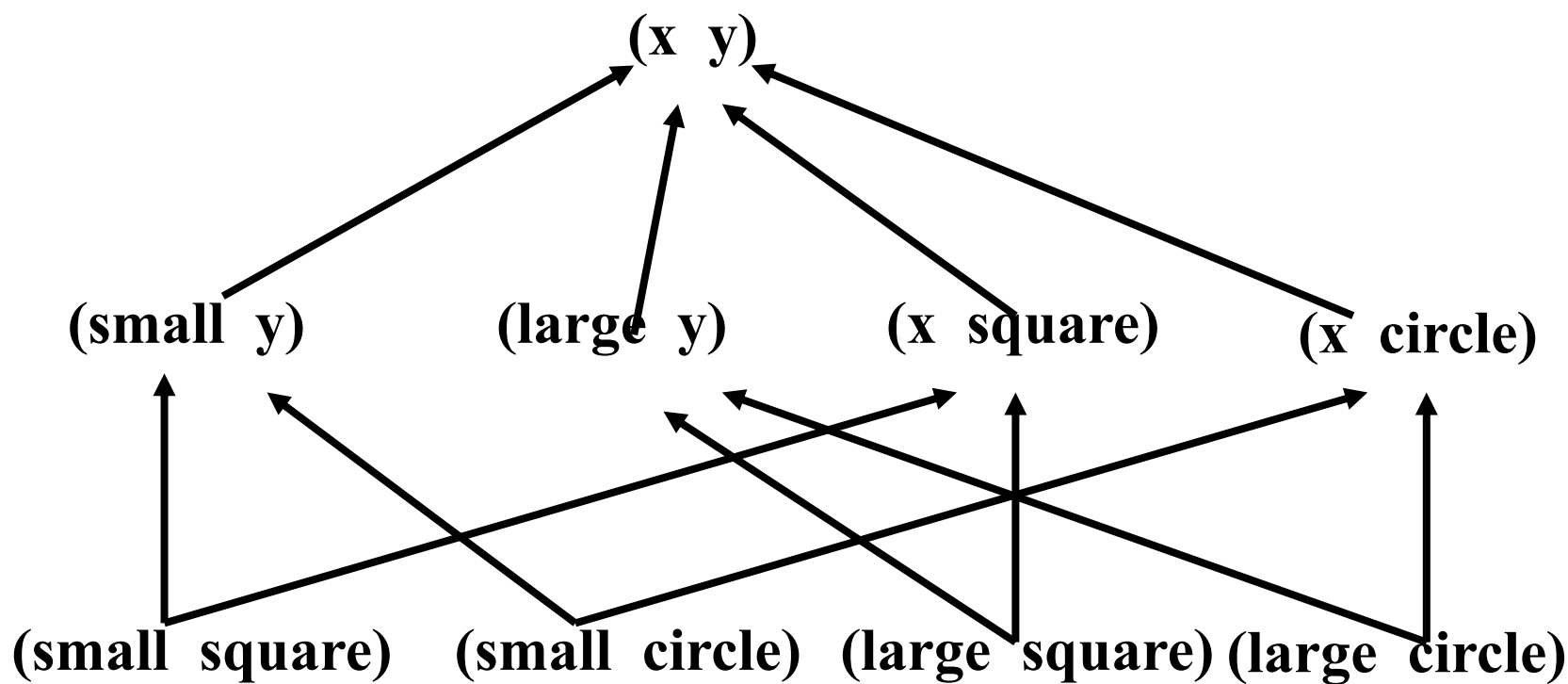


图4 概念空间的偏序集表示

消除候选者算法

- (1) 初始化概念空间 $H = \langle C, S \rangle$ ，其中， C 为零描述集合， S 为训练集合，初始化时， S 中只含第一个正例。
- (2) 接受一个新的训练实例。如果实例为**正例**，则从 C 中删去所有与该例**不相容**的概念，并更新集合 S ，**尽可能小地对 S 进行一般化**，以**相容**这个新的正实例；如果实例为反例，则首先从 S 中删去所有与该例相容的概念，并更新集合 C ，**尽可能小地特殊 C 中的元素**，以便它们**不相容**这个反例。
- (3) 如果 $C \neq S$ ，则重复步骤 (2)；否则输出 H 。

[例2]对例1给出的概念空间，应用消除候选者学习算法，说明学习得到一般概念为“圆”的学习步骤和学习结果。

注意以下术语：

- 概念空间
- 最大元素、最小元素
- 训练实例、正例、反例
- 相容的概念、不相容的概念
- 概念的一般化、特殊化

解: 由例1给出的概念空间可得知, “圆”的一般概念显然表示为(x circle)。这一概念可采用消除候选者算法获得, 过程如下:

1) 初始化概念空间 $H = \langle G, S \rangle$, 其中

$$G = \{(x \ y)\}$$

$$S = \{(\text{small square}), (\text{small circle}), \\ (\text{larges, square}), (\text{large, circle})\}$$

2) 若首先提供一个正例(small circle), 由消除候选者算法步骤(2)对 $H=\langle G, S \rangle$ 进行更新; $G=\{(x\ y)\}$ 中的概念 $(x\ y)$ 是最一般概念, 与提供的正例(small circle)相容, 故仍有 $G=\{(x\ y)\}$ 。删除 S 中与提供正例不相容的元素(large square), 对 S 中的另外2个元素尽可能小地一般化, 即使这2个元素分别为(small y)和(x circle), 从而使它们与提供的正例(small circle)相容, 故有

$$S=\{(\text{small } y), (\text{small circle}), (x\ \text{circle})\}$$

3) 此时有 $G \neq S$ ，若再提供一个反例 (large square)，继续对 H 更新；由于 S 中的元素与提供的反例不相容，故 S 不变。对 G 中的元素 $(x\ y)$ 尽可能小地特殊化，可由 $(x\ y)$ 得到与提供反例不相容的2个元素 (small y) 和 $(x\ \text{circle})$ ，故有 $G = \{(\text{small } y), (x\ \text{circle})\}$

4) 此时有 $G \neq S$ ，若再提供一个正例 (large circle)，继续对 H 更新：删除 G 中与提供的正例 (large circle) 不相容元素 (small y)，故有 $G = \{(x\ \text{circle})\}$ ，删除 S 中与提供的正例 (large circle) 不相容元素 (small y) 和 (small circle)，故有 $S = \{(x\ \text{circle})\}$ 。

5) 此时有 $G = S = \{(x\ \text{circle})\}$ ，故算法终止，由概念空间 $H = \langle G, S \rangle$ 得出经过学习后获得的一般概念 $(x\ \text{circle})$ 。

基于变形空间的归纳学习1

设实例空间是一组实体，每一个实例用“大小”和“形状”两个特征描述。“大小”的属性值有：small、large，“形状”的属性值有：square、circle、triangle。针对这个实例空间，请回答下面问题：

- (1) 给出其概念空间偏序集的最大元素和最小元素。
- (2) $(x \text{ square})$ 和 $(\text{large } y)$ 分别表示什么概念？
- (3) 给出 $(x \text{ square})$ 的一个正例和一个反例。
- (4) 对 $(x \text{ } y)$ 尽可能小地特殊化后，可以得到哪些概念？
- (5) 给出所有与 $(x \text{ square})$ 不相容的实例。

基于变形空间的归纳学习2

- 考虑积木世界的以下属性和属性值：

颜色={黄色, 蓝色, 绿色}

形状={圆锥形, 球形, 长方形}

硬度={硬, 软}

尺寸={大, 小}

- 通过以下实例学习概念：黄色。

实例集为：{(黄色, 圆锥形, 软, 大, +),
(蓝色, 长方形, 软, 小, -),
(黄色, 球形, 软, 小, +),
(黄色, 圆锥形, 硬, 大, +),
(黄色, 长方形, 软, 大, +),
(绿色, 球形, 硬, 大, -),
(蓝色, 圆锥形, 软, 大, -)}

10.3 The ID3 Decision Tree Induction Algorithm

ID3, like candidate elimination, **induces concepts** from examples.

ID3 represents concepts as **decision trees**, a representation that allows us to determine the **classification** of an object by testing its values for certain **properties**.

10.3.1 Top-Down Decision Tree Induction

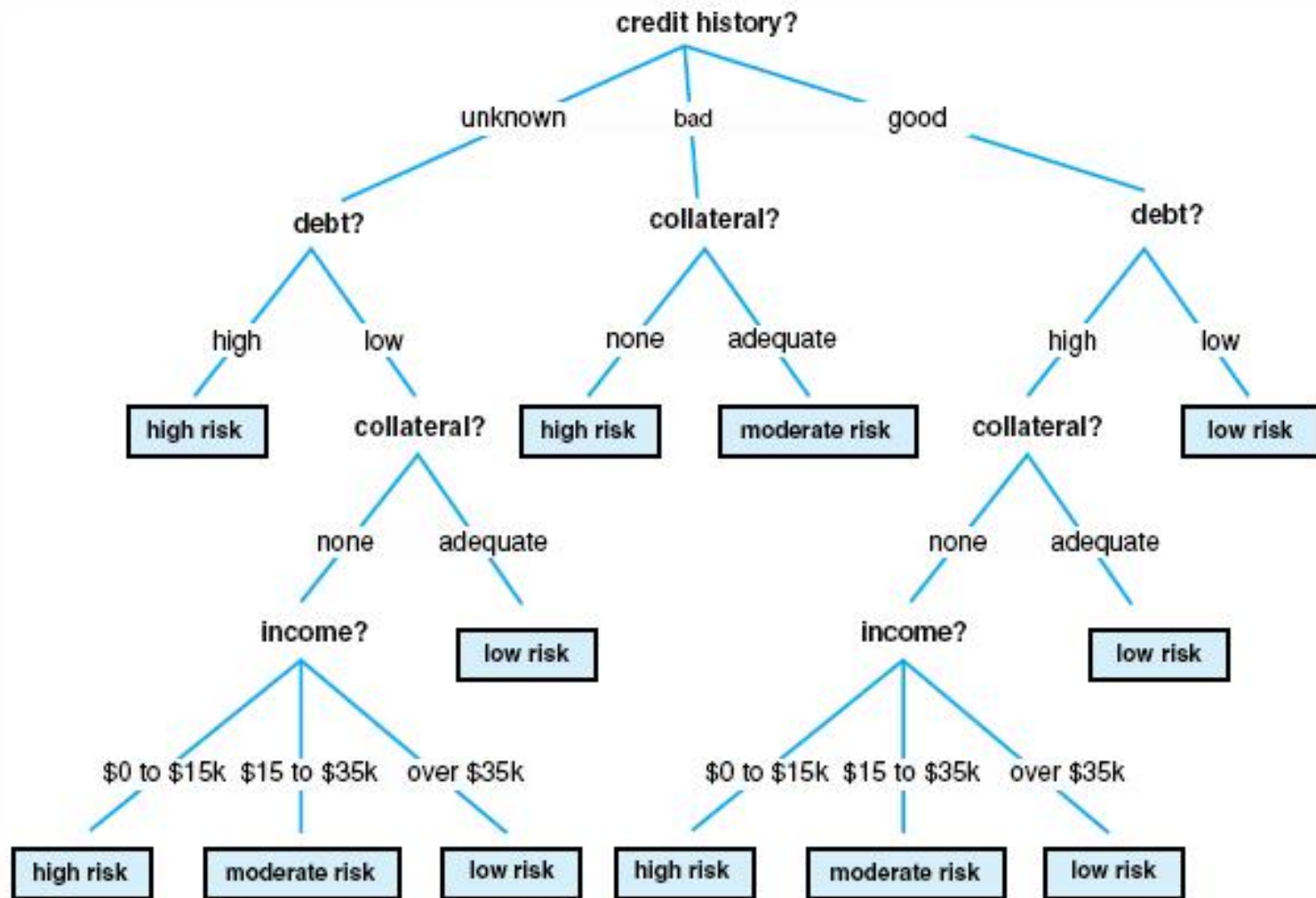
Table 10.1 Data from credit history of loan applications

NO.	RISK	CREDIT HISTORY	DEBT	COLLATERAL	INCOME
1.	high	bad	high	none	\$0 to \$15k
2.	high	unknown	high	none	\$15 to \$35k
3.	moderate	unknown	low	none	\$15 to \$35k
4.	high	unknown	low	none	\$0 to \$15k
5.	low	unknown	low	none	over \$35k
6.	low	unknown	low	adequate	over \$35k
7.	high	bad	low	none	\$0 to \$15k
8.	moderate	bad	low	adequate	over \$35k
9.	low	good	low	none	over \$35k
10.	low	good	high	adequate	over \$35k
11.	high	good	high	none	\$0 to \$15k
12.	moderate	good	high	none	\$15 to \$35k
13.	low	good	high	none	over \$35k
14.	high	bad	high	none	\$15 to \$35k

Table 10.1 Data from credit history of loan applications applications

NO.	风险	信用历史	债务	抵押	收入
1	高	坏	高	无	0-15000美元
2	高	未知	高	无	15000-35000美元
3	中等	未知	低	无	15000-35000美元
4	高	未知	低	无	0-15000美元
5	低	未知	低	无	超过35000美元
6	低	未知	低	无	超过35000美元
7	高	坏	坏	无	0-15000美元
8	中等	坏	坏	充分	超过35000美元
9	低	好	好	无	超过35000美元
10	低	好	好	充分	超过35000美元
11	高	好	好	无	0-15000美元
12	中等	好	好	无	15000-35000美元
13	低	好	好	无	超过35000美元
14	高	坏	坏	无	15000-35000美元

Fig 10.13 A decision tree for credit risk assessment.



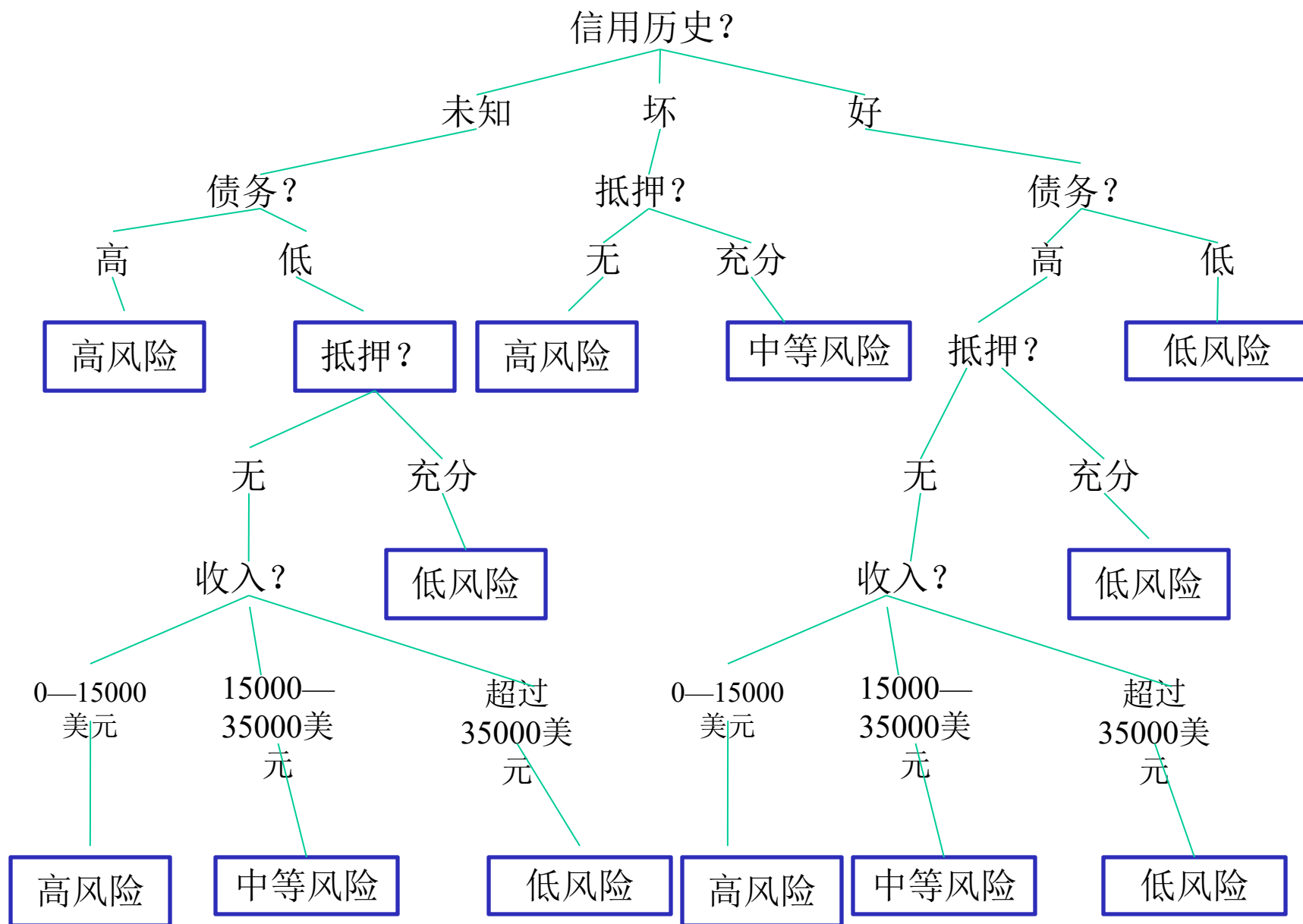
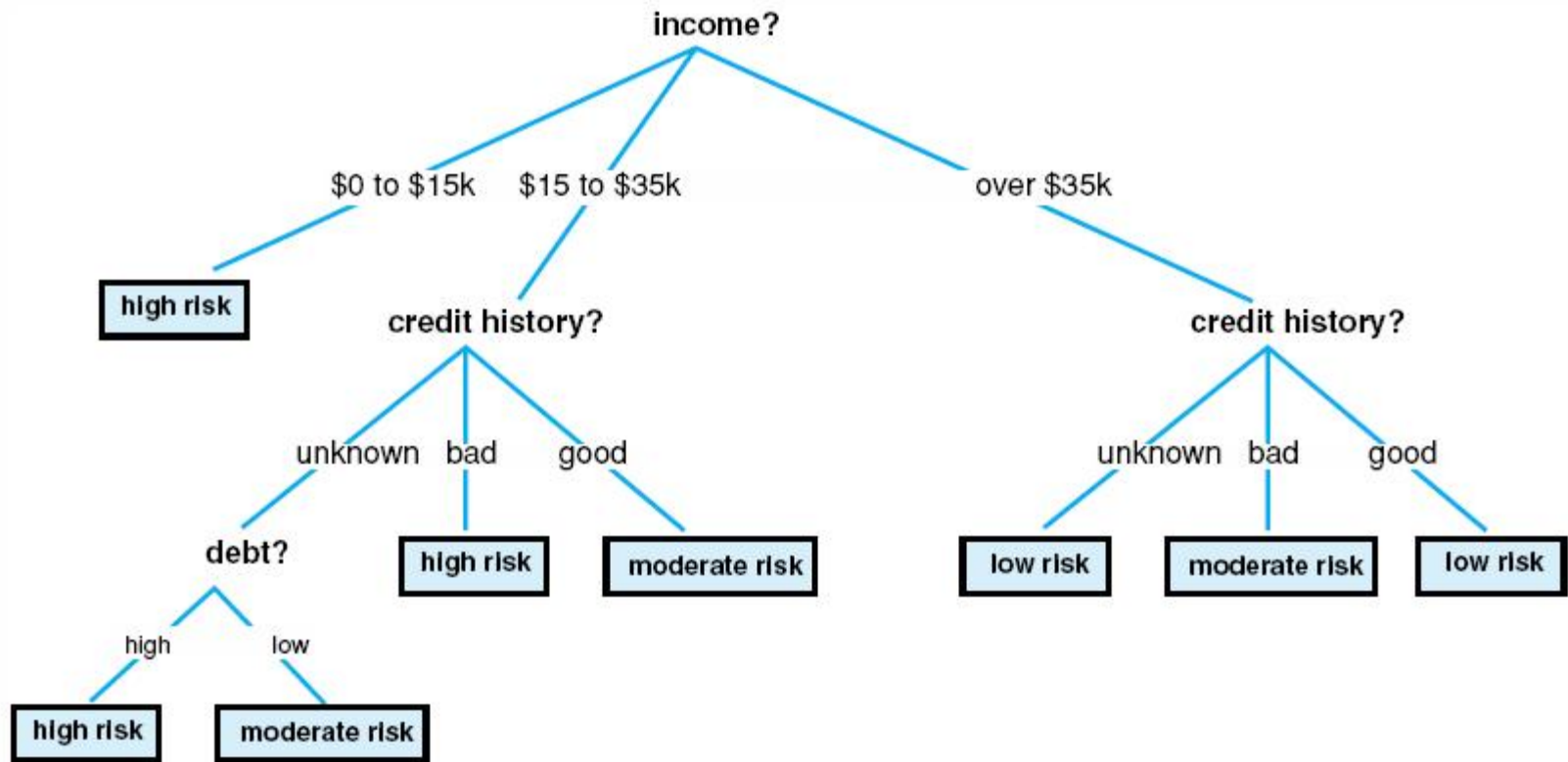


图10-13 信用风险评定的决策树

Fig 10.14 a simplified decision tree for credit risk assessment.



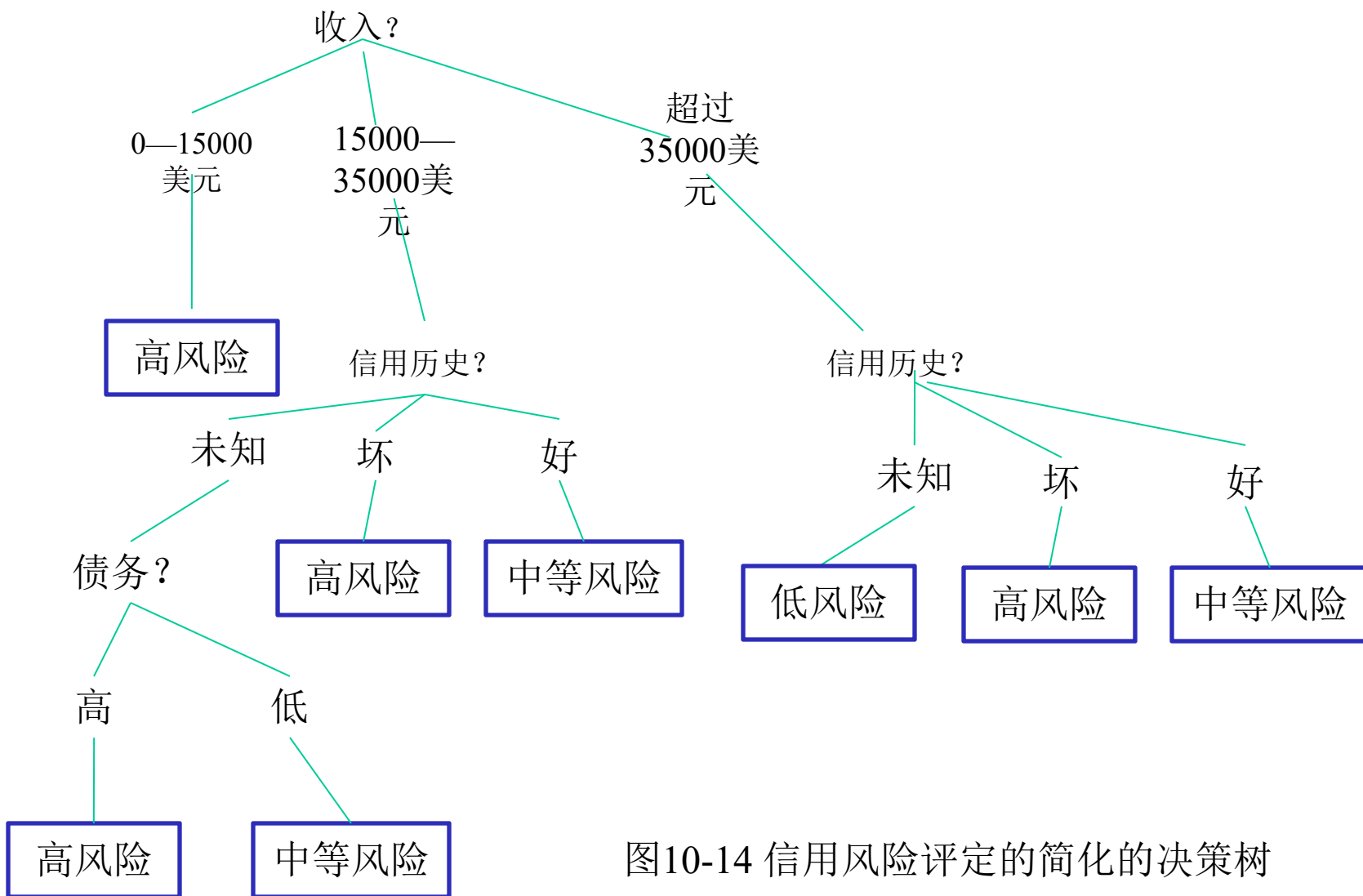
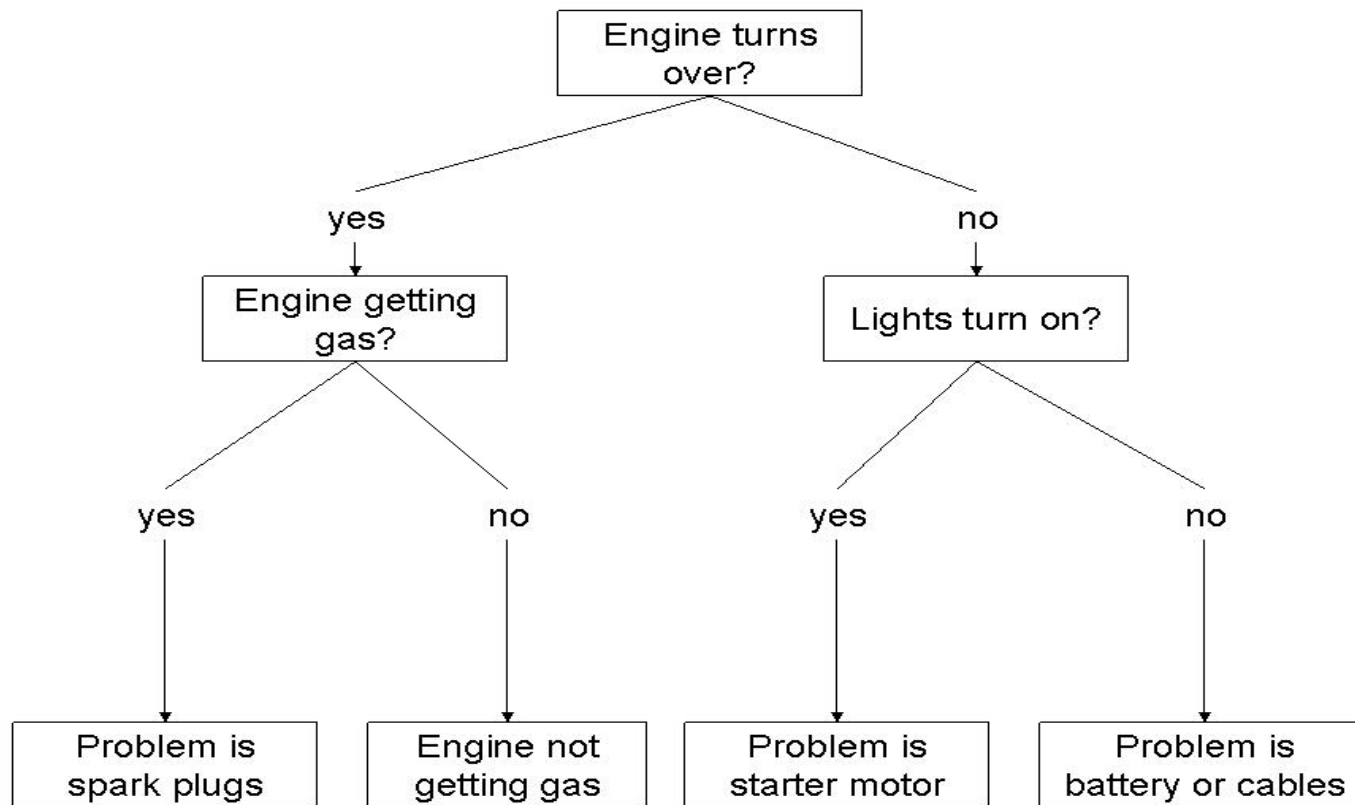


图10-14 信用风险评定的简化的决策树

Construct a decision tree such as is shown in Figures 10.13 and 10.14 for rules 1, 2, and 3 from the automobile diagnosis example from Chapter 8 Section 8.2.1.



The ID3 algorithm assumes that this is the simplest decision tree that covers all the training examples.

This principle, known as Occam's Razor:

It is vain to do with more what can be done with less..Entities should not be multiplied beyond necessity.

The induction algorithm begins with a sample of correctly classified members of the target categories. ID3 constructs a decision tree according to the algorithm:

```
function induce_tree (example_set, Properties)  
  
begin  
if all entries in example_set are in the same class  
then return a leaf node labeled with that class  
else if Properties is empty  
then return leaf node labeled with disjunction of all classes in example_set  
else begin  
    select a property, P, and make it the root of the current tree;  
    delete P from Properties;  
    for each value, V, of P,  
        begin  
            create a branch of the tree labeled with V;  
            let partitionv be elements of example_set with values V for property P;  
            call induce_tree(partitionv, Properties), attach result to branch V  
        end  
    end  
end  
end
```

基于决策树的归纳学习方法

- Hunt提出的概念学习系统CLS (Concept Learning System)是一种基于决策树的归纳学习系统。1979年, J.R.Quinlan在CLS基础上进行了发展, 提出了ID3算法, 该算法不仅能方便地表示概念的“属性-值”信息结构, 而且能从大量实例中有效地生成相应的决策树模型。

1.决策树及决策树构造算法CLS

在CLS的决策树中, 节点对应于待分类对象的属性, 由某一节点引出的弧对应于这个属性的可能值, 叶节点对应于分类的结果。

[例如]: 假设天气预测分类系统确定用于描述天气的分类有3种属性: outlook、humidity和windy, 这3种属性的值域分别为:

$\text{ValueType}(\text{outlook}) = \{\text{sunny}, \text{overcast}, \text{rain}\}$

$\text{ValueType}(\text{humidity}) = \{\text{high}, \text{normal}\}$

$\text{ValueType}(\text{windy}) = \{\text{true}, \text{false}\}$

若对天气实例只需给出2种分类结果N和P, 并认为P是正实例的分类结果, N是反实例的分类结果, 则天气预测的CLS决策树如图5。

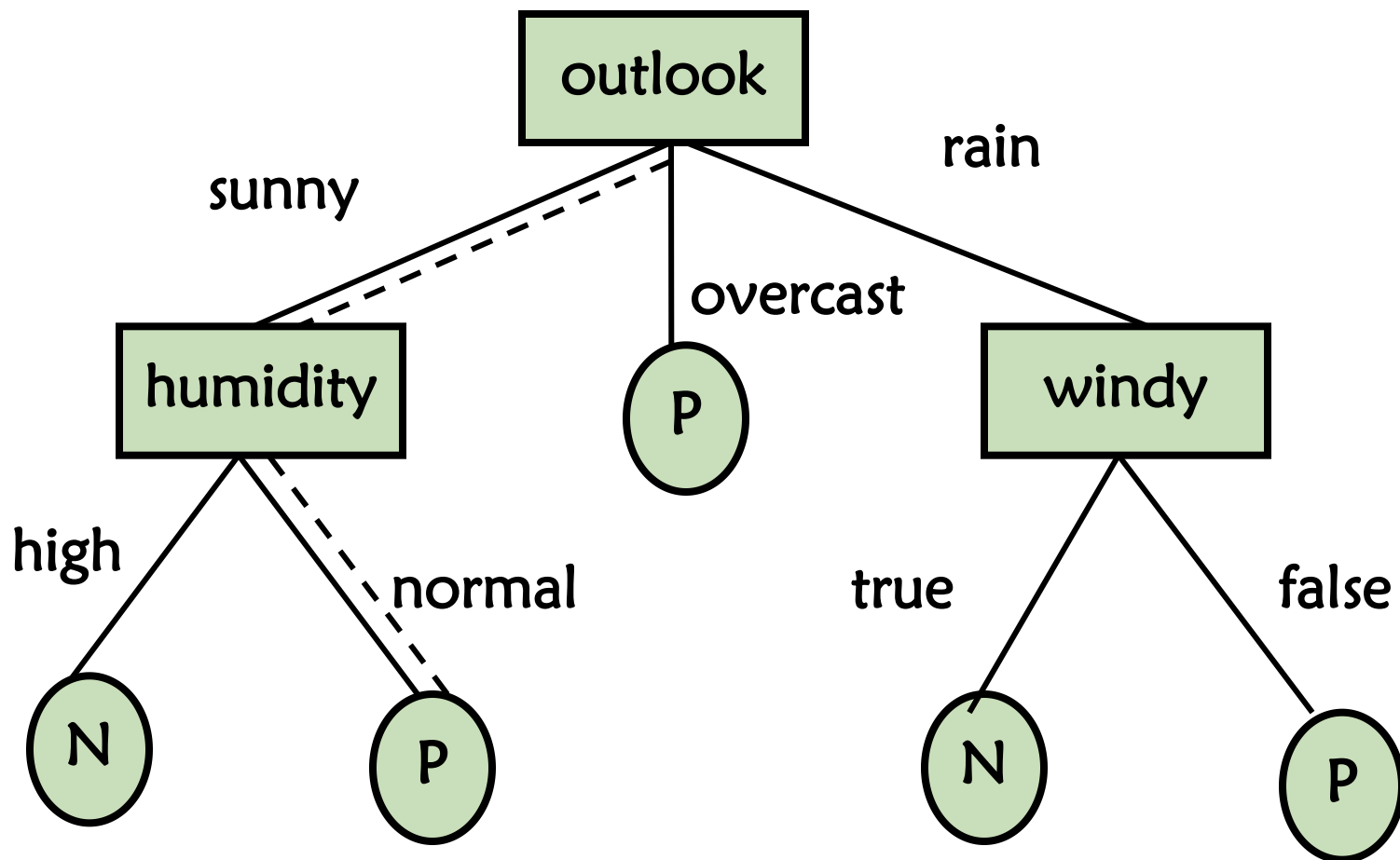


图5 天气预测的CLS决策树

术语：

- 属性：用节点表示。
- 属性值：用弧表示。
- 分类的结果：用叶子节点表示。
- 决策树：根据大量已知实例构造出的一般模型。

由此决策树可对任何的天气实例进行分类

- 设分类对象的属性表为 $\text{AttrList}=\{A_1, A_2, \dots, A_n\}$, 每个属性 A_i 的值域记为 $\text{Value Type}(A_i)$, 值域可是离散的, 也可能是连续的。
- 把分类结果组成分类结果表 $\text{Class}=\{C_1, C_2, \dots, C_m\}$, 一般应有 $n \geq 1, m \geq 2$ 。训练实例集 $T=\{\langle X, C \rangle\}$ 中的一个元素就是一个训练实例 $\langle X, C \rangle$, 实例 $\langle X, C \rangle$ 的特征向量 $X=(a_1, a_2, \dots, a_n)$, 其中 a_i 这个实例的第 i 个属性的取值, C 为这个实例的分类结果, $C \in \text{Class}$ 。可递归地描述下述决策树构造算法 CLS。

决策树构造算法CLS

- (1) 如果训练实例集 T 中所有实例的分类结果均为 C ，则返回 C 。
- (2) 从属性表 $AttList$ 中**任选**一个属性作为检测属性。
- (3) 若属性 A_i 的值域 $ValueType(A_i)$ 中有 s 个不同的取值，则将 T 分成 s 个子集 T_1, T_2, \dots, T_s ，每个子集的所有元素的属性 A_i 的取值相同。
- (4) 更新属性表，从属性表 $Attrlist$ 中删除检测属性 A_i 。
- (5) 对每个子集 T_k ，若子集 T_k 中的所有实例的分类结果均为 C_j ，则生成叶节点 C_j ；若子集中 T_k 所有实例的分类结果有两个或两个以上，则对子集 T_k 和更新后的属性表转到步骤2)，递归调用CLS构造算法，生成 T_k 的子树。

[例3] 假设需要根据人员的外貌特征对人员进行分类，用于人员分类的外貌特征有3个，它们组成人员分类属性表为：

$\text{AttrList} = \{\text{height}, \text{hair}, \text{eyes}\}$

个属性值域分别为：

$\text{ValueType}(\text{height}) = \{\text{short}, \text{tall}\}$

$\text{ValueType}(\text{hair}) = \{\text{blond}, \text{red}, \text{dark}\}$

$\text{ValueType}(\text{eyes}) = \{\text{blue}, \text{brown}\}$

(接下页)

若对人员进行分类的结果有两种，分别用“+”和“-”表示，组成分类结果表为：

Class={+,-}

若提供学习的训练实例集为：

$T = \{ \langle \text{short, blond, blue} \rangle, + \rangle, \langle \text{tall, blond, brown} \rangle, - \rangle, \langle \text{tall, red, blue} \rangle, + \rangle, \langle \text{short, dark, blue} \rangle, - \rangle, \langle \text{tall, dark, blue} \rangle, - \rangle, \langle \text{tall, blond, blue} \rangle, + \rangle, \langle \text{tall, dark, brown} \rangle, - \rangle, \langle \text{short, blond, brown} \rangle, - \rangle \}$

应用CLS构造算法构造决策树，说明构造过程，画出构造的决策树。

解: 执行CLS算法构造决策树的构造过程如下:

1) 若选取属性表AttrList={height,hair,eyes}中的属性height作为第一个检测属性, 则将T分为2个子集 T_1 和 T_2 , 其中, T_1 是属性height取值为short的训练实例子集, T_2 是属性height取值为tall的训练实例子集。故有:

$T_1 = \{ \langle \text{short, blond, blue} \rangle, + \rangle, \langle \text{short, dark, blue} \rangle, - \rangle, \langle \text{short, blond, brown} \rangle, - \rangle \}$

$T_2 = \{ \langle \text{tall, blond, brown} \rangle, - \rangle, \langle \text{tall, red, blue} \rangle, + \rangle, \langle \text{tall, dark, blue} \rangle, - \rangle, \langle \text{tall, blond, blue} \rangle, + \rangle, \langle \text{tall, dark, brown} \rangle, - \rangle \}$

从属性表中删除检测属性height, 故新的属性表为:

$\text{AttrList} = \{\text{hair, eyes}\}$

生产决策树如图6(a)所示

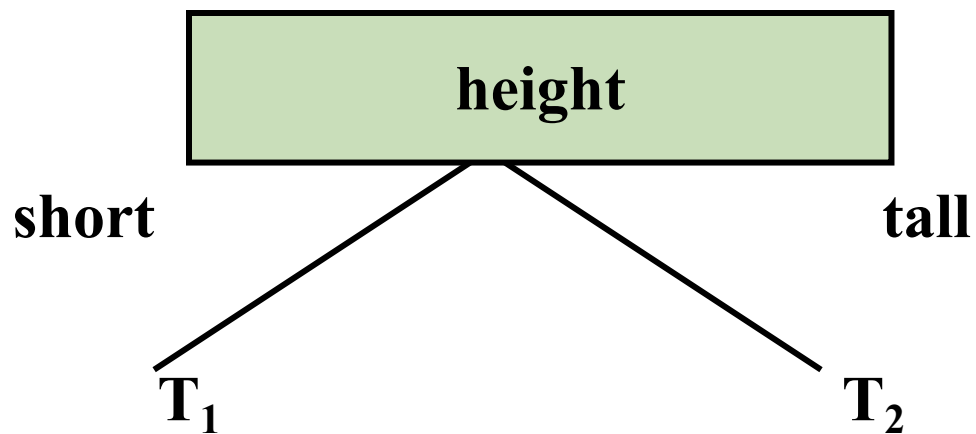


图6 生成的决策树(a)

2) 由于 T_1 和 T_2 中的实例的分类结果仍有两种，故需按新的属性表递归构造 T_1 和 T_2 的子树。若选取新属性hair作为第2个检测属性，且由于hair的值域 $\text{ValueType}(\text{hair}) = \{\text{blond}, \text{red}, \text{dark}\}$ 有3个取值，故将 T_1 分为3个子集 T_{11} 、 T_{12} 、 T_{13} ，将 T_2 也分为3个子集 T_{21} 、 T_{22} 和 T_{23} ，它们分别为：

$T_{11} = \{ \langle \text{short, blond, blue} \rangle, + \rangle, \langle \text{short, blond, brown} \rangle, - \rangle \}$

$T_{12} = \{ \}$

$T_{13} = \{ \langle \text{short, dark, blue} \rangle, - \rangle \}$

$T_{21} = \{ \langle \text{tall, blond, brown} \rangle, - \rangle, \langle \text{tall, blond, blue} \rangle, - \rangle \}$

$T_{22} = \{ \langle \text{tall, red, blue} \rangle, + \rangle \}$

$T_{23} = \{ \langle \text{tall, dark, blue} \rangle, - \rangle, \langle \text{tall, dark, brown} \rangle, - \rangle \}$

从属性表中删除检测属性hair，故新属性表为：

$\text{AttrList} = \{\text{eyes}\}$

至此，生产决策树如图6(b)所示。

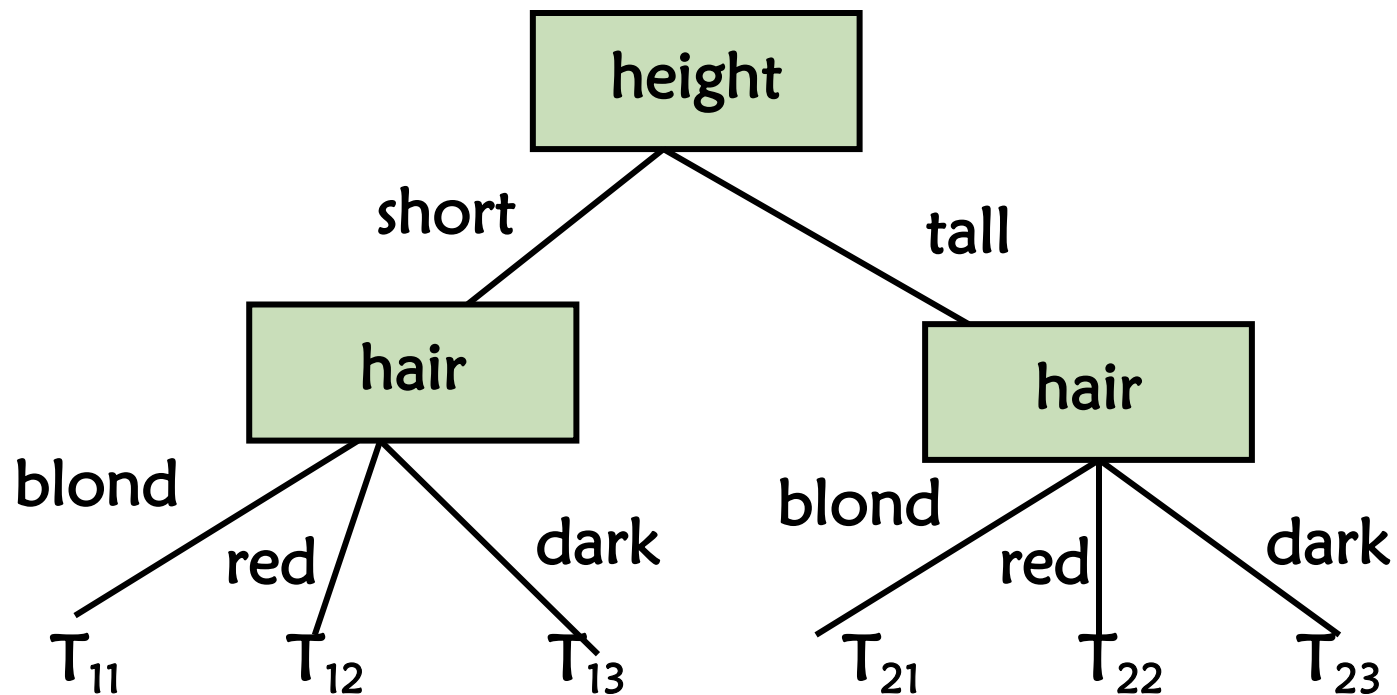


图6 生成的决策树(b)

3) 由于 T_{13} 、 T_{22} 和 T_{23} 中的实例的分类结果只有一种，故无需对 T_{13} 、 T_{22} 和 T_{23} 中的实例划分子集，并用 T_{13} 的实例分类结果“-”作为叶节点代替图6(b)所示决策树中节点的 T_{13} ，同样，用 T_{22} 的实例分类结果“+”作为叶节点代替 T_{22} ，用 T_{23} 的实例分类结果“-”为叶节点代替 T_{23} 。

T_{12} 是一个空集，故从决策树中删除节点 T_{12} 。

由于 T_{11} 和 T_{12} 中的实例的分类结果仍有两种，故还需按新的属性表递归构造 T_{11} 和 T_{21} 的子树。

以新属性表中的属性eyes作为第3个检测属性，且eyes的值域 $\text{ValueType}(\text{eyes}) = \{\text{blue}, \text{blown}\}$ 有2个取值，故将 T_{11} 分为2个子集 T_{111} 和 T_{112} ，将 T_{21} 也分为2个子集 T_{211} 和 T_{212} ，它们分别为：

$$T_{111} = \{ \langle (\text{short}, \text{blond}, \text{blue}), + \rangle \}$$

$$T_{112} = \{ \langle (\text{short}, \text{blond}, \text{brown}), - \rangle \}$$

$$T_{211} = \{ \langle (\text{tall}, \text{blond}, \text{blue}), + \rangle \}$$

$$T_{212} = \{ \langle (\text{tall}, \text{blond}, \text{brown}), - \rangle \}$$

从属性表中删除检测属性eyes，故新属性表已成为一个空表。

至此生产决策树如图6(c)所示。

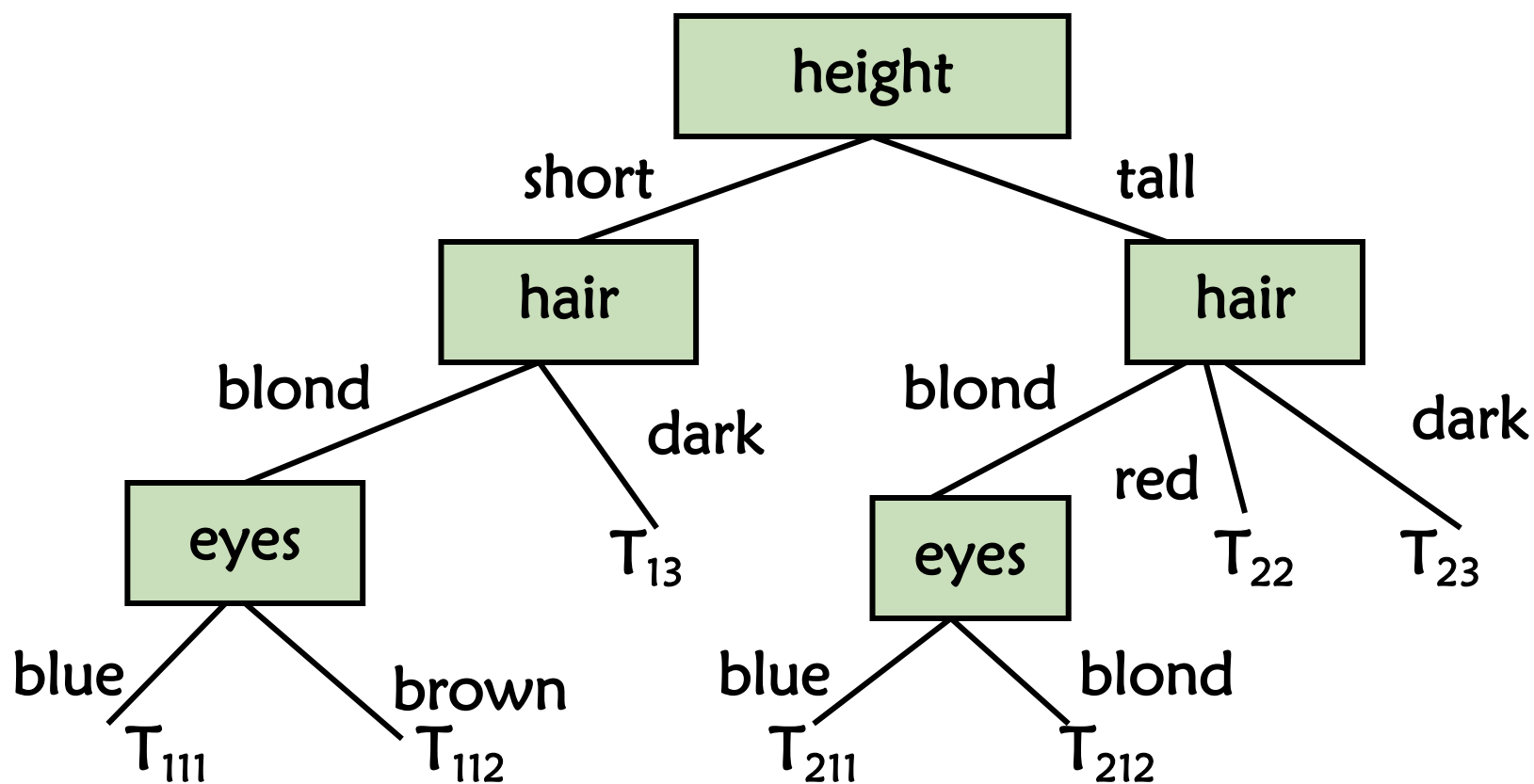


图6 生成的决策树(c)

4) 由于 T_{111} 、 T_{112} 、 T_{211} 和 T_{212} 中的实例的分类结果只有一种，则用它们的实例分类结果作为叶节点来分别代替它们，得到最后完整的决策树如图6(d)所示。

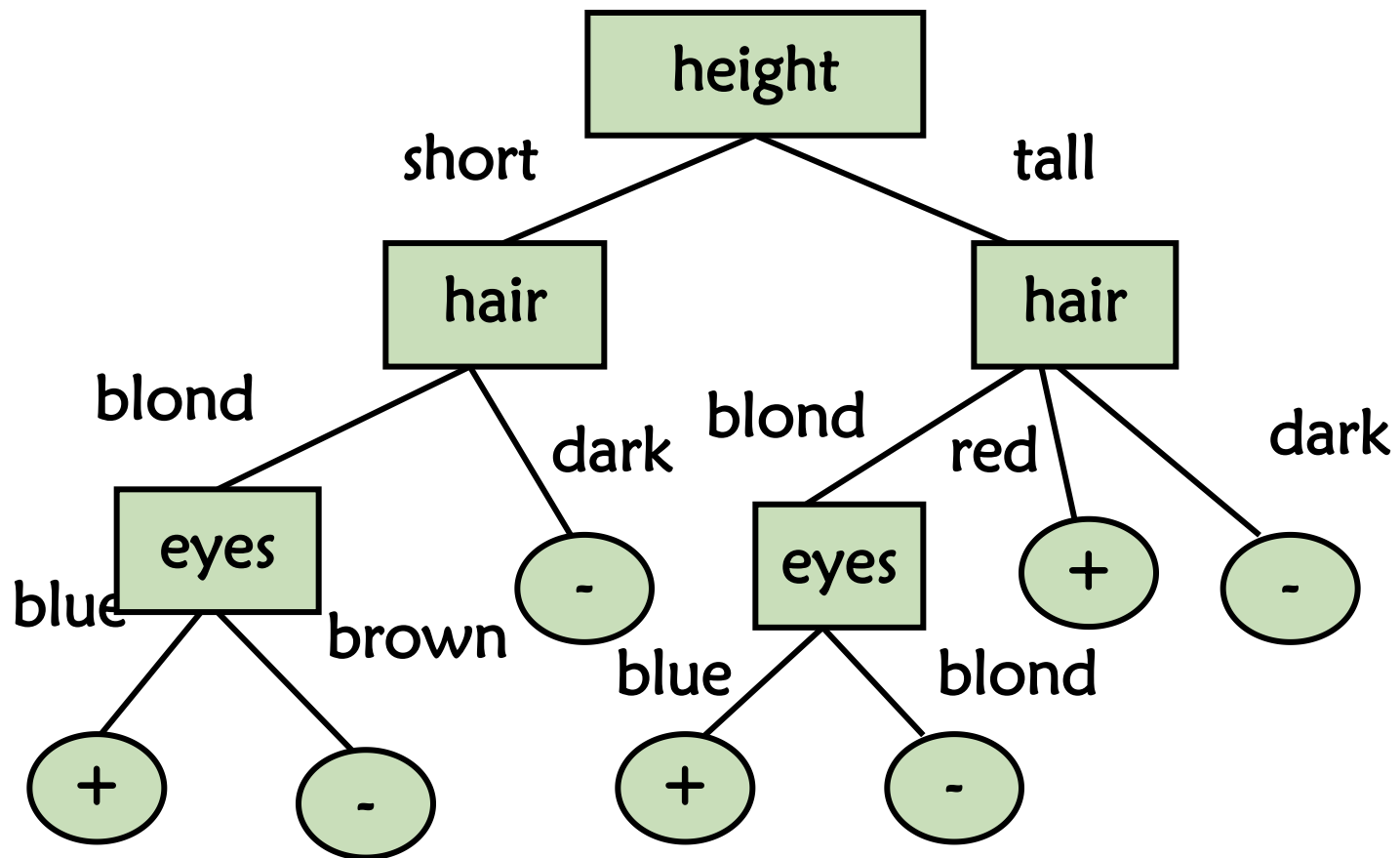


图6 生成的决策树(d)

基于决策树的归纳学习

假设需要根据人员的外貌特征对人员进行分类, 用于人员分类的外貌特征有3个, 它们组成人员分类属性表为: $\text{AttrList}=\{\text{height, hair, eyes}\}$

各属性的值域分别为:

$\text{ValueType}(\text{height})=\{\text{short, tall}\},$

$\text{ValueType}(\text{hair})=\{\text{blond, red, dark}\},$

$\text{ValueType}(\text{eyes})=\{\text{blue, brown}\}。$

对人员进行分类的结果有两种, 分别用“+”和“—”表示, 组成分类结果表为:

$\text{Class} = \{ +, - \},$ 提供学习的训练实例集为:

$T = \{ \langle \text{short, blond, blue, } + \rangle, \langle \text{(tall, red, brown), } - \rangle, \\ \langle \text{(tall, red, blue), } + \rangle, \langle \text{(short, dark, blue), } + \rangle, \\ \langle \text{(tall, dark, blue), } - \rangle, \langle \text{(tall, blond, blue), } - \rangle \}$

应用CLS算法构造决策树, 画出中间结果图和最终构造的决策树。

10.3.2 Information Theoretic Test Selection

We may think of each **property** of an instance as contributing a certain **amount of information** to its **classification**.

Information theory provides a mathematical basis for **measuring** the information content of a message(消息) .

- (1) The number of possible messages is important.
- (2) The influence of the probability of each message on the amount of information is important.

Shannon, definition: $-\log_2 P$

Given: $T = \{m_1, m_2, \dots, m_m\}$, the expected information content of a message is given by:

$$I(T) = -\sum_{j=1}^m P_j \log_2 P_j = -\sum_{j=1}^m \frac{e_j}{|T|} \log_2 \frac{e_j}{|T|} = -\left(\sum_{j=1}^m e_j \log_2 \frac{e_j}{|T|}\right) / |T|$$

The information in a message is measured in **bits**. For example, the information content of a message telling the outcome of the flip of an honest coin is:

$$\begin{aligned} I[\text{Coin toss}] &= -p(\text{heads})\log_2(p(\text{heads})) - p(\text{tails})\log_2(p(\text{tails})) \\ &= -1/2\log_2(1/2) - 1/2\log_2(1/2) \\ &= 1 \text{ bit} \end{aligned}$$

- If the coin has been rigged to come up heads 75 per cent of the time, then the information content of a message is:

$$\begin{aligned} I[\text{Coin toss}] &= -3/4\log_2(3/4) - 1/4\log_2(1/4) \\ &= -3/4 * (-0.145) - 1/4 * (-2) \\ &= 0.811 \text{ bits} \end{aligned}$$

Example label10-1:

$P(\text{risk is high})=6/14, p(\text{risk is moderate})=3/14, p(\text{risk is low})=5/14$

It follows that the distribution described in Table10.1,D9.1

$$I[D_{9.1}] = -6/14 \log_2(6/14) - 3/14 \log_2(3/14) - 5/14 \log_2(5/14)$$

$$= -6/14 * (-1.222) - 3/14 * (-2.222) - 5/14 * (-1.485)$$

$$= 1.531 \text{ bits}$$

$$E[p] = \sum_{i=1}^n \frac{|C_i|}{|C|} I[C_i]$$

$$\text{gain}(P) = I[C] - E[p]$$

Example: we make income the property tested at the root of the tree, this partitions the table of example into the partitions $C1=\{1,4,7,11\}$, $C2=\{2,3,12,14\}$, and $C3=\{5,6,8,9,10,13\}$ The expected information needed to complete the tree is:

$$\begin{aligned}
 E[\text{income}] &= 4/14 * I[C1] + 4/14 * I[C2] + 6/14 * I[C3] \\
 &= 4/14 * 0.0 + 4/14 * 1.0 + 6/14 * 0.650 \\
 &= 0.564 \text{ bits}
 \end{aligned}$$

The information gain for the distribution of Table 10.1 is :

$$\begin{aligned}
 \text{gain}(\text{income}) &= |[D9.1] - E[\text{income}] \\
 &= 1.531 - 0.564 \\
 &= 0.967 \text{ bits}
 \end{aligned}$$

$$\text{gain}(\text{credit history}) = 0.266$$

$$\text{gain}(\text{debt}) = 0.063$$

$$\text{gain}(\text{collateral}) = 0.206$$

Because income provides the greatest information gain, ID3 will select it as the root of the tree.

ID3 算法

- 由于CLS决策树构造算法每次从属性表AttrList中**任选**一个属性 A_i 作为检测属性来扩展生成决策树，按照训练集 T 中所有实例的属性 A_i 的取值来划分实例子集，因此，CLS算法的开销较大，效率较低。
- Quinlan提出了**ID3算法**，对检测属性的选择给出一种启发式规则，这个规则选择平均信息量（熵）最小的属性检测 A_i ，因此又称为**最小熵原理**。

若选择属性表AttrList={A₁,A₂,...,A_n}中的属性任选一个属性A_i为检测属性,则由A_i的值域
Value Type(A_i)={V₁,...,V_s}的s个取值把训练实例集T分为s个子集:

$$T = \bigcup_{K=1}^S T_K^{(i)}$$

子集T_K⁽ⁱ⁾中的所有实例的属性A_i的取值为V_k。

T中的实例的分类结果组成了分类结果表
Class={C₁,C₂,...,C_j,...,C_m} , 若分类结果为C_j的实例数为e_j,
1<=j<=m,且

$$\sum_{j=1}^m e_j = |T|$$

|T|表示训练实例集T中的实例总数,则实例分类结果为C_j的概率为:P_j=e_j/|T|

定义训练实例集T的**平均信息量**为：

$$I(T) = -\sum_{j=1}^m P_j \log_2 P_j = -\sum_{j=1}^m \frac{e_j}{|T|} \log_2 \frac{e_j}{|T|} = -\left(\sum_{j=1}^m e_j \log_2 \frac{e_j}{|T|}\right) / |T|$$

类似可定义 $T_K^{(i)}$ 的实例**平均信息量**为：

$$I\left(T_k^{(i)}\right) = -\left(\sum_{j=1}^m e_j^{(i)} \log_2 \frac{e_j^{(i)}}{|T_K^{(i)}|}\right) / |T_K^{(i)}|$$

其中, $e_j^{(i)}$ 为子集 $T_K^{(i)}$ 中分类结果 C_j 的实例数, $T_K^{(i)}$ 为子集中实例总数,

$$\sum_{j=1}^m e_j^{(i)} = |T_k^{(i)}|$$

• 若选择属性 A_i 作为检测属性将训练集 T 分为 s 个子集后，可由各实例子集的 **实例总信息量** $|T_{K(i)}| \cdot I(T_{K(i)})$ 之和对实例集 $|T|$ 的实例总数 T 的平均值来表示为实例集 T 的 **平均信息量**。

即：

$$I(T, A_i) = \left(\sum_{K=1}^s |T_{K(i)}| \cdot I(T_{K(i)}) \right) / |T|$$

选择属性 A_i 作为检测属性来将实例集 T 划分子集的启发式规则

选择属性作为检测属性，应由属性 A_i 的不同值把实例集划分为若干子集之前与之后的实例平均信息量的**差最大**，即使得 $GI(T, A_i) = I(T) - I(T, A_i)$ 的值最大。 $GI(T, A_i)$ 可认为是属性 A_i 对训练实例集 T 的信息量，因此，启发式规则实际是选择**信息量最小**的属性作为检测属性 A_i 来划分实例集。

例4：

**应用ID3算法构造例3 要求构造的决策树，
说明构造过程决策树。**

解: 应用ID3算法构造决策树的构造过程如下:

- 1) 实例集T分类结果为“+”的实例数为3, 分类结果为“-”的实例数为5, 故T的实例平均信息量为:

$$I(T) = -(5\log_2 5/8 + 3\log_2 3/8)/8 = 0.954 \text{ (bit)}$$

若选属性表AttrList={height,hair,eyes}中的第1个属性height为检测属性, 其值域ValueType(height)={short,tall}有2个值, 则将T分为2个子集 $T_1^{(1)}$ 和 $T_2^{(1)}$:

$$T_1^{(1)} = \{ \langle \text{short, blond, blue} \rangle, + \rangle, \langle \text{short, dark, blue} \rangle, - \rangle, \langle \text{short, blond, brown} \rangle, - \rangle \}$$

$$T_2^{(1)} = \{ \langle \text{tall, blond, brown} \rangle, - \rangle, \langle \text{tall, red, blue} \rangle, + \rangle, \langle \text{tall, dark, blue} \rangle, - \rangle, \langle \text{tall, blond, blue} \rangle, + \rangle, \langle \text{tall, dark, brown} \rangle, - \rangle \}$$

$T_1^{(1)}$ 中的分类结果为“+”的实例数为1，分类结果为“-”的实例数为2，可得到 $T_1^{(1)}$ 的实例平均信息量为：

$$I(T_1^{(1)}) = -(\log_2 1/3 + 2\log_2 2/3)/3 = 0.918 \text{ (bit)}$$

$T_2^{(1)}$ 中的分类结果为“+”的实例数为2，分类结果为“-”的实例数为3，可得到 $T_2^{(1)}$ 的实例平均信息量为

$$I(T_2^{(1)}) = -(2\log_2 2/5 + 3\log_2 3/5)/5 = 0.971 \text{ (bit)}$$

故而，若以属性height为检测属性扩展后的T的实例平均信息量为：

$$I(T, \text{height}) = (3 * I(T_1^{(1)}) + 5 * I(T_2^{(1)})) / 8 = 0.951 \text{ (bit)}$$

由此可得属性height对实例T信息量为：

$$GI(T, height) = I(T) - I(T, height) = 0.003 \text{ (bit)}$$

$$GI(T, hair) = 0.5 \text{ (bit)}$$

$$GI(T, eyes) = 0.347 \text{ (bit)}$$

由于 $GI(T, hair) > GI(T, eyes) > GI(T, height)$ ，故选取第2个属性 hair 作为检测属性，并从属性表删除属性 hair，更新属性表为：

$$AttrList = \{height, eyes\}$$

至此生成的决策树如图7(a)所示

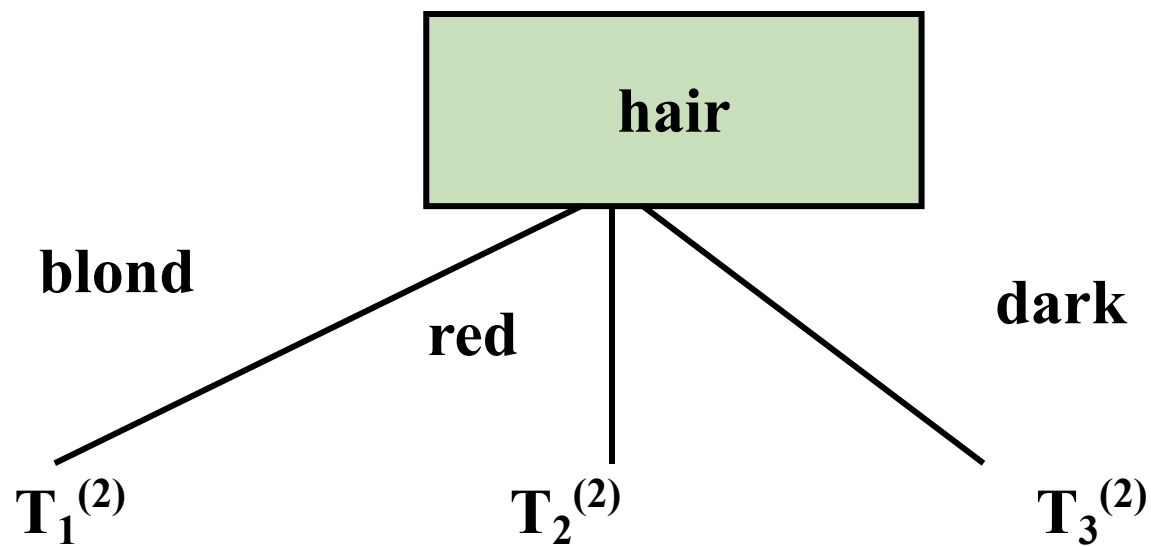


图7 生成的决策树(a)

图中的3个子节点 $T_1^{(2)}$ 、 $T_2^{(2)}$ 、 $T_3^{(2)}$ 是由检测属性hair的值域 $\text{ValueType}(\text{hair})=\{\text{blond},\text{red},\text{dark}\}$ 的3个取值将T划分成3个子集，分别是：

$$T_1^{(2)} = \{ \langle (\text{short}, \text{blond}, \text{blue}), + \rangle, \langle (\text{tall}, \text{blond}, \text{brown}), - \rangle, \\ \langle (\text{tall}, \text{blond}, \text{blue}), - \rangle, \langle (\text{short}, \text{blond}, \text{brown}), - \rangle \}$$

$$T_2^{(2)} = \{ \langle (\text{tall}, \text{red}, \text{blue}), + \rangle \}$$

$$T_3^{(2)} = \{ \langle (\text{short}, \text{dark}, \text{blue}), - \rangle, \langle (\text{tall}, \text{dark}, \text{blue}), - \rangle, \\ \langle (\text{tall}, \text{dark}, \text{brown}), - \rangle \}$$

2) 由于 $T_2^{(2)}$ 中的决策树分类结果只有“+”，故用叶节点“+”代替节点 $T_2^{(2)}$ ； $T_3^{(2)}$ 中的实例的分类结果都是“-”，故用叶节点“-”代替节点 $T_3^{(2)}$ 。由于 $T_1^{(2)}$ 中的分类结果仍有两种，故需按新属性表启发式选择检测属性来构造子树。

在实例集 $T_1^{(2)}$ 中，分类结果为“+”的实例数为2，分类结果为“-”的实例数也为2，可得 $T_1^{(2)}$ 的实例平均信息量为：

$$I(T_1^{(2)}) = -(2\log_2 2/4 + 2\log_2 2/4)/4 = 1 \text{ (bit)}$$

若选新属性表AttrList={height,eyes}中的第1个属性height为检测属性，则由值域ValueType(height)={short,tall}的2个取值将 $T_1^{(2)}$ 分为2个子集 $T_{11}^{(21)}$ 和 $T_{12}^{(21)}$ ：

$$T_{11}^{(21)} = \{ \langle (\text{short}, \text{blond}, \text{blue}), + \rangle, \langle (\text{short}, \text{blond}, \text{brown}), - \rangle \}$$

$$T_{12}^{(21)} = \{ \langle (\text{tall}, \text{blond}, \text{brown}), + \rangle, \langle (\text{tall}, \text{blond}, \text{brown}), - \rangle \}$$

$T_{11}^{(21)}$ 和 $T_{12}^{(21)}$ 中的分类结果为“+”的实例数都为1，分类结果为“-”的实例数也都为1，可得 $T_{11}^{(21)}$ 和 $T_{12}^{(21)}$ 的实例平均信息量分别为：

$$I(T_{11}^{(21)}) = -(\log_2 1/2 + \log_2 1/2)/2 = 1 \text{ (bit)}$$

$$I(T_{12}^{(21)}) = -(\log_2 1/2 + \log_2 1/2)/2 = 1 \text{ (bit)}$$

故而，若以属性height为检测属性扩展后的实例平均信息量为：

$$I(T_1^{(2)}, \text{height}) = (2 * I(T_{11}^{(21)}) + 2 * I(T_{12}^{(21)})) / 4 = 1 \text{ (bit)}$$

由此可得属性height对 $T_1^{(2)}$ 实例集的信息量为：

$$GI(T_1^{(2)}, \text{height}) = I(T_1^{(2)}) - I(T_1^{(2)}, \text{height}) = 0$$

若选新属性表AttrList={height,eyes} 中的第2个属性eyes为检测属性, 由其值域
 $\text{ValueType}(\text{eyes})=\{\text{blue},\text{brown}\}$ 的2个取值将 $T_1^{(2)}$
分为2个子集 $T_{11}^{(22)}$ 和 $T_{12}^{(22)}$:

$$T_{11}^{(22)} = \{ \langle (\text{short},\text{blond},\text{blue}), + \rangle, \langle (\text{tall},\text{blond},\text{blue}), - \rangle \}$$

$$T_{12}^{(22)} = \{ \langle (\text{tall},\text{blond},\text{brown}), + \rangle, \langle (\text{short},\text{blond},\text{brown}), + \rangle \}$$

由于height对实例集 $T_1^{(2)}$ 的信息量 $\text{GI}(T_1^{(2)}, \text{height})=0$,因此选择属性eyes作为检测属性, 由eyes的2个取值扩展生成2个子集节点 $T_{11}^{(22)}$ 和 $T_{12}^{(22)}$ 。

至此, 生成决策树如图7(b)。

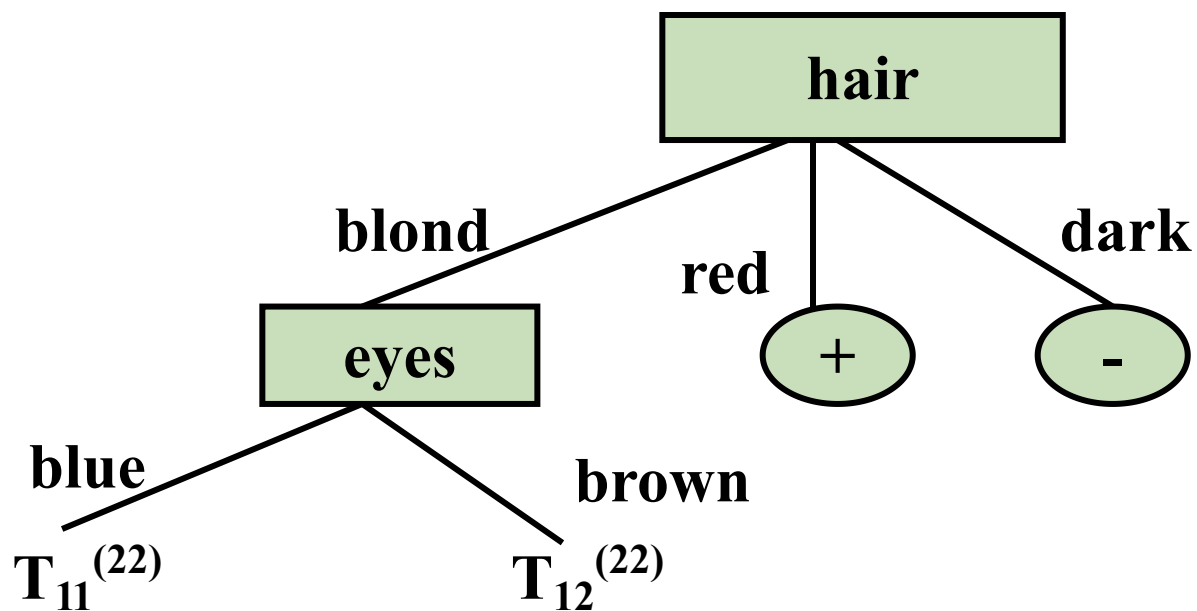


图7生成的决策树(b)

3) 由于 $T_{11}^{(22)}$ 中的决策树分类结果只有“+”，故用叶节点“+”代替节点 $T_{11}^{(22)}$ ； $T_{12}^{(22)}$ 中的实例的分类结果都是“-”，故用叶节点“-”代替节点 $T_{12}^{(22)}$ 。

至此，得到最后的完整决策树如图7(c)。

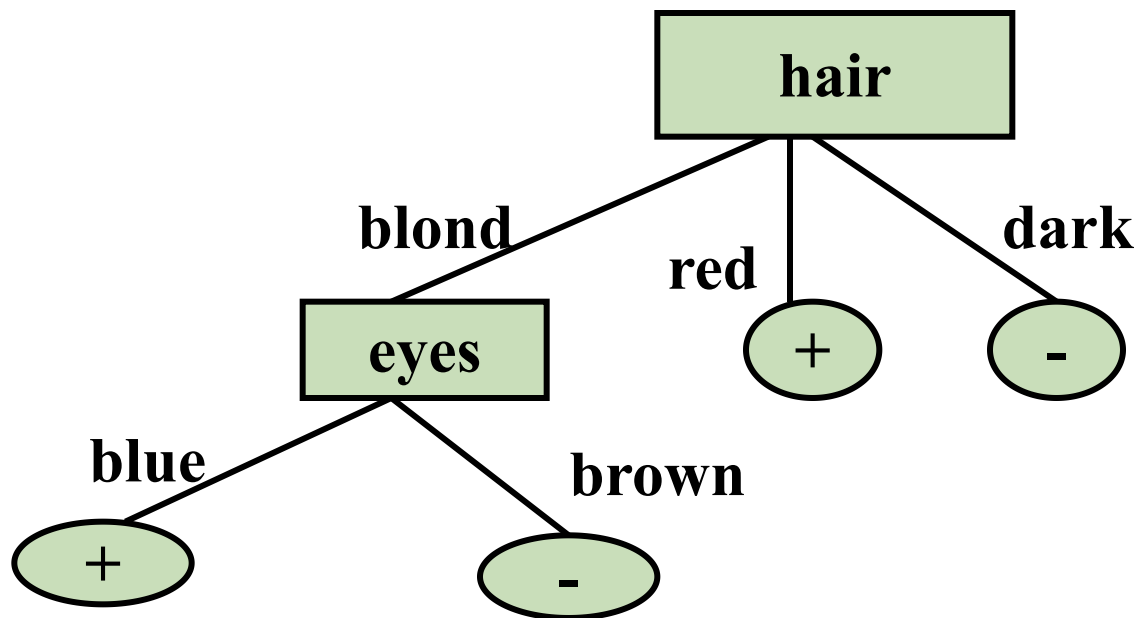


图7 生成的决策树(c)

练习

假设需要根据人员的外貌特征对人员进行分类, 用于人员分类的外貌特征有3个, 它们组成人员分类属性表为: AttrList={height, hair, eyes}

各属性的值域分别为:

ValueType(height)={short, tall},

ValueType(hair)={blond, red, dark},

ValueType(eyes)={blue, brown}。

对人员进行分类的结果有两种, 分别用“+”和“—”表示, 组成分类结果表为:

Class = { +, — }, 提供学习的训练实例集为:

$T = \{ \langle \text{short, blond, blue, +} \rangle, \langle \text{(tall, red, brown), —} \rangle, \langle \text{(tall, red, blue), +} \rangle, \langle \text{(short, dark, blue), +} \rangle, \langle \text{(tall, dark, blue), —} \rangle, \langle \text{(tall, blond, blue), —} \rangle \}$

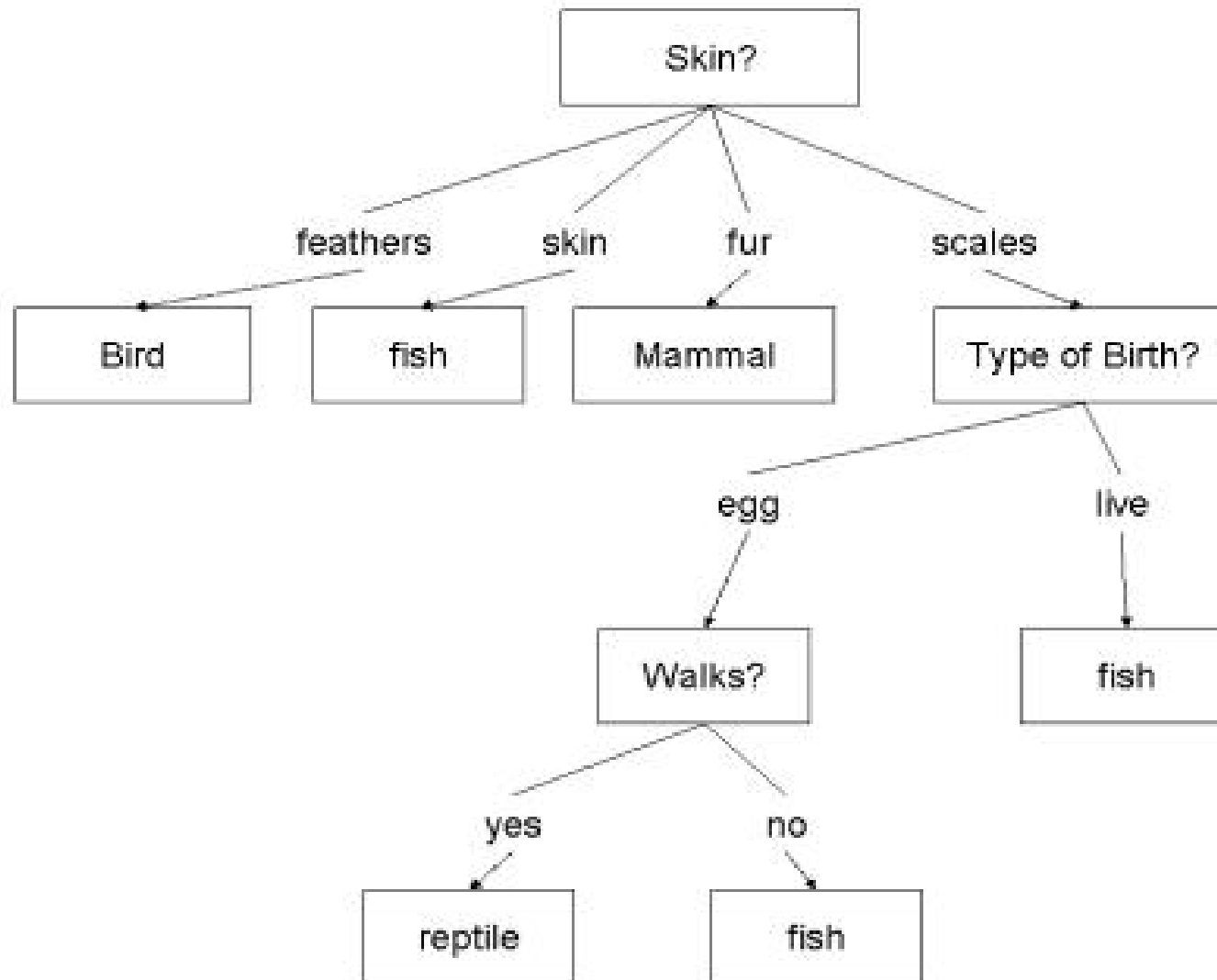
• 应用ID3算法构造决策树, 请回答下面有关问题:

1. 计算划分前 T 的实例平均信息量 $I(T)$
2. 若选第3个属性eyes为检测属性, 则将 T 划分为哪几个子集?
3. 每个子集的的实例平均信息量是多少?
4. 每个子集的的实例总信息量是多少?
5. 划分后 T 的实例平均信息量 $I(T, \text{eyes})$ 是多少?

Develop a simple table of examples to classify an animal as a fish, a bird, a reptile, or a mammal. Trace the construction of a decision tree by the ID3 algorithm.

Class	No.	Skin	Birth	Cares for young	Walks	Swims	Flies
Bird	1	Feathers	Egg	Yes	Yes	No	Yes
Bird	2	Feathers	Egg	Yes	Yes	No	No
Mammal	3	Fur	Live	Yes	Yes	Yes	No
Mammal	4	Fur	Live	Yes	Yes	Yes	No
Reptile	5	Scales	Egg	No	Yes	Yes	No
Fish	6	Skin	Live	No	No	Yes	No
Bird	7	Feathers	Egg	Yes	Yes	Yes	Yes
Mammal	8	Fur	Live	Yes	Yes	Yes	No
Fish	9	Skin	Egg	No	No	Yes	No
Mammal	10	Fur	Live	Yes	No	No	Yes
Fish	11	Scales	Egg	No	No	Yes	No
Bird	12	Feathers	Egg	Yes	Yes	No	No
Fish	13	Scales	Live	Yes	No	Yes	No
Mammal	14	Fur	Egg	Yes	Yes	Yes	No

the resulting decision tree:



5. Using Shannon's formula, show whether or not a message about the outcome of a fair six-sided die has more information than one about the outcome of a coin toss. What if the die message is "not three"?

Solution:

In Section 10.3.2 the information gain from a toss of a coin is shown to be one bit.

$$I[M] = \left(\sum_{i=1}^n -p(m_i) \log_2 p(m_i) \right)$$

$$I[M] = \left(\sum_{i=1}^6 -\frac{1}{6} \log_2 \frac{1}{6} \right)$$

$$= 6 * \left(-\frac{1}{6} \log_2 \frac{1}{6} \right)$$

$$= 6 * \left(-\frac{1}{6} \log_2 \frac{1}{6} \right)$$

$$= \left(-\log_2 \frac{1}{6} \right)$$

$$= 2.59$$

Thus, information about the roll of a six-sided fair die contains more information than the toss of a coin.

If we know the die cannot roll a three, then Shannon's formula gives us:

$$I[M] = \left(\sum_{i=1}^5 -\frac{1}{5} \log_2 \frac{1}{5} \right)$$

$$= \left(-\log_2 \frac{1}{5} \right)$$

$$= 2.32$$

which is consistent with the intuition that if something is less likely, it contains more information.

作业

- 2
- 4
- 6