

# 数据挖掘与机器学习

潘斌

panbin@nankai.edu.cn

范孙楼227

1

# 上节回顾

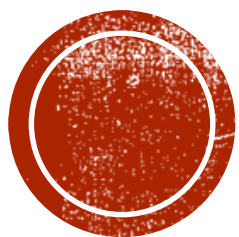
- 数据预处理的方法
  - 汇总统计
  - 可视化
  - 降维（特征提取、特征选择）



# 本节提要

- 概念学习
- 模型的评估与选择
- 实验课推迟至第8周开始，8-13周共5次





# 概念学习

Concept learning

- 概念（concept）：一个对象或事件集合，它是从更大的集合中选取的子集，或者是在这个较大集合中定义的布尔函数。
  - 如，从动物的集合中选取鸟类
    - 在动物集合中定义的函数，它对鸟类产生true并对其他动物产生false
- 概念学习：从样例中逼近布尔值函数。是指从有关某个布尔函数的输入输出训练样例中，推断出该布尔函数。



## ■ 一个概念学习的例子：Aldo Enjoy Sport

---

Example	<i>Sky</i>	<u><i>AirTemp</i></u>	<i>Humidity</i>	<i>Wind</i>	<i>Water</i>	<i>Forecast</i>	<u><i>EnjoySport</i></u>
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

---

- 输入某天的各种属性Sky、AirTemp、Humidity...等条件，Aldo是否进行水上运动？



## ■ 术语定义

- 概念定义在一个实例（instance）集合之上，这个集合表示为 $X$ 。
  - 例中， $X$ 是所有的日子，每个日子由Sky、AirTemp、Humidity、Wind、Water和Forecast六个属性表示。
- 待学习的概念或函数称为目标概念（target concept），记作 $c$ 。
  - $c$ 可以是定义在实例 $X$ 上的任意布尔函数，即 $c:X \rightarrow \{0, 1\}$ 。
  - 例中，目标概念对应于属性EnjoySport的值，当EnjoySport=Yes时 $c(x)=1$ ，当EnjoySport=No时 $c(x)=0$ 。



- **训练样例**（ training examples ）：  $X$  中的一个实例  $x$  以及它的目标概念值  $c(x)$  。
  - 用序偶  $\langle x, c(x) \rangle$  来描述训练样例
  - 符号  $D$  用来表示训练样例的集合。
- **正例**（ positive example ）： 对于  $c(x)=1$  的实例， 也称为目标概念的成员。
- **反例**（ negative example ）： 对于  $c(x)=0$  的实例， 也称为非目标概念成员。





- 给定目标概念  $c$  的训练样例集，学习器面临的问题是假设或估计  $c$ 。
- 使用符号  $H$  来表示所有可能假设的集合，称之为假设空间。
  - 通常  $H$  依设计者所选择的假设表示而定。
  - $H$  中每个假设  $h$  表示  $X$  上定义的布尔函数，即  $h:X \rightarrow \{0,1\}$ 。
  - 目标：寻找  $H$  中的假设  $h$ ，使对于  $X$  中的所有  $x$ ， $h(x)=c(x)$ 。



- 获取 $h$ 的方法：归纳学习
- 归纳学习：从特殊的样例得到普遍的规律
  - 归纳学习算法最多只能保证输出的假设与训练样例（特殊）相拟合。
- 归纳学习假设（ The Inductive Learning Hypothesis ）：任一假设如果在足够大的训练样例集中很好地逼近目标函数，它也能在未见实例中很好地逼近目标函数。（数据同分布）



## ■ 独立同分布

- 独立：每一个样本的出现或者生成，都是独立事件，即任意两个样本之间**不相关**
- 同分布：抽样内样本**服从总体**的分布（文字识别不平衡）



## ■ 表示假设 ( Representing Hypotheses )

- 很多种可能的表示方法
- 一个简单的形式：实例的各属性约束的合取式  $\wedge$  ( Conjunction )
- EnjoySport例中，令每个假设为6个约束（或变量）的向量，每个约束对应一个属性可取值范围，为
  - 由“?”表示任意值（如，AirTemp=?）
  - 明确指定的属性值（如，AirTemp=Warm）
  - 由“ $\emptyset$ ”表示不接受任何值（如，AirTemp= $\emptyset$ ）
  - 如，
    - $\langle \text{Sunny}, ?, ?, \text{Strong}, ?, \text{Same} \rangle$
    - $\langle ?, ?, ?, ?, ?, ? \rangle$  所有的样例都是正例（最一般，任取皆正）
    - $\langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$  所有的样例都是反例（最特殊，任取皆零）



- 任何概念学习任务能被描述为
  - 实例的集合（总体）
  - 实例集合上的目标函数（目标）
  - 训练样例的集合（样本）
  - 候选假设的集合（假设）
- EnjoySport 概念学习任务



# ■ EnjoySport 概念学习任务

---

- 已知:
    - 实例集  $X$ : 可能的日子, 每个日子由下面的属性描述:
      - $Sky$  (可取值为 *Sunny*, *Cloudy* 和 *Rainy*)
      - $AirTemp$  (可取值为 *Warm* 和 *Cold*)
      - $Humidity$  (可取值为 *Normal* 和 *High*)
      - $Wind$  (可取值为 *Strong* 和 *Weak*)
      - $Water$  (可取值为 *Warm* 和 *Cool*)
      - $Forecast$  (可取值为 *Same* 和 *Change*)
    - 假设集  $H$ : 每个假设描述为 6 个属性  $Sky$ ,  $AirTemp$ ,  $Humidity$ ,  $Wind$ ,  $Water$  和  $Forecast$  的值约束的合取。约束可以为 “?” (表示接受任意值), “ $\emptyset$ ” (表示拒绝所有值), 或一特定值。
    - 目标概念  $c: EnjoySport: X \rightarrow \{0, 1\}$
    - 训练样例集  $D$ : 目标函数的正例和反例
  - 求解:
    - $H$  中的一假设  $h$ , 使对于  $X$  中任意  $x$ ,  $h(x)=c(x)$ 。
- 



- 当假设的表示形式选定后，那么就隐含地为学习算法确定了所有假设的空间
  - 这些假设是学习程序所能表示的
  - 也是它能够学习的
  - 例，Enjoy-Sport的假设空间
- 概念学习可以看作一个搜索的过程
  - 搜索范围：假设的表示所隐含定义是整个空间
  - 搜索目标：能够最好地拟合训练样例的假设



- 假设的一般到特殊序关系
  - 考虑下面两个假设
    - $h1 = \langle \text{sunny}, ?, ?, \text{Strong}, ?, ? \rangle$
    - $h2 = \langle \text{Sunny}, ?, ?, ?, ?, ? \rangle$
  - 任何被 $h1$ 划分为正例的实例都会被 $h2$ 划分为正例，因此 $h2$ 比 $h1$ 更一般。
- 利用这个关系，无需列举所有假设，就能在无限的假设空间中进行较彻底的搜索





- **more-general-than-or-equal-to** ( 更一般或相等 )
- 定义： 令  $h_j$  和  $h_k$  为在  $X$  上定义的布尔函数。定义一个 **more-general-than-or-equal-to** 关系，记做  $\geq_g$ 。称  $h_j \geq_g h_k$  当且仅当

$$(\forall \mathbf{x} \in X)[(h_k(\mathbf{x})=1) \rightarrow (h_j(\mathbf{x})=1)]$$



- 对  $\mathbf{X}$  中任意实例  $\mathbf{x}$  和  $\mathbf{H}$  中任意假设  $h$ ，我们说  $\mathbf{x}$  满足  $h$  当且仅当  $h(\mathbf{x})=1$ 。
- 给定假设  $h_j$  和  $h_k$ ， $h_j$  more-general-than-or-equal-to  $h_k$ ，当且仅当任意一个满足  $h_k$  的实例同时也满足  $h_j$ 。
- $h_j$  严格的 more-general-than  $h_k$ （写作  $h_j >_g h_k$ ），当且仅当  $(h_j \geq_g h_k) \wedge \neg(h_k \geq_g h_j)$ 。
- 逆向的关系“比……更特殊”： $h_j$  more-specific-than  $h_k$ ，当  $h_k$  more-general-than  $h_j$ 。



## ■ Instance, Hypotheses, and More-General-Than

$h_1 = \langle \text{Sunny}, ?, ?, \text{Strong}, ?, ? \rangle$

$h_2 = \langle \text{Sunny}, ?, ?, ?, ?, ? \rangle$

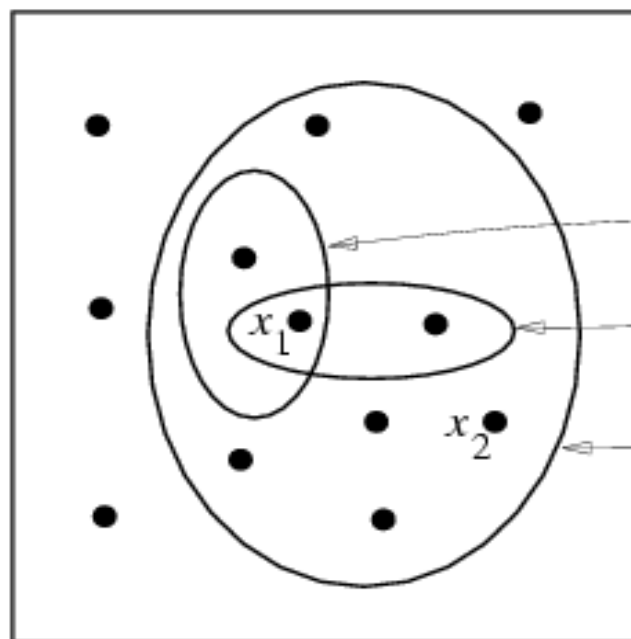
$h_3 = \langle \text{Sunny}, ?, ?, ?, \text{Cool}, ? \rangle$

$x_1 = \langle \text{Sunny}, \text{Warm}, \text{High}, \text{Strong}, \text{Cool}, \text{Same} \rangle$

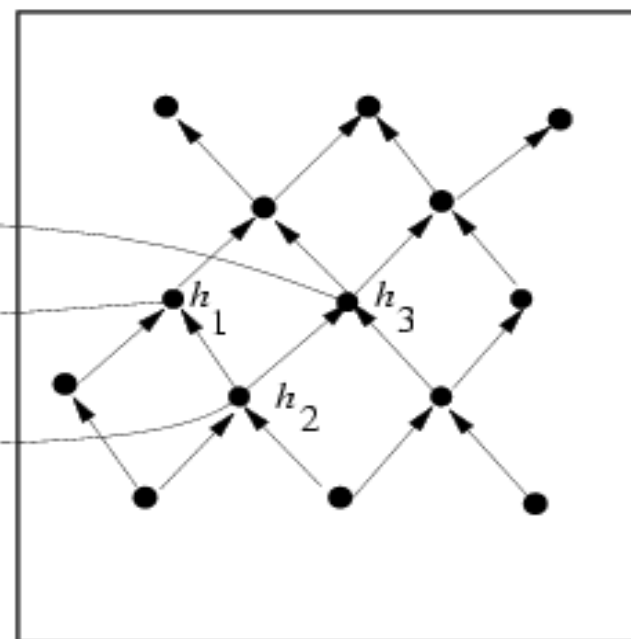
$x_2 = \langle \text{Sunny}, \text{Warm}, \text{High}, \text{Light}, \text{Warm}, \text{Same} \rangle$



*Instances X*



*Hypotheses H*



Specific

General

$x_1 = \langle \text{Sunny, Warm, High, Strong, Cool, Same} \rangle$

$x_2 = \langle \text{Sunny, Warm, High, Light, Warm, Same} \rangle$

$h_1 = \langle \text{Sunny, ?, ?, Strong, ?, ?} \rangle$

$h_2 = \langle \text{Sunny, ?, ?, ?, ?, ?} \rangle$

$h_3 = \langle \text{Sunny, ?, ?, ?, Cool, ?} \rangle$



- 练习：给出下列假设的偏序关系。

h1:      <Sunny, Warm, ?, Strong, ?, ?>

h2:      <Sunny, ?, ?, Strong, ?, ?>

h3:      <Sunny, Warm, ?, ?, ?, ?>

h4:      <?, Warm, ?, Strong, ?, ?>

h5:      <Sunny, ?, ?, ?, ?, ?>

h6:      <?, Warm, ?, ?, ?, ?>



- 一致 ( Consistent )

- 定义： 一个假设  $h$  与训练样例集合  $D$  一致，当且仅当对  $D$  中每一个样例  $\langle x, c(x) \rangle$ ， $h(x) = c(x)$ 。

$$\text{Consistent}(h, D) \equiv (\forall \langle x, c(x) \rangle \in D) h(x) = c(x)$$

- 与“满足”不同

- 一个样例  $x$  在  $h(x) = 1$  时称为满足假设  $h$ ，不论  $x$  是目标概念的正例还是反例。
- 这一样例是否与  $h$  一致与目标概念有关，即是否  $h(x) = c(x)$ 。



## ■ 求解 $h$ 的算法

### ■ Find-S Algorithm

#### ■ 基本思想：

- 使用more\_general\_than偏序的搜索算法
- 沿着偏序链，从较特殊的假设逐渐转移到较一般的假设。
- 在每一步，假设只在需要覆盖新的正例时被泛化。
- 每一步得到的假设，都是在那一点上与训练样例一致的最特殊的假设。



## ■ Find-S Algorithm

- 将  $h$  初始化为  $H$  中 **最特殊假设**
- 对每个 **正例  $x$** 
  - 对  $h$  的每个属性约束  $a_i$ 
    - 如果  $x$  满足  $a_i$
    - 那么不做事
    - 否则，将  **$h$  中  $a_i$**  替换为  $x$  满足的紧邻的更一般约束
- 输出假设  $h$





## ■ Find-S 算法实例 $h_0 \leftarrow \langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$

$x_1 = \langle \text{Sunny Warm Normal Strong Warm Same} \rangle, +$

$h_1 \leftarrow \langle \text{Sunny, Warm, Normal, Strong, Warm, Same} \rangle$

$x_2 = \langle \text{Sunny Warm High Strong Warm Same} \rangle, +$

$h_2 \leftarrow \langle \text{Sunny, Warm, ?, Strong, Warm, Same} \rangle$

$x_3 = \langle \text{Rainy Cold High Strong Warm Change} \rangle, -$

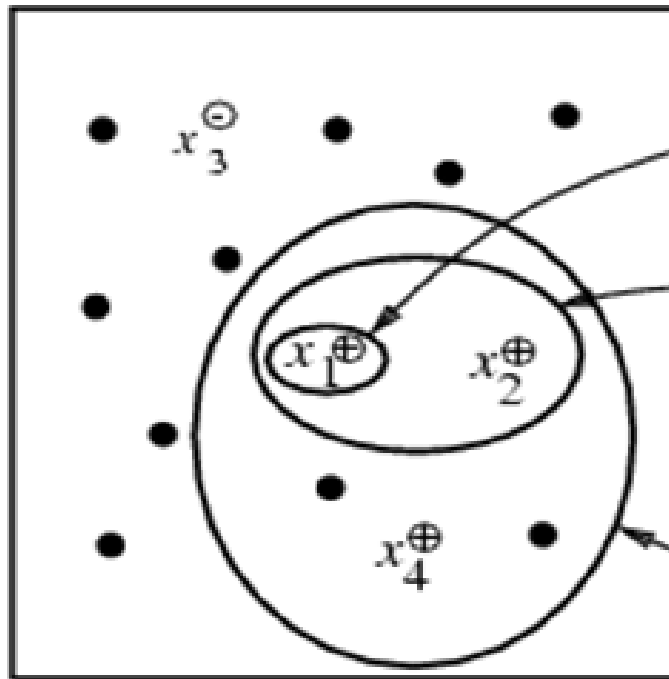
$h_3 \leftarrow \langle \text{Sunny, Warm, ?, Strong, Warm, Same} \rangle$

$x_4 = \langle \text{Sunny Warm High Strong Cool Change} \rangle, +$

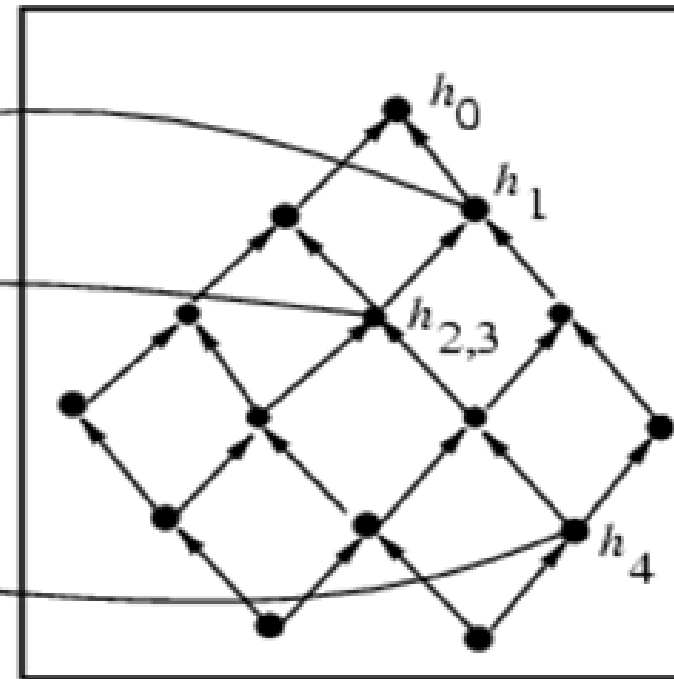
$h_4 \leftarrow \langle \text{Sunny, Warm, ?, Strong, ?, ?} \rangle$



*Instances  $X$*



*Hypotheses  $H$*



Specific

General



## ■ 练习

Outlook	Temperature	Humidity	Wind	PlayTennis
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Mild	High	Strong	No
Sunny	Cool	Normal	Weak	Yes
Overcast	Hot	Normal	Weak	Yes



## ■ Find-S 学习算法的特点及不足

- Find-S 的重要特点：对以属性约束的合取式<sup>^</sup>（conjunction）描述的假设空间  $H$ ，保证输出为  **$H$  中与正例一致的最特殊** 的假设。
- 存在的问题（反例；容错性）



- 变型空间（Version space，版本空间）

- 定义：关于假设空间 $H$ 和训练样例集 $D$ 的变型空间，标记为 $VS_{H,D}$ ，是 $H$ 中与训练样例集 $D$ 一致的所有假设构成的子集。

$$VS_{H,D} = \{h \in H \mid \text{Consistent}(h,D)\}$$

- 与训练样例集一致的所有假设组成的集合
- 包含的是目标概念的所有合理的变型



- 求解 $h$ 的算法
  - 列表后消除算法 ( The List-Then-Eliminate Algorithm )
    - 变型空间  $\text{VersionSpace} \leftarrow$  包含 $H$ 中所有假设的列表
    - 对每个训练样例  $\langle x, c(x) \rangle$  从变型空间中 **移除所有**  $h(x) \neq c(x)$  的假设  $h$
    - 输出  $\text{VersionSpace}$  中的假设列表
      - 只要假设空间是 **有限的**，就可使用
      - 保证得到 **所有** 与训练数据一致的假设
      - 非常繁琐地列出 $H$ 中的所有假设，大多数实际的假设空间 **无法做到**



## ■ 变型空间举例

### ■ EnjoySport

Example	<i>Sky</i>	<i>AirTemp</i>	<i>Humidity</i>	<i>Wind</i>	<i>Water</i>	<i>Forecast</i>	<i>EnjoySport</i>
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

结果如→

每来一个样本，压  
缩一次假设空间

h1:      <Sunny, Warm, ?, Strong, ?, ?>

h2:      <Sunny, ?, ?, Strong, ?, ?>

h3:      <Sunny, Warm, ?, ?, ?, ?>

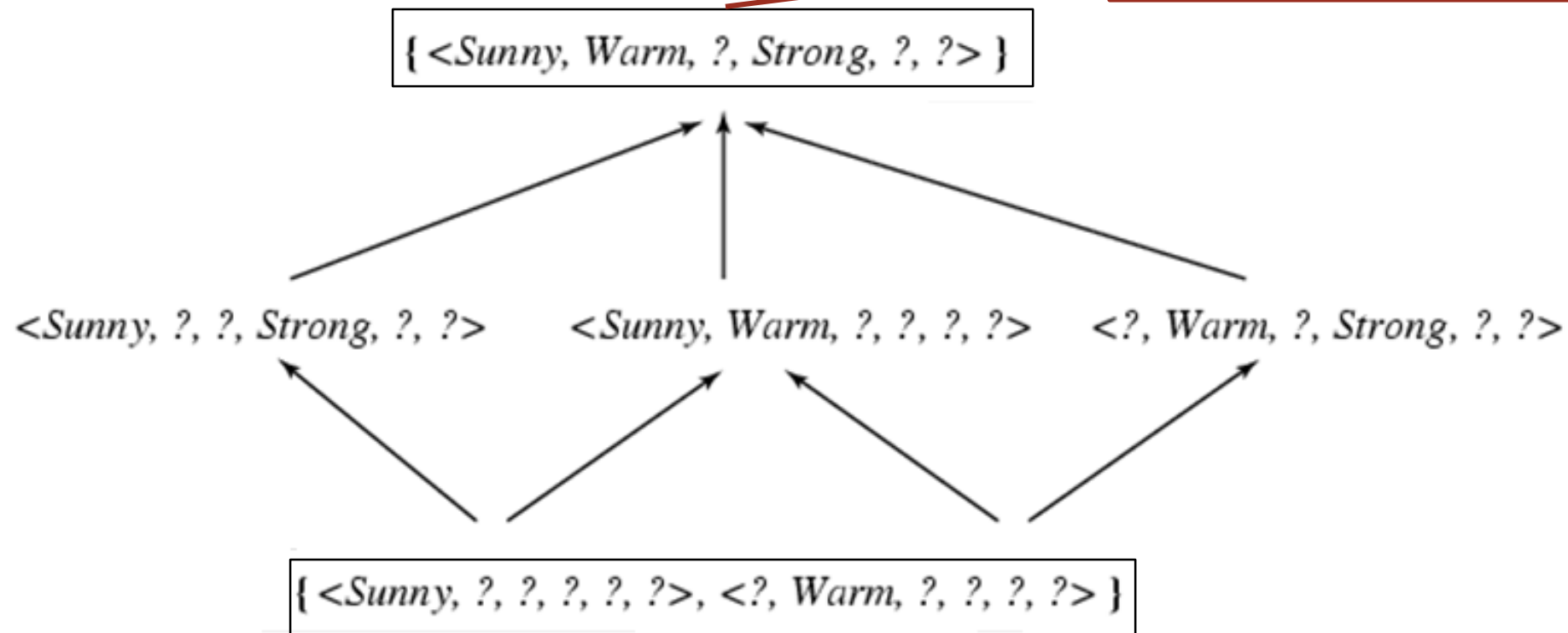
h4:      <?, Warm, ?, Strong, ?, ?>

h5:      <Sunny, ?, ?, ?, ?, ?>

h6:      <?, Warm, ?, ?, ?, ?>



Find-S输出的假设





- 求解**h**的算法： 候选消除算法

- 变型空间的两个边界定义

- 一般边界（ General Boundary ） G:

- H中与D（训练集）相一致的极大一般成员的集合。

$$G \equiv \{ g \in H \mid \text{Consistent}(g, D) \wedge (\neg \exists g' \in H)[(g' >_g g) \wedge \text{Consistent}(g', D)] \}$$

- 特殊边界（ Specific Boundary ） S:

- 在H中与D相一致的极大特殊成员的集合。

$$S \equiv \{ s \in H \mid \text{Consistent}(s, D) \wedge (\neg \exists s' \in H)[(s >_g s') \wedge \text{Consistent}(s', D)] \}$$

- 从变型空间的边界慢慢往中间找



## ■ 求解h的算法

将G集合初始化为H中极大一般假设

## ■ 候选消除算法

将S集合初始化为H中极大特殊假设

对每个训练例d，进行以下操作：

- 如果d是一正例

- 从G中移去所有与d不一致的假设

- 对S中每个与d不一致的假设s

- 从S中移去s

- 把s的所有的极小一般化式h加入到S中，其中h满足

- h与d一致，而且G的某个成员比h更一般

- 从S中移去所有这样的假设：它比S中另一假设更一般

- 如果d是一个反例

- 从S中移去所有d不一致的假设

- 对G中每个与d不一致的假设g

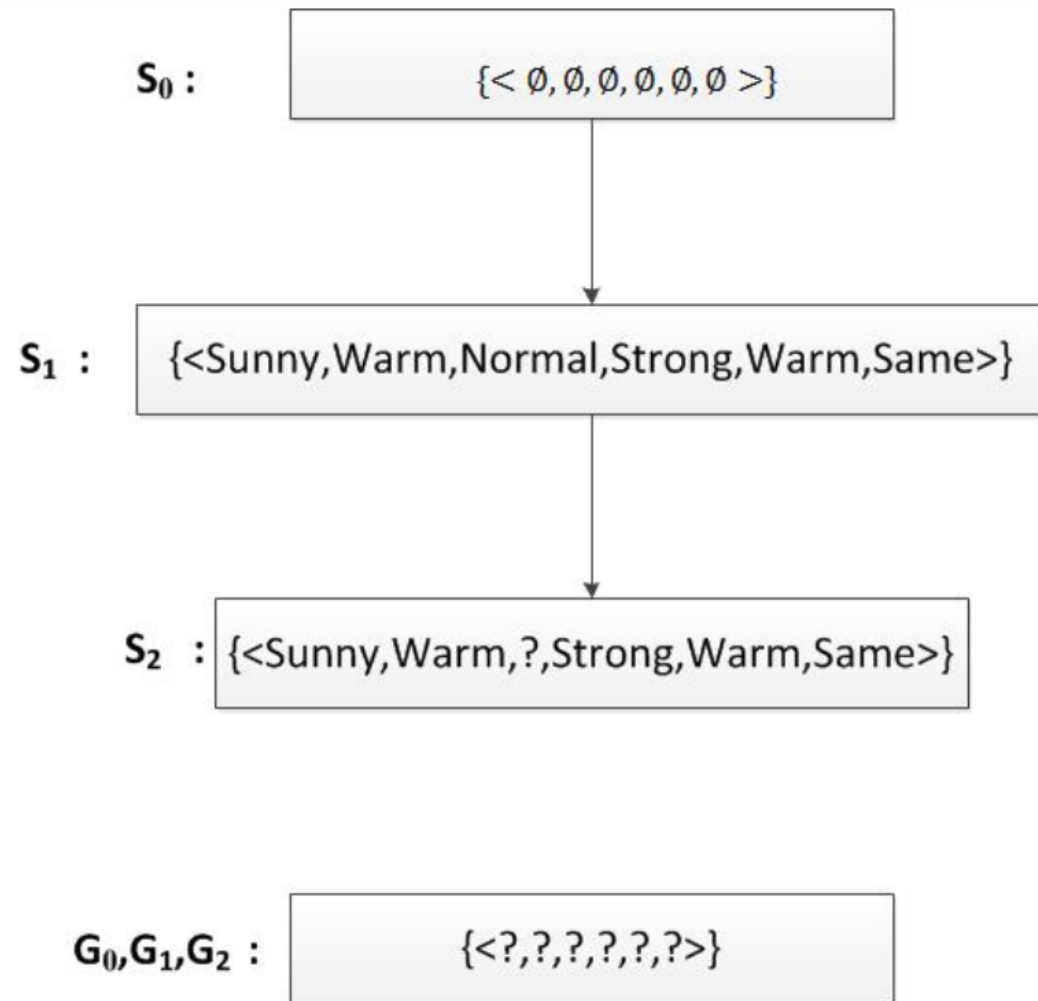
- 从G中移去g

- 把g的所有的极小特殊化式h加入到G中，其中h满足

- h与d一致，而且S的某个成员比h更特殊

- 从G中移去所有这样的假设：它比G中另一假设更特殊





▪ 分别看是否满足S, G

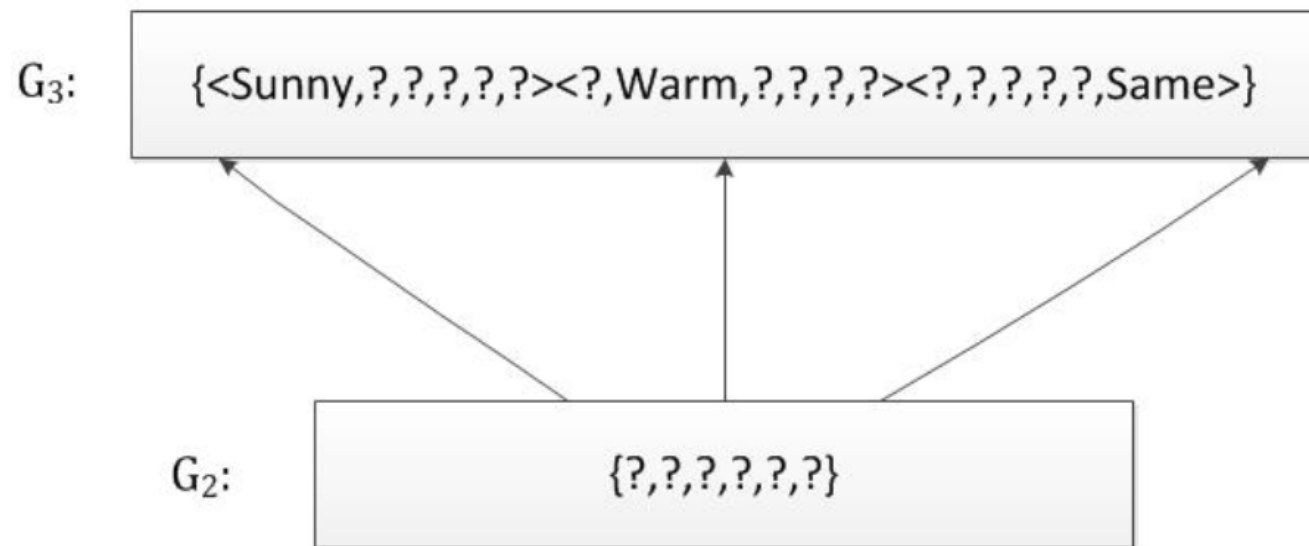
训练样例：

1.  $\langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same} \rangle, \text{EnjoySport} = \text{Yes}$

2.  $\langle \text{Sunny}, \text{Warm}, \text{High}, \text{Strong}, \text{Warm}, \text{Same} \rangle, \text{EnjoySport} = \text{Yes}$



$S_2, S_3$  {<Sunny,Warm,?,Strong,Warm,Same>}

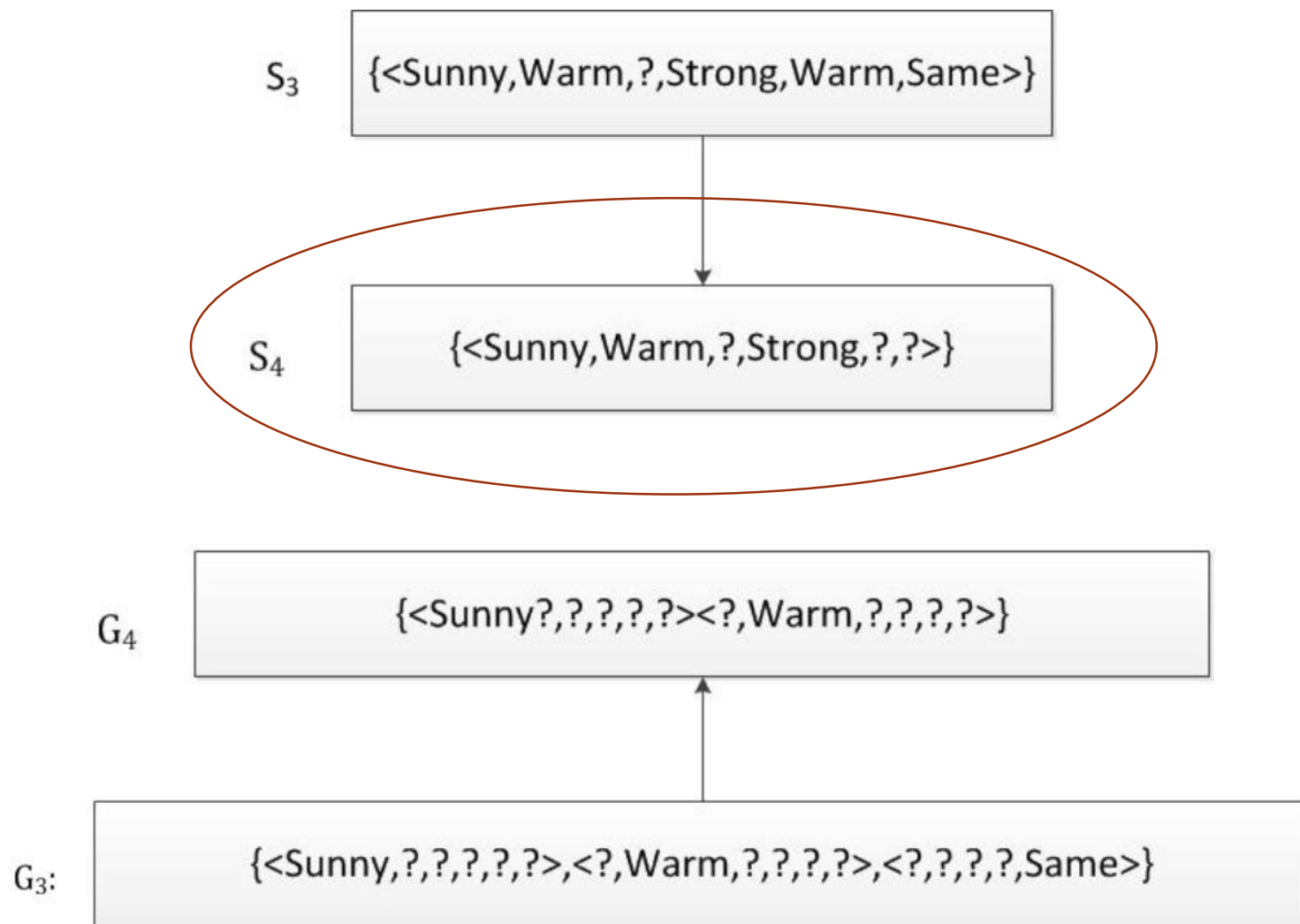


- 分别看是否满足S, G
- $G_3$ 中6选3
- 另外3个不在变型空间中
- $G_3$ 中有多个可选的极大一般假设

训练样例:

3.<Rainy,Cold,High,Strong,Warm,Change>,EnjoySport=No





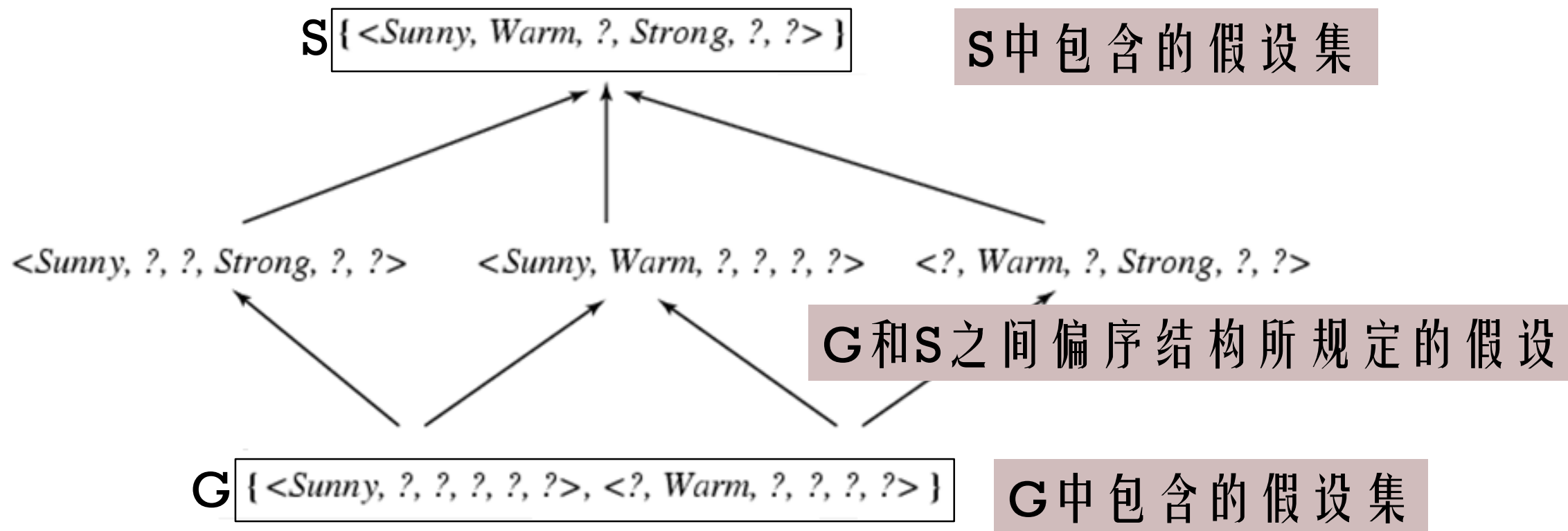
- 分别看是否满足S, G
- 红框为最终结果

训练样例:

4.<Sunny,Warm,High,Storage,Cool,Change>,EnjoySport=Yes



## ■ 本例中的变型空间



- 变型空间表示定理

- 令  $X$  为一任意的实例集合， $H$  为  $X$  上定义的布尔假设的集合。令  $c: X \rightarrow \{0,1\}$  为  $X$  上定义的任一目标概念，并令  $D$  为任一训练样例集合  $\{ \langle x, c(x) \rangle \}$ 。对所有的  $X, H, c, D$  以及定义好的  $S$  和  $G$ ，变型空间表示如下：

$$VS_{H,D} = \{ h \in H \mid (\exists s \in S)(\exists g \in G)(g \geq_g h \geq_g s) \}$$

- 上述定理表明：

- 变型空间由  $G$ 、 $S$ 、 $G$  和  $S$  之间的偏序结构所规定的假设  $h$  组成



## ■ 候选消除算法讨论

- 候选消除算法输出与训练样例一致的所有假设的集合
- 候选消除算法在描述这一集合时不需要明确列举所有成员
- 利用 `more_general_than` 偏序结构，可以得到一个一致假设集合的简洁表示
- 候选消除算法的缺点：同 **Find-S** 一样，容错性能差





- 候选消除算法收敛到正确目标概念的条件
  - 训练样例中没有错误
  - $H$ 中确实包含描述目标概念的正确假设



- 继续探讨：一个有偏的假设空间
  - 在EnjoySport这个例子中，假设空间限制为只包含属性值的合取（同时发生，交集）。（有偏）
  - 这一限制，导致假设空间不能够表示最简单的析取（发生一件，并集）形式的目标概念。

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Cool	Change	Yes
2	Cloudy	Warm	Normal	Strong	Cool	Change	Yes
3	Rainy	Warm	Normal	Strong	Cool	Change	No

“三个属性问题，Sky=Sunny或Sky=Cloudy”



- 归纳推理的一个重要的基本属性：
  - 学习器如果不对目标概念的形式做预先的假定，那么它从根本上无法对未见的样本（实例）进行分类（如：设定为合取形式、选最特殊 $h$ 的作为解）。
- 归纳偏置（**Inductive Bias**）：
  - 对归纳学习进行的某种形式的预先假定



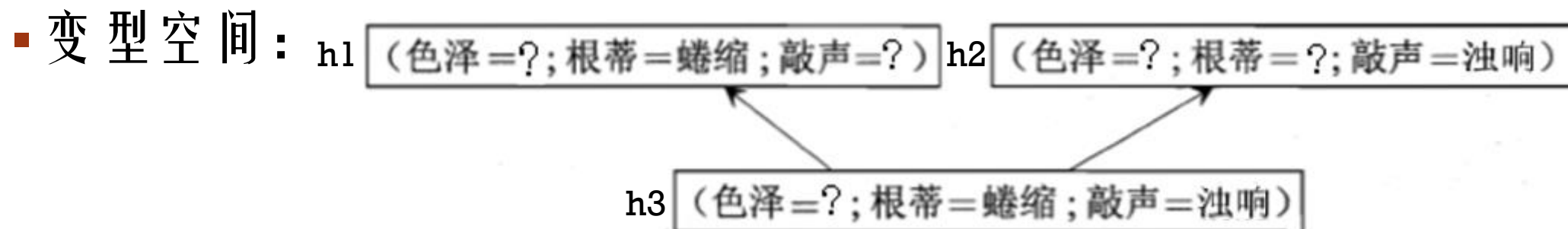
## ■ 归纳偏好（归纳偏置的偏好）：

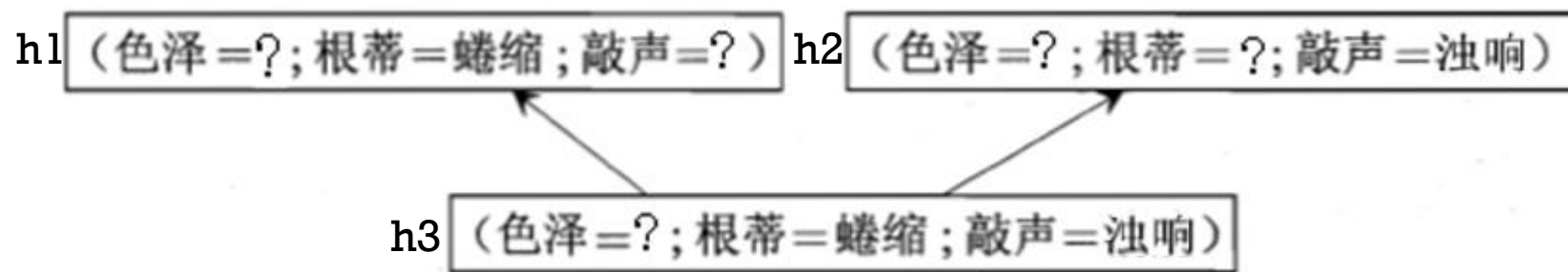
- 假如通过不同学习算法得到三个与训练集一致的假设，但是他们对应的模型在遇到相同的问题时，会产生不同的预测结果。那么，应该选择哪种模型？我们无法通过训练模型得知哪个模型“更好”。这时，学习算法本身的“偏好”就会起到决定性作用。机器学习算法在学习过程中对某种类型假设的偏好，称为：“归纳偏好”。



## ■ 西瓜数据集：学习概念“好瓜”

编号	色泽	根蒂	敲声	好瓜
1	青绿	蜷缩	浊响	是
2	乌黑	蜷缩	浊响	是
3	青绿	硬挺	清脆	否
4	乌黑	稍蜷	沉闷	否





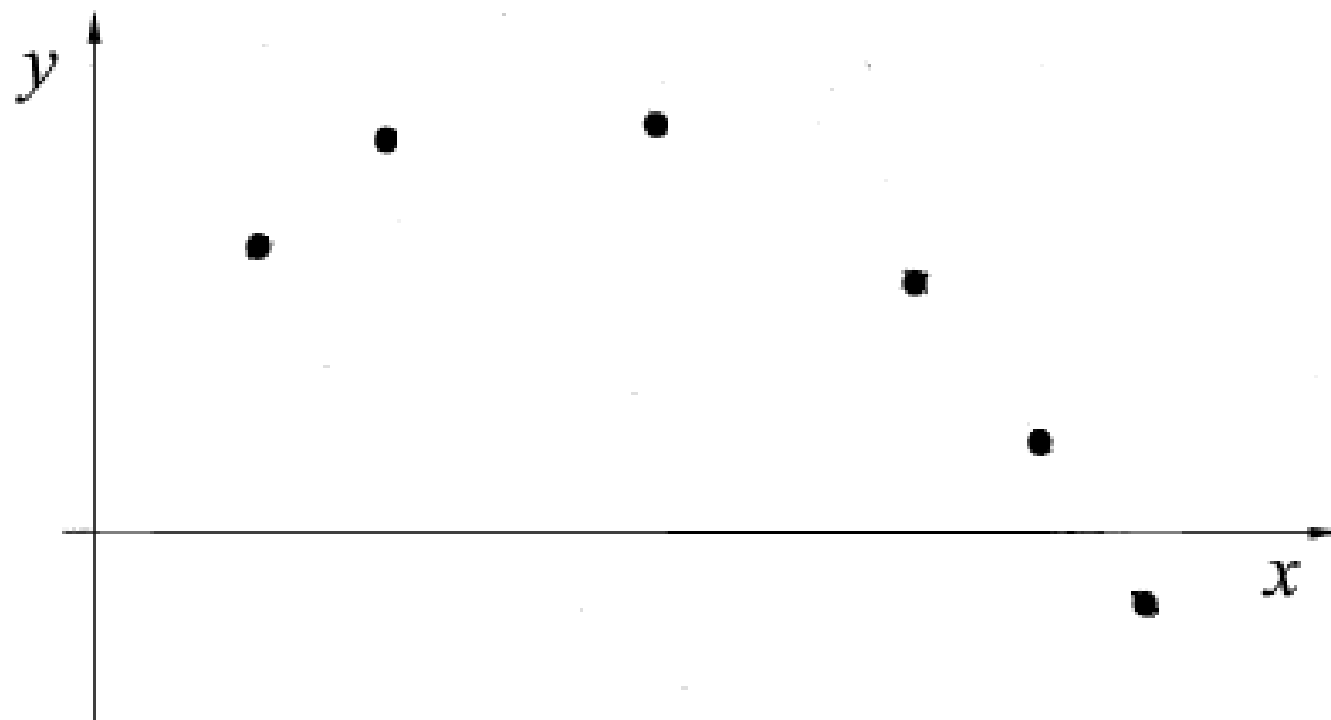
- 新瓜（好瓜）：色泽 = 青绿；根蒂 = 蜷缩；敲声 = 清脆
  - 算法偏好尽可能特殊的模型：否
  - 算法偏好尽可能一般的模型，且由于某种原因它更“相信”根蒂：是
  - 算法偏好多数表决：否



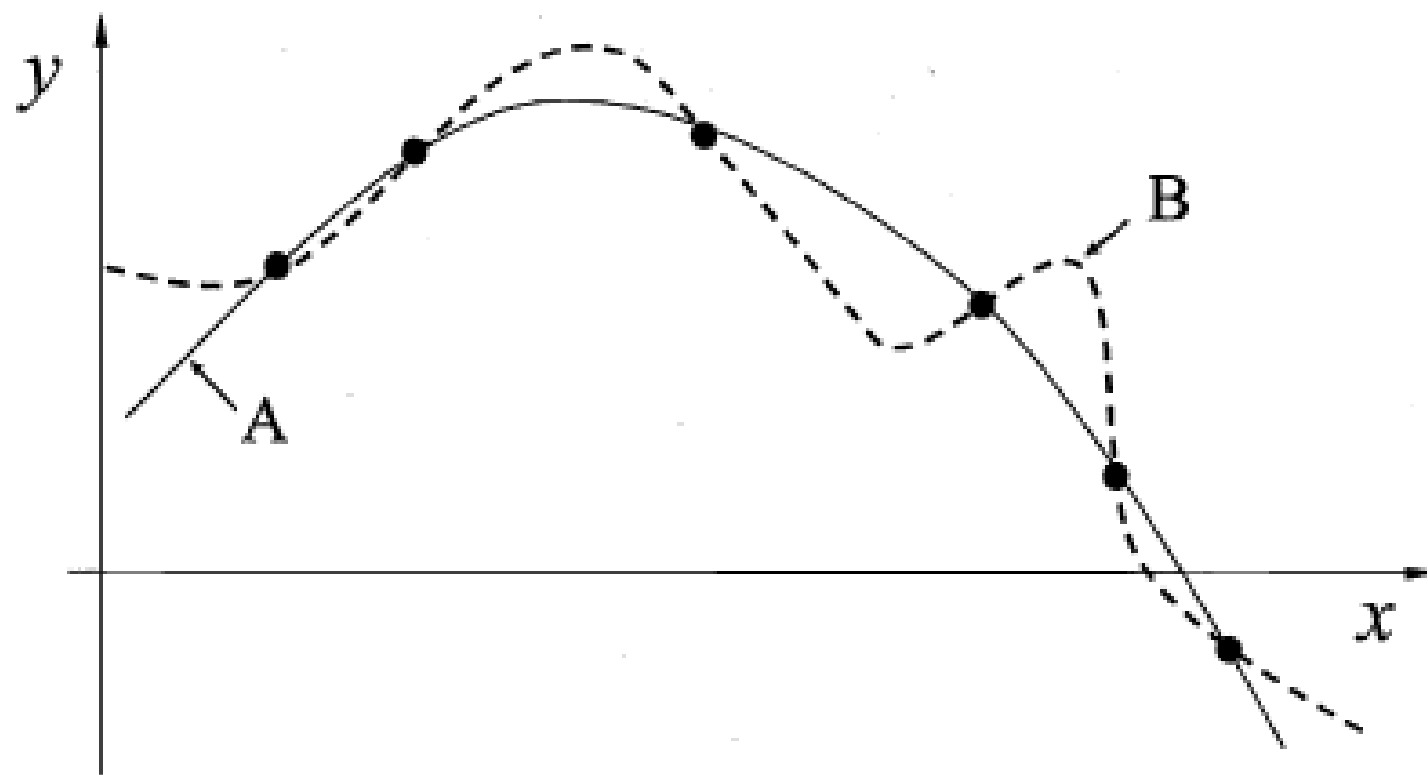
- 归纳偏好常表现为正则化项
- 可依据“奥卡姆剃刀”原理进行选择：若有多个假设与观察一致，则选择最简单的那个（如无必要，勿增实体）



## ■ 考察一个回归学习



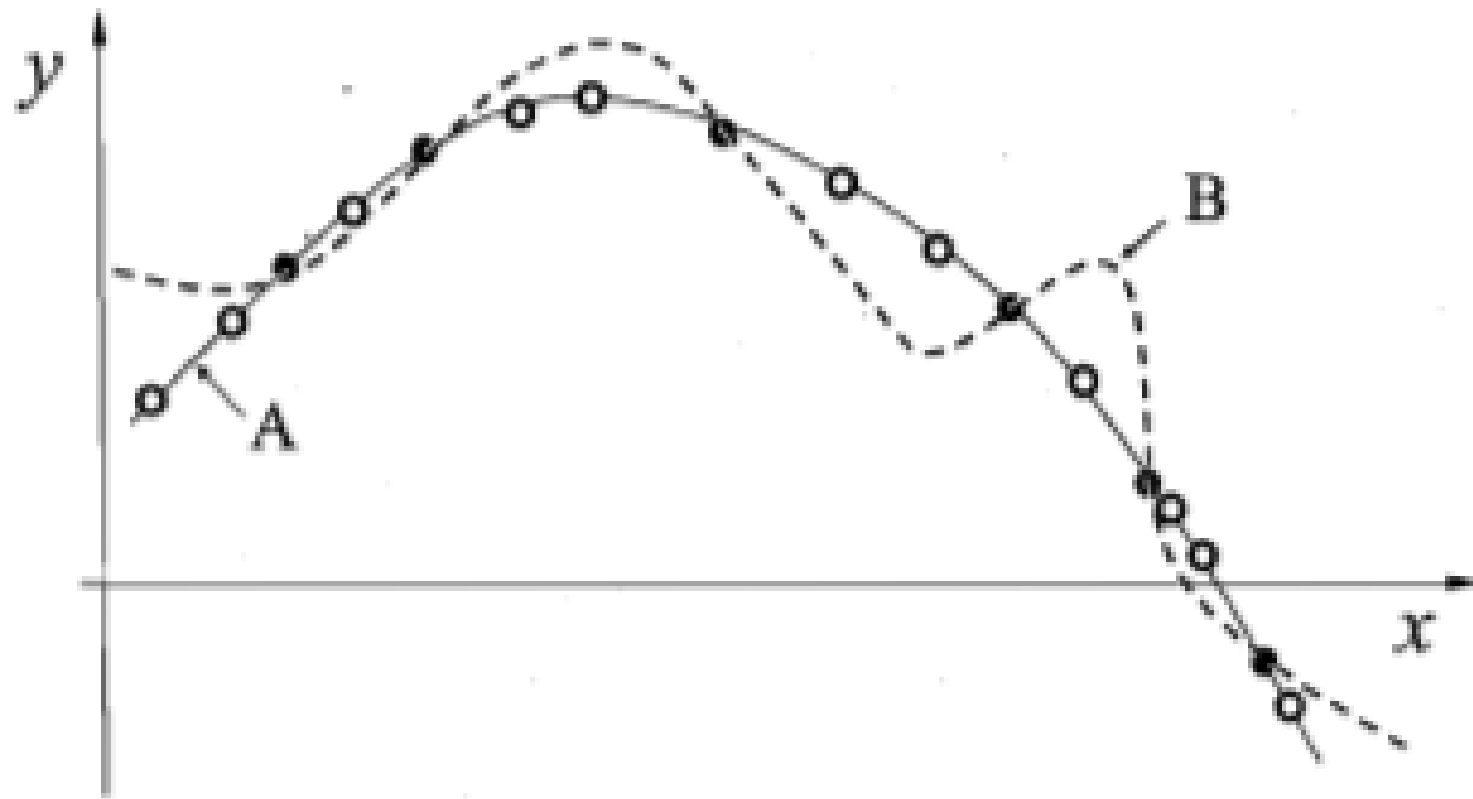




存在多条曲线与有限样本训练集一致

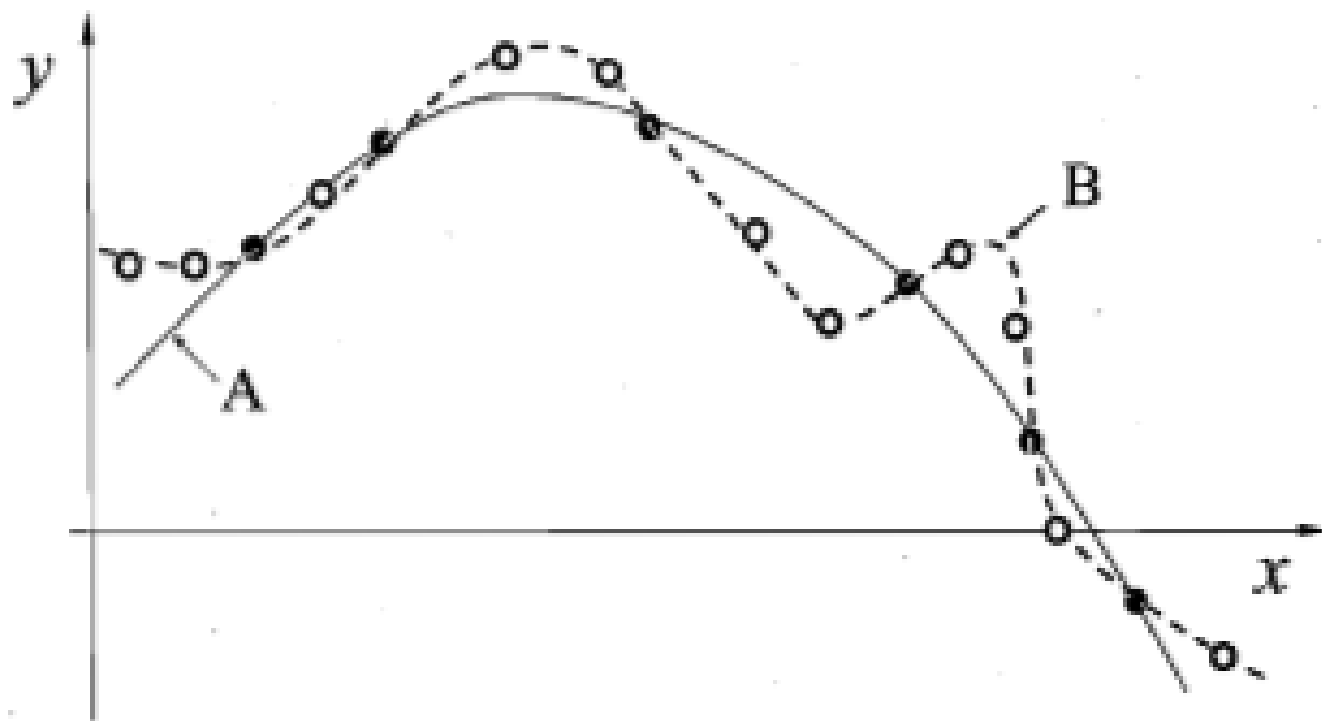
- 通过学习算法得到A、B两条拟合曲线





- 根据“奥卡姆剃刀”原理，A优于B





■ 但B优于A的情况也是完全可能存在的

对于一个学习算法fa，若它在某些问题上比学习算法fb好，则必然存在在另一些问题上，fb比fa好。这个结论对任何算法均成立。“**没有免费的午餐**”定理证实，无论学习算法fa多聪明、学习算法fb多笨拙，它们的期望性能竟然相同（训练集外误差）。



为简单起见, 假设样本空间  $\mathcal{X}$  和假设空间  $\mathcal{H}$  都是离散的. 令  $P(h|X, \mathcal{L}_a)$  代表算法  $\mathcal{L}_a$  基于训练数据  $X$  产生假设  $h$  的概率, 再令  $f$  代表我们希望学习的真实目标函数.  $\mathcal{L}_a$  的“训练集外误差”, 即  $\mathcal{L}_a$  在训练集之外的所有样本上的误差为

$$E_{ote}(\mathcal{L}_a|X, f) = \sum_h \sum_{\mathbf{x} \in \mathcal{X} - X} P(\mathbf{x}) \mathbb{I}(h(\mathbf{x}) \neq f(\mathbf{x})) P(h | X, \mathcal{L}_a) ,$$

其中  $\mathbb{I}(\cdot)$  是指示函数, 若  $\cdot$  为真则取值 1, 否则取值 0.



若  $f$  均匀分布, 则有一半的  $f$  对  $\mathbf{x}$  的预测与  $h(\mathbf{x})$  不一致.

对所有可能的  $f$  按均匀分布对误差求和, 有

$$\begin{aligned}\sum_f E_{ote}(\mathcal{L}_a | X, f) &= \sum_f \sum_h \sum_{\mathbf{x} \in \mathcal{X} - X} P(\mathbf{x}) \mathbb{I}(h(\mathbf{x}) \neq f(\mathbf{x})) P(h | X, \mathcal{L}_a) \\&= \sum_{\mathbf{x} \in \mathcal{X} - X} P(\mathbf{x}) \sum_h P(h | X, \mathcal{L}_a) \sum_f \mathbb{I}(h(\mathbf{x}) \neq f(\mathbf{x})) \\&= \sum_{\mathbf{x} \in \mathcal{X} - X} P(\mathbf{x}) \sum_h P(h | X, \mathcal{L}_a) \frac{1}{2} 2^{|\mathcal{X}|} \\&= \frac{1}{2} 2^{|\mathcal{X}|} \sum_{\mathbf{x} \in \mathcal{X} - X} P(\mathbf{x}) \sum_h P(h | X, \mathcal{L}_a) \\&= 2^{|\mathcal{X}|-1} \sum_{\mathbf{x} \in \mathcal{X} - X} P(\mathbf{x}) \cdot 1.\end{aligned}$$

总误差竟然与学习算法无关!



对于任意两个学习算法  $\mathcal{L}_a$  和  $\mathcal{L}_b$ ,

$$\sum_f E_{ote}(\mathcal{L}_a|X, f) = \sum_f E_{ote}(\mathcal{L}_b|X, f)$$

这就是“没有免费的午餐”定理 (No Free Lunch Theorem, 简称 NFL 定理)  
[Wolpert, 1996; Wolpert and Macready, 1995].

- 优化算法的等价性
- 任何优化算法都不比穷举法好
- 为什么还要研究最优化和机器学习算法呢?



- NFL定理有一个重要前提：所有“问题”出现的机会相同、或所有问题同等重要。但实际情形并不是这样的。很多时候，我们只关注自己正在试图解决的问题，希望为它找到一个解决方案，至于这个解决方案在别的问题上是否为好方案，我们并不关心。
- 脱离具体问题，空泛的谈论“什么学习算法更好”毫无意义
- 收敛速度





# 模型评估与选择

Model Assessment and Selection



# 损失函数 ( LOSS FUNCTION )

(1) 0-1 损失函数 (0-1 loss function)

$$L(Y, f(X)) = \begin{cases} 1, & Y \neq f(X) \\ 0, & Y = f(X) \end{cases}$$

(2) 平方损失函数 (quadratic loss function)

$$L(Y, f(X)) = (Y - f(X))^2$$

(3) 绝对损失函数 (absolute loss function)

$$L(Y, f(X)) = |Y - f(X)|$$

(4) 对数损失函数 (logarithmic loss function) 或对数似然损失函数 (log-likelihood loss function)

$$L(Y, P(Y | X)) = -\log P(Y | X)$$



# 风险函数 ( RISK FUNCTION )

- 理论上  $f(X)$  基于联合分布  $P(X,Y)$  的平均意义下的损失，即 **期望损失 ( Expected Loss )**、**风险函数**
- 风险函数可理解为样本总体的损失

$$R_{\text{exp}}(f) = E_P[L(Y, f(X))] = \int_{X \times Y} L(y, f(x)) P(x, y) dx dy$$

未知



# 经验风险 ( EMPIRICAL RISK )

- $f(\mathbf{x})$  关于训练数据集  $T = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$  的平均损失

$$R_{\text{emp}}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$



# 模型选择策略

- 经验风险最小化 ( Empirical Risk Minimization, ERM ) 求解最优化问题:

$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$

- 结构风险最小化 ( Structural Risk Minimization, SRM ) 求解最优化问题:

$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f)$$



# 模型性能度量方法

- 错误率（error rate）和准确度（accuracy）
  - 错误率：分类问题中，误分类样本的比例
  - 准确度：正确分类样本的比例，即正确率
- 对样本集 $\mathcal{D}$

$$E(f; D) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(\mathbf{x}_i) \neq y_i) .$$

$$\begin{aligned} \text{acc}(f; D) &= \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(\mathbf{x}_i) = y_i) \\ &= 1 - E(f; D) . \end{aligned}$$

$$E(f; \mathcal{D}) = \int_{\mathbf{x} \sim \mathcal{D}} \mathbb{I}(f(\mathbf{x}) \neq y) p(\mathbf{x}) d\mathbf{x} ,$$

$$\begin{aligned} \text{acc}(f; \mathcal{D}) &= \int_{\mathbf{x} \sim \mathcal{D}} \mathbb{I}(f(\mathbf{x}) = y) p(\mathbf{x}) d\mathbf{x} \\ &= 1 - E(f; \mathcal{D}) . \end{aligned}$$



- 训练误差（training error）和泛化误差（generalization error）

- 训练误差：关于训练数据集的平均损失

$$R_{\text{emp}}(\hat{f}) = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}(x_i))$$

- 泛化误差：对未知数据预测的误差

$$R_{\text{exp}}(\hat{f}) = E_P[L(Y, \hat{f}(X))] = \int_{\mathcal{X} \times \mathcal{Y}} L(y, \hat{f}(x)) P(x, y) dx dy$$



# 二分类问题的泛化误差上界

训练数据集  $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$  从联合概率分布  $P(X, Y)$  独立同分布产生的,  $X \in \mathbf{R}^n$ ,  $Y \in \{-1, +1\}$ . 设  $f$  是从  $\mathcal{F} = \{f_1, f_2, \dots, f_d\}$  中选取的函数. 损失函数是 0-1 损失.

关于  $f$  的期望风险和经验风险分别是

$$R(f) = E[L(Y, f(X))]$$

$$\hat{R}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$



经验风险最小化函数是

$$f_N = \arg \min_{f \in \mathcal{F}} \hat{R}(f)$$

$f_N$  的泛化能力

$$R(f_N) = E[L(Y, f_N(X))]$$





**定理（泛化误差上界）** 对二类分类问题，当假设空间是有限个函数的集合  $\mathcal{F} = \{f_1, f_2, \dots, f_d\}$  时，对任意一个函数  $f \in \mathcal{F}$ ，至少以概率  $1 - \delta$ ，以下不等式成立：

训练误差小的模型泛化误差也会小。

$$R(f) \leq \hat{R}(f) + \varepsilon(d, N, \delta)$$

$\mathcal{F}$ 中包含的函数越多，泛化误差上界越大。

样本容量  $N$  越大，训练误差与泛化误差越接近。

其中，

$$\varepsilon(d, N, \delta) = \sqrt{\frac{1}{2N} \left( \log d + \log \frac{1}{\delta} \right)}$$

**Hoeffding 不等式**

设  $S_n = \sum_{i=1}^n X_i$  是独立随机变量  $X_1, X_2, \dots, X_n$  之和， $X_i \in [a_i, b_i]$ ，则对任意  $t > 0$ ，

以下不等式成立：

随机变量的和与其期望值的偏差的概率上界

$$P(ES_n - S_n \geq t) \leq \exp \left( \frac{-2t^2}{\sum_{i=1}^n (b_i - a_i)^2} \right) \quad (1)$$

影响上界的因素：数量多，差异大

- 泛化误差的估计
  - 用训练误差（ training error ）估计
  - 其他理论方法？
- 测试集（ testing set ）与测试误差（ testing error ）
  - 假定从真实分布中独立采样得到
  - 与训练集互斥



ID	X1	X2	Y
1	A	7	T
2	A	7	T
3	B	5	T
4	A	3	F
5	A	2	F
6	A	6	F
7	A	7	F
8	A	7	T
9	B	2	T
10	A	3	F

ID	X1	X2	Y
1	A	7	T
2	A	7	T
4	A	3	F
6	A	6	F
9	B	2	T
10	A	3	F

IF X1 = B OR X2 > 6  
THEN Y = T

ID	X1	X2	Y
3	B	5	T
5	A	2	F
7	A	7	F
8	A	7	T

ID	X1	X2	Y
3	B	5	T
5	A	2	F
7	A	7	T
8	A	7	T

测试误差：0.25



# 准确率的局限性

- 考虑一个分类问题：
  - 类0中的样本数为 9990
  - 类1中的样本数为 10
  - 若模型将所有样例类别预测为类0，则分类准确率为  $9990/10000 = 99.9\%$ 
    - 模型没有发现任何类1的样例，因而准确率具有误导性（样本不均衡）。



# 混淆矩阵 ( CONFUSION MATRIX )

ACTUAL CLASS	PREDICTED CLASS	
	Class=Yes	Class=No
Class=Yes	a (TP)	b (FN)
	c (FP)	d (TN)

- a: TP (true positive) 真正，对应于被分类模型正确预测的正样本数
- b: FN (false negative) 假负，对应于被分类模型错误预测为负类的正样本数
- c: FP (false positive) 假正，对应于被分类模型错误预测为正类的负样本数
- d: TN (true negative) 真负，对应于被分类模型正确预测的负样本数



- 真正率（TPR）或称灵敏度（Sensitivity）
  - 模型正确预测的正样本的比例
  - $TPR = TP / (TP + FN) = a / (a + b)$
- 真负率（TNR）或称特指度（Specificity）
  - 模型正确预测的负样本的比例
  - $TNR = TN / (TN + FP) = d / (c + d)$
- 假正率（FPR）
  - 被预测为正类的负样本的比例
  - $FPR = FP / (TN + FP) = c / (c + d)$
- 假负率（FNR）
  - 被预测为负类的正样本的比例
  - $FNR = FN / (TP + FN) = b / (a + b)$

	PREDICTED CLASS		
		Class=Yes	Class=No
	Class=Yes	a (TP)	b (FN)
ACTUAL CLASS	Class=No	c (FP)	d (TN)



	PREDICTED CLASS	
	Class=Yes	Class=No
ACTUAL CLASS	Class=Yes a (TP)	b (FN)
	Class=No c (FP)	d (TN)

- **精度  $p$** :  $\text{Precision}(p) = \frac{a}{a + c}$ 
  - 也称查准率，确定在分类器断言为正类的那部分记录中实际为正类的记录所占的比率。
  - 精度越高，分类器的假正错误率就越低
- **召回率  $r$** :  $\text{Recall}(r) = \frac{a}{a + b}$ 
  - 度量被分类器正确预测的正样本的比例，亦称查全率。
  - 具有高召回率的分类器很少将正样本误分为负样本
- **$F_1$  度量**:  $F_1\text{-measure}(F) = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$ 
  - 精度和召回率的调和平均



# P-R曲线和P-R图-评估某分类器的能力

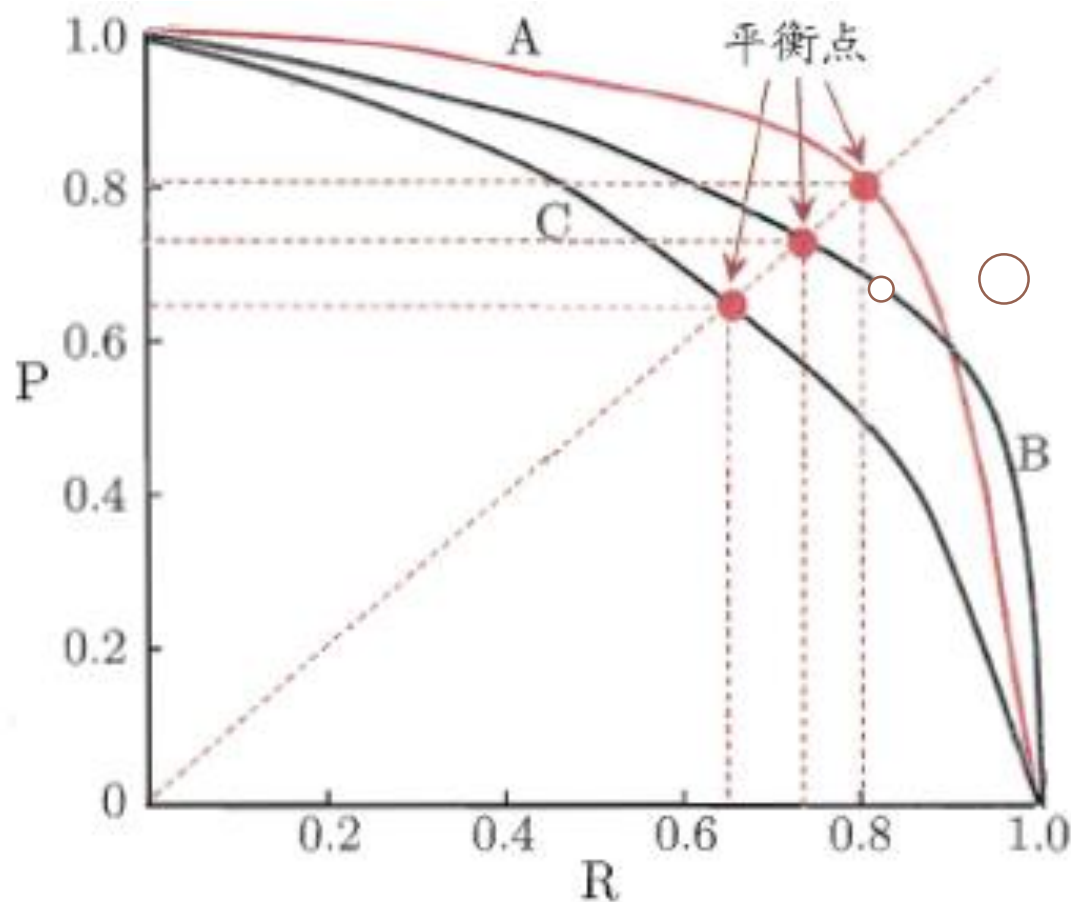
- 根据学习器的预测结果对样例进行排序（单学习器）
  - “最可能”是正例的样本排在最前面；
  - “最不可能”是正例的样本排在最后面；
- 按此顺序逐个将样本作为正例进行预测，每次计算 $p$ 、 $r$ 值
- 以 $p$ 为纵轴， $r$ 为横轴作图，即得“P-R曲线”
- 显示该曲线的图称为“P-R图”





# P-R曲线和P-R图

		PREDICTED CLASS	
ACTUAL CLASS		Class=Yes	Class=No
	Class=Yes	a (TP)	b (FN)
	Class=No	c (FP)	d (TN)



学习器A  
优于学  
习器C

- 精度 =  $a/(a+c)$
- 召回率 =  $a/(a+b)$

P-R曲线与平衡点示意图



# 多个二分类混淆矩阵

- 各混淆矩阵上度量的平均值

$$\text{macro-}P = \frac{1}{n} \sum_{i=1}^n P_i ,$$

$$\text{macro-}R = \frac{1}{n} \sum_{i=1}^n R_i ,$$

$$\text{macro-}F1 = \frac{2 \times \text{macro-}P \times \text{macro-}R}{\text{macro-}P + \text{macro-}R} .$$



- 将混淆矩阵对应元素平均后再求指标值

$$\text{micro-}P = \frac{\overline{TP}}{\overline{TP} + \overline{FP}} ,$$

$$\text{micro-}R = \frac{\overline{TP}}{\overline{TP} + \overline{FN}} ,$$

$$\text{micro-}F1 = \frac{2 \times \text{micro-}P \times \text{micro-}R}{\text{micro-}P + \text{micro-}R} .$$

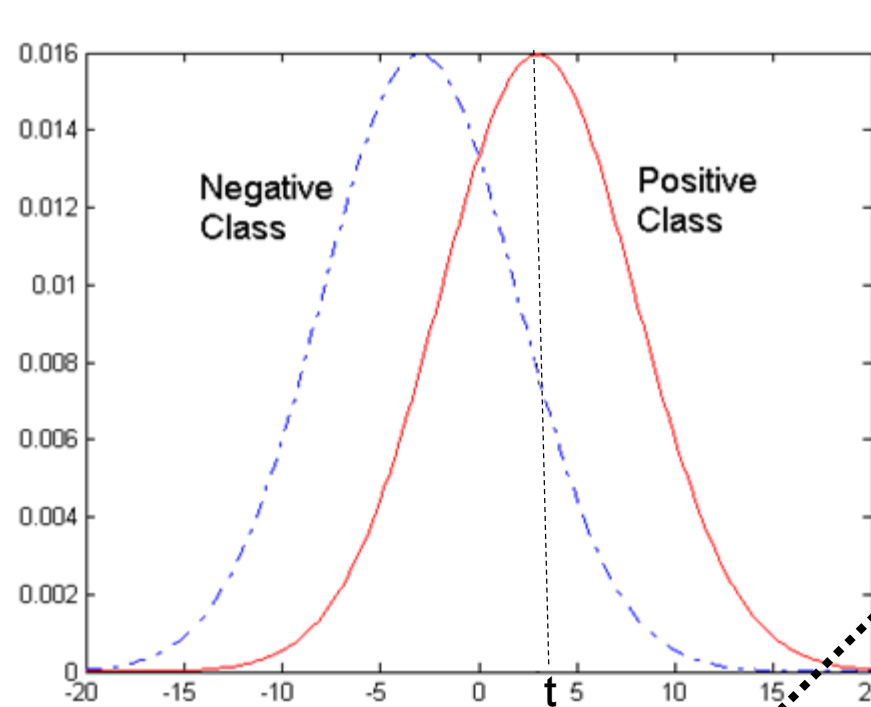


# 接受者操作特征 ( RECEIVER OPERATING CHARACTERISTIC, ROC ) 曲线—评估某模型能力

- 显示分类器真正率和假正率之间折中的一种图形化方法
- 真正率 ( TPR ) 沿y轴绘制 ( 召回率 )
- 假正率 ( FPR ) 沿x轴绘制
- 沿着曲线的每个点对应于一个分类器归纳的模型

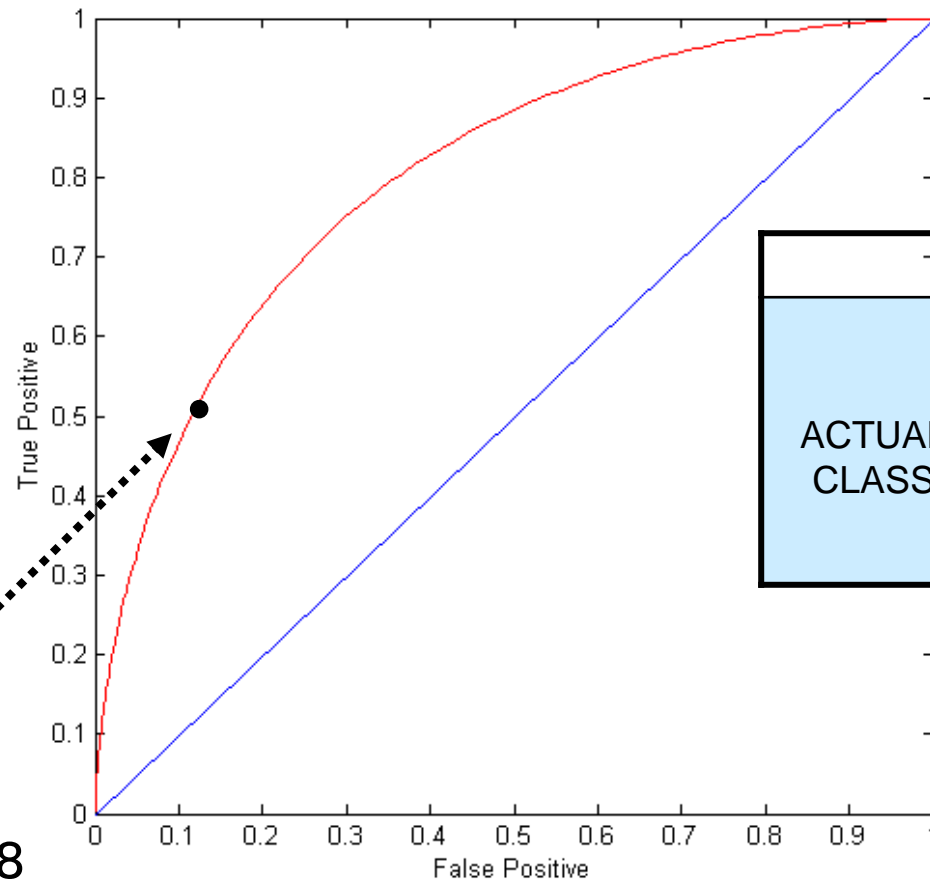


- 1-dimensional data set containing 2 classes (positive and negative)
- any points located at  $x > t$  is classified as positive



At threshold  $t$ :

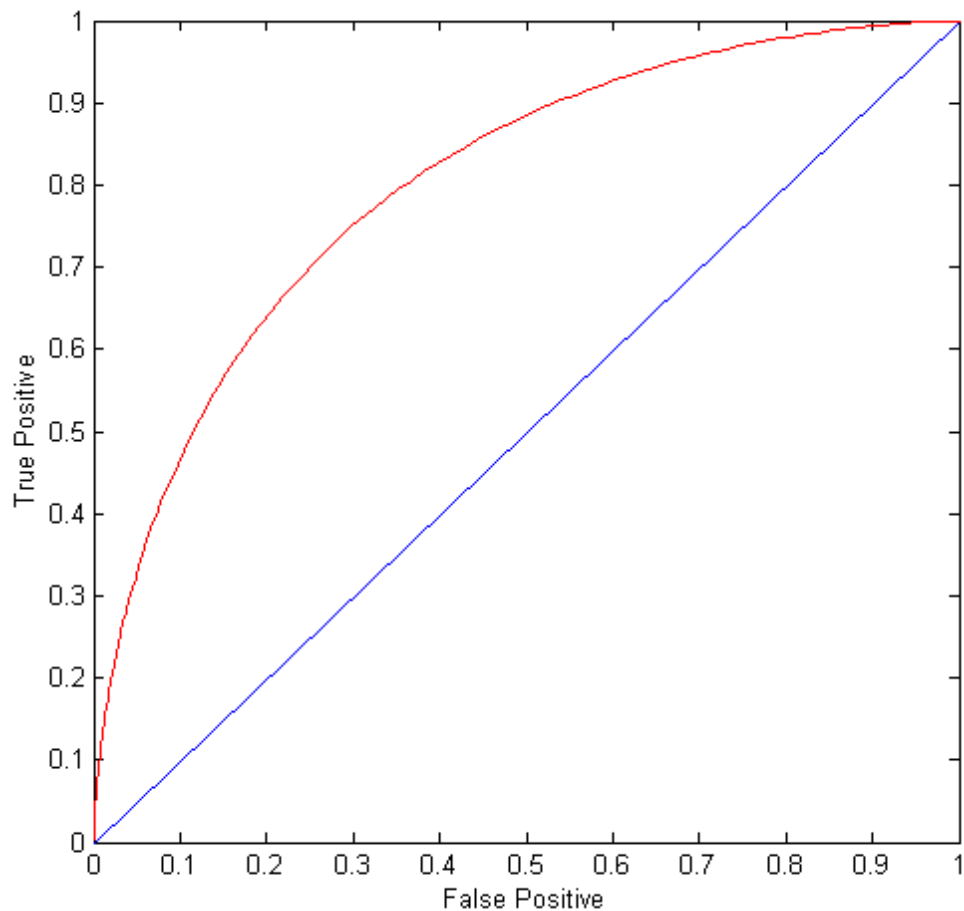
$TP=0.5$ ,  $FN=0.5$ ,  $FP=0.12$ ,  $TN=0.88$



		PREDICTED CLASS	
ACTUAL CLASS	Class=Yes	Class=Yes	Class=No
	Class=No	a (TP)	b (FN)
	Class=No	c (FP)	d (TN)

- $TPR = a / (a + b)$
- $FPR = c / (c + d)$





## ■ ROC 曲线

- $(\text{TPR}=0, \text{FPR}=0)$ : 把每个实例都预测为负类的模型
- $(\text{TPR}=1, \text{FPR}=1)$ : 把每个实例都预测为正类的模型
- $(\text{TPR}=1, \text{FPR}=0)$ : 理想模型
- 对角线: 随机猜测的模型
- 靠近图左上角的分类器是最优的



## ■ 如何绘制ROC曲线

(1) 假定为正类定义了连续值输出，对检验记录按它们的输出值递增排序。

(2) 选择秩最低的检验记录（即输出值最低的记录），把选择的记录以及那些秩高于它的记录指派为正类。这种方法等价于把所有的检验实例都分为正类。因为所有的正检验实例都被正确分类，而所有的负测试实例都被误分，因此  $TPR=FPR=1$ 。

(3) 从排序列表中选择下一个检验记录，把选择的记录以及那些秩高于它的记录指派为正类，而把那些秩低于它的记录指派为负类。通过考察前面选择的记录的实际类标号来更新  $TP$  和  $FP$  计数。如果前面选择的记录为正类，则  $TP$  计数减少而  $FP$  计数不变。如果前面选择的记录为负类，则  $FP$  计数减少而  $TP$  计数不变。

(4) 重复步骤 3 并相应地更新  $TP$  和  $FP$  计数，直到最高秩的记录被选择。

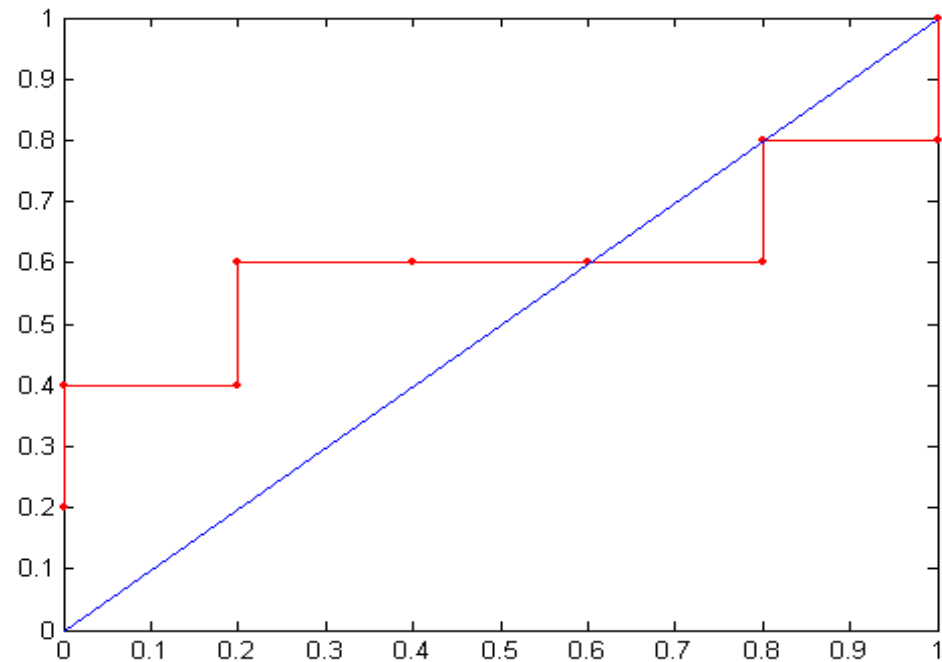
(5) 根据分类器的  $FPR$  画出  $TPR$  曲线。



- 对十个样本用不同阈值进行分类，5正5负

Class	+	-	+	-	-	-	+	-	+	+	
閾值	0.25	0.43	0.53	0.76	0.85	0.85	0.85	0.87	0.93	0.95	1.00
TP	5	4	4	3	3	3	3	2	2	1	0
FP	5	5	4	4	3	2	1	1	0	0	0
TN	0	0	1	1	2	3	4	4	5	5	5
FN	0	1	1	2	2	2	2	3	3	4	5
TPR	1	0.8	0.8	0.6	0.6	0.6	0.6	0.4	0.4	0.2	0
FPR	1	1	0.8	0.8	0.6	0.4	0.2	0.2	0	0	0

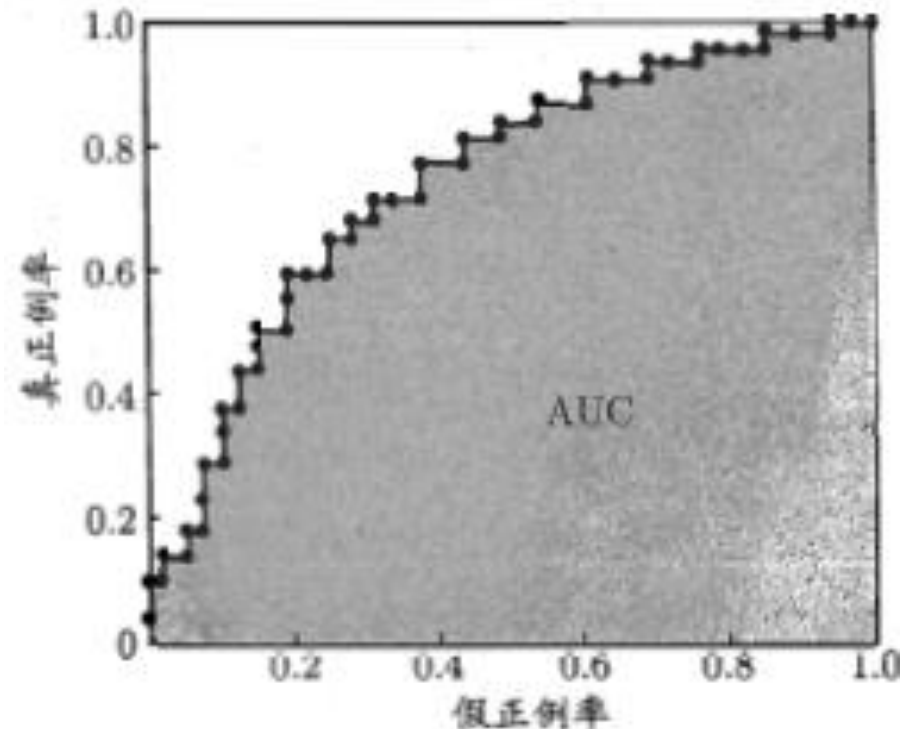
ROC Curve:





# AUC ( AREA UNDER ROC CURVE )

- 即 ROC 曲线下的面积
- 可以评价模型的性能



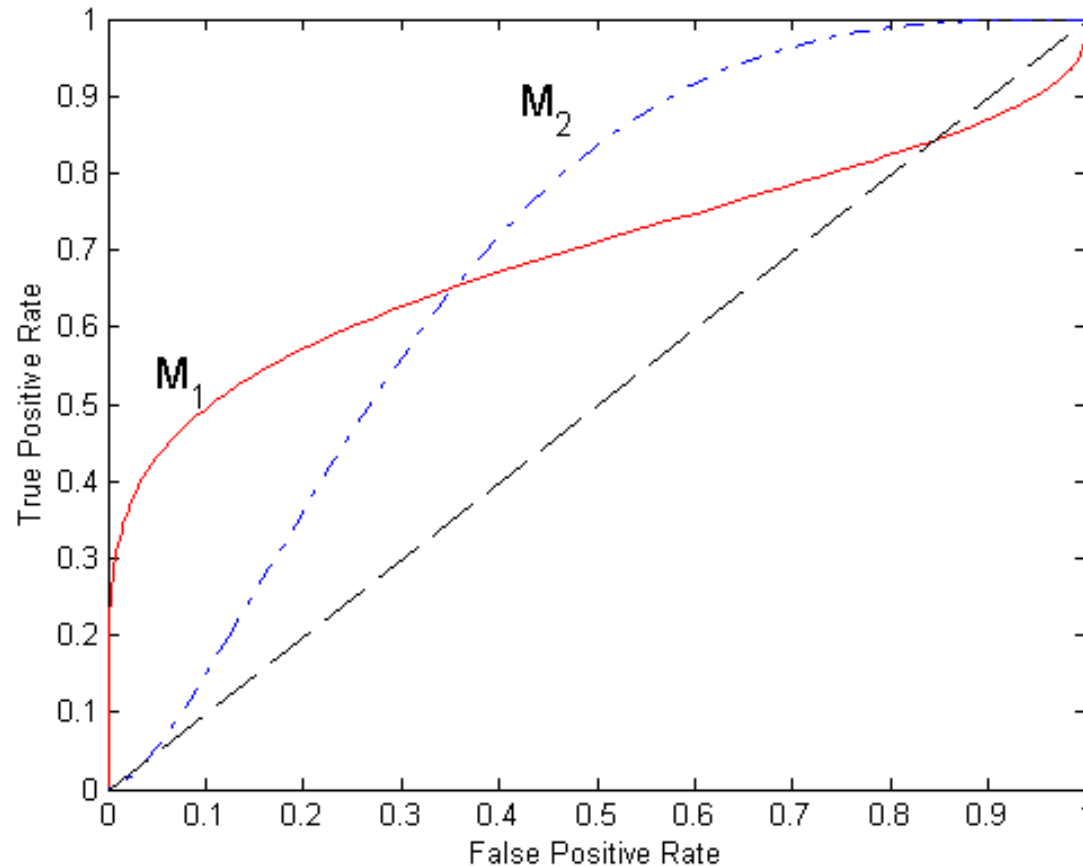
基于有限样例绘制的 ROC 曲线与 AUC

假定 ROC 曲线是由坐标为  $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$  的点按序连接而形成( $x_1 = 0, x_m = 1$ )

$$\text{AUC} = \frac{1}{2} \sum_{i=1}^{m-1} (x_{i+1} - x_i) \cdot (y_i + y_{i+1}) .$$



## ■ Using ROC for Model Comparison



■ No model consistently outperform the other

- $M_1$  is better for small FPR
- $M_2$  is better for large FPR

■ AUC



# PR 曲线与 ROC 曲线

- 例：计算下述二分类器的 TPR, FPR, P, R

M1	PREDICTED CLASS		
ACTUAL CLASS		Class=Yes	Class=No
	Class=Yes	90	10
	Class=No	10	1,999,890

M1: TPR=0.9,  
FPR=0.00000500025,  
P=0.9  
R=0.9

M2	PREDICTED CLASS		
ACTUAL CLASS		Class=Yes	Class=No
	Class=Yes	90	10
	Class=No	910	1,998,990

M2: TPR=0.9,  
FPR=0.00045502275,  
P=0.09  
R=0.9



# 代价矩阵

	PREDICTED CLASS		
	$C(i j)$	Class=Yes	Class=No
	Class=Yes	$C(\text{Yes} \text{Yes})$	$C(\text{No} \text{Yes})$
	Class=No	$C(\text{Yes} \text{No})$	$C(\text{No} \text{No})$

$C(i|j)$ : 预测（误分类）一个j类记录为i类的代价（医院诊断）



## ■ 计算分类代价

犯假负错误的代价是  
犯假正错误的100倍

尽管模型 $M_2$ 改善了准确率，  
但仍然较差。因这些改善  
是建立在增加代价更高的  
假负错误之上的。

Cost Matrix	PREDICTED CLASS		
	C(i j)	+	-
ACTUAL CLASS	+	-1	100
	-	1	0

标准的准确率度量趋向  
于 $M_2$ 优于 $M_1$

Model $M_1$	PREDICTED CLASS		
ACTUAL CLASS		+	-
	+	220	30
	-	70	180

Accuracy = 80%

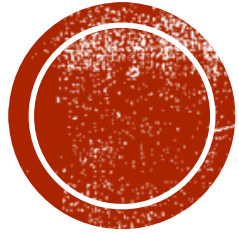
Cost = 2850

Model $M_2$	PREDICTED CLASS		
ACTUAL CLASS		+	-
	+	210	40
	-	10	240

Accuracy = 90%

Cost = 3800





**THE END!**

