

机器学习

李成龙

安徽大学人工智能学院

“多模态认知计算”安徽省重点实验室

- 什么是机器学习
- 机器如何学习
- 如何让机器学习的更好
- 为什么机器能学习

- 什么是机器学习
 - 机器学习的发展历史与背景
 - 机器学习的基本概念

- 必要性
- 定义
- 典型机器学习过程
- 基本术语
- 机器学习三要素

本节目录

- 必要性
- 定义
- 典型机器学习过程
- 基本术语
- 机器学习三要素

- 必要性

- 如何编程实现对“树”的判断？



- 编程显示的“定义”出树：很难

基本字义

树（樹） shù (尸乂)

- 1、木本植物的通称：树木。树林。树大根深（喻势力大，根基牢固）。
- 2、种植，培育：树艺（“艺”，种植）。树荆棘得刺，树桃李得荫。
- 3、立，建立：树立。树敌。
- 4、量词，相当于“株”、“棵”：一树梅花。
- 5、姓。

- 三岁的小孩却能从不断观察中学习识别，几乎不会出错
 - 机器学习：可以更“简单”的实现手工编程无法解决的复杂系统问题

- 必要性

- 系统太复杂，无法显示的编程解决：
 - 自动驾驶
- 无法明确定义出一个解决方案：
 - 图像识别
- 需要非常快速的判断和决策：
 - 高频交易
- 需要处理非常大量的数据：
 - 欧洲核子对撞机撞出来的数据

- 必要性

- 应用广泛

- 衣：淘宝同款、时尚搭配推荐
 - 食：餐馆推荐、特色菜品推荐
 - 住：建筑楼房预算估计、房价预测
 - 行：违章检测、自动驾驶

- 其他研究问题的基础

- 数据挖掘
 - 计算机视觉
 - 自然语言处理
 - 生物特征识别
 - 搜索引擎
 - 医学诊断
 - 检测信用卡欺诈
 - 证券市场分析
 - DNA序列测序
 - 语音和手写识别
 - 战略游戏
 - 机器人

本节目录

- 必要性
- **定义**
- 典型机器学习过程
- 基本术语
- 机器学习三要素

- 定义

- 五种实用的定义

最基本的机器学习是使用算法解析数据，从中学习，然后对世界上的一些事情做出决定或者是预测。

- Nvidia

机器学习是一门不需要明确编程就能让计算机运行的科学。

- 斯坦福大学

机器学习基于算法，可以从数据中进行学习而不依赖于基于规则的编程。

- 麦肯锡公司

机器学习算法可以通过例子从中挑选出执行最重要任务的方法。

- 华盛顿大学

机器学习领域旨在回答这样一个问题：我们如何建立能够根据经验自动改进的计算机系统，以及管理所有学习过程中的基本法则是什么？

- 卡内基梅隆大学

• 定义

“假设用 E 来评估计算机程序在某任务类 T 上的性能，若一个程序通过利用经验 P 在 T 中任务上获得了性能改善，则我们就说关于 T 和 E ，该程序对 P 进行了学习”

机器学习致力于研究如何通过计算的手段，利用经验来改善系统自身的性能，从而在计算机上从数据中产生“模型”，用于对新的情况给出判断

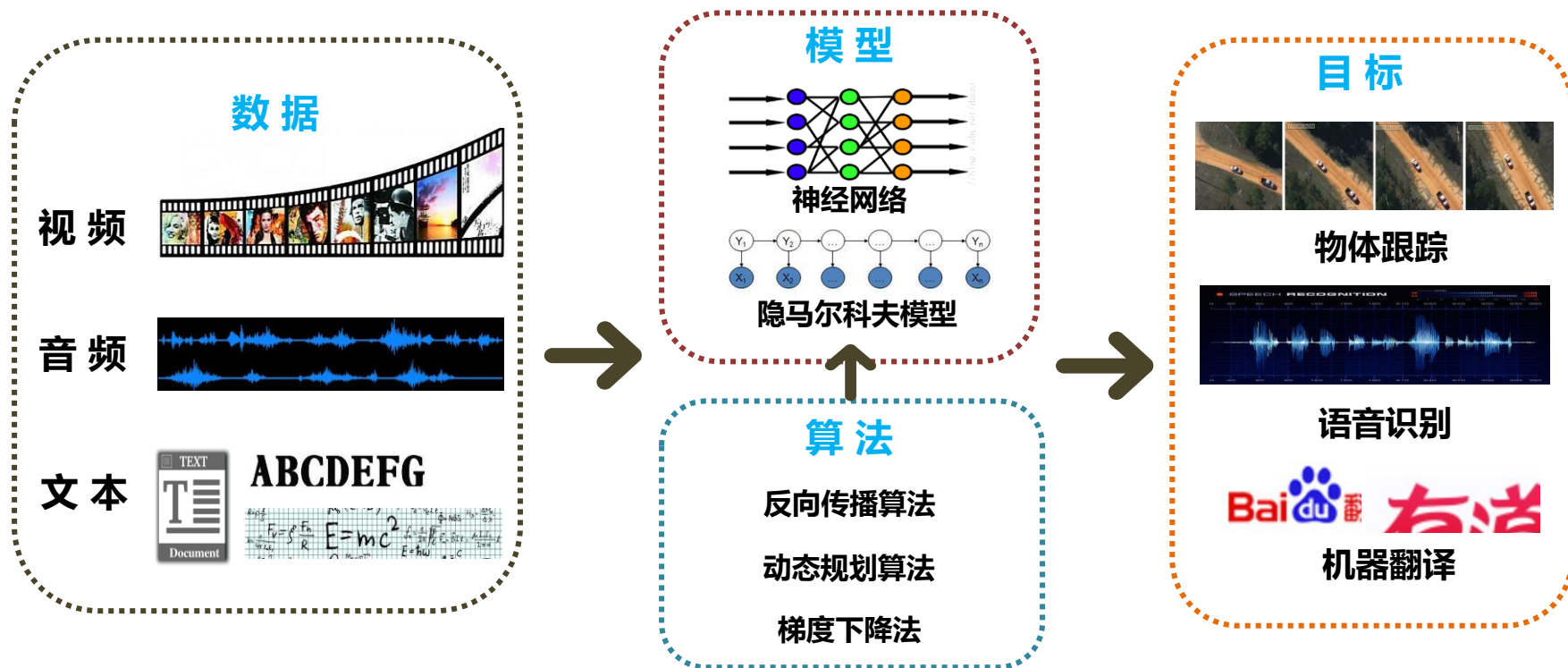
学习： 观察 → 学习 → 技能

机器学习： 数据 → 机器学习 → 提高某种性能指标

本节目录

- 必要性
- 定义
- **典型机器学习过程**
- 基本术语
- 机器学习三要素

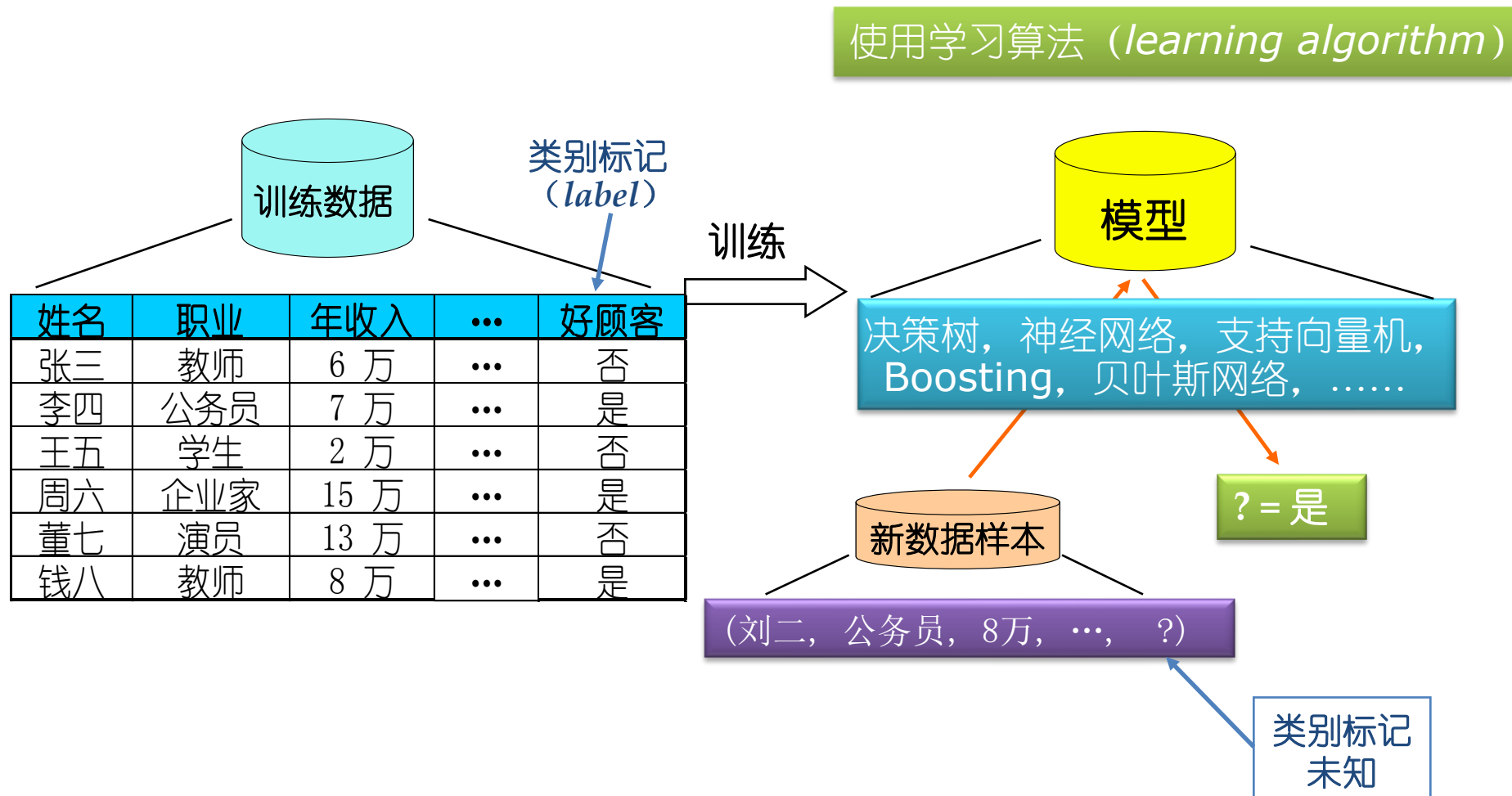
• 典型机器学习过程



学习/训练过程: $x \longrightarrow ? \longrightarrow y$

预测/测试过程: $x \longrightarrow f \longrightarrow ?$

• 典型机器学习过程



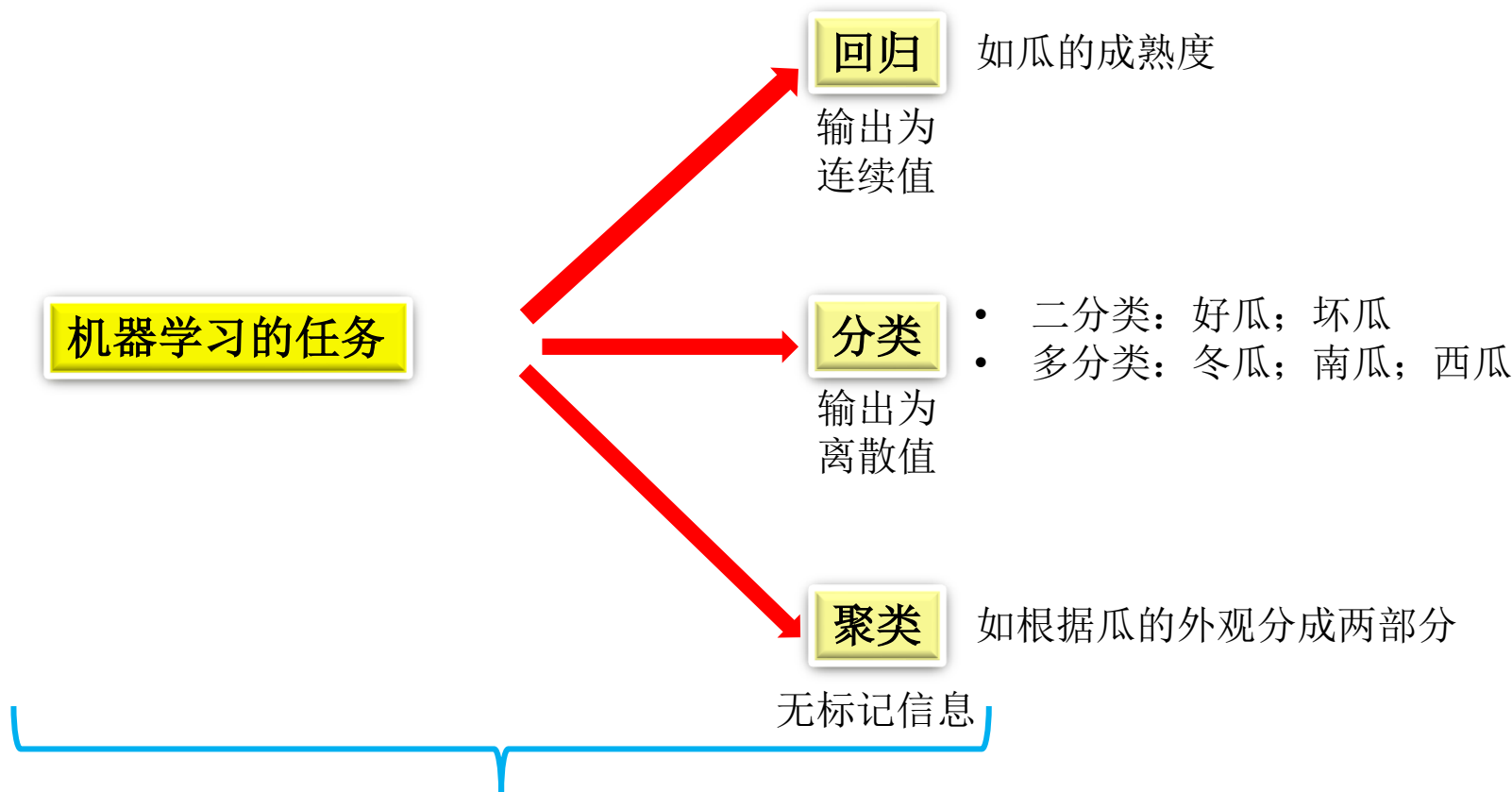
本节目录

- 必要性
- 定义
- 典型机器学习过程
- **基本术语**
- 机器学习三要素

- **基本术语**
 - 数据

		特征			标记
	编号	色泽	根蒂	敲声	好瓜
训练集	1	青绿	蜷缩	浊响	是
	2	乌黑	蜷缩	沉闷	是
	3	青绿	硬挺	清脆	否
	4	乌黑	稍蜷	沉闷	否
测试集	1	青绿	蜷缩	沉闷	?

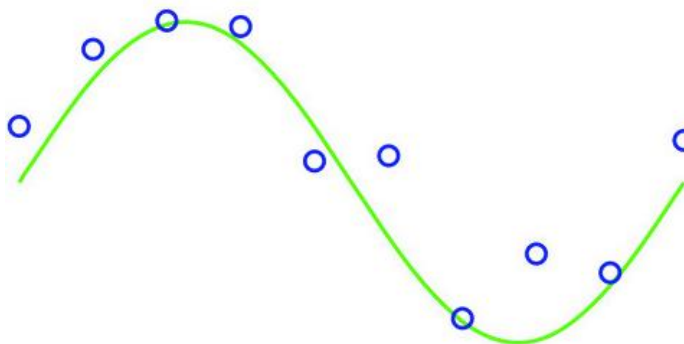
- 基本术语
 - 任务



所要解决的问题不同

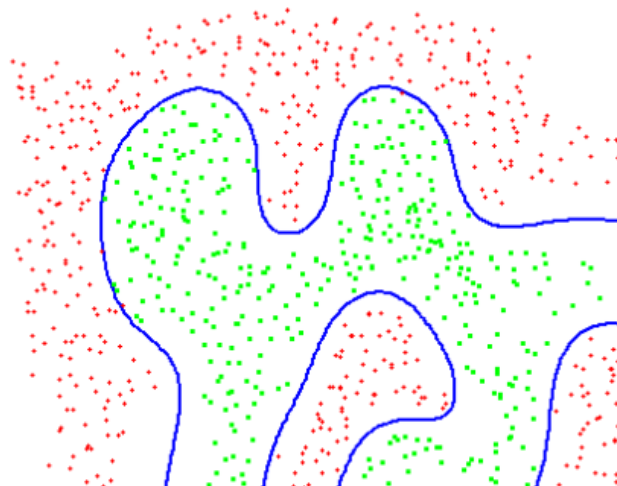
- 基本术语
 - 任务

回归任务是通过若干带有标注的样本数据构造出一个预测模型 $f(x)$ ，使得 $f(x)$ 的预测输出尽可能符合真实值，并称 $f(x)$ 为回归模型



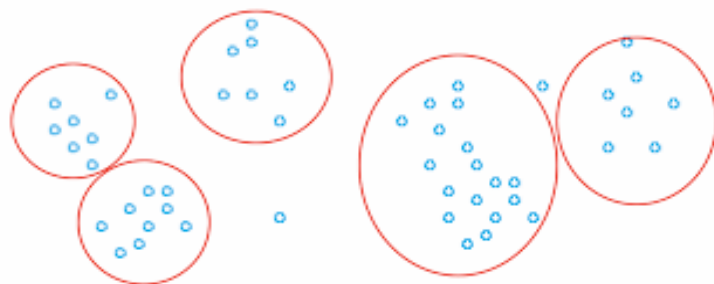
- 基本术语
 - 任务

分类任务的目标是通过训练样本构建合适的分类器 $f(x)$, 完成对目标的分类。用于分类任务的机器学习模型称为**分类模型**或**分类器**

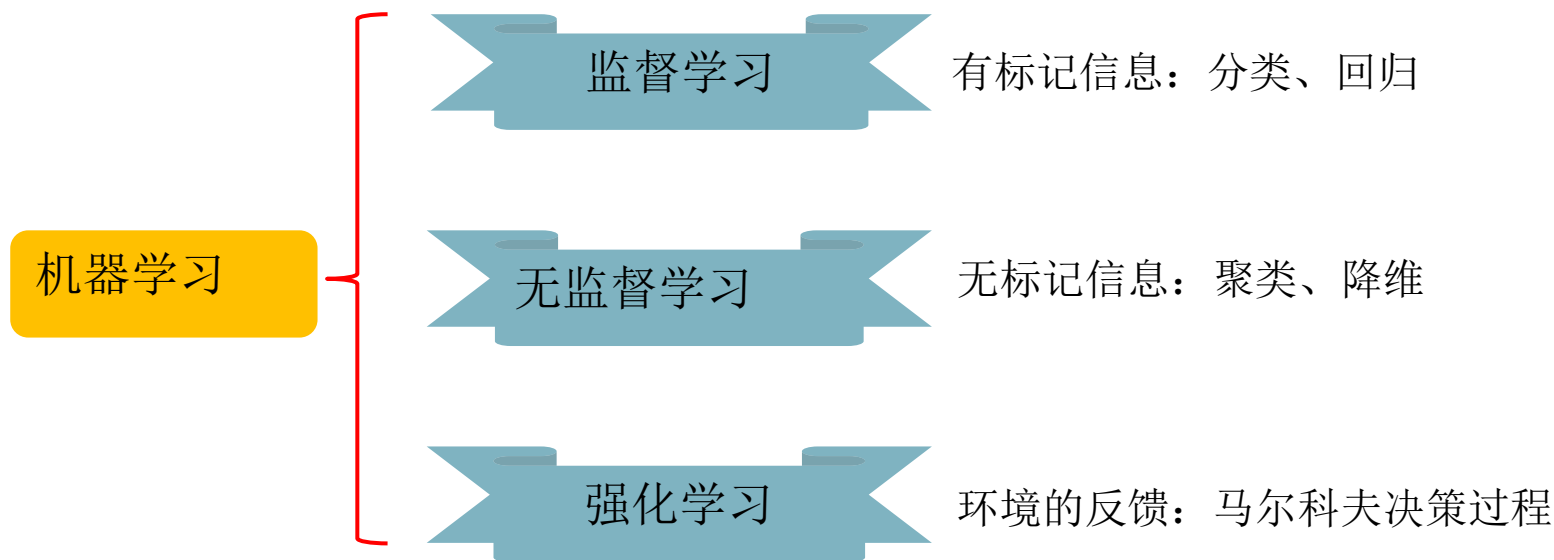


- 基本术语
 - 任务

聚类任务是对样本数据实现物以类聚的效果。聚类的类别由不同样本之间的某种相似性确定，因而聚类类别所表达的含义通常是不确定的



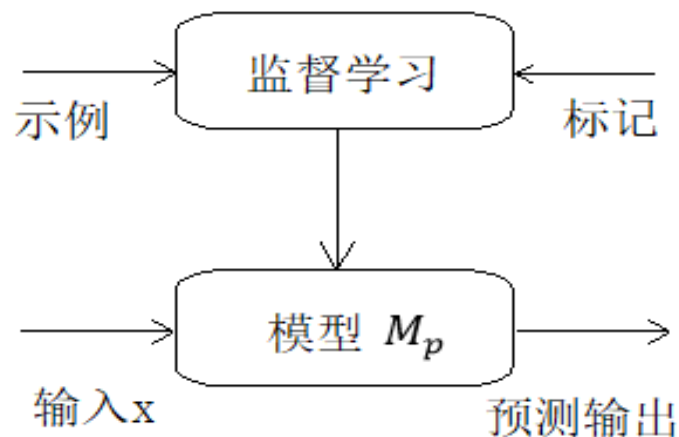
- 基本术语
 - 常见类型



依据先验信息的不同形式

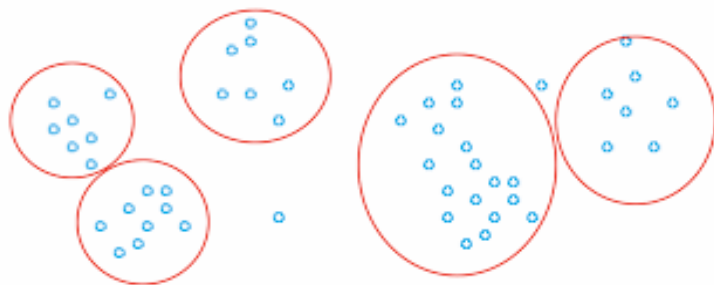
- 基本术语
 - 常见类型

有监督学习利用一组带标注样本调整模型参数，提升模型性能的学习方式。基本思想是通过标注值告诉模型在给定输入的情况下应该输出什么值，由此获得尽可能接近真实映射方式的优化模型

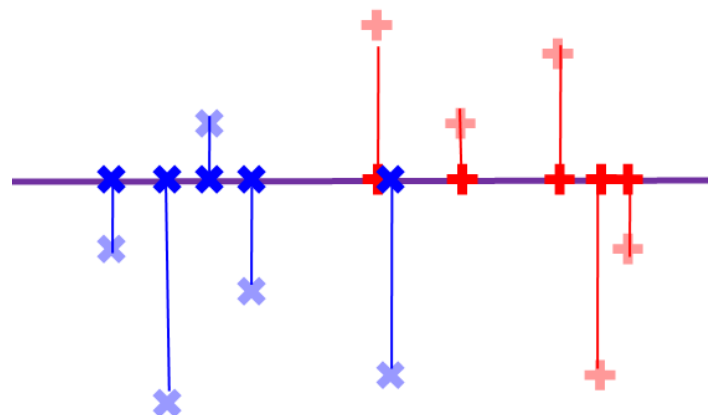


- 基本术语
 - 常见类型

无监督学习通过比较样本之间的某种联系实现对样本的数据分析。最大特点是学习算法的输入是无标记样本



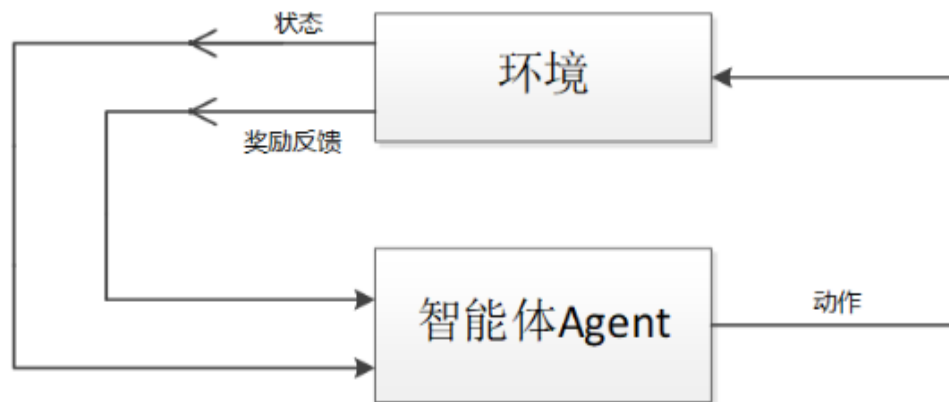
聚类



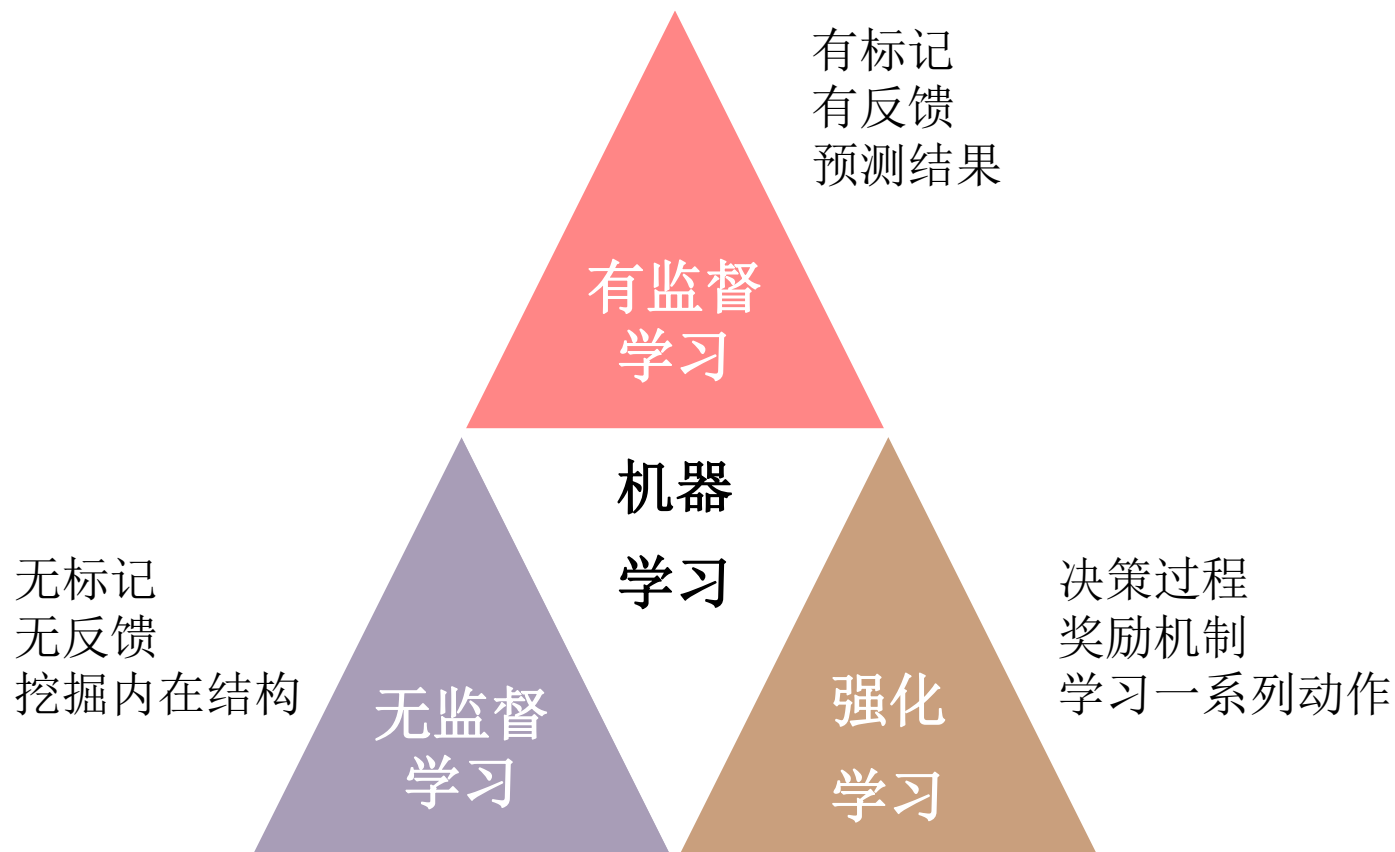
降维

- 基本术语
 - 常见类型

强化学习根据反馈信息来调整机器行为以实现自动决策的一种机器学习方式。强化学习主要由智能体和环境两个部分组成。智能体是行为的实施者，由基于环境信息的评价函数对智能体的行为做出评价，若智能体的行为正确，则由相应的回报函数给予智能体正向反馈信息以示奖励，反之则给予智能体负向反馈信息以示惩罚



- 基本术语
 - 常见类型



- 基本术语
 - 泛化能力

机器学习的目标是使得学到的模型能很好的适用于“新样本”，而不仅仅是训练集合，我们称模型适用于新样本的能力为泛化([generalization](#))能力

通常假设样本空间中的样本服从一个未知分布 \mathcal{D} ，样本从这个分布中独立获得，即“独立同分布”(i. i. d)。一般而言训练样本越多越有可能通过学习获得强泛化能力的模型

- 基本术语
 - 假设空间

对于一个具体的回归或分类任务，所有可能的模型输入数据组成的集合称为**输入空间**，所有可能的模型输出数据构成的集合称为**输出空间**

回归或分类机器学习任务的本质就是寻找一个从输入空间到输出空间的**映射**，并将该映射作为预测模型

从输入空间到输出空间的所有可能映射组成的集合称为**假设空间**

- 基本术语
 - 假设空间

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

$(\text{色泽}=\text{?}) \wedge (\text{根蒂}=\text{?}) \wedge (\text{敲声}=\text{?}) \leftrightarrow \text{好瓜}$

在模型空间中搜索不违背训练集的假设
假设空间大小： $3*4*4+1=49$

- 基本术语
 - 模型偏好

满足条件的映射通常不止一个，此时需要对多个满足条件的映射做出选择

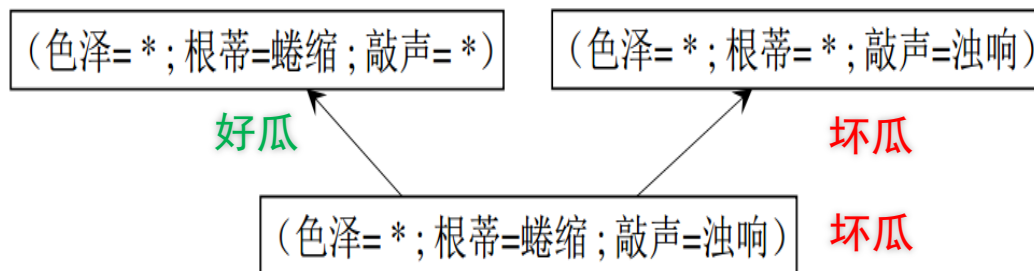
选择的主观倾向性称为机器学习算法的模型偏好

奥卡姆剃刀原则：在同等条件下选择简单事物的倾向性原则

- 基本术语
 - 模型偏好

满足条件的映射通常不止一个，此时需要对多个满足条件的映射做出选择

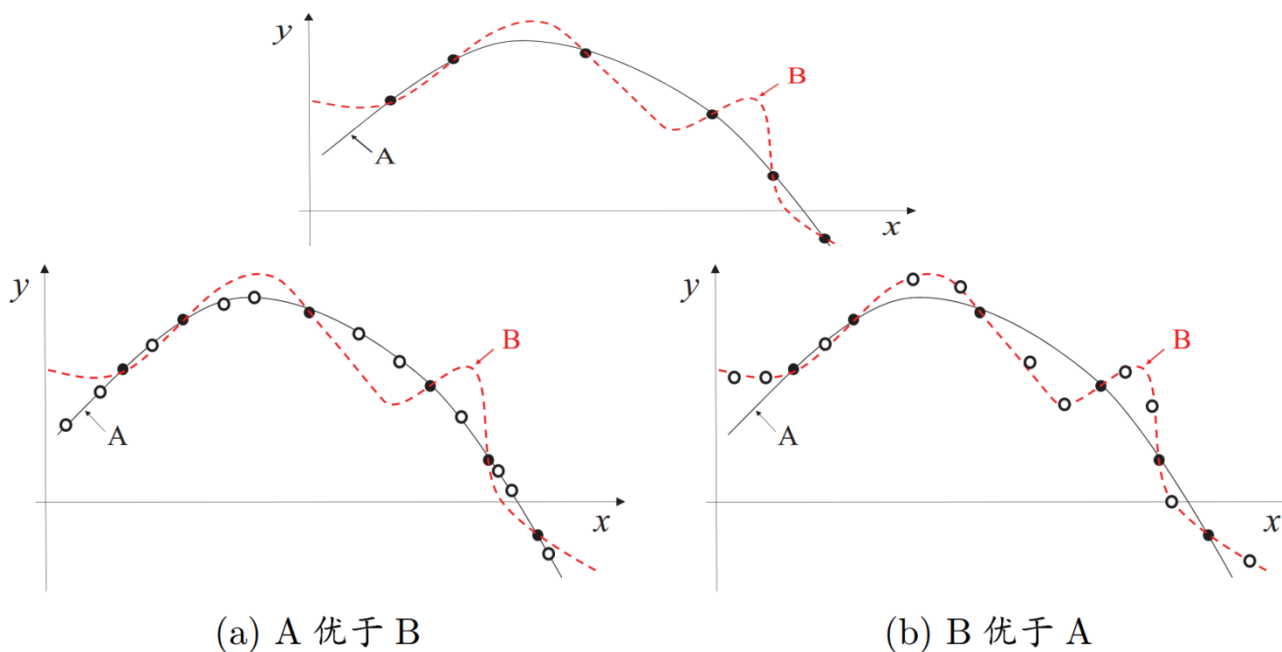
假设空间中有三个与训练集一致的假设，但他们对(色泽=青绿；根蒂=蜷缩；敲声=沉闷)的瓜会预测出不同的结果



选取哪个假设作为学习模型？

- 基本术语
 - 模型偏好

学习过程中对某种类型假设的偏好称作模型偏好



没有免费的午餐. (黑点: 训练样本; 白点: 测试样本)

- 基本术语
 - 模型偏好

模型偏好可看作学习算法自身在一个可能很庞大的假设空间中对假设进行选择的启发式或“价值观”

“奥卡姆剃刀”是一种常用的、自然科学研究中最基本的原则，即“若有多个假设与观察一致，选最简单的那个”

具体的现实问题中，学习算法本身所做的假设是否成立，也即算法的模型偏好是否与问题本身匹配，大多数时候直接决定了算法能否取得好的性能

- 基本术语
 - 模型偏好

一个算法 ξ_a 如果在某些问题上比另一个算法 ξ_b 好，必然存在另一些问题， ξ_b 比 ξ_a 好，也即没有免费的午餐定理

实际问题中，脱离具体问题，空谈“什么学习算法更好”毫无意义

- 基本术语

- 误差与损失函数

机器学习模型的输出结果与其对应的真实值之间往往会存在一定的差异，这种差异被称为模型的**输出误差**，简称为**误差**

通常需要构造**损失函数**用于度量模型对于单个样本的输出误差

对于给定的机器学习模型 f ，假设该模型对应于输入样本 x 的输出为 $\hat{y} = f(x)$ ，与 x 对应的实际真实值为 y ，则可用以 y 和 $f(x)$ 为自变量的某个函数 $L(y, f(x))$ 作为损失函数来度量模型 f 在输入样本 x 下的输出误差

例如 $L(y, f(x)) = (y - f(x))^2$ 和 $L(y, f(x)) = |y - f(x)|$

平方损失函数

绝对值损失函数

- 基本术语
 - 误差与损失函数

对于任意给定的一个 n 元样本集

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

模型 f 在 S 上的整体误差 $R_S(f)$ 定义为:

$$R_S(f) = E[L(y, f(x))] = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i))$$

即将 $R_S(f)$ 定义为 S 中所有单个样本所分别对应损失函数值的平均值

- 基本术语
 - 误差与损失函数

0-1损失函数:

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

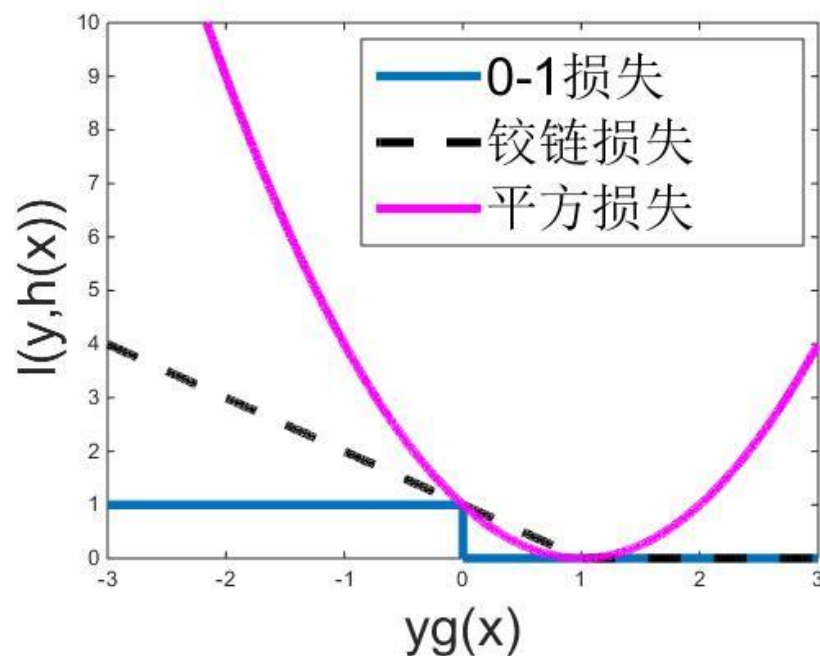
铰链损失(hinge loss):

$$L(y, f(x)) = \max(0, 1 - yf(x))$$

$$y \in \{-1, 1\}$$

平方损失(square loss):

$$L(y, f(x)) = (y - f(x))^2 = (1 - yf(x))^2$$



- 不同的任务需要不同的损失函数
- 损失函数能影响学习的好坏

- 基本术语
 - 泛化误差与经验风险

对于某个给定的机器学习任务，假设与该任务相关的所有样本构成的样本集合为 D ，则机器学习模型在样本集合 D 上的整体误差称为该模型关于该学习任务的泛化误差

$$R_{exp}(f) = E_{P(D)}[L(y, f(x))]$$

对于任意给定的 n 元训练样本集合 $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ ，假设模型 f 对输入样本 x 的预测输出为 $\hat{y} = f(x)$ ，则该模型关于训练样本集 G 的训练误差定义为：

$$R_{emp}(f) = \frac{1}{n} \sum_{k=1}^n L(y_k, f(x_k))$$

- 基本术语
 - 参数学习

根据经验风险最小化方法得到优化模型，即模型参数：

$$\hat{f} = \arg_{f \in F} \min R_{\text{emp}}(f)$$

- 基本术语
 - 测试误差

模型在测试样本集上的整体误差。对于任意给定的 v 元测试样本集合 $\{(x_1^t, y_1^t), (x_2^t, y_2^t), \dots, (x_v^t, y_v^t)\}$ ，该模型的测试误差定义为：

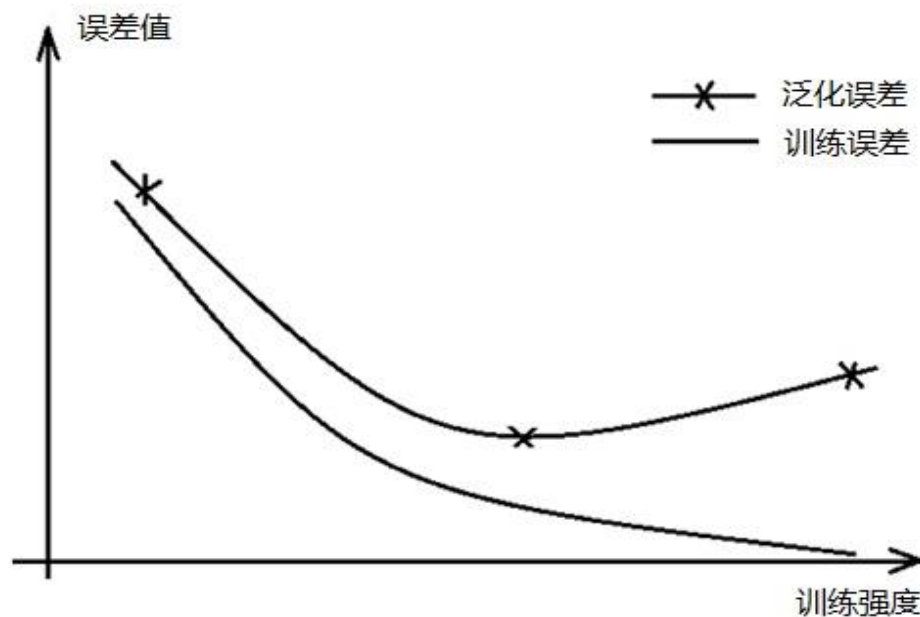
$$R_{test} = \frac{1}{v} \sum_{k=1}^v L(y_k^t, f(x_k^t))$$

- 基本术语

- 过拟合与欠拟合

过拟合是同时拟合训练样本的共性特征和个性特征（噪声）

欠拟合是未能充分拟合训练样本共性特征造成模型泛化误差较大而导致模型泛化能力较弱



- 基本术语

- 偏差与方差- “偏差-方差分解” 是解释学习算法泛化性能的重要工具

对于任意给定的一个初始模型 f ,假设 D_1, D_2, \dots, D_s 是 s 个不同的训练样本集合, 其中每个训练样本均采样自整个样本集合 D ,通过训练样本集合 D_i 训练初始模型 f_0 所得到的优化模型记为 $f_i, i \in (1, 2, \dots, s)$, $\hat{y}_i = f_i(x)$ 表示第 i 个模型对于输入样本 x 的模型输出, x 所对应的实际真实值为 y , 则这 s 个优化模型对于输入样本 x 的期望输出为:

$$E[f(x)] = \frac{1}{s} \sum_{i=1}^s f_i(x)$$

模型 $f(x)$ 在训练样本集 D_1, D_2, \dots, D_s 下所得优化模型 $f_1(x), f_2(x), \dots, f_s(x)$ 输出的方差为:

$$\begin{aligned} \text{var}[f(x)] &= E\{[f(x) - E[f(x)]]^2\} \\ &= \frac{1}{s} \sum_{i=1}^s [f_i(x) - E[f(x)]]^2 \end{aligned}$$

- 基本术语
 - 偏差与方差

模型的**学习能力**或**模型的容量**：机器学习模型这种适应训练数据变化的能力

使用模型输出在不同训练样本集合下的综合偏差对其进行度量，这种综合偏差称为模型输出的偏差，简称为**偏差**

对于模型 $f(x)$ 在训练样本集 D_1, D_2, \dots, D_s 下的优化模型 $f(x) = (f_1(x), f_2(x), \dots, f_s(x))^T$ ， $f(x)$ 作为一个离散随机变量与 x 所对应实际真实值 y 之间的偏差 $bias[f(x)]$ 为：

$$bias[f(x)] = E[f(x)] - y$$

- 基本术语
 - 偏差与方差

泛化误差：

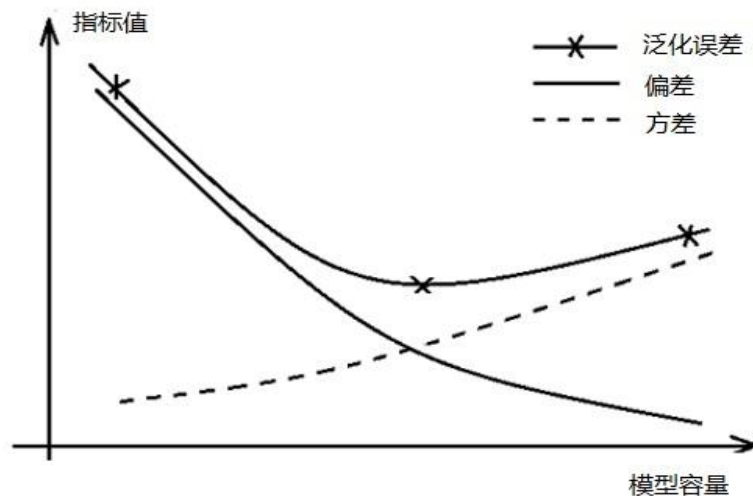
$$R_{exp}(f) = E[L(y, f(x))] = E[(f(x) - \hat{y})^2]$$

$$= \underbrace{\text{var}[f(x)]}_{\text{方差}} + \underbrace{[\text{bias}[f(x)]]^2}_{\text{偏差}} + \underbrace{\epsilon^2}_{\text{噪声}}$$

表达了同样大小训练集的变动导致的学习性能变化，刻画了数据扰动带来的影响

表达了期望预测与真实的偏离程度，刻画了学习算法对训练集的拟合能力

表达了期望泛化误差下界，刻画了学习问题的难度



本节目录

- 必要性
- 定义
- 典型机器学习过程
- 基本术语
- **机器学习三要素**

- 机器学习三要素

- 模型

- 感知机、朴素贝叶斯模型、支持向量机、决策树、随机森林...
 - 线性回归、逻辑回归、Softmax回归...
 - 神经网络...

- 学习准则

- 经验风险最小化
 - 损失函数

- 优化算法

- 梯度下降法
 - 反向传播算法
 - 动态规划算法
 - ...

- **必要性**

- 更“简单”的解决编程无法解决的复杂系统问题，应用广泛

- **定义**

- 给定任务和评价度量，程序对经验进行了学习

- **典型机器学习过程**

- 学习/训练过程与预测/测试过程

- **基本术语**

- 数据、任务、常见类型、泛化能力、假设空间、模型偏好、误差与损失函数、泛化误差与经验风险、参数学习、测试误差、过拟合与欠拟合、偏差与方差

- **机器学习三要素**

- 模型、准则、优化

思考题



1、机器学习的实质是

- A、根据现有数据,寻找输入数据和输出数据的映射关系/函数
- B、建立数据模型
- C、衡量输入数据和输出数据的映射关系/函数的好坏
- D、挑出输入数据和输出数据的最佳映射关系/函数

参考答案: A

3、以下说法正确的是

- A. 一个机器学习模型, 如果有较高准确率, 总是说明这个分类器是好的
- B. 如果增加模型复杂度, 那么模型的测试错误率总是会降低
- C. 如果增加模型复杂度, 那么模型的训练错误率总是会降低
- D. 如果降低模型复杂度, 那么模型的测试错误率总是会增加

参考答案: C

2、越复杂的模型,在训练集表现出越好的误差性能,但在测试集中并不总是表现出好的误差性能,这种现象叫?

- A、过拟合
- B、泛化性能
- C、欠拟合
- D、泛化能力

参考答案: A

4、一监狱人脸识别准入系统用来识别待进入人员的身份,此系统一共包括识别4种不同的人员:狱警、小偷、送餐员、其他,下面哪种学习方法最适合此种应用需求?

- A、二分类问题
- B、多分类问题
- C、聚类问题
- D、回归问题

参考答案: B