

监督式学习

维基百科，自由的百科全书

监督式学习（英语：Supervised learning），是机器学习的一种方法，可以由训练资料中学到或建立一个模式（函数 / learning model），并依此模式推测新的实例。^[1]训练资料是由输入物件（通常是向量）和预期输出所组成。函数的输出可以是一个连续的值（称为回归分析），或是预测一个分类标签（称作分类）。

一个监督式学习者的任务在观察完一些事先标记过的训练范例（输入和预期输出）后，去预测这个函数对任何可能出现的输入的输出。要达到此目的，学习者必须以"合理"（见归纳偏向）的方式从现有的资料中一般化到非观察到的情况。在人类和动物感知中，则通常被称为概念学习（concept learning）。

目录

回顾

经验风险最小化

主动式学习

策略和算法

应用

常见议题

参考文献

外部链接

回顾

监督式学习有两种形态的模型。最一般的，监督式学习产生一个全域模型，会将输入物件对应到预期输出。而另一种，则是将这种对应实作在一个区域模型。（如案例推论及最近邻居法）。为了解决一个给定的监督式学习的问题（手写辨识），必须考虑以下步骤：

1. 决定训练资料的范例的形态。在做其它事前，工程师应决定要使用哪种资料为范例。譬如，可能是一个手写字符，或一整个手写的辞汇，或一行手写文字。
2. 搜集训练资料。这资料须要具有真实世界的特征。所以，可以由人类专家或（机器或感测器的）测量中得到输入物件和其相对应输出。
3. 决定学习函数的输入特征的表示法。学习函数的准确度与输入的物件如何表示是有很大的关联度。传统上，输入的物件会被转成一个特征向量，包含了许多关于描述物件的特征。因为维数灾难的关系，特征的个数不宜太多，但也要足够大，才能准确的预测输出。
4. 决定要学习的函数和其对应的学习算法所使用的数据结构。譬如，工程师可能选择人工神经网络和决策树。
5. 完成设计。工程师接着在搜集到的资料上跑学习算法。可以借由将资料跑在资料的子集（称为验证集）或交叉验证（cross-validation）上来调整学习算法的参数。参数调整后，算法可以运行在不同于训练集的测试集上

另外对于监督式学习所使用的辞汇则是分类。现著有著各式的分类器，各自都有强项或弱项。分类器的表现很大程度上要跟要被分类的资料特性有关。并没有某一单一分类器可以在所有给定的问题上都表现最好，这被称为‘天下没有白吃的午餐理论’。各式的经验法则被用来比较分类器的表现及寻找会决定分类器表现的资料特性。决定适合某一问题的分类器仍旧是一项艺术，而非科学。

目前最广泛被使用的分类器有人工神经网络、支持向量机、最近邻居法、高斯混合模型、朴素贝叶斯方法、决策树和径向基函数分类。

经验风险最小化

监督式学习的目标是在给定一个 $(x, g(x))$ 的集合下，去找一个函数 g 。

假设符合 g 行为的样本集合是从某个更大甚至是无限的母体中，根据某种未知的概率分布 p ，以独立同分布随机变量方式来取样。则可以假设存在某个跟任务相关的损失函数 L

$$L: Y \times Y \rightarrow \mathbb{R}^+$$

其中， Y 是 g 的陪域，且 L 会对应到非负实数（ L 可能有其它限制）。如果预测出来 g 的值是 z ，但实际值是 y ，而 $L(z, y)$ 这个量是其间的损失。

某个函数 f 的风险是定义成损失函数的期望值。如果几率分布 p 是离散的（如果是连续的，则可采用定积分和几率密度函数），则定义如下：

$$R(f) = \sum_i L(f(x_i), g(x_i)) p(x_i)$$

现在的目标则是在一堆可能的函数中去找函数 f^* ，使其风险 $R(f^*)$ 是最小的。

然而，既然 g 的行为已知适用于此有限集合 $(x_1, y_1), \dots, (x_n, y_n)$ ，则我们可以求得出真实风险的近似值，譬如，其经验风险为：

$$\tilde{R}_n(f) = \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i)$$

选择会最小化经验风险的函数 f^* 就是一般所知的经验风险最小化原则。统计学习理论则是研究在什么条件下经验风险最小化才是可行的，且预斯其近似值将能多好？

主动式学习

一个情况是，有大量尚未标示的资料，但去标示资料则是很耗成本的。一种方法则是，学习算法会主动去向使用者或老师去询问标签。这种形态的监督式学习称为主动式学习。既然学习者可以选择例子，学习中要使用到的例子个数通常会比一般的监督式学习来得少。以这种策略则有一个风险是，算法可能会专注在于一些不重要或不合法的例子。

策略和算法

- 人工神经网络
- 案例推论
- 决策树学习
- 最近邻居法
- 支持向量机
- 随机森林
- 学习自动机

应用

- 生物资讯学

- 化学资讯学
 - 定量构效关系
- 手写辨识
- 资讯检索
- 电脑视觉中的物件识别
- 光学字元识别
- 侦测垃圾邮件
- 模式识别
- 语音识别
- 预测虚假财务报告

常见议题

- 计算学习理论
- 归纳偏向
- 过适现象
- 变形空间

参考文献

1. Stuart J. Russell, Peter Norvig (2010) *Artificial Intelligence: A Modern Approach, Third Edition*, Prentice Hall ISBN 9780136042594.

外部链接

- Matlab **SU**rrogate **MO**deling Toolbox – SUMO Toolbox (https://web.archive.org/web/20110830053048/http://sumo.intec.ugent.be/?q=SUMO_toolbox) – Matlab code for Active Learning + Model Selection + Supervised Learning (Surrogate Modeling)

取自“<https://zh.wikipedia.org/w/index.php?title=監督式學習&oldid=53185720>”

本页面最后修订于2019年2月14日 (星期四) 05:40。

本站的全部文字在知识共享 署名–相同方式共享 3.0协议之条款下提供，附加条款亦可能应用。（请参阅使用条款）
Wikipedia®和维基百科标志是维基媒体基金会的注册商标；维基™是维基媒体基金会的商标。
维基媒体基金会是按美国国内税收法501(c)(3)登记的非营利慈善机构。