

第一章 绪论

1.1 引言

1.何为机器学习？

在日常生活中经常涉及很多根据经验做出的预判。机器学习致力于研究如何通过计算的手段，利用经验改善系统自身的性能。

计算机中的经验通常以数据的形式存在，因此机器学习的主要内容就是从这些数据中产生模型。当遇到新的情况，我们可以用模型去帮助我们判断一些东西。

1.2 基本术语

1.数据集 **dataset**: 数据的集合

2.数据集中的每一条记录（也叫特征向量 **feature vector**）都有属性 **attribute**/特征 **feature**

3.特征张成的空间叫样本空间 **sample space**

4.学习 **learning**/训练 **training**: 从数据中学得模型的过程

5.训练过程中使用到的数据叫做训练数据 **training data**，他们组成训练集 **training set**

6.训练完成使用其进行预测的过程叫做测试 **test**

7.分类 **classification**: 预测离散值 回归 **regression**: 预测连续值

8.学习任务可以分为两大类：监督学习 **supervised learning** 和无监督学习 **unsupervised learning**

9.泛化 **generalization** 能力：学得模型适用于新样本的能力

1.3 归纳偏好

机器学习算法在学习过程中产生的对于某种类型假设的偏好

任何一个有效的机器学习算法必有归纳偏好

用什么样的规则去引导模型建立正确偏好？

1. 奥卡姆剃刀 **Occam's razor**: 如果有多个假设与观察一致，选择最简单的那个

举个例子。假如有一些连续点，可以用二次或更复杂的函数拟合，那么就用二次函数来拟合。问题是，怎么判断，哪一个假设更“简单”？这就要用其他机制来解决了，这个问题也一直困扰着研究者们，因此，对奥卡姆剃刀在机器学习领域的作用，一直存在争议。

2.没有免费的午餐（**No Free Lunch Theorem -NFL**）

如果简单的学习算法 **a**，它在某些问题上比算法 **b** 好，则必然存在另一些问题，**b** 比 **a** 的性能要好。

经过数学证明，这个结论对任何算法都成立。

也就是说，无论学习算法 **a** 有多聪明，**b** 有多笨拙，他们的期望性能是相同的。这就是 **NLF** 定理。

NLF 定理，让我们清楚认识到，脱离具体问题谈论什么“学习算法更好”是毫无意义的，一个算法无法在所有问题上都表现良好。