

Imitation Learning: A Survey of Learning Methods 阅读报告

韩馥光

南京大学人工智能学院

December 4, 2020

目录

引言

问题形式化

特征表示

直接模仿

间接学习

 强化学习

 优化

 迁移学习

 逆向强化学习

 主动学习

多智能体模仿学习

评估

未来方向



目录

引言

问题形式化

特征表示

直接模仿

间接学习

 强化学习

 优化

 迁移学习

 逆向强化学习

 主动学习

多智能体模仿学习

评估

未来方向



- ▶ 近年来，能够模仿人类行为的智能体的需求大大增加
 - ▶ 自动驾驶，辅助机器人和人机交互
- ▶ 复杂的应用中可能出现的场景数量太大，无法通过显式编程覆盖的情况，需要能够处理看不见的场景
 - ▶ 由专家提供的先验知识比从头开始寻找解决方案更为有效和高效
 - ▶ 为每个任务单独设计奖励是非常困难的



引言

- ▶ 一种更自然，更直观的方法是针对需要学习者模仿的所需行为进行演示
- ▶ 模仿学习的工作原理是：
 - ▶ 提取专家行为和周围环境的信息
 - ▶ 学习环境信息与演示行为之间的映射

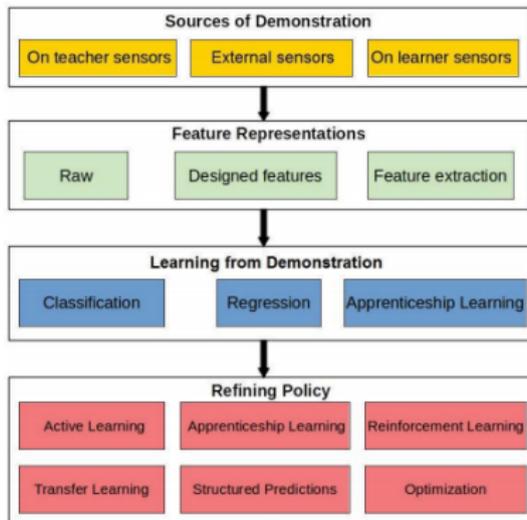


Figure: 模仿学习流程图



面临的挑战

- ▶ 传感器捕获的数据有噪声或错误
- ▶ teacher 和 learner 的匹配问题
- ▶ 可观察性问题
- ▶ 计算能力和内存的限制
- ▶ 泛化性问题

目录

引言

问题形式化

特征表示

直接模仿

间接学习

 强化学习

 优化

 迁移学习

 逆向强化学习

 主动学习

多智能体模仿学习

评估

未来方向



问题形式化

定义

- ▶ policy: 是将状态映射到行动的函数
- ▶ demonstration: (x, y) , x 是描述状态的特征向量, y 是演示者执行的动作
- ▶ experience: (s, a, r, s') , s 是状态, a 是在状态 s 采取的动作, r 是执行动作 a 的获得的奖赏, s' 是进行该操作产生的新状态
- ▶ non-stationary policy: 在学习策略参数时使用时间 t 的策略, 即该策略考虑了时间这一参数
- ▶ stationary policy: 在学习策略的过程中忽略时间这一参数的策略

除了奖赏函数外, 从演示中学习和从经验中学习的学习参数相似, 将两种方法进行结合比较常见
stationary policy 的一个优势是能够学习范围很大或未知的任务

目录

引言

问题形式化

特征表示

直接模仿

间接学习

 强化学习

 优化

 迁移学习

 逆向强化学习

 主动学习

多智能体模仿学习

评估

未来方向



挑战

- ▶ 之前介绍的匹配问题
- ▶ 演示不完整或不准确

原始特征

如果原始特征传达了足够的信息并且具有适当数量的维度，则无需进行进一步处理即可适合学习

例如一些 2D 游戏，可以直接将整个屏幕的图像作为输入而无需提取任何特征



手动设计特征

- ▶ 需要一些专业知识
- ▶ 在计算机视觉领域中很受欢迎
 - ▶ 特征从数值域到二进制图块的转换可显著改善学习



目录

引言

问题形式化

特征表示

直接模仿

间接学习

 强化学习

 优化

 迁移学习

 逆向强化学习

 主动学习

多智能体模仿学习

评估

未来方向



直接模仿

- ▶ 分类
 - ▶ 分类方法适用于学习者的行动可以被视为为一些离散类别
- ▶ 回归
 - ▶ 回归方法用于学习连续空间中的动作

直接模仿学习有一定的局限性，局限性归因于两个主要因素：演示错误、泛化性不强。



目录

引言

问题形式化

特征表示

直接模仿

间接学习

 强化学习

 优化

 迁移学习

 逆向强化学习

 主动学习

多智能体模仿学习

评估

未来方向



优化

- ▶ 定义：给定反映代理性能的成本函数 $f: A \rightarrow R$ ，其中 A 是一组输入参数， R 是一组实数，优化方法旨在查找输入参数 x_0 最小化成本函数，使得 $f(x_0) \leq f(x) \forall x \in A$
- ▶ 与强化学习类似，优化技术可用于通过从随机解开始并迭代改进以优化适应度函数来找到问题的解
- ▶ 进化算法是比较流行的优化算法，已广泛用于寻找机器人任务的运动轨迹
- ▶ 进化算法可以与模仿学习结合，以改善通过演示学习的轨迹或加快优化过程



迁移学习

- ▶ 迁移学习使用某一任务或领域的知识来加强对另一任务的学习
- ▶ 迁移学习与模仿学习的结合和机器人应用有关，因为获取样本既困难又昂贵，利用我们已经学习的一项任务的知识可能是效率高并且有效的



逆向强化学习

逆向强化学习（IRL）使用训练样本来学习奖赏函数，并使用它来改进训练后的模型



主动学习

- ▶ 主动学习能够向专家查询对给定状态的最佳响应，并使用这些主动样本来改进其策略
- ▶ 主动学习是使模型适应原始训练样本未包含的情况的有效方法
- ▶ 主动学习是使用置信度估计来识别模型中需要改进的部分，决定何时需要向专家查询

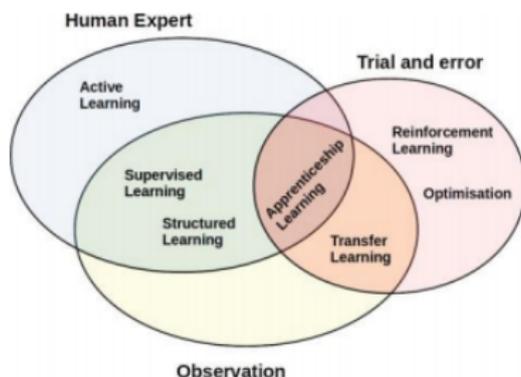


Figure: 不同来源的学习方法



目录

引言

问题形式化

特征表示

直接模仿

间接学习

 强化学习

 优化

 迁移学习

 逆向强化学习

 主动学习

多智能体模仿学习

评估

未来方向



多智能体模仿学习

尽管缺乏研究，但是模仿学习仍然很适合应用于多智能体领域：

- ▶ 在多智能体环境中，可以从演示中学习，因为可以在目标相似的智能体之间传递知识
- ▶ 智能体在一些任务中需要以人类角度来看比较切合实际的方式进行互动，在这些任务中模仿学习是有帮助的

目录

引言

问题形式化

特征表示

直接模仿

间接学习

 强化学习

 优化

 迁移学习

 逆向强化学习

 主动学习

多智能体模仿学习

评估

未来方向



目录

引言

问题形式化

特征表示

直接模仿

间接学习

 强化学习

 优化

 迁移学习

 逆向强化学习

 主动学习

多智能体模仿学习

评估

未来方向



未来方向

- ▶ General feature representations
- ▶ General task learning
- ▶ Benchmarking
- ▶ Multi-agent imitation
- ▶ Agent memory

Thank you

Thank you for listening!

