

Chapter1 Overview

智能体与环境 (Agent and Environment)

1. 核心思想与直觉 (Intuition)

核心概念：

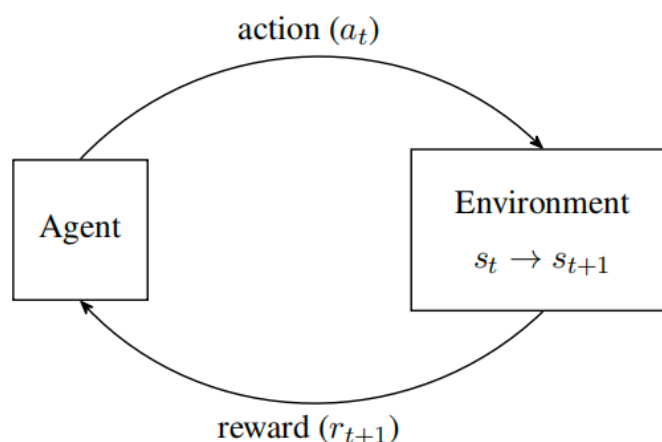
强化学习不是像监督学习那样，有一个老师告诉你“这张图是猫，那张图是狗”。RL 更像是“通过试错来学习” (Learning from interaction)。

生活案例：教小狗坐下

- 场景：你想教家里的小狗听到指令“坐下”时坐下来。
- 过程：
 1. 你发出指令“坐下”（这是**环境**给出的信号）。
 2. 小狗可能会迷茫，可能会转圈，也可能偶然坐下了（这是**智能体**尝试的动作）。
 3. 如果它转圈，你不给它吃的（**奖励**是0或负反馈）。
 4. 如果它坐下了，你立刻给它一块肉干（**奖励**是正数值）。
 5. 经过多次尝试，小狗通过**观察**（听到指令）、**行动**（坐下）和**反馈**（吃肉干）的循环，学会在听到指令时采取“坐下”这个最优动作。

在RL中，我们就是在通过数学语言描述这个循环过程。

2. 形式化定义 (Formalism)



让我们把图片中的英文定义转化为严谨的数学符号。

2.1 两个主角

- **Agent (智能体)**：决策者。它观察情况，做出选择。
 - 例子：下棋的人、扫地机器人、股票交易算法。
- **Environment (环境)**：除了智能体以外的一切。它接受动作，产生变化，并反馈结果。
 - 例子：棋盘局势、房间的布局、股市行情。

2.2 交互循环 (The Interaction Loop)

这是一个发生在离散时间步 (Discrete Time Steps) $t = 0, 1, 2, 3, \dots$ 的循环过程。

在每一个时刻 t ：

1. 观察 (Observation)：

智能体观察到当前环境的一个状态 (State)，记为 s_t 。

- $s_t \in \mathcal{S}$ ，其中 \mathcal{S} 是所有可能状态的集合。

2. 决策 (Action)：

基于状态 s_t ，智能体选择一个动作 (Action)，记为 a_t 。

- $a_t \in \mathcal{A}(s_t)$ ，其中 $\mathcal{A}(s_t)$ 是在状态 s_t 下所有可选动作的集合。

3. 演变与反馈 (Transition & Reward)：

环境接收到动作 a_t 后，发生两件事：

- **状态转移**：环境变了，进入下一个时刻的状态 s_{t+1} 。
- **奖励反馈**：环境给出一个数值反馈，称为奖励 (Reward)，记为 r_{t+1} 。
- 注意：奖励 $r_{t+1} \in \mathcal{R} \subset \mathbb{R}$ 是一个实数标量。

2.3 完整的轨迹 (Trajectory)

这样一个交互过程会产生一个序列，我们称之为**轨迹**或**历史**：

$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \dots$

符号解释与推导细节：

- 为什么奖励是 R_{t+1} 而不是 R_t ？
 - 这是一个常见的数学约定 (Convention)。
 - A_t 是在时刻 t 做出的。
 - 奖励和下一个状态是环境对 A_t 的**响应**，它们发生在 t 之后，因此在这个时间步结束、进入下一个时间步时 ($t + 1$) 才能被观测到。
 - 推导逻辑： $S_t, A_t \xrightarrow{\text{环境动力学}} S_{t+1}, R_{t+1}$

3. 关键推导：目标是什么？(The Goal)

图片中提到："The goal of the agent is to maximize the total reward throughout the entire process."

用数学公式表达，智能体的目标不是最大化当下的 r_{t+1} ，而是最大化**累积回报 (Return)**，通常记为 G_t 。

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$$

这里引出了图片中提到的两个挑战：

1. **Delayed Reward (延迟奖励)**：现在的动作 a_t 可能不会立即产生高 r_{t+1} ，但会导致未来获得巨大的 r_{t+100} 。
 - 例子：下围棋时，这一步弃子（短期损失，负奖励），是为了最终围杀对方（长期收益）。
2. **Sequential Decision Making (序列决策)**：现在的选择不仅影响当下的奖励，还会改变未来的状态 s_{t+1} ，进而限制未来可选的动作。

4. 深入剖析关键性质 (Key Properties)

4.1 状态 vs. 环境 (State represents the View)

Note 部分非常关键：

"The state is not necessarily a complete representation of the environment... states represent the agent's view of the environment."

- **完全可观测 (Fully Observable)**: 智能体能看到环境的所有细节（如：国际象棋，你看得见棋盘上所有棋子）。此时 $Observation = State$ 。
- **部分可观测 (Partially Observable)**: 智能体只能看到一部分（如：打扑克，你看不到对手的牌；或者机器人只有前置摄像头，看不到背后的障碍物）。
- **数学含义**：如果状态包含了一切用于预测未来的信息，我们称该状态具有 **马尔可夫性 (Markov Property)**。

4.2 案例解析 (Mapping Examples)

让我们用刚才定义的 S, A, R 拆解图片 Figure 1.1 下方的例子：

例子	智能体 (Agent)	环境 (Environment)	状态 (St)	动作 (At)	奖励 (R) [隐含]
移动机器人	机器人的控制软件	房间布局 + 自身电量	机器人位置、距离充电站的距离、电量百分比	移动到新房间 / 回去充电	完成清扫任务 (+)，没电关机 (-)

国际象棋	下棋的人/程序	棋盘上的棋子配置	当前棋盘局面 (Board Configuration)	移动某个棋子	赢棋(+1), 输棋(-1), 和棋(0)
股票交易	交易员/算法	股市市场 + 外部新闻 + 自身资金	股价走势、新闻情绪、账户余额	买入 / 卖出 / 持有	资产增值(+), 亏损(-)

注意：在机器人和交易员的例子中，“自身状态”（电量、资金）也被归类为环境的一部分，因为它们智能体决策时需要“观察”的外部约束。

5. 常见误区 (Common Pitfalls)

1. 误区：Agent 等于机器人身体。

- **纠正**：Agent 仅仅是那个**做决策的大脑/算法**。机器人的手臂、电机、传感器，甚至电量，通常被视为**环境**的一部分。
- **判断标准**：如果这个东西是你无法直接任意改变的（比如你不能凭空让电量变满，只能通过动作去充电），那它就是环境。

2. 误区：奖励是Agent自己定义的。

- **纠正**：奖励是**环境**给的。Agent 无法修改奖励函数，它只能想办法去获得更高的奖励。这就像如果你玩游戏，你不能自己修改代码让自己得分，你只能通过玩得更好来得分。

这部分内容深入到了强化学习系统的**内部构造**（Elements）以及它面临的**核心难题**（Challenges）。

如果说上一节讲的是“智能体”和“环境”这两个角色的**外部交互**，那么这一节就是打开“智能体”的大脑，看看里面装了什么，以及为什么它这么难学。

以下是详细的拆解与笔记整理：

1. 核心要素：智能体的大脑里有什么？

一个完整的强化学习系统包含四个主要要素：**策略 (Policy)**、**奖励信号 (Reward Signal)**、**价值函数 (Value Function)** 和 **环境模型 (Model of the Environment)**（可选）。

我们用一个**“走迷宫寻找宝藏”**的例子来直观理解这四个概念。

1.1 策略 (Policy, π) —— 行动指南

- **直觉**：这就是智能体的“锦囊妙计”或者“条件反射”。在迷宫的某个路口（状态），策略会告诉它：“往左走”或者“有80%的概率往左走”。
- **数学定义**：
 - 它是从**状态空间 (S)** 到 **动作空间 (A)** 的映射。
 - **确定性策略 (Deterministic)**: $a = \pi(s)$ 。在状态 s 下，一定做动作 a 。

- **随机策略 (Stochastic):** $\pi(a|s) = P(A_t = a|S_t = s)$ 。在状态 s 下，以一定概率选择动作 a 。

1.2 奖励信号 (Reward Signal, R) —— 即时反馈

- **直觉：**这是环境给的“糖果”或“电击”。在迷宫里，撞墙了扣1分，捡到金币加10分。它定义了任务的目标。
- **性质：**
 - 它是**即时 (Immediate)** 的。
 - 它是一个随机函数，取决于当前状态和采取的动作。
 - **目标：**最大化长期的总奖励，而不是单次的奖励。

1.3 价值函数 (Value Function, V) —— 长期眼光 (核心概念)

这是RL中最关键、也是初学者最容易混淆的概念。

- **直觉：**
 - **Reward** 是“我现在爽不爽”。
 - **Value** 是“我现在的处境好不好（未来能有多爽）”。
 - **例子：**在迷宫里，你走到了一个死胡同的尽头，虽然这里有一枚金币（**Reward高**），但拿完你就被困住了（**Value低**）；反之，你站在宝藏房间的门口，虽然现在手里没金币（**Reward低**），但你只要推门就能赢（**Value高**）。
- **数学定义：**
 - $V(s)$ 是一个实数值函数 $V : S \rightarrow \mathbb{R}$ 。
 - 它表示从状态 s 开始，未来能累积到的**总预期奖励 (Expected Total Reward)**。
- **关键关系：**

$Value \approx Immediate\ Reward + Value\ of\ Next\ State$

 - 根据图片解释，只要有了准确的价值函数，决策就变得很简单：选择那个能让你进入“最高价值下一状态”的动作。

1.4 环境模型 (Model) —— 想象力 (可选)

- **直觉：**这是智能体脑补的世界地图。如果它知道“在这个路口往东走会遇到陷阱”，它就不需要真的去踩陷阱（**规划/Planning**）。
- **分类：**
 - **Model-based (有模型)：**先在脑子里模拟推演，再行动。
 - **Model-free (无模型)：**像莽夫一样，在真实的试错中学习，不预测环境变化。

2. 核心挑战：为什么RL很难？

强化学习有两个监督学习（Supervised Learning）没有的独特难点。

2.1 探索与利用的困境 (Exploration-Exploitation Dilemma)

- **直觉**：你去一家餐厅吃饭。
 - **利用 (Exploit)**：点你以前吃过的最好吃的菜。这能保证你得到不错的体验（Reward），但你永远发现不了也许菜单背面有更好吃的菜。
 - **探索 (Explore)**：点一道从未尝试过的菜。这可能很难吃（低Reward），也可能那是你的新宠（发现更好的Action）。
- **定义**：
 - 为了获得高回报，必须**利用**已知的好动作。
 - 为了发现好动作，必须**探索**未知的动作。
 - **困境**：你不能同时既做利用又做探索，必须在两者间权衡。
- **注意**：这在监督学习里不存在，因为监督学习是训练时给答案，测试时照做，不需要去“探索”更好的答案。

2.2 延迟奖励 (Delayed Reward)

- **直觉**：下围棋。
 - 你在第10步下的一手好棋，可能当时看起来平平无奇（即时Reward为0），但它为你第50步的“绝杀”奠定了基础。
 - 当第50步赢了的时候，系统很难判断：是因为第49步下得好，还是第10步下得好？（这被称为**信用分配问题 Credit Assignment Problem**）。
 - **挑战**：智能体必须有**远见 (Foresight)**，不能只盯着眼前的利益。
-

1. 监督学习 (Supervised Learning) vs. 强化学习 (Reinforcement Learning)

通俗解释：老师教 vs. 自己悟

- **监督学习 (SL)** 就像是“在学校上课”：
 - **场景**：老师拿出一张卡片（输入），问你“这是什么？”你回答“猫”。老师立刻告诉你：“对，是猫”或者“错，是狗”（标签/Label）。

- **关键：有一个全知的“老师”**（数据集）直接告诉你每一步的**正确答案**是什么。你只需要照着学，目的是“模仿老师的答案”。
- **强化学习 (RL) 就像是“学骑自行车”：**
 - **场景：**没有老师告诉你这一秒你的左腿肌肉要收缩 30%，下一秒车把手要向右偏 5 度。你只能自己骑上去试。
 - **过程：**你歪了摔倒了（得到**负奖励**），你会痛，大脑记住了刚才那样做不行；你骑出去了 10 米（得到**正奖励**），你会爽，大脑记住了刚才那样做是对的。
 - **关键：没有正确答案。**你只能通过**试错 (Trial and Error)**，根据环境给你的反馈（痛或爽）来调整策略。

核心区别表

维度	监督学习 (SL)	强化学习 (RL)
数据来源	静态的历史数据 (Labeled Data)	动态的交互经验 (Interaction)
指导信号	告诉你是对是错，并给出 正确答案 (Correct Answer)	只告诉你得分多少 (Reward)，不告诉怎么做才更好
时间维度	每一个样本通常是独立的 (i.i.d.)	每一个决策都会改变未来 (Sequential)
核心难题	过拟合、泛化能力	探索与利用 (Exploration-Exploitation)

2. 强化学习解决了监督学习无法解决的什么问题？

监督学习很强大，但它有两个致命弱点，这正是 RL 存在的意义：

1. 没有“正确答案”的问题 (No Ground Truth)

- **例子：**在这个复杂的围棋局势下，这一步下哪里才是“绝对正确”的？
- 没有人知道。即使是人类最强棋手也不知道。监督学习无法训练，因为它需要标签。但 RL 可以，因为 RL 不需要知道哪一步最好，它只需要知道**最后赢没赢**。

2. 序列决策与长期后果 (Sequential Decision Making & Delayed Reward)

- **问题：**监督学习只看眼前。
- **RL 的特长：**RL 能处理**“现在的牺牲是为了未来的收益”**。
- **例子：**在股票交易或游戏中，有时候你必须先亏一点钱（监督学习会认为这是错的），才能在后面赚大钱。RL 通过**价值函数 (\$V\$)**能够看穿这一点，而监督学习通常做不到。

3. 马尔可夫性 (Markov Property)：完全可观测 vs. 部分可观测

这是一个非常容易混淆的概念。

直观定义：“历史无关性”

- 一句话定义：未来只取决于现在，与过去无关。
- 通俗例子：
 - **非马尔可夫**：你感冒了。医生问：“你昨天淋雨了吗？前天熬夜了吗？”（因为仅仅看你现在的体温，不足以判断病情，需要追溯历史）。
 - **马尔可夫**：一颗飞在空中的球。我们要预测它下一秒在哪里，只需要知道它**现在的【位置】和【速度】**。至于它是被乔丹扔出来的，还是被机器发射出来的（历史），完全不重要。只要现在的状态已知，历史就是多余的。

你的疑问：是完全可观测还是部分可观测？

答案是：马尔可夫性是对**“状态 (State)”** 的一种要求，而不是对环境的要求。

1. 完全可观测 (Fully Observable) \rightarrow 满足马尔可夫性

- 如果你能看到环境里的一切（如下围棋，棋盘上所有棋子都在你眼前），那么你看到的这个画面 (Observation) 就是**状态 (State)**。这个状态本身就包含了所有信息，所以它**满足**马尔可夫性。这被称为 **MDP (马尔可夫决策过程)**。

2. 部分可观测 (Partially Observable) \rightarrow 观测不满足马尔可夫性

- 如果你在玩FPS游戏（第一人称射击），你只能看到屏幕前的画面，看不到背后的敌人。
- 此时，你的**“观测 (Observation)”** 是**不满足**马尔可夫性的。因为仅仅凭现在的画面，你不知道背后有没有人追你，你必须记得“刚才听到了背后的脚步声”（依赖历史）。
- 这种情况被称为 **POMDP (部分可观测马尔可夫决策过程)**。
- **怎么解决？** 智能体必须自己在脑子里构建一个“信念状态 (Belief State)”，把历史记忆压缩进去，强行让它变得具有马尔可夫性。

总结：图片中的 Note 提到 “*state is not necessarily a complete representation... it only captures aspects... observable*”。这意味着在现实中，Agent 拿到的状态往往是不完美的（部分可观测的），但为了用数学解决它，我们通常**假设**它具有马尔可夫性，或者努力构造出一个满足马尔可夫性的状态。

4. 奖励函数 (Reward Function) 一般是什么？

奖励函数是 RL 系统的**指挥棒**。它定义了任务的**目标**，但不告诉智能体如何达成目标。

它是主观设定的

奖励是你（设计者）根据任务需求定义的**标量数值 (Scalar)**。你鼓励什么，就给正分；惩罚什么，就给负分。

常见的奖励设计模式：

1. Win/Loss 模式（最稀疏）

- **围棋/象棋**：赢了 +1，输了 -1，中间几百步全是 0。

- 难点：反馈太慢 (Delayed Reward)，智能体很难学。

2. Explicit Progress 模式 (稠密)

- 走迷宫：每走一步 -1 (为了鼓励它快点走出去)；找到出口 +100；撞墙 -10。
- 机器人走路：每前进一米 +1；摔倒 -100。

3. Human Feedback (RLHF) (ChatGPT 的核心)

- 对话系统：由于“好回答”很难用公式写出来，所以奖励是由**人类打分**的（或者由一个模仿人类打分的模型给出）。

关键点

奖励函数非常敏感。如果你设计得不好，智能体会**“钻空子” (Reward Hacking)**。

- 笑话例子：你让扫地机器人“把灰尘扫得越少越好”。结果机器人感应到灰尘后，选择原地关机（这样它眼里的灰尘就没有了），而不是去扫地。所以奖励必须精确对应你真正的意图。
-