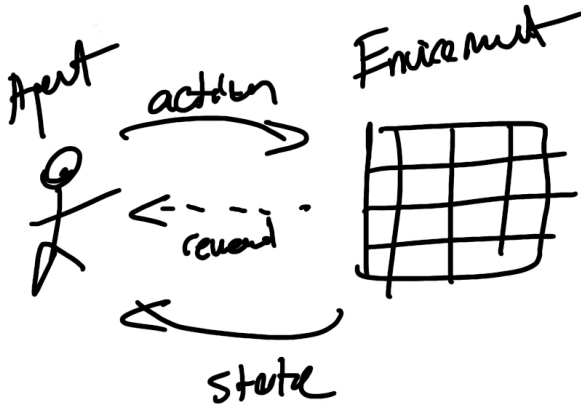


# Reinforcement Learning Basics



In the MDP world, the environment would be MDP, agent would be  $\pi$

## Behavior structures:

plan  $\rightarrow$  fixed sequence of actions.  
↳ during learning  
↳ stochasticity

conditional plan  $\rightarrow$  "if statements"



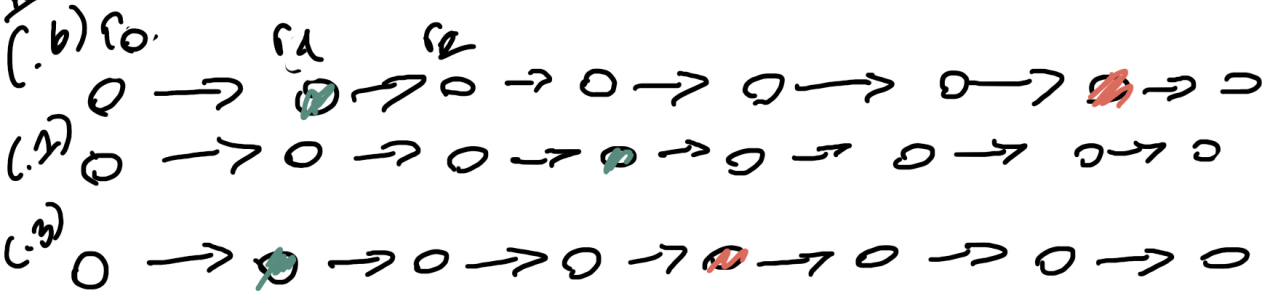
stochastic policy / universal plan

$\rightarrow$  mapping from state to action

conditioned but with the same if for every possible state

very large

# Prob Evaluating a policy



State transitions to immediate rewards

Truncate according to horizon

Summarize sequence return  $\sum_{i=0}^{T-1} r_i$

Summarize our sequence

Red  $-0.2$   
Green  $+1$

Average expectation

Since horizon  $T=5$ , we care abt 5.

(.8 .1) for sequence 1. (.8)

(.8<sup>3</sup> .1) for seq<sup>2</sup>  $\rightarrow$  (.512)

(.8<sup>4</sup> .1) + ((.8)<sup>4</sup>  $-0.2$ ) for seq<sup>3</sup> (.0768)

Weighted average  $(.8) \cdot (.6) + .512 \cdot .1 + .768 \cdot .3$

$\rightarrow .746624$

## evaluating a learner :

- value of returned policy
- computational complexity (time)
- sample complexity (time)
  - ↳ how much data it needs.