

# Reinforcement Learning

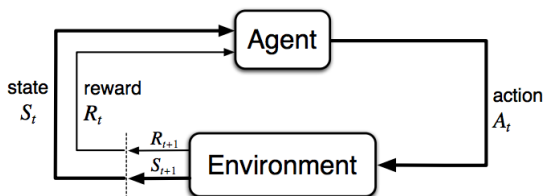
Jack Lanchantin

November 30, 2015

- 1 Reinforcement Learning
- 2 Markov Decision Processes
- 3 Second Section

- **Framing of the problem of learning from interaction to achieve a goal.**
- **Agent:** learner and decision maker
- **Environment:** what the learner interacts with (everything outside the agent)
- Agent selects actions and the environment responds to those actions and presents new situations

# Reinforcement Learning



- At each time step  $t$ , the agent receives the environment **state**  $S_t \in \mathcal{S}$ , and the agent then selects an **action**  $A_t \in \mathcal{A}(S_t)$ 
  - $\mathcal{S}$  is the set of possible states (whatever information is available to the agent).
  - $\mathcal{A}(S_t)$  is set of actions available in state  $S_t$
- One time step later, the agent receives a **reward**,  $R_{t+1} \in \mathcal{R} \subset \mathbb{R}$ , and ends up in a new state  $S_{t+1}$

- At each time step, the agent implements a mapping  $\pi_t$  from states to probabilities of selecting each possible action, where  $\pi_t$  is called a **policy**
  - $\pi_t(a|s)$  = probability that  $A_t = a$  if  $S_t = s$

## Reinforcement Learning Objective

The agent's goal is to maximize the total amount of reward it receives over the long run by changing its policy as a result of its experience

- Let the sequence of rewards after time step  $t$  is  $R_{t+1}, R_{t+2}, R_{t+3}, \dots$ , then we want to maximize the return  $G_t$
- The agent chooses  $A_t$  to maximize the discounted return:

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

where  $\gamma$  is the discount rate and  $0 \leq \gamma \leq 1$

- The closer  $\gamma$  is to 1, the more the agent accounts for future rewards

# Markov Property

- Probability of transitioning to new state  $s$ :

$$p(s', r|s, a) = Pr\{S_{t+1} = s', R_{t+1} = r | S_0, A_0, R_1, \dots, S_{t-1}, A_{t-1}, R_t, S_t, A_t\}$$

- Probability with Markov assumption

$$p(s', r|s, a) = Pr\{S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a\}$$

- Markov assumption allows us to predict the next state and expected rewards from knowledge of only the current state
  - Assume that the current state tells us everything we need to know for future (e.g. current state of checker board)

## Heading

- 1 Statement
- 2 Explanation
- 3 Example

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Integer lectus nisl, ultricies in feugiat rutrum, porttitor sit amet augue. Aliquam ut tortor mauris. Sed volutpat ante purus, quis accumsan dolor.



# Table

<b>Treatments</b>	<b>Response 1</b>	<b>Response 2</b>
Treatment 1	0.0003262	0.562
Treatment 2	0.0015681	0.910
Treatment 3	0.0009271	0.296

Table: Table caption

Richard S. Sutton and Andrew G. Barto. Reinforcement Learning: An Introduction. MIT Press 2015.