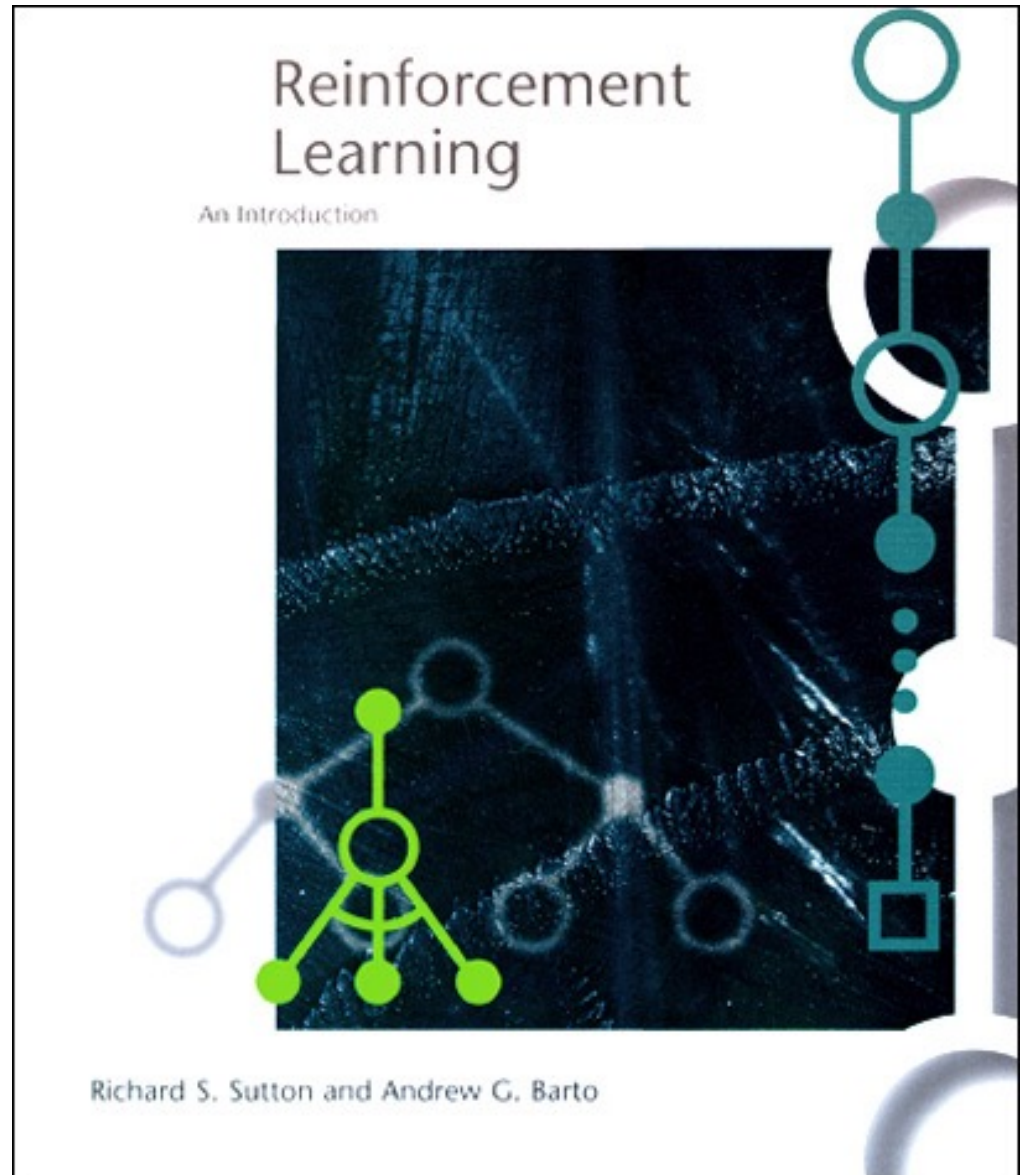


CMSC 478: Reinforcement Learning

There's an entire book!

<http://incompleteideas.net/book/the-book-2nd.html>



The Big Idea

- “Planning”: Find a sequence of steps to accomplish a goal.
 - Given start state, transition model, goal functions...
- This is a kind of **sequential decision making**.
 - Transitions are deterministic.
- What if they are stochastic (probabilistic)?
 - One time in ten, you drop your sock
- Probabilistic Planning: Make a plan that accounts for probability by **carrying it through the plan**.

Review: Formalizing Agents

- Given:
 - A state space S
 - A set of actions a_1, \dots, a_k including their results
 - Reward value at the end of each trial (series of action) (may be positive or negative)
- Output:
 - A **mapping from states to actions**

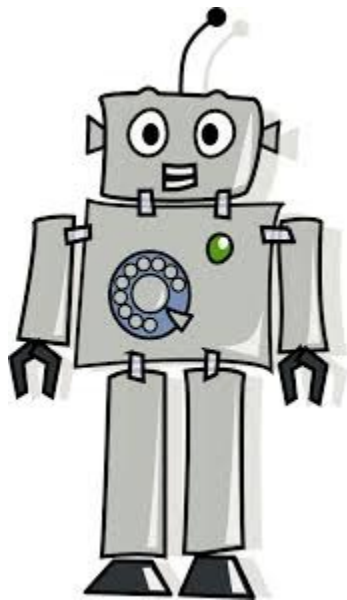
Review: Formalizing Agents

- Given:
 - A state space S
 - A set of actions a_1, \dots, a_k including their results
 - Reward value at the end of each trial (series of action) (may be positive or negative)
- Output:
 - A **mapping from states to actions**
 - Which is a **policy**, π

Reinforcement Learning

- We often have an agent which has a **task** to perform
 - It takes some actions in the world
 - At some later point, gets feedback on how well it did
 - The agent performs the same task repeatedly
- This problem is called **reinforcement learning**:
 - The agent gets positive reinforcement for tasks done well
 - And gets negative reinforcement for tasks done poorly
 - Must somehow figure out which actions to take next time

Reinforcement Learning



agent

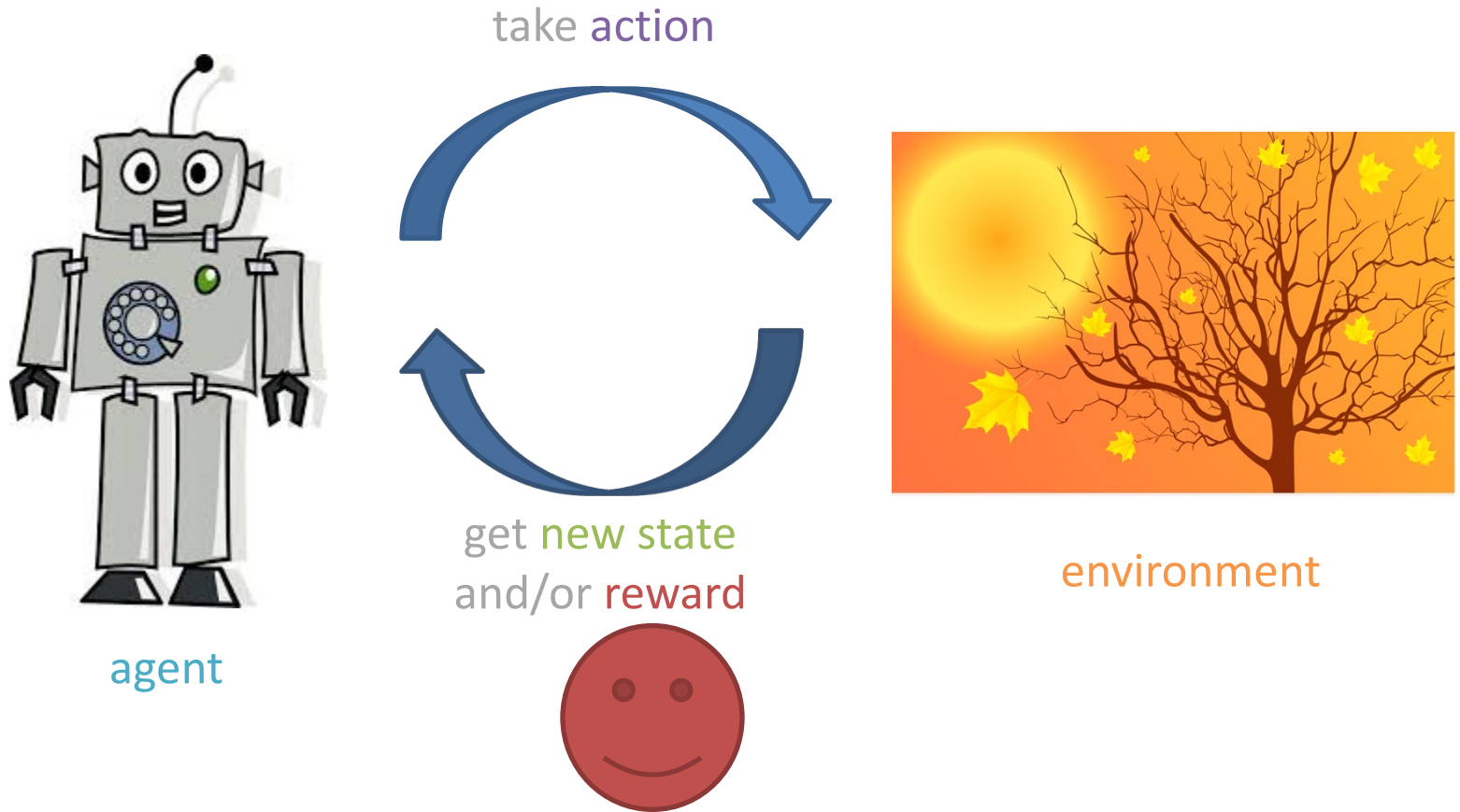


environment

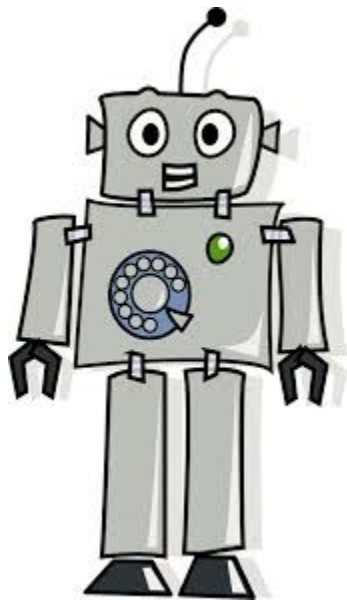
Reinforcement Learning



Reinforcement Learning

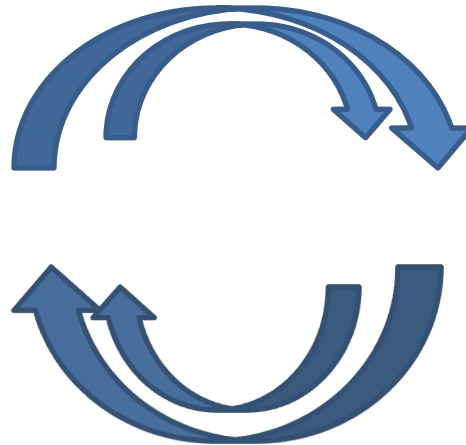


Reinforcement Learning



agent

take action



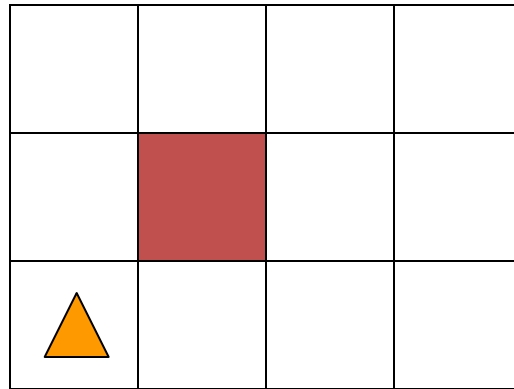
get new state
and/or reward



environment

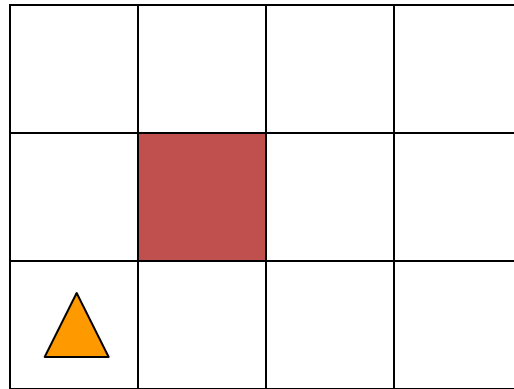


Simple Robot Navigation Problem



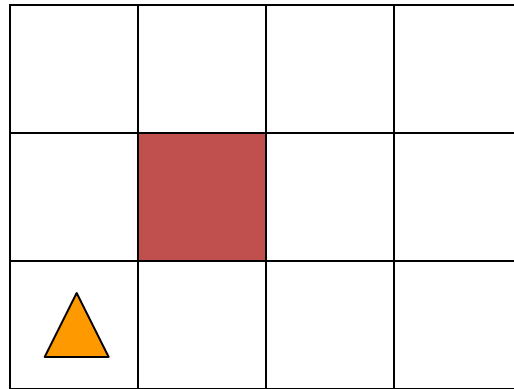
- In each state, the possible actions are **U**, **D**, **R**, and **L**

Probabilistic Transition Model



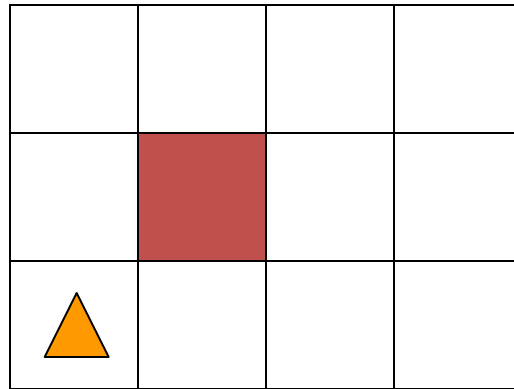
- In each state, the possible actions are **U**, **D**, **R**, and **L**
- The effect of **U** is as follows (**transition model**):
 - With probability 0.8, the robot moves up one square (if the robot is already in the top row, then it does not move)

Probabilistic Transition Model



- In each state, the possible actions are **U**, **D**, **R**, and **L**
- The effect of **U** is as follows (**transition model**):
 - With probability 0.8, the robot moves up one square (if the robot is already in the top row, then it does not move)
 - With probability 0.1, the robot moves right one square (if the robot is already in the rightmost row, then it does not move)

Probabilistic Transition Model



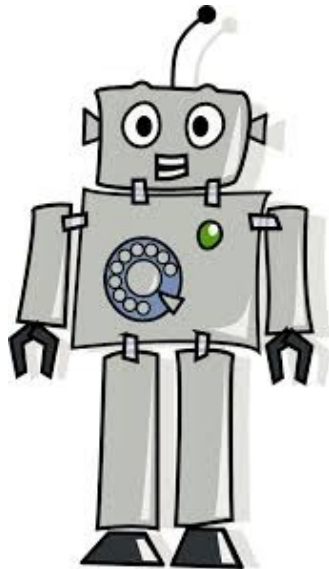
- In each state, the possible actions are **U**, **D**, **R**, and **L**
- The effect of **U** is as follows (**transition model**):
 - With probability 0.8, the robot moves up one square (if the robot is already in the top row, then it does not move)
 - With probability 0.1, the robot moves right one square (if the robot is already in the rightmost row, then it does not move)
 - With probability 0.1, the robot moves left one square (if the robot is already in the leftmost row, then it does not move)

Markov Property

The transition properties depend only on the current state, not on the previous history (how that state was reached)

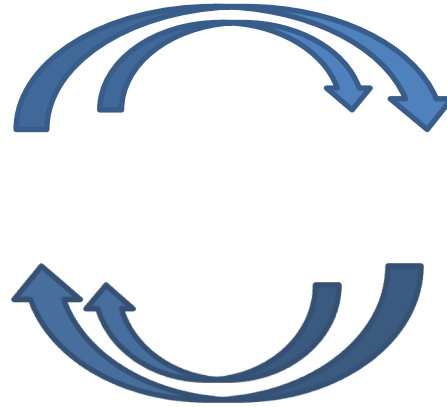
Markov assumption generally: current state only ever depends on previous state (or finite set of previous states).

Markov Decision Process: Formalizing Reinforcement Learning



agent

take action



get new state
and/or reward

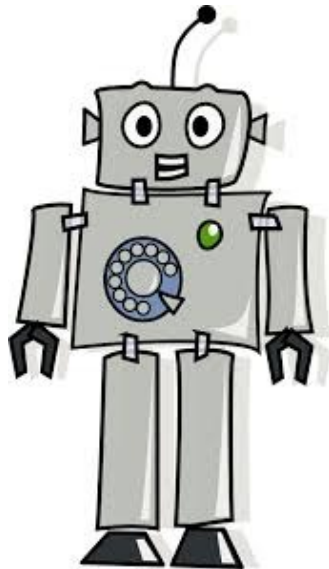


environment

Markov Decision
Process:

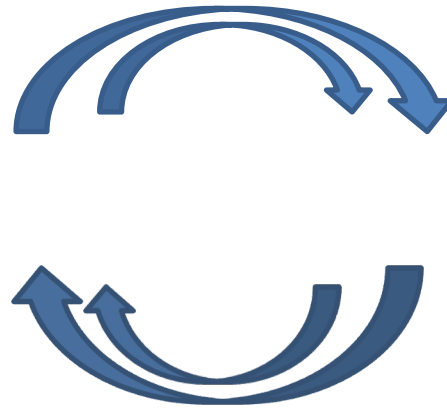
$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$

Markov Decision Process: Formalizing Reinforcement Learning



agent

take action



get new state
and/or reward



environment

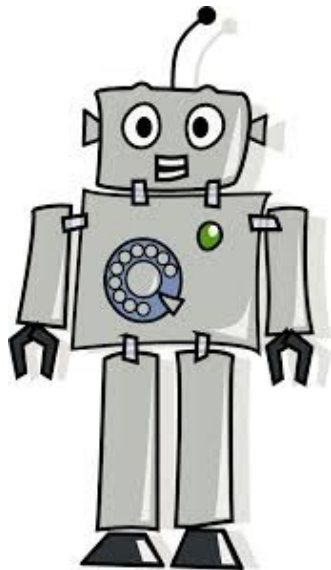
Markov Decision
Process:

set of
possible
actions

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

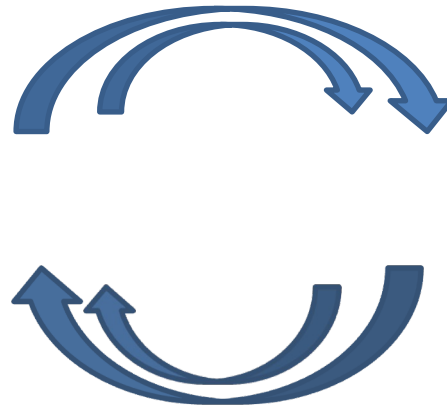
set of
possible
states

Markov Decision Process: Formalizing Reinforcement Learning



agent

take action



get new state
and/or reward



environment

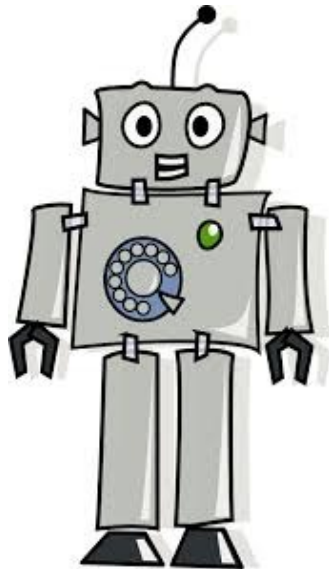
Markov Decision
Process:

set of possible actions

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

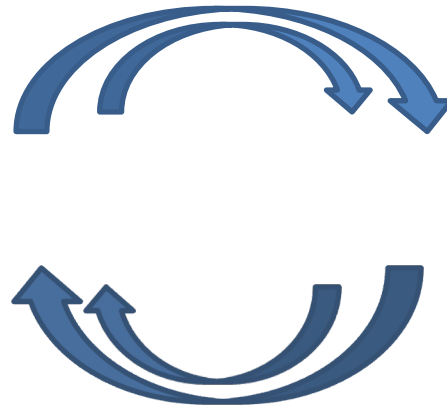
set of possible states reward of (state, action) pairs

Markov Decision Process: Formalizing Reinforcement Learning



agent

take action



get new state
and/or reward



environment

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

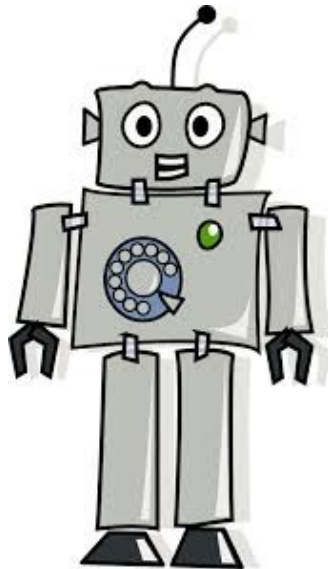
set of possible states

set of possible actions

state-action transition distribution

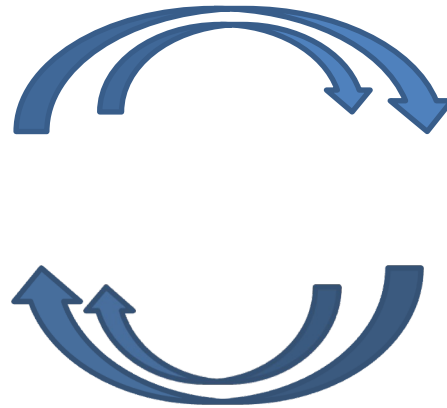
reward of (state, action) pairs

Markov Decision Process: Formalizing Reinforcement Learning



agent

take action



get new state
and/or reward



environment

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

set of possible states set of possible actions state-action transition distribution
 reward of (state, action) pairs discount factor

Robot in a room

			+1
			-1
START			

actions: UP, DOWN, LEFT, RIGHT

UP

80%

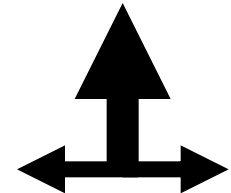
move UP

10%

move LEFT

10%

move RIGHT



reward +1 at [4,3], -1 at [4,2]
reward -0.04 for each step

Goal: what's the strategy to achieve the maximum reward?

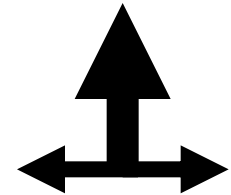
Robot in a room

			+1
			-1
START			

actions: UP, DOWN, LEFT, RIGHT

UP

80% move UP
10% move LEFT
10% move RIGHT



reward +1 at [4,3], -1 at [4,2]
reward -0.04 for each step

states: current location

actions: where to go next

rewards

what is the solution? Learn a mapping from (state, action) pairs to new states

Markov Decision Process: Formalizing Reinforcement Learning

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

set of possible states set of possible actions reward of (state, action) pairs state-action transition distribution discount factor

Start in initial state s_0

Markov Decision Process: Formalizing Reinforcement Learning

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

set of possible states set of possible actions state-action transition distribution reward of (state, action) pairs discount factor

Start in initial state s_0
for $t = 1$ to ...:
 choose action a_t

Markov Decision Process: Formalizing Reinforcement Learning

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

set of possible actions state-action transition distribution
set of possible states reward of (state, action) pairs discount factor

Start in initial state s_0

for $t = 1$ to ...:

choose action a_t

“move” to next state $s_t \sim \pi(\cdot | s_{t-1}, a_t)$

Policy
 $\pi: S \rightarrow A$

Markov Decision Process: Formalizing Reinforcement Learning

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

set of possible actions state-action transition distribution
set of possible states reward of (state, action) pairs discount factor

Start in initial state s_0

for $t = 1$ to ...:

choose action a_t

“move” to next state $s_t \sim \pi(\cdot | s_{t-1}, a_t)$

get reward $r_t = \mathcal{R}(s_t, a_t)$

Markov Decision Process: Formalizing Reinforcement Learning

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

set of possible actions state-action transition distribution
set of possible states reward of (state, action) pairs discount factor

Start in initial state s_0

for $t = 1$ to ...:

choose action a_t

“move” to next state $s_t \sim \pi(\cdot | s_{t-1}, a_t)$

get reward $r_t = \mathcal{R}(s_t, a_t)$

objective: choose
action over time
to maximize time-
discounted reward

Markov Decision Process: Formalizing Reinforcement Learning

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

set of possible actions state-action transition distribution
set of possible states reward of (state, action) pairs discount factor

Start in initial state s_0

for $t = 1$ to ...:

choose action a_t

“move” to next state $s_t \sim \pi(\cdot | s_{t-1}, a_t)$

get reward $r_t = \mathcal{R}(s_t, a_t)$

objective: choose action over
time to maximize discounted
reward

Consider all
possible future
times t

Reward at
time t

Markov Decision Process: Formalizing Reinforcement Learning

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

set of possible actions state-action transition distribution
set of possible states reward of (state, action) pairs discount factor

Start in initial state s_0

for $t = 1$ to ...:

choose action a_t

“move” to next state $s_t \sim \pi(\cdot | s_{t-1}, a_t)$

get reward $r_t = \mathcal{R}(s_t, a_t)$

objective: maximize
discounted reward

Consider all
possible future
times t

Discount at
time t

Reward at
time t

Markov Decision Process: Formalizing Reinforcement Learning

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

set of possible actions state-action transition distribution
set of possible states reward of (state, action) pairs discount factor

Start in initial state s_0

for $t = 1$ to ...:

choose action a_t

“move” to next state $s_t \sim \pi(\cdot | s_{t-1}, a_t)$

get reward $r_t = \mathcal{R}(s_t, a_t)$

objective: maximize
discounted reward

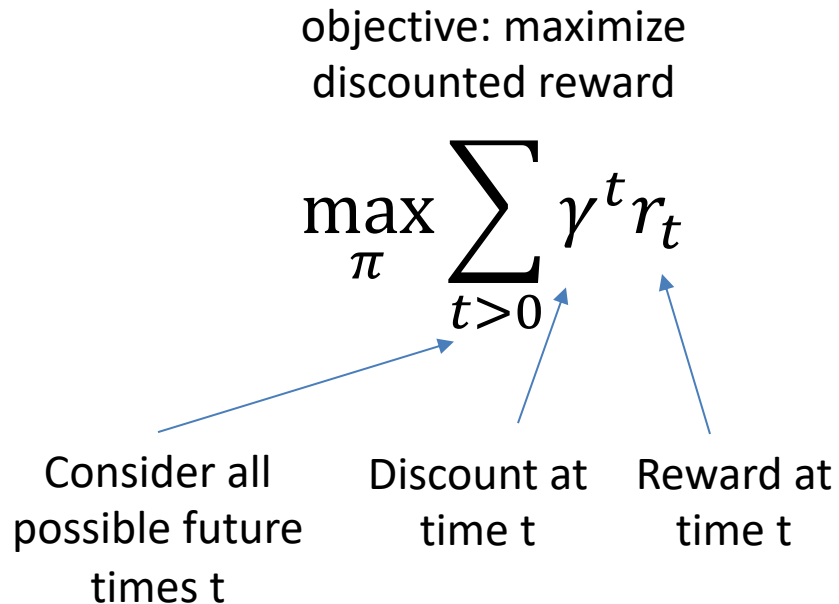
$$\max_{\pi} \sum_{t>0} \gamma^t r_t$$

Consider all
possible future
times t

Discount at
time t

Reward at
time t

Example of Discounted Reward



- If the discount factor $\gamma = 0.8$ then reward $0.8^0 r_0 + 0.8^1 r_1 + 0.8^2 r_2 + 0.8^3 r_3 + \dots + 0.8^n r_n + \dots$
- Allows you to consider all possible rewards in the future but preferring current vs. future self

Markov Decision Process: Formalizing Reinforcement Learning

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

set of possible actions state-action transition distribution
set of possible states reward of (state, action) pairs discount factor

Start in initial state s_0

for $t = 1$ to ...:

choose action a_t

“move” to next state $s_t \sim \pi(\cdot | s_{t-1}, a_t)$

get reward $r_t = \mathcal{R}(s_t, a_t)$

objective: maximize
discounted reward

$$\max_{\pi} \sum_{t>0} \gamma^t r_t$$

“solution”: the policy π^* that maximizes the
expected (average) time-discounted reward

Markov Decision Process: Formalizing Reinforcement Learning

Markov Decision
Process:

$$(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$$

set of possible actions state-action transition distribution
set of possible states reward of (state, action) pairs discount factor

Start in initial state s_0

for $t = 1$ to ...:

choose action a_t

“move” to next state $s_t \sim \pi(\cdot | s_{t-1}, a_t)$

get reward $r_t = \mathcal{R}(s_t, a_t)$

objective: maximize
discounted reward

$$\max_{\pi} \sum_{t>0} \gamma^t r_t$$

$$\text{“solution” } \pi^* = \operatorname{argmax}_{\pi} \mathbb{E} \left[\sum_{t>0} \gamma^t r_t ; \pi \right]$$

Markov Decision Process: Formalizing Reinforcement Learning

Mar

Here, r_t is a function of random variable s_t .

Start in initial

for $t = 1$ to ..

choose action a_t

“move” to next state $s_t \sim \pi(\cdot | s_{t-1}, a_t)$

get reward $r_t = \mathcal{R}(s_t, a_t)$

$$\max_{\pi} \sum_{t>0} \gamma^t r_t$$

“solution” $\pi^* = \operatorname{argmax}_{\pi} \mathbb{E} \left[\sum_{t>0} \gamma^t r_t ; \pi \right]$

Markov Decision Process: Formalizing Reinforcement Learning

Here, r_t is a function of random variable s_t . →

The expectation is over the different states s_t the agent could be in at time t (equiv. actions the agent could take).

Mar

Start in initial

for $t = 1$ to ..

choose action a_t

“move” to next state $s_t \sim \pi(\cdot | s_{t-1}, a_t)$

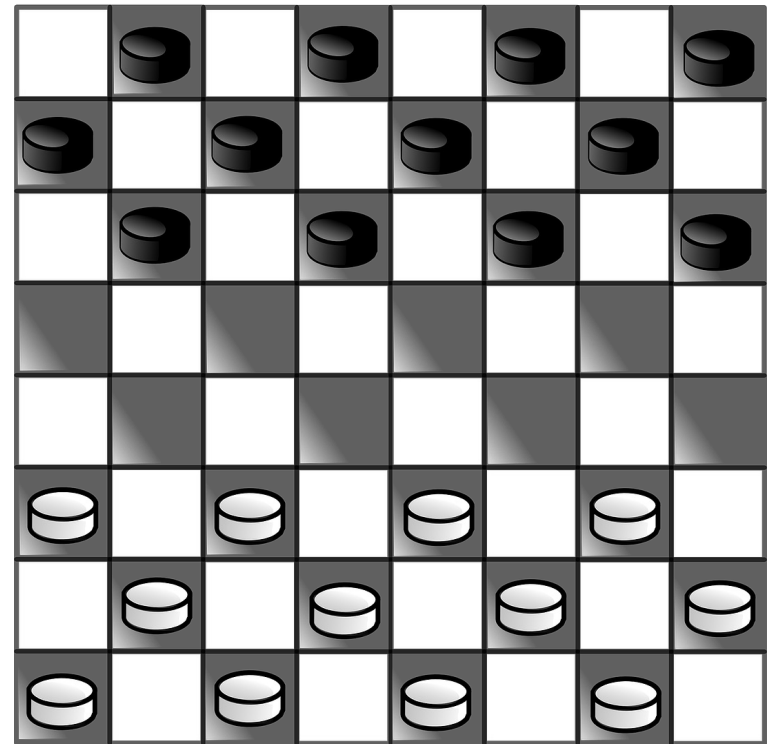
get reward $r_t = \mathcal{R}(s_t, a_t)$

$$\max_{\pi} \sum_{t>0} \gamma^t r_t$$

“solution” $\pi^* = \operatorname{argmax}_{\pi} \mathbb{E} \left[\sum_{t>0} \gamma^t r_t ; \pi \right]$

Simple Example

- Learn to play checkers
 - Two-person game
 - 8x8 boards, 12 checkers/side
 - relatively simple set of rules:
<http://www.darkfish.com/checkers/rules.html>
 - Goal is to eliminate all your opponent's pieces



Some Challenges

1. Representing states (and actions)
2. Defining our reward
3. Learning our policy

Overview: Learning Strategies


Dynamic Programming

Q-learning

Monte Carlo approaches

Reactive Agent Algorithm

Repeat:

- ◆ $s \leftarrow$ sensed state 
- ◆ If s is a terminal state then exit
- ◆ $a \leftarrow$ choose action (given s)
- ◆ Perform a

Policy (Reactive/Closed-Loop Strategy)

3	→	→	→	+1
2	↑		↑	-1
1	↑	←	←	←
	1	2	3	4

- In every state, we need to know what to do
- The **goal** doesn't change
- A **policy** (Π) is a complete mapping from *states* to *actions*
 - “If in [3,2], go up; if in [3,1], go left; if in...”

Optimal Policy

3	→	→	→	+1
2	↑		↑	-1
1	↑	←	←	←
	1	2	3	4

- A **policy** Π is a complete mapping from states to actions
- The **optimal policy** Π^* is the one that always yields a history (sequence of steps ending at a terminal state) with maximal ***expected*** utility

Optimal Policy

3	→	→	→	+1
2	↑		↑	-1
1	↑	←	←	←
	1	2	3	4

- A **policy** Π is a complete strategy that tells the agent what action to take in each state.
- The **optimal policy** Π^* is the policy that yields the highest expected utility for every state in the environment.

This problem is called a Markov Decision Problem (MDP)

How to compute Π^* ?

Defining Value Function

- Problem:
 - When making a decision, we only know the reward so far, and the possible actions

Defining Value Function

- Problem:
 - When making a decision, we only know the reward so far, and the possible actions
 - We've defined value function retroactively (i.e., the value function/utility of a history/sequence of states is known *once we finish it*)

Defining Value Function

- Problem:
 - When making a decision, we only know the reward so far, and the possible actions
 - We've defined value function retroactively (i.e., the value function/utility of a history/sequence of states is known *once we finish it*)
 - What is the value function of a particular ***state*** in the middle of decision making?

Defining Value Function

- Problem:
 - When making a decision, we only know the reward so far, and the possible actions
 - We've defined value function retroactively (i.e., the value function/utility of a history/sequence of states is known *once we finish it*)
 - What is the value function of a particular ***state*** in the middle of decision making?
 - Need to compute ***expected value function*** of possible future histories/states

Defining Value Function

$$V^\pi(s) = \mathbb{E} [R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots \mid s_0 = s, \pi].$$

$V^\pi(s)$ is simply the expected sum of discounted rewards upon starting in state s , and taking actions according to π .¹

Given a fixed policy π , its value function V^π satisfies the **Bellman equations**:

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in \mathcal{S}} P_{s\pi(s)}(s') V^\pi(s').$$



- What is the value function of a particular ***state*** in the middle of decision making?
- Need to compute ***expected value function*** of possible future histories/states

Dynamic programming

use value functions to structure the search for good policies





Dynamic programming

use value functions to structure the search for good policies

 policy evaluation: compute V^π from π
policy improvement: improve π based on V^π 

Dynamic programming

use value functions to structure the search for good policies

 policy evaluation: compute V^π from π 
 policy improvement: improve π based on V^π 

start with an arbitrary policy

repeat evaluation/improvement until convergence