

COMP 532

**Machine Learning and
BioInspired Optimization**

Lecture 4 Reinforcement Learning

Dr. Shan Luo

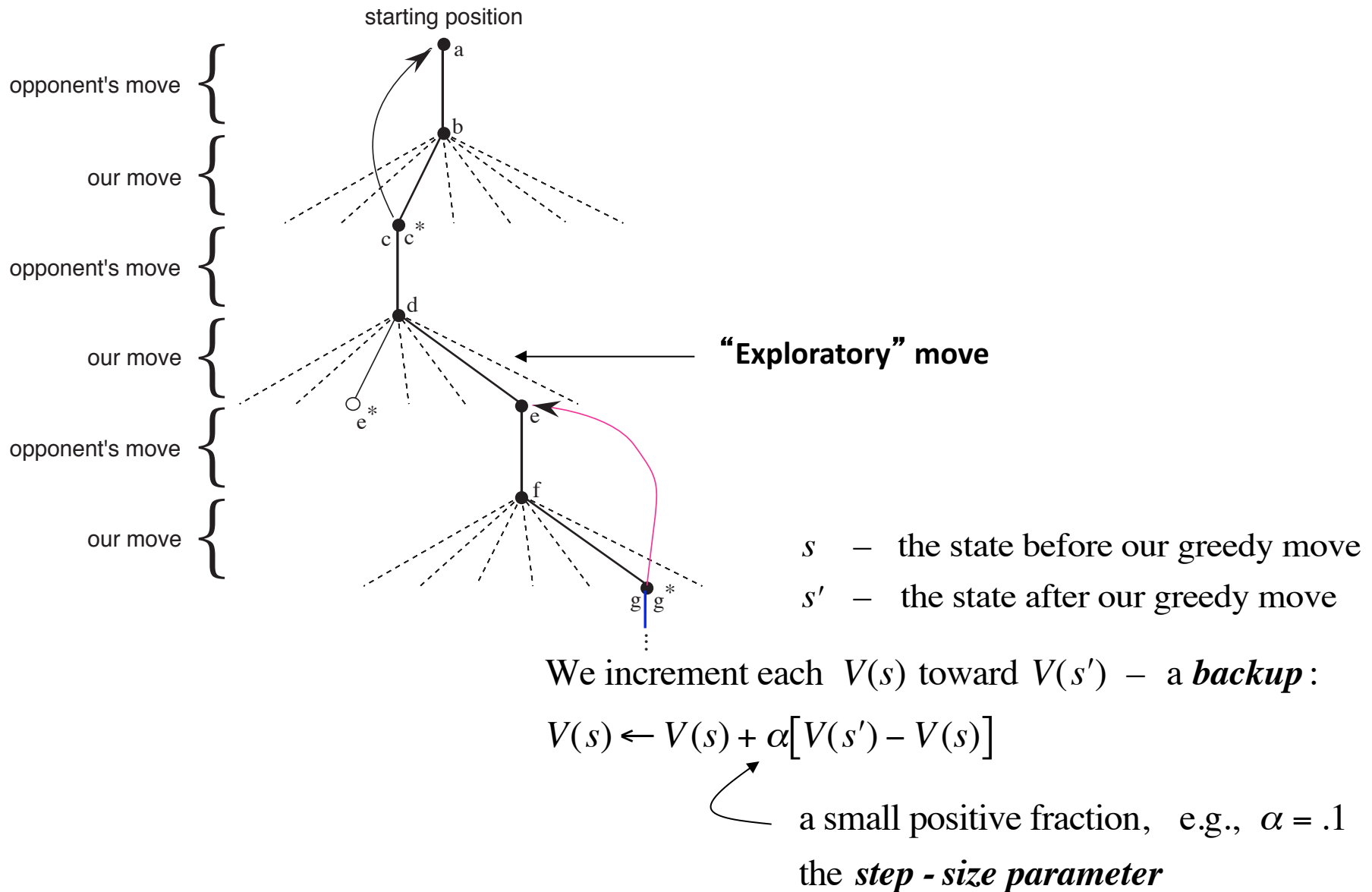
Department of computer science

shan.luo@liverpool.ac.uk

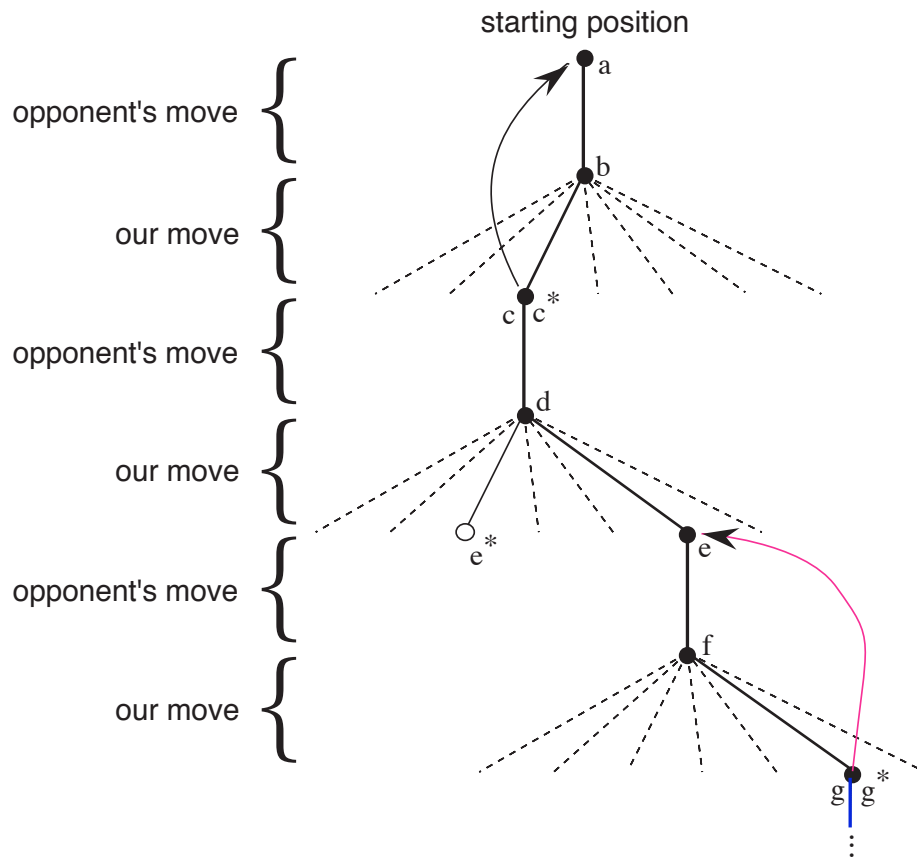
The Book

- Part I: The Problem
 - Introduction
 - Evaluative Feedback
 - The Reinforcement Learning Problem
- Part II: Elementary Solution Methods
 - Dynamic Programming
 - Monte Carlo Methods
 - Temporal Difference Learning
- Part III: A Unified View
 - Eligibility Traces
 - Generalization and Function Approximation
 - Planning and Learning
 - Dimensions of Reinforcement Learning
 - Case Studies

RL Learning Rule for Tic-Tac-Toe



RL Learning Rule for Tic-Tac-Toe



$$V(a) \leftarrow V(a) + \alpha[V(c) - V(a)]$$

$$V(e) \leftarrow V(e) + \alpha[V(g) - V(e)]$$

s — the state before our greedy move

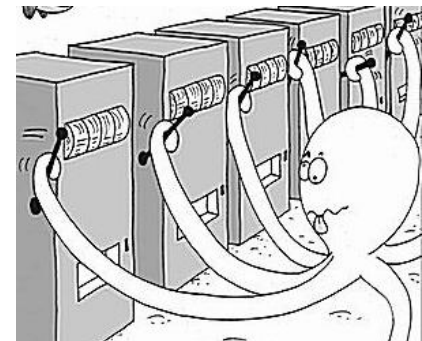
s' — the state after our greedy move

We increment each $V(s)$ toward $V(s')$ — a **backup**:

$$V(s) \leftarrow V(s) + \alpha[V(s') - V(s)]$$

Evaluative Feedback

- **Evaluating** actions vs. **instructing** by giving correct actions
- Pure evaluative feedback depends totally on the action taken.
Pure instructive feedback depends not at all on the action taken.
- Supervised learning is instructive; optimization is evaluative
- **Associative vs. Nonassociative:**
 - Associative: inputs mapped to outputs;
learn the best output **for each** input
 - Nonassociative: “learn” (find) one best output
- *n*-armed bandit (at least how we treat it) is:
 - Nonassociative
 - Evaluative feedback



The n -Armed Bandit Problem

- Choose repeatedly from one of n actions; each choice is called a **play**
- After each play a_t , you get a reward r_t , where

$$E \{ r_t \mid a_t \} = Q^*(a_t)$$

These are unknown ***action values***

Distribution of r_t depends only on a_t

- Objective is to maximize the reward in the long term, e.g., over 1000 plays

To solve the n -armed bandit problem,
you must **explore** a variety of actions
and **exploit** the best of them

The Exploration/Exploitation Dilemma

- Suppose you form estimates

$$Q_t(a) \approx Q^*(a) \quad \text{action value estimates}$$

- The **greedy** action at t is a_t^*

$$a_t^* = \arg \max_a Q_t(a)$$

$$a_t = a_t^* \Rightarrow \text{exploitation}$$

$$a_t \neq a_t^* \Rightarrow \text{exploration}$$

- You can't exploit all the time; you can't explore all the time
- You can never stop exploring; but you should always reduce exploring. Maybe.

Action-Value Methods

- Methods that adapt action-value estimates and nothing else, e.g.: suppose by the t -th play, action a had been chosen k_a times, producing rewards r_1, r_2, \dots, r_{k_a} , then

$$Q_t(a) = \frac{r_1 + r_2 + \dots + r_{k_a}}{k_a}$$

“sample average”

- $\lim_{k_a \rightarrow \infty} Q_t(a) = Q^*(a)$

ϵ -Greedy Action Selection

- Greedy action selection:

$$a_t = a_t^* = \arg \max_a Q_t(a)$$

- ϵ -Greedy:

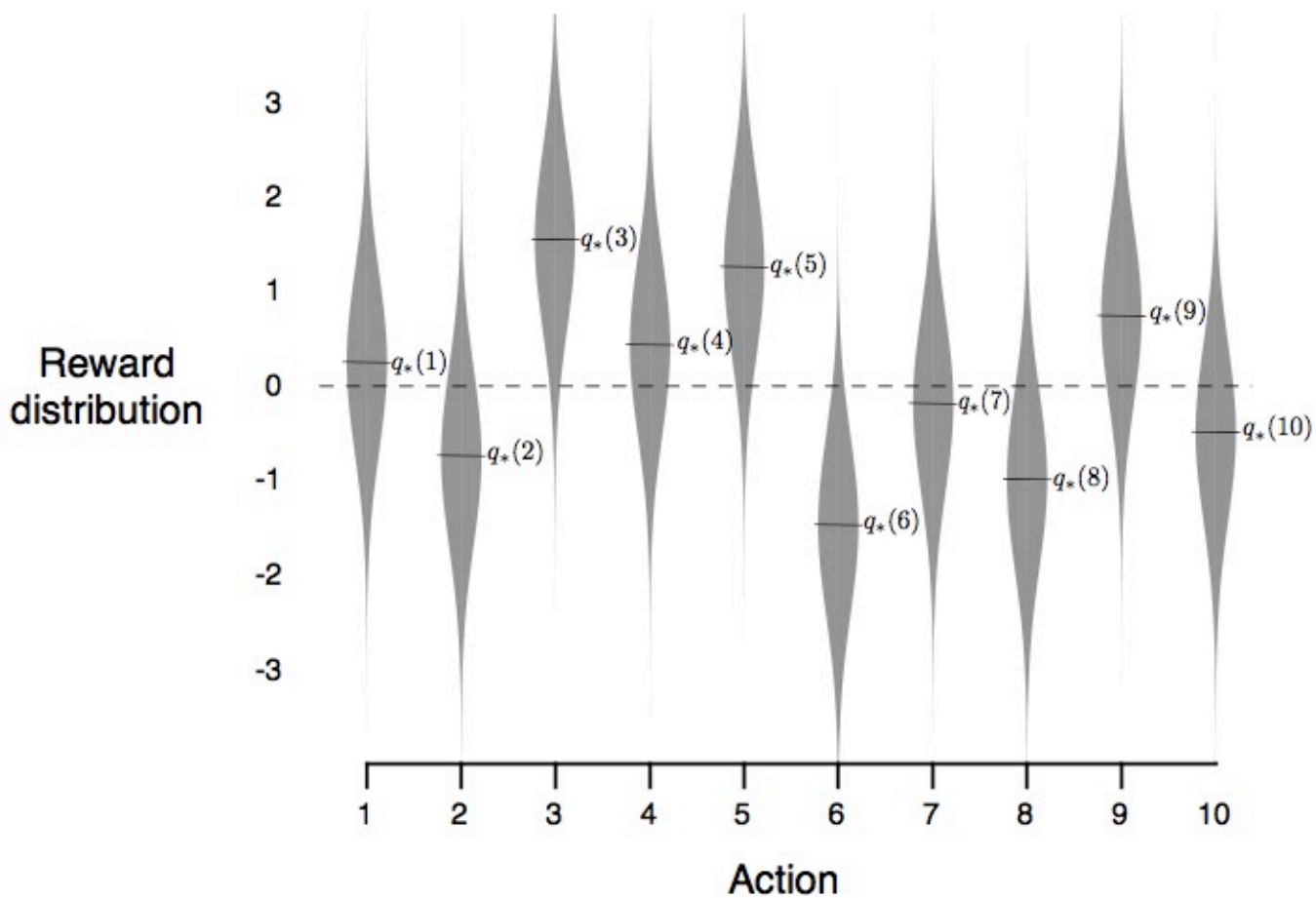
$$a_t = \begin{cases} a_t^* & \text{with probability } 1 - \epsilon \\ \text{random action} & \text{with probability } \epsilon \end{cases}$$

... the simplest way to balance exploration and exploitation

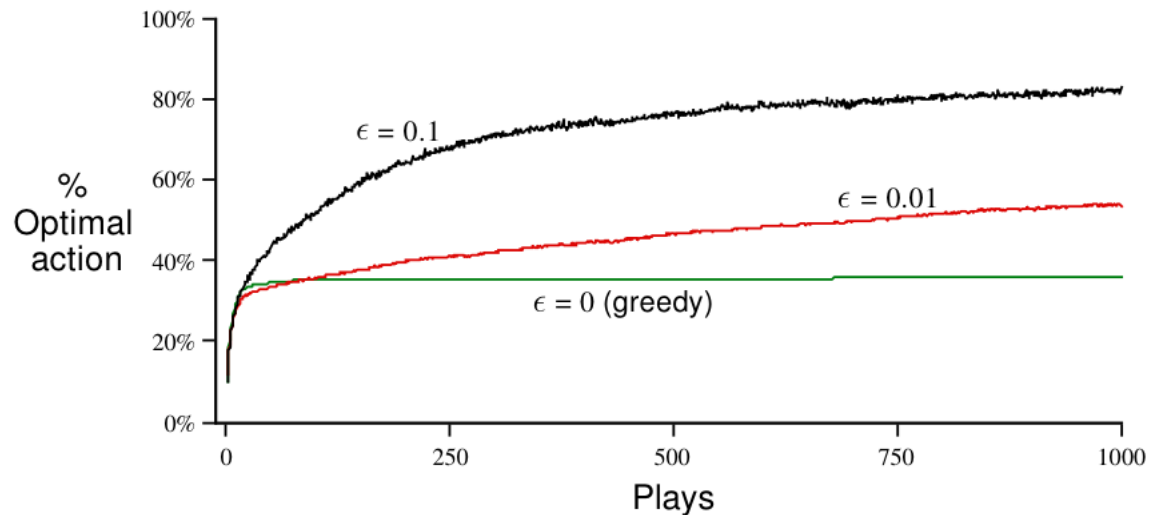
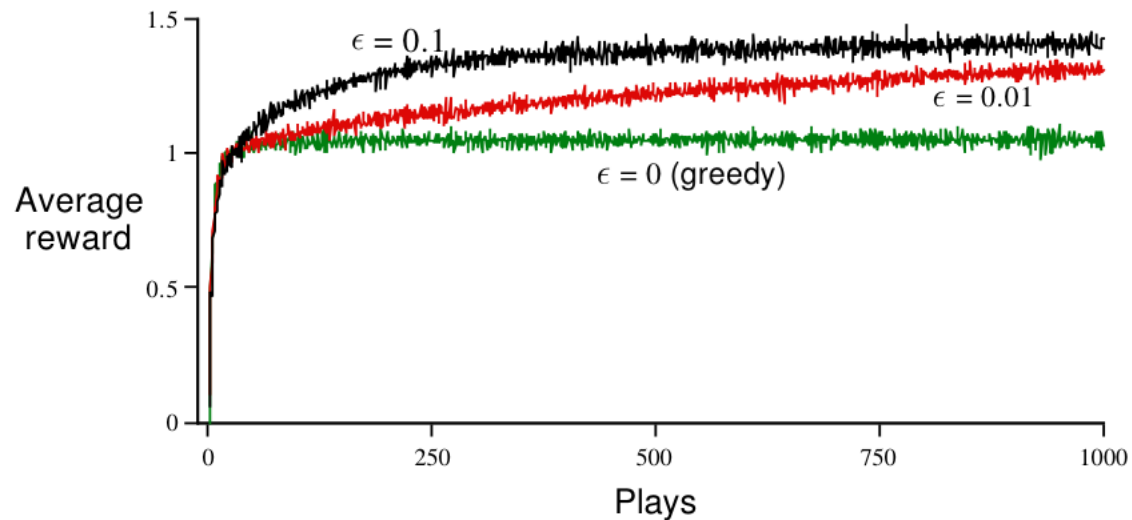
10-Armed Testbed

- $n = 10$ possible actions
- Each $Q^*(a)$ is chosen randomly from a normal distrib.: $\eta(0, 1)$
- each r_t is also normal: $\eta(Q^*(a_t), 1)$
- 1000 plays
- repeat the whole thing 2000 times and average the results

10-Armed Testbed



ϵ -Greedy Methods on the 10-Armed Testbed



Softmax Action Selection

- Softmax action selection methods define action probabilities by estimated values.
- The most common softmax uses a Gibbs, or Boltzmann, distribution:

Choose action a on play t with probability

$$\frac{e^{Q_t(a)/\tau}}{\sum_{b=1}^n e^{Q_t(b)/\tau}}$$

where τ is the

“computational temperature”

Incremental Implementation

Recall the sample average estimation method:

The average of the first k rewards is
(dropping the dependence on a):

$$Q_k = \frac{r_1 + r_2 + \cdots + r_k}{k}$$

Can we do this incrementally (without storing all the rewards)?

We could keep a running sum and count, or, equivalently:

$$Q_{k+1} = Q_k + \frac{1}{k+1} [r_{k+1} - Q_k]$$

This is a common form for update rules:

$$\textit{NewEstimate} = \textit{OldEstimate} + \textit{StepSize}[\textit{Target} - \textit{OldEstimate}]$$

Tracking a Nonstationary Problem

Choosing Q_k to be a sample average is appropriate in a stationary problem

(i.e., when none of the $Q^*(a)$ change over time)

But not in a nonstationary problem.

Better in the nonstationary case is:

$$Q_{k+1} = Q_k + \alpha [r_{k+1} - Q_k]$$

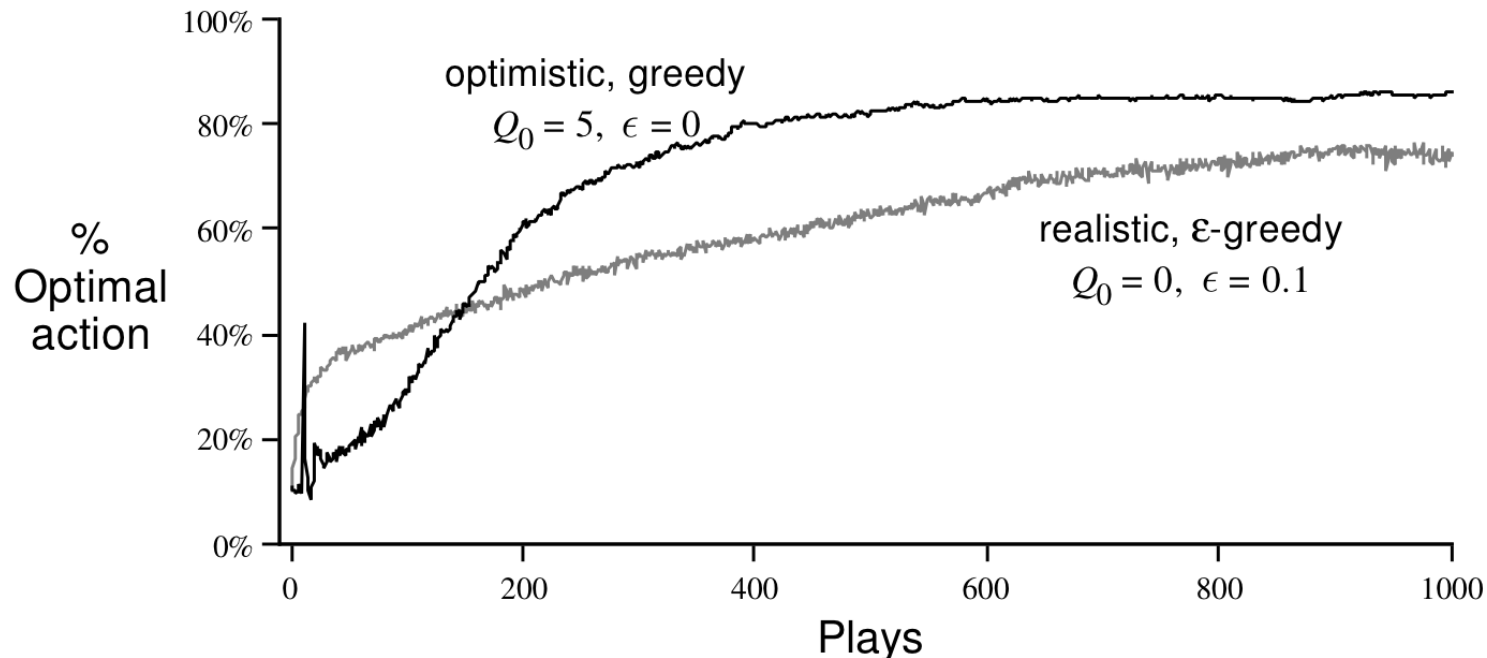
for *constant* α , $0 < \alpha \leq 1$

$$= (1 - \alpha)^k Q_0 + \sum_{i=1}^k \alpha (1 - \alpha)^{k-i} r_i$$

exponential, recency-weighted average

Optimistic Initial Values

- All methods so far depend on $Q_0(a)$, i.e., they are **biased**.
- Suppose instead we initialize the action values optimistically, i.e., on the 10-armed testbed, use $Q_0(a) = 5$ for all a



Conclusions

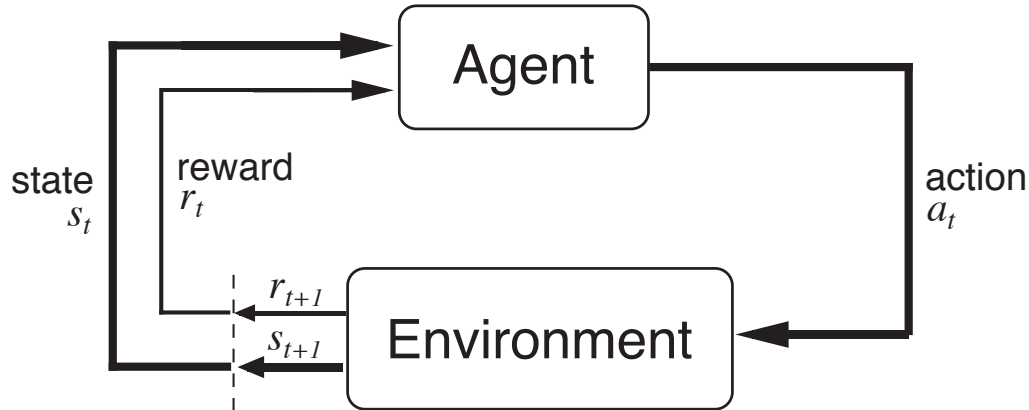
- These are all very simple methods
 - but they are complicated enough—we will build on them
 - we should understand them *completely*

The Reinforcement Learning Problem

Markov Decision Processes and Value Functions

- describe the RL problem we will be studying for the remainder of this part of the course
- present idealized form of the RL problem for which we have precise theoretical results;
- introduce key components of the mathematics: value functions and Bellman equations;

The Agent-Environment Interface



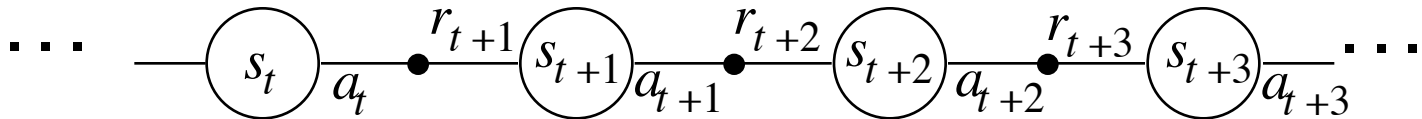
Agent and environment interact at discrete time steps : $t = 0, 1, 2, K$

Agent observes state at step t : $s_t \in S$

produces action at step t : $a_t \in A(s_t)$

gets resulting reward : $r_{t+1} \in \mathfrak{R}$

and resulting next state : s_{t+1}



The Agent Learns a Policy

Policy at step t , π_t :

a mapping from states to action probabilities

$\pi_t(s, a) =$ probability that $a_t = a$ when $s_t = s$

- Reinforcement learning methods specify how the agent changes its policy as a result of experience.
- Roughly, the agent's goal is to get as much reward as it can over the long run.

Getting the Degree of Abstraction Right

- Time steps need not refer to fixed intervals of real time.
- Actions can be low level (e.g., voltages to motors), or high level (e.g., accept a job offer), “mental” (e.g., shift in focus of attention), etc.
- States can be low-level “sensations”, or they can be abstract, symbolic, based on memory, or subjective (e.g., the state of being “surprised” or “lost”).
- An RL agent is not like a *whole* animal or robot.
- Reward computation is in the agent’s environment because the agent cannot change it arbitrarily.
- The environment is not necessarily unknown to the agent, only incompletely controllable.

Goals and Rewards

- Is a scalar reward signal an adequate notion of a goal?—maybe not, but it is surprisingly flexible.
- A goal should specify **what** we want to achieve, not **how** we want to achieve it.
- A goal must be outside the agent's direct control—thus outside the agent.
- The agent must be able to measure success:
 - explicitly;
 - frequently during its lifespan.

The reward hypothesis

- That all of what we mean by goals and purposes can be well thought of as the maximization of the cumulative sum of a received scalar signal (reward)

Returns

Suppose the sequence of rewards after step t is:

$$r_{t+1}, r_{t+2}, r_{t+3}, \dots$$

What do we want to maximize?

In general,

we want to maximize the **expected return**, $E\{R_t\}$, for each step t .

Episodic tasks: interaction breaks naturally into episodes, e.g., plays of a game, trips through a maze.

$$R_t = r_{t+1} + r_{t+2} + \dots + r_T,$$

where T is a final time step at which a **terminal state** is reached, ending an episode.

Returns for Continuing Tasks

Continuing tasks: interaction does not have natural episodes.

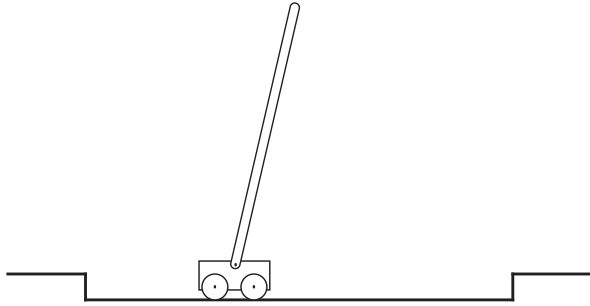
Discounted return:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1},$$

where γ , $0 \leq \gamma \leq 1$, is the **discount rate**.

shortsighted $0 \leftarrow \gamma \rightarrow 1$ farsighted

An Example



Avoid **failure**: the pole falling beyond a critical angle or the cart hitting end of track.

As an **episodic task** where episode ends upon failure:

reward = +1 for each step before failure

\Rightarrow return = number of steps before failure

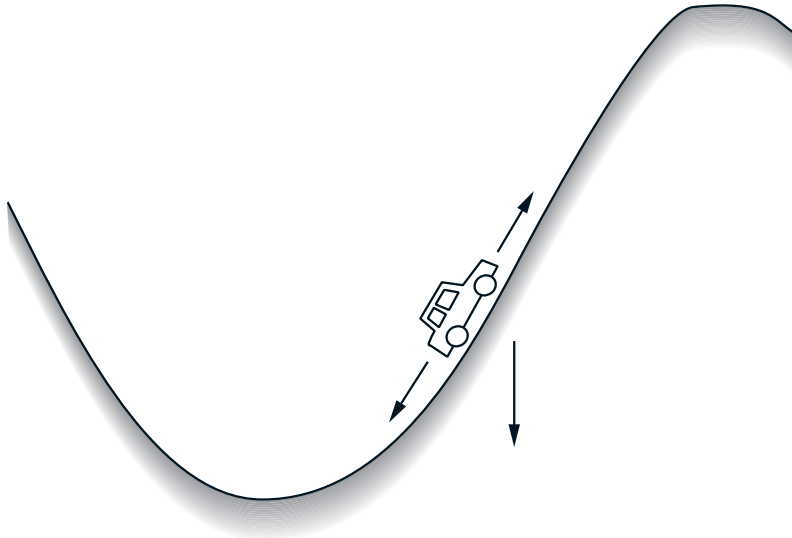
As a **continuing task** with discounted return:

reward = -1 upon failure; 0 otherwise

\Rightarrow return = $-\gamma^k$, for k steps before failure

In either case, return is maximized by avoiding failure for as long as possible.

Another Example



Get to the top of the hill
as quickly as possible.

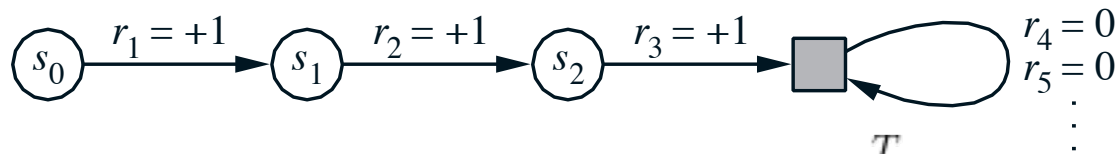
reward = -1 for each step where **not** at top of hill

\Rightarrow return = - number of steps before reaching top of hill

Return is maximized by minimizing
number of steps to reach the top of the hill.

A Unified Notation

- In episodic tasks, we number the time steps of each episode starting from zero.
- We usually do not have to distinguish between episodes, so we write s_t instead of $s_{t,j}$ for the state at step t of episode j .
- Think of each episode as ending in an absorbing state that always produces reward of zero:



- We can cover all cases by writing
$$R_t = \sum_{k=0}^T \gamma^k r_{t+k+1},$$

including the possibility that $T = \infty$ or $\gamma = 1$ (not both)

The Markov Property

- By “the state” at step t , the book means whatever information is available to the agent at step t about its environment.
- The state can include immediate “sensations”, highly processed sensations, and structures built up over time from sequences of sensations.
- Ideally, a state should summarize past sensations so as to retain all “essential” information, i.e., it should have the

Markov Property:

$$\begin{aligned} Pr \{s_{t+1} = s', r_{t+1} = r \mid s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_1, s_0, a_0\} \\ = Pr \{s_{t+1} = s', r_{t+1} = r \mid s_t, a_t\}, \end{aligned}$$

for all s' and r , and all histories $s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_1, s_0, a_0$.

Markov Decision Processes

- If a reinforcement learning task has the Markov Property, it is basically a **Markov Decision Process (MDP)**.
- If state and action sets are finite, it is a **finite MDP**.
- To define a finite MDP, you need to give:
 - **state and action sets**
 - one-step “dynamics” defined by **transition probabilities**:

$$\mathbf{P}_{ss'}^a = \Pr\{s_{t+1} = s' \mid s_t = s, a_t = a\} \quad \text{for all } s, s' \in S, a \in A(s).$$

- **reward probabilities**:

$$\mathbf{R}_{ss'}^a = E\{r_{t+1} \mid s_t = s, a_t = a, s_{t+1} = s'\} \quad \text{for all } s, s' \in S, a \in A(s).$$

An Example Finite MDP

Recycling Robot

- At each step, robot has to decide whether it should (1) actively search for a can, (2) wait for someone to bring it a can, or (3) go to home base and recharge.
- Searching is better but runs down the battery; if runs out of power while searching, has to be rescued (which is bad).
- Decisions made on basis of current energy level: high, low.
- Reward = number of cans collected

Recycling Robot MDP

$$S = \{\text{high}, \text{low}\}$$

$$A(\text{high}) = \{\text{search}, \text{wait}\}$$

$$A(\text{low}) = \{\text{search}, \text{wait}, \text{recharge}\}$$

$\mathbf{R}^{\text{search}}$ = expected no. of cans while searching

\mathbf{R}^{wait} = expected no. of cans while waiting

$$\mathbf{R}^{\text{search}} > \mathbf{R}^{\text{wait}}$$

