### **Reinforcement Learning**

Slides by Rich Sutton

Mods by Dan Lizotte

Refer to "Reinforcement Learning: An Introduction" by Sutton and Barto Alpaydin Chapter 16











































## The Exploration/Exploitation Dilemma

Suppose you form estimates

 $Q_t(a) \approx Q^*(a)$  action value estimates

• The greedy action at t is  $a_t$ 

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

$$a_t = a_t^* \Rightarrow$$
 exploitation

 $a_t \neq a_t^* \Rightarrow$  exploration

- If you need to learn, you can't exploit all the time; if you need to do well, you can't explore all the time
- You can never stop exploring; but you should always reduce exploring. Maybe.













### 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:

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

This is a common form for update rules:

*NewEstimate* = *OldEstimate* + *StepSize*[*Target* – *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  $\alpha$ ,  $0 < \alpha \le 1$  $= (1 - \alpha)^k Q_0 + \sum_{i=1}^k \alpha (1 - \alpha)^{k-i} r_i$ exponential, recency-weighted average



### Conclusions

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



### Chapter 3: The Reinforcement Learning Problem

Objectives of this chapter:

- describe the RL problem we will be studying for the remainder 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;
- describe trade-offs between applicability and mathematical tractability.



### The Agent Learns a Policy

**Policy** at step *t*,  $\pi_t$ :

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 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.



### 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**:

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

where  $\gamma, 0 \le \gamma \le 1$ , is the **discount rate** 

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

















# Bellman Equation for a Policy $\pi$ The basic idea: $R_{t} = r_{t+1} + \gamma r_{t+2} + \gamma^{2} r_{t+3} + \gamma^{3} r_{t+4} \cdots$ $= r_{t+1} + \gamma \left( r_{t+2} + \gamma r_{t+3} + \gamma^{2} r_{t+4} \cdots \right)$ $= r_{t+1} + \gamma R_{t+1}$ So: $V^{\pi}(s) = E_{\pi} \left\{ R_{t} \mid s_{t} = s \right\}$ $= E_{\pi} \left\{ r_{t+1} + \gamma V(s_{t+1}) \mid s_{t} = s \right\}$ Or, without the expectation operator:

$$V^{\pi}(s) = \sum_{a} \pi(s, a) \sum_{s'} \mathbf{P}_{ss'}^{a} \left[ \mathbf{R}_{ss'}^{a} + \gamma V^{\pi}(s') \right]$$



















### Solving the Bellman Optimality Equation

- Finding an optimal policy by solving the Bellman Optimality Equation requires the following:
  - accurate knowledge of environment dynamics;
  - we have enough space and time to do the computation;
  - the Markov Property.
- How much space and time do we need?
  - polynomial in number of states (via dynamic programming methods; Chapter 4),
  - BUT, number of states is often huge (e.g., backgammon has about 10<sup>20</sup> states).
- We usually have to settle for approximations.
- Many RL methods can be understood as approximately solving the Bellman Optimality Equation.

