

**CS 287: Advanced Robotics  
Fall 2009**

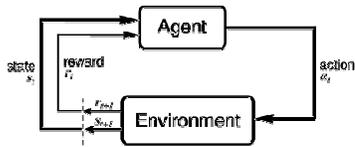
Lecture 9: Reinforcement Learning 1: Bandits

Pieter Abbeel  
UC Berkeley EECS

**Multi-armed bandits**

- Slot machines
- Clinical trials
- Advertising
- Merchandising

**Reinforcement Learning**



- Model: Markov decision process (S, A, T, R,  $\gamma$ )
  - Goal: Find  $\pi$  that maximizes expected sum of rewards
  - T and R might be unknown

[Drawing from Sutton and Barto, Reinforcement Learning: An Introduction, 1998]

**Multi-armed bandits**

Consider slot machines  $H_1, H_2, \dots, H_n$ .

Slot machine  $i$  has pay-off =  $\begin{cases} 0, & \text{with probability } 1 - \theta_i \\ 1, & \text{with probability } \theta_i \end{cases}$   
where  $\theta_i$  is unknown.

Now the objective to maximize is:

$$E \left[ \sum_{t=1}^{\infty} \gamma^t R(s_t, a_t) \right] \text{ (where } s_t \text{ is unchanged).}$$

**Exploration vs. exploitation**

- = classical dilemma in reinforcement learning
- A conceptual solution: Bayesian approach:
  - State space =  $\{x : x = \text{probability distribution over } T, R\}$ 
    - For known initial state --- tree of sufficient statistics could suffice
  - Transition model: describes transitions in new state space
  - Reward = standard reward
- Today: one particular setting in which the Bayesian solution is in fact computationally practical

**Information state**

$$\begin{pmatrix} \# \text{ of successes on } H_1 \\ \# \text{ of failures on } H_1 \\ \# \text{ of successes on } H_2 \\ \# \text{ of failures on } H_2 \\ \vdots \\ \# \text{ of successes on } H_n \\ \# \text{ of failures on } H_n \end{pmatrix}$$

## Semi-MDP

- Transition model:

$$P(s_{t+1} = s', t_{t+1} = t_k + \Delta | s_t = s, a_t = \alpha)$$

- Objective:

$$E[\sum_{k=0}^{\infty} \gamma^k R(s_{t_{k+1}}, \Delta_k, s_{t_k})]$$

- Bellman update:

$$V(s) = \max_{\alpha} \sum_{s', \Delta} P(s', \Delta | s, \alpha) [R(s, \Delta, s') + \gamma V(s')]$$

## Optimal stopping

- One approach:
  - Solve the optimal stopping problem for many values of  $g$ , and for each state keep track of the smallest value of  $g$  which causes stopping

## Optimal stopping

- A specialized version of the Semi-MDP is the "Optimal stopping problem". At each of the times, we have two choices:
  - continue
  - stop and accumulate reward  $g$  for current time and for all future times.
- The optimal stopping problem has the following Bellman update:

$$V(s) = \max\left\{\sum_{s', \Delta} P(s', \Delta | s) [R(s, \Delta, s') + \gamma V(s')], \frac{g}{1-\gamma}\right\}$$

## Reward rate

- Reward rate

$$\sum_{t=1}^{\Delta_{t_0}-1} \gamma^t r(s_{t_0}, \Delta_{t_0}, s_{t_0+t}) = R(s_{t_0}, \Delta_{t_0}, s_{t_0+t})$$

- Expected reward rate

$$\bar{r}(s) = \mathbb{E}_{\Delta, s'} [r(s, \Delta, s')] = \sum_{\Delta, s'} P(s', \Delta | s) r(s, \Delta, s')$$

## Optimal stopping

- Optimal stopping Bellman update:

$$V(s) = \max\left\{\sum_{s', \Delta} P(s', \Delta | s) [R(s, \Delta, s') + \gamma V(s')], \frac{g}{1-\gamma}\right\}$$

- Hence, for fixed  $g$ , we can find the value of each state in the optimal stopping problem by dynamic programming

- However, we are interested in  $g^*(s)$  for all  $s$ :

$$g^*(s) = \min\left\{g, \frac{g}{1-\gamma} \geq \max_{\tau} \mathbb{E}_{\tau} \left[ \sum_{k=0}^{\tau-1} \gamma^k R(s_{t_k}, \Delta_k, s_{t_{k+1}}) + \sum_{t=\tau}^{\infty} \gamma^t g \right] \right\}$$

- Note:  $\tau$  is a random variable, which denotes the stopping time. It is the policy in this setting.
- Any stopping policy can be represented as a set of states in which we decided to stop. The random variable  $\tau$  takes on the value = time when we first visit a state in the stopping set.

## Basic idea to find $g^*$

Now consider

$$s^* = \arg\max_s V(s)$$

Of all states, the state  $s^*$  would receive the highest payoff if we were willing to stop. Namely,

$$g^*(s^*) = r(s^*)$$

This means that when  $g < g^*(s^*)$ , the optimal stopping policy will choose to continue at  $s^*$ . Note that for  $s \neq s^*$ ,  $g^*(s) < g^*(s^*)$ .

To compute  $g^*(s)$  for the other states, we consider a new semi-MDP which differs from the existing one only in that we always continue when at state  $s^*$ . This is equivalent to letting the new state space  $\mathcal{S} = \mathcal{S} \cup \{s^*\}$ .

