#### **Chapter 5: Monte Carlo Methods**

- ☐ Monte Carlo methods learn from *complete* sample returns
  - Only defined for episodic tasks
- □ Monte Carlo methods learn directly from experience
  - *On-line:* No model necessary and still attains optimality
  - *Simulated:* No need for a *full* model

#### **Monte Carlo Policy Evaluation**

- **Goal:** learn  $V^{\pi}(s)$
- **Given:** some number of episodes under  $\pi$  which contain *s*
- **Idea:** Average returns observed after visits to s



- Every-Visit MC: average returns for every time s is visited in an episode
- □ *First-visit MC:* average returns only for *first* time *s* is visited in an episode
- **D** Both converge asymptotically

R. S. Sutton and A. G. Barto: Reinforcement Learning: An Introduction

### **First-visit Monte Carlo policy evaluation**

Initialize:

R. S. Sutton and A. G. Barto: Reinforcement Learning: An Introduction

 $\pi \leftarrow \text{policy to be evaluated}$ 

- $V \leftarrow$  an arbitrary state-value function
- $Returns(s) \leftarrow an \text{ empty list, for all } s \in S$

#### Repeat forever:

- (a) Generate an episode using  $\pi$
- (b) For each state s appearing in the episode:
  - $R \leftarrow$  return following the first occurrence of sAppend R to Returns(s)
    - $V(s) \leftarrow \operatorname{average}(Returns(s))$

### Blackjack example

- □ *Object:* Have your card sum be greater than the dealers without exceeding 21.
- **States** (200 of them):
  - current sum (12-21)
  - dealer's showing card (ace-10)



- do I have a useable ace?
- **Reward:** +1 for winning, 0 for a draw, -1 for losing
- Actions: stick (stop receiving cards), hit (receive another card)
- Delicy: Stick if my sum is 20 or 21, else hit

1

4

2





11

#### **Convergence of MC Control**

Greedified policy meets conditions for policy improvement:

$$Q^{\pi_k}(s, \pi_{k+1}(s)) = Q^{\pi_k}(s, \arg\max_a Q^{\pi_k}(s, a))$$
$$= \max_a Q^{\pi_k}(s, a)$$
$$\ge Q^{\pi_k}(s, \pi_k(s))$$
$$= V^{\pi_k}(s)$$

- $\square$  And thus must be  $\ge \pi_k$  by the policy improvement theorem
- □ This assumes exploring starts and infinite number of episodes for MC policy evaluation
- **T**o solve the latter:
  - update only to a given level of performance
  - alternate between evaluation and improvement per episode

R. S. Sutton and A. G. Barto: Reinforcement Learning: An Introduction

# **Monte Carlo Exploring Starts**

Initialize, for all $s \in S$ , $a \in \mathcal{A}(s)$ : $Q(s, a) \leftarrow \text{arbitrary}$ $\pi(s) \leftarrow \text{arbitrary}$	Fixed point is optimal policy $\pi^*$
$Returns(s, a) \leftarrow \text{empty list}$	Proof is open question
Repeat forever:	
(a) Generate an episode using exploring starts and $\pi$	
(b) For each pair $s, a$ appearing in the episode:	
$R \leftarrow$ return following the first occurrence of $s, a$	
Append $R$ to $Returns(s, a)$	
$Q(s, a) \leftarrow \operatorname{average}(Returns(s, a))$	
(c) For each $s$ in the episode:	
$\pi(s) \leftarrow \arg\max_a Q(s,a)$	

R. S. Sutton and A. G. Barto: Reinforcement Learning: An Introduction



#### Learning about $\pi$ while following $\pi'$

Suppose we have  $n_s$  returns,  $R_i(s)$ , from state s, each with probability  $p_i(s)$  of being generated by  $\pi$  and probability  $p_i'(s)$  of being generated by  $\pi'$ . Then we can estimate

$$V^{\pi}(s) \approx \frac{\sum_{i=1}^{n_s} \frac{p_i(s)}{p'_i(s)} R_i(s)}{\sum_{i=1}^{n_s} \frac{p_i(s)}{p'_i(s)}}$$

which depends on the environmental probabilities  $p_i(s)$  and  $p'_i(s)$ . However,

$$p_i(s_t) = \prod_{k=t}^{T_i(s)-1} \pi(s_k, a_k) \mathcal{P}_{s_k s_{k+1}}^{a_k}$$

and

$$\frac{p_i(s_t)}{p_i'(s_t)} = \frac{\prod_{k=t}^{T_i(s)-1} \pi(s_k, a_k) \mathcal{P}^{a_k}_{s_k s_{k+1}}}{\prod_{k=t}^{T_i(s)-1} \pi'(s_k, a_k) \mathcal{P}^{a_k}_{s_k s_{k+1}}} = \prod_{k=t}^{T_i(s)-1} \frac{\pi(s_k, a_k)}{\pi'(s_k, a_k)}$$

Thus the weight needed,  $p_i(s)/p'_i(s)$ , depends only on the two policies and not at all on the environmental dynamics.

R. S. Sutton and A. G. Barto: Reinforcement Learning: An Introduction

#### 17

#### **Off-policy MC control**



#### R. S. Sutton and A. G. Barto: Reinforcement Learning: An Introduction

**Incremental Implementation** 

**MC** can be implemented incrementally

saves memory

Compute the weighted average of each return

$$V_{n} = \frac{\sum_{k=1}^{n} w_{k} R_{k}}{\sum_{k=1}^{n} w_{k}} \qquad V_{n+1} = V_{n} + \frac{W_{n+1}}{W_{n+1}} \Big[ R_{n+1} - V_{n} - V$$

non-incremental

incremental equivalent



R. S. Sutton and A. G. Barto: Reinforcement Learning: An Introduction

18

## Summary

□ MC has several advantages over DP:

- Can learn directly from interaction with environment
- No need for full models
- No need to learn about ALL states
- Less harm by Markovian violations (later in book)
- MC methods provide an alternate policy evaluation process
- □ One issue to watch for: maintaining sufficient exploration
  - exploring starts, soft policies
- □ No bootstrapping (as opposed to DP)

R. S. Sutton and A. G. Barto: Reinforcement Learning: An Introduction

21