#### Chapter 4: Dynamic Programming

Objectives of this chapter:

- Overview of a collection of classical solution methods for MDPs known as dynamic programming (DP)
- Show how DP can be used to compute value functions, and hence, optimal policies
- Discuss efficiency and utility of DP

## Policy Evaluation Policy Evaluation: for a given policy $\pi$ , compute the state-value function $V^{\pi}$

Recall: State - value function for policy  $\pi$ :

$$V^{\pi}(s) = E_{\pi} \left\{ R_{t} \mid s_{t} = s \right\} = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} \mid s_{t} = s \right\}$$

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

- a system of |S| simultaneous linear equations



3

#### **Iterative Policy Evaluation**

Input  $\pi$ , the policy to be evaluated Initialize V(s) = 0, for all  $s \in S^+$ Repeat  $\Delta \leftarrow 0$ For each  $s \in S$ :  $v \leftarrow V(s)$   $V(s) \leftarrow \sum_a \pi(s, a) \sum_{s'} \mathcal{P}^a_{ss'} \left[ \mathcal{R}^a_{ss'} + \gamma V(s') \right]$   $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ until  $\Delta < \theta$  (a small positive number) Output  $V \approx V^{\pi}$ 

4





#### Policy Improvement

Suppose we have computed  $V^{\pi}$  for a deterministic policy  $\pi$ .

For a given state s, would it be better to do an action  $a \neq \pi(s)$ ?

The value of doing a in state s is :

$$Q^{\pi}(s, a) = E_{\pi} \left\{ r_{t+1} + \gamma V^{\pi}(s_{t+1}) \middle| s_{t} = s, a_{t} = a \right\}$$
$$= \sum_{s'} P_{ss'}^{a} \left[ R_{ss'}^{a} + \gamma V^{\pi}(s') \right]$$

It is better to switch to action a for state s if and only if

$$Q^{\pi}(s,a) > V^{\pi}(s)$$

7

#### Policy Improvement Cont.

Do this for all states to get a new policy  $\pi'$  that is greedy with respect to  $V^{\pi}$ :

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

8

# **Policy Improvement Cont.** What if $V^{\pi'} = V^{\pi}$ ? i.e., for all $s \in S$ , $V^{\pi'}(s) = \max_{a} \sum_{s'} P^{a}_{ss'} [R^{a}_{ss'} + \gamma V^{\pi}(s')]$ ? But this is the Bellman Optimality Equation. So $V^{\pi'} = V^*$ and both $\pi$ and $\pi'$ are optimal policies.





### Jack's Car Rental

- \$10 for each car rented (must be available when request received
- Two locations, maximum of 20 cars at each
- Cars returned and requested randomly
  - Poisson distribution, *n* returns/requests with prob  $\frac{\lambda}{n!}e^{-\lambda}$
  - 1st location: average requests = 3, average returns = 2
  - 2nd location: average requests = 4, average returns = 2

12

- Can move up to 5 cars between locations overnight
- States, Actions, Rewards?
- Transition probabilities?





#### Value Iteration

Recall the **full policy evaluation backup**:

$$V_{k+1}(s) \leftarrow \sum_{a} \pi(s, a) \sum_{s'} P^a_{ss'} \Big[ R^a_{ss'} + \gamma V_k(s') \Big]$$

Here is the **full value iteration backup**:

$$V_{k+1}(s) \leftarrow \max_{a} \sum_{s'} P^a_{ss'} \Big[ R^a_{ss'} + \gamma V_k(s') \Big]$$

15

Value Iteration Cont.

Initialize V arbitrarily, e.g., V(s) = 0, for all  $s \in S^+$ Repeat  $\Delta \leftarrow 0$ For each  $s \in S$ :  $v \leftarrow V(s)$   $V(s) \leftarrow \max_a \sum_{s'} \mathcal{P}^a_{ss'} [\mathcal{R}^a_{ss'} + \gamma V(s')]$   $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ until  $\Delta < \theta$  (a small positive number) Output a deterministic policy,  $\pi$ , such that  $\pi(s) = \arg \max_a \sum_{s'} \mathcal{P}^a_{ss'} [\mathcal{R}^a_{ss'} + \gamma V(s')]$ 

16

15



- Gambler can repeatedly bet \$ on a coin flip
- Heads he wins his stake, tails he loses it
- Initial capital: \$1, \$2, ..., \$99
- Gambler wins if his capital becomes \$100; loses if it becomes \$0
- Coin is unfair
  - Heads (gambler wins) with probability p = 0.4

17

• States, Actions, Rewards?

#### Herd Management

- You are a consultant to a farmer managing a herd of cows
- Herd consists of 5 kinds of cows:
  - Young
  - Milking
  - Breeding
  - Old
  - Sick
- Number of each kind is the State
- Number sold of each kind is the Action
- Cows transition from one kind to another
- Young cows can be born

#### Gambler's Problem Solution



#### Asynchronous DP

- All the DP methods described so far require exhaustive sweeps of the entire state set.
- Asynchronous DP does not use sweeps. Instead it works like this:
  - Repeat until convergence criterion is met:
    - Pick a state at random and apply the appropriate backup
- Still need lots of computation, but does not get locked into hopelessly long sweeps
- Can you select states to backup intelligently? YES: an agent's experience can act as a guide.

20

#### **Generalized Policy Iteration**

**Generalized Policy Iteration** (GPI):

any interaction of policy evaluation and policy improvement, independent of their granularity.



#### Efficiency of DP

- To find an optimal policy is polynomial in the number of states...
- BUT, the number of states is often astronomical, e.g., often growing exponentially with the number of state variables (what Bellman called "the curse of dimensionality").
- In practice, classical DP can be applied to problems with a few millions of states.
- Asynchronous DP can be applied to larger problems, and appropriate for parallel computation.

22

 It is surprisingly easy to come up with MDPs for which DP methods are not practical.

#### Summary

- Policy evaluation: backups without a max
- Policy improvement: form a greedy policy, if only locally
- Policy iteration: alternate the above two processes
- Value iteration: backups with a max
- Full backups (to be contrasted later with sample backups)
- Generalized Policy Iteration (GPI)
- Asynchronous DP: a way to avoid exhaustive sweeps
- **Bootstrapping**: updating estimates based on other estimates