

# Dynamic Programming

- Collection of algorithms used to compute optimal policies given a perfect model of the environment as an MDP
- DP provides an essential foundation for understanding reinforcement learning algorithms

# MDP recap

- Assume environment is a finite MDP (states, actions, rewards)
- Dynamics given by set probabilities  $p(s', r | s, a)$
- Bellman optimally equations:

$$v_*(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

$$= \max_a \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma v_*(s') \right], \text{ or}$$

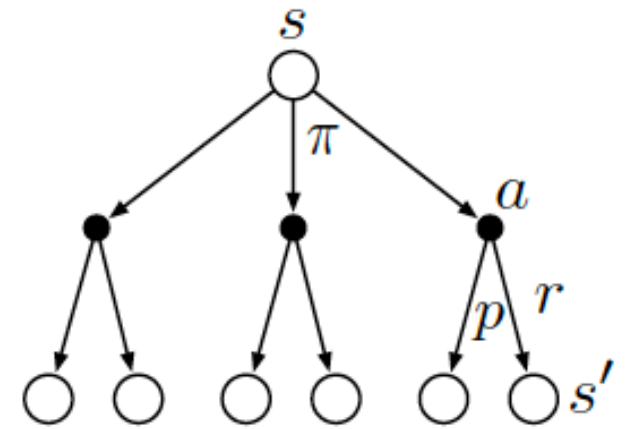
$$q_*(s, a) = \mathbb{E} \left[ R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right]$$

$$= \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma \max_{a'} q_*(s', a') \right],$$

# Policy Evaluation (Prediction)

- State-value function expected update rule (based on Bellman equation):

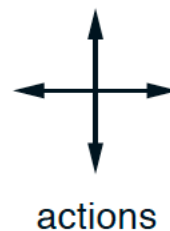
$$\begin{aligned} v_{k+1}(s) &\doteq \mathbb{E}_{\pi}[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_k(s')], \end{aligned}$$



Backup diagram for  $v_{\pi}$

- Bootstrapping: update estimates on the basis of estimates.
- Convergence (update no longer results in value changes) guaranteed as number of samples  $\rightarrow \infty$

# Gridworld example



	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

$R_t = -1$   
on all transitions

$v_k$  for the  
random policy

greedy policy  
w.r.t.  $v_k$

$k = 0$

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

	↕	↕	↕
↕	↕	↕	↕
↕	↕	↕	↕
↕	↕	↕	

random  
policy

$k = 3$

0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0

	←	←	↘
↑	↖	↘	↓
↑	↖	↘	↓
↙	→	→	

$k = 1$

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

	←	↕	↕
↑	↕	↕	↕
↕	↕	↕	↓
↕	↕	→	

$k = 10$

0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0

	←	←	↘
↑	↖	↘	↓
↑	↖	↘	↓
↙	→	→	

optimal  
policy

$k = 2$

0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7
-2.0	-2.0	-1.7	0.0

	←	←	↕
↑	↖	↕	↓
↑	↕	↘	↓
↕	→	→	

$k = \infty$

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0

	←	←	↘
↑	↖	↘	↓
↑	↖	↘	↓
↙	→	→	

## Iterative Policy Evaluation, for estimating $V \approx v_\pi$

Input  $\pi$ , the policy to be evaluated

Algorithm parameter: a small threshold  $\theta > 0$  determining accuracy of estimation

Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop:

$\Delta \leftarrow 0$

Loop for each  $s \in \mathcal{S}$ :

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until  $\Delta < \theta$

# Policy Improvement

- State-action value function  $q_{\pi}(s, a) \doteq \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a]$   
 $= \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_{\pi}(s')].$
- Policy improvement theorem: policy  $\pi'$  must be as good as, or better than,  $\pi$ .  
 $v_{\pi'}(s) \geq v_{\pi}(s).$
- Greedy policy  $\pi'(s) \doteq \arg\max_a q_{\pi}(s, a)$   
 $= \arg\max_a \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a]$   
 $= \arg\max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_{\pi}(s')],$

# Policy Iteration

Policy Iteration (using iterative policy evaluation) for estimating  $\pi \approx \pi_*$

1. Initialization

$V(s) \in \mathbb{R}$  and  $\pi(s) \in \mathcal{A}(s)$  arbitrarily for all  $s \in \mathcal{S}$

2. Policy Evaluation

Loop:

$\Delta \leftarrow 0$

Loop for each  $s \in \mathcal{S}$ :

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until  $\Delta < \theta$  (a small positive number determining the accuracy of estimation)

3. Policy Improvement

*policy-stable*  $\leftarrow$  *true*

For each  $s \in \mathcal{S}$ :

*old-action*  $\leftarrow \pi(s)$

$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

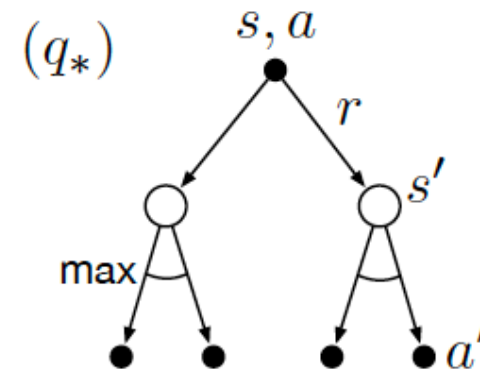
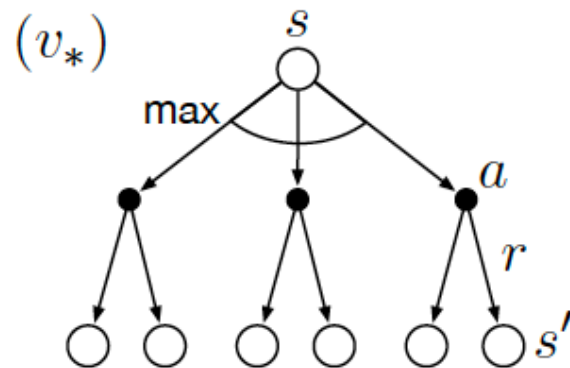
If *old-action*  $\neq \pi(s)$ , then *policy-stable*  $\leftarrow$  *false*

If *policy-stable*, then stop and return  $V \approx v_*$  and  $\pi \approx \pi_*$ ; else go to 2

# Value Iteration

- Truncate policy evaluation to a single sweep/update
- Update rule that combines policy improvement and evaluation (Bellman optimality equation):

$$\begin{aligned}v_{k+1}(s) &\doteq \max_a \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s, A_t = a] \\&= \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_k(s')],\end{aligned}$$





## Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold  $\theta > 0$  determining accuracy of estimation

Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop:

|  $\Delta \leftarrow 0$

| Loop for each  $s \in \mathcal{S}$ :

|      $v \leftarrow V(s)$

|      $V(s) \leftarrow \max_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$

|      $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until  $\Delta < \theta$

Output a deterministic policy,  $\pi \approx \pi_*$ , such that

$$\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$$

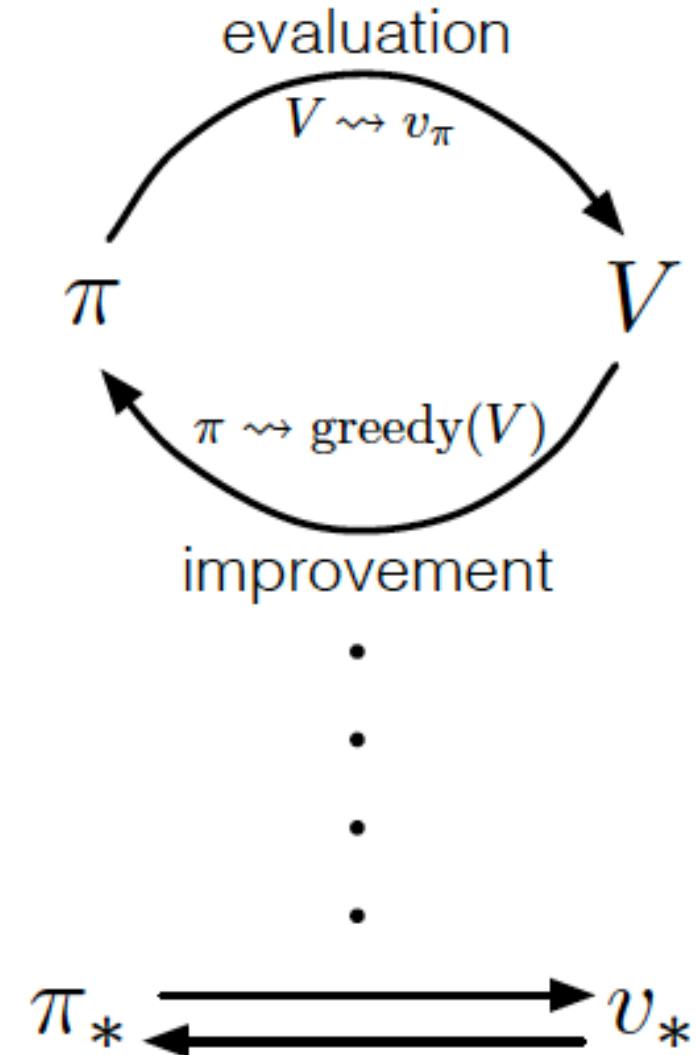
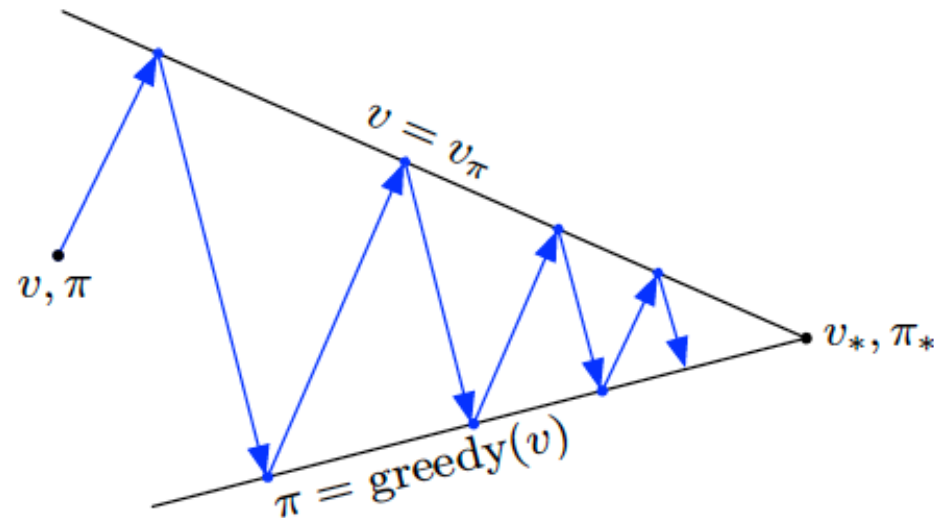
# Asynchronous Dynamic Programming

- Synchronous involves a complete sweep of the state space every iteration
- Asynchronous is preferred for problems with large state spaces where the values of states are updated in any order
- Allows flexibility in selecting states to update and the policy can be improved more frequently
- To converge correctly, asynchronous algorithms must continue to update the values of all the states
- Allows real-time interaction

# Generalized Policy Iteration

Two interacting processes:

- Policy evaluation - makes the value function consistent with the current policy
- Policy improvement - make the policy greedy with respect to the current value function



# Efficiency of Dynamic Programming

- Worst case time DP methods take to find an optimal policy is polynomial in the number of states ( $n$ ) and actions ( $k$ ).
- Faster convergence with good initial value functions or policies.
- However, direct search of policy space is exponential in state space
- Classical DP algorithms are of limited use in RL because of known dynamics assumption and great computational expense