# Reinforcement Learning

Fateme Taroodi

# Planning with DP

Fateme Taroodi

**From last lectures**

$$v_\pi(s) = \mathbb{E}_\pi[G_t|S_t = s] \; = \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a)[r + \gamma v_\pi(s')].$$

Fateme Taroodi

## Dynamic programming for planning MDPs

In reinforcement learning, we want to use dynamic programming to solve MDPs. So given an MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ and a policy $\pi$:

First, we want to find the value function $v_\pi$ for that policy:

- This is done by **policy evaluation** (the prediction problem)

Then, when we're able to evaluate the policy, we want find the best policy $v_*$ (the control problem). This is done with two strategies:

1. **Policy iteration**
2. **Value iteration**

Fateme Taroodi

Policy Evaluation is to compute the state-value $V_\pi$ for a given policy $\pi$:

$$V_{t+1}(s) = \mathbb{E}_\pi[r + \gamma V_t(s') | S_t = s] = \sum_a \pi(a|s) \sum_{s',r} P(s',r|s,a)(r + \gamma V_t(s'))$$

Fateme Taroodi

Recall: Bellman equation for $v_\pi$ is system of linear equations

$$v_\pi(s_1) = \sum_a \pi(a|s_1) \sum_{s',r} p(s',r|s_1,a)\left[r + \gamma v_\pi(s')\right]$$

$$v_\pi(s_2) = \sum_a \pi(a|s_2) \sum_{s',r} p(s',r|s_2,a)\left[r + \gamma v_\pi(s')\right]$$

$$\vdots$$

$$v_\pi(s_n) = \sum_a \pi(a|s_n) \sum_{s',r} p(s',r|s_n,a)\left[r + \gamma v_\pi(s')\right]$$

Could use this for policy evaluation step, but expensive

- Gauss elimination (de facto standard) has time complexity $O(n^3)$

Fateme Taroodi

## Iterative Policy Evaluation

We can use Bellman equation as operator to *iteratively* compute $v_\pi$:

- Initialise $v_0(s) = 0$

- Then repeatedly perform updates:

$$v_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[r + \gamma v_k(s')\right] \quad \text{for all } s \in \mathcal{S}$$

- Sequence converges to fixed point $v_\pi$, so stop when no more changes to $v_k$

*Updating estimates based on other estimates is called* *bootstrapping*

Fateme Taroodi

Input $\pi$, the policy to be evaluated

Initialize an array $V(s) = 0$, for all $s \in \mathcal{S}^+$

Repeat

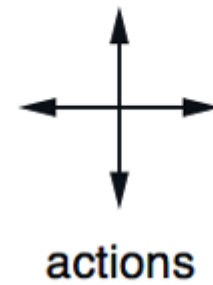$\quad \Delta \leftarrow 0$

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

$\quad\quad v \leftarrow V(s)$

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

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

until $\Delta < \theta$ (a small positive number)

Output $V \approx v_\pi$

Fateme Taroodi

# Example: Policy evaluation



$$R_t = -1$$
on all transitions

actions

## Policy evaluation

$k = 0$

| 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 |

$k = 1$

| 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 |

$k = 2$

| 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 |

Fateme Taroodi

# Policy evaluation

$k = 3$

| 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 |

$k = 10$

| 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 |

$k = \infty$

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

Fateme Taroodi

# Policy evaluation + Greedy Improvement



random policy

Fateme Taroodi

# Policy evaluation + Greedy Improvement

$k = 3$

| 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 |

$k = 10$

| 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

$k = \infty$

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

Fateme Taroodi

▶ The example already shows we can use evaluation to then improve our policy

▶ In fact, just being greedy with respect to the values of the random policy sufficed! (That is not true in general)

## Algorithm

Iterate, using

انتخاب حریصانه ← اینتخاب حریصانه

$$\forall s: \quad \pi_{new}(s) = \text{argmax}_a \, q_\pi(s, a)$$

$$= \text{argmax}_a \, \mathbb{E}\left[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a\right]$$

Then, evaluate $\pi_{new}$ and repeat

▶ Claim: One can show that $v_{\pi_{new}}(s) \geq v_\pi(s)$, for all $s$

$$v_{\pi_{new}} \geq v_\pi$$

Policy **Evaluation**:
Estimate $v_\pi$

Policy **Improvement**:
Generate better $\pi'$

Fateme Taroodi

$$\pi_0 \xrightarrow{\text{evaluation}} V_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluation}} V_{\pi_1} \xrightarrow{\text{improve}} \pi_2 \xrightarrow{\text{evaluation}} \cdots \xrightarrow{\text{improve}} \pi_* \xrightarrow{\text{evaluation}} V_*$$

$k = 3$

| 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 |

$k = 10$

| 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 |

$k = \infty$

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

optimal policy

## 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* $\leftarrow \pi(s)$
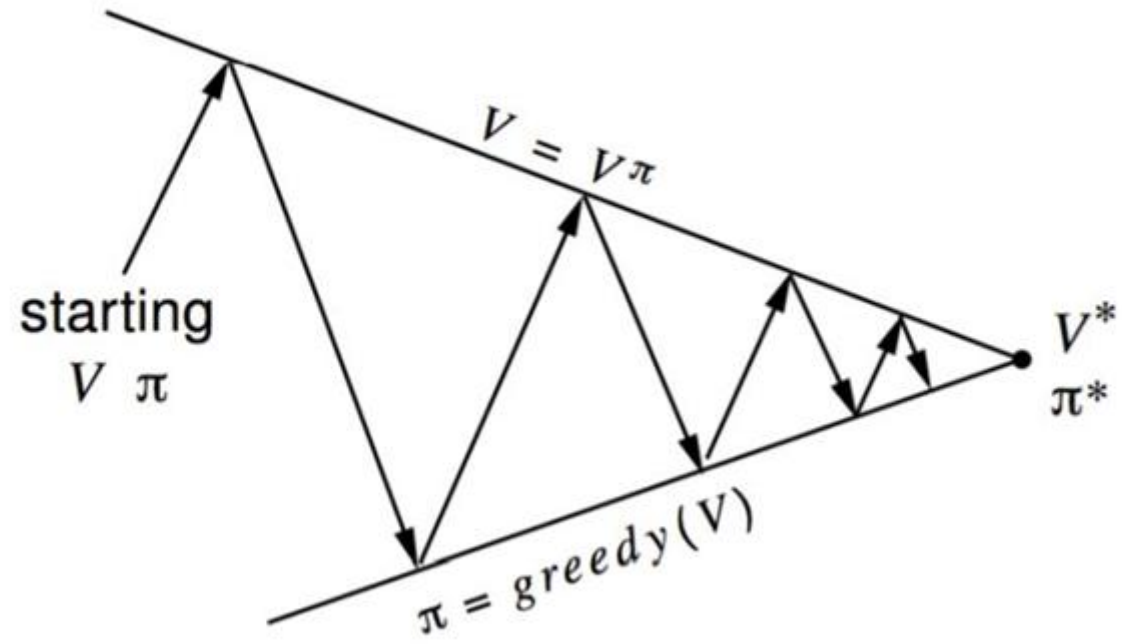       $\pi(s) \leftarrow \arg\max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$
       If *old-action* $\neq \pi(s)$, then *policy-stable* $\leftarrow$ *false*
   If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Fateme Taroodi

# Value Iteration

Fateme Taroodi

starting $V$ $\pi$

$V = V\pi$

$\pi = greedy(V)$

$V^*$
$\pi^*$

Fateme Taroodi

## 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 $s \in \mathcal{S}$:

           $v \leftarrow V(s)$

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

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

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

   <span style="color:red">**Big loop**</span>
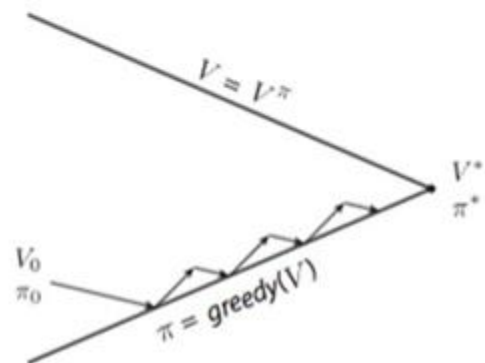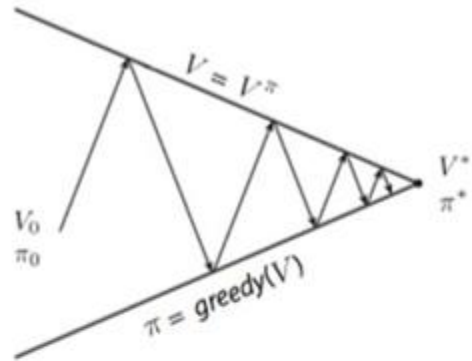
3. Policy Improvement

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

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

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

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

       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
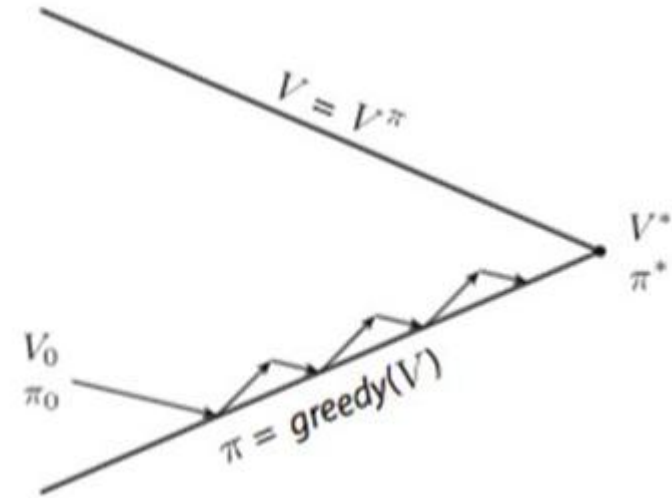
Fateme Taroodi

**Can we make it faster?!**

**Value Iteration**

Fateme Taroodi

## Value Iteration

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

$$= \max_a \sum_{s',r} p(s',r \mid s,a)\big[r + \gamma v_k(s')\big],$$

## 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 $V(terminal) = 0$

Loop:
$\quad | \quad \Delta \leftarrow 0$
$\quad | \quad$ Loop for each $s \in \mathcal{S}$:
$\quad | \qquad v \leftarrow V(s)$
$\quad | \qquad V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$
$\quad | \qquad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that
$\quad \pi(s) = \arg\max_a \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$

Fateme Taroodi

# Generalized Policy Iteration

letting policy-evaluation and policy improvement processes interact, independent of the granularity and other details of the two processes



Fateme Taroodi