

Probabilistic Decision Making

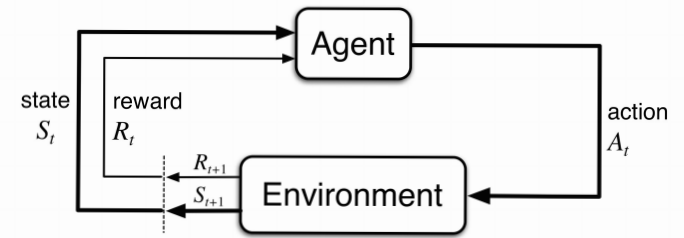
Lecture 15

Reinforcement Learning 2

Robert Peharz

Institute of Machine Learning and Neural Computation
Graz University of Technology
Winter Term 2025/26

Recap



Markov Decision Process (MDP) $(\mathcal{S}, \mathcal{A}, p)$

state space \mathcal{S} , action space \mathcal{A} , dynamics $p(s', r | s, a)$

Policy $\pi(a|s)$

Discounted Return

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}, \quad \gamma \in [0, 1]$$

Value Function

$$V_{\pi}(s) := \mathbb{E}_{\pi}[G_t | S_t = s] = \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s\right]$$

Recap

Bellman Expectation Equation



Image: wikipedia.org

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

$$\underline{v_{\pi}}(s) = r(s) + \gamma \sum_{s'} p(s'|s) \underline{v_{\pi}}(s') \quad \forall s \in \mathcal{S}$$

$$p(s'|s) = \sum_{r,a} p(s', r | s, a) \pi(a|s)$$

Matrix-Vector form:

$$\underline{v} = \begin{pmatrix} v_{\pi}(s_1) \\ v_{\pi}(s_2) \\ \vdots \\ v_{\pi}(s_N) \end{pmatrix} \quad \underline{r} = \begin{pmatrix} r(s_1) \\ r(s_2) \\ \vdots \\ r(s_N) \end{pmatrix} \quad P = \begin{pmatrix} p(s_1|s_1) & \dots & p(s_N|s_1) \\ \vdots & \ddots & \vdots \\ p(s_1|s_N) & \dots & p(s_N|s_N) \end{pmatrix}$$

$$\underline{v} = \underline{r} + \gamma P \underline{v}$$

Analytic solution: $\underline{v_{\pi}} = (I - \gamma P)^{-1} \underline{r}$

Recap

Iterative Policy Evaluation

Bellman equation
(interpreted as update rule)

Iterative Policy Evaluation

- initialize $v(s)$ arbitrarily (e.g. all 0)
- repeat
 - for $s \in \mathcal{S}$ do $V_{\text{new}}(s) \leftarrow r(s) + \gamma \sum_{s'} p(s'|s) v(s')$
 - if $\forall s: V_{\text{new}}(s) \approx v(s) \rightarrow$ break
 - $V \leftarrow V_{\text{new}}$

Random walk policy $\hat{\pi}(a|s) = 0.25$ ($\leftarrow, \rightarrow, \uparrow, \downarrow$)

$k=0$

| | | | |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

$k=1$

| | | | |
|----|----|----|----|
| 0 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 0 |

$k=2$

| | | | |
|------|------|------|------|
| 0 | -1.7 | -2 | -2 |
| -1.7 | -2 | -2 | -2 |
| -2 | -2 | -2 | -1.7 |
| -2 | -2 | -1.7 | 0 |

...

$k=10$

| | | | |
|------|------|------|------|
| 0 | -6.1 | -8.3 | -8.9 |
| -6.1 | -7.7 | -8.4 | -8.3 |
| -8.3 | -8.4 | -7.7 | -6.1 |
| -8.9 | -8.3 | -6.1 | 0 |

...

$k=199$

| | | | |
|-------|-------|-------|-------|
| 0 | -13.8 | -19.6 | -21.6 |
| -13.8 | -17.7 | -19.6 | -19.6 |
| -19.6 | -19.6 | -17.7 | -13.8 |
| -21.6 | -19.6 | -13.8 | 0 |

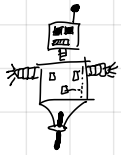
$k=200$

| | | | |
|-------|-------|-------|-------|
| 0 | -13.8 | -19.6 | -21.6 |
| -13.8 | -17.7 | -19.6 | -19.6 |
| -19.6 | -19.6 | -17.7 | -13.8 |
| -21.6 | -19.6 | -13.8 | 0 |

$\approx V_{\pi}$

How to Improve Policies?

Policy:



$$\pi(a|s) \equiv 0.25$$

| | | | |
|---|---|---|---|
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |

Value function v_π

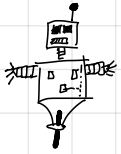
| | | | |
|-------|-------|-------|-------|
| 0 | -13.8 | -19.6 | -21.6 |
| -13.8 | -17.7 | -19.6 | -19.6 |
| -19.6 | -19.6 | -17.7 | -13.8 |
| -21.6 | -19.6 | -13.8 | 0 |

Let's help the robot !

- the value function tells us "how good" a state s is under current policy π
- assume any state s
- assume you could "hack" the agent's policy for one time step, i.e. you can select the agent's next action
- after that, the agent resumes its policy
- which action do you pick?

How to Improve Policies?

Policy:



$$\pi(a|s) \equiv 0.25$$

| | | | |
|---|---|---|---|
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |

Value function V_π

| | | | |
|-------|-------|-------|-------|
| 0 | -13.8 | -19.6 | -21.6 |
| -13.8 | -17.7 | -19.6 | -19.6 |
| -19.6 | -19.6 | -17.7 | -13.8 |
| -21.6 | -19.6 | -13.8 | 0 |

better: larger value

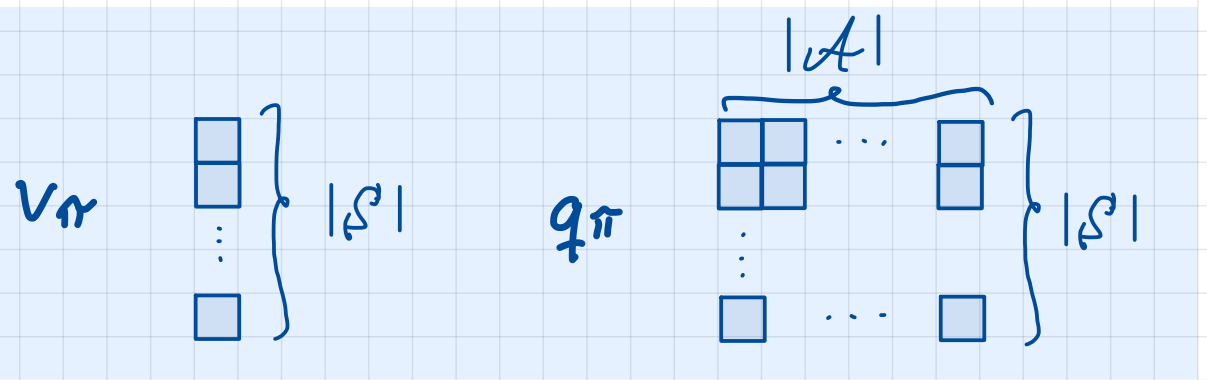
- we should move the agent to better states

| | | | |
|---|---|---|---|
| ↕ | ← | ← | ↙ |
| ↑ | ↖ | ↙ | ↓ |
| ↑ | ↗ | ↘ | ↓ |
| ↘ | → | → | ↕ |

(multiple arrows denote equally good actions)

- we actually just came up with
a new policy ! (greedy policy)
- is it better than the old one?

The q Function



- to formalize greedy policies we introduce the q -function
- we have previously defined the value function

$$V_\pi(s) = \mathbb{E}_\pi [G_t \mid S_t = s] = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

the q -function is the value function where we first pick an action of our choice:

$$q_\pi(s, a) := \mathbb{E}_\pi [G_t \mid S_t = s, A_t = a] = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

also called state-action value function.

How to Compute q ?

Method 1 v and q are related via

$$v_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(s, a) \quad (1)$$

$$q_{\pi}(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\pi}(s') \quad (2)$$

Thus, (2) allows to compute q_{π} from v_{π}

Method 2 Use Bellman equations in q_{π}

$$q_{\pi}(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) \sum_{a'} \pi(a'|s') q_{\pi}(s', a')$$

- Solve system of $|S| \times |A|$ equations, or
- do iterative Bellman updates

Greedy Policy

- let π be any policy and q_π its q -function
- define a new deterministic policy π' :

$$\pi'(a|s) := \begin{cases} 1 & \text{if } a = \operatorname{argmax}_{a'} q_\pi(s, a') \\ 0 & \text{otherwise} \end{cases}$$

(pick arbitrary $a' \in \operatorname{argmax}$ if maximizer is not unique)

$\pi'(a|s)$ is called the greedy policy w.r.t. q_π (v_π)

Example: Grid World

Value function

| | | | |
|-------|-------|-------|-------|
| 0 | -13.8 | -19.6 | -21.6 |
| -13.8 | -17.7 | -19.6 | -19.6 |
| -19.6 | -19.6 | -17.7 | -13.8 |
| -21.6 | -19.6 | -13.8 | 0 |

Greedy Policy

| | | | |
|----|---|---|----|
| ↕↔ | ← | ← | ↖ |
| ↑ | ↖ | ↖ | ↓ |
| ↑ | ↗ | ↗ | ↓ |
| ↗ | → | → | ↕↔ |

(multiple arrows : non-unique argmax)

is it better than the old one?

Policy Improvement Theorem

Let $\hat{\pi}$ be any policy and $\hat{\pi}'$ the (a) greedy policy w.r.t. $v_{\hat{\pi}}(q_{\hat{\pi}})$.
Then

$$v_{\hat{\pi}}(s) \leq v_{\hat{\pi}'}(s) \quad \forall s \in \mathcal{S}$$

Recall:

$$v_{\hat{\pi}}(s) = \sum_a \hat{\pi}(a|s) q_{\hat{\pi}}(s, a)$$

$$q_{\hat{\pi}}(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\hat{\pi}}(s')$$

$$v_{\hat{\pi}}(s) = \sum_a \hat{\pi}(a|s) q_{\hat{\pi}}(s, a)$$

$$\leq \max_a q_{\hat{\pi}}(s, a)$$

// use $\hat{\pi}'$ once, then $\hat{\pi}$

$$= \max_a r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\hat{\pi}}(s')$$

// repeat argument

$$\leq \max_a r(s, a) + \gamma \sum_{s'} p(s'|s, a) \max_{a'} q_{\hat{\pi}}(s', a')$$

// use $\hat{\pi}'$ twice, then $\hat{\pi}$

$$= \max_a r(s, a) + \gamma \sum_{s'} p(s'|s, a) \left(\max_{a'} r(s', a') + \gamma \sum_{s''} p(s''|s', a') v_{\hat{\pi}}(s'') \right)$$

$$\leq \dots \leq v_{\hat{\pi}'}(s)$$

// by induction, $\hat{\pi}'$ is at least as good as $\hat{\pi}$



Policy Iteration

We might apply the greedy policy iteratively, leading to policy iteration

- initialize $\hat{\pi}_0$, $K \leftarrow 0$

- repeat

- (1) compute $q_{\hat{\pi}_K}$ // e.g. with iterative policy evaluation

- (2) let $\hat{\pi}_{K+1}(a|s) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} q_{\hat{\pi}_K}(a', s) \\ 0 & \text{otherwise} \end{cases}$ // greedy policy

- if $q_{\hat{\pi}_{K+1}} \approx q_{\hat{\pi}_K} \rightarrow \text{break}$

- $K \leftarrow K + 1$

Example: Grid World

\hat{V}_K

$k=0$

| | | | |
|---|---|---|---|
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |

$$V_{\pi_K} \quad (q_{\pi_K} = r(s,a) + \gamma E_{S_{t+1}}[V_{\pi_K}(S_{t+1})])$$

| | | | |
|-------|-------|-------|-------|
| 0 | -13.8 | -19.6 | -21.6 |
| -13.8 | -17.7 | -19.6 | -19.6 |
| -19.6 | -19.6 | -17.7 | -13.8 |
| -21.6 | -19.6 | -13.8 | 0 |

$k=1$

| | | | |
|---|---|---|---|
| ↕ | ← | ← | ← |
| ↑ | ↑ | ↓ | ↓ |
| ↑ | ↑ | → | ↓ |
| ↑ | → | → | ↑ |

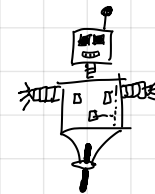
| | | | |
|----|----|----|----|
| 0 | -1 | -2 | -3 |
| -1 | -2 | -3 | -2 |
| -2 | -3 | -2 | -1 |
| -3 | -2 | -1 | 0 |

$k=2$

| | | | |
|---|---|---|---|
| ↕ | ← | ← | ← |
| ↑ | ↑ | ↓ | ↓ |
| ↑ | ↑ | → | ↓ |
| ↑ | → | → | ↑ |

| | | | |
|----|----|----|----|
| 0 | -1 | -2 | -3 |
| -1 | -2 | -3 | -2 |
| -2 | -3 | -2 | -1 |
| -3 | -2 | -1 | 0 |

Converged!



This policy seems globally optimal!

Does this always happen?

Optimal Policies

A policy $\hat{\pi}^*$ is optimal if $v_{\hat{\pi}^*}(s) \geq v_{\hat{\pi}}(s)$, $\forall \hat{\pi}, s \in \mathcal{S}$

Theorem: Policy Iteration Converges to Optimal Policy

For any initial policy $\hat{\pi}_0$ policy iteration converges to an optimal policy

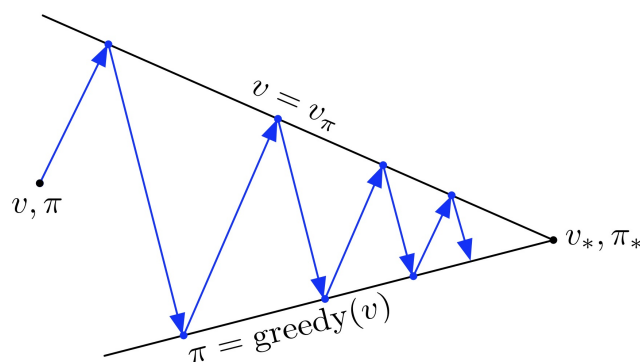
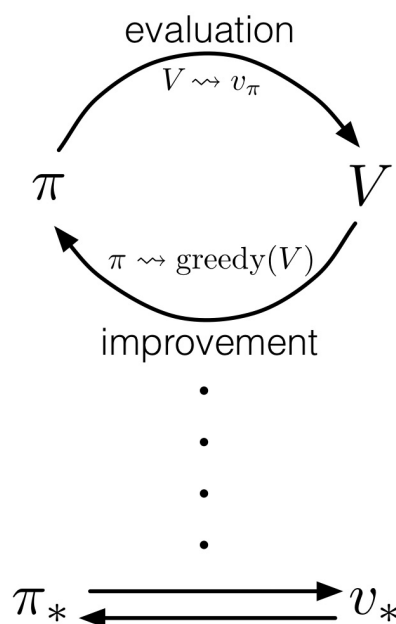
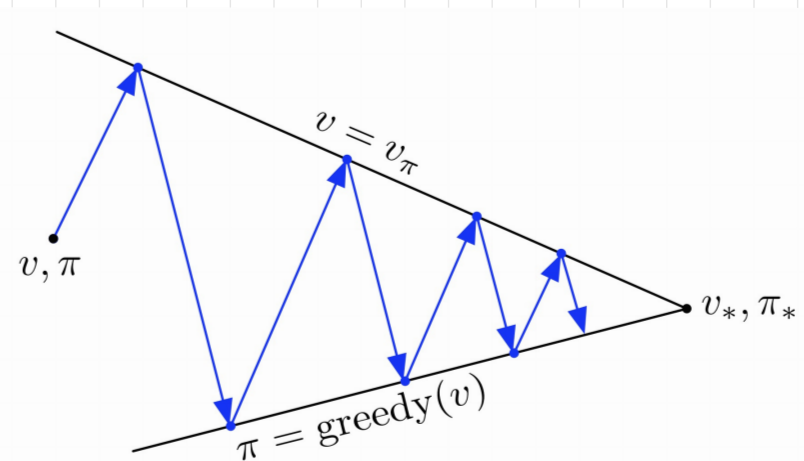
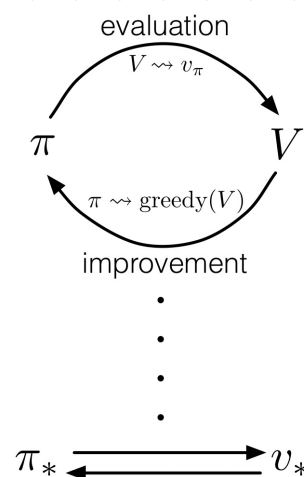


Image: Sutton & Barto

Efficiency of Policy Iteration?

Image: Sutton & Barto



- policy evaluation is expensive (in particular closed form)
- Iterative Policy Evaluation only converges for $k \rightarrow \infty$
- before that, intermediate results are typically not even a value function of any policy
- how many iterations of Iterative Policy Evaluation do we really need?

Grid World

$$\hat{\pi}(a|s) \equiv 0.25$$

Iterative Policy Evaluation

$k=0$

| | | | |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

$k=1$

| | | | |
|----|----|----|----|
| 0 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 0 |

$k=2$

| | | | |
|------|------|------|------|
| 0 | -1.7 | -2 | -2 |
| -1.7 | -2 | -2 | -2 |
| -2 | -2 | -2 | -1.7 |
| -2 | -2 | -1.7 | 0 |

$k=3$

| | | | |
|------|------|------|------|
| 0 | -2.4 | -2.9 | -3 |
| -2.4 | -2.8 | -3 | -2.9 |
| -2.9 | -3 | -2.8 | -2.4 |
| -3 | -2.9 | -2.4 | 0 |

...

$k=10$

| | | | |
|------|------|------|------|
| 0 | -6.1 | -8.3 | -8.9 |
| -6.1 | -7.7 | -8.4 | -8.3 |
| -8.3 | -8.4 | -7.7 | -6.1 |
| -8.9 | -8.3 | -6.1 | 0 |

...

$k=\infty$

| | | | |
|-------|-------|-------|-------|
| 0 | -13.8 | -19.6 | -21.6 |
| -13.8 | -17.7 | -19.6 | -19.6 |
| -19.6 | -19.6 | -17.7 | -13.8 |
| -21.6 | -19.6 | -13.8 | 0 |

Greedy Policy

| | | | |
|---|---|---|---|
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↕ |

| | | | |
|---|---|---|---|
| ↕ | ← | ↕ | ↕ |
| ↑ | ↕ | ↕ | ↕ |
| ↕ | ↕ | ↕ | ↓ |
| ↕ | ↕ | → | ↕ |

| | | | |
|---|---|---|---|
| ↕ | ← | ← | ↕ |
| ↑ | ↕ | ↕ | ↓ |
| ↑ | ↕ | ↕ | ↓ |
| ↕ | → | → | ↕ |

| | | | |
|---|---|---|---|
| ↕ | ← | ← | ↕ |
| ↑ | ↕ | ↕ | ↓ |
| ↑ | ↕ | ↕ | ↓ |
| ↕ | → | → | ↕ |

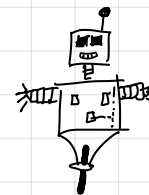
...

| | | | |
|---|---|---|---|
| ↕ | ← | ← | ↕ |
| ↑ | ↕ | ↕ | ↓ |
| ↑ | ↕ | ↕ | ↓ |
| ↕ | → | → | ↕ |

...

| | | | |
|---|---|---|---|
| ↕ | ← | ← | ↕ |
| ↑ | ↕ | ↕ | ↓ |
| ↑ | ↕ | ↕ | ↓ |
| ↕ | → | → | ↕ |

Optimal policy!



Policy Iteration with Truncated Policy Evaluation

initialize $\hat{\pi}, v$

repeat until convergence

repeat for K steps

$$\forall s: v(s) \leftarrow r(s) + \gamma \sum_a \hat{\pi}(a|s) \sum_{s'} p(s'|s, a) v(s')$$

$$\hat{\pi} \leftarrow \text{greedy}(v)$$

- this algorithm also converges, for any K , to an optimal policy — no need for exact policy evaluation
- which K will work best, in order to minimize total compute? Unknown, open problem...
- special case $K=1$ is called value iteration