

Reinforcement Learning

Lecture 4

Planning: Control

Greedy Policy, Policy Improvement Theorem

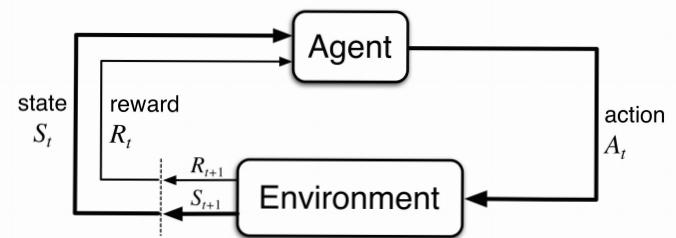
Policy Iteration

Robert Peharz

Institute of Theoretical Computer Science
Graz University of Technology

Winter Term 2023/24

Recap



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

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

Policy $\tilde{\pi}(a|s)$

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

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

Recap

Bellman Expectation Equation

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



Image: wikipedia.org

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

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

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

Recap

Iterative Policy Evaluation

Bellman equation
(interpreted as update rule)

Iterative Policy Evaluation

- initialize $v(s)$ arbitrarily (e.g. all 0)
- repeat

- for $s \in S$ 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 |

$\approx V_{\pi}$

$k=200$

| | | | |
|-------|-------|-------|-------|
| 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 |

Recap

Contraction (Contraction mapping)

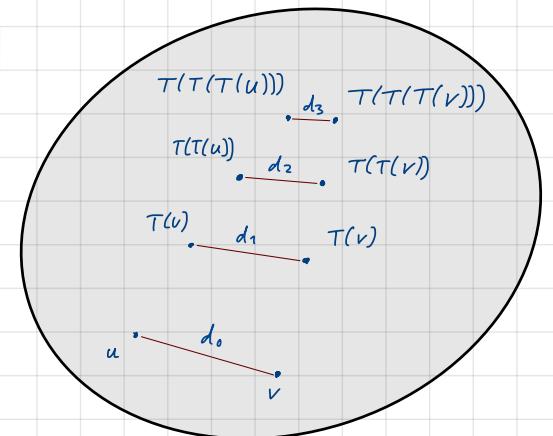
Let V be any vector space. An L -Lipschitz function $T: V \rightarrow V$ with $L < 1$ is called a contraction (L -contraction).

Fixed point

Let V be some vector space and let $T: V \rightarrow V$ be a function.

A vector $v \in V$ is called a fixed point of T if

$v = T(v)$, i.e. if T maps v onto itself.



Banach's Fixed Point Theorem

Let V be some vector space and let $T: V \rightarrow V$ be a contraction. Then T has a unique fixed point $v^* \in V$.

Moreover, for any $v_0 \in V$, the sequence defined via

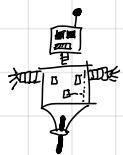
$$v_n = T(v_{n-1})$$

converges to v^* , i.e. $v_n \rightarrow v^*$ as $n \rightarrow \infty$.

→ theoretical
underpinning
for many RL
algorithms!
(needed for your homework)

How to Improve Policies?

Policy:



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

| | | | |
|---|---|---|---|
| ↔ | ↔ | ↔ | ↔ |
| ↔ | ↔ | ↔ | ↔ |
| ↔ | ↔ | ↔ | ↔ |
| ↔ | ↔ | ↔ | ↔ |

Value function $v_{\hat{\pi}}$

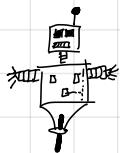
| | | | |
|-------|-------|-------|-------|
| 0 | -13.8 | -19.6 | -21.6 |
| -13.8 | -17.7 | -19.6 | -19.6 |
| -19.6 | -19.6 | -17.7 | -13.8 |
| -21.6 | -19.6 | -13.8 | 0 |

Let's help the robot! 🚀

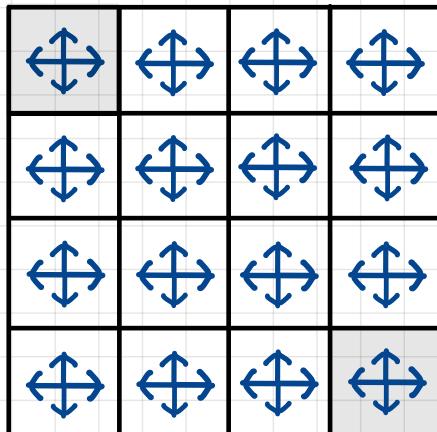
- the value function tells us "how good" a state s is under current policy $\hat{\pi}$
- assume any state s
- assume you could "hack" the agent's policy for one time step, i.e. you can select the agent's next action
- after that, the agent resumes its policy
- which action do you pick?

How to Improve Policies?

Policy:



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

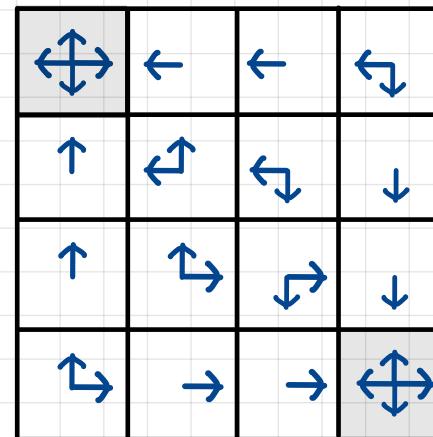


Value function v_{π}

| | | | |
|-------|-------|-------|-------|
| 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

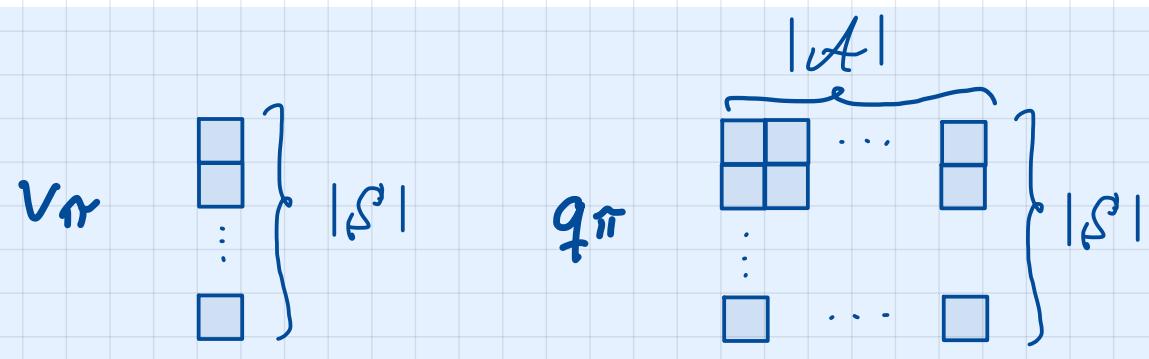
- we should move the agent to better states



(multiple arrows denote equally good actions)

- we actually just came up with a new policy! **(greedy policy)**
- is it always 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.

Relation Between v_{π} and q_{π}

We can express v_{π} via q_{π} :

$$\underline{v_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s] = \mathbb{E}_{A_t} \left[\underbrace{\mathbb{E}_{\pi}[G_t | s, A_t]}_{q(s, A_t)} \mid S_t = s \right] = \sum_a \hat{\pi}(a | s) q_{\pi}(s, a)}$$

(1)

We can express q_{π} via v_{π} :

$$\underline{q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a]}$$

$$= \mathbb{E}_{\pi} \left[\underbrace{R_{t+1} + \gamma G_{t+1}}_{= G_t} \mid S_t = s, A_t = a \right]$$

expected reward

$$= r(s, a) + \gamma \mathbb{E}_{S_{t+1}} \left[\underbrace{\mathbb{E}_{\pi}[G_{t+1} | S_t = s, A_t = a, S_{t+1}]}_{= v_{\pi}(S_{t+1})} \mid S_t = s, A_t = a \right]$$

Markov!

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

(2)

Bellman Expectation Equation in q

Plugging (1) into (2) yields the Bellman equation in q :

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

How to compute q_{π} :

Option 1: Compute V_{π} and let $q_{\pi}(s, a) = r(s, a) + \gamma \sum_{s'} p(s' | s, a) V_{\pi}(s')$

Option 2: Directly solve Bellman equation in q

- closed form (matrix inversion)
- iterating Bellman equation with arbitrary initial q_0

Greedy Policy

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

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

(pick arbitrary $a' \in \arg \max$
if maximizer is not unique)

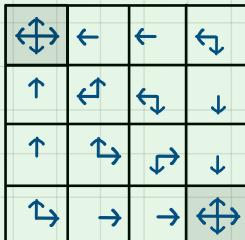
$\pi'(a|s)$ is called the greedy policy w.r.t. q_π (v_π)

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

$$V_{\tilde{\pi}}(s) \leq V_{\tilde{\pi}'}(s) \quad \forall s \in S$$

Recall:

$$V_{\tilde{\pi}}(s) = \sum_a \tilde{\pi}(a|s) q_{\tilde{\pi}}(s, a)$$

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

$$V_{\tilde{\pi}}(s) = \sum_a \tilde{\pi}(a|s) q_{\tilde{\pi}}(s, a)$$

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

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

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

// repeat argument

$$\leq \max_a r(s, a) + \gamma \sum_{s'} p(s'|s, a) \max_{a'} q_{\tilde{\pi}}(s', a') // \text{use } \tilde{\pi}' \text{ twice, then } \tilde{\pi}$$

$$\leq \max_a r(s, a) + \gamma \sum_{s'} p(s'|s, a) \left(\max_{a''} r(s', a'') + \gamma \sum_{s''} p(s''|s', a'') V_{\tilde{\pi}}(s'') \right)$$

$$\leq \dots \leq V_{\tilde{\pi}'}(s)$$

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

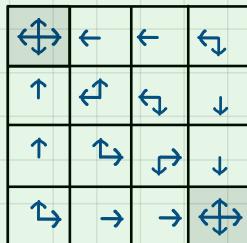


Example: Grid World

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?

Yes!

A greedy policy is always at least as good as the original policy.

- Note that a greedy policy is deterministic
- Thus, deterministic policies are "enough" *
- We can now iterate the "greedyfication" → policy iteration

* For planning (known MDP) – for general RL, stochastic policies will be important for exploration.

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

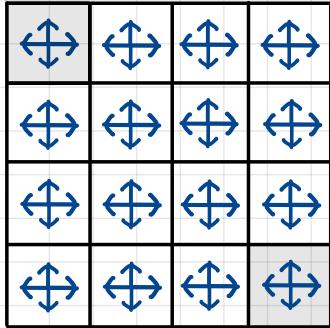
- if $q_{\hat{\pi}_{K+1}} \approx q_{\hat{\pi}_K} \rightarrow \text{break}$

- $K \leftarrow K + 1$

Example: Grid World

π_k

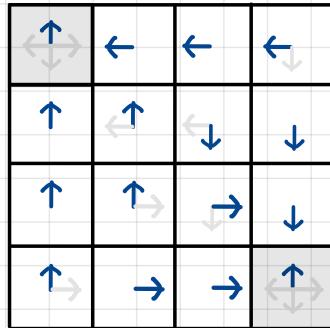
$k=0$



$$V_{\pi_k} \quad (q_{\pi_k} = r(s,a) + \gamma \mathbb{E}_{s_{t+1}}[v_{\pi_k}(s_{t+1})])$$

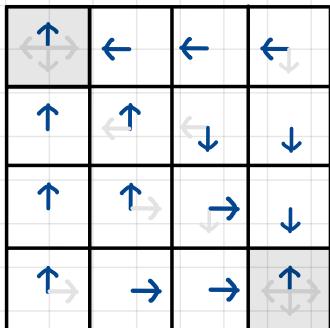
| | | | |
|-------|-------|-------|-------|
| 0 | -13.8 | -19.6 | -21.6 |
| -13.8 | -17.7 | -19.6 | -19.6 |
| -19.6 | -19.6 | -17.7 | -13.8 |
| -21.6 | -19.6 | -13.8 | 0 |

$k=1$



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

$k=2$

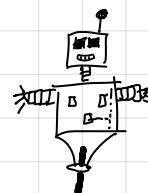


| | | | |
|----|----|----|----|
| 0 | -1 | -2 | -3 |
| -1 | -2 | -3 | -2 |
| -2 | -3 | -2 | -1 |
| -3 | -2 | -1 | 0 |

This policy seems globally optimal!

Does this always happen?

Converged!



Optimal Value Function

- What does it mean for a policy to be optimal?

Intuitively, a policy π^* is optimal if $V_{\pi^*}(s) \geq V_{\pi}(s)$, $\forall \pi, s \in S$

- Does an optimal policy always exist?

(Not evident: They might be incomparable, meaning for any two π_1, π_2 ,
there could be $s, s' \in S$, such that $V_{\pi_1}(s) > V_{\pi_2}(s)$
 $V_{\pi_1}(s') < V_{\pi_2}(s')$)

- First, let's define optimal value functions

$$V^*(s) := \max_{\pi} V_{\pi}(s)$$

$$q^*(s, a) := \max_{\pi} q_{\pi}(s, a)$$

— upper bound for all V_{π}, q_{π} ,
since \max_{π} is taken separately
for all s, a .

- If there was a π^* with $V_{\pi^*} = V^*$ then π^* would be optimal.

Bellman Optimality Equation

$$\hat{\pi}'(a|s) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} q_{\pi}(s, a') \\ 0 & \text{otherwise} \end{cases}$$

When $V_{\pi} = V_{\pi'}$, then policy iteration has converged.

$$\begin{aligned} \underline{V_{\pi}(s)} &= E_{\pi}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s] \\ &= \underline{E_{\pi'}[R_{t+1} + \gamma V_{\pi'}(S_{t+1}) | S_t = s]} = V_{\pi'}(s) \end{aligned}$$

$$\begin{aligned} \underline{V_{\pi}(s)} &= E_{\pi'}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s] \\ &= \sum_a \hat{\pi}'(a|s) \left[r(s, a) + \gamma \sum_{s'} p(s'|s, a) \overset{q_{\pi}(s, a)}{\circ} V_{\pi}(s') \right] \\ &= \max_a r(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s') \end{aligned}$$

Bellman Optimality Equation

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

Holds, when
policy iteration
has converged

Bellman Optimality Equation (in q)

Policy iteration converged: $\hat{q}_{\pi} = \hat{q}_{\pi'}$

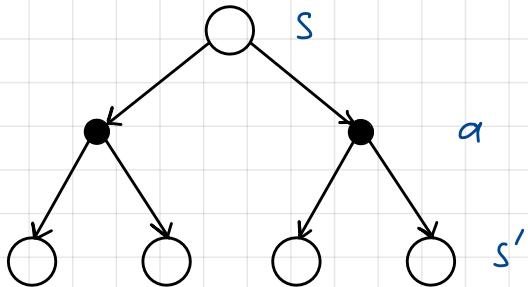
$$\begin{aligned}\underline{q_{\pi}(s, a)} &= r(s, a) + \gamma \sum_{s'} p(s' | s, a) \sum_{a'} \hat{\pi}(a' | s') q_{\pi}(s', a') \\ &= \underline{r(s, a) + \gamma \sum_{s'} p(s' | s, a) \sum_{a'} \hat{\pi}'(a' | s') q_{\pi'}(s', a')} = q_{\pi'}(s, a)\end{aligned}$$

Bellman Optimality Equation (in q)

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

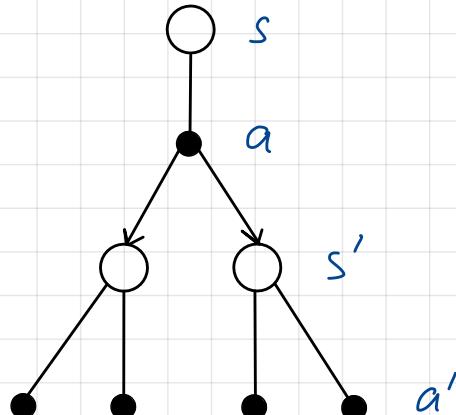
Backup Diagrams: Graphical Descriptions of RL Algorithms

Bellman Expectation Equation in v



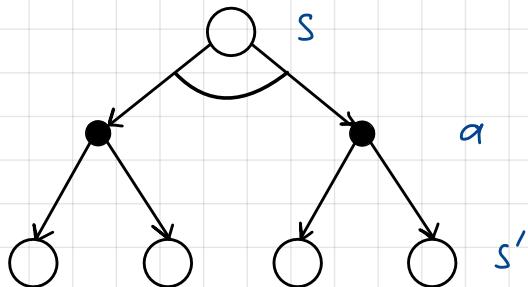
$$v_{\pi}(s) = r(s) + \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) v_{\pi}(s')$$

Bellman Expectation Equation in q



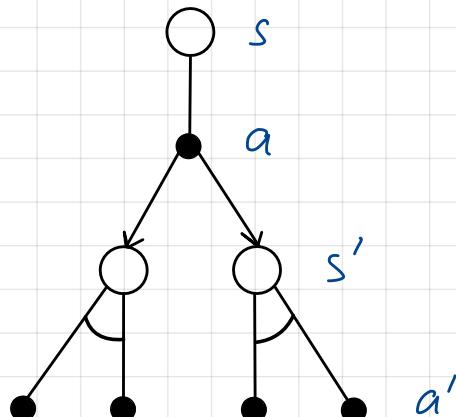
$$q(s,a) = r(s,a) + \gamma \sum_{s'} p(s'|s,a) \sum_{a'} \pi(a'|s') q(s',a')$$

Bellman Optimality Equation in v



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

Bellman Optimality Equation in q



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

Theorem: Global Optimal Policy

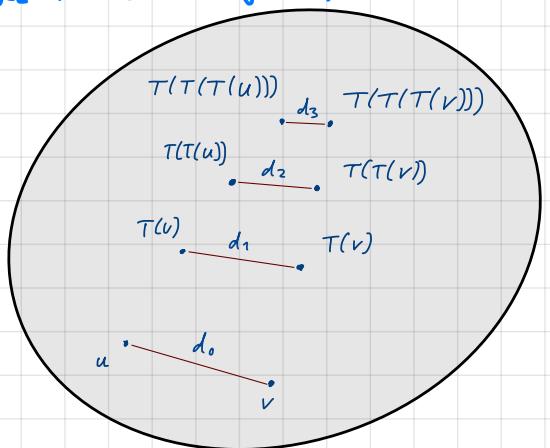
For any initial policy π_0 , policy iteration converges to an optimal policy π^* with value function v^*

Proof sketch

(see Szepesvári, page 78)

— (why will it converge?)

- When Policy Iteration has converged, the Bellman Optimality Equation holds
- Show that Bellman Optimality Equation is a contraction
- Hence, it converges to a unique fixed point
- Show that this fixed point is an upper bound of v^*
- But the fixed point is also $v_{\pi'}$ of the learned π' why?
- Hence v^* must be the value function of π'



Schematic of Policy Iteration

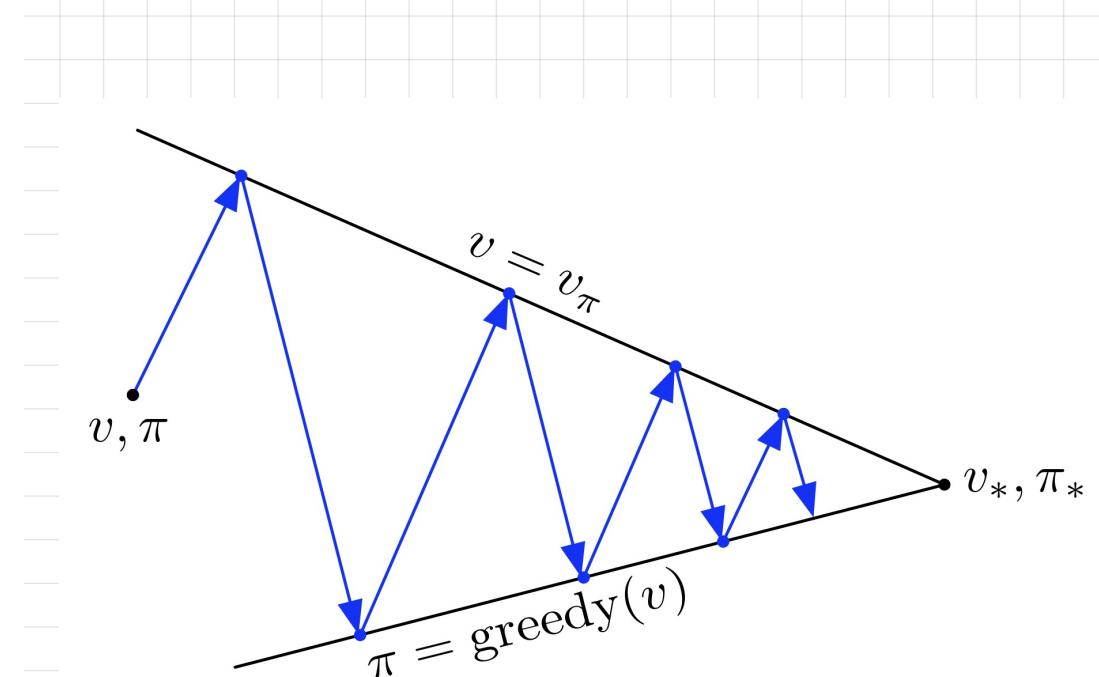
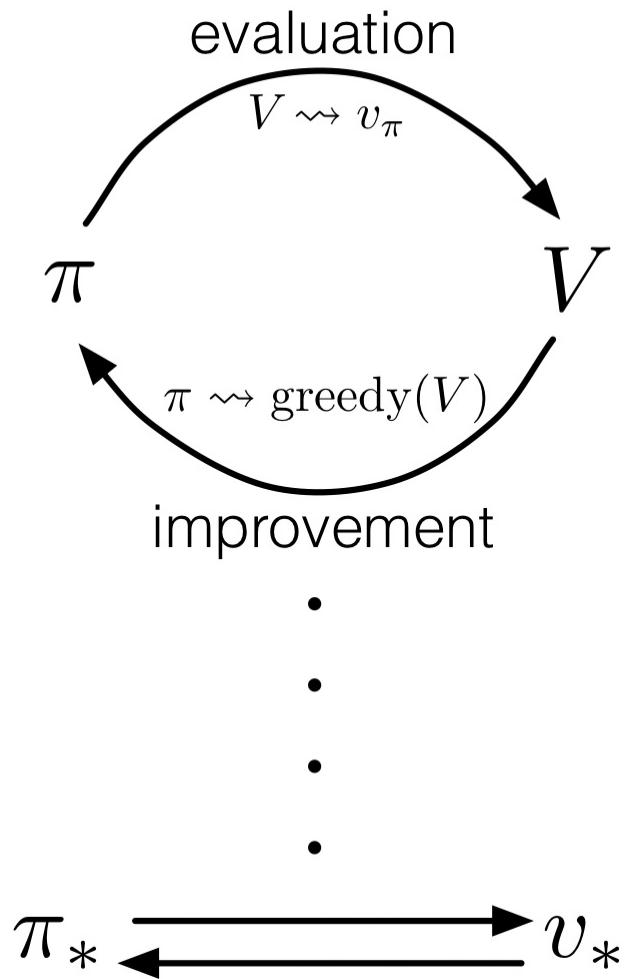


Image: Sutton & Barto

Car Rental

(Sutton & Barto, p. 81)

- a manager is in charge of two car rental stations
- if a customer requests a car at a station:
 - if car available \rightarrow rented out for 10 €
 - otherwise \rightarrow 0 €
- returned cars get available the next day
- maximal 20 cars at each station (additional cars returned elsewhere)
- numbers of requested and returned cars are Poisson distributed

$$p(n) = \frac{\lambda^n}{n!} e^{-\lambda}$$

station 1: $\lambda_{\text{request}} = 3$ $\lambda_{\text{return}} = 3$

station 2: $\lambda_{\text{request}} = 4$ $\lambda_{\text{return}} = 2$

(i.e., typically, more cars are requested at station 2 than returned)

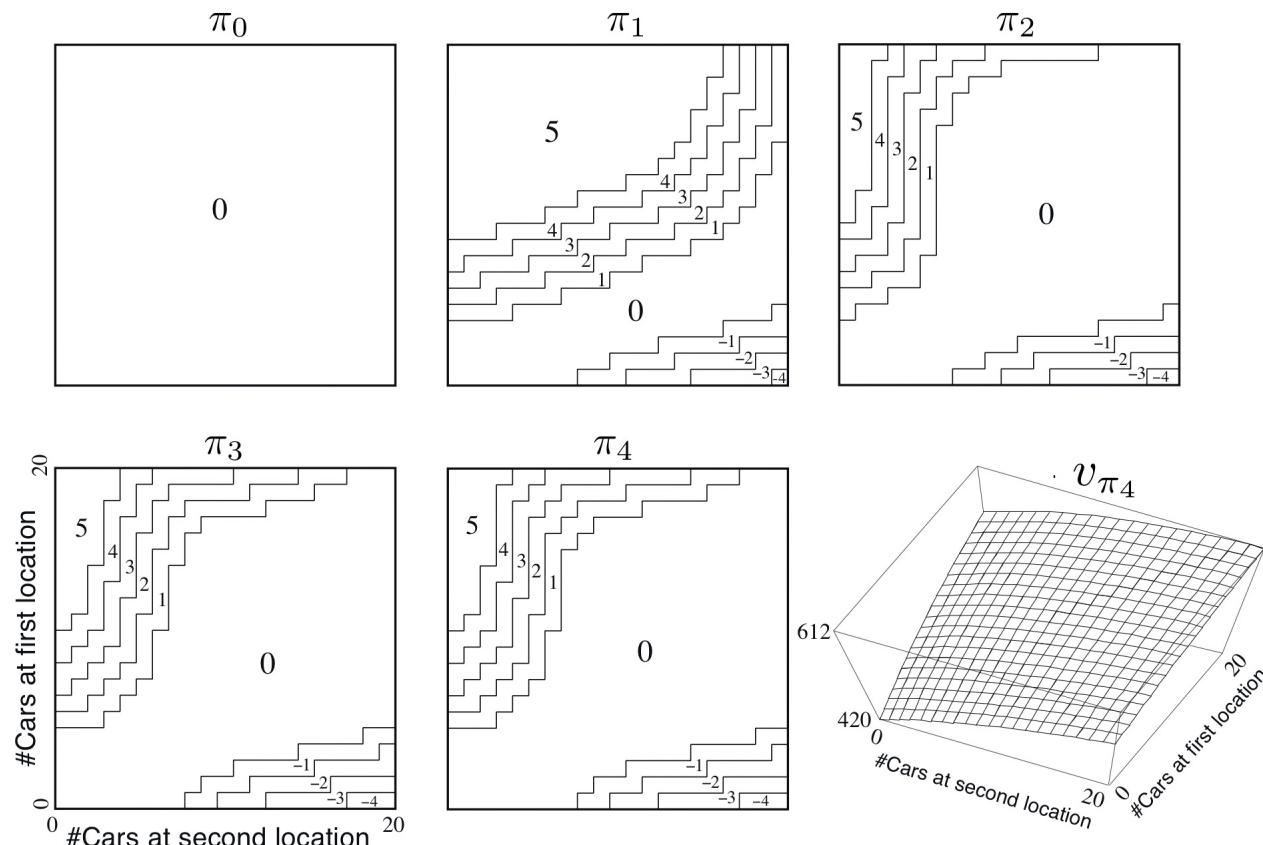
- manager can transfer up to 5 cars between stations each day for the cost of 2 €

Car Rental cont'd

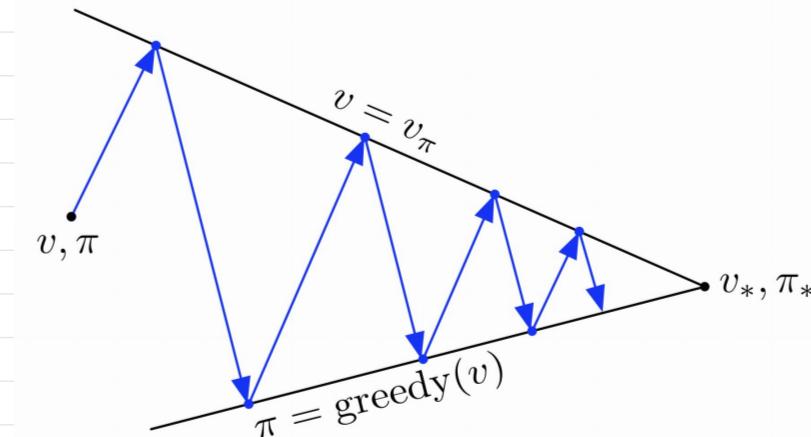
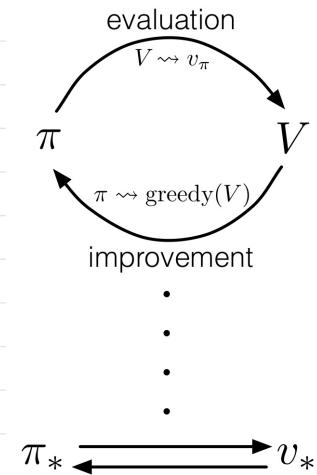
(Sutton & Barto, p. 81)

- state space $S = \{0, 1, 2, \dots, 20\} \times \{0, 1, 2, \dots, 20\}$
- action space $A = \{-5, -4, -3, -2, -1, 0, 1, 2, 3, 4, 5\}$
number of cars moved from station 1 to station 2
- $p(s'|s, a)$ and $r(s, a)$ governed by Poisson distributions

Policy iteration:



Efficiency of Policy Iteration?



- policy evaluation is expensive (in particular closed form)
- Iterative Policy Evaluation only converges for $n \rightarrow \infty$

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

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

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

| | | | |
|---|---|---|---|
| ↗ | ↗ | ↗ | ↗ |
| ↗ | ↗ | ↗ | ↗ |
| ↗ | ↗ | ↗ | ↗ |
| ↗ | ↗ | ↗ | ↗ |

| | | | |
|---|---|---|---|
| ↗ | ↖ | ↗ | ↗ |
| ↑ | ↗ | ↗ | ↗ |
| ↗ | ↗ | ↗ | ↓ |
| ↗ | ↗ | → | ↗ |

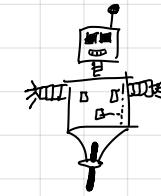
| | | | |
|---|---|---|---|
| ↗ | ↖ | ↖ | ↗ |
| ↑ | ↖ | ↗ | ↓ |
| ↑ | ↗ | ↗ | ↓ |
| ↗ | → | ↗ | ↗ |

| | | | |
|---|---|---|---|
| ↗ | ↖ | ↖ | ↖ |
| ↑ | ↖ | ↖ | ↓ |
| ↑ | ↗ | ↖ | ↓ |
| ↗ | → | ↗ | ↗ |

...

| | | | |
|---|---|---|---|
| ↗ | ↖ | ↖ | ↖ |
| ↑ | ↖ | ↖ | ↓ |
| ↑ | ↗ | ↖ | ↓ |
| ↗ | → | ↗ | ↗ |

Optimal policy!



Policy Iteration with Truncated Policy Evaluation

initialize $\hat{\pi}$, v

repeat until convergence

repeat for K steps

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

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

Value Iteration

Special case for $k=1$

initialize $\hat{\pi}$, v

repeat until convergence

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

$$\hat{\pi}'(a|s) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} \underbrace{r(s,a') + \gamma \sum_{s'} p(s'|s,a) v(s')}_{q(s,a)} \\ 0 & \text{o.w.} \end{cases}$$

$$\begin{aligned} v^{new}(s) &= \underbrace{r(s)}_{\check{r}(s)} + \gamma \sum_a \hat{\pi}'(a|s) \sum_{s'} p(s'|s,a) v^{old}(s') \\ &= \sum_a \hat{\pi}'(a|s) \left(r(s,a) + \gamma \sum_{s'} p(s'|s,a) v^{old}(s') \right) \\ &= \max_a r(s,a) + \gamma \sum_{s'} p(s'|s,a) v^{old}(s') \end{aligned}$$

This is the iterating the Bellman Optimality Equation !

Theorem: Global Optimal Policy

showed that this converges to v^* and $\hat{\pi}^*$.

Value Iteration

Policy Iteration with
Truncated Policy Evaluation

Policy Iteration

1 number k of policy evaluation steps $\dots \infty$

Converges to v^* and π^* for all K !

Optimal K , with least computation? Not easy to determine.

Asynchronous Updates

...
repeat until convergence

repeat for K steps

$\forall s:$

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

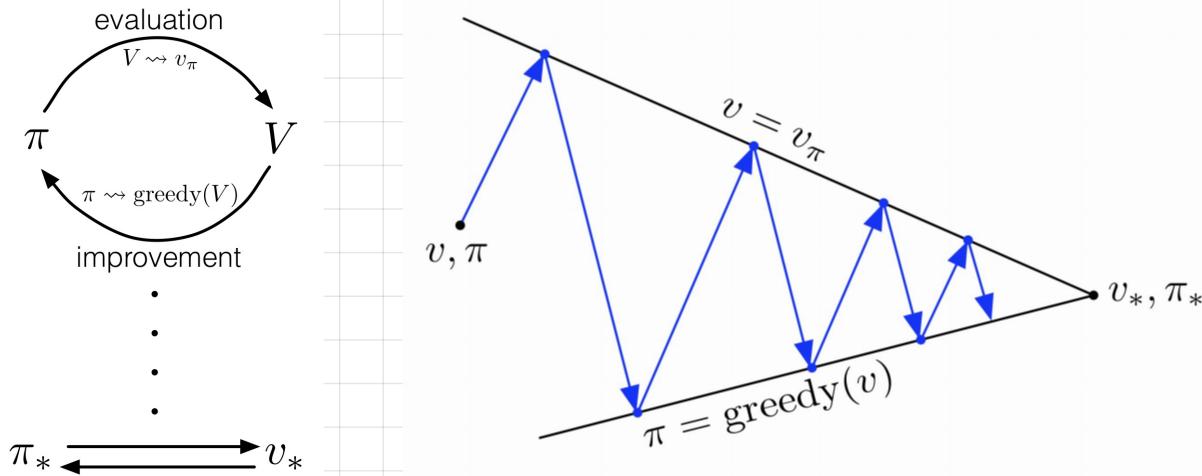
Do we have to always update all state values?

- If number of states is large, one iteration takes very long
- We could use updated values immediately

No.

- We can update subsets of S (also single s), in any order
- Convergence to v^* , π^* if
 - all states are updated eventually
 - infinitely often

Generalised Policy Iteration



- Truncated policy evaluation & asynchronous updates
approximate exact policy evaluation
- Generalized policy iteration: any scheme which moves
 - v closer to v_π
 - π closer to $\text{greedy}(v)$

Typically converges to v^*, π^* (shown via contraction arguments)