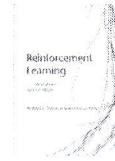


Outline and References

- Outline
 - ▶ Policy Evaluation
 - ▶ Policy Improvement
 - ▶ Policy Iteration
 - ▶ Generalized Policy Iteration
 - ▶ Efficiency of DP

- References
 - ▶ Reinforcement Learning: An Introduction [RL Chapter 4]
 - ▶ Fundamentals of Reinforcement Learning (Coursera)

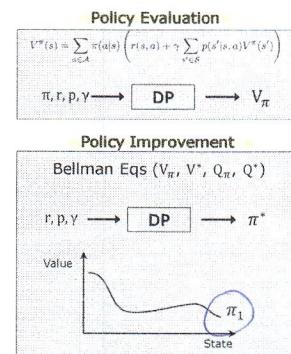


Machine Learning - Daniele Lolliano

2

Why Dynamic Programming?

- To solve an MDP we need to find an optimal policy
- Unfortunately we cannot use a brute-force approach:
 - ▶ $|\mathcal{A}|^{|S|}$ deterministic policies to evaluate
 - ▶ $|S|$ linear equations to solve for each policy
- Dynamic Programming (DP) is a method that allow to solve a complex problem by breaking it down into simpler sub-problems in a recursive manner
- We will see how to use DP to solve an MDP thanks to the Bellman Equations



because of size of the state space we cannot proceed using linear systems solver

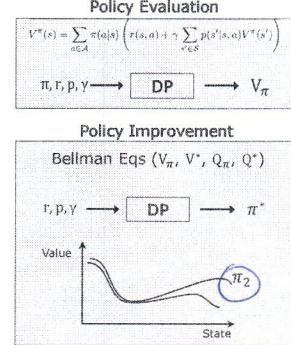
thanks to the dynamic programming we'll improve iteratively the policy (from $\pi_1 \rightarrow \pi_2 \rightarrow \dots \rightarrow \pi^*$)

Machine Learning - Daniele Lolliano

3

Why Dynamic Programming?

- To solve an MDP we need to find an optimal policy
- Unfortunately we cannot use a brute-force approach:
 - ▶ $|\mathcal{A}|^{|S|}$ deterministic policies to evaluate
 - ▶ $|S|$ linear equations to solve for each policy
- Dynamic Programming (DP) is a method that allow to solve a complex problem by breaking it down into simpler sub-problems in a recursive manner
- We will see how to use DP to solve an MDP thanks to the Bellman Equations

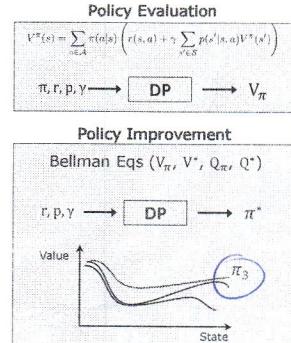


Machine Learning - Daniele Lolliano

4

Why Dynamic Programming?

- To solve an MDP we need to find an optimal policy
- Unfortunately we cannot use a brute-force approach:
 - ▶ $|\mathcal{A}|^{|S|}$ deterministic policies to evaluate
 - ▶ $|S|$ linear equations to solve for each policy
- Dynamic Programming (DP) is a method that allow to solve a complex problem by breaking it down into simpler sub-problems in a recursive manner
- We will see how to use DP to solve an MDP thanks to the Bellman Equations

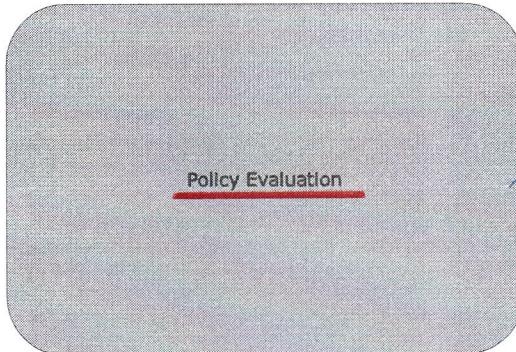
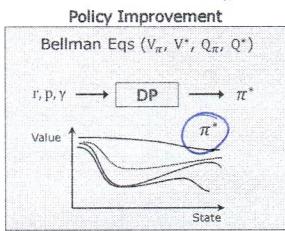
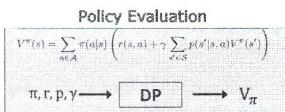


Machine Learning - Daniele Lolliano

5

Why Dynamic Programming?

- To solve an MDP we need to find an optimal policy
- Unfortunately we cannot use a brute-force approach:
 - $|\mathcal{A}|^{|S|}$ deterministic policies to evaluate
 - $|S|$ linear equations to solve for each policy
- Dynamic Programming (DP) is a method that allows to solve a complex problem by breaking it down into simpler sub-problems in a recursive manner
- We will see how to use DP to solve an MDP thanks to the Bellman Equations



we have a given policy π , we need to evaluate it

Iterative Policy Evaluation

- We search the solution of the Bellman expectation equation:

$$V_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left(r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V_\pi(s') \right)$$

- DP solves this problem through iterative application of Bellman equation:

$$V_{k+1}(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \left(r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V_k(s') \right)$$



- At each iteration k , the value-function V_k is updated for all state $s \in S$

initialization of the value function

$V_0 \rightarrow V_1 \rightarrow \dots \rightarrow V_k \rightarrow V_{k+1} \rightarrow V_\pi$ sweep

It can be proved that V_k converge to V_π as $k \rightarrow \infty$ for any V_0

Iterative Policy Evaluation (2)

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

Input π , the policy to be evaluated

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

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

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in S$:

$v \leftarrow V(s)$

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

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

until $\Delta < \theta$

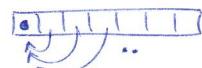
«in place» update

the alternative is:

old:

new:

we use the vector $V(i)$ to update itself:



This brings the property that the final result depends on the order of the spaces in the vector $V(\cdot)$.

Moreover, this cannot be parallelized. However, it converges fast.

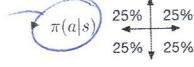
Iterative Policy Evaluation: a Small Gridworld Example

- Let consider gridworld environment with two terminal states, where

$$\gamma = 1$$

$$r(s, a) = -1 \quad \forall s, a$$

we want to evaluate this policy



terminal state

$$V_0 = \underline{\underline{0}}$$

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

$$V_1$$

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	-1.0

$$V_2$$

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

$$V_3$$

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

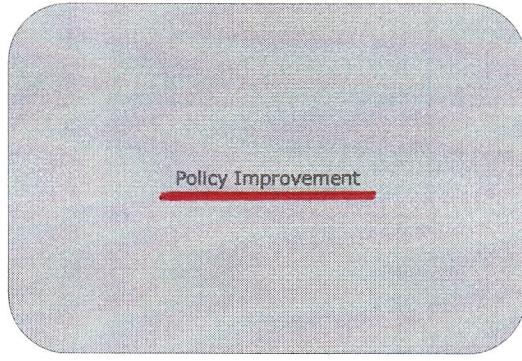
$$V_{10}$$

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

$$V_\infty$$

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

(Here we're not using the in-place update, but the *)



Policy Improvement

- Do you remember how to derive optimal policy from optimal value functions?

$$\pi^*(s) = \arg \max_a \left\{ r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^*(s') \right\} = \arg \max_a Q^*(s, a)$$

- What happens if we act greedy with respect to non optimal value function?

$$\pi'(s) = \arg \max_a \left\{ r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V_\pi(s') \right\} = \arg \max_a Q_\pi(s, a), \quad \forall s \in S$$

- Is π' different from π ?

Yes, they'll be equal only when $\pi = \pi' = \pi^*$ (optimal)

If not, it means that π is already the optimal policy π^* (as it satisfies the Bellman Optimality equations)

Otherwise, is π' better or as good as π ?

If $\pi'(s) \neq \pi(s)$ then:
 $Q_{\pi}(s, \pi(s)) \leq Q_{\pi'}(s, \pi'(s))$

Policy Improvement Theorem

- For any pair deterministic policies π' and π such that:

$$Q_\pi(s, \pi'(s)) \geq Q_\pi(s, \pi(s)), \quad \forall s \in S$$

then π' is better or as good as π

$$\pi' \geq \pi$$

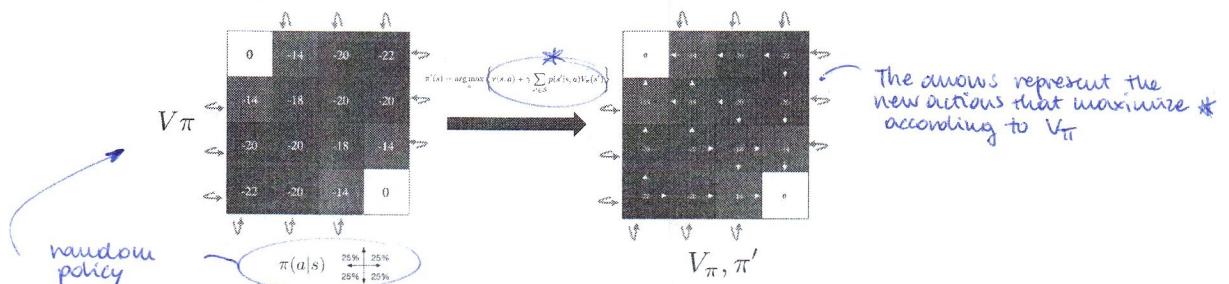
If $\exists s \in S$ s.t. $Q_\pi(s, \pi'(s)) > Q_\pi(s, \pi(s))$ then $\pi' > \pi$

- Proof

$$\begin{aligned} V_\pi(s) &\leq Q_\pi(s, \pi'(s)) = \mathbb{E}_{\pi'} [R_{t+1} + \gamma V_\pi(S_{t+1}) | S_t = s] \\ &\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma Q_\pi(S_{t+1}, \pi'(S_{t+1})) | S_t = s] \\ &\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 Q_\pi(S_{t+2}, \pi'(S_{t+2})) | S_t = s] \\ &\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \dots | S_t = s] = V_{\pi'}(s) \end{aligned}$$

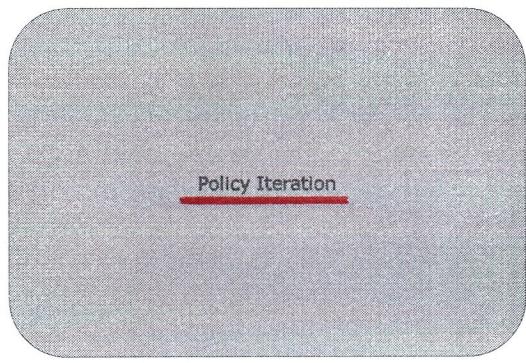
Policy Improvement: a Small Gridworld Example

- Let's go back to the value function we found iteratively for a random policy in the Small Gridworld example



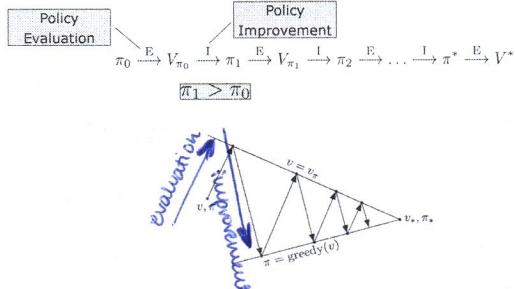
Policy Improvement: a Small Gridworld Example

V_0	$\begin{array}{ c c c c } \hline 0.0 & 0.0 & 0.0 & 0.0 \\ \hline 0.0 & 0.0 & 0.0 & 0.0 \\ \hline 0.0 & 0.0 & 0.0 & 0.0 \\ \hline 0.0 & 0.0 & 0.0 & 0.0 \\ \hline \end{array}$	\rightarrow	$\begin{array}{ c c c c } \hline \uparrow & \downarrow & \leftarrow & \rightarrow \\ \hline \uparrow & \downarrow & \leftarrow & \rightarrow \\ \hline \uparrow & \downarrow & \leftarrow & \rightarrow \\ \hline \uparrow & \downarrow & \leftarrow & \rightarrow \\ \hline \end{array}$	V_2	$\begin{array}{ c c c c } \hline 0.0 & -1.7 & 2.0 & -2.0 \\ \hline -1.7 & -2.0 & 2.0 & -2.0 \\ \hline 2.0 & -2.0 & 2.0 & -1.7 \\ \hline -2.0 & -2.0 & 1.7 & 0.0 \\ \hline \end{array}$	\rightarrow	$\begin{array}{ c c c c } \hline \leftarrow & \uparrow & \downarrow & \rightarrow \\ \hline \uparrow & \downarrow & \leftarrow & \rightarrow \\ \hline \uparrow & \downarrow & \leftarrow & \rightarrow \\ \hline \uparrow & \downarrow & \leftarrow & \rightarrow \\ \hline \end{array}$
V_1	$\begin{array}{ c c c c } \hline 0.0 & -1.0 & -1.0 & -1.0 \\ \hline -1.0 & -1.0 & -1.0 & -1.0 \\ \hline -1.0 & -1.0 & -1.0 & -1.0 \\ \hline -1.0 & -1.0 & -1.0 & 0.0 \\ \hline \end{array}$	\rightarrow	$\begin{array}{ c c c c } \hline \uparrow & \downarrow & \leftarrow & \rightarrow \\ \hline \uparrow & \downarrow & \leftarrow & \rightarrow \\ \hline \uparrow & \downarrow & \leftarrow & \rightarrow \\ \hline \uparrow & \downarrow & \leftarrow & \rightarrow \\ \hline \end{array}$	V_3	$\begin{array}{ c c c c } \hline 0.0 & -2.4 & 2.5 & -3.0 \\ \hline -2.4 & -2.8 & 3.0 & -2.8 \\ \hline 2.5 & -3.0 & 2.4 & -2.4 \\ \hline -3.0 & -2.4 & 0.0 & 0.0 \\ \hline \end{array}$	\rightarrow	π^*

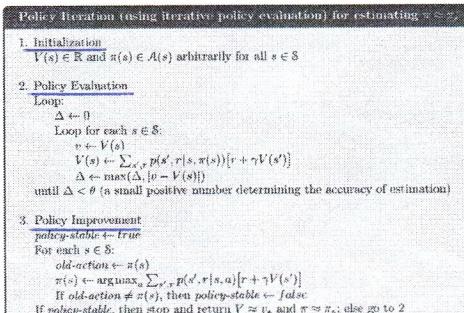


Policy Iteration

- We can exploit the policy improvement theorem to find the optimal policy:

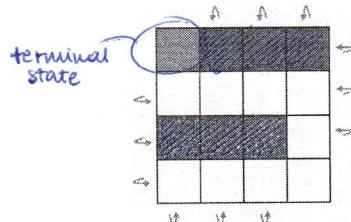


Policy Iteration (2)



Policy Iteration: Example

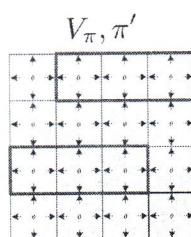
Goal: we want to reach the terminal state with the smaller price to pay



for this reason
it's convenient to
avoid the blue states

$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

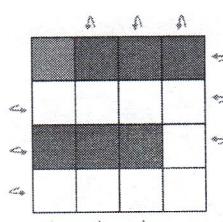
$$\gamma = 1$$



value function V_{π}
initialized to 0's
and initial policy
completely random

Evaluation
Improvement

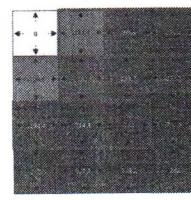
Policy Iteration: Example



$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

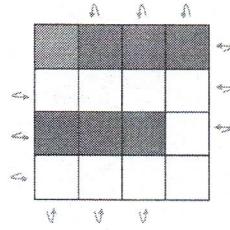
$$\gamma = 1$$

We perform: evaluation (of the completely random policy)
 V_{π}, π'



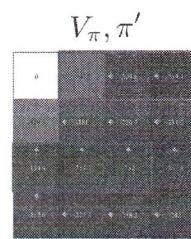
Evaluation
Improvement

Policy Iteration: Example



$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

$$\gamma = 1$$



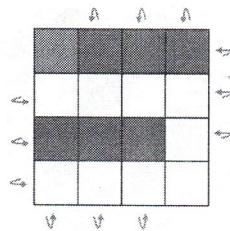
Evaluation
Improvement

(using the knowledge discovered evaluating the random policy)

The arrows (choices/actions) are now modified: we created a new policy

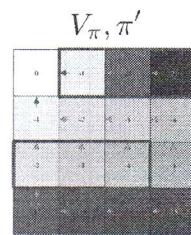
Machine Learning - Daniela Lolli

Policy Iteration: Example



$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

$$\gamma = 1$$

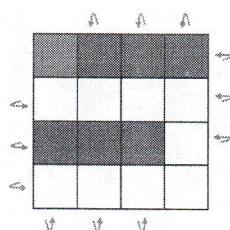


Evaluation
Improvement

The new policy is already meaningful

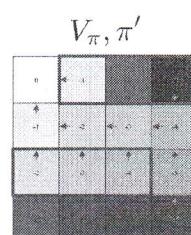
Machine Learning - Daniela Lolli

Policy Iteration: Example



$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

$$\gamma = 1$$

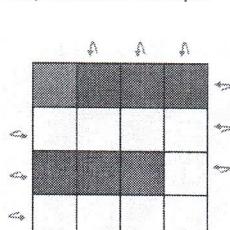


Evaluation
Improvement

We create a 3rd policy

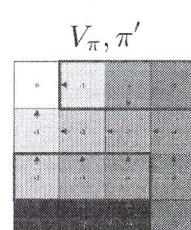
Machine Learning - Daniela Lolli

Policy Iteration: Example



$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

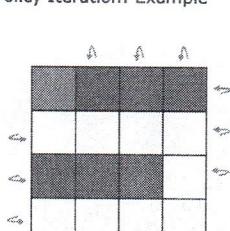
$$\gamma = 1$$



Evaluation
Improvement

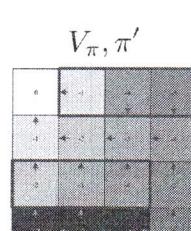
Machine Learning - Daniela Lolli

Policy Iteration: Example



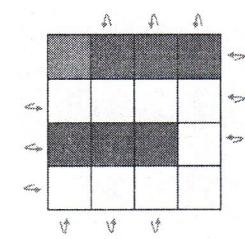
$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

$$\gamma = 1$$



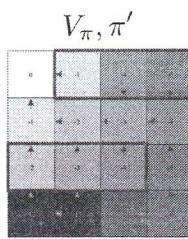
Evaluation
Improvement

Policy Iteration: Example



$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

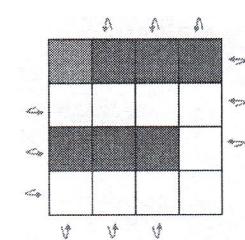
$$\gamma = 1$$



Evaluation
Improvement

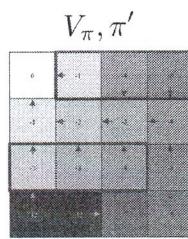
26

Policy Iteration: Example



$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

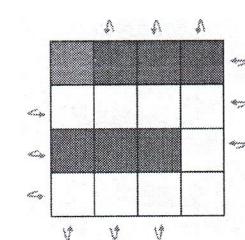
$$\gamma = 1$$



Evaluation
Improvement

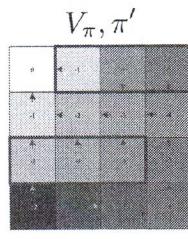
27

Policy Iteration: Example



$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

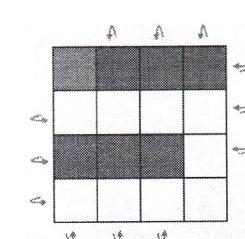
$$\gamma = 1$$



Evaluation
Improvement

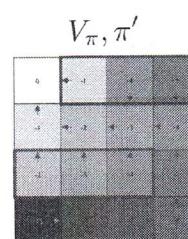
28

Policy Iteration: Example



$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

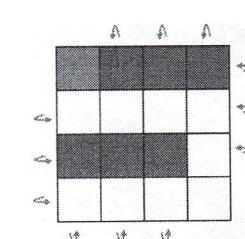
$$\gamma = 1$$



Evaluation
Improvement

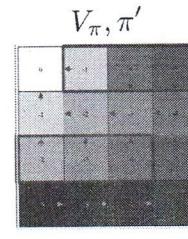
29

Policy Iteration: Example



$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

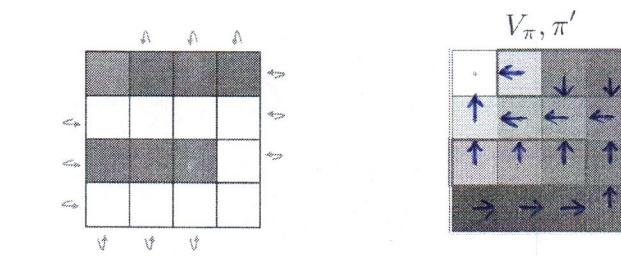
$$\gamma = 1$$



Evaluation
Improvement

30

Policy Iteration: Example

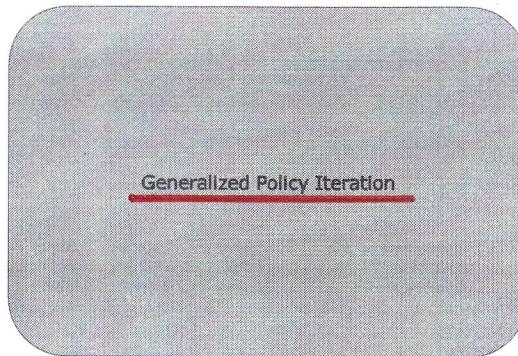


$$R_t = \begin{cases} -10 & \text{in blue states} \\ -1 & \text{in other states} \end{cases}$$

$$\gamma = 1$$

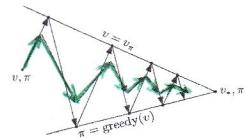
Optimal policy and value function

Reasonably it tries to avoid the blue zones, if, however, we end up in a blue zone we go out of it immediately



Generalized Policy Iteration

- Policy iteration alternates complete policy evaluation and improvement up to the convergence:



- Policy iteration framework allows also to find the optimal policy interleaving partial evaluation and improvement steps
- In particular, Value Iteration is one of the most popular GPI method

Value Iteration

- In the policy evaluation step, only a single sweep of updates is performed:

improvement:

$$\pi'(s) = \arg \max_a \left\{ r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V_\pi(s') \right\}, \forall s \in S$$

evaluation:

$$V_{k+1}(s) \leftarrow \sum_{a \in A} \pi'(a|s) \left(r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V_k(s') \right), \forall s \in S$$

} we combine the one-step improvement and the one-step evaluation: we do them together at once

- Combining them, we simply need to iterate the update of the value function using the Bellman optimality equation:

$$V_{k+1}(s) \leftarrow \max_a \left[r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V_k(s') \right], \forall s \in S$$

- It can be proved that $\lim_{k \rightarrow \infty} V_k = V^*$

Value Iteration (2)

This does not try to compute explicitly the policy, only the value function. At the very end we still need to compute the optimal policy: once we have the optimal value function we derive it (through arg max)

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

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in S^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

```

    | Δ ← 0
    | Loop for each  $s \in S$ :
    |   |  $v \leftarrow V(s)$ 
    |   |  $V(s) \leftarrow \max_a \sum_{s', r} p(s', r|s, a) [r + \gamma V(s')]$ 
    |   | Δ ← max(Δ, |v - V(s)|)
    | until Δ < θ
  
```

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

Efficiency of DP

Asynchronous Dynamic Programming

- All the DP methods described so far require exhaustive sweeps of the entire state set.
- Asynchronous DP does not use sweeps:
 - Pick a state at random
 - Apply the appropriate backup (= update)
 - Repeat until convergence criterion is met
- Can you select states to backup intelligently?
 - An agent's experience can act as a guide.

Efficiency of DP

- To find an optimal policy is polynomial in the number of states and actions
 - Value Iteration: $O(|S|^2|A|)$ (instead of the brute force $O(|A|^{|S|})$)
 - Policy Iteration: iterative evaluation $O\left(\frac{|S|^2 \log(1/\epsilon)}{\log(1/\gamma)}\right)$, improvement $O\left(\frac{|A|}{1-\gamma} \log\left(\frac{|S|}{1-\gamma}\right)\right)$
- Unfortunately, the number of states is often astronomical, e.g., often growing exponentially with the number of state variables (**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 is appropriate for parallel computation
 - But it is **easy** to come up with MDPs for which DP methods are not practical
- Linear programming approaches can be also used instead of DP but they do not typically scale well on larger problems