

Solving the Bellman Optimality Equations: Basic Methods

AI & Agents for IET
Lecturer: Liliana Mamani Sanchez

<http://www.scss.tcd.ie/~mamanisl/teaching/cs7032/>

November 30, 2015

The basic background & assumptions

- ▶ Environment is a **finite MDP** (i.e. A and S are finite).
- ▶ MDP's dynamics defined by **transition probabilities**:
- ▶ and expected **immediate rewards**,
- ▶ Goal: to **search for good policies** π
- ▶ Strategy: use **value functions** to structure search:

The basic background & assumptions

- ▶ Environment is a **finite MDP** (i.e. A and S are finite).
- ▶ MDP's dynamics defined by **transition probabilities**:

$$\mathcal{P}_{ss'}^a = P(s_{t+1} = s' | s_t = s, a_t = a)$$

- ▶ and expected **immediate rewards**,
- ▶ Goal: to **search for good policies** π
- ▶ Strategy: use **value functions** to structure search:

The basic background & assumptions

- ▶ Environment is a **finite MDP** (i.e. A and S are finite).
- ▶ MDP's dynamics defined by **transition probabilities**:

$$\mathcal{P}_{ss'}^a = P(s_{t+1} = s' | s_t = s, a_t = a)$$

- ▶ and expected **immediate rewards**,

$$\mathcal{R}_{ss'}^a = E\{r_{t+1} | s_t = s, a_t = a, s_{t+1} = s'\}$$

- ▶ Goal: to **search for good policies** π
- ▶ Strategy: use **value functions** to structure search:

The basic background & assumptions

- ▶ Environment is a **finite MDP** (i.e. A and S are finite).
- ▶ MDP's dynamics defined by **transition probabilities**:

$$\mathcal{P}_{ss'}^a = P(s_{t+1} = s' | s_t = s, a_t = a)$$

- ▶ and expected **immediate rewards**,

$$\mathcal{R}_{ss'}^a = E\{r_{t+1} | s_t = s, a_t = a, s_{t+1} = s'\}$$

- ▶ Goal: to **search for good policies** π
- ▶ Strategy: use **value functions** to structure search:

$$V^*(s) = \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^*(s')] \quad \text{or}$$

$$Q^*(s, a) = \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma \max_{a'} Q^*(s', a')]$$

Overview of methods for solving Bellman's equations

- ▶ Dynamic programming:
 - ▶ well-understood mathematical properties...
 - ▶ ...but require a complete and accurate model of the environment

Overview of methods for solving Bellman's equations

- ▶ **Dynamic programming:**
 - ▶ well-understood mathematical properties...
 - ▶ ...but require a complete and accurate model of the environment
- ▶ **Monte Carlo** (simulation methods):
 - ▶ conceptually simple
 - ▶ no model required...
 - ▶ ...but unsuitable for incremental computation

Overview of methods for solving Bellman's equations

- ▶ **Dynamic programming:**
 - ▶ well-understood mathematical properties...
 - ▶ ...but require a complete and accurate model of the environment
- ▶ **Monte Carlo** (simulation methods):
 - ▶ conceptually simple
 - ▶ no model required...
 - ▶ ...but unsuitable for incremental computation
- ▶ **Temporal difference** methods
 - ▶ also require no model;
 - ▶ suitable for incremental computation...
 - ▶ ... but mathematically complex to analyse

Dynamic programming

- ▶ Basic Idea: “sweep” through S performing a full backup operation on each s .
- ▶ A few different methods exist. E.g.:
 - ▶ Policy Iteration and
 - ▶ Value Iteration.
- ▶ The building blocks:
 - ▶ Policy Evaluation: how to compute V^π for an arbitrary π .
 - ▶ Policy Improvement: how to compute an improved π given V^π .

Policy Evaluation

- ▶ The task of computing V^π for an arbitrary π is known as the **prediction problem**.
- ▶ As we have seen, a state-value function is given by

$$\begin{aligned} V^\pi(s) &= E_\pi\{R_t | s_t = s\} = E_\pi\{r_{t+1} + \gamma V^\pi(s_{t+1}) | s_t = s\} \\ &= \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')] \end{aligned}$$

- ▶ a system of $|S|$ linear equations in $|S|$ unknowns (the state values $V^\pi(s)$)

Iterative Policy evaluation

- ▶ Consider the sequence of approximations $V_0, \dots V^\pi$.
- ▶ Choose V_0 arbitrarily and set each successive approximation accommodation to the Bellman equation:

$$\begin{aligned} V_{k+1}(s) &\leftarrow E_\pi\{r_{t+1} + \gamma V_k(s_{t+1}) | s_t = s\} \\ &\leftarrow \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s')] \quad (1) \end{aligned}$$

- ▶ “Sweeps”:

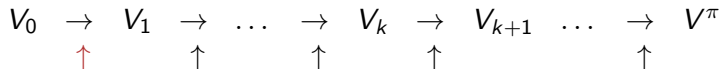
$$\begin{array}{ccccccccccc} V_0 & \rightarrow & V_1 & \rightarrow & \dots & \rightarrow & V_k & \rightarrow & V_{k+1} & \dots & \rightarrow & V^\pi \\ & \uparrow & & \uparrow & & \uparrow & & \uparrow & & & \uparrow & \end{array}$$

Iterative Policy evaluation

- ▶ Consider the sequence of approximations $V_0, \dots V^\pi$.
- ▶ Choose V_0 arbitrarily and set each successive approximation accommodation to the Bellman equation:

$$\begin{aligned} V_{k+1}(s) &\leftarrow E_\pi\{r_{t+1} + \gamma V_k(s_{t+1}) | s_t = s\} \\ &\leftarrow \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s')] \end{aligned} \quad (1)$$

- ▶ “Sweeps”:

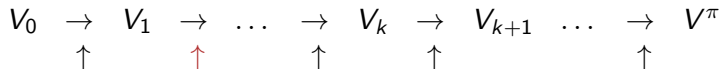


Iterative Policy evaluation

- ▶ Consider the sequence of approximations $V_0, \dots V^\pi$.
- ▶ Choose V_0 arbitrarily and set each successive approximation accommodation to the Bellman equation:

$$\begin{aligned} V_{k+1}(s) &\leftarrow E_\pi\{r_{t+1} + \gamma V_k(s_{t+1}) | s_t = s\} \\ &\leftarrow \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s')] \end{aligned} \quad (1)$$

- ▶ “Sweeps”:

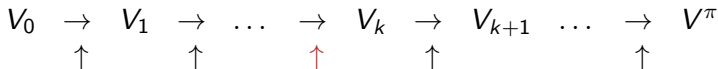


Iterative Policy evaluation

- ▶ Consider the sequence of approximations $V_0, \dots V^\pi$.
- ▶ Choose V_0 arbitrarily and set each successive approximation accommodation to the Bellman equation:

$$\begin{aligned} V_{k+1}(s) &\leftarrow E_\pi\{r_{t+1} + \gamma V_k(s_{t+1}) | s_t = s\} \\ &\leftarrow \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s')] \end{aligned} \quad (1)$$

- ▶ “Sweeps”:

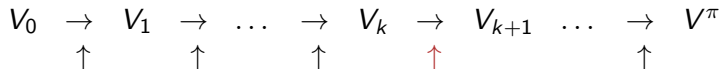


Iterative Policy evaluation

- ▶ Consider the sequence of approximations $V_0, \dots V^\pi$.
- ▶ Choose V_0 arbitrarily and set each successive approximation accommodation to the Bellman equation:

$$\begin{aligned} V_{k+1}(s) &\leftarrow E_\pi\{r_{t+1} + \gamma V_k(s_{t+1}) | s_t = s\} \\ &\leftarrow \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s')] \quad (1) \end{aligned}$$

- ▶ “Sweeps”:

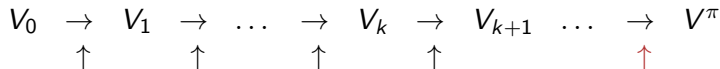


Iterative Policy evaluation

- ▶ Consider the sequence of approximations $V_0, \dots V^\pi$.
- ▶ Choose V_0 arbitrarily and set each successive approximation accommodation to the Bellman equation:

$$\begin{aligned} V_{k+1}(s) &\leftarrow E_\pi\{r_{t+1} + \gamma V_k(s_{t+1}) | s_t = s\} \\ &\leftarrow \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s')] \end{aligned} \quad (1)$$

- ▶ “Sweeps”:



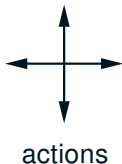
Iterative Policy Evaluation Algorithm

```
1  Initialisation :
2    for (each  $s \in S$ )
3       $V(s) \leftarrow 0$ 
4
5  IPE( $\pi$ )                                /*  $\pi$ : policy to be evaluated */
6    repeat
7       $\Delta \leftarrow 0$ 
8       $V_k \leftarrow V$ 
9      for (each  $s \in S / \{s_{terminal}\}$ )
10        $V(s) \leftarrow \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s')]$ 
11        $\Delta \leftarrow \max(\Delta, |V_k(s) - V(s)|)$ 
12    until  $\Delta < \theta$                       /*  $\theta > 0$ : a small constant */
13    return  $V$                              /*  $V \approx V^\pi$  */
```

- NB: alternatively one could evaluate V^π *in place* (i.e. using a single vector V to store all values and update it directly).

An example: an episodic GridWorld

- ▶ Rewards of -1 until terminal state (shown in grey) is reached
- ▶ Undiscounted episodic task:



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

$r = -1$
on all transitions

Policy Evaluation for the GridWorld

- Iterative evaluation of V_k for equiprobable random policy π :

Policy Evaluation for the GridWorld

- Iterative evaluation of V_k for equiprobable random policy π :

$V_0 \rightarrow$

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

Policy Evaluation for the GridWorld

- Iterative evaluation of V_k for equiprobable random policy π :

$V_0 \rightarrow$	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

$V_1 \rightarrow$	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

Policy Evaluation for the GridWorld

- Iterative evaluation of V_k for equiprobable random policy π :

$V_0 \rightarrow$	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

$V_1 \rightarrow$	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

$V_2 \rightarrow$	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

Policy Evaluation for the GridWorld

- Iterative evaluation of V_k for equiprobable random policy π :

$V_0 \rightarrow$

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

$V_1 \rightarrow$

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

$V_2 \rightarrow$

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

$V_3 \rightarrow$

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

Policy Evaluation for the GridWorld

- Iterative evaluation of V_k for equiprobable random policy π :

$V_0 \rightarrow$

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

$V_1 \rightarrow$

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

$V_2 \rightarrow$

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

$V_3 \rightarrow$

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

$V_{10} \rightarrow$

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

Policy Evaluation for the GridWorld

- Iterative evaluation of V_k for equiprobable random policy π :

$V_0 \rightarrow$

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

$V_1 \rightarrow$

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

$V_2 \rightarrow$

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

$V_3 \rightarrow$

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

$V_{10} \rightarrow$

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

$V_\infty \rightarrow$

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

Policy Improvement

- ▶ Consider the following: how would the expected return change for a policy π if instead of following $\pi(s)$ for a given state s we choose an action $a \neq \pi(s)$?
- ▶ For this setting, the value would be:

Policy Improvement

- ▶ Consider the following: how would the expected return change for a policy π if instead of following $\pi(s)$ for a given state s we choose an action $a \neq \pi(s)$?

- ▶ For this setting, the value would be:

$$\begin{aligned} Q^\pi(s, a) &= E_\pi\{r_{t+1} + \gamma V^\pi(s_{t+1}) | s_t = s, a_t = a\} \\ &= \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s_{t+1})] \end{aligned}$$

- ▶ So, a should be preferred iff $Q^\pi(s, a) > V^\pi(s)$

Policy improvement theorem

If choosing $a \neq \pi(s)$ implies $Q^\pi(s, a) \geq V^\pi(s)$ for a state s , then the policy π' obtained by choosing a every time s is encountered (and following π otherwise) is at least as good as π (i.e. $V^{\pi'}(s) \geq V^\pi(s)$).
If $Q^\pi(s, a) > V^\pi(s)$ then $V^{\pi'}(s) > V^\pi(s)$

- ▶ If we apply this strategy to all states to get a new greedy policy $\pi'(s) = \arg \max_a Q^\pi(s, a)$, then $V^{\pi'} \geq V^\pi$
- ▶ $V^{\pi'} = V^\pi$ implies that

$$V^{\pi'}(s) = \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')]$$

which is...

Policy improvement theorem

If choosing $a \neq \pi(s)$ implies $Q^\pi(s, a) \geq V^\pi(s)$ for a state s , then the policy π' obtained by choosing a every time s is encountered (and following π otherwise) is at least as good as π (i.e. $V^{\pi'}(s) \geq V^\pi(s)$).
If $Q^\pi(s, a) > V^\pi(s)$ then $V^{\pi'}(s) > V^\pi(s)$

- ▶ If we apply this strategy to all states to get a new greedy policy $\pi'(s) = \arg \max_a Q^\pi(s, a)$, then $V^{\pi'} \geq V^\pi$
- ▶ $V^{\pi'} = V^\pi$ implies that

$$V^{\pi'}(s) = \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')]$$

which is... a form of the Bellman optimality equation.

Policy improvement theorem

If choosing $a \neq \pi(s)$ implies $Q^\pi(s, a) \geq V^\pi(s)$ for a state s , then the policy π' obtained by choosing a every time s is encountered (and following π otherwise) is at least as good as π (i.e. $V^{\pi'}(s) \geq V^\pi(s)$). If $Q^\pi(s, a) > V^\pi(s)$ then $V^{\pi'}(s) > V^\pi(s)$

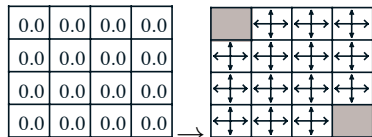
- ▶ If we apply this strategy to all states to get a new greedy policy $\pi'(s) = \arg \max_a Q^\pi(s, a)$, then $V^{\pi'} \geq V^\pi$
- ▶ $V^{\pi'} = V^\pi$ implies that

$$V^{\pi'}(s) = \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')]$$

which is... a form of the Bellman optimality equation.

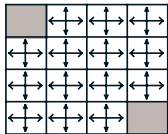
- ▶ Therefore $V^\pi = V^{\pi'} = V^*$

Improving the GridWorld policy

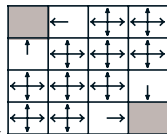


Improving the GridWorld policy

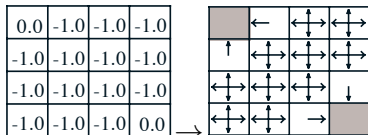
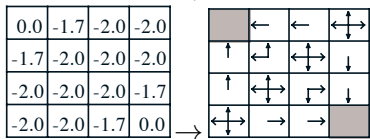
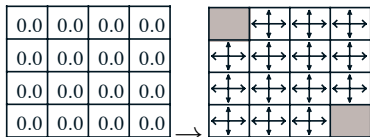
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0



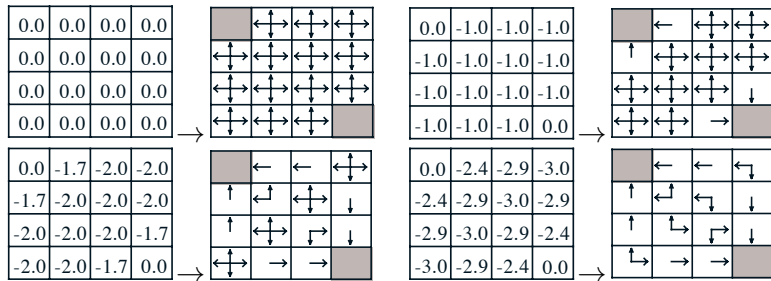
0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0



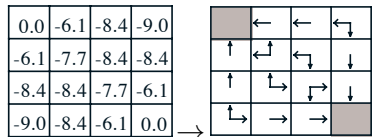
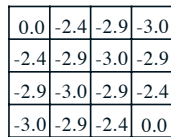
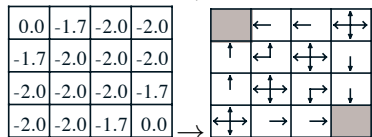
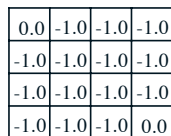
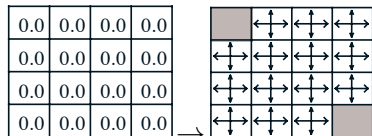
Improving the GridWorld policy



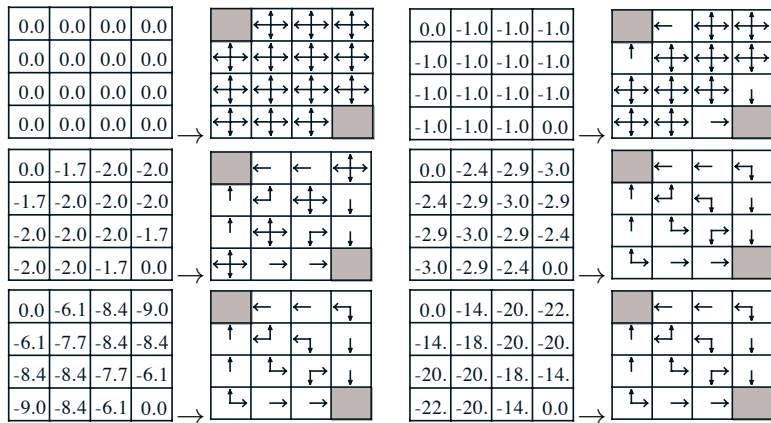
Improving the GridWorld policy



Improving the GridWorld policy



Improving the GridWorld policy



Putting them together: Policy Iteration

π_0

Putting them together: Policy Iteration

$$\pi_0 \xrightarrow{\text{eval}}$$

Putting them together: Policy Iteration

$$\pi_0 \xrightarrow{\text{eval}} V_{\pi_0}$$

Putting them together: Policy Iteration

$$\pi_0 \xrightarrow{\text{eval}} V_{\pi_0} \xrightarrow{\text{improve}}$$

Putting them together: Policy Iteration

$$\pi_0 \xrightarrow{\text{eval}} V_{\pi_0} \xrightarrow{\text{improve}} \pi_1$$

Putting them together: Policy Iteration

$$\pi_0 \xrightarrow{\text{eval}} V_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{e}$$

Putting them together: Policy Iteration

$$\pi_0 \xrightarrow{\text{eval}} V_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{e} V_{\pi_1}$$

Putting them together: Policy Iteration

$$\pi_0 \xrightarrow{\text{eval}} V_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{e} V_{\pi_1} \xrightarrow{i} \dots \xrightarrow{i}$$

Putting them together: Policy Iteration

$$\pi_0 \xrightarrow{\text{eval}} V_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{e} V_{\pi_1} \xrightarrow{i} \dots \xrightarrow{i} \pi^*$$

Putting them together: Policy Iteration

$$\pi_0 \xrightarrow{\text{eval}} V_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{e} V_{\pi_1} \xrightarrow{i} \dots \xrightarrow{i} \pi^* \xrightarrow{e}$$

Putting them together: Policy Iteration

$$\pi_0 \xrightarrow{\text{eval}} V_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{e} V_{\pi_1} \xrightarrow{i} \dots \xrightarrow{i} \pi^* \xrightarrow{e} V^*$$

Putting them together: Policy Iteration

$$\pi_0 \xrightarrow{\text{eval}} V_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{e} V_{\pi_1} \xrightarrow{i} \dots \xrightarrow{i} \pi^* \xrightarrow{e} V^*$$

1 Initialisation :

2 for all $s \in S$

3 $V(s) \leftarrow$ an arbitrary $v \in \mathbb{R}$

5 Policy_Improvement(π):

6 do

7 $\text{stable}(\pi) \leftarrow \text{true}$

8 $V \leftarrow \text{IPE}(\pi)$

9 for each $s \in S$

10 $b \leftarrow \pi(s)$

11 $\pi(s) \leftarrow \arg \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V(s')]$

12 if ($b \neq \pi(s)$)

13 $\text{stable}(\pi) \leftarrow \text{false}$

14 while (not $\text{stable}(\pi)$)

15 return π

Other DP methods

- ▶ **Value Iteration**: evaluation is stopped after a **single sweep** (one backup of each state). The backup rule is then:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s')]$$

- ▶ **Asynchronous DP**: back up the values of states **in any order**, using whatever values of other states happen to be available.
 - ▶ On problems with **large state spaces**, asynchronous DP methods are often preferred

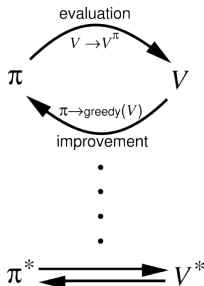
Other DP methods

- **Value Iteration**: evaluation is stopped after a **single sweep** (one backup of each state). The backup rule is then:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s')]$$

- **Asynchronous DP**: back up the values of states **in any order**, using whatever values of other states happen to be available.
 - On problems with **large state spaces**, asynchronous DP methods are often preferred

Generalised policy iteration

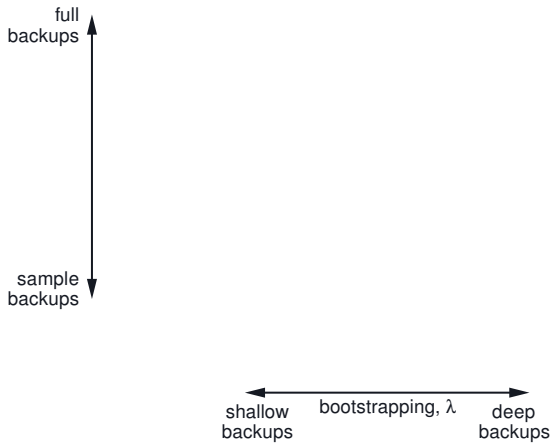


Value Iteration

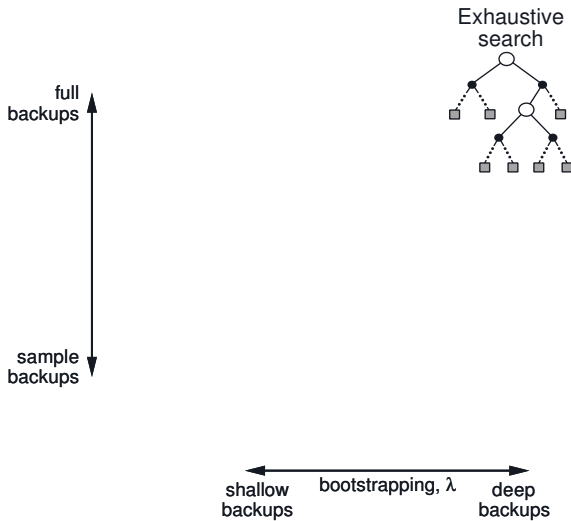
- Potential computational savings over Policy Iteration in terms of policy evaluation

```
1  Initialisation :
2    for all  $s \in S$ 
3       $V(s) \leftarrow$  an arbitrary  $v \in \mathbb{R}$ 
4
5  Value Iteration( $\pi$ ):
6    repeat
7       $\Delta \leftarrow 0$ 
8      for each  $s \in S$ 
9         $v \leftarrow V(s)$ 
10        $V(s) \leftarrow \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V(s')]$ 
11        $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
12    until  $\Delta < \theta$ 
13    return deterministic  $\pi$  s.t.
14       $\pi(s) = \arg \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V(s')]$ 
```

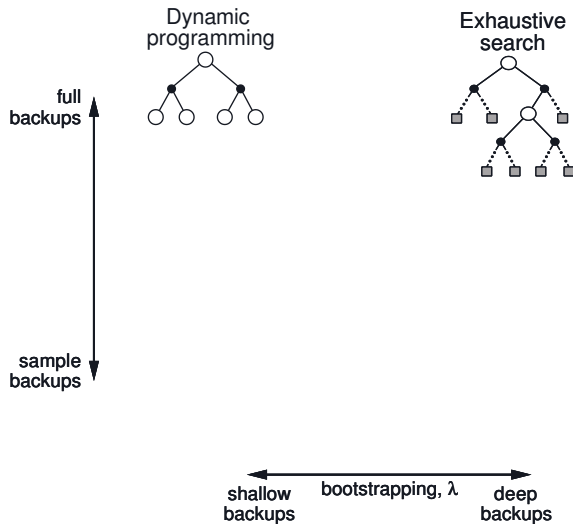
Summary of methods



Summary of methods



Summary of methods



Monte Carlo Methods

- ▶ Complete knowledge of environment is not necessary
- ▶ Only **experience** is required
- ▶ Learning can be **on-line** (no model needed) or through **simulated experience** (model only needs to generate sample transitions).
 - ▶ In both cases, learning is based on **averaged sample returns**.
- ▶ As in DP, one can use an **evaluation-improvement** strategy.
- ▶ Evaluation can be done by keeping averages of:
 - ▶ **Every Visit** to a state in an episode, or
 - ▶ of the **First Visit** to a state in an episode.

Estimating value-state functions in MC

The **first visit policy evaluation** method:

```
1  FirstVisitMC( $\pi$ )
2  Initialisation:
3       $V \leftarrow$  arbitrary state values
4       $Returns(s) \leftarrow$  empty list of size  $|S|$ 
5
6  Repeat
7      Generate an episode  $E$  using  $\pi$ 
8      For each  $s$  in  $E$ 
9           $R \leftarrow$  return following the first occurrence of  $s$ 
10         Append  $R$  to  $Returns(s)$ 
11          $V(s) \leftarrow mean(Returns(s))$ 
```

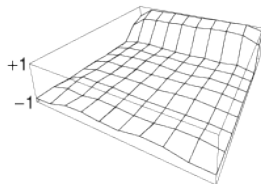
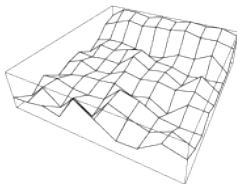
Example

Evaluate the policy described below for **blackjack**
[Sutton and Barto, 1998, section 5.1]

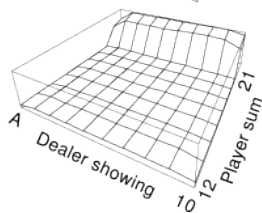
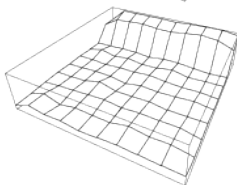
- ▶ Actions: stick (stop receiving cards), hit (receive another card)
- ▶ Play against dealer, who has a fixed strategy ('hit' if **sum** < 17 ; 'stick' otherwise).
- ▶ You win if your card sum is greater than the dealer's without exceeding 21.
- ▶ States:
 - ▶ current sum (12-21)
 - ▶ dealers showing card (ace-10)
 - ▶ do I have a useable ace (can be 11 without making sum exceed 21)?
- ▶ Reward: +1 for winning, 0 for a draw, -1 for losing
- ▶ Policy: Stick if my sum is 20 or 21, else hit

After 500,000 episodes

+1
—



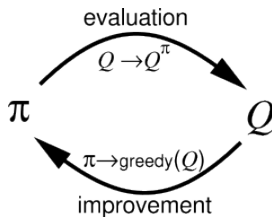
No
usable
ace



- Question: compare MC to DP in estimating the value function. Do we know the the **environment**? The **transition probabilities**? The **expected returns** given each state and action? What does the MC backup diagram look like?

Monte Carlo control

- ▶ Monte Carlo **version** of DP policy iteration:



- ▶ Policy **improvement** theorem applies:

$$\begin{aligned} Q^{\pi_k}(s, \pi_{k+1}(s)) &= Q^{\pi_k}(s, \arg \max_a Q^{\pi_k}(s, a)) \\ &= \max_a Q^{\pi_k}(s, a) \\ &\geq Q^{\pi_k}(s, a) \\ &= V^{\pi_k}(s) \end{aligned}$$

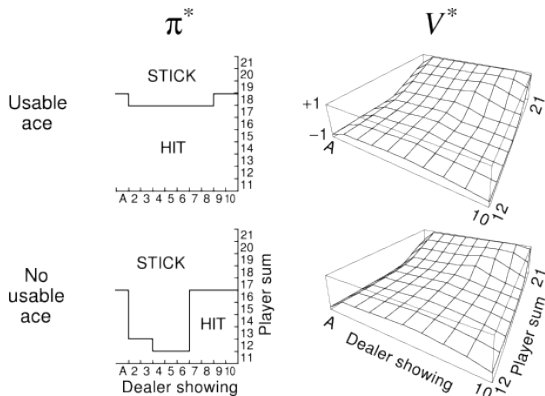
MC policy iteration (exploring starts)

- ▶ As with DP, we have **evaluation-improvement** cycles.
- ▶ Learn Q^* (if no model is available)
- ▶ One must make sure that each state-action pair can be a **starting pair** (with probability > 0).

```
1  MonteCarloES()
2  Initialisation ,  $\forall s \in S, a \in A$ :
3       $Q(s,a) \leftarrow \text{arbitrary}$ ;  $\pi(s) \leftarrow \text{arbitrary}$ 
4       $Returns(s,a) \leftarrow \text{empty list of size } |S|$ 
5
6  Repeat until stop-criterion met
7      Generate an episode  $E$  using  $\pi$  and exploring starts
8      For each  $(s,a)$  in  $E$ 
9           $R \leftarrow \text{return}$  following the first occurrence of  $s,a$ 
10         Append  $R$  to  $Returns(s,a)$ 
11          $Q(s,a) \leftarrow \text{mean}(Returns(s,a))$ 
12     For each  $s$  in  $E$ 
13          $\pi(s) \leftarrow \arg \max_a Q(s,a)$ 
```

Optimal policy for blackjack example

- Optimal **policy** found by **MonteCarloES** for the blackjack example, and its **state-value** function:



On- and off- policy MC control

- ▶ MonteCarloES assumes that all states are observed an infinite number of times and episodes are generated with exploring starts
 - ▶ For an analysis of convergence properties, see [Tsitsiklis, 2003]
- ▶ On-policy and off-policy methods relax these assumptions to produce practical algorithms
- ▶ On-policy methods use a given policy and ϵ -greedy strategy (see lecture on Evaluative Feedback) to generate episodes.
- ▶ Off-policy methods evaluate a policy while generating an episode through a different policy

On-policy control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$Q(s, a) \leftarrow$ arbitrary

$Returns(s, a) \leftarrow$ empty list

$\pi \leftarrow$ an arbitrary ε -soft policy

Repeat forever:

(a) Generate an episode using π

(b) For each pair s, a appearing in the episode:

$R \leftarrow$ return following the first occurrence of s, a

Append R to $Returns(s, a)$

$Q(s, a) \leftarrow \text{average}(Returns(s, a))$

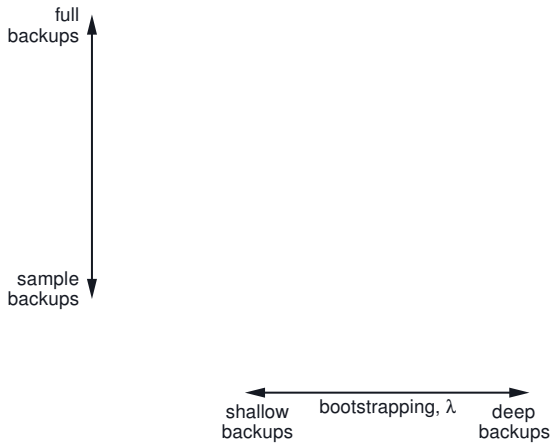
(c) For each s in the episode:

$a^* \leftarrow \arg \max_a Q(s, a)$

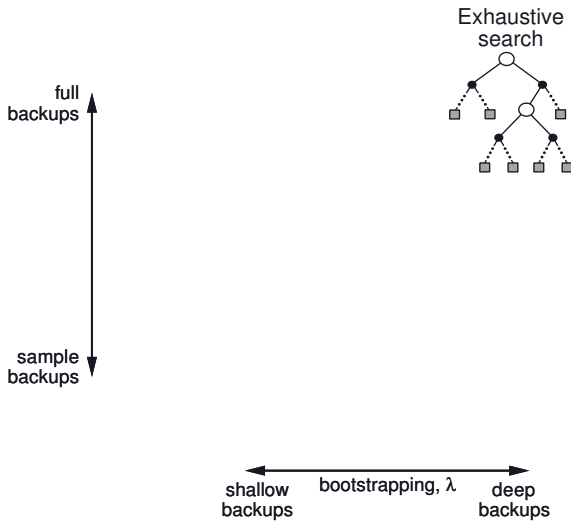
For all $a \in \mathcal{A}(s)$:

$$\pi(s, a) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(s)| & \text{if } a = a^* \\ \varepsilon / |\mathcal{A}(s)| & \text{if } a \neq a^* \end{cases}$$

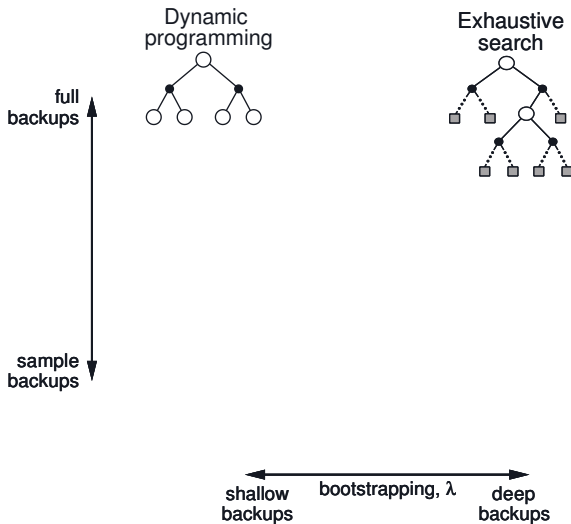
Summary of methods



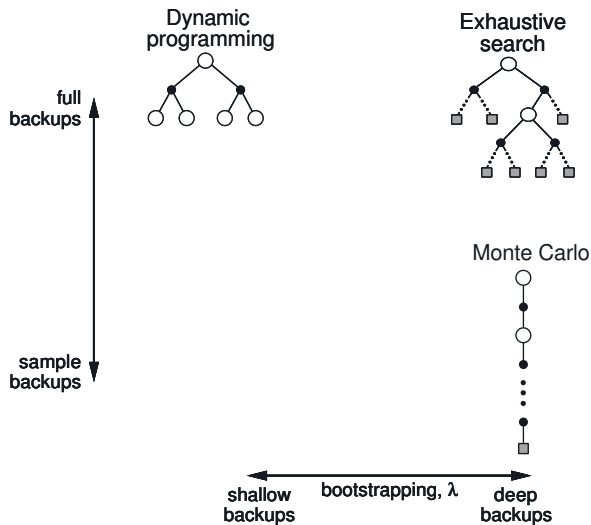
Summary of methods



Summary of methods



Summary of methods



References

Notes based on [Sutton and Barto, 1998, ch 4-6].
Convergence results for several MC algorithms are given by [Tsitsiklis, 2003].



Sutton, R. S. and Barto, A. G. (1998).
Reinforcement Learning: An Introduction.
MIT Press, Cambridge, MA.



Tsitsiklis, J. N. (2003).
On the convergence of optimistic policy iteration.
The Journal of Machine Learning Research, 3:59–72.