

# Lecture 03: Dynamic Programming

Oliver Wallscheid



# Table of contents

- 1 Introduction
- 2 Policy evaluation
- 3 Policy improvement
- 4 Policy and value iteration
- 5 Further aspects

# What is dynamic programming (DP)?

## Basic DP definition

- ▶ **Dynamic**: sequential or temporal problem structure
- ▶ **Programming**: mathematical optimization, i.e., numerical solutions

## Further characteristics:

- ▶ DP is a collection of algorithms to solve MDPs and neighboring problems.
  - ▶ **We will focus only on finite MDPs.**
  - ▶ In case of continuous action/state space: apply quantization.
- ▶ Use of value functions to organize and structure the search for an optimal policy.
- ▶ Breaks problems into subproblems and solves them.

# Requirements for DP

DP can be applied to problems with the following characteristics.

- ▶ Optimal substructure:
  - ▶ Principle of optimality applies.
  - ▶ Optimal solution can be derived from subproblems.
- ▶ Overlapping subproblems:
  - ▶ Subproblems recur many times.
  - ▶ Hence, solutions can be cached and reused.

How is that connected to MDPs?

- ▶ MDPs satisfy above's properties:
  - ▶ Bellman equation provides recursive decomposition.
  - ▶ Value function stores and reuses solutions.

## Example: DP vs. exhaustive search (1)

Fig. 3.1: Shortest path problem to travel from Paderborn to Bielefeld: Exhaustive search requires 14 travel segment evaluations since every possible travel route is evaluated independently.

## Example: DP vs. exhaustive search (2)

Fig. 3.2: Shortest path problem to travel from Paderborn to Bielefeld: DP requires only 10 travel segment evaluations in order to calculate the optimal travel policy due to the reuse of subproblem results.

# Utility of DP in the RL context

DP is used for iterative **model-based** prediction and control in an MDP.

▶ Prediction:

- ▶ Input: MDP  $\langle \mathcal{X}, \mathcal{U}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  and policy  $\pi$
- ▶ Output: (estimated) value function  $\hat{v}_\pi \approx v_\pi$

▶ Control:

- ▶ Input: MDP  $\langle \mathcal{X}, \mathcal{U}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
- ▶ Output: (estimated) optimal value function  $\hat{v}_\pi^* \approx v_\pi^*$  or policy  $\hat{\pi}^* \approx \pi^*$

In both applications **DP requires full knowledge of the MDP** structure.

- ▶ Feasibility in real-world engineering applications (model vs. system) is therefore limited.
- ▶ But: **following DP concepts are largely used in modern data-driven RL algorithms.**

# Table of contents

- 1 Introduction
- 2 Policy evaluation**
- 3 Policy improvement
- 4 Policy and value iteration
- 5 Further aspects

# Policy evaluation background (1)

- ▶ Problem: evaluate a given policy  $\pi$  to predict  $v_\pi$ .
- ▶ Recap: Bellman expectation equation for  $x_k \in \mathcal{X}$  is given as

$$\begin{aligned}v_\pi(x_k) &= \mathbb{E}_\pi [G_k | X_k = x_k], \\ &= \mathbb{E}_\pi [R_{k+1} + \gamma G_{k+1} | X_k = x_k], \\ &= \mathbb{E}_\pi [R_{k+1} + \gamma v_\pi(X_{k+1}) | X_k = x_k].\end{aligned}$$

- ▶ Or in matrix form:

$$\begin{aligned}\mathbf{v}_\pi^\pi &= \mathbf{r}_\pi^\pi + \gamma \mathbf{P}_{xx'}^\pi \mathbf{v}_\pi^\pi, \\ \begin{bmatrix} v_1^\pi \\ \vdots \\ v_n^\pi \end{bmatrix} &= \begin{bmatrix} \mathcal{R}_1^\pi \\ \vdots \\ \mathcal{R}_n^\pi \end{bmatrix} + \gamma \begin{bmatrix} p_{11}^\pi & \cdots & p_{1n}^\pi \\ \vdots & & \vdots \\ p_{n1}^\pi & \cdots & p_{nn}^\pi \end{bmatrix} \begin{bmatrix} v_1^\pi \\ \vdots \\ v_n^\pi \end{bmatrix}.\end{aligned}$$

- ▶ Solving the Bellman expectation equation for  $v_\pi$  requires handling a linear equation system with  $n$  unknowns (i.e., number of states).

## Policy evaluation background (2)

- ▶ Problem: directly calculating  $v_\pi$  is numerically costly for high-dimensional state spaces (e.g., by matrix inversion).
- ▶ General idea: **apply iterative approximations**  $\hat{v}_i(x_k) = v_i(x_k)$  of  $v_\pi(x_k)$  with decreasing errors:

$$\|v_i(x_k) - v_\pi\|_\infty \rightarrow 0 \quad \text{for } i = 1, 2, 3, \dots \quad (3.1)$$

- ▶ The Bellman equation in matrix form can be rewritten as:

$$\underbrace{(\mathbf{I} - \gamma \mathcal{P}_{xx'}^\pi)}_A \underbrace{\mathbf{v}_\mathcal{X}^\pi}_\zeta = \underbrace{\mathbf{r}_\mathcal{X}^\pi}_b. \quad (3.2)$$

- ▶ To iteratively solve this linear equation  $A\zeta = b$ , one can apply numerous methods such as
  - ▶ General gradient descent,
  - ▶ Richardson iteration,
  - ▶ Kyrlov subspace methods.

# Richardson iteration (1)

In the MDP context, the Richardson iteration became the default solution approach to iteratively solve:

$$\mathbf{A}\zeta = \mathbf{b}.$$

The **Richardson iteration** is

$$\zeta_{i+1} = \zeta_i + \omega(\mathbf{b} - \mathbf{A}\zeta_i) \quad (3.3)$$

with  $\omega$  being a scalar parameter that has to be chosen such that the sequence  $\zeta_i$  converges. To choose  $\omega$  we inspect the series of approximation errors  $\mathbf{e}_i = \zeta_i - \zeta$  and apply it to (3.3):

$$\mathbf{e}_{i+1} = \mathbf{e}_i - \omega\mathbf{A}\mathbf{e}_i = (\mathbf{I} - \omega\mathbf{A})\mathbf{e}_i. \quad (3.4)$$

To evaluate convergence we inspect the following norm:

$$\|\mathbf{e}_{i+1}\|_\infty = \|(\mathbf{I} - \omega\mathbf{A})\mathbf{e}_i\|_\infty. \quad (3.5)$$

## Richardson iteration (2)

Since any induced matrix norm is sub-multiplicative, we can approximate (3.5) by the inequality:

$$\|e_{i+1}\|_\infty \leq \|(\mathbf{I} - \omega\mathbf{A})\|_\infty \|e_i\|_\infty. \quad (3.6)$$

Hence, the series converges if

$$\|(\mathbf{I} - \omega\mathbf{A})\|_\infty < 1. \quad (3.7)$$

Inserting from (3.2) leads to:

$$\|(\mathbf{I}(1 - \omega) + \omega\gamma\mathcal{P}_{xx'}^\pi)\|_\infty < 1. \quad (3.8)$$

For  $\omega = 1$  we receive:

$$\gamma \|(\mathcal{P}_{xx'}^\pi)\|_\infty < 1. \quad (3.9)$$

Since the row elements of  $\mathcal{P}_{xx'}^\pi$  always sum up to 1,

$$\gamma < 1 \quad (3.10)$$

follows. Hence, **when discounting the Richardson iteration always converges for MDPs** even if we assume  $\omega = 1$ .

# Iterative policy evaluation by Richardson iteration (1)

Applying the Richardson iteration (3.3) with  $w = 1$  to the Bellman equation (2.17) for any  $x_k \in \mathcal{X}$  at iteration  $i$  results in:

$$v_{i+1}(x_k) = \sum_{u_k \in \mathcal{U}} \pi(u_k | x_k) \left( \mathcal{R}_x^u + \gamma \sum_{x_{k+1} \in \mathcal{X}} p_{xx'}^u v_i(x_{k+1}) \right). \quad (3.11)$$

Matrix form based on (2.19) then is:

$$\mathbf{v}_{\mathcal{X}, i+1}^{\pi} = \mathbf{r}_{\mathcal{X}}^{\pi} + \gamma \mathbf{P}_{xx'}^{\pi} \mathbf{v}_{\mathcal{X}, i}^{\pi}. \quad (3.12)$$

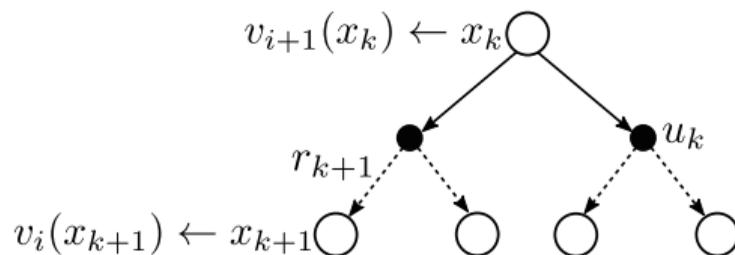


Fig. 3.3: Backup diagram for iterative policy evaluation

## Iterative policy evaluation by Richardson iteration (2)

- ▶ During one Richardson iteration the 'old' value of  $x_k$  is replaced with a 'new' value from the 'old' values of the successor state  $x_{k+1}$ .
  - ▶ Update  $v_{i+1}(x_k)$  from  $v_i(x_{k+1})$ , see Fig. 3.3.
  - ▶ Updating estimates ( $v_{i+1}$ ) on the basis of other estimates ( $v_i$ ) is often called **bootstrapping**.
- ▶ The Richardson iteration can be interpreted as a gradient descent algorithm for solving (3.2).
- ▶ This leads to **synchronous, full backups** of the entire state space  $\mathcal{X}$ .
- ▶ Also called **expected update** because it is based on the expectation over all possible next states (utilizing full model knowledge).
- ▶ In subsequent lectures, the expected update will be supplemented by data-driven samples from the environment.

# Iterative policy evaluation example: forest tree MDP

Let's reuse the forest tree MDP example from Fig. 2.10 with *fifty-fifty policy*:

$$\mathcal{P}_{xx'}^\pi = \begin{bmatrix} 0 & \frac{1-\alpha}{2} & 0 & \frac{1+\alpha}{2} \\ 0 & 0 & \frac{1-\alpha}{2} & \frac{1+\alpha}{2} \\ 0 & 0 & \frac{1-\alpha}{2} & \frac{1+\alpha}{2} \\ 0 & 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{r}_{\mathcal{X}}^\pi = \begin{bmatrix} 0.5 \\ 1 \\ 2 \\ 0 \end{bmatrix}.$$

$i$	$v_i(x=1)$	$v_i(x=2)$	$v_i(x=3)$	$v_i(x=4)$
0	0	0	0	0
1	0.5	1	2	0
2	0.82	1.64	2.64	0
3	1.03	1.85	2.85	0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\infty$	1.12	1.94	2.94	0

Tab. 3.1: Policy evaluation by Richardson iteration (3.12) for forest tree MDP with  $\gamma = 0.8$  and  $\alpha = 0.2$

## Variant: in-place updates

Instead of applying (3.12) to the entire vector  $v_{\mathcal{X},i+1}^{\pi}$  in 'one shot' (synchronous backup), an elementwise **in-place** version of the policy evaluation can be carried out:

**input:** full model of the MDP, i.e.,  $\langle \mathcal{X}, \mathcal{U}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  including policy  $\pi$

**parameter:**  $\delta > 0$  as accuracy termination threshold

**init:**  $v_0(x) \forall x \in \mathcal{X}$  arbitrary except  $v_0(x) = 0$  if  $x$  is terminal

**repeat**

$\Delta \leftarrow 0;$

**for**  $\forall x_k \in \mathcal{X}$  **do**

$\tilde{v} \leftarrow \hat{v}(x_k);$

$\hat{v}(x_k) \leftarrow \sum_{u_k \in \mathcal{U}} \pi(u_k | x_k) \left( \mathcal{R}_x^u + \gamma \sum_{x_{k+1} \in \mathcal{X}} p_{xx'}^u \hat{v}(x_{k+1}) \right);$

$\Delta \leftarrow \max(\Delta, |\tilde{v} - \hat{v}(x_k)|);$

**until**  $\Delta < \delta;$

**Algo. 3.1:** Iterative policy evaluation using in-place updates (output: estimate of  $v_{\mathcal{X}}^{\pi}$ )

# In-place policy evaluation updates for forest tree MDP

- ▶ In-place algorithms allow to update states in a beneficial order.
- ▶ May converge faster than regular Richardson iteration if state update order is chosen wisely (sweep through state space).
- ▶ For forest tree MDP: reverse order, i.e., start with  $x = 4$ .
- ▶ As can be seen in Tab. 3.2 the in-place updates especially converge faster for the 'early states'.

$i$	$v_i(x = 1)$	$v_i(x = 2)$	$v_i(x = 3)$	$v_i(x = 4)$
0	0	0	0	0
1	1.03	1.64	2	0
2	1.09	1.85	2.64	0
3	1.11	1.91	2.85	0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\infty$	1.12	1.94	2.94	0

Tab. 3.2: In-place updates for forest tree MDP

# Table of contents

- 1 Introduction
- 2 Policy evaluation
- 3 Policy improvement**
- 4 Policy and value iteration
- 5 Further aspects

# General idea on policy improvement

- ▶ If we know  $v_\pi$  of a given MDP, how to improve the policy?
- ▶ The simple idea of policy improvement is:
  - ▶ Consider a new (non-policy conform) action  $u \neq \pi(x_k)$ .
  - ▶ Follow thereafter the current policy  $\pi$ .
  - ▶ Check the action value of this 'new move'. If it is better than the 'old' value, take it:

$$q_\pi(x_k, u_k) = \mathbb{E} [R_{k+1} + \gamma v_\pi(X_{k+1}) | X_k = x_k, U_k = u_k] . \quad (3.13)$$

## Theorem 3.1: Policy improvement

If for any deterministic policy pair  $\pi$  and  $\pi'$

$$q_\pi(x, \pi'(x)) \geq v_\pi(x) \quad \forall x \in \mathcal{X} \quad (3.14)$$

applies, then the policy  $\pi'$  must be as good as or better than  $\pi$ . Hence, it obtains greater or equal expected return

$$v_{\pi'}(x) \geq v_\pi(x) \quad \forall x \in \mathcal{X}. \quad (3.15)$$

# Proof of policy improvement theorem

Start with (3.14) and recursively reapply (3.13):

$$\begin{aligned}v_{\pi}(x_k) &\leq q_{\pi}(x_k, \pi'(x_k)), \\&= \mathbb{E} [R_{k+1} + \gamma v_{\pi}(X_{k+1}) | X_k = x_k, U_k = \pi'(x_k)], \\&= \mathbb{E}_{\pi'} [R_{k+1} + \gamma v_{\pi}(X_{k+1}) | X_k = x_k], \\&\leq \mathbb{E}_{\pi'} [R_{k+1} + \gamma q_{\pi}(x_{k+1}, \pi'(x_{k+1})) | X_k = x_k], \\&= \mathbb{E}_{\pi'} [R_{k+1} + \gamma \mathbb{E}_{\pi'} [R_{k+2} + \gamma v_{\pi}(X_{k+2}) | X_{k+1}, \pi'(x_{k+1})] | X_k = x_k], \\&= \mathbb{E}_{\pi'} [R_{k+1} + \gamma R_{k+2} + \gamma^2 v_{\pi}(X_{k+2}) | X_k = x_k], \\&\leq \mathbb{E}_{\pi'} [R_{k+1} + \gamma R_{k+2} + \gamma^2 R_{k+3} + \gamma^3 v_{\pi}(X_{k+3}) | X_k = x_k], \\&\vdots \\&\leq \mathbb{E}_{\pi'} [R_{k+1} + \gamma R_{k+2} + \gamma^2 R_{k+3} + \gamma^3 R_{k+4} + \dots | X_k = x_k], \\&= v_{\pi'}(x_k).\end{aligned}\tag{3.16}$$

# Greedy policy improvement (1)

- ▶ So far, policy improvement addressed only changing the policy at a single state.
- ▶ Now, extend this scheme to all states by selecting the best action according to  $q_\pi(x_k, u_k)$  in every state (**greedy policy improvement**):

$$\begin{aligned}\pi'(x_k) &= \arg \max_{u_k \in \mathcal{U}} q_\pi(x_k, u_k), \\ &= \arg \max_{u_k \in \mathcal{U}} \mathbb{E} [R_{k+1} + \gamma v_\pi(X_{k+1}) | X_k = x_k, U_k = u_k], \\ &= \arg \max_{u_k \in \mathcal{U}} \mathcal{R}_x^u + \gamma \sum_{x_{k+1} \in \mathcal{X}} p_{xx'}^u v_\pi(x_{k+1}).\end{aligned}\tag{3.17}$$

## Greedy policy improvement (2)

- ▶ Each greedy policy improvement takes the best action in a one-step look-ahead search and, therefore, satisfies Theo. 3.1.
- ▶ If after a policy improvement step  $v_\pi(x_k) = v_{\pi'}(x_k)$  applies, it follows:

$$\begin{aligned} v_{\pi'}(x_k) &= \max_{u_k \in \mathcal{U}} \mathbb{E} [R_{k+1} + \gamma v_{\pi'}(X_{k+1}) | X_k = x_k, U_k = u_k], \\ &= \max_{u_k \in \mathcal{U}} \mathcal{R}_x^u + \gamma \sum_{x_{k+1} \in \mathcal{X}} p_{xx'}^u v_{\pi'}(x_{k+1}). \end{aligned} \tag{3.18}$$

- ▶ This is the Bellman optimality equation, which guarantees that  $\pi' = \pi$  must be optimal policies.
- ▶ Although proof for policy improvement theorem was presented for deterministic policies, transfer to stochastic policies  $\pi(u_k | x_k)$  is possible.
- ▶ Takeaway message: **policy improvement theorem guarantees finding optimal policies in finite MDPs** (e.g., by DP).

# Table of contents

- 1 Introduction
- 2 Policy evaluation
- 3 Policy improvement
- 4 Policy and value iteration**
- 5 Further aspects

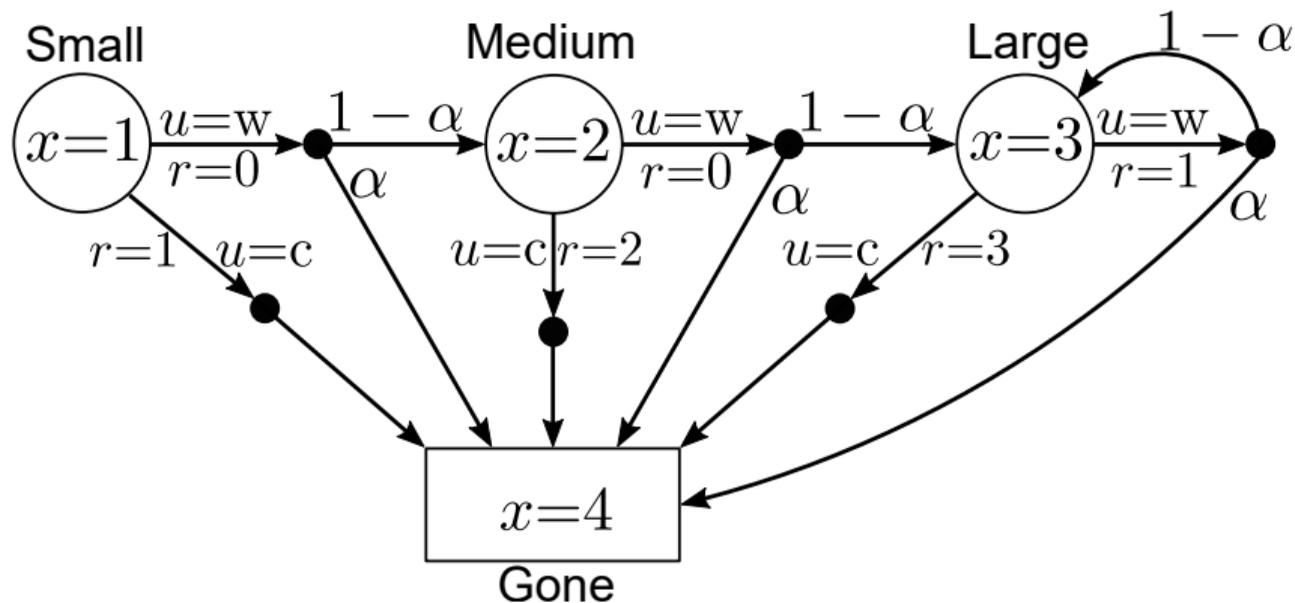
# Concept of policy iteration

- ▶ Policy iteration **combines the previous policy evaluation and policy improvement** in an iterative sequence:

$$\pi_0 \rightarrow v_{\pi_0} \rightarrow \pi_1 \rightarrow v_{\pi_1} \rightarrow \cdots \pi^* \rightarrow v_{\pi^*} \quad (3.19)$$

- ▶ Evaluate  $\rightarrow$  improve  $\rightarrow$  evaluate  $\rightarrow$  improve ...
- ▶ In the 'classic' policy iteration, each policy evaluation step in (3.19) is fully executed, i.e., for each policy  $\pi_i$  an exact estimate of  $v_{\pi_i}$  is provided either by iterative policy evaluation with a sufficiently high number of steps or by any other method that fully solves (3.2).

# Policy iteration example: forest tree MDP (1)



- ▶ Two actions possible in each state:
  - ▶ Wait  $u = w$ : let the tree grow.
  - ▶ Cut  $u = c$ : gather the wood.

## Policy iteration example: forest tree MDP (2)

Assume  $\alpha = 0.2$  and  $\gamma = 0.8$  and start with 'tree hater' initial policy:

- 1  $\pi_0 = \pi(u_k = \mathbf{c} | x_k) \quad \forall x_k \in \mathcal{X}$ .
- 2 Policy evaluation:  $v_{\mathcal{X}}^{\pi_0} = [1 \quad 2 \quad 3 \quad 0]^T$
- 3 Greedy policy improvement:

$$\begin{aligned}\pi_1(x_k) &= \arg \max_{u_k \in \mathcal{U}} \mathbb{E} [R_{k+1} + \gamma v_{\pi_0}(X_{k+1}) | X_k = x_k, U_k = u_k], \\ &= \{\pi(u_k = \mathbf{w} | x_k = 1), \pi(u_k = \mathbf{c} | x_k = 2), \pi(u_k = \mathbf{c} | x_k = 3)\}\end{aligned}$$

- 4 Policy evaluation:  $v_{\mathcal{X}}^{\pi_1} = [1.28 \quad 2 \quad 3 \quad 0]^T$
- 5 Greedy policy improvement:

$$\begin{aligned}\pi_2(x_k) &= \arg \max_{u_k \in \mathcal{U}} \mathbb{E} [R_{k+1} + \gamma v_{\pi_1}(X_{k+1}) | X_k = x_k, U_k = u_k], \\ &= \{\pi(u_k = \mathbf{w} | x_k = 1), \pi(u_k = \mathbf{c} | x_k = 2), \pi(u_k = \mathbf{c} | x_k = 3)\}, \\ &= \pi_1(x_k) \\ &= \pi^*\end{aligned}$$

# Policy iteration example: forest tree MDP (3)

Assume  $\alpha = 0.2$  and  $\gamma = 0.8$  and start with 'tree lover' initial policy:

1  $\pi_0 = \pi(u_k = \mathbf{w} | x_k) \quad \forall x_k \in \mathcal{X}.$

2 Policy evaluation:  $v_{\mathcal{X}}^{\pi_0} = [1.14 \quad 1.78 \quad 2.78 \quad 0]^T$

3 Greedy policy improvement:

$$\begin{aligned}\pi_1(x_k) &= \arg \max_{u_k \in \mathcal{U}} \mathbb{E} [R_{k+1} + \gamma v_{\pi_0}(X_{k+1}) | X_k = x_k, U_k = u_k], \\ &= \{\pi(u_k = \mathbf{w} | x_k = 1), \pi(u_k = \mathbf{c} | x_k = 2), \pi(u_k = \mathbf{c} | x_k = 3)\}\end{aligned}$$

4 Policy evaluation:  $v_{\mathcal{X}}^{\pi_1} = [1.28 \quad 2 \quad 3 \quad 0]^T$

5 Greedy policy improvement:

$$\begin{aligned}\pi_2(x_k) &= \arg \max_{u_k \in \mathcal{U}} \mathbb{E} [R_{k+1} + \gamma v_{\pi_1}(X_{k+1}) | X_k = x_k, U_k = u_k], \\ &= \{\pi(u_k = \mathbf{w} | x_k = 1), \pi(u_k = \mathbf{c} | x_k = 2), \pi(u_k = \mathbf{c} | x_k = 3)\}, \\ &= \pi_1(x_k) \\ &= \pi^*\end{aligned}$$

## Value iteration (1)

- ▶ Policy iteration involves full policy evaluation steps between policy improvements.
- ▶ In large state-space MDPs the full policy evaluation may be numerically very costly.
- ▶ **Value iteration**: One step iterative policy evaluation followed by policy improvement.
- ▶ Allows simple update rule which **combines policy improvement with truncated policy evaluation in a single step**:

$$\begin{aligned} v_{i+1}(x_k) &= \max_{u_k \in \mathcal{U}} \mathbb{E} [R_{k+1} + \gamma v_i(X_{k+1}) | X_k = x_k, U_k = u_k], \\ &= \max_{u_k \in \mathcal{U}} \mathcal{R}_x^u + \gamma \sum_{x_{k+1} \in \mathcal{X}} p_{xx'}^u v_i(x_{k+1}). \end{aligned} \tag{3.20}$$

## Value iteration (2)

**input:** full model of the MDP, i.e.,  $\langle \mathcal{X}, \mathcal{U}, \mathcal{P}, \mathcal{R}, \gamma \rangle$

**parameter:**  $\delta > 0$  as accuracy termination threshold

**init:**  $v_0(x) \forall x \in \mathcal{X}$  arbitrary except  $v_0(x) = 0$  if  $x$  is terminal

**repeat**

$\Delta \leftarrow 0;$

**for**  $\forall x_k \in \mathcal{X}$  **do**

$\tilde{v} \leftarrow \hat{v}(x_k);$

$\hat{v}(x_k) \leftarrow \max_{u_k \in \mathcal{U}} \left( \mathcal{R}_x^u + \gamma \sum_{x_{k+1} \in \mathcal{X}} p_{xx'}^u \hat{v}(x_{k+1}) \right);$

$\Delta \leftarrow \max(\Delta, |\tilde{v} - \hat{v}(x_k)|);$

**until**  $\Delta < \delta;$

**output:** deterministic policy  $\pi \approx \pi^*$ , such that

$\pi(x_k) \leftarrow \arg \max_{u_k \in \mathcal{U}} \left( \mathcal{R}_x^u + \gamma \sum_{x_{k+1} \in \mathcal{X}} p_{xx'}^u \hat{v}(x_{k+1}) \right);$

**Algo. 3.2:** Value iteration (note: compared to policy iteration, value iteration does not require an initial policy but only a state-value guess)

# Value iteration example: forest tree MDP

- ▶ Assume again  $\alpha = 0.2$  and  $\gamma = 0.8$ .
- ▶ Similar to in-place update policy evaluation, reverse order and start value iteration with  $x = 4$ .
- ▶ As shown in Tab. 3.3 value iteration converges in one step (for the given problem) to the optimal state value.

$i$	$v_i(x = 1)$	$v_i(x = 2)$	$v_i(x = 3)$	$v_i(x = 4)$
0	0	0	0	0
1	1.28	2	3	0
*	1.28	2	3	0

Tab. 3.3: Value iteration for forest tree MDP

# Table of contents

- 1 Introduction
- 2 Policy evaluation
- 3 Policy improvement
- 4 Policy and value iteration
- 5 Further aspects**

# Summarizing DP algorithms

- ▶ All DP algorithms are based on the state value  $v(x)$ .
  - ▶ Complexity is  $\mathcal{O}(m \cdot n^2)$  for  $m$  actions and  $n$  states.
  - ▶ Evaluate all  $n^2$  state transitions while considering up to  $m$  actions per state.
- ▶ Could be also applied to action values  $q(x, u)$ .
  - ▶ Complexity is inferior with  $\mathcal{O}(m^2 \cdot n^2)$ .
  - ▶ There are up to  $m^2$  action values which require  $n^2$  state transition evaluations each.

Problem	Relevant Equations	Algorithm
prediction	Bellman expectation eq.	policy evaluation
control	Bellman expectation eq. & greedy policy improvement	policy iteration
control	Bellman optimality eq.	value iteration

Tab. 3.4: Short overview addressing the treated DP algorithms

# Curse of dimensionality

- ▶ DP is much more efficient than an exhaustive search over all  $n$  states and  $m$  actions in finite MDPs in order to find an optimal policy.
  - ▶ Exhaustive search for deterministic policies:  $m^n$  evaluations.
  - ▶ DP results in polynomial complexity regarding  $m$  and  $n$ .
- ▶ Nevertheless, DP uses full-width backups:
  - ▶ For each state update, every successor state and action is considered.
  - ▶ While utilizing full knowledge of the MDP structure.
- ▶ Hence, DP is can be effective up to medium-sized MDPs (i.e., million finite states)
- ▶ For large problems DP suffers from the **curse of dimensionality**:
  - ▶ Single update step may become computational infeasible.
  - ▶ Also: if continuous states need quantization, number of finite states  $n$  grows exponentially with the number of state variables (assuming fixed number of discretization levels).

# Generalized policy iteration (GPI)

- ▶ Almost all RL methods are well-described as GPI.
- ▶ **Push-pull**: Improving the policy will deteriorate value estimation.
- ▶ Well balanced **trade-off between evaluating and improving** is required.

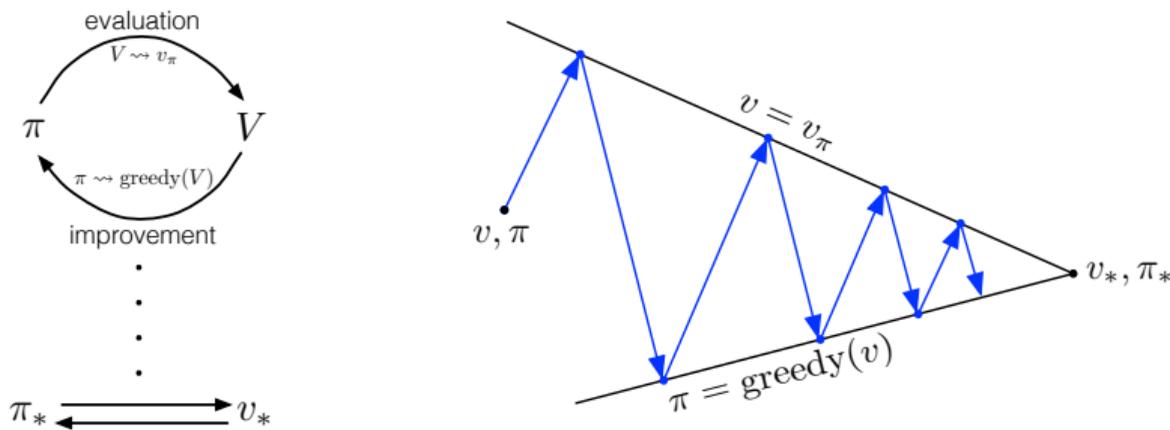


Fig. 3.4: Interpreting generalized policy iteration to switch back and forth between (arbitrary) evaluations and improvement steps (source: R. Sutton and G. Barto, Reinforcement learning: an introduction, 2018, [CC BY-NC-ND 2.0](#))

## Summary: what you've learned in this lecture

- ▶ DP is applicable for prediction and control problems in MDPs.
- ▶ But requires always full knowledge about the environment (i.e., it is a model-based solution).
- ▶ DP is more efficient than exhaustive search.
- ▶ But suffers from the curse of dimensionality for large MDPs.
- ▶ (Iterative) policy evaluations and (greedy) improvements solve MDPs.
- ▶ Both steps can be combined via value iteration.
- ▶ This idea of (generalized) policy iteration is a basic scheme of RL.
- ▶ Implementing DP algorithms comes with many degrees of freedom regarding the update order.