

Reinforcement Learning Part I: Discrete State-Action Optimal Control ...with Learned Models

Jan Peters
Gerhard Neumann
Guilherme Maeda

Motivation for optimal decision making in robotics

Typically, **imitation is not enough**

Imperfect demonstrations

Correspondance problem

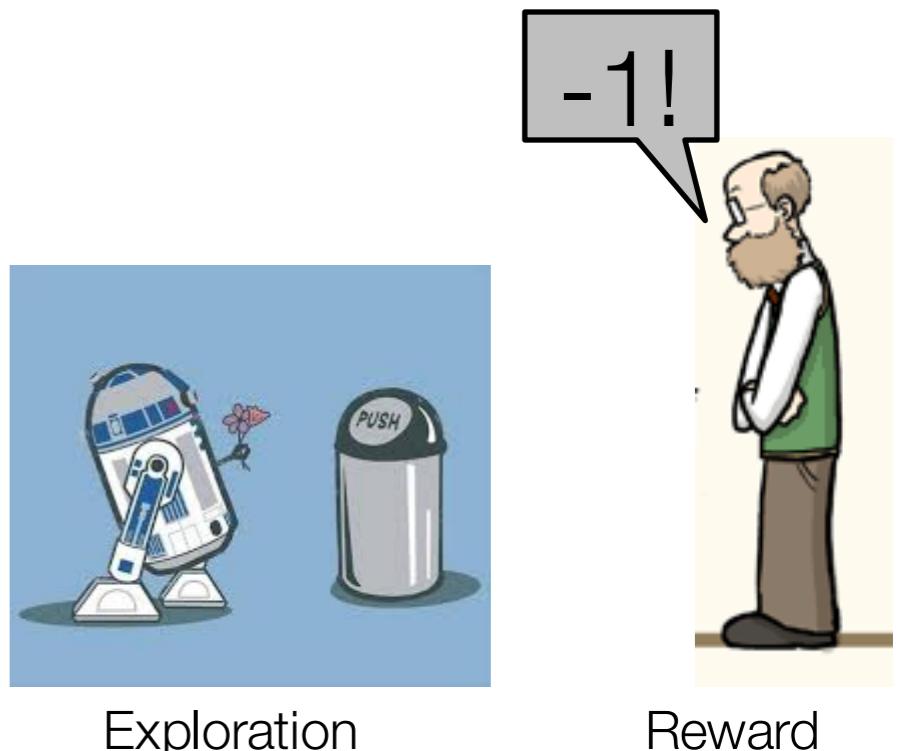
We can not demonstrate everything

Hence, we need **self-improvement!**

The robot explores by trial and error

We give evaluative feedback → reward

Today, we are going to look at the problem of how to **take optimal decision that maximize the reward**



Exploration

Reward



Outline of the Lecture

1. Introduction

- Example of Discrete State-Action Control
- Formalization of Optimal Control as Markov Decision Process

2. Finite-Horizon Optimal Control

- Value Iteration with a For-Loop

3. Infinite Horizon Value Iteration

- Value Iteration with a Repeat-Until-Loop

4. Infinite Horizon Policy Iteration

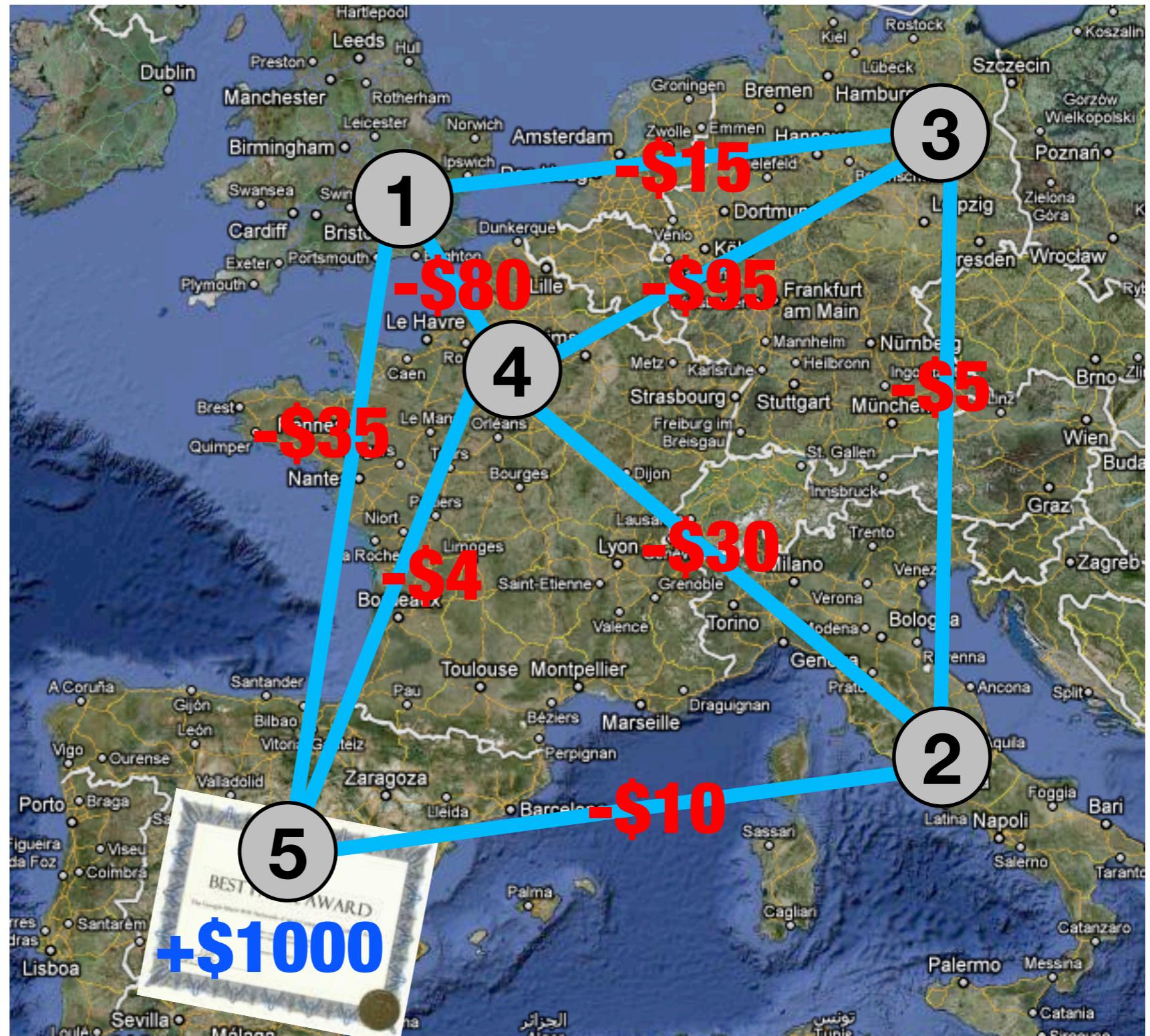
- Policy Evaluation: Generate the value function for a fixed policy
- Policy Improvement: Compute a better policy

Illustration of basic idea...



You have won
a Best-Paper
Award in
Madrid!

What is the
Optimal
Policy to
Collect it?



Dynamic Programming



“An optimal sequence of controls in a multistage optimization problem has the property that **whatever the initial stage, state and controls are, the remaining controls must constitute an optimal sequence of decisions for the remaining problem** with stage and state resulting from previous controls considered as initial conditions.”

Dynamic Programming

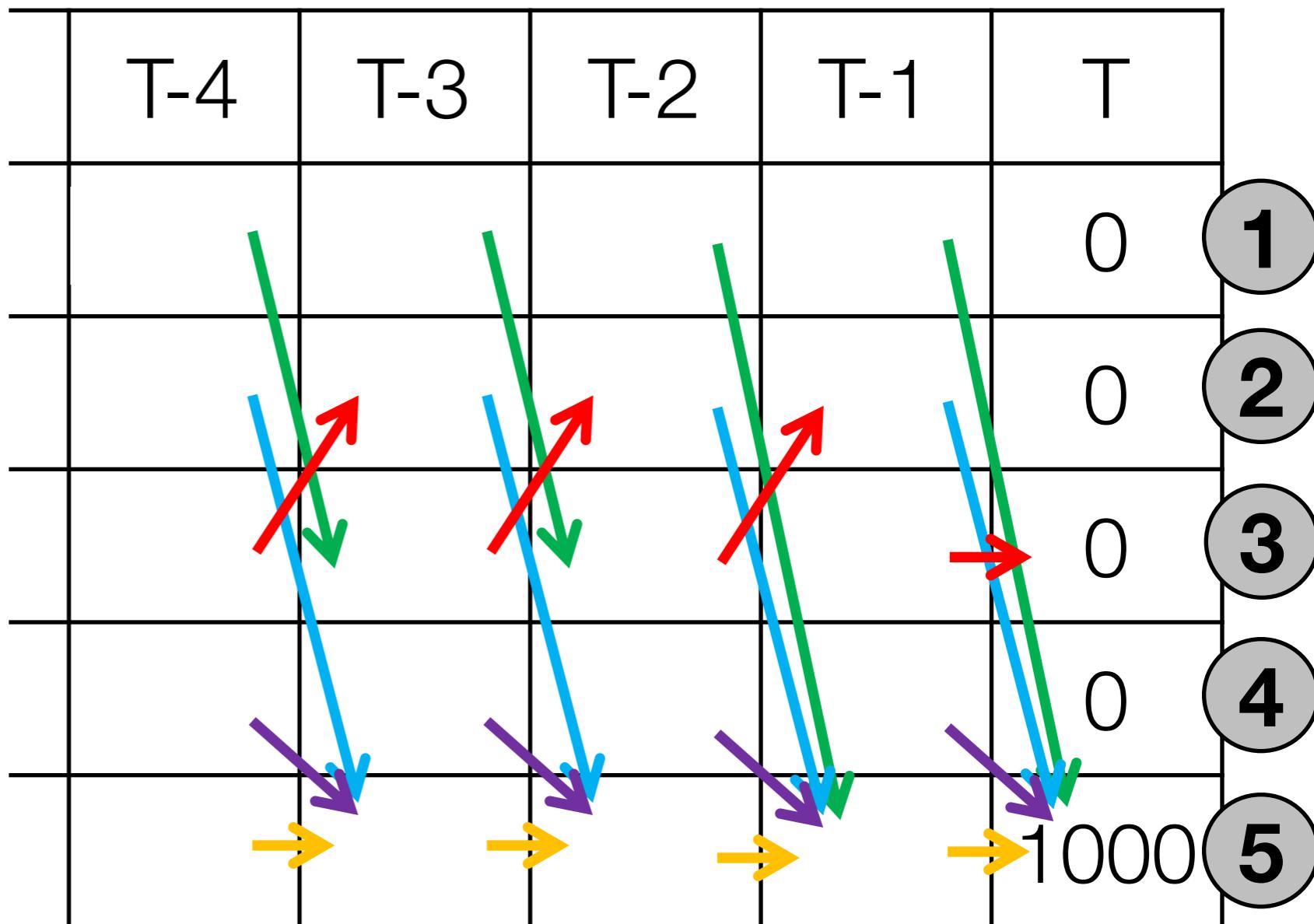
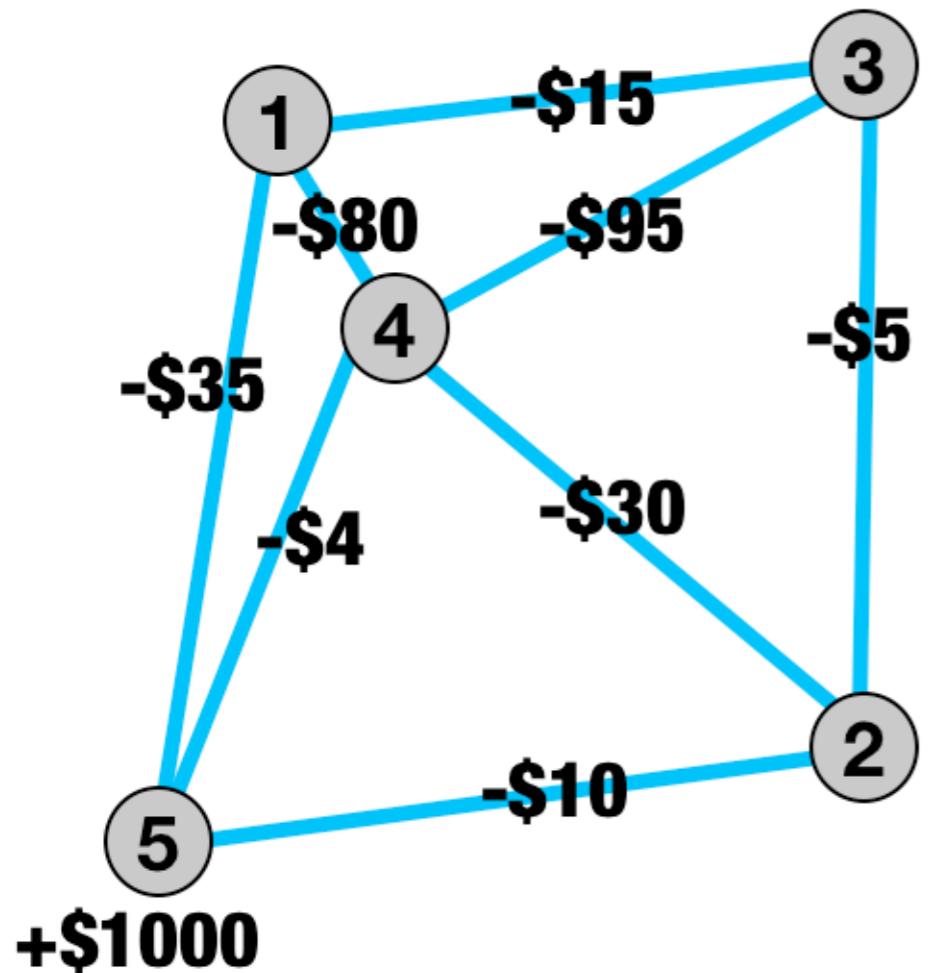


Idea 1: If the optimal solution is broken in small parts,
each part must be optimal

Idea 2: reuse solutions that were already computed (we
will store them as V or Q)

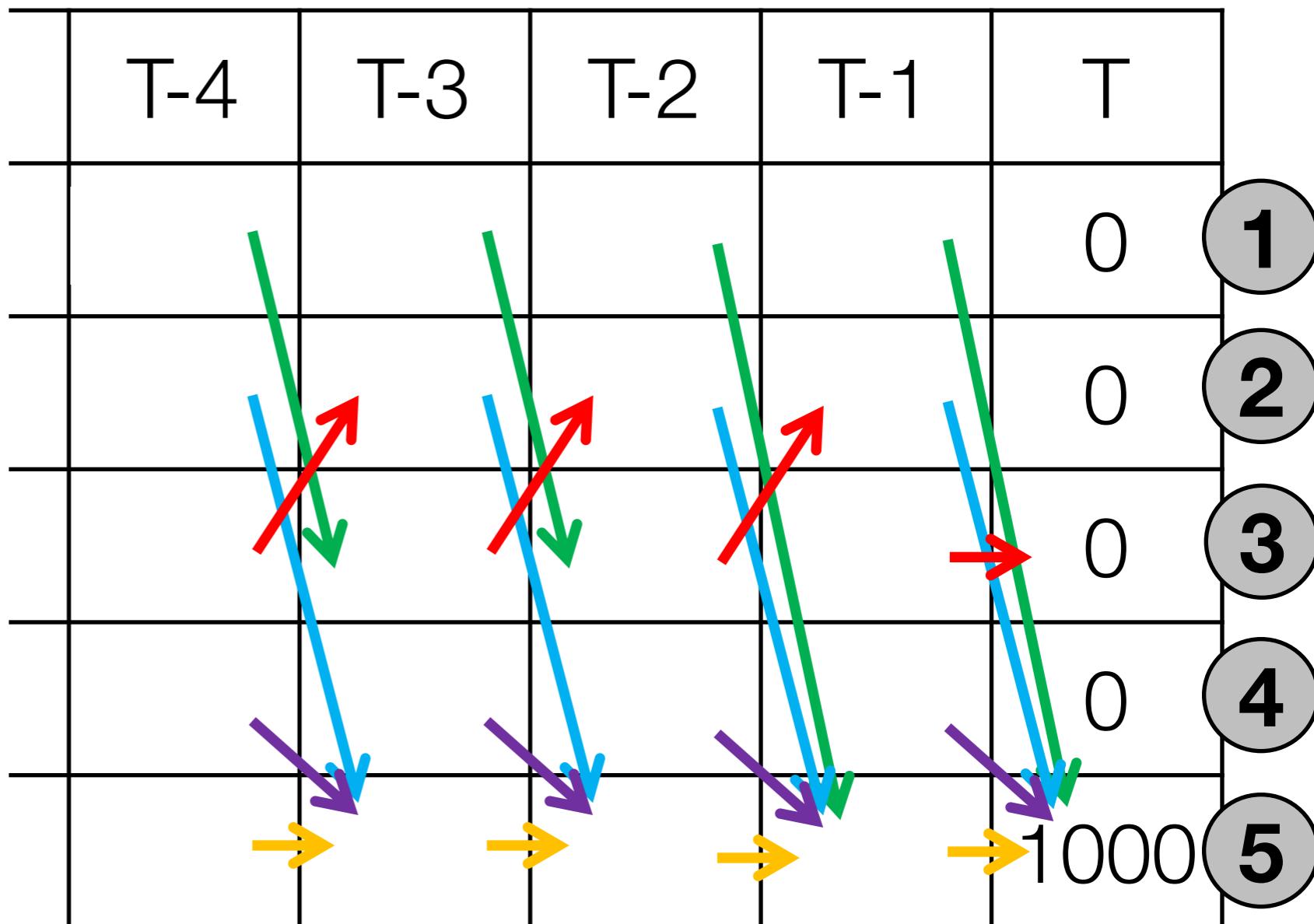
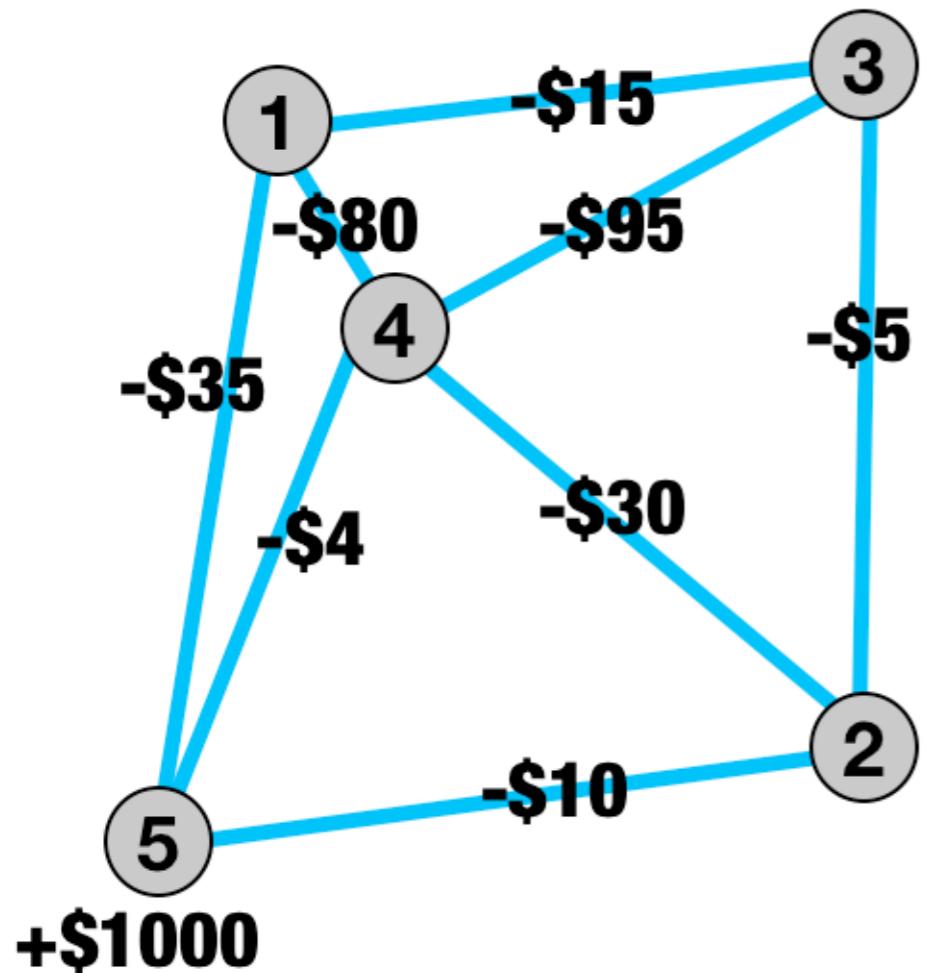


Let's Try this Example!





Let's Try this Example!

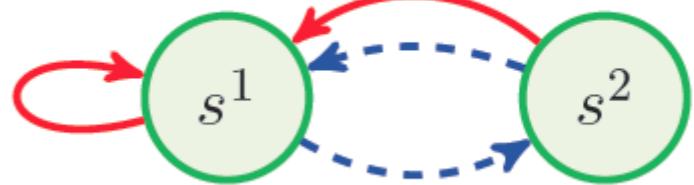




Markov Decision Problems (MDP)

A (non-stationary) **MDP** is defined by:

- its state space $s \in \mathcal{S}$
- its action space $a \in \mathcal{A}$
- its transition dynamics $\mathcal{P}_t(s_{t+1}|s_t, a_t)$
- its reward function $r_t(s, a)$
- and its initial state probabilities $\mu_0(s)$



Markov property:

$$\mathcal{P}_t(s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}, \dots) = \mathcal{P}_t(s_{t+1}|s_t, a_t)$$

- Transition dynamics depends on only on the current time step



Outline of the Lecture

1. Introduction

- Example of Discrete State-Action Control
- Formalization of Optimal Control as Markov Decision Process

2. Finite-Horizon Optimal Control

- Value Iteration with a For-Loop

3. Infinite Horizon Value Iteration

- Value Iteration with a Repeat-Until-Loop

4. Infinite Horizon Policy Iteration

- Policy Evaluation: Generate the value function for a fixed policy
- Policy Improvement: Compute a better policy



Finite Horizon Objective

The goal of the agent is to find a policy $\pi(a|s)$ that maximizes its expected return J_π **for a finite time horizon**

Finite Horizon T: Accumulated expected reward for T steps

$$J_\pi = \mathbb{E}_{\mu_0, \mathcal{P}, \pi} \left[\sum_{t=1}^{T-1} r_t(s_t, a_t) + r_T(s_T) \right]$$

$r_T(s_T)$... final reward



Algorithmic Description of Value Iteration

Init: $V_T^*(s) \leftarrow r_T(s)$, $t = T$

Repeat $t = t - 1$

Compute Q-Function for time step t (for each state action pair)

$$Q_t^*(s, a) = r_t(s, a) + \sum_{s'} P_t(s'|s, a) V_{t+1}^*(s')$$

Compute V-Function for time step t (for each state)

$$V_t^*(s) = \max_a Q_t^*(s, a)$$

Until $t = 1$

Return: Optimal policy for each time step

$$\pi_t^*(s) = \operatorname{argmax}_a Q_t^*(s, a)$$

[check animation]



Value Iteration for Finite Horizon

So how does **dynamic programming** work now?

- Start with last layer... (no transition)

$$V_T^*(\mathbf{s}) = r_T(\mathbf{s})$$

- Iterate **backwards** in time

$$V_t^*(\mathbf{s}) = \max_{\mathbf{a}} (r_t(\mathbf{s}_t, \mathbf{a}_t) + \mathbb{E}_{\mathcal{P}} [V_{t+1}^*(\mathbf{s}_{t+1}) | \mathbf{s}_t, \mathbf{a}_t])$$

- The optimal value function/policy for time step t is obtained after $T - t + 1$ iterations

$$V_T^*(\mathbf{s}_T) \rightarrow V_{T-1}^*(\mathbf{s}_{T-1}) \rightarrow \dots \rightarrow V_1^*(\mathbf{s}_1)$$



What does a finite life-time T imply?

In the finite horizon case, the **time index is part of the state**

- It matters, how many time steps are left
- We can only visit **each state (including time index) once!**
- We have **a layered / multi stage decision** problem
- **The optimal policy is time-dependent**

$$\pi_t^*(\mathbf{a}|\mathbf{s}) = \pi^*(\mathbf{a}|\mathbf{s}, t)$$

- The reward function and the transition model can be time-dependent, i.e.,

$$r_t(\mathbf{s}, \mathbf{a}) \text{ and } \mathcal{P}_t(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$$



Outline of the Lecture

1. Introduction

- Example of Discrete State-Action Control
- Formalization of Optimal Control as Markov Decision Process

2. Finite-Horizon Optimal Control

- Value Iteration with a For-Loop

3. Infinite Horizon Value Iteration

- Value Iteration with a Repeat-Until-Loop

4. Infinite Horizon Policy Iteration

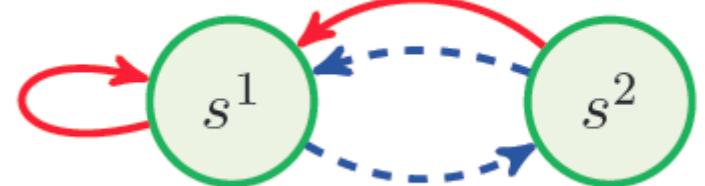
- Policy Evaluation: Generate the value function for a fixed policy
- Policy Improvement: Compute a better policy



Markov Decision Problems (MDP)

A stationary **MDP** is defined by:

- its state space $s \in \mathcal{S}$
- its action space $a \in \mathcal{A}$
- its transition dynamics $\mathcal{P}(s_{t+1}|s_t, a_t)$
- its reward function $r(s, a) = r_t(s, a)$ and $\mathcal{P}_t(s_{t+1}|s_t, a_t)$
- and its initial state probabilities $\mu_0(s)$



$$r_t(s, a) \text{ and } \mathcal{P}_t(s_{t+1}|s_t, a_t)$$

Markov property:

$$\mathcal{P}(s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}, \dots) = \mathcal{P}(s_{t+1}|s_t, a_t)$$

- Transition dynamics depends on only on the current time step

So what's special for an infinite horizon,
i.e., when you live forever?



In the infinite horizon case, the **time index is not part of the state**

- Your robot lives for all it does, it does not matter, how many time steps are left
- optimal policy is time-independent $\pi_t^*(a|s) = \pi^*(a|s)$
- The reward function and the transition model can no longer be time-dependent.
- There is a single, stationary value function for all times.
- We have two different approaches to learning:
 1. **Value Iteration:** Use value iteration from before, choose a large T for which the value function converges.
 2. **Policy Iteration:** Learn a value function for the current policy. Then update the policy and learn a new value function. Repeat.



Optimality Objective

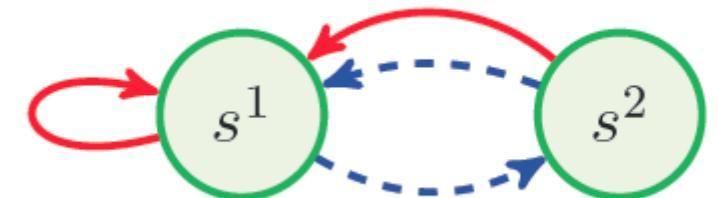
The goal of the agent is to find an optimal policy π^* that maximizes its **expected long term reward** J_π

$$\pi^* = \operatorname{argmax}_\pi J_\pi, \quad J_\pi = \mathbb{E}_{\mu_0, \mathcal{P}, \pi} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

- $0 \leq \gamma < 1$... discount factor
- Discount Factor **trades-off long term vs. immediate reward**
- Time Horizon: Infinite

Example: Two State Problem

States: s^1, s^2



Actions: red (a^1) and blue (a^2) edges

Transition:

$$\mathcal{P}(s^1|s^1, a^1) = 1, \mathcal{P}(s^2|s^1, a^1) = 0, \mathcal{P}(s^1|s^1, a^2) = 0, \mathcal{P}(s^2|s^1, a^2) = 1$$

$$\mathcal{P}(s^1|s^2, a^1) = 1, \mathcal{P}(s^2|s^2, a^1) = 0, \mathcal{P}(s^1|s^2, a^2) = 1, \mathcal{P}(s^2|s^2, a^2) = 0$$

Rewards: $r(s^1) = 1, r(s^2) = 0$

Policy: What is the optimal infinite horizon policy?



How do we find an optimal policy?

Typically done iteratively:

- **Policy Evaluation:**

Policy Evaluation:
Estimate the Value Function V^π

Policy Improvement:
Update the Policy π

Estimate quality of states (and actions) with current policy

- **Policy Improvement:**

Improve policy by taking actions with the highest quality

Such iterations are called **Policy Iteration**

→ will lead to the Bellman Equation → solved by Value Iteration



Value functions and State-Action Value Functions

Value function $V^\pi(s)$:

Long-term reward for state s when following policy $\pi(a|s)$

$$V^\pi(s) = E_{\mathcal{P}, \pi} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s \right]$$

→ **Quality measure** for state s

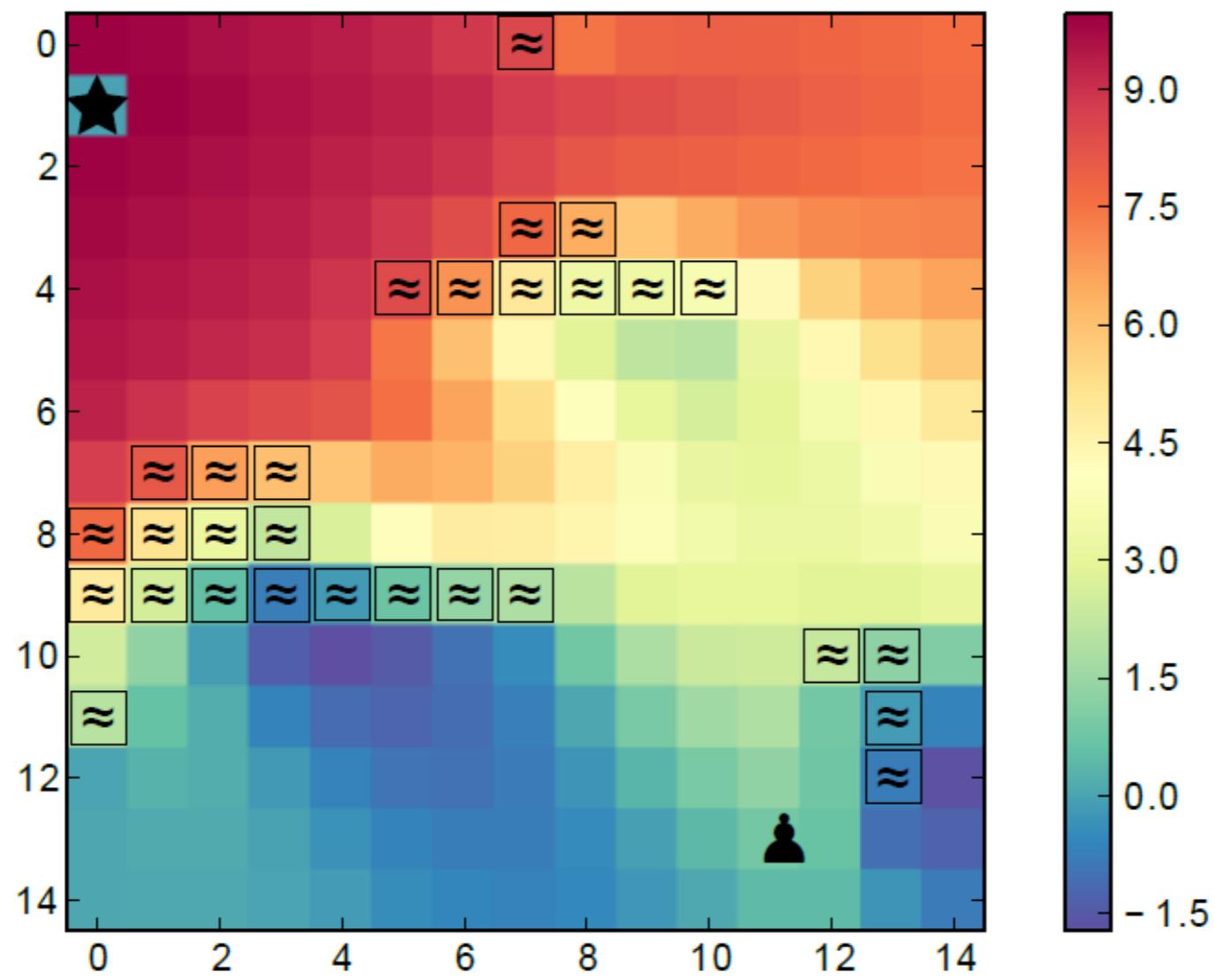
„How good“ is it to be in state s under policy $\pi(a|s)$?



Value functions

An Illustration...

Policy always goes directly to the star
Going through puddles is punished





Value functions and State-Action Value Functions

Q-function $Q^\pi(s, a)$:

Long-term reward for taking action a in state s and subsequently following policy $\pi(a|s)$

$$Q^\pi(s, a) = \mathbb{E}_{\mathcal{P}, \pi} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s, a_0 = a \right]$$

→ **Quality measure** for taking action a in state s

„How good“ is it to take action a in state s under policy $\pi(a|s)$?



Value functions and State-Action Value Functions

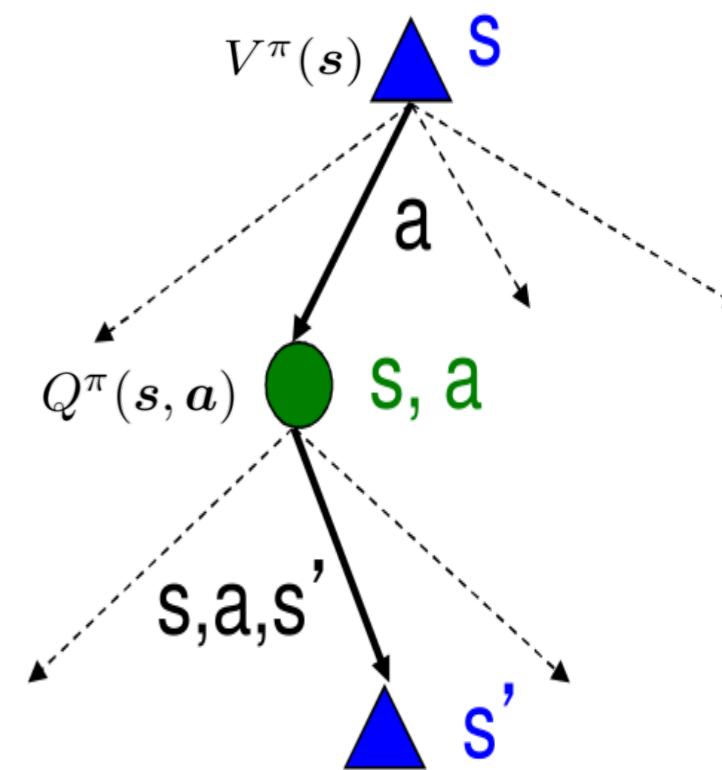
... and can be **easily computed from each other**

Computing V-Function from Q-Function

$$V^\pi(s) = \mathbb{E}_\pi [Q^\pi(s, a)|s] = \int \pi(a|s)Q^\pi(s, a)da$$

Computing Q-Function from V-Function

$$\begin{aligned} Q^\pi(s, a) &= r(s, a) + \gamma \mathbb{E}_{\mathcal{P}} [V^\pi(s')|s, a] \\ &= r(s, a) + \gamma \int \mathcal{P}(s'|s, a)V^\pi(s')ds' \end{aligned}$$



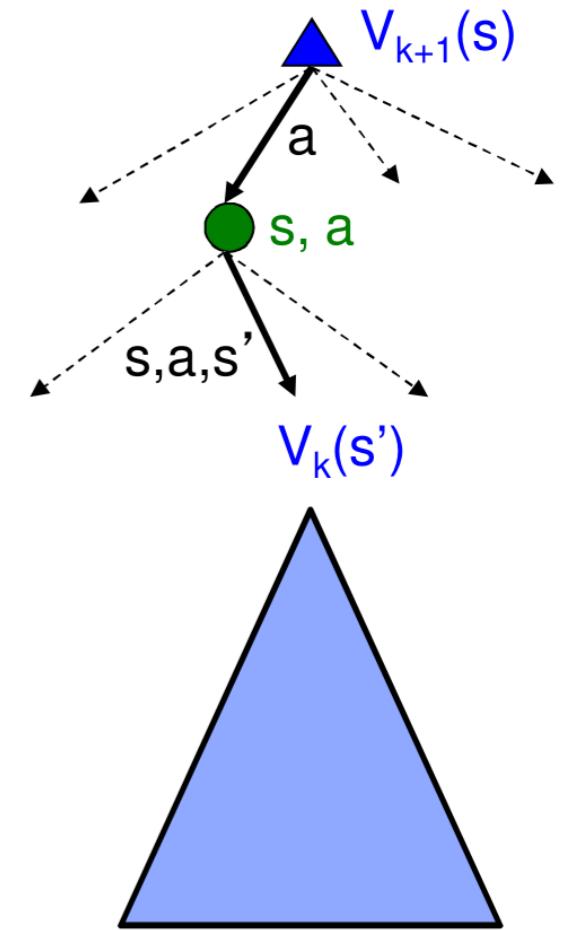


Value functions and State-Action Value Functions

... both functions can also be **estimated recursively**

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi \left[r(s, a) + \gamma \mathbb{E}_{\mathcal{P}} [V^\pi(s')] \mid s \right] \\ &= \int \pi(a|s) \left(r(s, a) + \gamma \int \mathcal{P}(s'|s, a) V^\pi(s') ds' \right) da \end{aligned}$$

$$\begin{aligned} Q^\pi(s, a) &= r(s, a) + \gamma \mathbb{E}_{\mathcal{P}, \pi} \left[Q^\pi(s', a') \mid s, a \right] \\ &= r(s, a) + \gamma \int \mathcal{P}(s'|s, a) \int \pi(a'|s') Q^\pi(s', a') da' ds' \end{aligned}$$



→ If I know the value of the next state s' , I can compute the value of the current state

Iterating these equations converges to the true V or Q function



Algorithmic Description of Policy Evaluation

Simplification: For discrete states....

Init: $V_0^\pi(s) \leftarrow 0, \forall s$ and $k = 0$

Repeat

Compute Q-Function (for each state action pair)

$$Q_{k+1}^\pi(s, a) = r(s, a) + \gamma \sum_{s'} P(s'|s, a) V_k^\pi(s')$$

Compute V-Function (for each state)

$$V_{k+1}^\pi(s) = \sum_a \pi(a|s) Q_{k+1}^\pi(s, a)$$

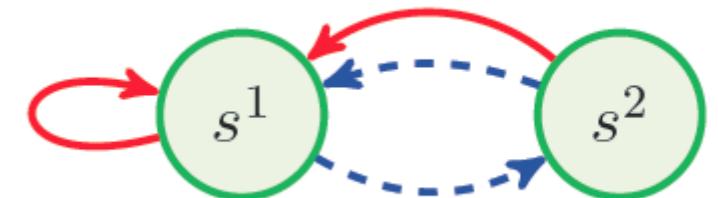
$$k = k + 1$$

until convergence

This algorithm is called Dynamic Programming!

Example: Two State Problem

States: s^1, s^2



Actions: red (a^1) and blue (a^2) edges

Transition:

$$\mathcal{P}(s^1|s^1, a^1) = 1, \mathcal{P}(s^2|s^1, a^1) = 0, \mathcal{P}(s^1|s^1, a^2) = 0, \mathcal{P}(s^2|s^1, a^2) = 1$$

$$\mathcal{P}(s^1|s^2, a^1) = 1, \mathcal{P}(s^2|s^2, a^1) = 0, \mathcal{P}(s^1|s^2, a^2) = 1, \mathcal{P}(s^2|s^2, a^2) = 0$$

Rewards: $r(s^1) = 1, r(s^2) = 0$

Policy Evaluation: What is the value function of the uniform policy?

→ HOMEWORK!



Outline of the Lecture

1. Introduction

- Example of Discrete State-Action Control
- Formalization of Optimal Control as Markov Decision Process

2. Finite-Horizon Optimal Control

- Value Iteration with a For-Loop

3. Infinite Horizon Value Iteration

- Value Iteration with a Repeat-Until-Loop

4. Infinite Horizon Policy Iteration

- Policy Evaluation: Generate the value function for a fixed policy
- Policy Improvement: Compute a better policy



How do we find an optimal policy?

Typically done iteratively:

- **Policy Evaluation:**

Estimate quality of states (and actions) with current policy

- **Policy Improvement:**

Improve policy by taking actions with the highest quality

For all states:

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

Iterating Policy Evaluation and Policy Improvement converges to the **optimal policy** and is called **Policy Iteration**



Algorithmic Description of Policy Iteration

Init: $V_0^\pi(s) \leftarrow 0, \pi \leftarrow \text{uniform}$

Repeat

Repeat $k = k + 1$

Compute Q-Function (for each state action pair)

$$Q_{k+1}^\pi(s, a) = r(s, a) + \gamma \sum_{s'} P(s'|s, a) V_k^\pi(s')$$

Compute V-Function (for each state)

$$V_{k+1}^\pi(s) = \sum_a \pi(a|s) Q_{k+1}^\pi(s, a)$$

until convergence of V

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

until convergence of policy

These are the optimal actions based on values of a suboptimal policy!



Value iteration

Can we also **stop policy evaluation before convergence** and perform a policy update?

Yes! We will still converge to the **optimal policy** !

„Extreme“ case: Stop policy evaluation **after 1 iteration**

$$V^*(s) = \max_a \left(r(s, a) + \gamma \mathbb{E}_{\mathcal{P}} [V^*(s') | s, a] \right)$$

This equation is called the **Bellman Equation**

Iterating this equation computes the **value function** $V^*(s)$ of the **optimal policy**



Value Iteration

Alternatively we can also **iterate Q-functions...**

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{\mathcal{P}} [\max_{a'} Q^*(s', a') | s, a]$$

Small side note:

Computing optimal V-Function from optimal Q-Function

$$V^*(s) = \max_a Q^*(s, a)$$

Computing optimal Q-Function from optimal V-Function

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{\mathcal{P}} [V^*(s') | s, a]$$



Algorithmic Description of Value Iteration

Init: $V_0^*(s) \leftarrow 0$

Repeat $k = k + 1$

Compute Q-Function (for each state action pair)

$$Q_{k+1}^*(s, a) = r(s, a) + \gamma \sum_{s'} P(s'|s, a) V_k^*(s')$$

Compute V-Function (for each state)

$$V_{k+1}^*(s) = \max_a Q_{k+1}^*(s, a)$$

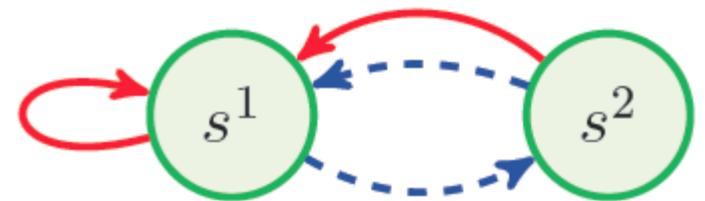
until convergence of V

[check animation]

Example: Value Iteration

- The Two state example.

→ HOMEWORK!



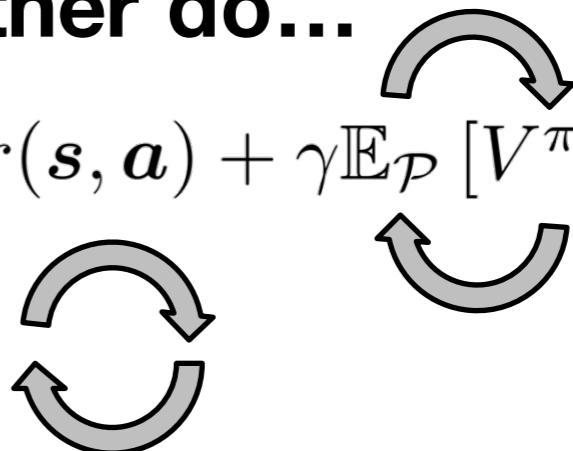


Wrap-Up: Dynamic Programming

To compute an **optimal policy** we can either do...

Policy Iteration:

$$V^\pi(s) = \mathbb{E}_\pi \left[r(s, a) + \gamma \mathbb{E}_{\mathcal{P}} [V^\pi(s')] \mid s \right]$$



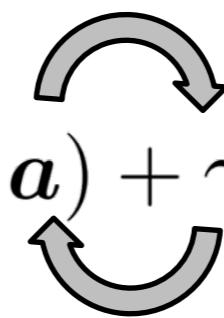
Policy Evaluation:

Policy Improvement:

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

Value Iteration:

Iterate: $V^*(s) = \max_a \left(r(s, a) + \gamma \mathbb{E}_{\mathcal{P}} [V^*(s') \mid s, a] \right)$



Get optimal policy after convergence:

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



Outline of the Lecture

1. Introduction

- Example of Discrete State-Action Control
- Formalization of Optimal Control as Markov Decision Process

2. Finite-Horizon Optimal Control

3. Infinite Horizon Value-Functions

- Policy Evaluation for a fixed policy

4. Computing an Optimal Policy for Any Value Function

- Policy Improvement
- Value iteration



Wrap-Up: Dynamic Programming

We now know how to compute optimal policies for both objectives (finite and infinite horizon)

Cool, that's all we need. Lets go home...

Wait, there is a catch!

Unfortunately, we can only do this in 2 cases

- Discrete Systems
 - Easy: integrals turn into sums
 - ...but the world is not discrete!
- Linear Systems, Quadratic Reward, Gaussian Noise (LQR) (next lecture)
 - ... but the world is not linear!



Wrap-Up: Dynamic Programming

In all other cases, we have to use approximations!

Why?

1. Representation of the V-function:

How to represent V in continuous state spaces?

2. We need to solve:

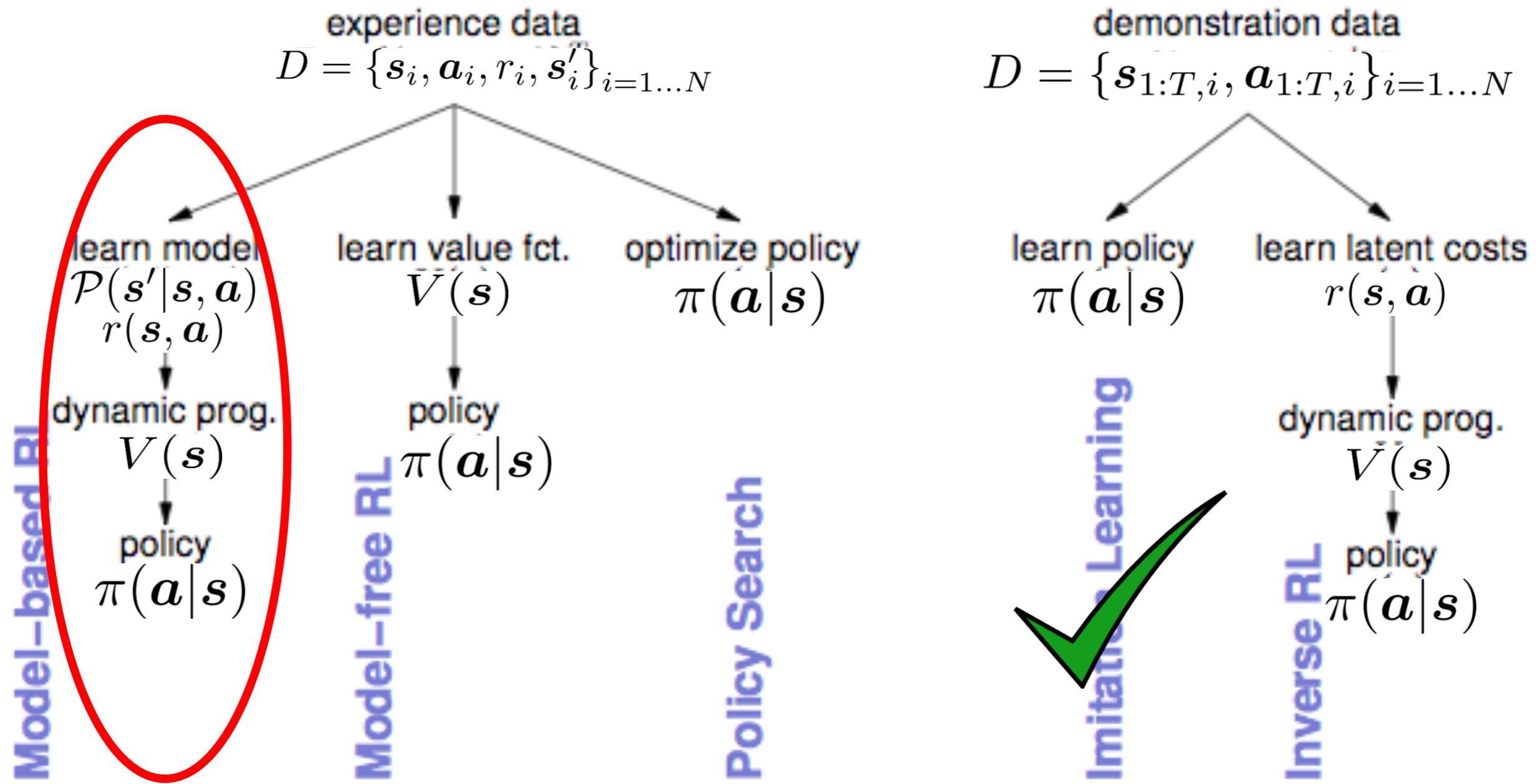
$\max_a Q^*(s, a)$: difficult in **continuous action spaces**

$\mathbb{E}_{\mathcal{P}} [V^*(s') | s, a]$: difficult for **arbitrary functions V and models \mathcal{P}**

We will hear about that in the next lectures....!



The Bigger Picture: How to learn policies



1. Next Lecture

2.

3.

4.

Discrete State-Action Optimal Control: Summary



What you should know...

- What is a MDP, a value function and a state-action value function...
- What is policy evaluation, policy improvement, policy iteration and value iteration
- The Bellman equation
- Differences of finite and infinite horizon objectives
- Why is it difficult?