# Foundations for Optimal Decision Making

**Jan Peters**
**Gerhard Neumann**

# Motivation for optimal decision making in robotics

Typically, **imitation is not enough**

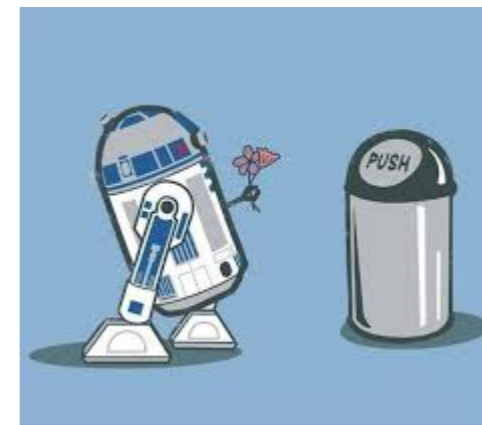    Imperfect demonstrations

    Correspondance problem

    We can not demonstrate everything

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

    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. <span style="color:red">Introduction to MDPs</span>**

**2. Value-Functions**

- Policy Evaluation for a fixed policy

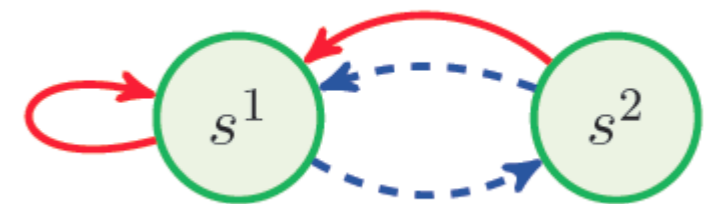**3. Computing an Optimal Policy**

- Policy Improvement

- Value iteration

**4. Infinite vs Finite Horizon**

3

# Markov Decision Processes (MDP)

A **MDP** is defined by:

- its state space $\boldsymbol{s} \in \mathcal{S}$

- its action space $\boldsymbol{a} \in \mathcal{A}$

- its transition dynamics $\mathcal{P}(\boldsymbol{s}_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t)$

- its reward function $r(\boldsymbol{s}, \boldsymbol{a})$

- and its initial state probabilities $\mu_0(\boldsymbol{s})$

Markov property:

$$\mathcal{P}(\boldsymbol{s}_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t, \boldsymbol{s}_{t-1}, \boldsymbol{a}_{t-1}, \dots) = \mathcal{P}(\boldsymbol{s}_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t)$$

- Transition dynamics depends on only of current time step

# Optimality Objective

The goal of the agent is to find an optimal policy $\pi^*$ that maximizes its <span style="color:red">expected long term reward</span> $J_\pi$

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

- $0 \leq \gamma < 1$ … discount factor

- Discount Factor <span style="color:red">trades-off long term vs. immediate reward</span>

- Time Horizon: Infinite

5

# Example: Two State Problem

**States:** $s^1, s^2$
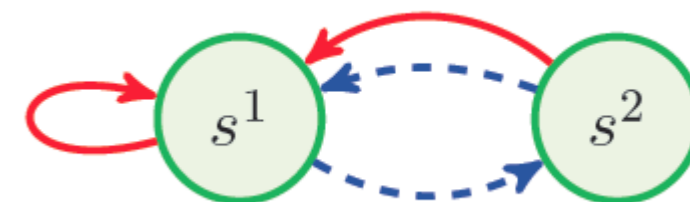


**Actions:** red ( $a^1$) and blue ( $a^2$) edges

**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 policy?**
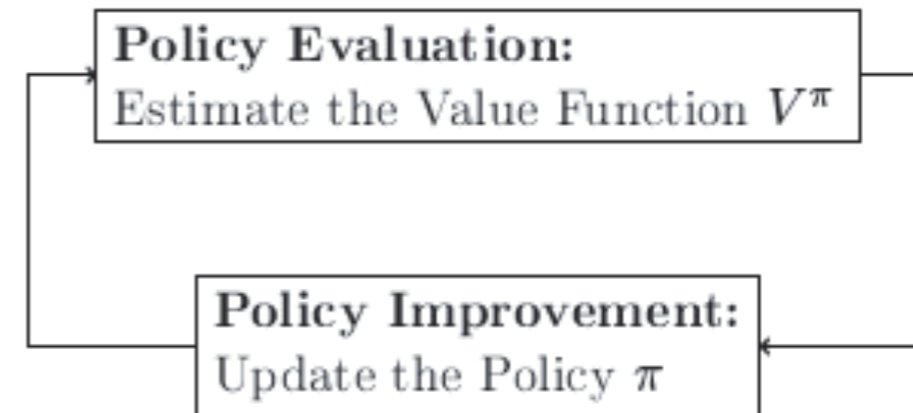
6

# How do we find an optimal policy?

**Typically done iteratively:**



**Policy Evaluation:**
Estimate the Value Function $V^{\pi}$

**Policy Improvement:**
Update the Policy $\pi$

- Policy Evaluation:

  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**

7

# Outline of the Lecture

**1. Introduction to MDPs**

**2. <span style="color:red">Value-Functions</span>**

- Policy Evaluation for a fixed policy

**3. Computing an Optimal Policy**

- Policy Improvement

- Value iteration

**4. Infinite vs Finite Horizon**

8

**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** $\pi(a|s)$ **?**

9

**An Illustration…**

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



10

**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** $\pi(a|s)$ **?**

11

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

Computing V-Function from Q-Function

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

Computing Q-Function from V-Function

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

$$= r(\boldsymbol{s}, \boldsymbol{a}) + \gamma\int \mathcal{P}(\boldsymbol{s}'|\boldsymbol{s}, \boldsymbol{a})V^\pi(\boldsymbol{s}')d\boldsymbol{s}'$$

**… both functions can also be <span style="color:red">estimated recursively</span>**

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

$$= \int \pi(a|s)\left(r(s,a) + \gamma \int \mathcal{P}(s'|s,a)V^\pi(s')ds'\right)da$$

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

$$= r(s,a) + \gamma \int \mathcal{P}(s'|s,a)\int \pi(a'|s')Q^\pi(s',a')da'ds'$$

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

**<span style="color:red">Iterating these equations</span>** converges to the true V or Q function

13

# Algorithmic Description of Policy Evaluation

**Simplification:** For discrete states….

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

**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 to MDPs**

**2. Value-Functions**

- Policy Evaluation for a fixed policy

**3. <span style="color:red">Computing an Optimal Policy</span>**

- Policy Improvement

- Value iteration

**4. Infinite vs Finite Horizon**

# How do we find an optimal policy?

**Typically done iteratively:**



Policy Evaluation:
Estimate the Value Function $V^\pi$

Policy Improvement:
Update the Policy $\pi$

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

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

17

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

**Repeat**

    **Repeat** $k = k + 1$

<span style="color:red">Compute Q-Function (for each state action pair)</span>

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

<span style="color:red">Compute V-Function (for each state)</span>

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

    **until convergence of V**

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

**until convergence of policy**

18

# 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

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

This equation is called the Bellman Equation

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

# Value Iteration

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

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

**Small side note:**

Computing optimal V-Function from optimal Q-Function

$$V^*(\boldsymbol{s}) = \max_{\boldsymbol{a}} Q^*(\boldsymbol{s}, \boldsymbol{a})$$

Computing optimal Q-Function from optimal V-Function

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

20

# Algorithmic Description of Value Iteration

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

**Repeat** $k = k + 1$

<span style="color:red">Compute Q-Function (for each state action pair)</span>

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

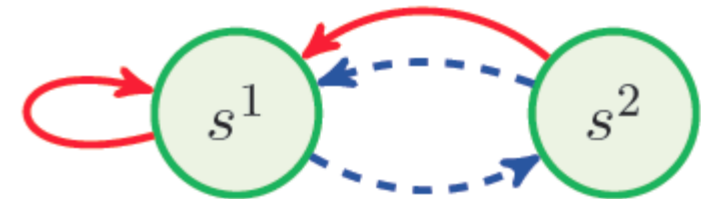<span style="color:red">Compute V-Function (for each state)</span>

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

**until convergence of V**

# Example: Value Iteration

- The Two state example.

➡ HOMEWORK!

# Wrap-Up: Dynamic Programming

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

**Policy Iteration:**

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

**Policy Evaluation:**

**Policy Improvement:**

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

**Value Iteration:**

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

**Get optimal policy after convergence:**

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

23

# Outline of the Lecture

**1. Introduction to MDPs**

**2. Value-Functions**

- Policy Evaluation for a fixed policy

**3. <span style="color:red">Computing an Optimal Policy</span>**

- Policy Improvement

- Value iteration

**4. Infinite vs Finite Horizon**

24

# Finite Horizon Objectives

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

**Finite Horizon T:** Accumulated expected reward for T steps

$$J_{\boldsymbol{\pi}} = \mathbb{E}_{\mu_0, \mathcal{P}, \pi}\left[\sum_{t=1}^{T-1} r_t(\boldsymbol{s}_t, \boldsymbol{a}_t) + r_T(\boldsymbol{s}_T)\right]$$

$r_T(\boldsymbol{s}_T)$ … final reward

# Bellman again…



"*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.*"

Richard Bellman, Dynamic Programming, 1957

# Illustration of basic idea…

**You have won a Best-Paper Award in Madrid!**

**What is the Optimal Policy to Collect it?**

# So what changes to the infinite horizon case?

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

→ It matters, how many time steps are left

→ We can only visit each state (including time index) once!

→ We get a layered / multi stage decision problem

→ optimal policy becomes time-dependent

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

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

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

# Value Iteration for Finite Horizon

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

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

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

➡ Iterate backwards in time

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

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

$$V_T^*(\boldsymbol{s}_T) \implies V_{T-1}^*(\boldsymbol{s}_{T-1}) \implies \ldots \implies V_1^*(\boldsymbol{s}_1)$$

30

# 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

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

# Wrap-Up: Dynamic Programming

**We now know how to compute <span style="color:red">optimal policies for both objectives</span> (finite and infinite horizon)**

Cool, thats all we need. Lets go home…

**Wait, <span style="color:red">there is a catch!</span>**

**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:**

$$\max_{\boldsymbol{a}} Q^*(\boldsymbol{s}, \boldsymbol{a})$$ : difficult in continuous action spaces

$$\mathbb{E}_{\mathcal{P}}\left[V^*(\boldsymbol{s}')\middle|\boldsymbol{s}, \boldsymbol{a}\right]$$ : 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

experience data
$$D = \{\boldsymbol{s}_i, \boldsymbol{a}_i, r_i, \boldsymbol{s}_i'\}_{i=1\ldots N}$$

demonstration data
$$D = \{\boldsymbol{s}_{1:T,i}, \boldsymbol{a}_{1:T,i}\}_{i=1\ldots N}$$

**Model–based RL**

learn model
$\mathcal{P}(\boldsymbol{s}'|\boldsymbol{s}, \boldsymbol{a})$
$r(\boldsymbol{s}, \boldsymbol{a})$

dynamic prog.
$V(\boldsymbol{s})$

policy
$\pi(\boldsymbol{a}|\boldsymbol{s})$

**Model–free RL**

learn value fct.
$\bar{V}(\boldsymbol{s})$

policy
$\pi(\boldsymbol{a}|\boldsymbol{s})$

**Policy Search**

optimize policy
$\pi(\boldsymbol{a}|\boldsymbol{s})$

**Imitation Learning**

learn policy
$\pi(\boldsymbol{a}|\boldsymbol{s})$

**Inverse RL**

learn latent costs
$r(\boldsymbol{s}, \boldsymbol{a})$

dynamic prog.
$\bar{V}(\boldsymbol{s})$

policy
$\pi(\boldsymbol{a}|\boldsymbol{s})$

1. Next Lecture    2.    3.                4.

# Optimal Decision Making: 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?