

Reinforcement Learning 1

Kintan Saha

January 6, 2026

Outline

Overview of ML Theory Series

Reinforcement Learning

Next Steps

Why Machine Learning Theory?

- ▶ Machine Learning Theory seeks to answer a fundamental question: *When, why, and how do learning algorithms work?*
- ▶ It provides **provable guarantees** on performance, generalization, and stability of learning algorithms.
- ▶ When algorithms fail in practice, theory helps explain *why failure is unavoidable* under certain assumptions.
- ▶ Through this seminar series, we aim to showcase key ideas, techniques, and results across modern ML theory.

Why Not Just Take Courses?

- ▶ A full, rigorous treatment of ML theory typically requires multiple advanced courses.
- ▶ Completing these courses often takes 1–2 years of sustained effort.
- ▶ This seminar series is designed as a **gentle, structured introduction** to several subfields of ML theory.
- ▶ The goal is to provide intuition, key results, and clear pointers for deeper independent study.

Format of the Seminar Series

- ▶ The series will be organized as a sequence of **mini-courses** on different topics in ML theory.
- ▶ Example themes include:
 - ▶ Reinforcement Learning and Stochastic Approximation
 - ▶ Over-parameterized Models
 - ▶ Diffusion Models
- ▶ Each mini-course will consist of 3–4 lectures.
- ▶ These will be followed by faculty or industry talks showcasing state-of-the-art research.
- ▶ Where possible, we will include hands-on components such as paper reproductions or small projects.

Tentative Schedule

- ▶ **Jan 6, 8, 10, 12, 14:** Reinforcement Learning + Stochastic Approximation
- ▶ **Jan 15–18:** Reinforcement Learning Workshop (CSA)
- ▶ **Jan 24/25:** Talks by Prof. Shubhada Agrawal / Prof. Gugan Thoppe
- ▶ **Jan 31, Feb 7, 14, 21:** Over-parameterized Models
- ▶ **Feb 28:** Talk by Prof. Anant Raj
- ▶ **March:** Diffusion Models + Industry Talks (Microsoft Research)

Note

The continuation of the series depends on audience interest and participation. If you want to interact with leading researchers, make sure to attend!

RL Workshop

Timings are as per the Indian Standard Time

Day	8:45am to 9:00am	9:00am to 10:00am	10:00am to 10:30am	10:30am to 11.30am	11:30 am to 12:30 pm	12:30 pm to 2:00 pm	2:00 pm to 3:00 pm	3:00 pm to 4:00 pm	4:00 pm to 4:30 pm	4:30 pm to 5:30 pm
Thursday 15th January 2026		Tutorial: RL Theory Gugan Thoppe (IISc)	Coffee Break	Tutorial: RL Theory Gugan Thoppe (IISc)	Tutorial: RL Theory Gugan Thoppe (IISc)	Lunch	Tutorial: RL Theory Gugan Thoppe (IISc)	Tutorial: RL Simulations Gugan Thoppe (IISc)	Coffee Break	Tutorial: RL Simulations Gugan Thoppe (IISc)
Friday 16th January 2026	Inauguration by CSA Chair	Furong Huang <i>(UMD)</i> <input type="button" value="(Online)"/>	Coffee Break	Vivek Borkar <i>(IIT Bombay)</i>	M Vidyasagar <i>(IIT Hyderabad)</i>	Lunch	Dheeraj Nagraj <i>(Google)</i>	Raunak B <i>(IIT Delhi)</i>	Coffee Break	Arunselvan Ramaswamy <i>(Karistad University)</i>
Saturday 17th January 2026		Serdar Yuksel <i>(Queens)</i> <input type="button" value="(Online)"/>	Coffee Break	Sridharan Devarajan <i>(CMU)</i>	Prashanth L A <i>(IIT Madras)</i>	Lunch	Paras Chopra <i>(LossFunk)</i>	Karthik Sundaresan <i>(Walmart)</i>	Coffee Break	Pranay Sharma <i>(IIT Bombay)</i>
Sunday 18th January 2026		Aviral Kumar <i>(IISc)</i> <input type="button" value="(Online)"/>	Coffee Break	Shubhada Agarwal <i>(IISc)</i>	Swanand Khare <i>(IIT Madras)</i>	Lunch	Networking and Casual Discussions	Networking and Casual Discussions	Coffee Break	Networking and Casual Discussions

Outline

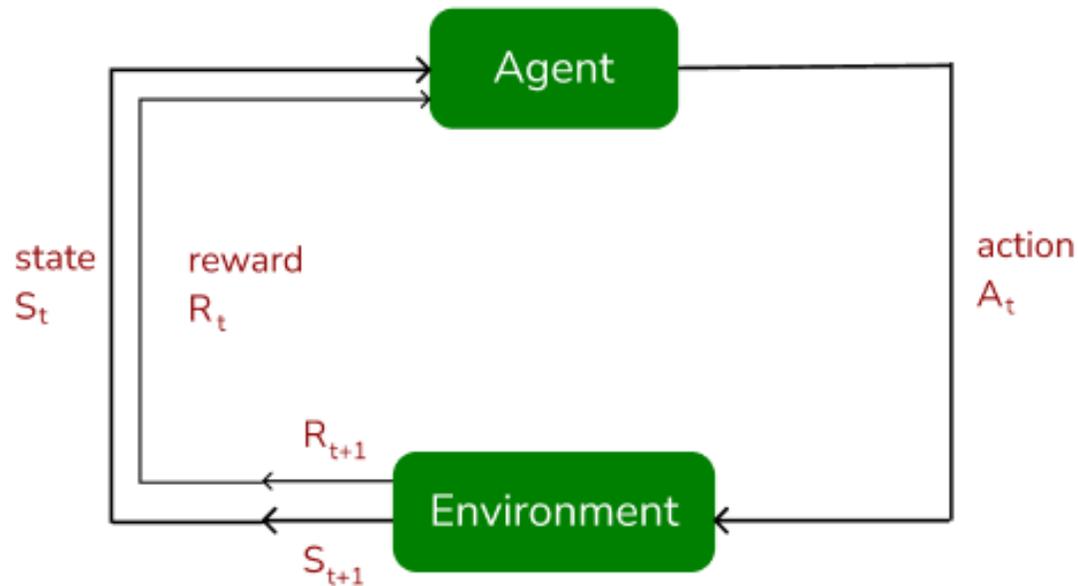
Overview of ML Theory Series

Reinforcement Learning

Next Steps

Setup: Agent–Environment Interaction

Reinforcement Learning (RL) studies *sequential decision making under uncertainty*.



Key idea

An **agent** repeatedly interacts with an **environment**, observing states, choosing actions, and receiving rewards.

What is Sequential and What is Uncertain?

Sequential decision making

At time t :

$$S_t \xrightarrow{A_t} S_{t+1}$$

$$\mathbb{E} \left[\sum_{t=0}^{\infty} R_t \right]$$

Actions affect not only immediate reward but also *future states*.

Uncertainty

- ▶ The reward function $R(s, a)$ is unknown
- ▶ The transition kernel $P(\cdot | s, a)$ is unknown

$$\lim_{T \rightarrow \infty} \mathbb{E} \left[\frac{\sum_{t=0}^T R_t}{T} \right] \rightarrow \text{average reward criterion}$$

Objective

Maximize expected discounted return:

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t \right], \quad \gamma \in (0, 1)$$

Markov Decision Processes (MDPs)

\mathcal{S}

Definition

A Markov Decision Process (MDP) is a tuple

$$(\mathcal{S}, \mathcal{A}, R, P, \gamma)$$

where:

s, a

$$s' \sim P(\cdot | s, a)$$

- ▶ \mathcal{S} : state space
- ▶ \mathcal{A} : action space
- ▶ $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_+$: reward function
- ▶ $P(\cdot | s, a)$: transition kernel (probability simplex over \mathcal{S})
- ▶ $\gamma \in (0, 1)$: discount factor

$(s_0, \underbrace{A_0, s_1, A_1, s_2, \dots}_{\text{seq. of random variables}})$

Markov property

$$\mathbb{P}(S_{t+1} | S_0, A_0, \dots, S_t, A_t) = P(S_{t+1} | S_t, A_t)$$

Policies and Value Functions

Policy

A (stationary) policy π is a mapping

$$\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$$

↳ set of prob. distributions on \mathcal{A}
 $s \quad a \sim \pi(\cdot | s)$

$$s \quad a_0 \quad \hookrightarrow P(a_0) = 1$$

$$P(a) = 0 \quad \forall a \neq a_0$$

where $\pi(a | s)$ is the probability of selecting action a in state s .

State-value function

For a policy π ,

$$V^\pi(s) := \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s \right]$$

$$= \underbrace{\sum_a \pi(a|s) r(s,a)}_{\mathbb{E}[R(S_0, A_0) | S_0 = s]} + \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) \sum_{a'} \pi(a'|s') r(s', a')$$

Action-value function

$$Q^\pi(s, a) := \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s, A_0 = a \right]$$

$$V^\pi(s) = \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_0 = s \right]$$

$$= \underbrace{\mathbb{E}^\pi [R(S_0, A_0) \mid S_0 = s]} + \gamma \underbrace{\mathbb{E}^\pi [R(S_1, A_1) \mid S_0 = s]}$$

$$\sum_{a \in A} \pi(a|s) r(s, a)$$

$\downarrow S_1, A_1$
 $S_1 \sim P(\cdot \mid S_0, A_0)$
 $A_0 \sim \pi(\cdot \mid S_0)$

$$\sum_a \pi(a|s) \sum_{s'} P(s' \mid s, a) \sum_{a'} \pi(a'|s') r(s', a')$$

$$= \mathbb{E} [R(S_1, A_1)]$$

$$= \sum_{s' \in S} R(s', A_1) \cdot \underbrace{P(s')}_{\parallel}$$

$$\left[\sum_{a \in A} P(s' \mid s, a) \pi(a|s) \right]$$

Optimality and Bellman Optimality Equation

Optimal value function

$$V^*(s) := \sup_{\pi} V^{\pi}(s)$$

Bellman optimality equation

$$V^*(s) = \max_{a \in \mathcal{A}} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

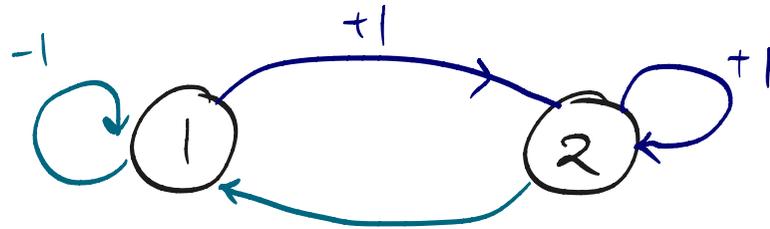


Bellman optimality operator

$$(TV)(s) := \max_a \left[R(s, a) + \gamma \mathbb{E} V(s') \right]$$

MDP Example $\gamma = 0.9$

2 actions: right left



Policy π_1 : "always right" $\rightarrow V_{\pi_1}(1) = \sum_{t=0}^{\infty} \gamma^t \cdot 1 = \frac{1}{1-\gamma} = 10$ $V_{\pi_1}(2) = \frac{1}{1-\gamma} = 10$

π_2 : "always left" $\rightarrow V_{\pi_2}(1) = \sum_{t=0}^{\infty} \gamma^t \cdot (-1) = -10$ $V_{\pi_2}(2) = -10$

Optimal policy: $V^*(2) = \max \left\{ \underset{\text{right} \downarrow}{+1 + \gamma V^*(2)}, \underset{\text{left} \downarrow}{-1 + \gamma V^*(1)} \right\}$
 $V^*(1) = \max \left\{ \underset{\text{right} \downarrow}{+1 + \gamma V^*(2)}, \underset{\text{left} \downarrow}{-1 + \gamma V^*(1)} \right\}$

if right is optimal in both states then $V^*(2) = V^*(1) = 10$

CHECK: at state 1, $1 + \gamma V^*(2) = 10 > -1 + \gamma V^*(1) = 8$

at state 2, $1 + \gamma V^*(2) = 10 > -1 + \gamma V^*(1) = 8$

\Rightarrow right is optimal in both states

Derivation of the Bellman Optimality Equation (I)

Definition of the optimal value function

For any state s ,

$$V^*(s) := \sup_{\pi} \mathbb{E}^{\pi} \left[\sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s \right]$$

Unroll the return (first step)

Separate the immediate reward:

$$\begin{aligned} V^*(s) &= \sup_{\pi} \mathbb{E}^{\pi} \left[R(S_0, A_0) + \sum_{t=1}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s \right] \\ &= \sup_{\pi} \mathbb{E}^{\pi} \left[R(s, A_0) + \gamma \sum_{t=0}^{\infty} \gamma^t R(S_{t+1}, A_{t+1}) \mid S_0 = s \right] \end{aligned}$$

Derivation of the Bellman Optimality Equation (II)

Apply the tower property

Condition on the next state S_1 :

$$\mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_{t+1}, A_{t+1}) \mid S_0 = s \right] = \mathbb{E}^\pi \left[\mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_{t+1}, A_{t+1}) \mid S_1 \right] \mid S_0 = s \right]$$

Optimal substructure (principle of optimality)

For any s' ,

$$\sup_{\pi} \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_{t+1}, A_{t+1}) \mid S_1 = s' \right] = \sup_{\pi} \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s' \right] = V^*(s')$$

Derivation of the Bellman Optimality Equation (III)

Substitute the tail by V^*

Combining the previous steps gives

$$V^*(s) = \sup_{\pi} \mathbb{E}^{\pi} [R(s, A_0) + \gamma V^*(S_1) \mid S_0 = s]$$

Reduce \sup_{π} to \max_a (first decision)

At time 0 the policy selects a distribution over actions; the supremum is attained by a deterministic choice:

$$V^*(s) = \max_{a \in \mathcal{A}} [R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s')]$$

Bellman optimality equation and operator form

$$V^*(s) = (TV^*)(s) := \max_a [R(s, a) + \gamma \mathbb{E} V^*(s')]$$

Bellman Evaluation Operator

Bellman expectation equation

For a fixed policy π :

$$V^\pi(s) = \sum_a \pi(a|s) \left[R(s, a) + \gamma \mathbb{E} V^\pi(s') \right]$$

Bellman evaluation operator

$$(T^\pi V)(s) := \sum_a \pi(a|s) \left[R(s, a) + \gamma \mathbb{E} V(s') \right]$$

Fixed point

$$V^\pi = T^\pi V^\pi$$

Derivation of the Bellman Expectation Equation (I)

Definition of the value function

For a fixed (stationary) policy π and state s ,

$$V^\pi(s) := \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s \right]$$

Unroll the return

Separate the immediate reward:

$$\begin{aligned} V^\pi(s) &= \mathbb{E}^\pi \left[R(S_0, A_0) + \sum_{t=1}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s \right] \\ &= \mathbb{E}^\pi \left[R(s, A_0) + \gamma \sum_{t=0}^{\infty} \gamma^t R(S_{t+1}, A_{t+1}) \mid S_0 = s \right] \end{aligned}$$

Derivation of the Bellman Expectation Equation (II)

$$\mathbb{E} X = \mathbb{E} \left[\mathbb{E} [X | Y] \right]$$

Apply the tower property

Condition on the next state S_1 :

$$\mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_{t+1}, A_{t+1}) \mid S_0 = s \right] = \mathbb{E}^\pi \left[\mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_{t+1}, A_{t+1}) \mid S_1 \right] \mid S_0 = s \right]$$

Identify the tail return

By the Markov property and stationarity of π , for any s' ,

$$\mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_{t+1}, A_{t+1}) \mid S_1 = s' \right] = V^\pi(s')$$

Derivation of the Bellman Expectation Equation (III)

Substitute back

$$V^\pi(s) = \mathbb{E}^\pi[R(s, A_0) + \gamma V^\pi(S_1) \mid S_0 = s]$$

$$\underbrace{f(x) = x}$$

Expand the expectation

$$V^* = T V^*$$

Using $A_0 \sim \pi(\cdot \mid s)$ and $S_1 \sim P(\cdot \mid s, A_0)$,

$$V^\pi(s) = \sum_a \pi(a \mid s) \left[R(s, a) + \gamma \sum_{s'} P(s' \mid s, a) V^\pi(s') \right]$$

Bellman evaluation fixed-point equation

$$\boxed{V^\pi = T^\pi V^\pi}$$

$$(T^\pi V)(s) := \sum_a \pi(a \mid s) \left[R(s, a) + \gamma \mathbb{E}_{s'} V(s') \right]$$

What Do We Want to Solve?

Core RL questions

1. **Policy evaluation:** Given π , compute V^π
2. **Policy improvement:** Improve π using V^π

Goal

Construct a sequence of policies converging to an optimal policy π^* .

Solving Bellman Equations Directly

Observation

Bellman equations are systems of linear equations:

$$V^\pi = R^\pi + \gamma P^\pi V^\pi$$

Why not solve directly?

- ▶ State space can be huge
- ▶ Matrix inversion is expensive

A Detour: Fixed Points

Operator form

$$V^\pi = T^\pi V^\pi, \quad V^* = TV^*$$

Key question

Can we compute fixed points using repeated application of operators?

Banach Fixed Point Theorem

Contraction mapping

An operator T is a contraction on $(\mathcal{X}, \|\cdot\|)$ if

$$\|Tx - Ty\| \leq c\|x - y\|, \quad c < 1$$

Banach Fixed Point Theorem

If T is a contraction on a complete normed space, then:

- ▶ T has a unique fixed point
- ▶ Iteration $x_{k+1} = Tx_k$ converges to it

$$\begin{aligned} f(a) &= x \\ f & \\ a, f(a), f^2(a) \\ \dots f^n(a) \\ \lim_{n \rightarrow \infty} f^n(a) &= x \end{aligned}$$

Choosing a Norm

$\|\cdot\|$

Common norms

$$\textcircled{1} \|x\| + \|y\| \geq \|x+y\| \quad \textcircled{2} \|\alpha x\| = |\alpha| \|x\| \quad \forall \alpha \in \mathbb{R}$$

$$S \quad \|\cdot\| : S \rightarrow \mathbb{R}^+$$

$$\textcircled{3} \|x\| = 0$$

$$\|V\|_\infty = \max_s |V(s)|, \quad \|V\|_2, \quad \|V\|_1 \quad \Leftrightarrow x=0$$

Equivalence of norms

In finite-dimensional spaces, all norms are equivalent:

$$c_1 \|x\|_a \leq \|x\|_b \leq c_2 \|x\|_a$$

$$x \in \mathbb{R}^d$$

$$\|x\|_2 = \sqrt{\sum x_i^2}$$

$$\|x\|_1 = \sum |x_i|$$

$$\|x\|_\infty = \max_{i \in [d]} |x_i|$$

Why l_∞ ?

It interacts cleanly with max and expectations.

Contraction of Bellman Operators

Theorem

Both T^π and T are γ -contractions in $\|\cdot\|_\infty$:

$$\|TV - TW\|_\infty \leq \gamma \|V - W\|_\infty$$

$$\|T^\pi V - T^\pi W\|_\infty \leq \gamma \|V - W\|_\infty$$

Contraction of the Bellman Optimality Operator T (I)

Recall

$$(TV)(s) := \max_{a \in \mathcal{A}} \left[R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$$

Goal

Show that T is a γ -contraction in the ℓ_∞ norm:

$$\|TV - TW\|_\infty \leq \gamma \|V - W\|_\infty$$

Key inequality

For any real vectors (x_a) and (y_a) ,

$$\left| \max_a x_a - \max_a y_a \right| \leq \max_a |x_a - y_a|$$

Contraction of the Bellman Optimality Operator T (II)

Proof

Fix a state s . Define, for each action a ,

$$x_a := R(s, a) + \gamma \mathbb{E}_{s'} V(s'), \quad y_a := R(s, a) + \gamma \mathbb{E}_{s'} W(s')$$

Then

$$\begin{aligned} |(TV)(s) - (TW)(s)| &= \left| \max_a x_a - \max_a y_a \right| \\ &\leq \max_a |x_a - y_a| \\ &= \max_a \gamma \left| \mathbb{E}[V(s') - W(s')] \right| \\ &\leq \gamma \|V - W\|_\infty \end{aligned}$$

Taking the supremum over all states s ,

$$\|TV - TW\|_\infty \leq \gamma \|V - W\|_\infty$$

Value Iteration Algorithm

Algorithm

Initialize V_0 arbitrarily. For $k \geq 0$:

$$V_{k+1} = TV_k$$

V^*
 ~~V~~
 TV
 $T^2V \dots$ $\lim_{n \rightarrow \infty} T^n V = V^*$

Policy extraction

$$\pi_k(s) \in \arg \max_a \left[R(s, a) + \gamma \mathbb{E} V_k(s') \right]$$

When should the algorithm stop btw?

$TV^* = V^*$

Value Iteration: Guarantees

Main guarantees

- ▶ Monotonic value improvement
- ▶ Convergence to V^*
- ▶ Geometric rate γ^k

$$V_0 \leq TV_0 \leq T^2V_0 \dots$$
$$\neq V_0 \leq V_1 \leq V_2 \dots$$

Monotonic Value Improvement

Assumption

Rewards are nonnegative and $V_0 \equiv 0$.

Claim

$$V_0 \leq V_1 \leq V_2 \leq \dots \leq V^*$$

Reason

Bellman operator is monotone:

$$V \leq W \Rightarrow TV \leq TW$$

Monotonicity of the Bellman Operator (I)

Claim (Monotonicity)

If $V(s) \leq W(s)$ for all $s \in \mathcal{S}$, then

$$(TV)(s) \leq (TW)(s) \quad \text{for all } s.$$

Proof

Fix a state s and any action a . Since $V \leq W$ pointwise,

$$\mathbb{E}_{s'}[V(s')] \leq \mathbb{E}_{s'}[W(s')].$$

Therefore,

$$R(s, a) + \gamma \mathbb{E}V(s') \leq R(s, a) + \gamma \mathbb{E}W(s').$$

Monotonicity of the Bellman Operator (II)

Preservation under maximization

Since the above inequality holds for every action a ,

$$\max_a \left[R(s, a) + \gamma \mathbb{E} V(s') \right] \leq \max_a \left[R(s, a) + \gamma \mathbb{E} W(s') \right].$$

That is,

$$(TV)(s) \leq (TW)(s).$$

Conclusion

$$\boxed{V \leq W \Rightarrow TV \leq TW}$$

Hence, the Bellman optimality operator is *order-preserving* (monotone).

Convergence to the Optimal Value

$$\|V_{k+1} - V^*\|_\infty$$

$$= \|TV_k - TV^*\|_\infty$$

$$\leq \gamma \|V_k - V^*\|_\infty$$

Contraction argument

$$\|V_k - V^*\|_\infty \leq \gamma^k \|V_0 - V^*\|_\infty$$

Conclusion

Value Iteration converges geometrically fast to V^* .

Outline

Overview of ML Theory Series

Reinforcement Learning

Next Steps

Remarks and Outlook

$$T^\pi V = \sum_{\pi(a|s)} \left[\underbrace{r(s,a)}_{???} + \gamma \sum_{\pi(s'|s,a)} \underbrace{P(s'|s,a)}_{???} V(s') \right]$$

- ▶ Recall the MDP setup and the agent–environment interaction framework.
- ▶ In practice, the reward function R and transition kernel P are *unknown*.
- ▶ Consequently, the Bellman operators T^π and T are not available explicitly.
- ▶ As a result, classical value iteration cannot be implemented directly.
- ▶ This motivates the study of **model-free** algorithms, which learn from data without explicit knowledge of R and P .
- ▶ In the next session, we will study such algorithms and introduce **stochastic approximation** as a key tool for their analysis.

Thank You

Thank you for your attention!

Questions, comments, or clarifications?