

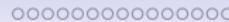
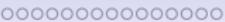
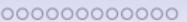
# Introduction to Reinforcement Learning

## Part 1: MDP and Bellman Equations

Eric J. Pauwels

CWI & VU

Version, November 25, 2021





## What is Reinforcement Learning (RL)?

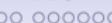
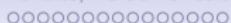
Markov Decision Processes (MDP)

Policies, Value Functions and the Bellman Equation

Bellman equations

Bellman equations for optimality

Summary and Outlook



# Outline

What is Reinforcement Learning (RL)?

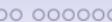
Markov Decision Processes (MDP)

Policies, Value Functions and the Bellman Equation

Bellman equations

Bellman equations for optimality

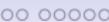
Summary and Outlook



# What is Reinforcement Learning (RL)

- RL is one of **three main learning problems** studied in AI:
- Learning by **trial and error** and **improving** over time;
- **Natural form of learning**, ubiquitous in biology and psychology:





# What is Reinforcement Learning?

- **Supervised learning** needs **teacher** (labeled data!)
  - **RL learns on the job:** interact, explore, exploit experience!





## RL Example 1: Playing a song by ear



- **Actions and transitions:** Moving from key to key...
- **Immediate reward:** correct key or not;
- **Cumulative reward:** Being able to **play whole song**
- Relatively easy, one gets **immediate valuable feedback**...

## RL Example 2: Cooking

- **States:**
- **Actions:** e.g. cooking, seasoning, stirring, etc.
- **Transitions:** e.g. from raw to cooked
- **Immediate reward**  
e.g. smell, look
- **Final reward:** enjoyable dinner
- **Problem:** what's the best sequence of actions (i.e. **optimal policy**)?



## RL Examples: Playing Chess

- **States:** board configurations
- **Actions:** legal moves
- **Value of state/action:** very valuable to determine next move (action),
- **Final reward:**  
Win(2), lose(0), draw(1)
- **Problem:** what's the optimal action to take in each state?  
**(policy)**



# What makes Reinforcement Learning challenging?

## 1. No Oracle

- **Supervised learning:** compare estimate to correct result (provided by oracle);
- **RL: no information on optimal action**, just indicative *reward*
  - i.e. **evaluative feedback, not prescriptive!**
- **RL: only data on states/actions that have been experienced:** (might be very **small subset** of entire state space);

## 2. Sparsity of feedback

- For many RL problems, **no feedback until target** is reached!
- Initial stages: just random search!

## 3. Data generated during training

- **Online learning:** improve policies while interacting w. environment
- **No independent data** generation!

policy  $\longleftrightarrow$  data generation

## Reinforcement Learning: more abstract view

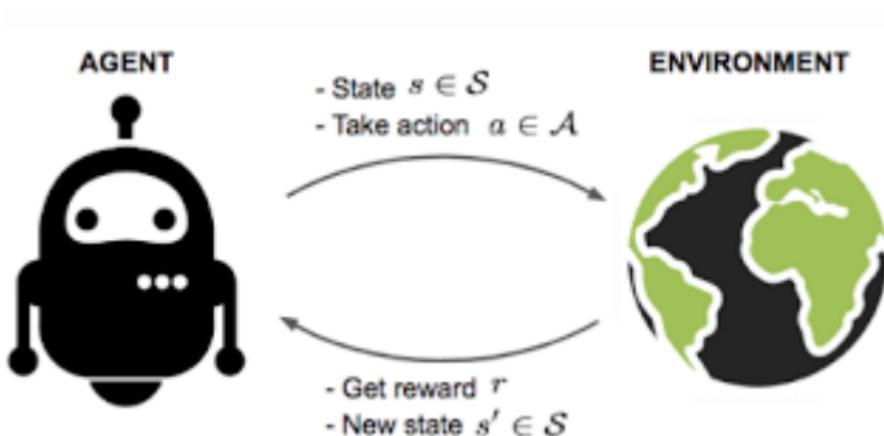
- **RL:** Learning from (sequential) interaction:  
take **actions** to move from **state** to (better!?) **state**;
- No labeled data, but **feedback** (**reward**, possibly delayed)
- **Goal:** optimise **end-result** (i.e. long-term/cumulative reward)



## RL: Even more abstract version

### Agent interacting with Environment:

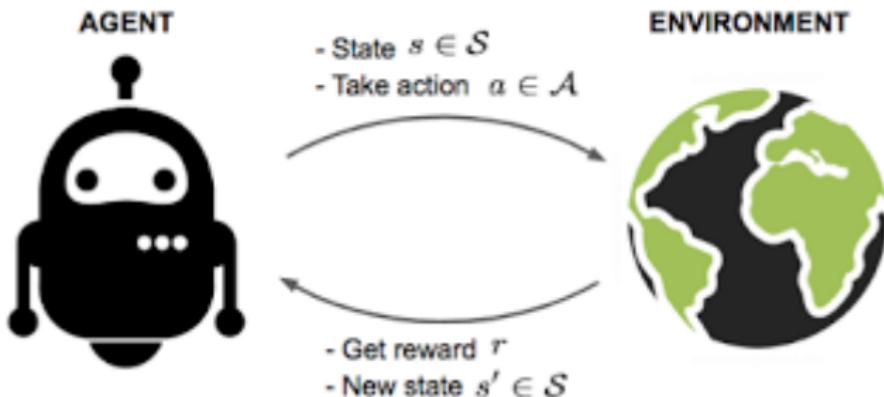
- states( $s$ ), actions ( $a$ ) and transitions:  $s \xrightarrow{a} s'$  ( $p(s' | s, a)$ )
- (immediate) reward  $r(s, a, s')$



## RL GOAL: Find optimal policy!

**Policy  $\pi$** , tells for each state ( $s$ ), which action ( $a$ ) to take:

$$s \xrightarrow{\pi} a$$



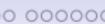
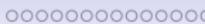
**RL Goal:** Given environment, determine **optimal policy  $\pi^*$** :

**what action to choose in each state to achieve best FINAL reward?**



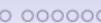
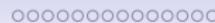
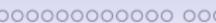
# Overview

MDP
states, actions, transitions, rewards

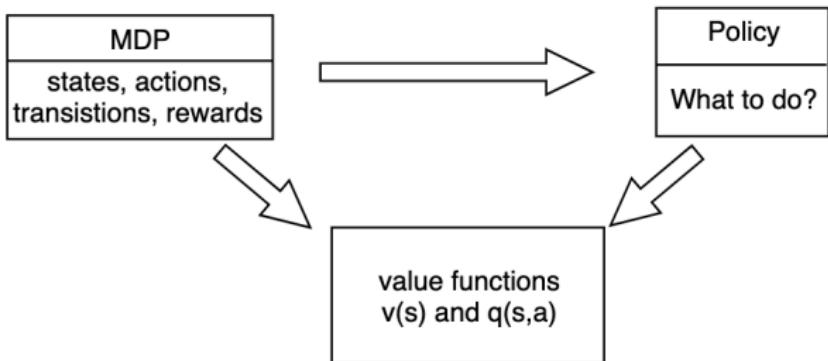


# Overview

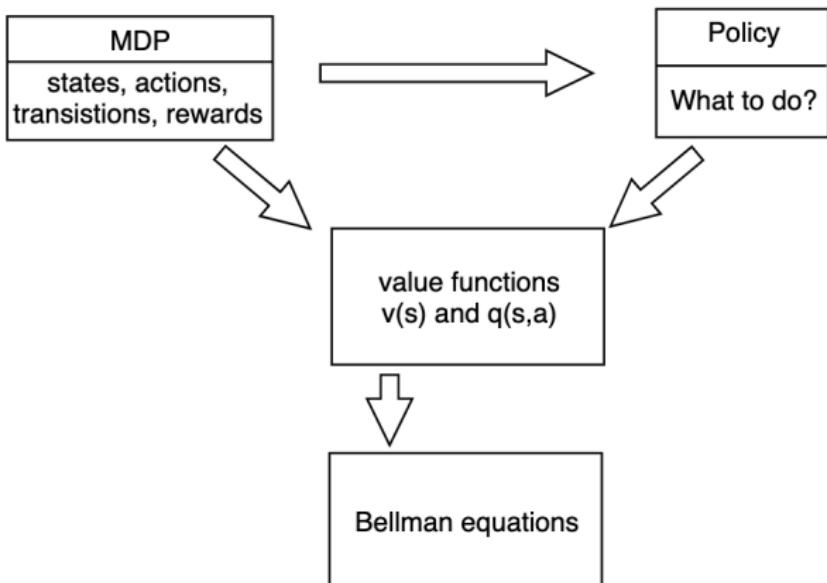


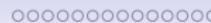


## Overview

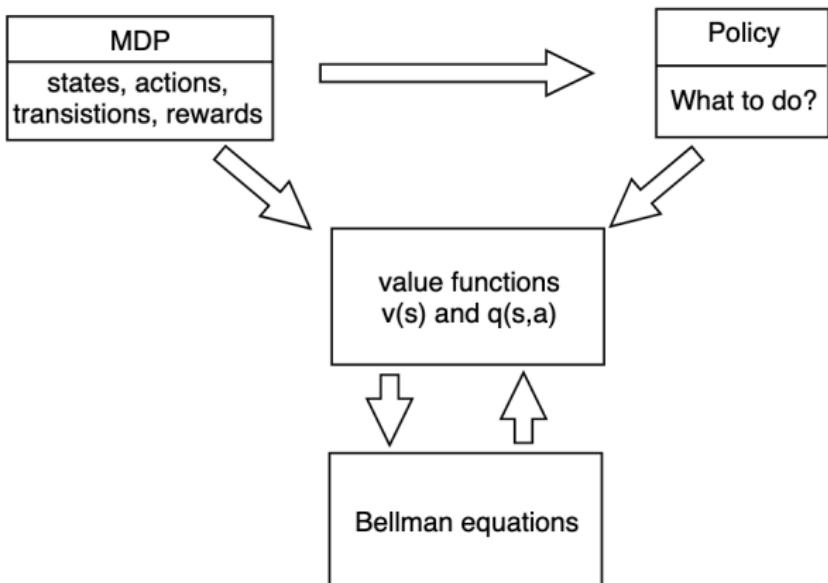


# Overview

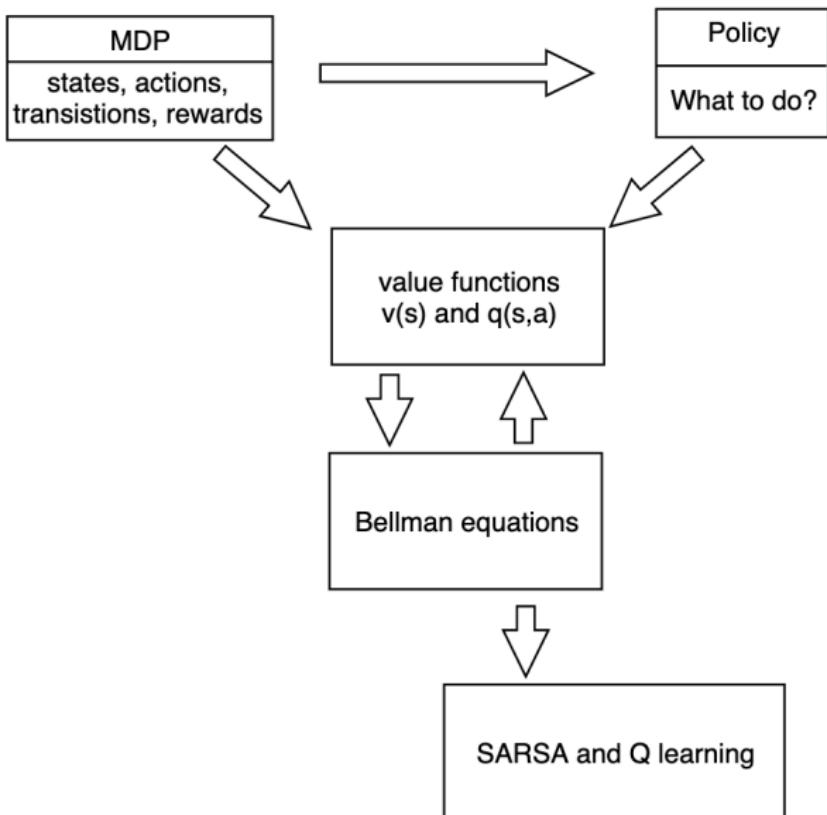




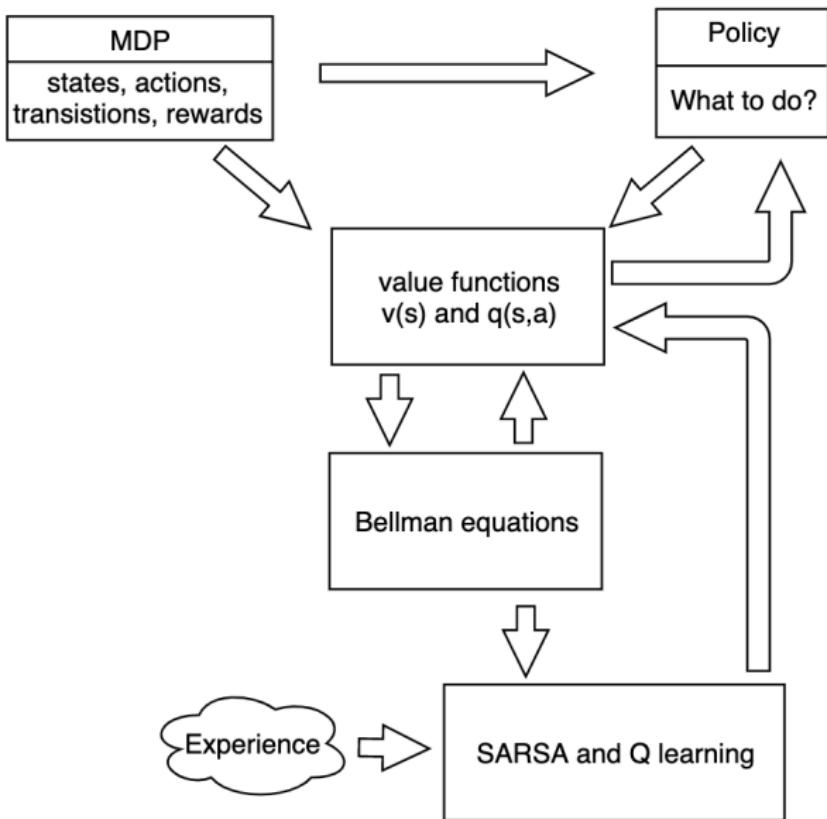
## Overview



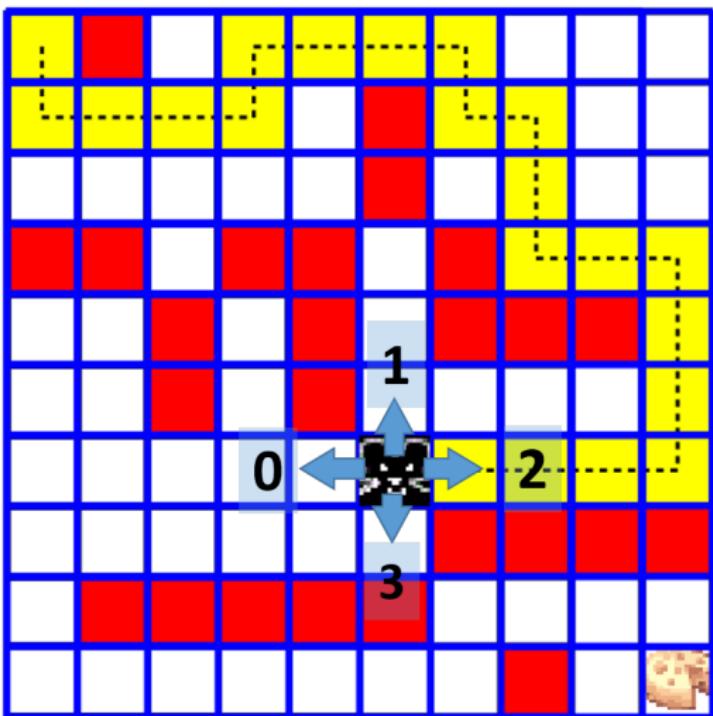
# Overview



# Overview



## Gridworld Maze as Prototypical RL Example



# Outline

What is Reinforcement Learning (RL)?

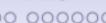
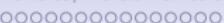
Markov Decision Processes (MDP)

Policies, Value Functions and the Bellman Equation

Bellman equations

Bellman equations for optimality

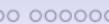
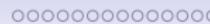
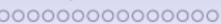
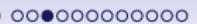
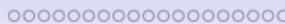
Summary and Outlook



# Markov Decision Process (MDP)

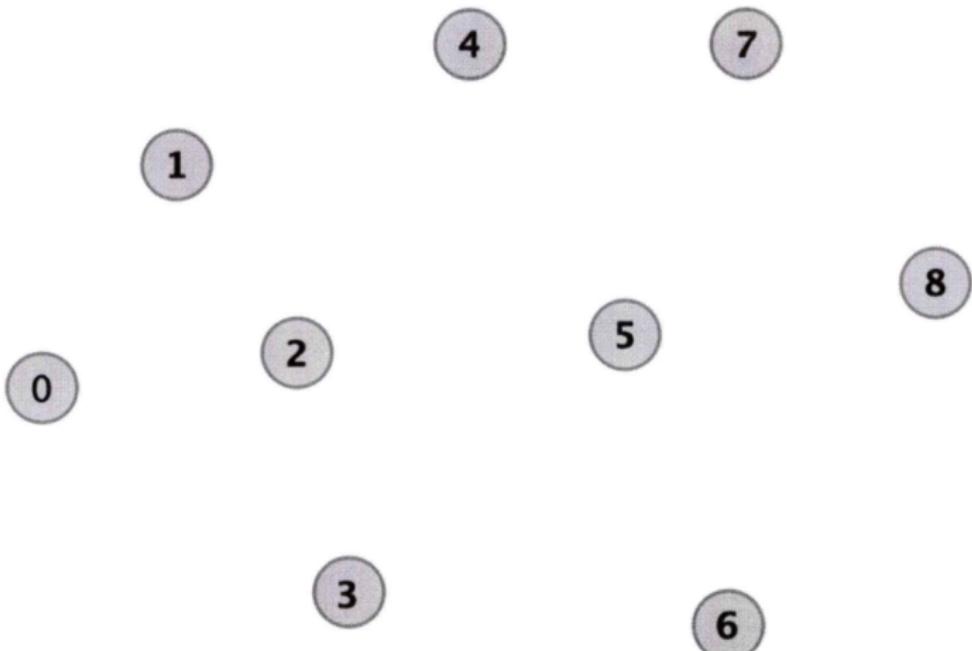
**Mathematical formalism** to capture the **essential characteristics** of the RL problem:

- MDP = states, actions, transitions, rewards, (discount factor)



# Markov Decision Process (MDP) States

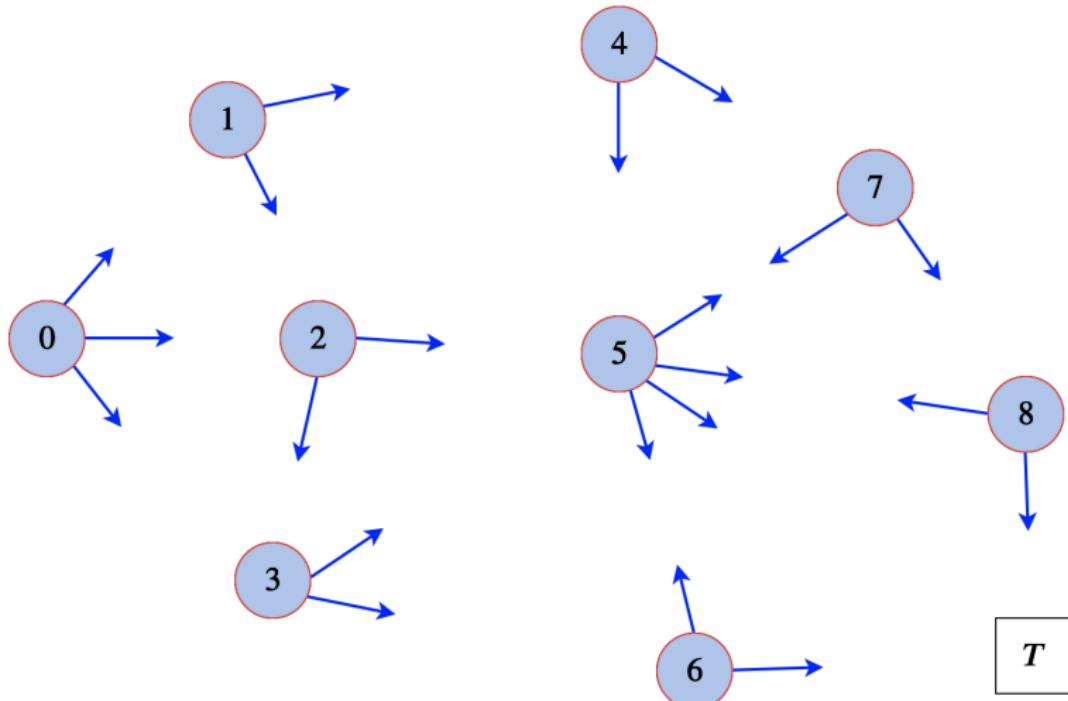
STATES





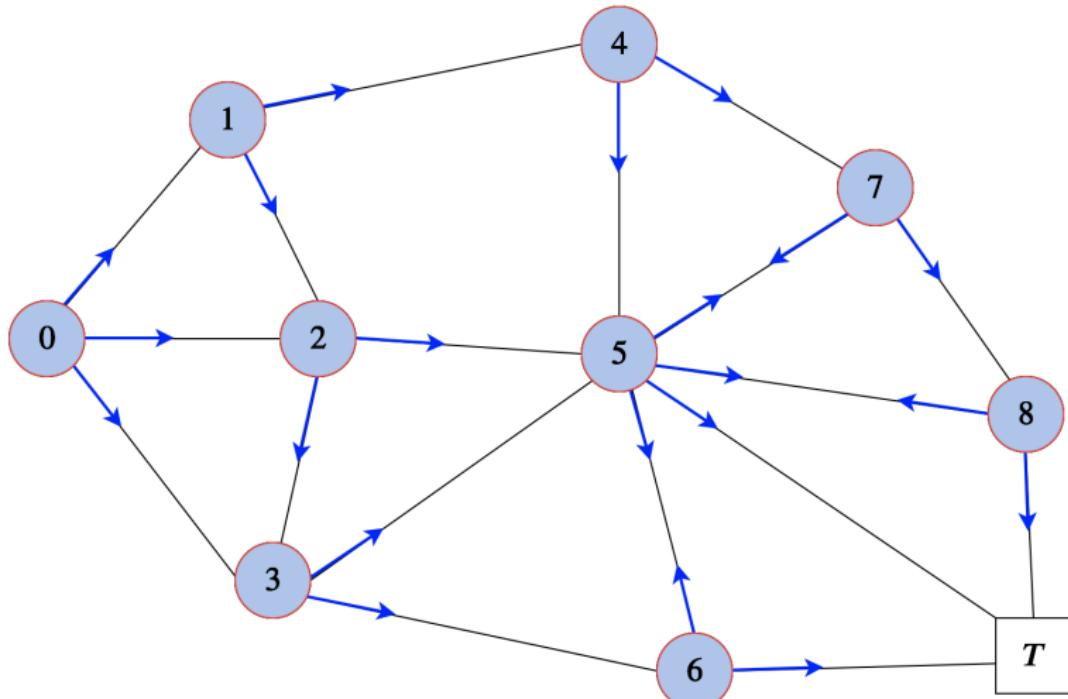
# Markov Decision Process (MDP)

## States + Actions



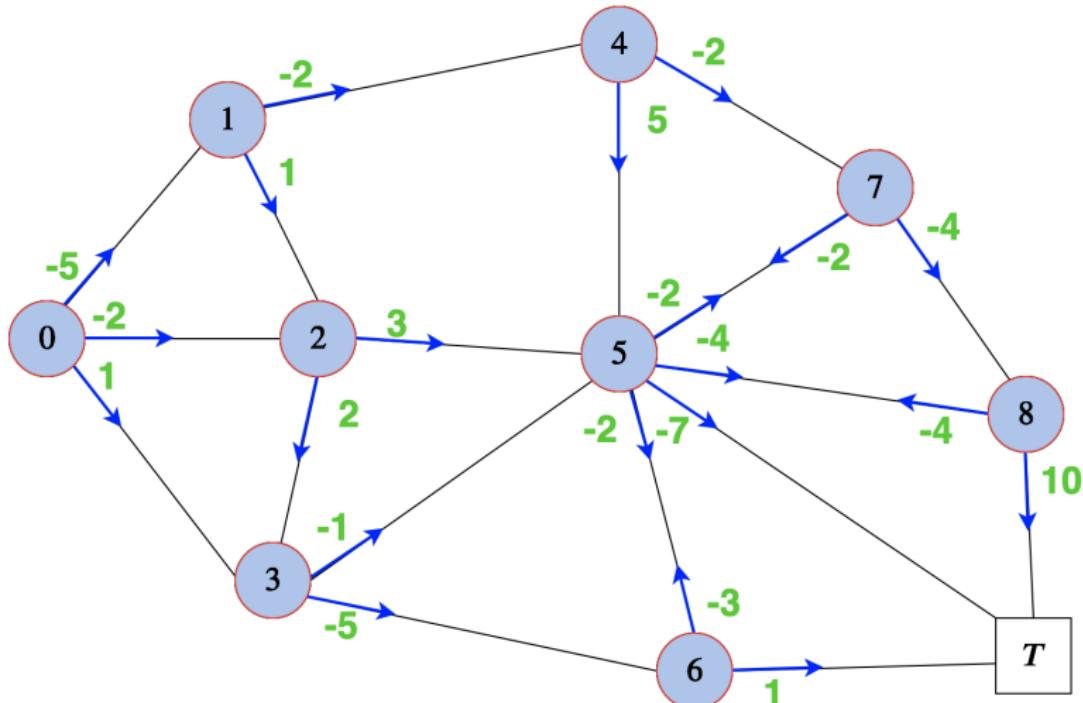
# Markov Decision Process (MDP)

## States + Actions + Transitions



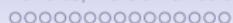
# Markov Decision Process (MDP)

## States + Actions + Transitions + Rewards



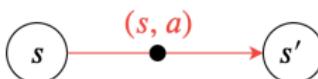
# Markov decision processes (MDP)

- **Markov decision processes (MDP)** provide a **formal model** for a **sequential decision problem**;
- A **finite MDP**  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$  consists of:
  - **Discrete time**  $t = 0, 1, 2, \dots$
  - A discrete set of **states**  $s \in \mathcal{S}$
  - A discrete set of **actions**  $a \in \mathcal{A}(s)$  for each  $s$
  - A **transition function**  $p(s'|s, a)$ : probability of transitioning to state  $s'$  when taking action  $a$  at state  $s$
  - A **reward function**  $r(s, a, s') = E[R | s, a, s']$ : expected reward when taking action  $a$  at state  $s$  and transitioning to  $s'$
  - A planning horizon  $H$  or **discount factor**  $\gamma$ :
    - How important are future rewards? planning horizon: 规划范围
    - shortsighted  $0 \leftarrow \gamma \rightarrow 1$  farsighted  
**短视**



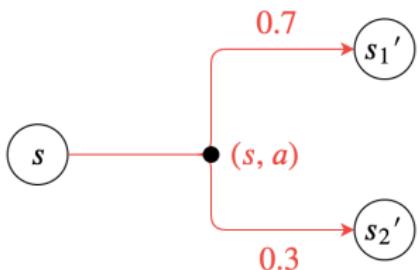
## Detailed Definition and Notation (1)

Deterministic



$$p(s' \mid s, a) = 1$$

Probabilistic



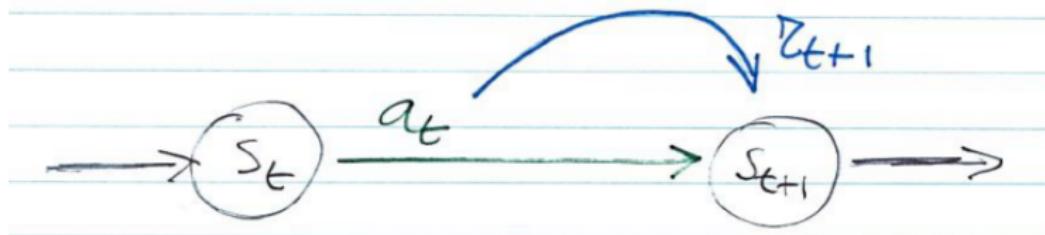
$$p(s_1' \mid s, a) = 0.7$$

$$p(s_2' \mid s, a) = 0.3$$

**Transition probability** to state  $s'$  when taking action  $a$  in state  $s$  :

$$p(s' \mid s, a) := P(S_{t+1} = s' \mid S_t = s, A_t = a)$$

## Detailed Definition and Notation (2)



- **Expected immediate reward** when transitioning from state  $s$  to state  $s'$  under action  $a$ :

$$r(s, a, s') := E(R_{t+1} = r \mid S_t = s, A_t = a, S_{t+1} = s')$$

# The Markov property

## Markov property: informally

- The present (state) has all the information necessary to predict the future: no need to keep track of the past.
- "*The future is independent of the past, given the present*"

## Mathematical formulation

Joint distribution:

$$P(S_{t+1}, R_{t+1} | s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_1, s_0, a_0) = P(S_{t+1}, R_{t+1} | s_t, a_t)$$

If reward  $R_{t+1}$  does NOT depend on successor state  $S_{t+1}$ :

$$P(S_{t+1} | s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_1, s_0, a_0) = P(S_{t+1} | s_t, a_t)$$

$$P(R_{t+1} | s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_1, s_0, a_0) = P(R_{t+1} | s_t, a_t)$$

## Markov Property

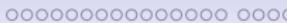
- MP simplifies the mathematics and makes it tractable;
- MP is sufficiently powerful as **states** can be expanded to include all the information necessary to predict future;
  - E.g.: to estimate the speed of a vehicle, we need at least two positions!

## Example of Markov game

- **Game of chess**

- The **current state** of the board provides all information needed to determine **next (optimal) move**;
- **No need** to know the **history** (i.e. the path leading up to the current state)





# Outline

What is Reinforcement Learning (RL)?

Markov Decision Processes (MDP)

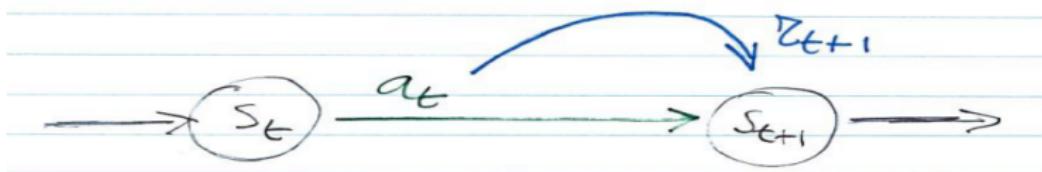
Policies, Value Functions and the Bellman Equation

Bellman equations

Bellman equations for optimality

Summary and Outlook

## Policies for a MDP



- **Policy** specifies what actions should be taken in a given state
- Hence, a **policy  $\pi$  maps states into actions**;
  - **Deterministic policy:**  $a = \pi(s)$
  - **Stochastic policy:**

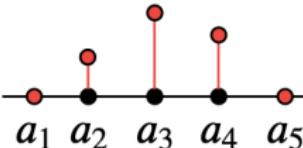
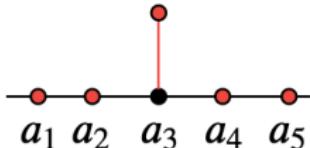
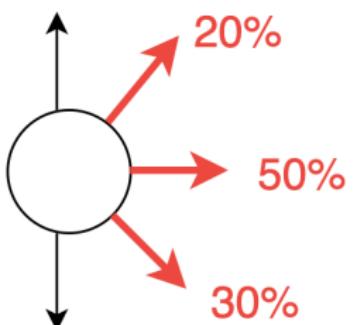
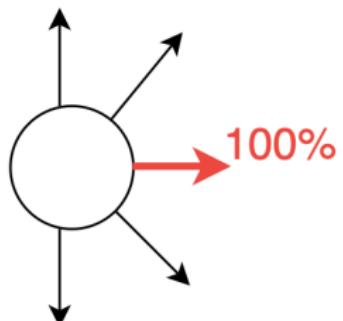
$$\pi(a | s) := P(A_{t+1} = a | S_t = s)$$

- **Important: policies are not part of the MDP!**



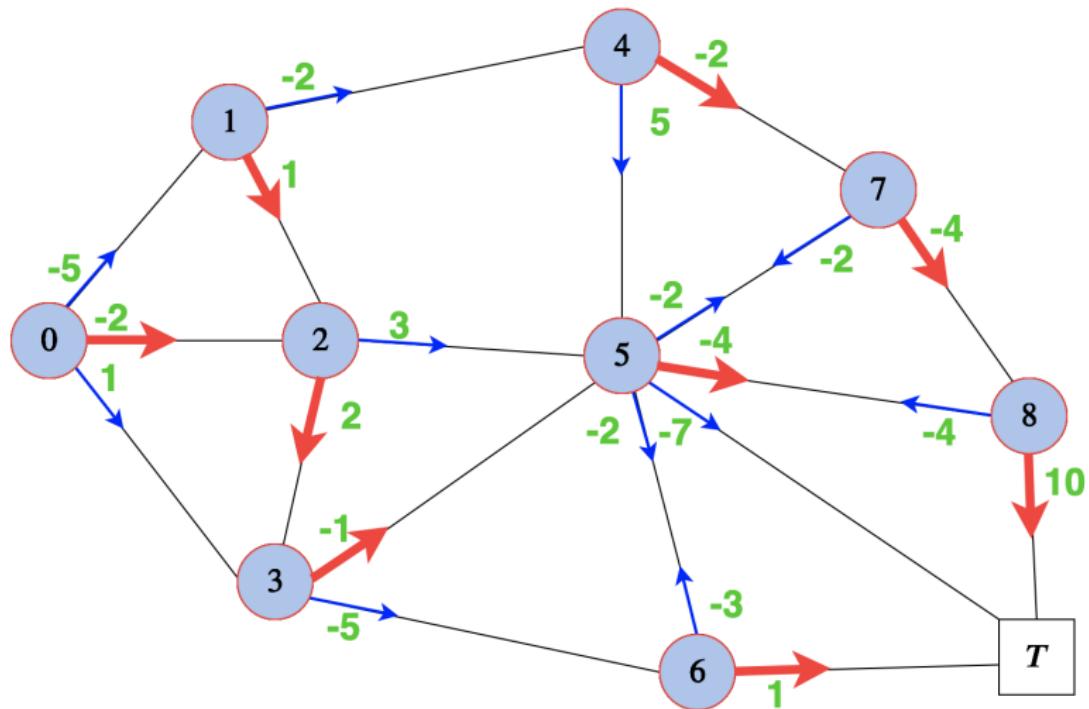
## Deterministic vs. Probabilistic Policy

deterministic                            probabilistic



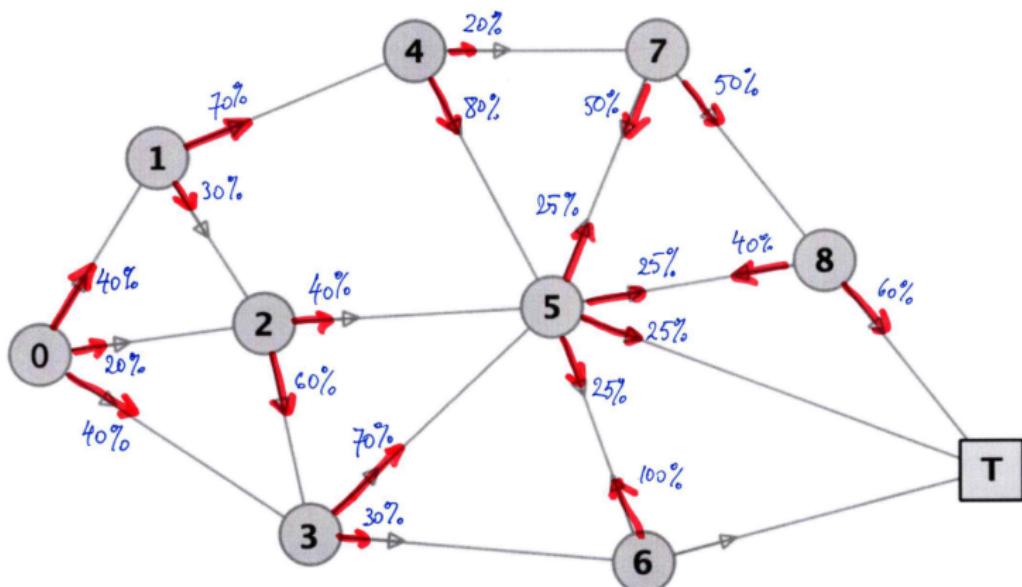
$$s \rightarrow \pi(a \mid s)$$

## MDP + Policy (deterministic)

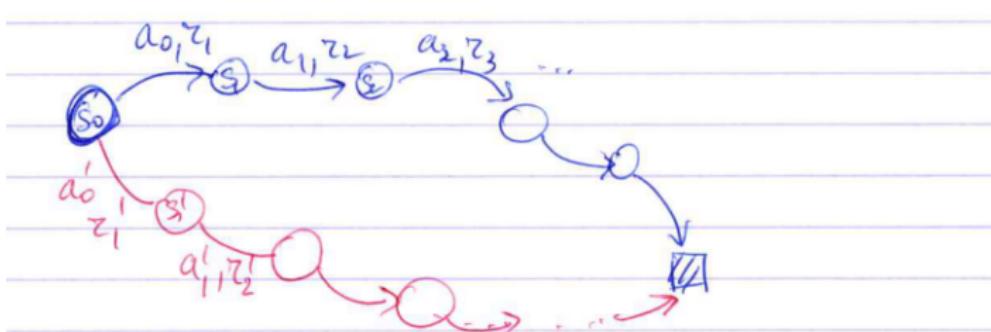
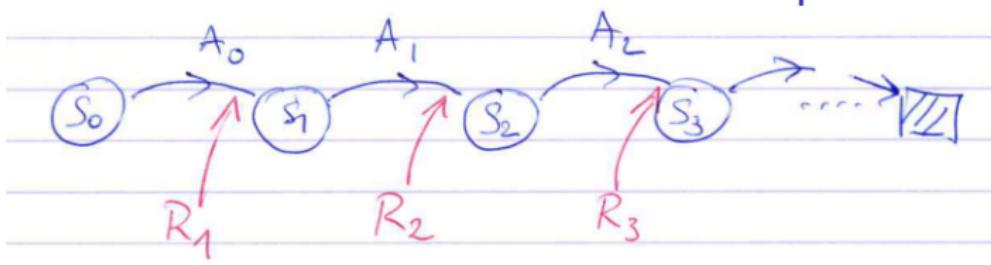


# MDP + Policy (Probabilistic/Stochastic)

PROBABILISTIC POLICY :  $\pi(a|s)$



## Stochastic actions and returns: interpretation



Long-term return:

$$G_0(s_0) = R_1 + R_2 + \dots + R_T = \sum_{t=1}^T R_t$$

## Return and Rewards

- The goal of the agent is to maximize the expected **(long-term) return**, a **(discounted)** sum over the immediate rewards received:

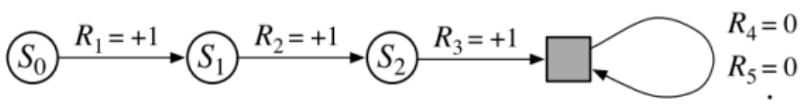
$$G_0 = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots = \sum_{k=1}^{\infty} \gamma^{k-1} R_k$$

or more generally:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=1}^{\infty} \gamma^{k-1} R_{t+k}$$

无限范围

- Episodic tasks** can be interpreted as infinite-horizon, if we represent episode termination as transition to an **absorbing state** with self-transitions and zero reward.



# Long-term (cumulative) return: some remarks

G<sub>t</sub> -> return  
R<sub>t</sub> -> rewards

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=1}^{\infty} \gamma^{k-1} R_{t+k}$$

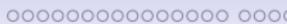
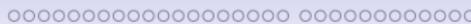
## Remarks

- Notice that  $G$  depends on both the **starting state** ( $s$ ) and the **policy applied** ( $\pi$ ):

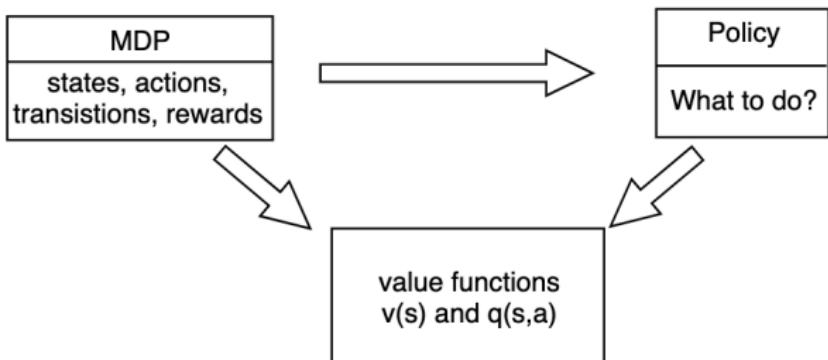
$$G \equiv G_{\pi}(s)$$

- Recursion relation:

$$G_t = R_{t+1} + \gamma G_{t+1}.$$



## Overview

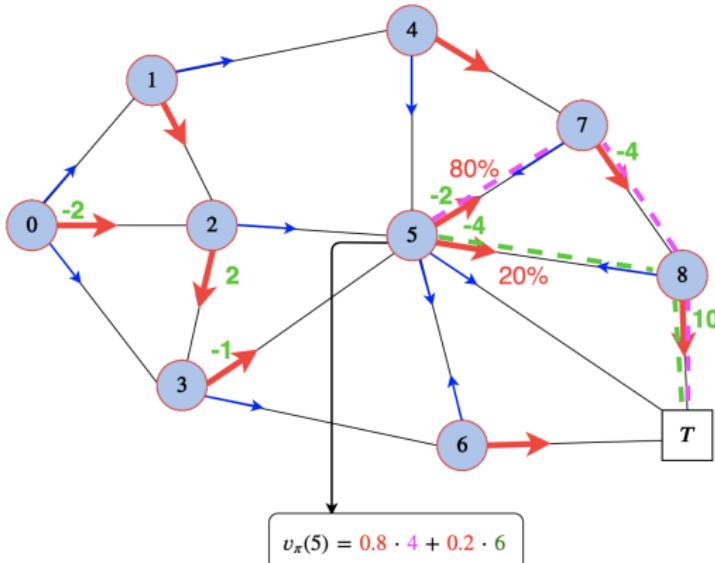




**MDP + Policy  $\rightarrow$  state value function  $v_\pi(s)$**

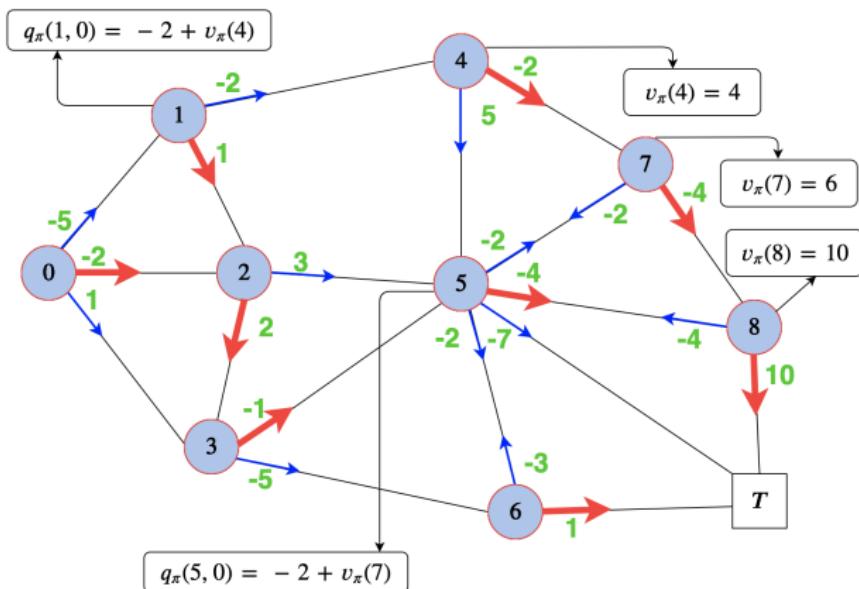
$v_\pi(s)$ : expected value of state  $s$  when actions are specified by  $\pi$ :

$$v_\pi(s) = E_\pi \left( \sum_{t=1}^T R_t \mid s_0 = s \right)$$



MDP + Policy  $\longrightarrow$  state-action value function  $q_{\pi}(s, a)$

$$q_{\pi}(s, a) = E_{\pi} \left( \sum_{t=1}^T R_t \mid s_0 = s, a_0 = a \right)$$



## Value functions: Tools for reasoning about future reward

- The **state value function**  $v_\pi(s)$  assigns to each state  $s$  the **expected total return** when **starting** in that state  $s$  and **applying policy**  $\pi$ ;

$$v_\pi(s) := E_\pi \left[ G_0 \mid S_0 = s \right] = E_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{k+1} \mid S_0 = s \right]$$

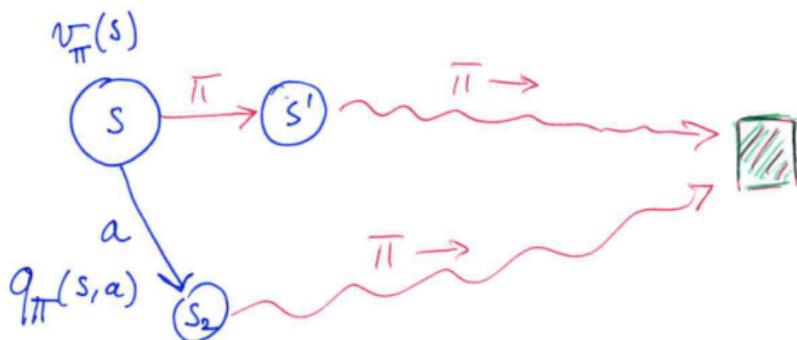
- The **state-action value function**  $q_\pi(s, a)$  of a policy  $\pi$  is:

$$q_\pi(s, a) := E_\pi \left[ G_0 \mid S_0 = s, A_0 = a \right] = E_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{k+1} \mid S_0 = s, A_0 = a \right]$$

- $q(s, a)$  state space typically much larger than  $v(s)$ -state space!  
Hence, more difficult to learn.

# Difference between value functions $v_\pi(s)$ and $q_\pi(s, a)$

$\pi = \text{policy}$

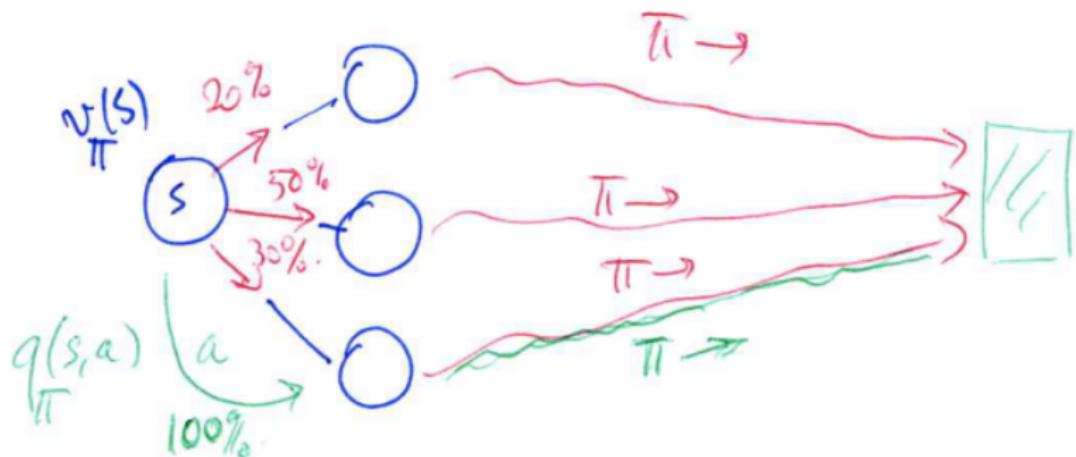


TODO: what does  $q(s, a)$  used for??

- $v_\pi(s)$ : policy  $\pi$  dictates **every action along the path**;
- $q_\pi(s, a)$ : **first action  $a$  is taken independently of the policy  $\pi$ ,** from then onward,  $\pi$  dictates the remaining actions taken along the rest of the path.

## Difference between value functions $v_{\pi}(s)$ and $q_{\pi}(s, a)$

Similar interpretation for stochastic policy.





## Relationship between value and value-action function

Notice that

$$v_\pi(s) = \sum_a \pi(a | s) q_\pi(s, a)$$

Indeed:

$$\begin{aligned} v_\pi(s) &= E_\pi \left[ G_0 \mid S_0 = s \right] \\ &= \sum_a E_\pi \left[ G_0 \mid S_0 = s, A_0 = a \right] \pi(a | s) \\ &= \sum_a q_\pi(s, a) \pi(a | s) \end{aligned}$$

# Outline

What is Reinforcement Learning (RL)?

Markov Decision Processes (MDP)

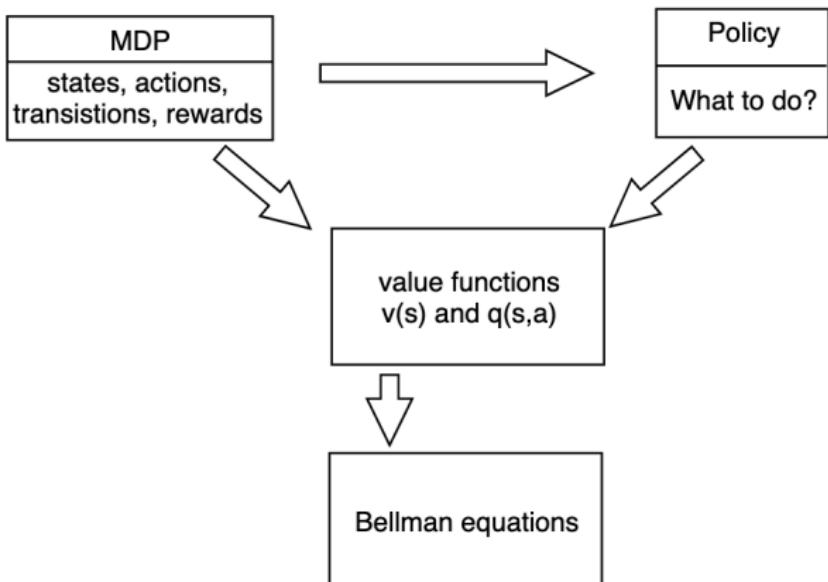
Policies, Value Functions and the Bellman Equation

Bellman equations

Bellman equations for optimality

Summary and Outlook

# Overview



## Bellman equation for value functions

**Value function**  $v_\pi(s)$ :  $v_\pi(s) = \sum_a \pi(a|s)q_\pi(s, a)$

**Value function**  $q_\pi(s, a)$       GO - long-term return, refer to p40

p - transition probability

$$q_\pi(s, a) = E_\pi \left[ G_0 \mid S_0 = s, A_0 = a \right]$$

$$= \sum E_\pi \left[ G_0 \mid S_0 = s, A_0 = a, S_1 = s' \right] p(s' \mid s, a)$$

TODO: ?? Bellman equation: The value of your current starting point is the reward you expect to get from being there, plus the value of wherever you land next.

$$= \sum_{s'} E_\pi \left[ R_1 + \gamma G_1 \mid S_0 = s, A_0 = a, S_1 = s' \right] p(s' \mid s, a)$$

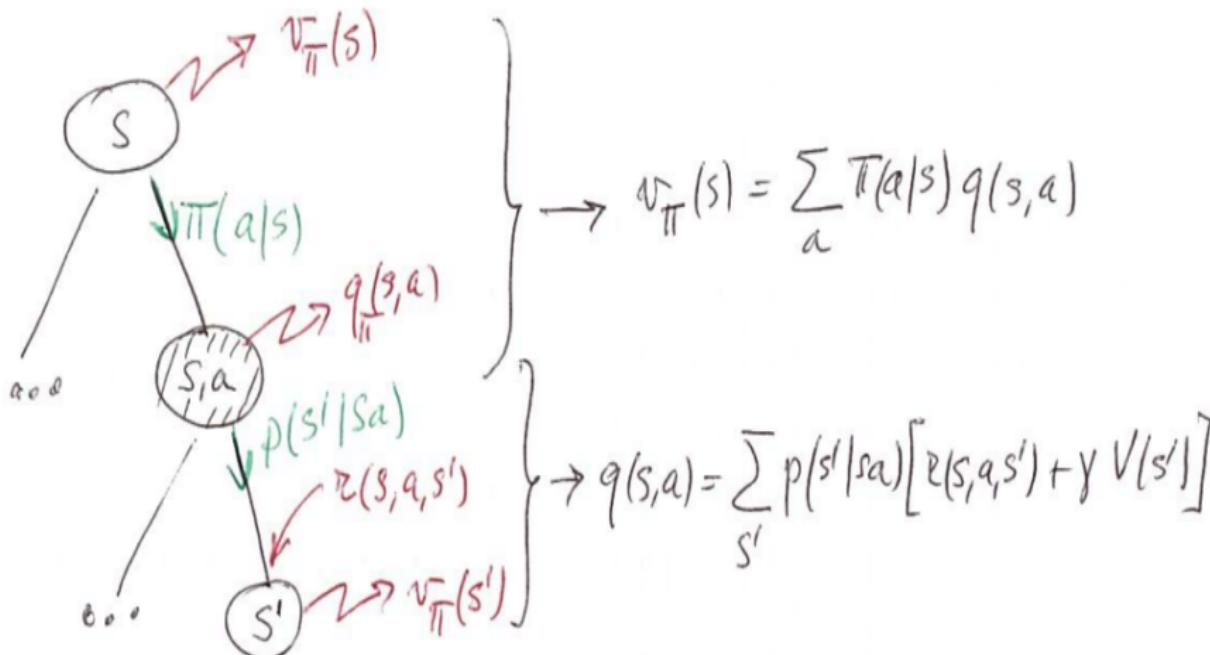
$$= \sum_{s'} \left[ r(s, a, s') + \gamma E_\pi(G_1 \mid S_1 = s') \right] p(s' \mid s, a)$$

$$= \sum_{s'} \left[ r(s, a, s') + \gamma E_\pi G_0(s') \right] p(s' \mid s, a)$$

$$= \sum_{s'} p(s' \mid s, a) \left[ r(s, a, s') + \gamma v_\pi(s') \right]$$

r(s,a,s'): expected immediate value

Bellman equations (schematically): Backup Diagram



## Bellman equation: Summary

- **Back-up**

$$v_\pi(s) = \sum_a \pi(a | s) q(s, a)$$

$$q(s, a) = \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma v(s')]$$

- **Combined into recursion equations**

$$v_\pi(s) = \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma v_\pi(s')]$$

$$q_\pi(s, a) = \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma \sum_{a'} \pi(a' | s') q_\pi(s', a')]$$

## Bellman equation: Summary

- The definition of  $v_\pi$  can be rewritten recursively by making use of the transition model, **yielding the Bellman equation:**

$$v_\pi(s) = \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma v_\pi(s')]$$

- This is a set of **linear equations**, one for each state, the solution of which defines the value of  $\pi$
- A similar recursive relation holds for Q-values:

$$q_\pi(s, a) = \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma \sum_{a'} \pi(a' | s') q_\pi(s', a')]$$

## Matrix form of Bellman equation

$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma v_{\pi}(s')]$$

We can rewrite this as:

$$\begin{aligned} v_{\pi}(s) &= \gamma \underbrace{\left( \sum_{s'} \underbrace{\left( \sum_a \pi(a | s) p(s' | s, a) \right)}_{P(s, s')} \right)}_{R(s, a)} v_{\pi}(s') \\ &\quad + \underbrace{\left( \sum_a \left( \sum_{s'} p(s' | s, a) r(s, a, s') \right) \right)}_{R(s, a)} \pi(a | s) \end{aligned}$$

The second term is immediate reward, and it's independent of a

In matrix notation:

P: transition function

$$\mathbf{v} = \gamma P \mathbf{v} + \mathbf{r}$$

In matrix form,

or again

$$(I - \gamma P) \mathbf{v} = \mathbf{r}$$

## Matrix form of Bellman equation

$$\mathbf{v} = \gamma P\mathbf{v} + \mathbf{r} \quad \text{or again} \quad (I - \gamma P)\mathbf{v} = \mathbf{r}$$

where

- square matrix  $P(s, s')$  is **transition probability**  $s \rightarrow s'$ :

$$P(s, s') := \sum_a \pi(a | s) p(s' | s, a)$$

- $\mathbf{r}(s)$  is the **expected (immediate) reward** in  $s$ :

$$\mathbf{r}(s) = \sum_a \pi(a | s) R(s, a) \quad \text{where} \quad R(s, a) = \sum_{s'} p(s' | s, a) r(s, a, s')$$

- Notice that  $\mathbf{v} = (I - \gamma P)^{-1} \mathbf{r} = (I + \gamma P + \gamma^2 P^2 + \dots) \mathbf{r}$   
based on taylor expansion:  $1/(1+x) = 1+x+x^{**}2+x^{**}3\dots$

## Example: MDP + value functions

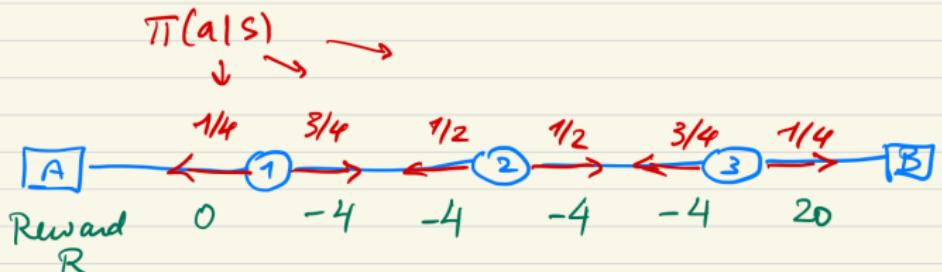
TODO: discount???

Example: - MDP + policy: see fig.

Assumptions:

① No discounting!  $\gamma = 1$

② Deterministic transitions:  $s \xrightarrow{a} s'$   
 $\text{i.e.: } p(s'|s,a) = \text{degenerate.}$   
 $\hookrightarrow \text{either 1 or 0.}$



## Example: MDP + value functions

$$P_{\pi}(s, s') = \sum_a \pi(a|s) p(s'|s, a)$$

$= P(\text{transition } s \rightarrow s' \text{ in 1 time step})$

		s' \ s				
		A	1	2	3	B
P =	A	1	0	0	0	0
	1	1/4	0	3/4	0	0
	2	0	1/2	0	1/2	0
	3	0	0	3/4	0	1/4
	B	0	0	0	0	1

A, B is absorbing point, so  $P(A,A)=P(B,B)=1$ ,  $P(1,1)=P(2,2)=\dots=0$

## Example: MDP + value functions

$$R(s,a) = \sum_{s'} p(s'|s,a) r(s,a,s')$$

= expected immediate reward when taking action  $a$  in state  $s$ .

$$R = \begin{array}{c|cc} & L & R \\ \hline A & 0 & 0 \\ 1 & 0 & -4 \\ 2 & -4 & -4 \\ 3 & -4 & 20 \\ B & 0 & 0 \end{array}$$

## Example: MDP + value functions

$$\mathcal{V}(s) = \sum_a R(s,a) \pi(a|s)$$

$$\mathcal{V} = \begin{bmatrix} \mathcal{V}(A) \\ \mathcal{V}(1) \\ \mathcal{V}(2) \\ \mathcal{V}(3) \\ \mathcal{V}(B) \end{bmatrix} = \text{diag} \left( \begin{bmatrix} L & R \\ 0 & 0 \\ 0 & -4 \\ -4 & -4 \\ -4 & 20 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} A & 1 & 2 & 3 & B \\ L & * & 1/4 & 1/2 & 3/4 & * \\ R & * & 3/4 & 1/2 & 1/4 & * \end{bmatrix} \right) \mathcal{V}(a|s)$$

## Example: MDP + value functions

$$\mathcal{R}(s) = \sum_a R(s, a) \pi(a|s)$$

= expected immediate return in  $s$

$$\mathcal{R} = \begin{bmatrix} 0 & (\frac{1}{4} \cdot 0 + \frac{3}{4}(-4)) = -3 & -4 & \frac{3}{4}(-4) + \frac{1}{4}20 = 2 & 0 \\ A & 1 & 2 & 3 & B \end{bmatrix}$$

## Example: MDP + value functions

Bellman eqs in matrix form:

$$\begin{aligned} \mathbf{v} &= \begin{bmatrix} v(A) \\ v(1) \\ v(2) \\ v(3) \\ v(B) \end{bmatrix} = \begin{bmatrix} 4 & 1 & 2 & 3 & B \\ 1 & 0 & 0 & 0 & 0 \\ 1/4 & 0 & 3/4 & 0 & 0 \\ 0 & 1/2 & 0 & 1/2 & 0 \\ 0 & 0 & 3/4 & 0 & 1/4 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \mathbf{v} + \begin{bmatrix} 0 \\ -3 \\ -4 \\ 2 \\ 0 \end{bmatrix} \\ \mathbf{v} &= P \cdot \mathbf{v} + \mathbf{z} \end{aligned}$$

$$\text{Example: } v(2) = \frac{1}{2} v(1) + \frac{1}{2} v(3) + (-4)$$

$$v(3) = \frac{3}{4} v(2) + \frac{1}{4} \underbrace{v(B)}_0 + 2$$

不动点定理？？

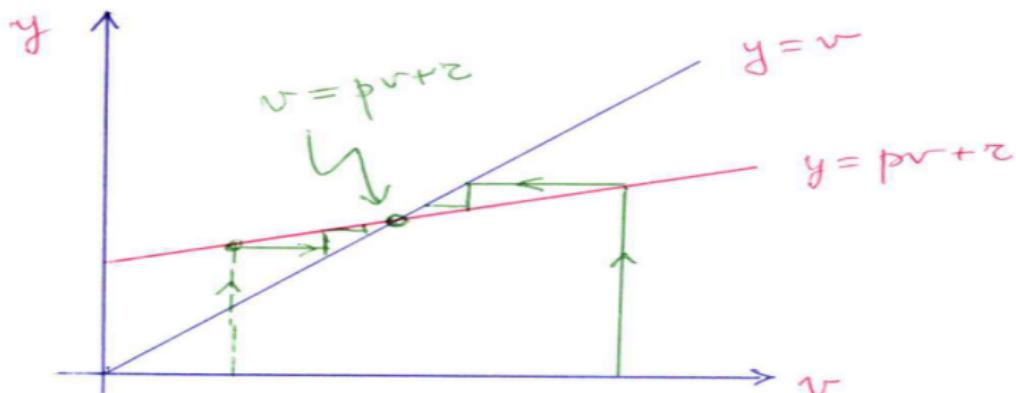
## Solving Bellman eqs: Iteration to Fix-Point

Matrix form of **Bellman equation** corresponds to fix-point:

$$\mathbf{v}_\pi = \gamma P_\pi \mathbf{v}_\pi + \mathbf{r}_\pi \quad \Rightarrow \quad \mathbf{v}_\pi = (I - \gamma P_\pi)^{-1} \mathbf{r}_\pi$$

**Iterative solution (fix-point solution):** update rule:

$$\mathbf{v}^{k+1} = \gamma P \mathbf{v}^k + \mathbf{r}$$



# Outline

What is Reinforcement Learning (RL)?

Markov Decision Processes (MDP)

Policies, Value Functions and the Bellman Equation

Bellman equations

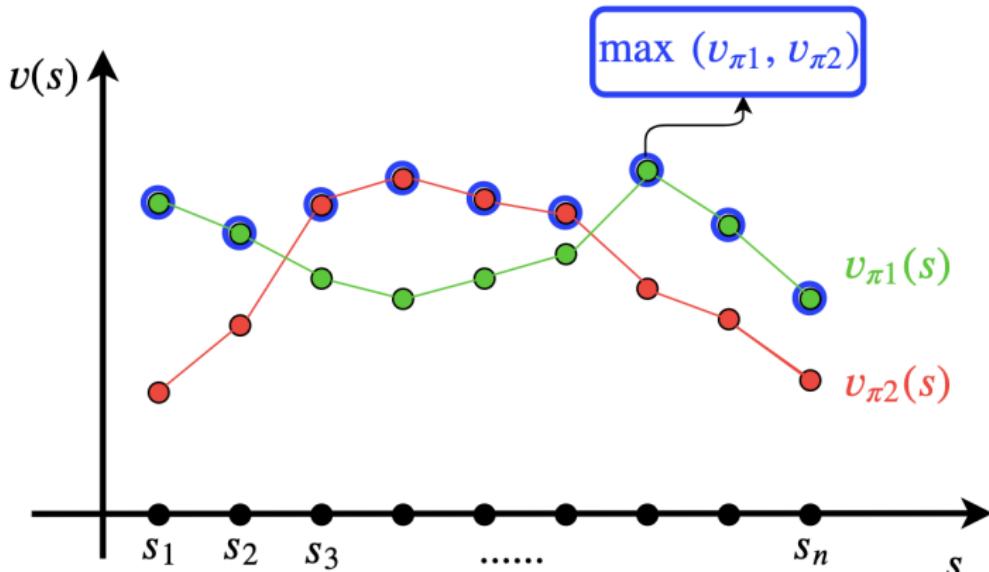
Bellman equations for optimality

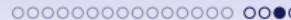
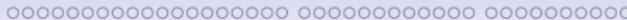
Summary and Outlook

## Optimal value functions

The optimal value functions are defined by the **pointwise maximum over policies**:

$$\forall s, a : \quad v^*(s) := \max_{\pi} v_{\pi}(s) \quad \text{and} \quad q^*(s, a) := \max_{\pi} q(s, a)$$



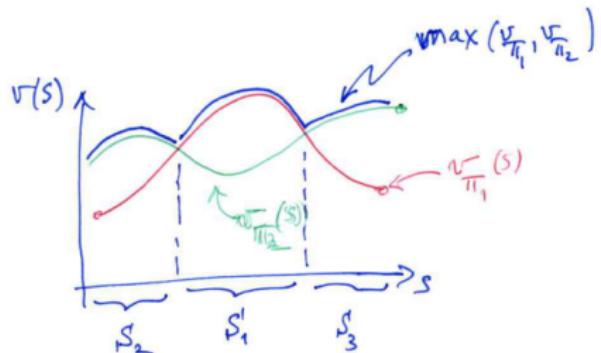


## Optimal policy $\pi^*$

There exists an **optimal policy**  $\pi^*$  such that its value functions corresponds to the optimal value functions:

$$v_{\pi^*}(s) = v^*(s) \quad \text{and} \quad q_{\pi^*}(s, a) = q^*(s, a)$$

## Intuition:



$$\pi^*(s) = \begin{cases} \pi_1(s) & \text{if } s \in S_2 \\ \pi_2(s) & \text{if } s \in S_1 \cup S_3 \end{cases}$$

## Optimal value functions

- Value functions define a partial ordering over policies:

$$\pi \succ \pi' \Rightarrow v_{\pi(s)} \geq v_{\pi'}(s), \forall s \in S$$

- There can be multiple optimal policies but they all share the same **optimal state-value function**:

$$v^*(s) = \max_{\pi} v_{\pi}(s), \quad \forall s \in S$$

- They also share the same **optimal action-value function**:

$$q^*(s, a) = \max_{\pi} q_{\pi}(s, a), \quad \forall s \in S, a \in A$$

# Optimal value functions: from weighted mean to max

- State-value function

- General:

$$v_{\pi}(s) = \sum_a \pi(a | s) q_{\pi}(s, a)$$

- Optimal

$$v^*(s) = v_{\pi^*}(s) = \max_a q_{\pi^*}(s, a) = \max_a q^*(s, a)$$

$$v^*(s) = \max_a q^*(s, a)$$

# Optimal value functions: from weighted mean to max

- State-action value function
  - General:

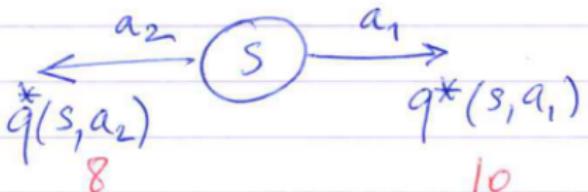
$$q_{\pi}(s, a) = \sum_{s'} p(s' | s, a)(r(s, a, s') + \gamma v_{\pi}(s))$$

Optimal:

$$q_{\pi^*}(s, a) = \sum_{s'} p(s' | s, a)(r(s, a, s') + \gamma v_{\pi^*}(s))$$

$$q^*(s, a) = \sum_{s'} p(s' | s, a)(r(s, a, s') + \gamma v^*(s))$$

## Bellman Optimality Equations



$q^*(s, a_1)$  = best possible expected return  
 - if taking action  $\underline{a_1}$  in  $s$ .  $\underline{= 10}$

$q^*(s, a_2)$  = best possible expected return  
 - if taking action  $\underline{a_2}$  in  $s$ .  $\underline{= 8}$

$v^*(s)$  = best possible expected return in  $s$ .

$$\boxed{v^*(s) = \max_a q^*(s, a)} (= 10)$$

## Bellman Optimality Equations

Deterministic transition:  $s \xrightarrow{a} s'$

$$q^*(s, a) = r(s, a, s') + \gamma v^*(s')$$

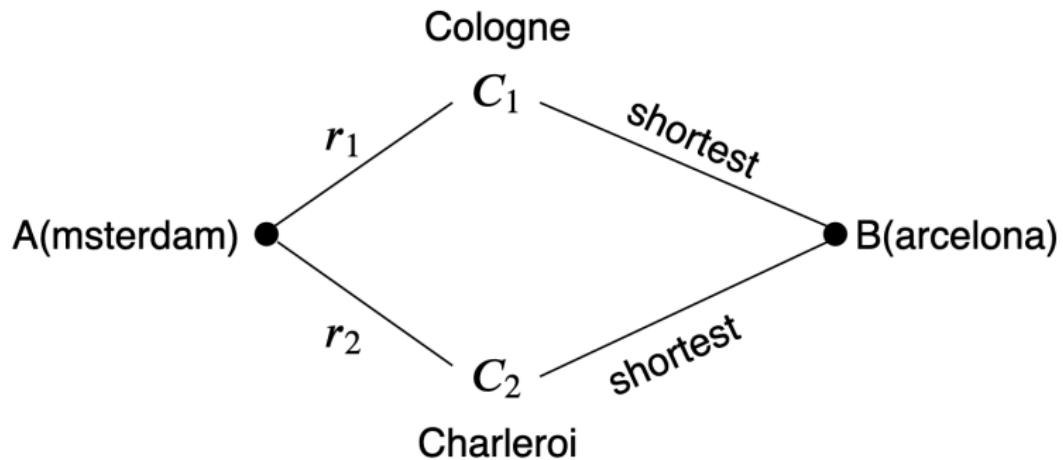
Best possible expected  
Return when  
Committed to action a

Best possible  
expected return  
in this state  $s'$

$$q^*(s, a) = \sum_{s'} p(s'|sa) [r(s, a, s') + \gamma v^*(s')]$$

General

## Bellman Optimality equation: Travel Distance Analogy



$$d^*(A, B) = \min_{C_i} \{r_i + d^*(C_i, B)\}$$

## Bellman optimality conditions (for deterministic transitions)

- Travel Analogy: Shortest distance paths:

$$d^*(A, B) = \min_{C_i} \{r_i + d^*(C_i, B)\}$$

- Deterministic transitions:  $s \xrightarrow{a} s_a$

$$v^*(s) = \max_a \{r(s, a, s_a) + v^*(s_a)\}$$

$$q^*(s, a) = r(s, a, s_a) + v^*(s) = r(s, a, s_a) + \max_{a'} q^*(s_a, a')$$

## Bellman optimality equations (General form)

- **Optimal state value function**

$$v^*(s) = \max_a q^*(s, a)$$

- **Optimal state-action value function**

$$q^*(s, a) = \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma v^*(s')]$$

- **Combined**

$$v^*(s) = \max_{a \in A} \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma v^*(s')]$$

$$q^*(s, a) = \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma \max_{a' \in A} q^*(s', a')]$$

# Backup Diagram for Bellman Optimality Equations

**Optimize over actions!**

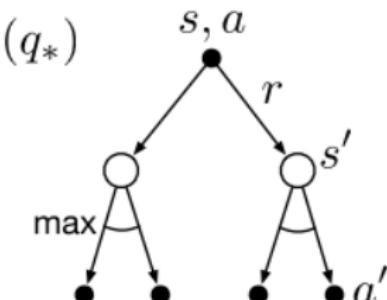
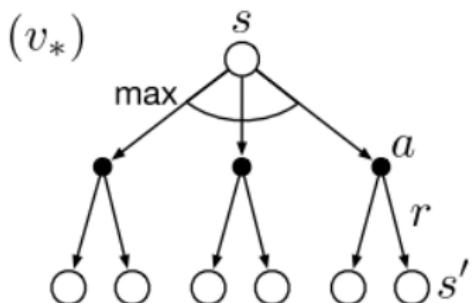


Figure 3.5: Backup diagrams for  $v_*$  and  $q_*$

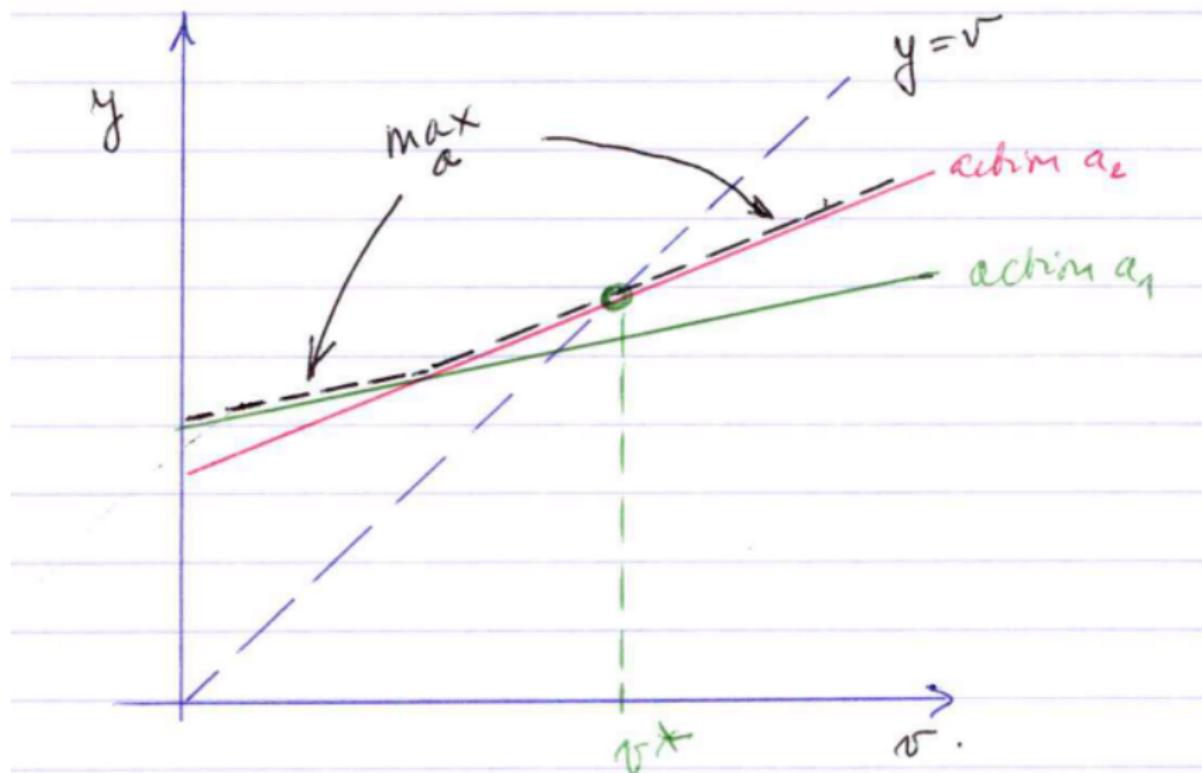
## Bellman optimality equation in matrix form

$$\begin{aligned} v^*(s) &= \max_{a \in A} \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma v^*(s')] \\ &= \max_a \left( R(s, a) + \gamma \sum_{s'} \underbrace{p(s' | s, a)}_{T_a(s, s')} v^*(s') \right) \\ &= \max_a \left( R(s, a) + \gamma \sum_{s'} T_a(s, s') v^*(s') \right) \end{aligned}$$

or in matrix notation:

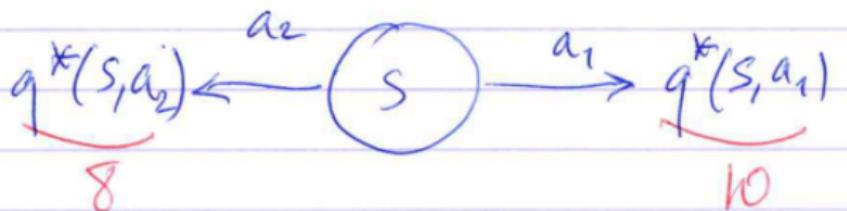
$$\mathbf{v}^* = \max_a (R_a + \gamma T_a \mathbf{v}^*)$$

## Bellman optimality equation



## Importance of $v^*$ and $q^*$

- Markov property: Optimal decision only depends on current state:
- **Greedification:** In state  $s$  pick action  $a$  with highest  $q^*(s, a)$



$$\pi_g : s \mapsto a_1$$

## Bellman Optimality Conditions

- $v^*$  and  $q^*$  satisfy a set of (non-linear) equations analogous to Bellman equations for  $v_\pi$  and  $q_\pi$ ;
- these non-linear equations **do not refer** explicitly to the optimal policy;
- As a consequence we can find the optimal policy  $\pi^*$  by
  1. first solving these equations to find  $v^*(s)$  (and  $q^*(s, a)$ ) and
  2. then use **greedification** to determine the corresponding optimal policy  $\pi^*$ :

$$\forall s : \quad a_{opt} = \arg \max_a q^*(s, a)$$

## Example: Bellman optimality conditions

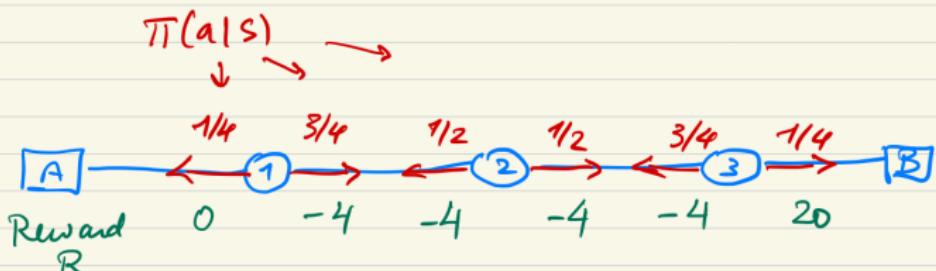
MDP

Example: - MDP + policy : see fig.

Assumptions:

① No discounting,  $\gamma = 1$

② Deterministic transitions,  $s \xrightarrow{a} s'$   
 i.e.:  $p(s'|s,a) = \text{degenerate}$ .  
 ↳ either 1 or 0.



## Example: Bellman optimality conditions

### Bellman Optimality Conditions

$$\begin{aligned} v^*(s) &= \max_a \sum_{s'} p(s'|sa) [r(s,a,s') + \gamma v^*(s')] \\ &= \max_a \left\{ R(s,a) + \gamma \sum_{s'} p(s'|sa) v^*(s') \right\} \end{aligned}$$

Figure: General expression for  $v^*$

# Bellman optimality conditions: Deterministic transition

Deterministic transition :  $s \xrightarrow{a} s_a$

$$\pi^*(s) = \max_a [R(s,a) + \gamma v^*(s_a)]$$

In this example, optimal policy is simple

$\pi^*$ : always move right

$$v^* = 12 \quad v^* = 16 \quad v^* = 20$$



Eg.  $s=2$

$a=L \rightarrow s_a = 1, \pi^*(s_a) = 12, R(s,a) = -4$   
 $R(s,a) + v^*(s_a) = -4 + 12 = 8$

$a=R \rightarrow s_a = 3, \pi^*(s_a) = 20, R(s,a) = -4$   
 $R(s,a) + v^*(s_a) = -4 + 20 = 16$

## Bellman optimality conditions: Deterministic transition

Computing  $v^*$ , Worked example.

$$R =$$

0	0
0	-4
-4	-4
-4	20
0	0

Deterministic:  $s \xrightarrow{a} s_a$

$$v^*(s) = \max_a [R(s_a) + \gamma v^*(s_a)]$$

↓

$\gamma = 1$

Figure:  $R(s, a)$  and  $v^*(s)$  for deterministic transitions

## Bellman optimality condition: Computing $v^*$ using iteration

Initialise

$t=0$

$\rightarrow t=1$

$$v^* = 0$$

$$v^* = \max_a \left( \begin{array}{c|cc} & L & R \\ \hline a & 0 & 0 \\ & 0 & -4 \\ & -4 & -4 \\ & -4 & 20 \\ \hline 0 & 0 \end{array} \right) + \left( \begin{array}{c|cc} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ \hline 0 & 3 \end{array} \right) = \left( \begin{array}{c} 0 \\ 0 \\ -4 \\ 20 \\ 0 \end{array} \right)$$

$$\underline{t=2:} \quad v^* = \max_a$$

$$\left\{ \left( \begin{array}{c|cc} 0 & 0 \\ 0 & -4 \\ -4 & -4 \\ -4 & 20 \\ \hline 0 & 0 \end{array} \right) + \left( \begin{array}{c|cc} 0 & 0 \\ 0 & -4 \\ 0 & 20 \\ -4 & 0 \\ \hline 0 & 0 \end{array} \right) \right\} = \left( \begin{array}{c} 0 \\ 0 \\ 16 \\ 20 \\ 0 \end{array} \right)$$

## Example: Bellman optimality conditions

$$t=3: \quad v^k = \max_a \left\{ \begin{array}{c} \left[ \begin{array}{cc} 0 & 0 \\ 0 & -4 \\ -4 & -4 \\ -4 & 20 \\ 0 & 0 \end{array} \right] + \left[ \begin{array}{cc} 0 & 0 \\ 0 & 16 \\ 0 & 20 \\ 16 & 0 \\ 0 & 0 \end{array} \right] \end{array} \right\} = \left[ \begin{array}{c} 0 \\ 12 \\ 16 \\ 20 \\ 0 \end{array} \right]$$

Figure:  $t = 3$

## Example: Bellman optimality conditions

$t=4$

$$v^* = \max_a \left\{ \begin{array}{c} \left[ \begin{array}{cc} 0 & 0 \\ 0 & -4 \\ -4 & -4 \\ -4 & 20 \\ 0 & 0 \end{array} \right] + \left[ \begin{array}{cc} 0 & 0 \\ 0 & 16 \\ 12 & 20 \\ 16 & 0 \\ 0 & 0 \end{array} \right] \end{array} \right\} = \left[ \begin{array}{c} 0 \\ 12 \\ 16 \\ 20 \\ 0 \end{array} \right]$$

Converged!

Figure: Convergence for  $t = 4$

For known MDP: compute  $q^*(s, a)$  based on  $v^*(s)$

$$s \xrightarrow{a} s_a$$

$$q^*(s, a) = z(s, a, s_a) + \gamma v^*(s_a)$$

$$v^* = \begin{pmatrix} 0 \\ 12 \\ 16 \\ 20 \\ 0 \end{pmatrix} \quad q^* = R(s, a) + \gamma v^*(s_a) = \begin{matrix} & L & R \\ A & \boxed{0 \ 0} \\ 1 & 0 \ -4 \\ 2 & -4 \ -4 \\ 3 & -4 \ 20 \\ B & 0 \ 0 \end{matrix} + \begin{matrix} & L & R \\ A & \boxed{0 \ 0} \\ 1 & 0 \ 16 \\ 2 & 12 \ 20 \\ 3 & 16 \ 0 \\ B & 0 \ 0 \end{matrix} = \begin{pmatrix} 0 & 0 \\ 0 & 12 \\ 8 & 16 \\ 12 & 20 \\ 0 & 0 \end{pmatrix}$$

$$a^*(s) = \arg \max a q^*(s, a)$$

Hence  $a^*(s) = R$  for  $s=1, 2, 3$

Figure: Computing  $q^*(s, a)$  from  $v^*(s)$

## Bellman optimality equations

If MDP is fully known, it suffices to compute  $v^*$ , as  $q^*$  can then be derived:

- **Optimal state-action value function**

$$q^*(s, a) = \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma v^*(s')]$$

- **Recall: Optimal state value function**

$$v^*(s) = \max_a q^*(s, a)$$

## Model-based vs model-free

- **Model-based:** the MDP =  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$  is completely specified;
  - Solve the Bellman equations!
- **Model-free:** only **direct experience**, i.e. sample paths (states, actions and rewards) are given. Put differently, only experience-based information is given!
  - Random search but Bellman equations allow to propagate values!

# Outline

What is Reinforcement Learning (RL)?

Markov Decision Processes (MDP)

Policies, Value Functions and the Bellman Equation

Bellman equations

Bellman equations for optimality

Summary and Outlook

# Bellman equations: General versus Optimal

- **State value functions:**

- General:

$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma v_{\pi}(s')]$$

- Optimal:

$$v^*(s) = \max_a \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma v^*(s')]$$

- **State-action value function**

- General:

$$q_{\pi}(s, a) = \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a')]$$

- Optimal:

$$q^*(s, a) = \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma \max_{a'} q^*(s', a')]$$

# Taxonomy of RL problems

	<b>Prediction</b> <i>Estimation:</i> <i>Given <math>\pi</math>, what is <math>v</math>?</i>	<b>(Optimal) Control</b> <i>Optimisation:</i> <i>What is optimal <math>\pi</math>?</i>
model-based (MDP given)	Policy evaluation using Dyn. Programming (DP)	Policy improvement (+ Policy evaluation) = Policy iteration
model-free (MDP unknown)	Monte Carlo (MC) Temporal Diff <sup>ing</sup> (TD) = "impatient MC" <i>bootstrapping!</i>	Generalized Policy Iteration <i>"simultaneous"</i>