# Deep Learning

## Lecture 8

from really good course in AI masters (https://ozonmasters.ru/reinforcementlearning).

# Recap

- Semantic segmentation problem
- Upsampling
- Architectures
- Panoptic / Instance segmentation

# What is Reinforcement Learning?

Let's start from...

# Supervised Learning Problem

**Supervised Learning case:**

Given Dataset $D := \{(X_i, y_i)\}$

Learn a function that will predict y from X: $f_\theta: X \to y$

e.g. find parameters theta that will minimize: $L(f_\theta(X_i), y_i)$ where $L$ is a loss function

**Standard Assumptions:**
- Samples in Dataset are I.I.D
- We have ground truth labels $y$

# No ground truth answers
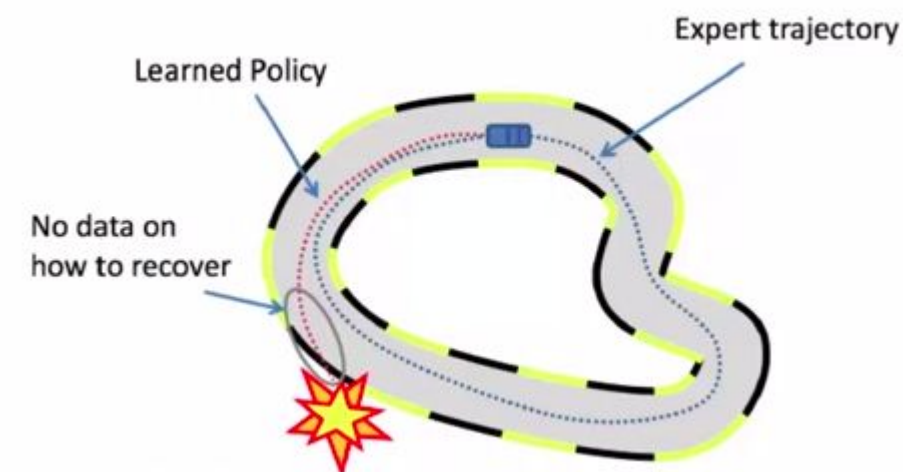
**You don't have answers at all**

**Your answers are not good enough**

# Choice matters

Assume that we have expert trajectories, i.e. sufficiently good answers:

- Treat trajectories as a dataset:

  $D = \{(x_1, a_1), .. (x_N, a_N)\}$
- Train with Supervised Learning
- Done?:)



HOW TO BE AN EXPERT DRIVER
by Al Esper, Chief Test Driver, Ford Motor Company





Expert trajectory

Learned Policy

No data on how to recover

# Choice matters

New Plan ([DAGGER algorithm](#)):

1. Train a model from human trajectories :
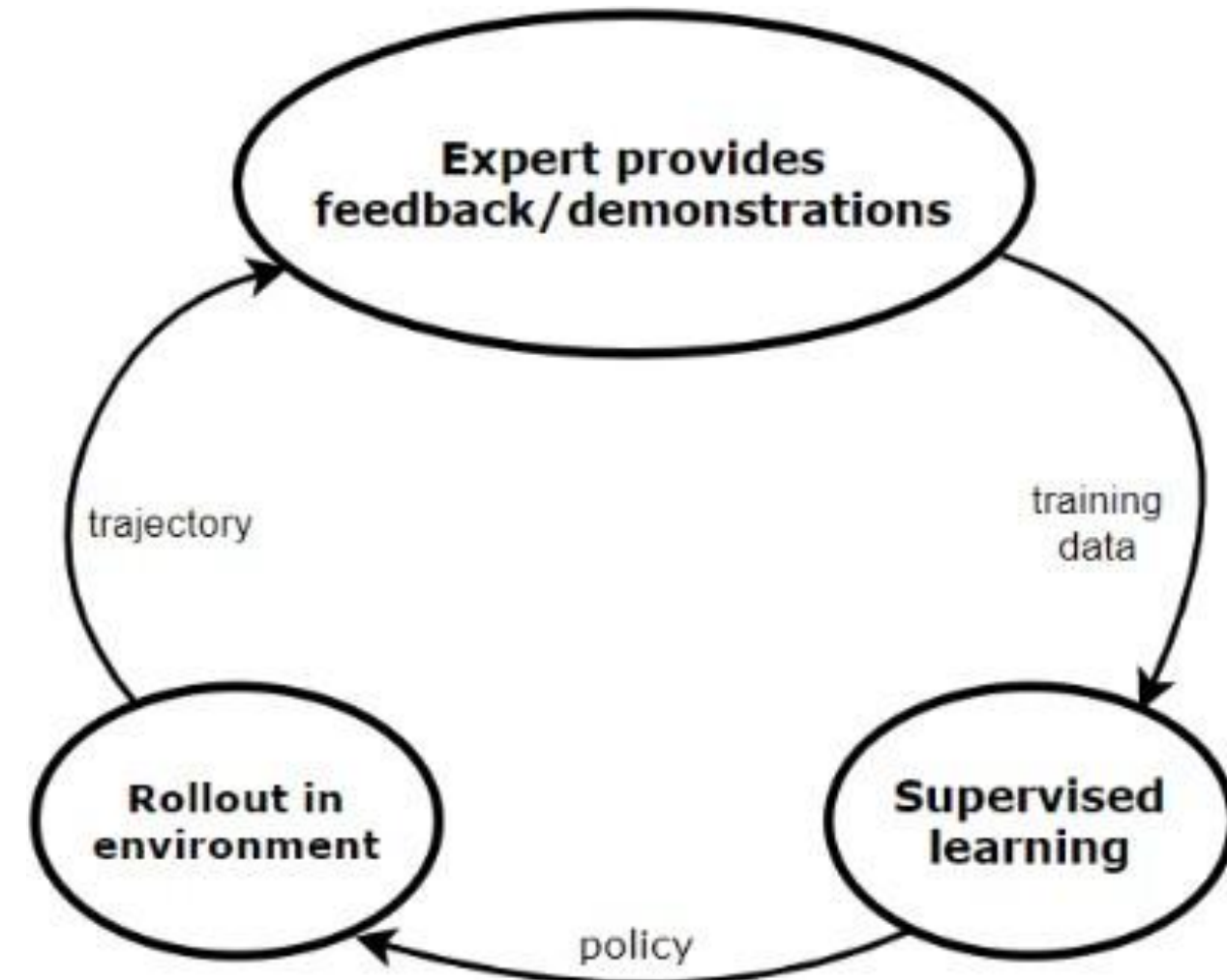
   $D_0 = \{(x_1, a_1), ..(x_N, a_N)\}$

2. Run the model to get new trajectories:

   $D' = \{(x_1, ?), ..(x_N, ?)\}$

3. Ask humans to label $D'$ with actions $a_t$

4. Aggregate: $D_1 \leftarrow D_0 \cup D'$

5. Repeat

# Choice matters

But this is really hard to do: 3. Ask humans to label $D'$ with actions $a_t$

# Reinforcement learning

**If You know what you want, but don't know how to do it...**    USE REWARDS!



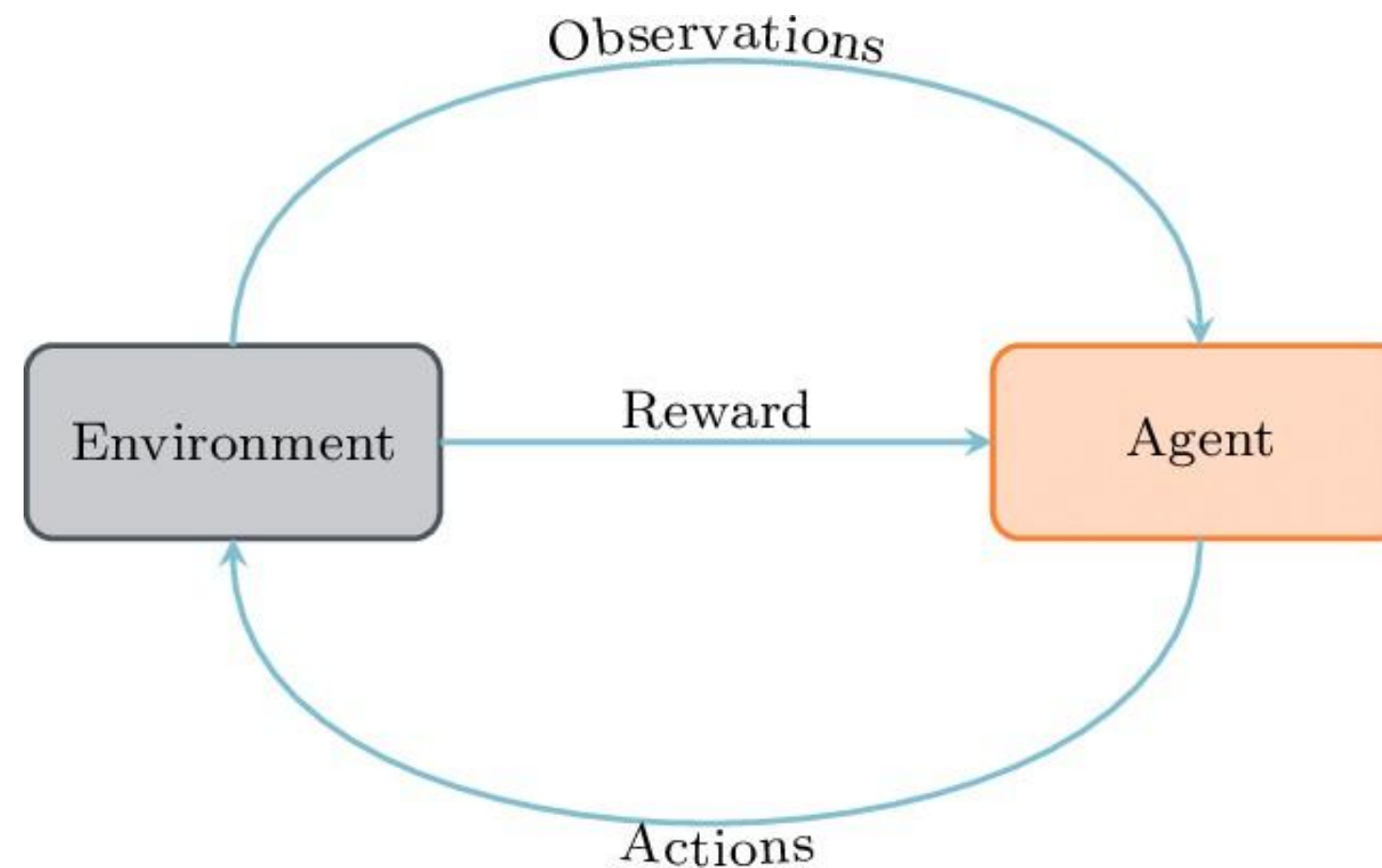**WIN: +1**

**LOSE: -1**

**Assumptions:**

- **It's easy to compute reward**
- **You can express your goals with rewards!**

# Reinforcement Learning Problem

You have **Agent** and **Environment** that interact with each other:

- Agent's actions change the state environment
- After each action agent receives new state and reward

Interaction with environment  is typically divided into **episodes.**

# Reinforcement Learning Problem

Agent has a policy:  $\pi(\text{action}\,|\,\text{observations from env})$

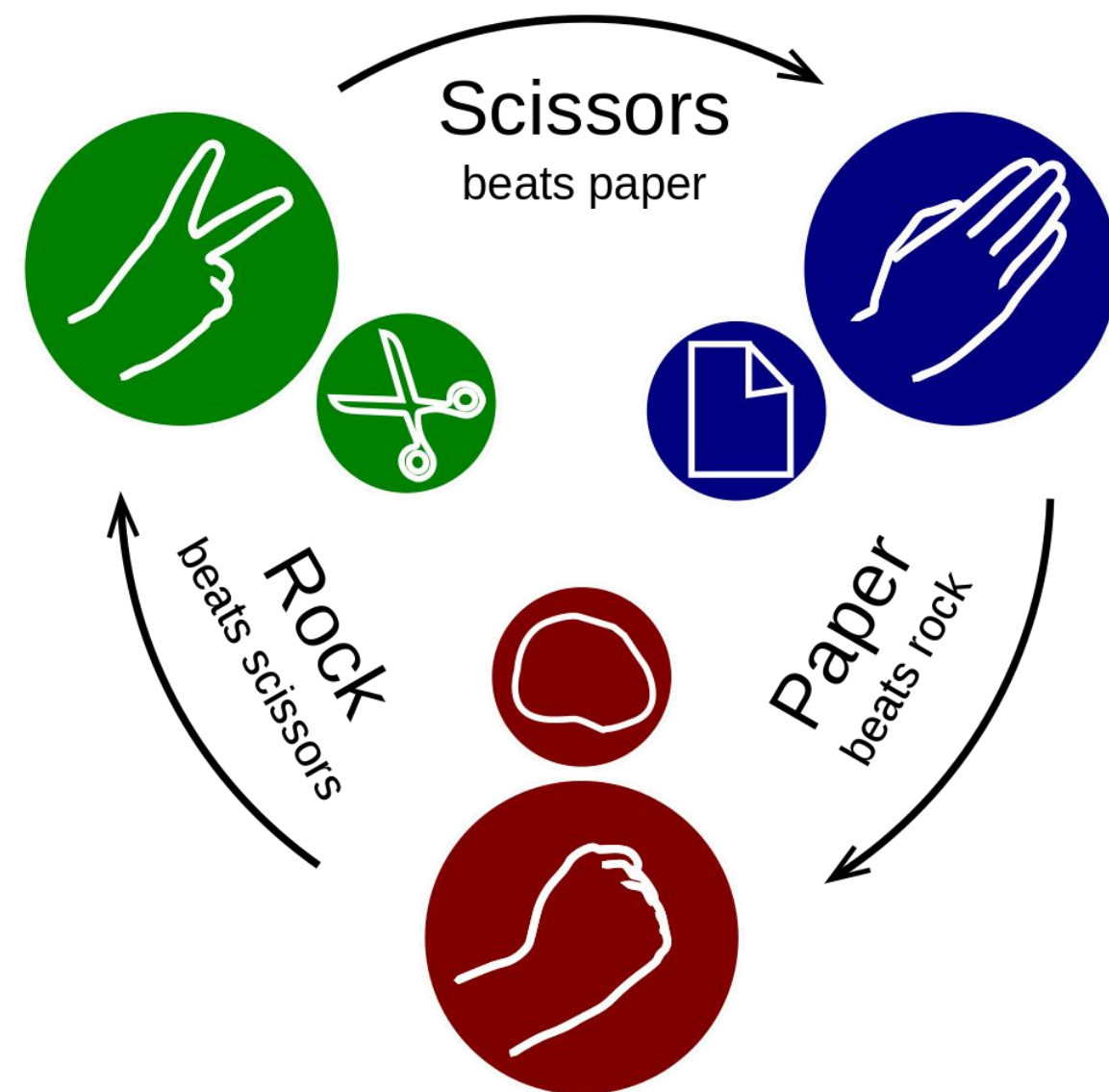Agent learns its policy via **Trial and Error**!

The goal is to find a policy that maximizes **total expected reward:**
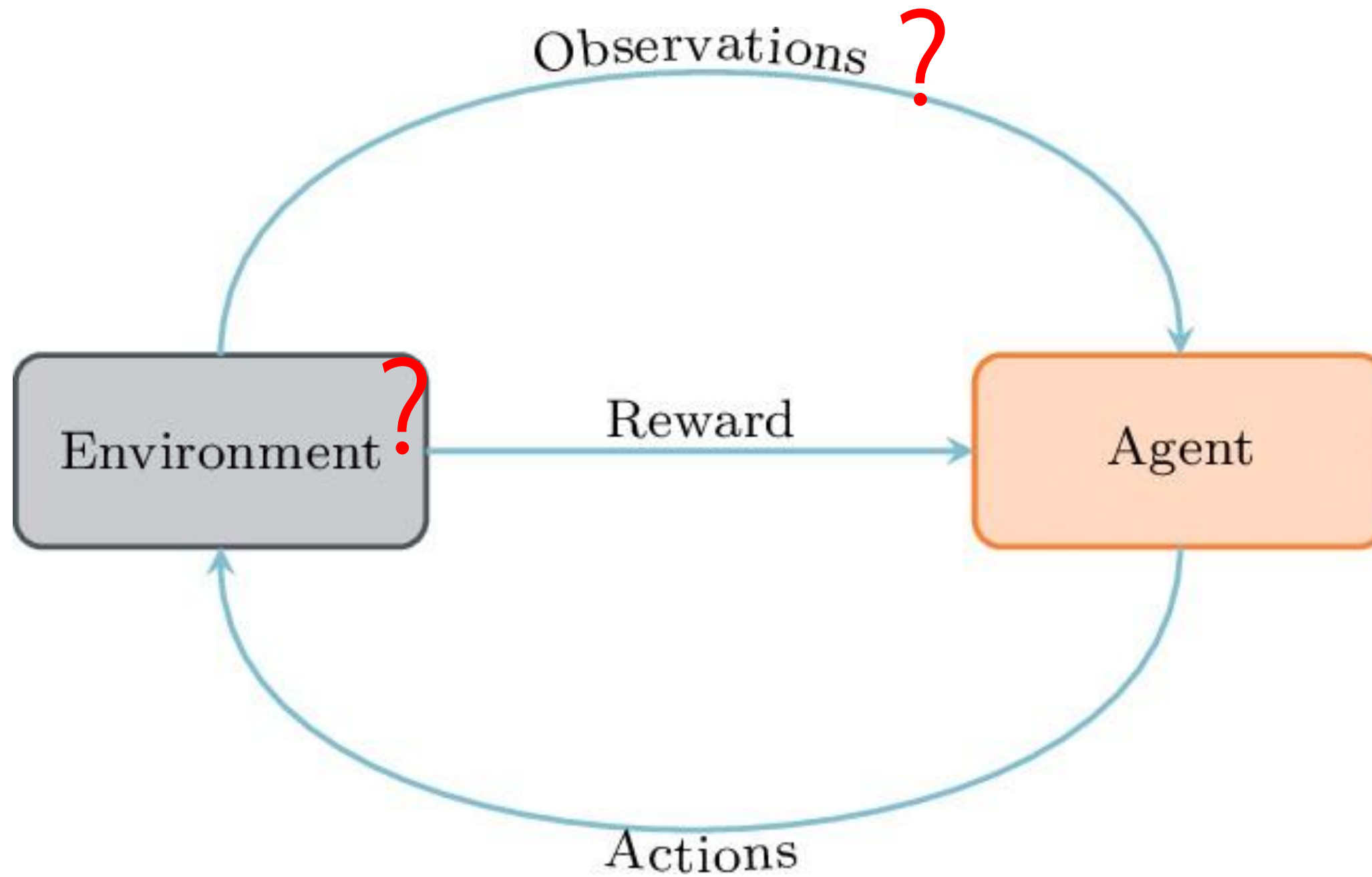
$$\text{maximize}_\pi \text{E}_\pi[\textstyle\sum^{\text{T}}_{t=0} r_t]$$

Why we need $\text{E}_\pi$ ?

A non-deterministic policy or environment lead to

a distribution of total rewards!

Why not use $\text{max}_\pi[\textstyle\sum^{\text{T}}_{t=0} r_t]$, $\text{min}_\pi[\textstyle\sum^{\text{T}}_{t=0} r_t]$ ?



Scissors
beats paper

Rock
beats scissors

Paper
beats rock

# Reinforcement Learning Problem

# Environment and Observation

What should an agent observe?

- Wheel speed
- Acceleration
- LiDAR
- Battery

- Map of the apartment
- Location

Is this enough?
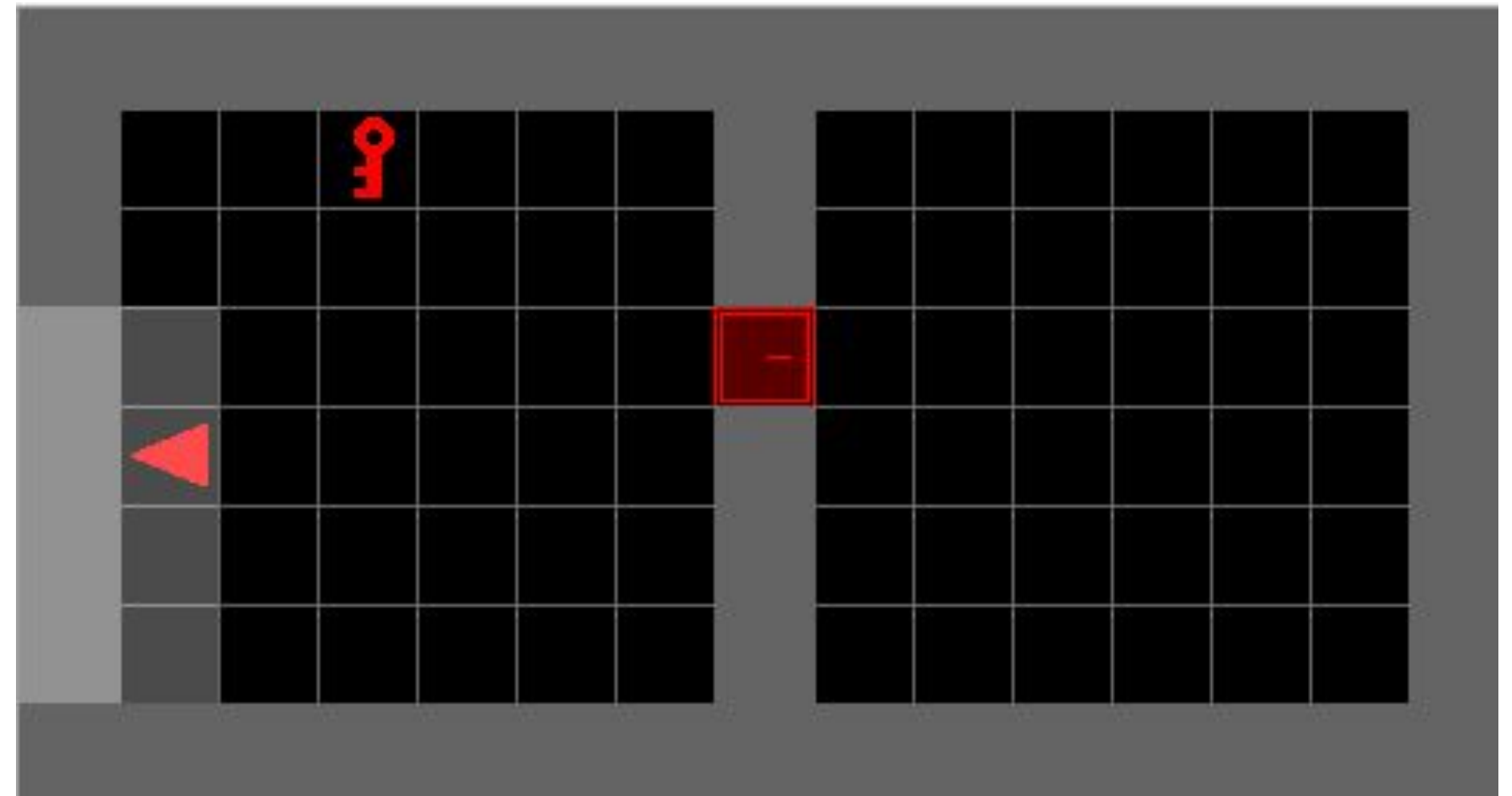
Does agent need past observations?

# Markovian Property

**Task:** Open the red door with the key

**Details:** Agent starts at random location
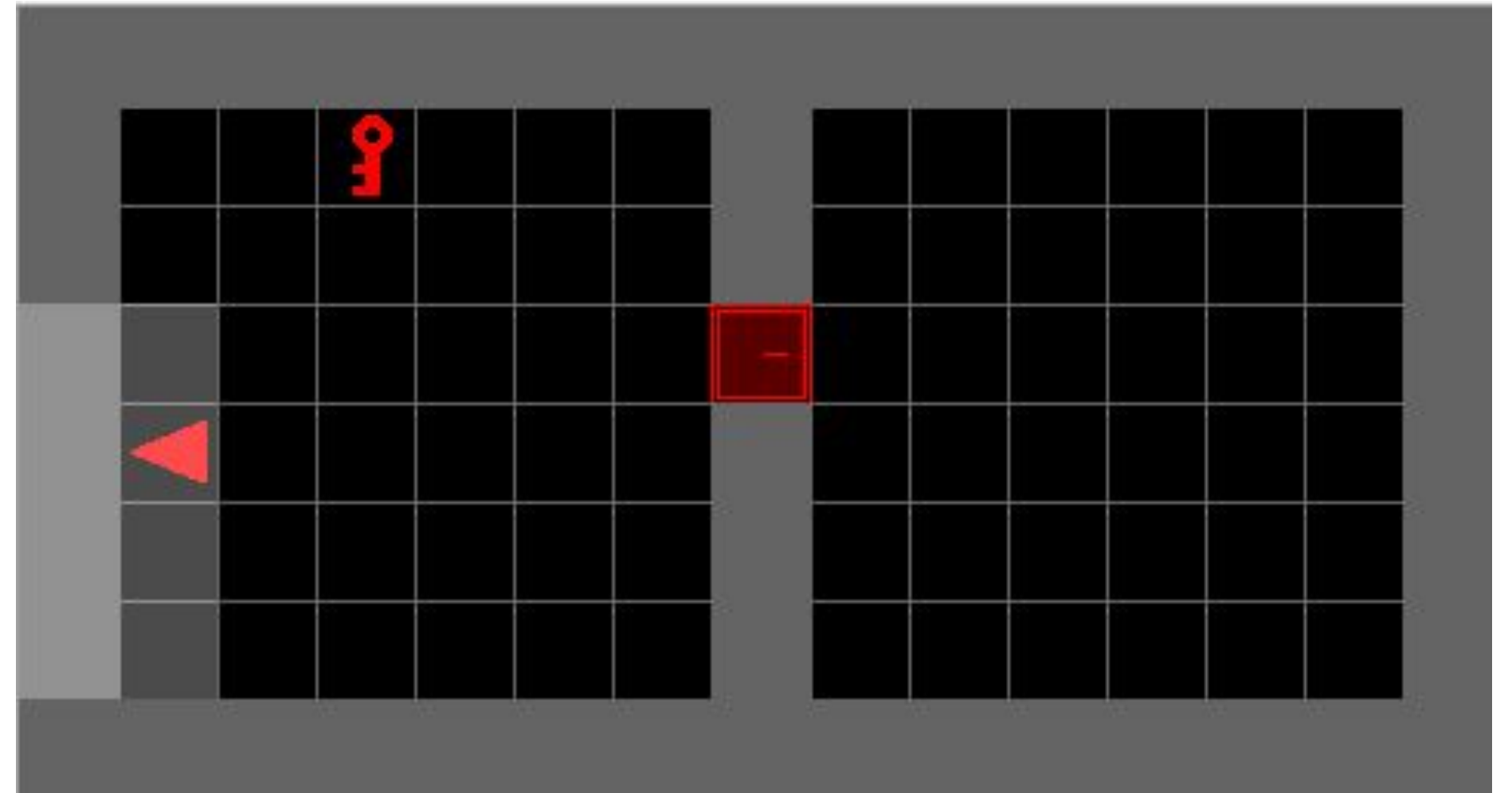
**Actions:**

- move up/left/right/down
- pick up an object
- apply an object to the door
  (when near the door)

# Markovian Property

Which observations are enough to
learn the optimal policy?

1. Agent's coordinates, and previous action
2. Full image of the maze
3. Agent's coordinates and does it has key



For 2 and 3 agent doesn't need to remember it's history:

$$P(o_{t+1}, r_{t+1} | o_t, a_t) = P(o_{t+1}, r_{t+1} | o_t, a_t, ..., o_1, a_1, o_0, a_0)$$

Markovian property: **"The future is independent of the past given the present."**

# Markov Decision Process

MDP is a 5-tuple $< S, A, R, T, \gamma >$:

- $S$ is a set of states
- $A$ is a set of actions
- $R : S \times A \to$ R is a reward function
- $T : S \times A \times S \to [0, 1]$ is a transition function
  $T(s, a, s') = P(S_{t+1} = s' | S_t = s, A_t =$
- $a)$ $\gamma \in [0, 1]$ is a discount factor

Discount factor $\gamma$ determines how much we should care about the future!

Given Agent's policy $\pi$, RL objective become: $E_\pi[\sum_{t=0}^{T} \gamma^t r_t]$
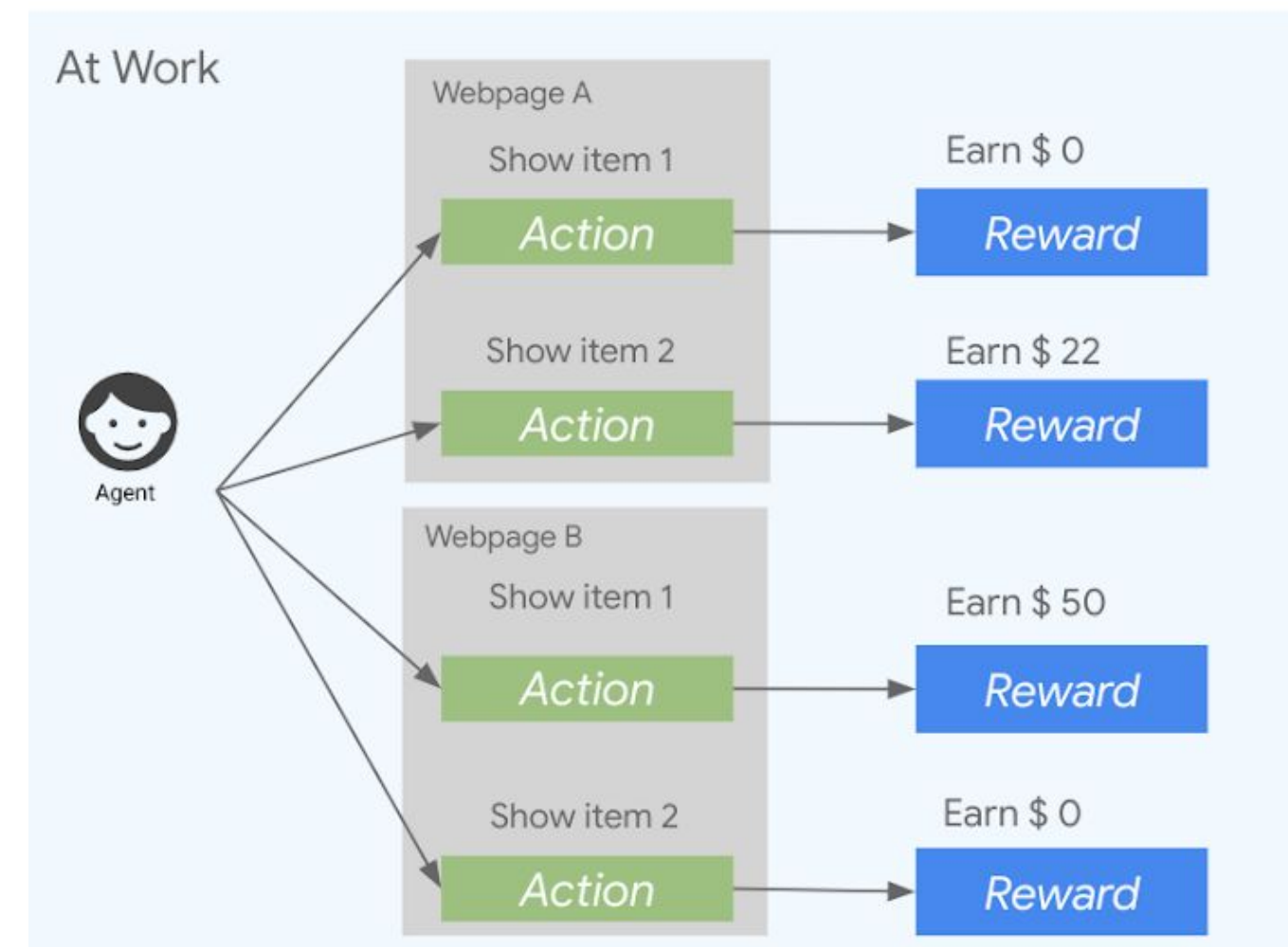
# Multi-Armed Bandits

MDP for Multi-Armed Bandits:

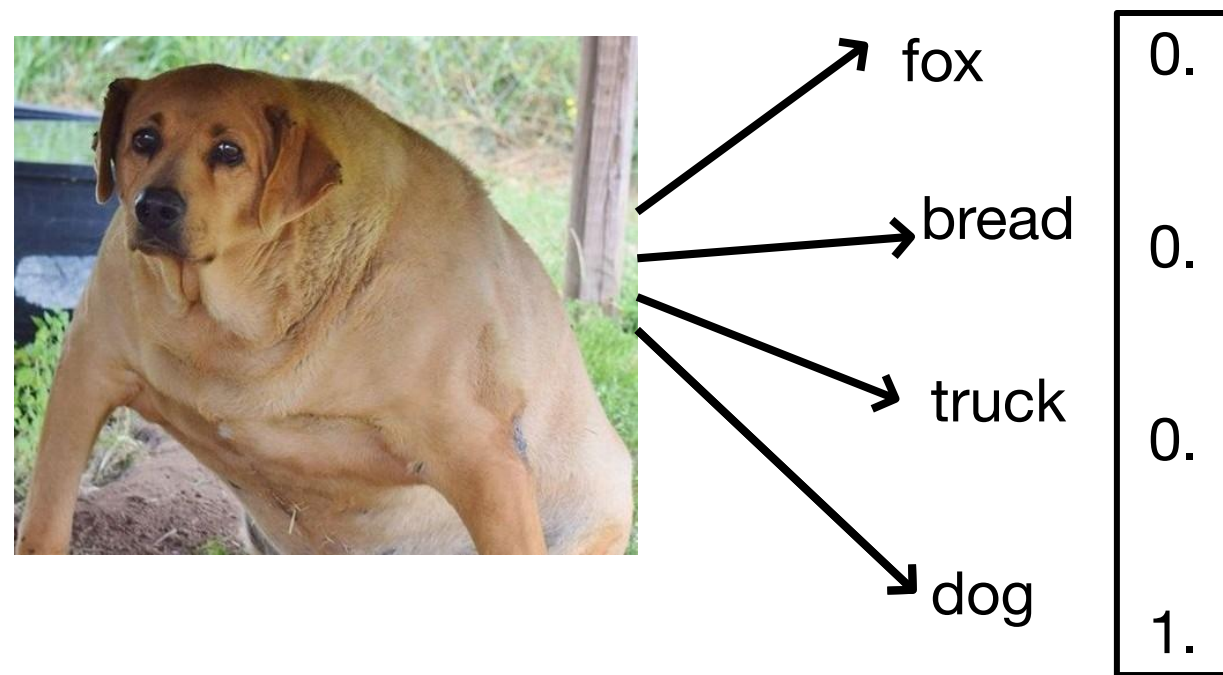1. Only one state
2. Rewards are immediate (follows from 1.)



MDP for Contextual Multi-Armed Bandits:

1. $$P(S' \mid S, A) = P(S)$$

2. Rewards are immediate (follows from 1.)

# Exploration-Exploitation Dilemma

We have ground truth labels



fox

bread

truck

dog

| 0. |
| 0. |
| 0. |
| 1. |

We have rewards



fox **?**

bread  -3

truck  **?**

dog  **?**

**?**

- Was that action optimal?
- Should you explore other actions?
- When you need to stop exploration?

Reward contains **less information** than a correct answer!

# Reinforcement Learning: SL as RL

You can formulate SL problem as RL problem!

Given Dataset: $D := \{(X_i, y_i)\}$

We consider X_i as states, and y_i as correct actions!

Then the reward function will be $R(X_i, a_i) = 1.$ if $a_i = y_i$ else 0.

Why don't we use Reinforcement learning every where?

Because Reinforcement learning is a harder problem!





Boston Dynamics

# Reward Specification Problem

Goal: Train a bot to win the game!

Rewards:

- +100 for the first place
- +5 for additional targets along the course

https://www.youtube.com/embed/tlOlHko8ySg?enablejsapi=1

**Reward is a proxy for you goal, but they are not the same!**

# Credit Assignment Problem



- Agent makes a move at step 8
- At step 50 agent loses: R = -1
- Was it a good move?

Your data is not i.i.d. Previous actions affect future states and rewards.

Credit Assignment Problem:

**How to determine which actions are responsible for the outcome?**

# Distributional shift: In case of Deep RL



The training dataset is changing with the policy.

This can lead to a **catastrophic forgetting problem**:

Agent unlearns it's policy in some parts of the State Space

# What we have discussed:

- What is RL?
- When do we need it?
- State, action, policy, reward, markovian property
- Why don't we use it everywhere?

Let's estimate policies.

# Basics

$s \sim S$; $a \sim A$ - state/action spaces (can be infinite)

$p(s_{t+1} | s_t, a_t)$ - dynamics of transitions in the environment (Markovian)

$r(s, a)$ - reward for action a in state s (can be random or depends on other variables)

$\pi(a | s)$ - agent policy

$p(\tau | \pi) = p(s_0) \prod^T_{t=0} \pi(a_t | s_t) p(s_{t+1} | a_t, s_t)$ - agent policy

where $\tau = (s_0, a_0, s_1, a_1, \ldots, s_T, a_T)$ - agent trajectory

$R_t = \sum^\infty_{\tau=t} \gamma^{\tau-t} r(s_{t+\tau}, a_{t+\tau})$
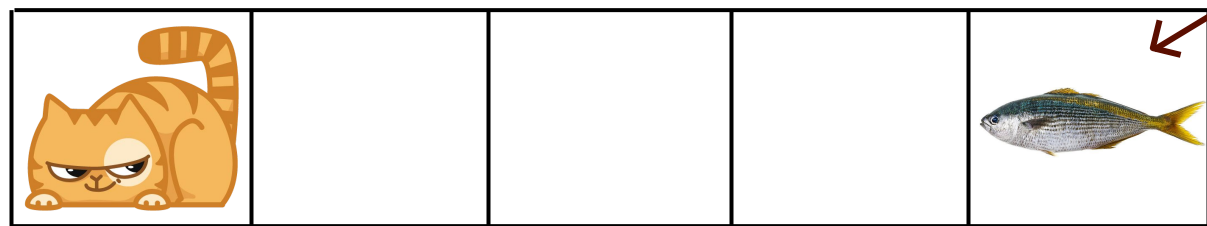
Return – random variable. Why?

$R_t = r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \gamma^2 r(s_{t+2}, a_{t+2}) + \ldots$ - reward to go **or** return

# Rate policy

How good is the policy $\pi$, if we start in state $s$?

$$V_t^\pi(s) = E_{\tau \sim \pi} [R_t | s_t = s]$$

Policy $\forall s$: $\longrightarrow$



Terminal state

+1 for fish

$V$ - value function:

| $\gamma^4$ | $\gamma^3$ | $\gamma^2$ | $\gamma$ | 1 |
|---|---|---|---|---|
| $s_0$ | $s_1$ | $s_2$ | | $s_3$ |
| $s_4$ | | | | |

# Rate policy

How good is the policy $\pi$, if we start in state $s$?

$$V_t^\pi(s) = E_{\tau \sim \pi}[R_t \mid s_t = s]$$

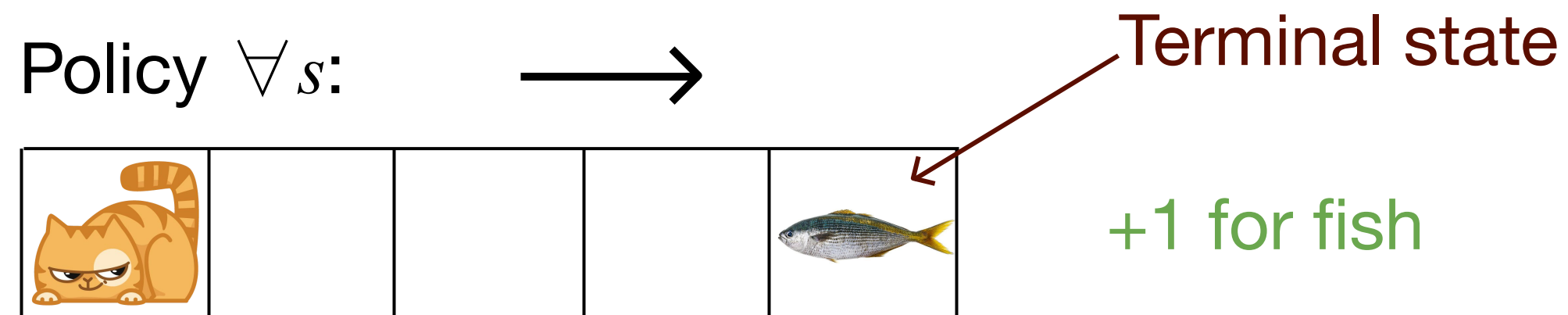Policy $\forall s$: $\longrightarrow$

Terminal state



+1 for fish

What if we "force" to choose the action $a$ in $s$, and only then follow the policy $\pi$?

$$Q_t^\pi(s, a) = E_{\tau \sim \pi}[R_t \mid s_t = s, a_t = a]$$

In complex environments, it is inconvenient to count!

$V$ - value function:

| $\gamma^4$ | $\gamma^3$ | $\gamma^2$ | $\gamma$ | 1 |
|---|---|---|---|---|
| $s_0$ | $s_1$ | $s_2$ | $s_3$ | |

$s_4$

$Q$ - value function:

$\longrightarrow$
$\longleftarrow$

| $\gamma^4$ | $\gamma^3$ | $\gamma^2$ | $\gamma$ | 1 |
|---|---|---|---|---|
| $\gamma^5$ | $\gamma^5$ | $\gamma^4$ | $\gamma^3$ | 1 |
| $s_0$ | $s_1$ | $s_2$ | $s_3$ | $s_4$ |

# Finite and infinite

in time MDP

Let the length of the episode $T = \infty$, then it is easy to see that:

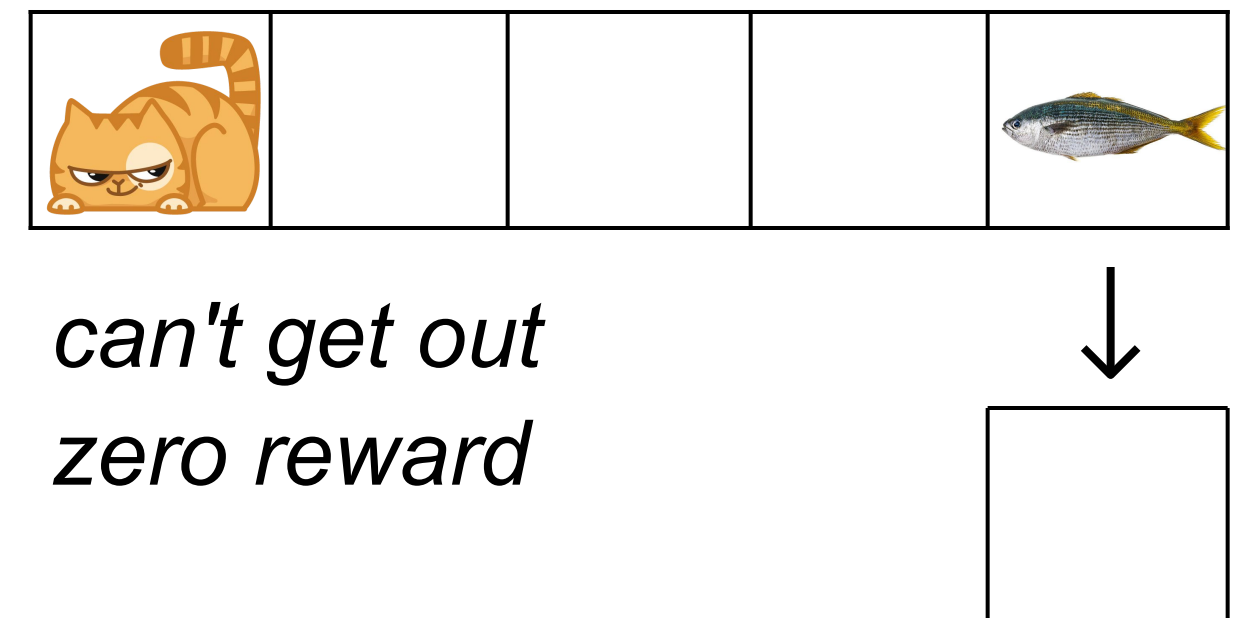$$V_t^\pi(x) = \mathbf{E}_{\tau \sim \pi}\left[\sum_{i=t}^{\infty} \gamma^{i-t} r_i | s_t = s\right] = \left[\sum_{i=0}^{\infty} \gamma^i r_i | s_0 = s\right]$$

That mean,

$$V^\pi(s) = V_0^\pi(s) = V_t^\pi \qquad \text{does not depend on time!}$$

Such MDPs are called **infinite.**

MDPs with **terminal** states are reduced to infinite by adding an **absorbing** state.
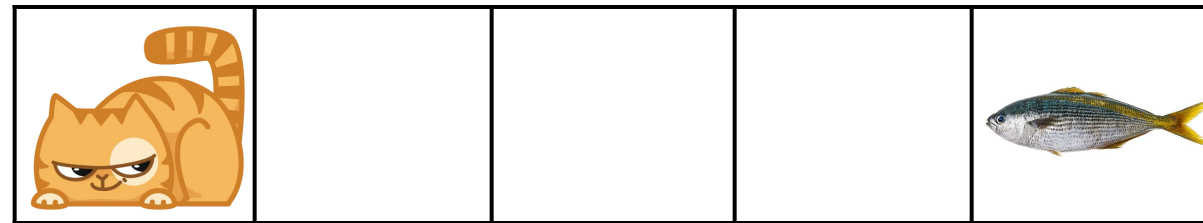


*can't get out*
*zero reward*

# Finite and infinite

in time MDP

If the episode length $T < \infty$, then in general case V depends on time.

For example, if $T = 4$:



$$V_0^\pi(s_0) = \gamma^4$$

$$V_1^\pi(s_0) = 0$$

In theory, however, t is omitted here as well.

To do this, it suffices to assume that the state contains t:

$$s \to (s, t)$$

Such MDP are called **finite**.

# Dynamic programming

Reformulation of a complex problem as a recursive sequence of simpler problems.

Get the recursive ratio for the cumulative reward $R_t$ :

$$R_t = r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \gamma^2 r(s_{t+2}, a_{t+2}) + \ldots = r(s_t, a_t) + \gamma\, (r(s_{t+1}, a_{t+1}) + \gamma r(s_{t+2}, a_{t+2}) + \ldots ) = r(s_t, a_t) + \gamma R_{t+1}$$
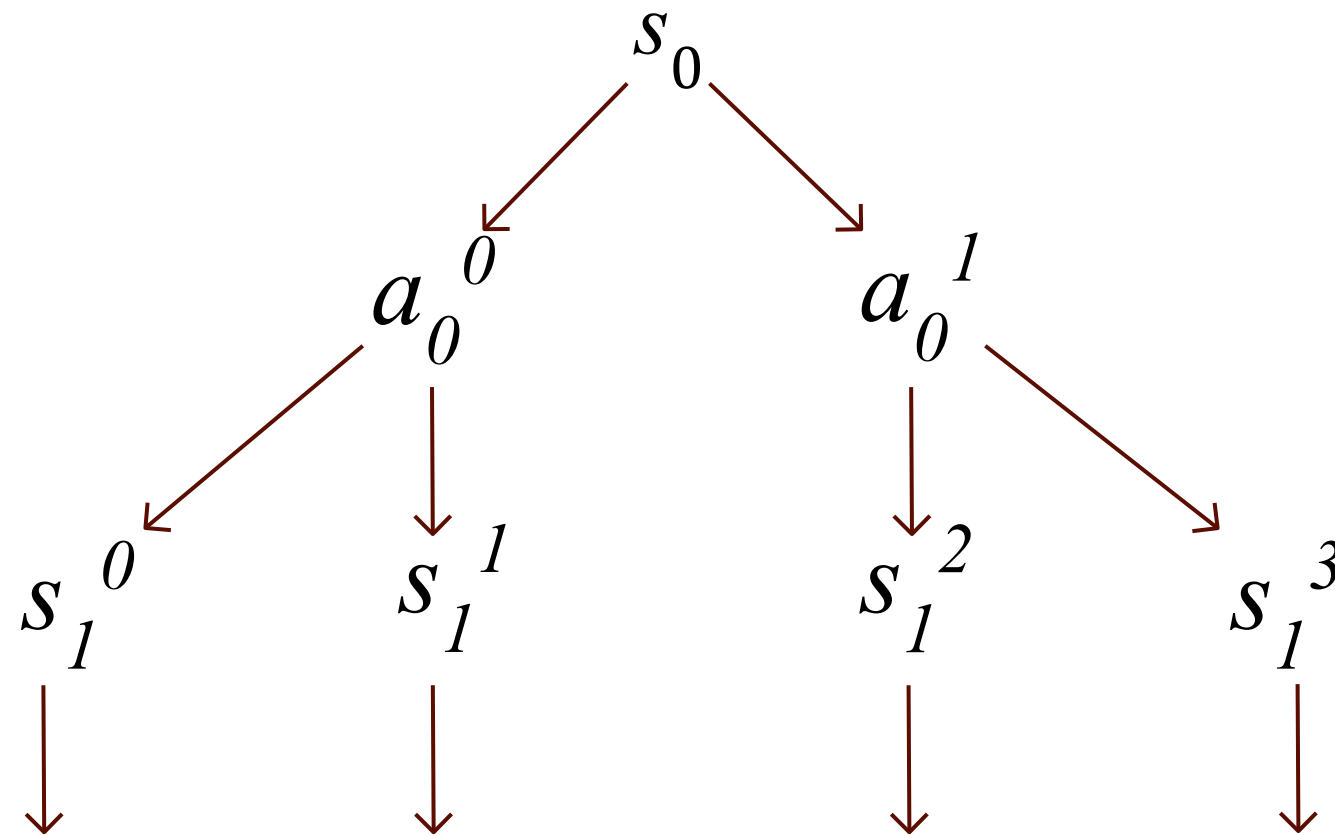
<u>For $V$ - function:</u>

$$V^\pi(s) = \mathbf{E}\,[R_t | s_t = s] = \mathbf{E}\,[r(s_t, a_t) + \gamma R_{t+1} | s_t = s] = \mathbf{E}_{a \sim \pi(\cdot | s)}\,[r(s, a) + \gamma \mathbf{E}_{s' \sim p(s' | s, a)} \mathbf{E}[R_{t+1} | s_{t+1} = s']] =$$
$$\mathbf{E}_{a \sim \pi(\cdot | s)}\,[r(s, a) + \gamma \mathbf{E}_{s' \sim p(s' | s, a)} V^\pi(s')]$$

<u>For Q - function:</u>

$$Q^\pi(s, a) = r(s, a) + \gamma\, \mathbf{E}_{s' \sim p(\cdot | s,a)} \mathbf{E}_{a' \sim \pi(\cdot | s')} Q^\pi(s', a')$$

# Dynamic programming

If states never repeat in the environment, the graph of this MDP will be a tree

$$V^{\pi}(s_T) = \mathsf{E}a{\sim}\pi(\cdot \mid s_T) \, r(s_T, a)$$



$$V^{\pi}(s) = \mathsf{E}_{a{\sim}\pi(\cdot \mid s)} \left[ r(s, a) + \gamma \mathsf{E}_{s'{\sim}p(s' \mid s,a)} V^{\pi}(s') \right]$$

Bellman's Equations tell you how to calculate value "backwards".

# Relationship of Q and V functions

Expressing V in terms of Q:

$$V^\pi(s) = \mathsf{E}_{a \sim \pi(\cdot \,|\, s)} \, Q^\pi(s, a)$$

V - this is Q, in which the action from the policy was substituted

Expressing Q in terms of V:

$$Q^\pi(s, a) = r(s, a) + \mathsf{E}_{s' \sim p(\cdot \,|\, s, a)} \, V^\pi(s')$$

Q - is the instant reward for (s, a) plus future state value

# How to solve the Bellman equation?

Like SLAE:

$$V^{\pi}(s) = \mathsf{E}_{a\sim\pi(\cdot\,|s)}\left[r(s, a) + \gamma\mathsf{E}_{s'\sim p(s'\,|s,a)}V^{\pi}(s')\right] =$$

$$= \mathsf{E}_{a\sim\pi(\cdot\,|s)}r(s, a) + \gamma\mathsf{E}_{s'\sim p(s'\,|s)}V^{\pi}(s') = u(s) + \gamma\mathsf{E}_{s'\sim p(s'\,|s)}V^{\pi}(s')$$
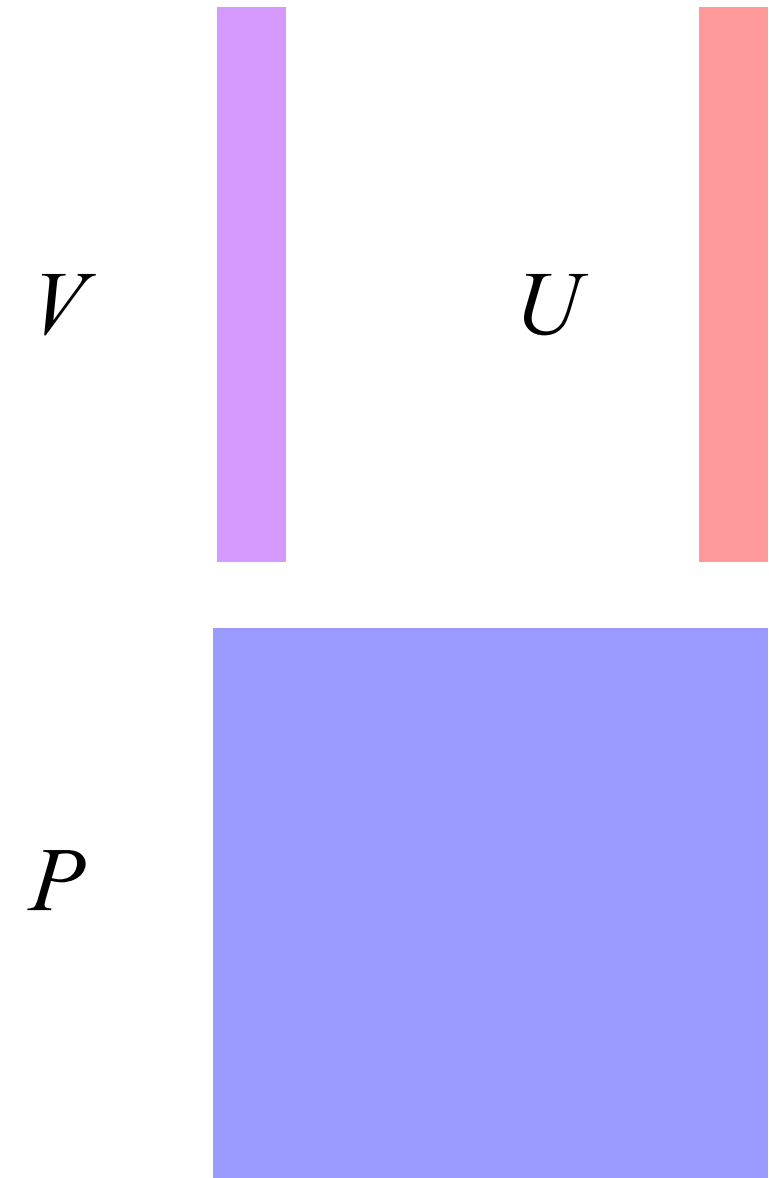
Everything is linear with respect to V

$$V = U + \gamma PV$$

$$(I - \gamma P)V = U$$

$$V = (I - \gamma P)^{-1}U$$

It will be expensive!
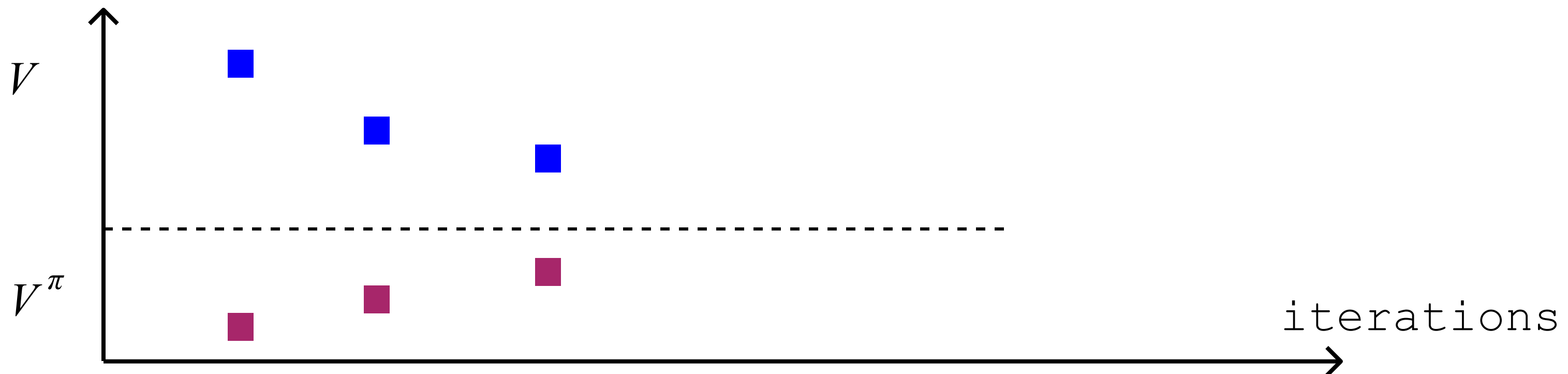
Without taking into account $|A|$ - already $O(|S|^3)$

$V$

$U$

$P$

# How to solve the Bellman equation?

Simple iteration method:

$$V_{new} = F(V_{old})$$

$$F(V_s) = \mathsf{E}_{a \sim \pi(\cdot \,|\, s)} \left[ r(s, a) + \gamma \mathsf{E}_{s' \sim p(s' \,|\, s, a)} \, V_s \right]$$

Will the algorithm converge? *It will be if the mapping F is contractive*
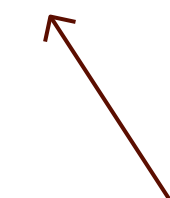
# Is the Bellman operator a contraction?

By the infinite norm:

reminder: $||x||_\infty = max_i |x_i|$

$$||F(V) - F(W)||_\infty = ||U + \gamma PV - U - \gamma PW||_\infty =$$

$$= ||\gamma P(V - W)||_\infty \leq \gamma ||P||_\infty ||V - W||_\infty$$

where is the matrix norm:

$$||P||_\infty = \max_{x: ||x||_\infty = 1} ||Px||_\infty = \max_{x: ||x||_\infty = 1} \max_i |\sum_j P_{ij} x_j|$$

$$= \max_i |\sum_j P_{ij}| = 1$$

$$x_j = \text{sign}(P_{ij})$$

Q.E.D.: $||F(V) - F(W)||_\infty \leq \gamma ||V - W||_\infty$

# Algorithm Policy Evaluation

- Initialize $V(s) \ \forall s$

- Repeat:

  - $\Delta = 0$

  - For all $s$:

    - $v = V(s)$
    - $V(s) = \mathsf{E}_{a \sim \pi(\cdot \mid s)} \left[ r(s, a) + \gamma \mathsf{E}_{s' \sim p(\cdot \mid s, a)} V(s') \right]$
    - $\Delta = \max \left( \Delta, \left| v - V(s) \right| \right)$

- while $\Delta > \epsilon$

# Policy improvement

# Optimal Bellman Equations
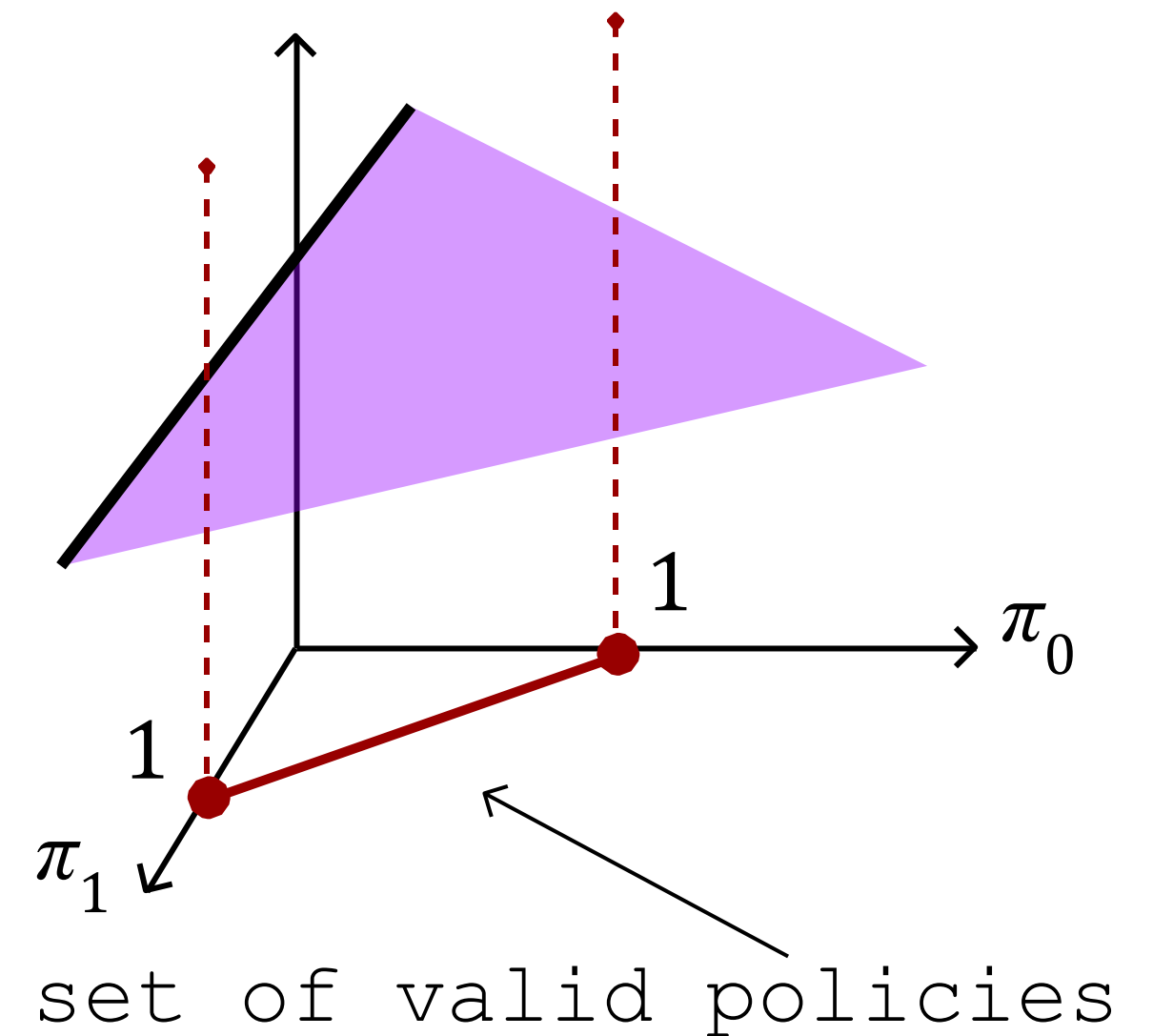
V - function for optimal policy:

$$V^*(s) = \max_{\pi^0} (\mathsf{E}_a [r(s, a) + \gamma \mathsf{E}_{s'} \max_{\pi^1,\dots} V^{\pi_1,\dots}(s')])$$

With respect to $\pi_0$ the following problem is solved:

$$\begin{cases} \Sigma_i \pi_i y_i \rightarrow \max \\ \pi_i \geq 0 \\ \Sigma_i \pi_i = 1 \end{cases}$$

LP task

Optimal Bellman equations:

$$V^*(s) = \max_a [r(s, a) + \gamma \mathsf{E}_{s'} V^*(s')]$$



`set of valid policies`

Among the optimal policies there is always a deterministic (greedy)

# Optimality Bellman Equations

V - function for optimal policy:

$$V^*(s) = \max_\pi V^\pi(s) = \max_\pi (\mathsf{E}_a [r(s, a) + \gamma \mathsf{E}_{s'} V^\pi(s')]) = \max_{\pi^0, \pi^1, \dots} (\mathsf{E}_a [r(s, a) + \gamma \mathsf{E}_{s'} V^{\pi_1, \dots}(s')]) =$$

$$= \max_{\pi^0} (\mathsf{E}_a [r(s, a) + \gamma \mathsf{E}_{s'} \max_{\pi^1, \dots} V^{\pi_1, \dots}(s')])$$

With respect to $\pi_0$ the following problem is solved:

$$
\begin{cases}
\Sigma_i \pi_i\, y_i \;\to \max \\
\pi_i \geq 0 \\
\Sigma_i \pi_i = 1
\end{cases}
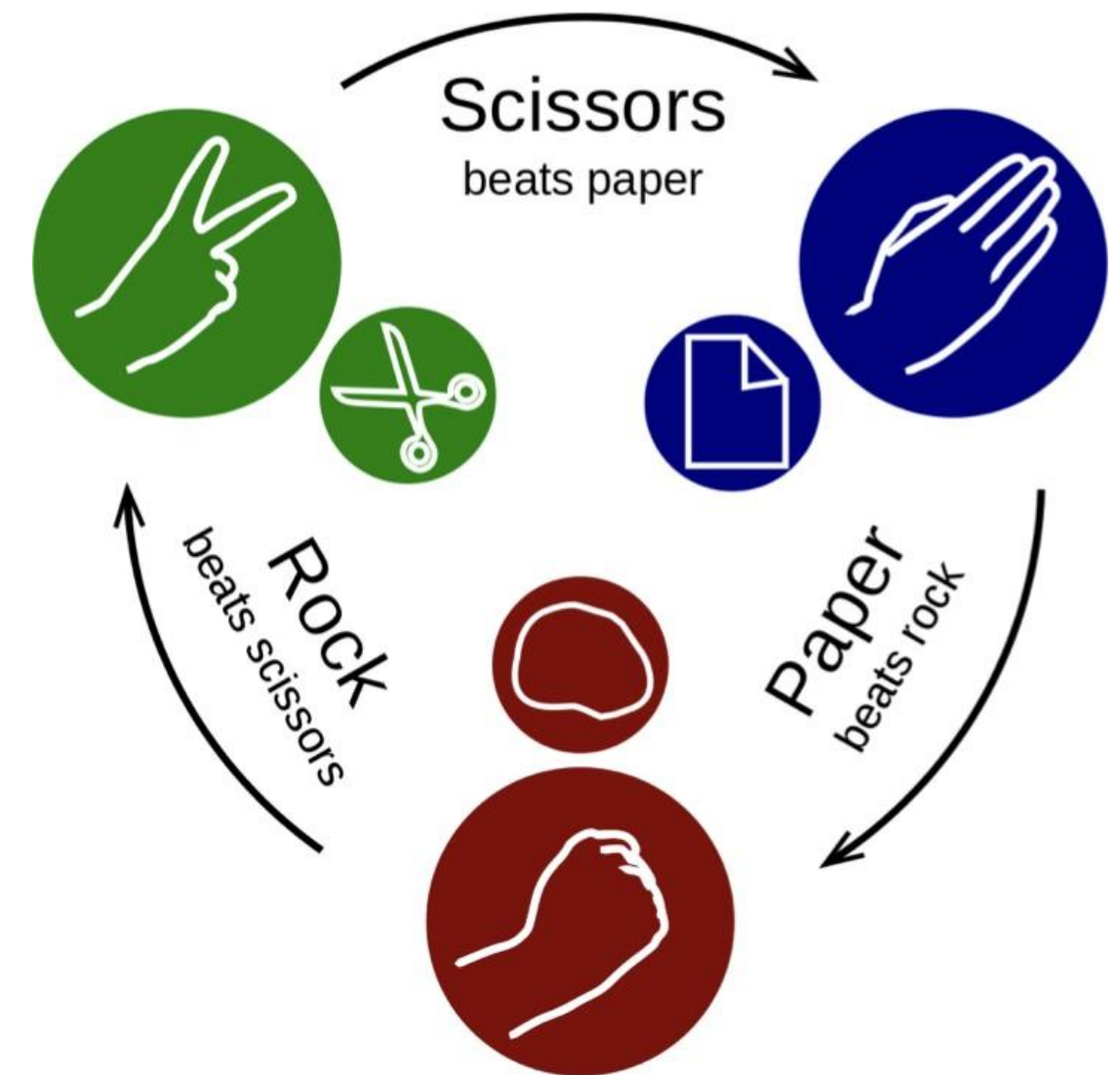$$

# Is this true?

Rock - Paper - Scissors

What is the optimal strategy?

- there are no states - it is unique
- cenvironment == random reward



- If your opponent has a fixed policy, there is an optimal deterministic one for you!
- If the opponent adjusts, then it is not MDP.

Reward in this environment:

$p(r\,|\,history)$ **or** $p(r\,|\,\pi)$

# Optimal Bellman Equations

Optimal Bellman equation for V-function:

$$V^*(s) = \max_a [r(s, a) + \gamma \mathsf{E}_{s'} V^*(s')])$$

Optimal Bellman equation for Q-function:

$$Q^*(s, a) = r(s, a) + \gamma \mathsf{E}_{s'} \max_{a'} Q^*(s', a')$$

Expressing $V^*$ in terms of $Q^*$:

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

Expressing $Q^*$ in terms of $V^*$:

$$Q^*(s, a) = r(s, a) + \gamma \mathsf{E}_{s'} V^*(s')$$

# Policy improvement

**Def.**    $\pi' \geq \pi$ if
$V^{\pi'}(s) \geq V^{\pi}(s) \quad \forall s$

Our Policy Update Strategy:

- let $\exists s$ be such that:

  $\exists a : Q^{\pi}(s, a) > V^{\pi}(s)$

- then $\pi'(s) := a$,

  In all $s^{\sim} \mathrel{!=} s$ define $\pi'(s^{\sim}) = \pi(s^{\sim})$

In this case, $\pi' \geq \pi$     **CHECK IT**

# Policy improvement

Our Policy Update Strategy:

- пусть $\exists s$ такой, что:

  $$\exists a : Q^\pi(s, a) > V^\pi(s)$$

- then $\pi'(s) := a,$

  In all $\tilde{s} \mathrel{!}= s$ define $\pi'(\tilde{s}) = \pi(\tilde{s})$

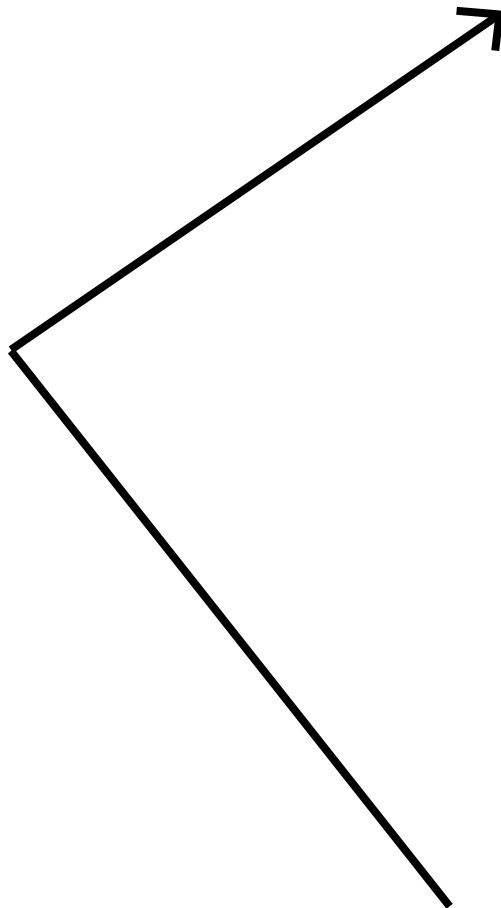$$V^\pi(s) \leq Q^\pi(s, \pi'(s))$$

$$= r(s, \pi'(s)) + \mathsf{E}_{s'} V^\pi(s') \leq$$

$$\leq r(s, \pi'(s)) + \mathsf{E}_{s'} Q^\pi(s', \pi'(s')) \leq$$

$$\leq \cdots \leq V^{\pi'}(s)$$

Q.E.D.

# Algorithm Policy Iteration

- **Initialize** $V(s), \pi(s) \quad \forall s$
- estimate $V$ for policy $\pi$ by method PE
- stop = *True*
- **For all** $s$:
  - $a = \pi(s)$
  - $\pi(s) = \arg\max_a [r(s, a) + \mathsf{E}_{s'} V(s')]$
  - **if** $a =/= \pi(s)$:
    - stop = *False*
- **if not** stop

43

# Algorithm Value Iteration

- **Initialize** $V(s) \ \forall \ s$
- **Repeat**:
  - $\Delta = 0$
  - **For all** $s$:
    - $v = V(s)$
    - $V(s) = \max_a [r(s, a) + \gamma \mathrm{E}_{s' \sim p(\cdot \,|\, s, a)} V(s')]$
    - $\Delta = \max (\Delta, |v - V(s)|)$
- **while** $\Delta > \epsilon$

# Recap

- What is RL?
- When do we need it?
- State, action, policy, reward, markovian property, MDP
- Why don't we use it everywhere?
- V-function
- Q-function
- Value Iteration