#### Reinforcement learning

Tomáš Svoboda, Petr Pošík

Vision for Robots and Autonomous Systems, Center for Machine Perception
Department of Cybernetics
Faculty of Electrical Engineering, Czech Technical University in Prague

April 15, 2025

### (Multi-armed) Bandits







Think about not one but 10 arms you may choose to pull.

### (Multi-armed) Bandits



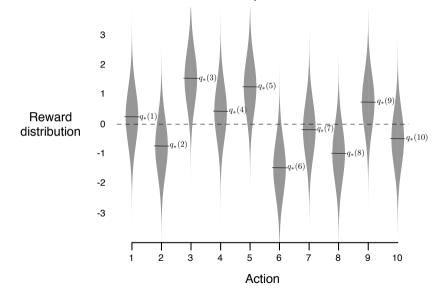




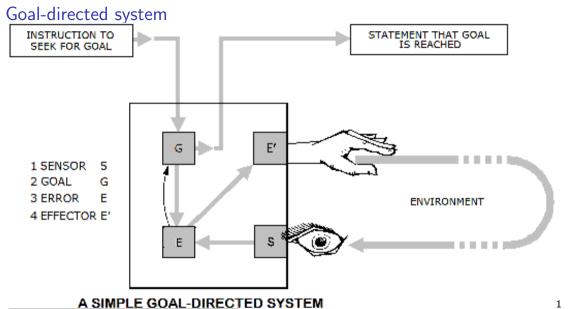
Think about not one but 10 arms you may choose to pull.

p(s'|s, a) and r(s, a, s') not known!

#### 10 armed bandit, what arm to pull?

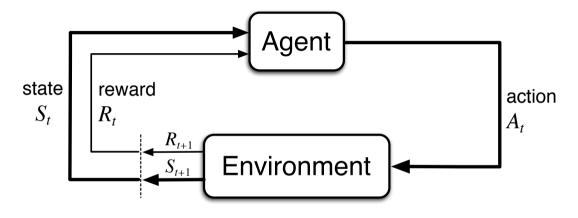


See chapters 2.2 and 2.3 in [4] for more detailed discussion



<sup>1</sup>Figure from http://www.cybsoc.org/gcyb.htm

## Reinforcement Learning - performing actions, learning from rewards



- Feedback in form of Rewards
- Learn to act so as to maximize expected rewards.

<sup>2</sup>Scheme from [4]

/ 39

### Reinforcement Learning - performing actions, learning from rewards



### Reinforcement Learning - performing actions, learning from rewards





#### Examples, robot learning, Atari games, ...

#### **Autonomous Flipper Control with Safety Constraints**

Martin Pecka, Vojtěch Šalanský, Karel Zimmermann, Tomáš Svoboda

experiments utilizing
Constrained Relative Entropy Policy Search

Video: Learning safe policies<sup>3</sup>

<sup>&</sup>lt;sup>3</sup>M. Pecka, V. Salansky, K. Zimmermann, T. Svoboda. Autonomous flipper control with safety constraints. In Intelligent Robots and Systems (IROS), 2016, https://youtu.be/\_oUMbBtoRcs

## From off-line (MDPs) to on-line (RL)

Markov decision process – MDPs. Off-line search, we know:

- ▶ A set of states  $s \in \mathcal{S}$  (map)
- ▶ A set of actions per state.  $a \in A$
- A transition model T(s, a, s') or p(s'|s, a) (robot)
- ▶ A reward function r(s, a, s') (map, robot)

Looking for the optimal policy  $\pi(s)$ . We can plan/search before the robot enters the environment.

## From off-line (MDPs) to on-line (RL)

Markov decision process – MDPs. Off-line search, we know:

- ▶ A set of states  $s \in S$  (map)
- ▶ A set of actions per state.  $a \in A$
- A transition model T(s, a, s') or p(s'|s, a) (robot)
- A reward function r(s, a, s') (map, robot)

Looking for the optimal policy  $\pi(s)$ . We can plan/search before the robot enters the environment.

#### On-line problem:

- ightharpoonup Transition model p and reward function r not known.
- Agent/robot must act and learn from experience.

## (Transition) Model-based learning

The main idea: Do something, and:

- Learn an approximate model from experiences.
- ► Solve as if the model was correct.

### (Transition) Model-based learning

#### The main idea: Do something, and:

- Learn an approximate model from experiences.
- Solve as if the model was correct.

#### Learning MDP model:

- ln s try a, observe s', count (s, a, s').
- Normalize to get and estimate of  $p(s' \mid s, a)$ .
- ▶ Discover (by observation) each r(s, a, s') when experienced.

### (Transition) Model-based learning

The main idea: Do something, and:

- ► Learn an approximate model from experiences.
- Solve as if the model was correct.

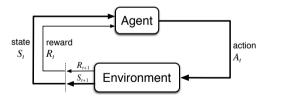
#### Learning MDP model:

- In s try a, observe s', count (s, a, s').
- Normalize to get and estimate of  $p(s' \mid s, a)$ .
- ▶ Discover (by observation) each r(s, a, s') when experienced.

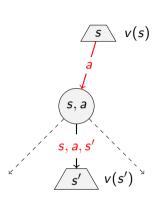
Solve the learned MDP.

## Reward function r(s, a, s')

- ightharpoonup r(s, a, s') reward for taking a in s and landing in s'.
- ▶ In Grid world, we assumed r(s, a, s') to be the same everywhere.
- ▶ In the real world, it is different (going up, down, ...)



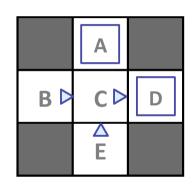
In ai-gym env.step(action) returns s', r(s, action, s').



## Model-based learning: Grid example

Input Policy  $\pi$ 

**Observed Episodes (Training)** 



Assume:  $\gamma = 1$ 

## Episode 1

B, east, C, -1 C, east, D, -1 D, exit, x, +10

## Episode 2

B, east, C, -1 C, east, D, -1 D, exit, x, +10

# Episode 3

E, north, C, -1 C, east, D, -1 D, exit, x, +10

## Episode 4

E, north, C, -1 C, east, A, -1 A, exit, x, -10

1

#### Learning transition model

$$\hat{p}(D \mid C, east) = ?$$



### Episode 1

B, east, C, -1 C, east, D, -1 D, exit, x, +10

### Episode 2

B, east, C, -1 C, east, D, -1 D, exit, x, +10

## Episode 3

E, north, C, -1 C, east, D, -1 D, exit, x, +1

### Episode 4

E, north, C, -1 C, east, A, -1 A, exit, x, -10

### Learning reward function

$$\hat{r}(C, east, D) = ?$$



B, east, C, -1 C, east, D, -1 D, exit, x, +10

## Episode 2

B, east, C, -1 C, east, D, -1 D, exit, x, +10

## Episode 3

E, north, C, -1 C, east, D, -1 D, exit, x, +10

### Episode 4

E, north, C, -1 C, east, A, -1 A, exit, x, -10

## Model based vs model-free: Expected age E [A]

Random variable age A.

$$\mathsf{E}\left[A\right] = \sum_{a} P(A = a)a$$

We do not know P(A = a). Instead, we collect N samples  $[a_1, a_2, \dots a_N]$ .

## Model based vs model-free: Expected age E [A]

Random variable age A.

$$\mathsf{E}\left[A\right] = \sum_{a} P(A=a)a$$

We do not know P(A = a). Instead, we collect N samples  $[a_1, a_2, \dots a_N]$ .

#### Model based

$$\hat{P}(a) = \frac{\mathsf{num}(a)}{N}$$

$$E[A] \approx \sum_{a} \hat{P}(a)a$$

## Model based vs model-free: Expected age E[A]

Random variable age A.

$$\mathsf{E}\left[A\right] = \sum_{a} P(A=a)a$$

We do not know P(A = a). Instead, we collect N samples  $[a_1, a_2, \dots a_N]$ .

Model based

Model free

$$\hat{P}(a) = \frac{\mathsf{num}(a)}{N}$$

$$P(a) = \frac{1}{N}$$

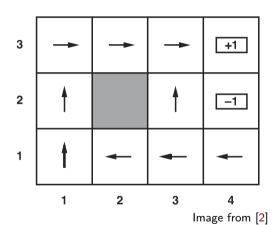
$$\mathsf{E}\left[A\right] pprox \sum_{a} \hat{P}(a)a$$

$$\mathsf{E}\left[A\right] \approx \frac{1}{N} \sum_{i} a_{i}$$

# Model-free learning

## Passive learning (evaluating given policy)

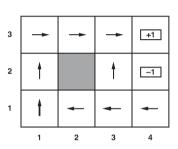
- ▶ **Input:** a fixed policy  $\pi(s)$
- ▶ We want to know how good it is.
- ightharpoonup r, p not known.
- ► Execute policy . . .
- and learn on the way.
- ▶ **Goal:** learn the state values  $v^{\pi}(s)$



### Direct evaluation from episodes

Value of s for  $\pi$  – expected sum of discounted rewards – expected return

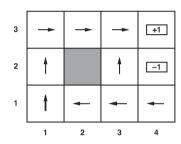
$$v^{\pi}(S_t) = \mathsf{E}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}
ight]$$
 $v^{\pi}(S_t) = \mathsf{E}\left[G_t\right]$ 



### Direct evaluation from episodes

Value of s for  $\pi$  – expected sum of discounted rewards – expected return

$$v^{\pi}(S_t) = \mathsf{E}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}
ight]$$
  $v^{\pi}(S_t) = \mathsf{E}\left[G_t
ight]$ 



$$(1,1)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (2,3)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (2,3)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (3,2)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \rightsquigarrow (2,1)_{\textbf{-.04}} \rightsquigarrow (3,1)_{\textbf{-.04}} \rightsquigarrow (3,2)_{\textbf{-.04}} \rightsquigarrow (4,2)_{\textbf{-1}} .$$

## Direct evaluation from episodes, $v^{\pi}(S_t) = E[G_t], \gamma = 1$

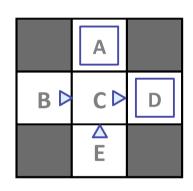
$$\begin{array}{l} (1,1)_{\textbf{-.04}} \leadsto (1,2)_{\textbf{-.04}} \leadsto (1,3)_{\textbf{-.04}} \leadsto (1,2)_{\textbf{-.04}} \leadsto (1,3)_{\textbf{-.04}} \leadsto (2,3)_{\textbf{-.04}} \leadsto (3,3)_{\textbf{-.04}} \leadsto (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \leadsto (1,2)_{\textbf{-.04}} \leadsto (1,3)_{\textbf{-.04}} \leadsto (2,3)_{\textbf{-.04}} \leadsto (3,3)_{\textbf{-.04}} \leadsto (3,2)_{\textbf{-.04}} \leadsto (3,3)_{\textbf{-.04}} \leadsto (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \leadsto (2,1)_{\textbf{-.04}} \leadsto (3,1)_{\textbf{-.04}} \leadsto (3,2)_{\textbf{-.04}} \leadsto (4,2)_{\textbf{-1}} \end{array}.$$

What is v(3,2) after these episodes?

### Direct evaluation: Grid example

#### Input Policy $\pi$





Assume:  $\gamma = 1$ 

## Episode 1

B, east, C, -1 C, east, D, -1 D, exit, x, +10

## Episode 2

B, east, C, -1 C, east, D, -1 D, exit, x, +10

## Episode 3

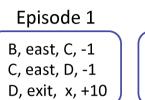
E, north, C, -1 C, east, D, -1 D, exit, x, +10

## Episode 4

E, north, C, -1 C, east, A, -1 A, exit, x, -10

### Direct evaluation: Grid example, $\gamma = 1$

What is v(C) after the 4 episodes?



# Episode 2

B, east, C, -1 C, east, D, -1 D, exit, x, +10

## Episode 3

E, north, C, -1 C, east, D, -1 D, exit, x, +10 E, north, C, -1 C, east, A, -1 A, exit, x, -10

### Episode 4

### Direct evaluation: Grid example, $\gamma=1$

What is v(C) after the 4 episodes?

What is the relation of M and N?

Let M be the number of recorded episodes. Let N be the number of samples used to compute the averages.

- A N = M
- $\mathbf{B} \ N \leq M$
- $C N \geq M$
- D N has no relation to M

#### Episode 1

B, east, C, -1 C, east, D, -1 D, exit, x, +10

### Episode 2

B, east, C, -1 C, east, D, -1 D, exit, x, +10

#### Episode 3

E, north, C, -1 C, east, D, -1 D, exit, x, +10

### Episode 4

E, north, C, -1 C, east, A, -1 A, exit, x, -10

## Direct evaluation algorithm (every-visit version)

$$\begin{array}{l} (1,1)_{\textbf{-.04}} \leadsto (1,2)_{\textbf{-.04}} \leadsto (1,3)_{\textbf{-.04}} \leadsto (1,2)_{\textbf{-.04}} \leadsto (1,3)_{\textbf{-.04}} \leadsto (2,3)_{\textbf{-.04}} \leadsto (3,3)_{\textbf{-.04}} \leadsto (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \leadsto (1,2)_{\textbf{-.04}} \leadsto (1,3)_{\textbf{-.04}} \leadsto (2,3)_{\textbf{-.04}} \leadsto (3,3)_{\textbf{-.04}} \leadsto (3,2)_{\textbf{-.04}} \leadsto (3,3)_{\textbf{-.04}} \leadsto (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \leadsto (2,1)_{\textbf{-.04}} \leadsto (3,1)_{\textbf{-.04}} \leadsto (3,2)_{\textbf{-.04}} \leadsto (4,2)_{\textbf{-1}} \end{array}.$$

Input: a policy  $\pi$  to be evaluated Initialize:

$$V(s) \in \mathbb{R}$$
, arbitrarily, for all  $s \in \mathcal{S}$ 

$$Returns(s) \leftarrow$$
 an empty list, for all  $s \in \mathcal{S}$ 

Loop forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 

$$G \leftarrow 0$$

Loop backwards for each step of episode,  $t = T - 1, T - 2, \dots, 0$ :

$$G \leftarrow R_{t+1} + \gamma G$$

Append G to  $Returns(S_t)$ 

$$V(S_t) \leftarrow average(Returns(S_t))$$

# Direct evaluation algorithm (first-visit version)

```
(1,1)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (2,3)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (2,3)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (3,2)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \rightsquigarrow (2,1)_{\textbf{-.04}} \rightsquigarrow (3,1)_{\textbf{-.04}} \rightsquigarrow (3,2)_{\textbf{-.04}} \rightsquigarrow (4,2)_{\textbf{-1}} .
```

Input: a policy  $\pi$  to be evaluated Initialize:

$$V(s) \in \mathbb{R}$$
, arbitrarily, for all  $s \in \mathcal{S}$ 

 $Returns(s) \leftarrow ext{an empty list, for all } s \in \mathcal{S}$ 

Loop forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ 

$$G \leftarrow 0$$

Loop backwards for each step of episode,  $t = T - 1, T - 2, \dots, 0$ :

$$G \leftarrow R_{t+1} + \gamma G$$

If  $S_t$  does not appear in  $S_0, S_1, \ldots, S_{t-1}$ : // Use the return for the first visit only

Append G to  $Returns(S_t)$ 

 $V(S_t) \leftarrow \operatorname{average}(Returns(S_t))$ 

#### The good:

- ► Simple, easy to understand and implement.
- ▶ Does not need p, r and eventually it computes the true  $v^{\pi}$ .

#### The good:

- Simple, easy to understand and implement.
- ▶ Does not need p, r and eventually it computes the true  $v^{\pi}$ .

#### The bad:

$$(1,1)_{\textbf{-.04}} \rightarrow (1,2)_{\textbf{-.04}} \rightarrow (1,3)_{\textbf{-.04}} \rightarrow (1,2)_{\textbf{-.04}} \rightarrow (1,3)_{\textbf{-.04}} \rightarrow (2,3)_{\textbf{-.04}} \rightarrow (3,3)_{\textbf{-.04}} \rightarrow (4,3)_{\textbf{+1}}$$

$$(1,1)_{\textbf{-.04}} \rightarrow (1,2)_{\textbf{-.04}} \rightarrow (1,3)_{\textbf{-.04}} \rightarrow (2,3)_{\textbf{-.04}} \rightarrow (3,3)_{\textbf{-.04}} \rightarrow (3,2)_{\textbf{-.04}} \rightarrow (3,3)_{\textbf{-.04}} \rightarrow (4,3)_{\textbf{+1}}$$

$$(1,1)_{\textbf{-.04}} \rightarrow (2,1)_{\textbf{-.04}} \rightarrow (3,1)_{\textbf{-.04}} \rightarrow (3,2)_{\textbf{-.04}} \rightarrow (4,2)_{\textbf{-1}}.$$

#### The good:

- Simple, easy to understand and implement.
- ▶ Does not need p, r and eventually it computes the true  $v^{\pi}$ .

#### The bad:

$$(1,1)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (2,3)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (2,3)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (3,2)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \rightsquigarrow (2,1)_{\textbf{-.04}} \rightsquigarrow (3,1)_{\textbf{-.04}} \rightsquigarrow (3,2)_{\textbf{-.04}} \rightsquigarrow (4,2)_{\textbf{-1}}.$$

Each state value learned in isolation.

#### The good:

- Simple, easy to understand and implement.
- ▶ Does not need p, r and eventually it computes the true  $v^{\pi}$ .

#### The bad:

$$(1,1)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (2,3)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \rightsquigarrow (1,2)_{\textbf{-.04}} \rightsquigarrow (1,3)_{\textbf{-.04}} \rightsquigarrow (2,3)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (3,2)_{\textbf{-.04}} \rightsquigarrow (3,3)_{\textbf{-.04}} \rightsquigarrow (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \rightsquigarrow (2,1)_{\textbf{-.04}} \rightsquigarrow (3,1)_{\textbf{-.04}} \rightsquigarrow (3,2)_{\textbf{-.04}} \rightsquigarrow (4,2)_{\textbf{-1}}.$$

- Each state value learned in isolation.
- State values are not independent
- $ightharpoonup v^{\pi}(s) = \sum_{s'} p(s' \mid s, \pi(s)) [r(s, \pi(s), s') + \gamma v^{\pi}(s')]$

### (on-line) Policy evaluation?

#### In MDP, we did:

- lnitialize the values:  $V_0^{\pi}(s) = 0$
- ▶ In each iteration, replace V with a one-step-look-ahead:

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) [r(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

# (on-line) Policy evaluation?

In MDP, we did:

- lnitialize the values:  $V_0^{\pi}(s) = 0$
- ▶ In each iteration, replace *V* with a one-step-look-ahead:

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) \left[ r(s, \pi(s), s') + \gamma V_k^{\pi}(s') \right]$$

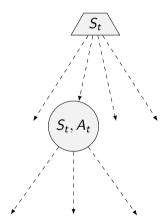
Problem: both  $p(s' | s, \pi(s))$  and  $r(s, \pi(s), s')$  unknown!

MDP (p, r known): Update V estimate by a weighted average:

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) [r(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

MDP (p, r known): Update V estimate by a weighted average:

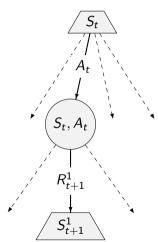
$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) [r(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$



MDP (p, r known): Update V estimate by a weighted average:

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) [r(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

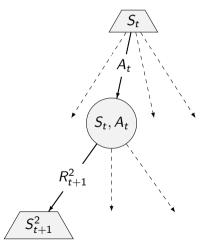
$$\mathsf{sample}^1 = R^1_{t+1} + \gamma V(S^1_{t+1})$$



MDP (p, r known): Update V estimate by a weighted average:

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) [r(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

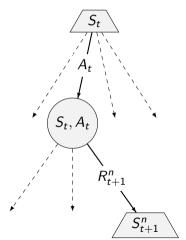
sample<sup>1</sup> = 
$$R_{t+1}^1 + \gamma V(S_{t+1}^1)$$
  
sample<sup>2</sup> =  $R_{t+1}^2 + \gamma V(S_{t+1}^2)$ 



MDP (p, r known): Update V estimate by a weighted average:

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) [r(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

$$\begin{array}{lll} \mathsf{sample}^1 & = & R_{t+1}^1 + \gamma \ V(S_{t+1}^1) \\ \mathsf{sample}^2 & = & R_{t+1}^2 + \gamma \ V(S_{t+1}^2) \\ & \vdots & = & \vdots \\ \mathsf{sample}^n & = & R_{t+1}^n + \gamma \ V(S_{t+1}^n) \end{array}$$

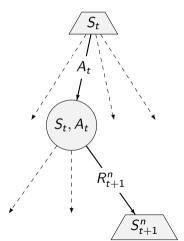


MDP (p, r known): Update V estimate by a weighted average:

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) [r(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

$$\begin{aligned} \mathsf{sample}^1 &=& R_{t+1}^1 + \gamma \ V(S_{t+1}^1) \\ \mathsf{sample}^2 &=& R_{t+1}^2 + \gamma \ V(S_{t+1}^2) \\ &\vdots &=& \vdots \\ \mathsf{sample}^n &=& R_{t+1}^n + \gamma \ V(S_{t+1}^n) \end{aligned}$$

$$V(S_t) \leftarrow \frac{1}{n} \sum_i \mathsf{sample}^i$$



MDP (p, r known): Update V estimate by a weighted average:

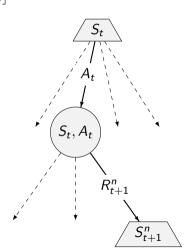
$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) [r(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

What about stop, try, try, ..., and average? Samples at time t.  $\pi(S_t) \to A_t$ , repeat  $A_t$ .

$$\begin{aligned} \mathsf{sample}^1 &=& R_{t+1}^1 + \gamma \ V(S_{t+1}^1) \\ \mathsf{sample}^2 &=& R_{t+1}^2 + \gamma \ V(S_{t+1}^2) \\ &\vdots &=& \vdots \\ \mathsf{sample}^n &=& R_{t+1}^n + \gamma \ V(S_{t+1}^n) \\ \end{aligned}$$

$$V(S_t) \leftarrow \frac{1}{n} \sum \mathsf{sample}^i$$

Problem: We cannot re-set to  $S_t$  easily.



$$\begin{array}{l} (1,1)_{\textbf{-.04}} \leadsto (1,2)_{\textbf{-.04}} \leadsto (1,3)_{\textbf{-.04}} \leadsto (1,2)_{\textbf{-.04}} \leadsto (1,3)_{\textbf{-.04}} \leadsto (2,3)_{\textbf{-.04}} \leadsto (3,3)_{\textbf{-.04}} \leadsto (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \leadsto (1,2)_{\textbf{-.04}} \leadsto (1,3)_{\textbf{-.04}} \leadsto (2,3)_{\textbf{-.04}} \leadsto (3,3)_{\textbf{-.04}} \leadsto (3,2)_{\textbf{-.04}} \leadsto (3,3)_{\textbf{-.04}} \leadsto (4,3)_{\textbf{+1}} \\ (1,1)_{\textbf{-.04}} \leadsto (2,1)_{\textbf{-.04}} \leadsto (3,1)_{\textbf{-.04}} \leadsto (3,2)_{\textbf{-.04}} \leadsto (4,2)_{\textbf{-1}} \ . \end{array}$$

$$\gamma = 1$$

$$\begin{array}{l} (1,1) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (2,3) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (4,3) \text{+1} \\ (1,1) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (2,3) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (3,2) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (4,3) \text{+1} \\ (1,1) \text{-.04} \leadsto (2,1) \text{-.04} \leadsto (3,1) \text{-.04} \leadsto (3,2) \text{-.04} \leadsto (4,2) \text{-1} \end{array}.$$

$$\gamma=1$$
 From first trial (episode):  $V(2,3)=$  ,  $V(1,3)=$  ,  $\dots$ 

$$\begin{array}{l} (1,1) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (2,3) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (4,3) \text{+1} \\ (1,1) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (2,3) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (3,2) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (4,3) \text{+1} \\ (1,1) \text{-.04} \leadsto (2,1) \text{-.04} \leadsto (3,1) \text{-.04} \leadsto (3,2) \text{-.04} \leadsto (4,2) \text{-1} \end{array}.$$

$$\gamma=1$$
 From first trial (episode):  $V(2,3)=0.92,\ V(1,3)=0.84,\ldots$ 

$$\begin{array}{l} (1,1) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (2,3) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (4,3) \text{+1} \\ (1,1) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (2,3) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (3,2) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (4,3) \text{+1} \\ (1,1) \text{-.04} \leadsto (2,1) \text{-.04} \leadsto (3,1) \text{-.04} \leadsto (3,2) \text{-.04} \leadsto (4,2) \text{-1} \end{array}.$$

$$\gamma = 1$$

From first trial (episode): V(2,3) = 0.92, V(1,3) = 0.84, ... In second episode, going from  $S_t = (1,3)$  to  $S_{t+1} = (2,3)$  with reward  $R_{t+1} = -0.04$ , hence:

$$V(1,3) = R_{t+1} + V(2,3) = -0.04 + 0.92 = 0.88$$

$$\begin{array}{l} (1,1) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (2,3) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (4,3) \text{+1} \\ (1,1) \text{-.04} \leadsto (1,2) \text{-.04} \leadsto (1,3) \text{-.04} \leadsto (2,3) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (3,2) \text{-.04} \leadsto (3,3) \text{-.04} \leadsto (4,3) \text{+1} \\ (1,1) \text{-.04} \leadsto (2,1) \text{-.04} \leadsto (3,1) \text{-.04} \leadsto (3,2) \text{-.04} \leadsto (4,2) \text{-1} \end{array}.$$

$$\gamma = 1$$

From first trial (episode): V(2,3) = 0.92, V(1,3) = 0.84,... In second episode, going from  $S_t = (1,3)$  to  $S_{t+1} = (2,3)$  with reward  $R_{t+1} = -0.04$ , hence:

$$V(1,3) = R_{t+1} + V(2,3) = -0.04 + 0.92 = 0.88$$

First estimate 0.84 is a bit lower than 0.88.  $V(S_t)$  is different than  $R_{t+1} + \gamma V(S_{t+1})$ 

$$\begin{array}{l} (1,1)_{\text{-.04}} \leadsto (1,2)_{\text{-.04}} \leadsto (1,3)_{\text{-.04}} \leadsto (1,2)_{\text{-.04}} \leadsto (1,3)_{\text{-.04}} \leadsto (2,3)_{\text{-.04}} \leadsto (3,3)_{\text{-.04}} \leadsto (4,3)_{\text{+1}} \\ (1,1)_{\text{-.04}} \leadsto (1,2)_{\text{-.04}} \leadsto (1,3)_{\text{-.04}} \leadsto (2,3)_{\text{-.04}} \leadsto (3,3)_{\text{-.04}} \leadsto (3,2)_{\text{-.04}} \leadsto (3,3)_{\text{-.04}} \leadsto (4,3)_{\text{+1}} \\ (1,1)_{\text{-.04}} \leadsto (2,1)_{\text{-.04}} \leadsto (3,1)_{\text{-.04}} \leadsto (3,2)_{\text{-.04}} \leadsto (4,2)_{\text{-1}} \ . \end{array}$$

$$\gamma = 1$$

From first trial (episode): V(2,3) = 0.92, V(1,3) = 0.84,... In second episode, going from  $S_t = (1,3)$  to  $S_{t+1} = (2,3)$  with reward  $R_{t+1} = -0.04$ , hence:

$$V(1,3) = R_{t+1} + V(2,3) = -0.04 + 0.92 = 0.88$$

- First estimate 0.84 is a bit lower than 0.88.  $V(S_t)$  is different than  $R_{t+1} + \gamma V(S_{t+1})$
- ▶ Update ( $\alpha$ × difference):  $V(S_t) \leftarrow V(S_t) + \alpha \Big( [R_{t+1} + \gamma V(S_{t+1})] V(S_t) \Big)$
- ightharpoonup lpha is the learning rate.

$$\begin{array}{l} (1,1) \textbf{-.04} \leadsto (1,2) \textbf{-.04} \leadsto (1,3) \textbf{-.04} \leadsto (1,2) \textbf{-.04} \leadsto (1,3) \textbf{-.04} \leadsto (2,3) \textbf{-.04} \leadsto (3,3) \textbf{-.04} \leadsto (4,3) \textbf{+1} \\ (1,1) \textbf{-.04} \leadsto (1,2) \textbf{-.04} \leadsto (1,3) \textbf{-.04} \leadsto (2,3) \textbf{-.04} \leadsto (3,3) \textbf{-.04} \leadsto (3,2) \textbf{-.04} \leadsto (3,3) \textbf{-.04} \leadsto (4,3) \textbf{+1} \\ (1,1) \textbf{-.04} \leadsto (2,1) \textbf{-.04} \leadsto (3,1) \textbf{-.04} \leadsto (3,2) \textbf{-.04} \leadsto (4,2) \textbf{-1} \end{array}.$$

$$\gamma = 1$$

From first trial (episode): V(2,3) = 0.92, V(1,3) = 0.84,... In second episode, going from  $S_t = (1,3)$  to  $S_{t+1} = (2,3)$  with reward  $R_{t+1} = -0.04$ , hence:

$$V(1,3) = R_{t+1} + V(2,3) = -0.04 + 0.92 = 0.88$$

- First estimate 0.84 is a bit lower than 0.88.  $V(S_t)$  is different than  $R_{t+1} + \gamma V(S_{t+1})$
- ▶ Update ( $\alpha$ × difference):  $V(S_t) \leftarrow V(S_t) + \alpha \Big( [R_{t+1} + \gamma V(S_{t+1})] V(S_t) \Big)$
- $ightharpoonup \alpha$  is the learning rate.
- $V(S_t) \leftarrow (1-\alpha)V(S_t) + \alpha \text{ (new sample)}$

#### Exponential moving average

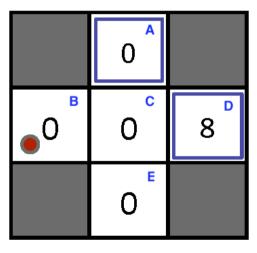
$$\overline{x}_n = (1 - \alpha)\overline{x}_{n-1} + \alpha x_n$$

What does it remember about the past? Try to derive:

$$\overline{x}_n = f(\alpha, x_n, x_{n-1}, x_{n-2}, x_{n-3}, \dots)$$

#### Example: TD Value learning

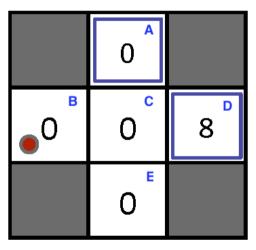
$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



- ightharpoonup Values represent initial V(s)
- Assume:  $\gamma = 1, \alpha = 0.5, \pi(s) = \rightarrow$

#### Example: TD Value learning

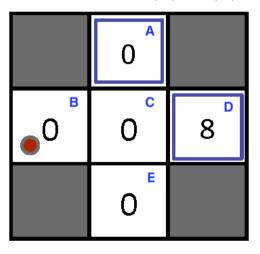
$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



- $\triangleright$  Values represent initial V(s)
- ightharpoonup Assume:  $\gamma = 1, \alpha = 0.5, \pi(s) = \rightarrow$
- $\triangleright$   $(B, \rightarrow, C), -2, \Rightarrow V(B)$ ?

# Example: TD Value learning

$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



- $\triangleright$  Values represent initial V(s)
- ightharpoonup Assume:  $\gamma = 1, \alpha = 0.5, \pi(s) = \rightarrow$
- $\triangleright$   $(B, \rightarrow, C), -2, \Rightarrow V(B)$ ?
- $\blacktriangleright$   $(C, \rightarrow, D), -2, \Rightarrow V(C)$ ?

# Temporal difference value learning: algorithm

Input: the policy  $\pi$  to be evaluated Algorithm parameter: step size  $\alpha \in (0,1]$ Initialize V(s), for all  $s \in \mathbb{S}^+$ , arbitrarily except that V(terminal) = 0Loop for each episode: Initialize SLoop for each step of episode:  $A \leftarrow \text{action given by } \pi \text{ for } S$ Take action A, observe R, S' $V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$  $S \leftarrow S'$ until S is terminal

What is wrong with the temporal difference Value learning?

The Good: Model-free value learning by mimicking Bellman updates.

# What is wrong with the temporal difference Value learning?

The Good: Model-free value learning by mimicking Bellman updates.

The Bad: How to turn values into a (new) policy?

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

# What is wrong with the temporal difference Value learning?

The Good: Model-free value learning by mimicking Bellman updates.

The Bad: How to turn values into a (new) policy?

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

$$\pi(s) = \arg\max_{a} Q(s, a)$$

# Q-learning

#### Reminder: V, Q-value iteration for MDPs

Value/Utility iteration (depth limited evaluation):

- ▶ Start:  $V_0(s) = 0$
- ▶ In each step update V by looking one step ahead:  $V_{k+1}(s) \leftarrow \max \sum_{s'} p(s' \mid s, a) [r(s, a, s') + \gamma V_k(s')]$

Q values more useful (think about updating  $\pi$ )

- ► Start:  $Q_0(s, a) = 0$
- ightharpoonup In each step update Q by looking one step ahead:

$$Q_{k+1}(s,a) \leftarrow \sum_{s'} p(s' \mid s,a) \left[ r(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

MDP update: 
$$Q_{k+1}(s,a) \leftarrow \sum_{s'} p(s'\mid s,a) \left[ r(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

MDP update: 
$$Q_{k+1}(s, a) \leftarrow \sum_{s'} p(s' \mid s, a) \left[ r(s, a, s') + \gamma \max_{a'} Q_k(s', a') \right]$$

Learn Q values as the robot/agent goes (temporal difference)

▶ Drive the robot (take action a in state s) and fetch rewards (s, a, s', R)

MDP update: 
$$Q_{k+1}(s,a) \leftarrow \sum_{s'} p(s' \mid s,a) \left[ r(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

Learn Q values as the robot/agent goes (temporal difference)

- ▶ Drive the robot (take action a in state s) and fetch rewards (s, a, s', R)
- ▶ We know old estimates Q(s, a) (and Q(s', a')), if not, initialize.

MDP update: 
$$Q_{k+1}(s, a) \leftarrow \sum_{s'} p(s' \mid s, a) \left[ r(s, a, s') + \gamma \max_{a'} Q_k(s', a') \right]$$

Learn Q values as the robot/agent goes (temporal difference)

- ▶ Drive the robot (take action a in state s) and fetch rewards (s, a, s', R)
- We know old estimates Q(s, a) (and Q(s', a')), if not, initialize.
- A new sample estimate (of Q(s, a)) at time t sample =  $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$

MDP update: 
$$Q_{k+1}(s,a) \leftarrow \sum_{s'} p(s'\mid s,a) \left[ r(s,a,s') + \gamma \max_{a'} Q_k(s',a') 
ight]$$

Learn Q values as the robot/agent goes (temporal difference)

- ▶ Drive the robot (take action a in state s) and fetch rewards (s, a, s', R)
- ▶ We know old estimates Q(s, a) (and Q(s', a')), if not, initialize.
- A new sample estimate (of Q(s, a)) at time t sample =  $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$
- ▶  $\alpha$  update  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\mathsf{sample} Q(S_t, A_t))$  or (the same)  $Q(S_t, A_t) \leftarrow (1 \alpha)Q(S_t, A_t) + \alpha \,\mathsf{sample}$

MDP update: 
$$Q_{k+1}(s, a) \leftarrow \sum_{s'} p(s' \mid s, a) \left[ r(s, a, s') + \gamma \max_{a'} Q_k(s', a') \right]$$

Learn Q values as the robot/agent goes (temporal difference)

- ▶ Drive the robot (take action a in state s) and fetch rewards (s, a, s', R)
- ▶ We know old estimates Q(s, a) (and Q(s', a')), if not, initialize.
- A new sample estimate (of Q(s, a)) at time t sample =  $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$
- $\begin{array}{l} \color{red} \alpha \; \text{update} \\ Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (\text{sample} Q(S_t, A_t)) \\ \text{or (the same)} \\ Q(S_t, A_t) \leftarrow (1 \alpha) Q(S_t, A_t) + \alpha \, \text{sample} \end{array}$

In each step Q approximates the optimal  $q^*$  function.

# Q-learning: algorithm (repeating episodes, until terminal or exhausted)

```
step size 0 < \alpha < 1
initialize Q(s, a) for all s \in \mathcal{S}, a \in \mathcal{A}(s)
repeat episodes:
    initialize S
    for for each step of episode: do
        choose A from A(S)
        take action A. observe R, S'
        Q(S,A) \leftarrow Q(S,A) + \alpha [R + \gamma \max_{a} Q(S',a) - Q(S,A)]
      S \leftarrow S'
until S is terminal
until Time is up, ...
```

- ▶ Drive the robot and fetch rewards. (s, a, s', R)
- ▶ We know old estimates Q(s, a) (and Q(s', a')), if not, initialize.
- ▶ A new sample estimate: sample =  $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$
- ightharpoonup lpha update:  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + lpha(\mathsf{sample} Q(S_t, A_t))$

- ▶ Drive the robot and fetch rewards. (s, a, s', R)
- ▶ We know old estimates Q(s, a) (and Q(s', a')), if not, initialize.
- ▶ A new sample estimate: sample =  $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$
- ightharpoonup lpha update:  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + lpha(\mathsf{sample} Q(S_t, A_t))$

#### Technicalities for the Q-learning agent

- ▶ Drive the robot and fetch rewards. (s, a, s', R)
- ▶ We know old estimates Q(s, a) (and Q(s', a')), if not, initialize.
- ▶ A new sample estimate: sample =  $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$
- ightharpoonup lpha update:  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + lpha(\mathsf{sample} Q(S_t, A_t))$

#### Technicalities for the Q-learning agent

► How to represent the *Q*-function?

- ▶ Drive the robot and fetch rewards. (s, a, s', R)
- ▶ We know old estimates Q(s, a) (and Q(s', a')), if not, initialize.
- ▶ A new sample estimate: sample =  $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$
- ightharpoonup lpha update:  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + lpha(\mathsf{sample} Q(S_t, A_t))$

#### Technicalities for the Q-learning agent

- ► How to represent the *Q*-function?
- ▶ What is the value for terminal? Q(s, Exit) or Q(s, None)

# From Q-learning to Q-learning agent

- ▶ Drive the robot and fetch rewards. (s, a, s', R)
- ▶ We know old estimates Q(s, a) (and Q(s', a')), if not, initialize.
- ▶ A new sample estimate: sample =  $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$
- ightharpoonup lpha update:  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + lpha(\mathsf{sample} Q(S_t, A_t))$

### Technicalities for the Q-learning agent

- ► How to represent the *Q*-function?
- ▶ What is the value for terminal? Q(s, Exit) or Q(s, None)
- How to drive? Where to drive next? Does it change over the course?







▶ Drive the known road or try a new one?







- ▶ Drive the known road or try a new one?
- ▶ Go to the university menza or try a nearby restaurant?







- ▶ Drive the known road or try a new one?
- ▶ Go to the university menza or try a nearby restaurant?
- ▶ Use the SW (operating system) I know or try a new one?







- ▶ Drive the known road or try a new one?
- ▶ Go to the university menza or try a nearby restaurant?
- ▶ Use the SW (operating system) I know or try a new one?
- Go to bussiness or study a demanding program?







- ▶ Drive the known road or try a new one?
- ► Go to the university menza or try a nearby restaurant?
- ▶ Use the SW (operating system) I know or try a new one?
- Go to bussiness or study a demanding program?
- **.** . . .

Random ( $\epsilon$ -greedy):

### Random ( $\epsilon$ -greedy):

Flip a coin every step.

### Random ( $\epsilon$ -greedy):

- ► Flip a coin every step.
- ightharpoonup With probability  $\epsilon$ , act randomly.

### Random ( $\epsilon$ -greedy):

- ► Flip a coin every step.
- $\blacktriangleright$  With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 \epsilon$ , use the policy.

#### Random ( $\epsilon$ -greedy):

- Flip a coin every step.
- $\blacktriangleright$  With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 \epsilon$ , use the policy.

Problems with randomness?

### Random ( $\epsilon$ -greedy):

- Flip a coin every step.
- $\blacktriangleright$  With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 \epsilon$ , use the policy.

#### Problems with randomness?

► Keeps exploring forever.

#### Random ( $\epsilon$ -greedy):

- Flip a coin every step.
- $\blacktriangleright$  With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 \epsilon$ , use the policy.

#### Problems with randomness?

- ► Keeps exploring forever.
- ▶ Should we keep  $\epsilon$  fixed (over learning)?

### Random ( $\epsilon$ -greedy):

- ► Flip a coin every step.
- $\blacktriangleright$  With probability  $\epsilon$ , act randomly.
- ▶ With probability  $1 \epsilon$ , use the policy.

#### Problems with randomness?

- ► Keeps exploring forever.
- ▶ Should we keep  $\epsilon$  fixed (over learning)?
- ightharpoonup  $\epsilon$  same everywhere?

### What we have learned

- ► Agent/robot may learn by acting an getting rewards
- Model based vs. model-free methods
- ▶ Direct learning vs. temporal-difference learning
- From learning state values to Q-learning

### References I

Further reading: Chapter 21 of [2] (chapter 23 of [3]). More detailed discussion in [4], chapters 2, 5, and 6.

- Dan Klein and Pieter Abbeel.
   UC Berkeley CS188 Intro to AI course materials. http://ai.berkeley.edu/.
   Used with permission of Pieter Abbeel.
- [2] Stuart Russell and Peter Norvig. Artificial Intelligence: A Modern Approach. Prentice Hall, 3rd edition, 2010. http://aima.cs.berkeley.edu/.
- [3] Stuart Russell and Peter Norvig.

  Artificial Intelligence: A Modern Approach.

  Prentice Hall, 4th edition, 2021.

### References II

[4] Richard S. Sutton and Andrew G. Barto. Reinforcement Learning; an Introduction.

MIT Press, 2nd edition, 2018.

http://www.incompleteideas.net/book/the-book-2nd.html.