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

Notes -

In Czech: Posilované učení

(Multi-armed) Bandits







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

(Multi-armed) Bandits







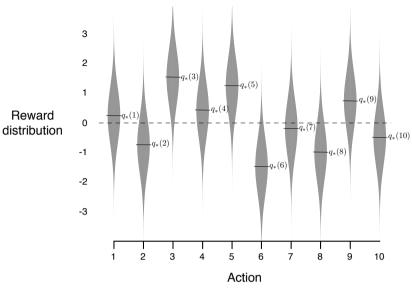
Think about not one but 10 arms you may choose to pull.

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

2/39

Notes -

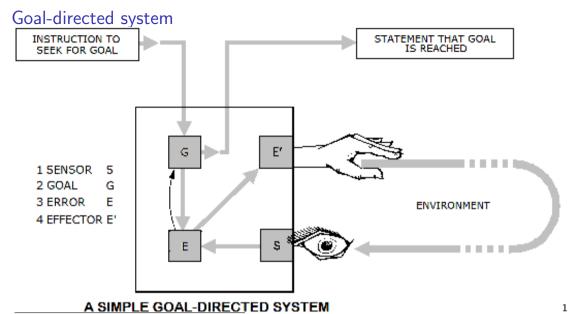
10 armed bandit, what arm to pull?



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

Notes -

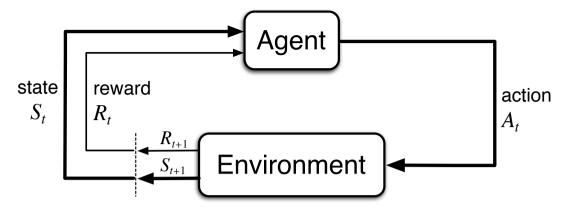
- 10 different arms
- action pulling k—th arm
- value of the action, i.e. q(a) is stochastic (Gaussian around $q^*(a)$) and uknown!
- Think about expectimax tree. How would it look like?
- Playing (pulling) many times, what is the policy? Think also that each sample costs something.



Notes -

¹Figure 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.

²Scheme from [4]

5 / 39

2

Notes

Reinforcement Learning - performing actions, learning from rewards



6/39

Notes -

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 policies³

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

7/39

Notes -

Policy search is a more advanced topic, only touched by this course. Later in master programme. Reinforement learning beating humans in playing Atari games: https://deepmind.google/discover/blog/agent57-outperforming-the-human-atari-benchmark/

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.

Notes

8/39

For MDPs, we know p, r for all possible states and actions.

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:

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

Notes —

8/39

For MDPs, we know p, r for all possible states and actions.

(Transition) Model-based learning

The main idea: Do something, and:

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

Notes -

- Where to start?
- When does it end?
- How long does it take?
- When to stop (the learning phase)?

(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.

Notes

- Where to start?
- When does it end?
- How long does it take?
- When to stop (the learning phase)?

(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.

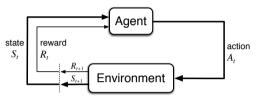
Solve the learned MDP.

Notes

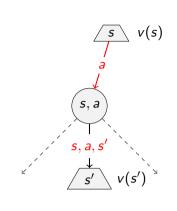
- Where to start?
- When does it end?
- How long does it take?
- When to stop (the learning phase)?

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').



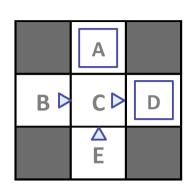
Notes

In ai-gym env.step(action) returns s', r(s, action, s'), It is defined by the environment (robot simulator, system, ...) not by the (algorithms)

Model-based learning: Grid example

Input Policy π

Observed Episodes (Training)



Assume: γ = 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

⁴Figure from [1]

Notes

Learned Model

$$\widehat{T}(s,a,s')$$

T(B, east, C) = 1.00 T(C, east, D) = 0.75 T(C, east, A) = 0.25

$$\widehat{R}(s,a,s')$$

R(B, east, C) = -1 R(C, east, D) = -1 R(D, exit, x) = +10 ...

Learning transition model

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

Episode 1 Episode 2

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 Episode 4

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

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

12 / 39

Notes

(C, east) combination performed 4 times, 3 times landed in D, once in A. Hence, $\hat{p}(D \mid C, east) = 0.75$.

Learning reward function

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

Episode 1

Episode 2 B, east, C, -1

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

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

Episode 3

Episode 4

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

Notes

Whenever (C, east, D) performed, received reward was -1. Hence, $\hat{r}(C, east, D) = -1$.

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 $[a_1, a_2, \dots a_N]$.

14 / 39

Notes -

Just to avoid confusion. There are many more samples than possible ages (positive integer). Think about $N\gg 100$.

- Model based eventually, we learn the correct model.
- Model free no need for weighting; this is achieved through the frequencies of different ages within the samples (most frequent and hence most probable ages simply come up many times).

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

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

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

Notes

Just to avoid confusion. There are many more samples than possible ages (positive integer). Think about $N\gg 100$.

- Model based eventually, we learn the correct model.
- Model free no need for weighting; this is achieved through the frequencies of different ages within the samples (most frequent and hence most probable ages simply come up many times).

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

Random variable age A.

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

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

Model based

Model free

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

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

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

Notes

Just to avoid confusion. There are many more samples than possible ages (positive integer). Think about $N\gg 100$.

- Model based eventually, we learn the correct model.
- Model free no need for weighting; this is achieved through the frequencies of different ages within the samples (most frequent and hence most probable ages simply come up many times).

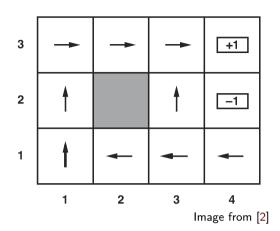
Model-free learning

15 / 39

Notes -

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 $v^{\pi}(s)$



16 / 39

Notes -

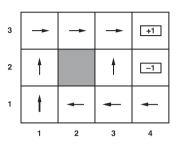
Executing policies - training, then learning from the observations. We want to do the policy evaluation but the necessary model is not known.

The word passive means we just follow a prescribed policy $\pi(s)$.

Direct evaluation from episodes

Value of s for π – expected sum of discounted rewards – expected return

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



Notes

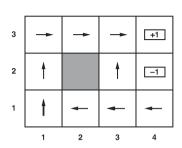
- Act according to the policy.
- When visiting a state, remember what the sum of discounted rewards (returns) turned out to be.
- Compute average of the returns.
- Each trial episode provides a sample of v^{π} .

What is $v^{\pi}(3,2)$ after these episodes?

Direct evaluation from episodes

Value of s for π – expected sum of discounted rewards – expected return

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



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

Notes

- Act according to the policy.
- When visiting a state, remember what the sum of discounted rewards (returns) turned out to be.
- Compute average of the returns.
- Each trial episode provides a sample of v^{π} .

What is $v^{\pi}(3,2)$ after these episodes?

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

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

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

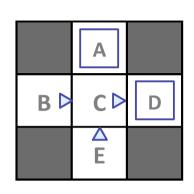
Notes

- Not visited during the first episode.
- Visited once in the second, gathered return G = -0.04 0.04 + 1 = 0.92.
- Visited once in the third, return G = -0.04 1 = -1.04.
- Value, average return is (0.92 1.04)/2 = -0.06.

Direct evaluation: Grid example

Input Policy π

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

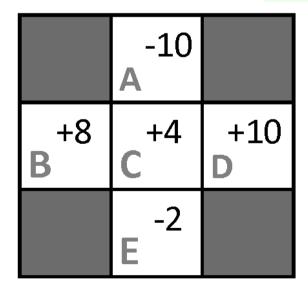
Episode 3

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

Notes



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

What is v(C) after the 4 episodes?

Episode 1

Episode 2

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

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

B, east, C, -1

Episode 3

Episode 4

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

Notes -

- Episode 1, G = -1 + 10 = 9
- Episode 2, G = -1 + 10 = 9
- Episode 3, G = -1 + 10 = 9
- Episode 4, G = -1 10 = -11
- Average return v(C) = (9+9+9-11)/4 = 4

For first-visit variant, B is correct. For every-visit variant, D is correct.

N can be lower than M (state does not have to be attended in every episode). For every-visit variant, N can be higher than M (a state can be visited several times in one episode).

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

What is v(C) after the 4 episodes?

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

What is the relation of M and N?

- A N = M
- $\mathbf{B} \ N \leq M$
- C N > 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

Notes -

- Episode 1, G = -1 + 10 = 9
- Episode 2, G = -1 + 10 = 9
- Episode 3, G = -1 + 10 = 9
- Episode 4, G = -1 10 = -11
- Average return v(C) = (9+9+9-11)/4 = 4

For first-visit variant, B is correct. For every-visit variant, D is correct.

N can be lower than M (state does not have to be attended in every episode). For every-visit variant, N can be higher than M (a state can be visited several times in one episode).

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 $\boldsymbol{\pi}$ to be evaluated

Loop forever (for each episode):

Initialize:

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

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

Genera

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

$$G \leftarrow 0$$

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

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

Append G to $Returns(S_t)$

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

Notes -

21 / 39

The algorithm can be easily expanded to $Q(S_t, A_t)$. Instead of visiting S_t we consider visiting of a pair S_t, A_t .

Direct evaluation algorithm (first-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 π 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 followi

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

$$G \leftarrow 0$$

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 $Returns(S_t)$

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

Notes -

21 / 39

The algorithm can be easily expanded to $Q(S_t, A_t)$. Instead of visiting S_t we consider visiting of a pair S_t, A_t .

The good:

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

22 / 39

Notes -

In second trial, we visit (3,2) for the first time. We already know that the successor (3,3) has probably a high value but the method does not use until the end of the trial episode.

Before updating V(s) we have to wait until the training episode ends.

The good:

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

The bad:

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

22 / 39

Notes -

In second trial, we visit (3,2) for the first time. We already know that the successor (3,3) has probably a high value but the method does not use until the end of the trial episode.

Before updating V(s) we have to wait until the training episode ends.

The good:

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

The bad:

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

► Each state value learned in isolation.

22 / 39

Notes -

In second trial, we visit (3,2) for the first time. We already know that the successor (3,3) has probably a high value but the method does not use until the end of the trial episode.

Before updating V(s) we have to wait until the training episode ends.

The good:

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

The bad:

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

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

Notes

In second trial, we visit (3,2) for the first time. We already know that the successor (3,3) has probably a high value but the method does not use until the end of the trial episode.

Before updating V(s) we have to wait until the training episode ends.

(on-line) Policy evaluation?

In MDP, we did:

- Initialize the values: $V_0^{\pi}(s) = 0$
- ▶ In each iteration, replace V with a one-step-look-ahead: $V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) \left[r(s, \pi(s), s') + \gamma V_k^{\pi}(s') \right]$

(on-line) Policy evaluation?

In MDP, we did:

- ▶ Initialize the values: $V_0^{\pi}(s) = 0$
- ▶ In each iteration, replace V with a one-step-look-ahead: $V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) [r(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$

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:

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

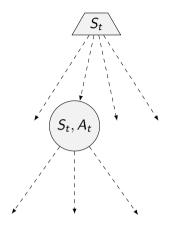
24 / 39

Notes

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 .



24 / 39

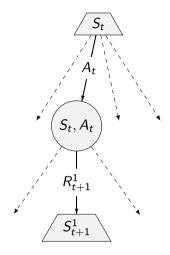
Notes

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 .

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



24 / 39

Notes

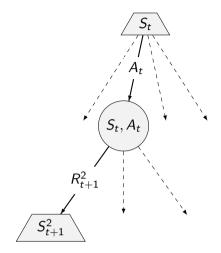
MDP (p, r known): Update V estimate by a weighted average:

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

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

sample¹ =
$$R_{t+1}^1 + \gamma V(S_{t+1}^1)$$

sample² = $R_{t+1}^2 + \gamma V(S_{t+1}^2)$



24 / 39

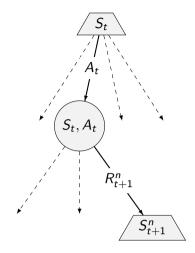
Notes

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 .

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



24 / 39

Notes

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

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

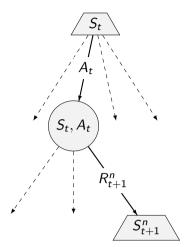
What about stop, try, try, ..., and average? Samples at time t. $\pi(S_t) \to A_t$, repeat A_t .

sample¹ =
$$R_{t+1}^1 + \gamma V(S_{t+1}^1)$$

sample² = $R_{t+1}^2 + \gamma V(S_{t+1}^2)$
 \vdots = \vdots

sampleⁿ =
$$R_{t+1}^n + \gamma V(S_{t+1}^n)$$

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



24 / 39

Notes

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

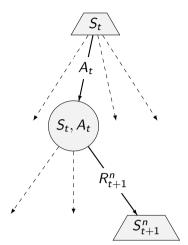
$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} p(s' \mid s, \pi(s)) [r(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

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

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

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

Problem: We cannot re-set to S_t easily.



Notes

It looks promising. Unfortunately, we cannot do it that way. After an action, the robot is in a next state and cannot go back to the very same state where it was before. Energy was consumed and some actions may be irreversible; think about falling into a hole. We have to utilize the s, a, s' experience anytime when performed/visited.

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

 $\gamma = 1$

Notes -

Trial episode: acting, observing, until it stops (in a terminal state or by a limit).

We visit S(1,3) twice during the first episode. Its value estimate is the average of two returns. Note the main difference. In *Direct evaluation*, we had to wait until the end of the episode, compute G_t for each

t on the way, and then we update $V(S_t)$. We can do it α incrementally

$$V(S_t) \leftarrow V(S_t) + \alpha \Big(G_t - V(S_t)\Big)$$

In TD learning, we update as we go.

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

$$\gamma = 1$$

From first trial (episode): V(2,3) =, V(1,3) =,...

25 / 39

Notes -

Trial episode: acting, observing, until it stops (in a terminal state or by a limit).

We visit S(1,3) twice during the first episode. Its value estimate is the average of two returns. Note the main difference. In *Direct evaluation*, we had to wait until the end of the episode, compute G_t for each

t on the way, and then we update $V(S_t)$. We can do it α incrementally

$$V(S_t) \leftarrow V(S_t) + \alpha \Big(G_t - V(S_t)\Big)$$

In TD learning, we update as we go.

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

$$\gamma = 1$$

From first trial (episode): V(2,3) = 0.92, V(1,3) = 0.84,...

25 / 39

Notes -

Trial episode: acting, observing, until it stops (in a terminal state or by a limit).

We visit S(1,3) twice during the first episode. Its value estimate is the average of two returns. Note the main difference. In *Direct evaluation*, we had to wait until the end of the episode, compute G_t for each t on the way, and then we update $V(S_t)$. We can do it α incrementally

$$V(S_t) \leftarrow V(S_t) + \alpha \Big(G_t - V(S_t)\Big)$$

In TD learning, we update as we go.

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

$$\gamma = 1$$

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:

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

25 / 39

Notes -

Trial episode: acting, observing, until it stops (in a terminal state or by a limit).

We visit S(1,3) twice during the first episode. Its value estimate is the average of two returns. Note the main difference. In *Direct evaluation*, we had to wait until the end of the episode, compute G_t for each t on the way, and then we update $V(S_t)$. We can do it α incrementally

$$V(S_t) \leftarrow V(S_t) + \alpha \Big(G_t - V(S_t)\Big)$$

In TD learning, we update as we go.

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

$$\gamma = 1$$

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 $R_{t+1} + \gamma V(S_{t+1})$

Notes -

Trial episode: acting, observing, until it stops (in a terminal state or by a limit).

We visit S(1,3) twice during the first episode. Its value estimate is the average of two returns.

Note the main difference. In *Direct evaluation*, we had to wait until the end of the episode, compute G_t for each t on the way, and then we update $V(S_t)$. We can do it α incrementally

$$V(S_t) \leftarrow V(S_t) + \alpha \Big(G_t - V(S_t)\Big)$$

In TD learning, we update as we go.

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

$$\gamma = 1$$

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 (α × 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.

Notes -

Trial episode: acting, observing, until it stops (in a terminal state or by a limit).

We visit S(1,3) twice during the first episode. Its value estimate is the average of two returns. Note the main difference. In *Direct evaluation*, we had to wait until the end of the episode, compute G_t for each

t on the way, and then we update $V(S_t)$. We can do it α incrementally

$$V(S_t) \leftarrow V(S_t) + \alpha \Big(G_t - V(S_t)\Big)$$

In TD learning, we update as we go.

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

$$\gamma = 1$$

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 (α × 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.
- $V(S_t) \leftarrow (1-\alpha)V(S_t) + \alpha \text{ (new sample)}$

Notes -

Trial episode: acting, observing, until it stops (in a terminal state or by a limit).

We visit S(1,3) twice during the first episode. Its value estimate is the average of two returns.

Note the main difference. In *Direct evaluation*, we had to wait until the end of the episode, compute G_t for each t on the way, and then we update $V(S_t)$. We can do it α incrementally

$$V(S_t) \leftarrow V(S_t) + \alpha \Big(G_t - V(S_t)\Big)$$

In TD learning, we update as we go.

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:

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

Notes

Recursively insetring we end up with

$$\overline{x}_n = \alpha \left[x_n + (1 - \alpha)x_{n-1} + (1 - \alpha)^2 x_{n-2} + \cdots \right]$$

We already know the sum of geometric series for r < 1

$$1 + r + r^2 + r^3 + \dots = \frac{1}{1 - r}$$

Putting $r = 1 - \alpha$, we see that

$$\frac{1}{\alpha} = 1 + (1 - \alpha) + (1 - \alpha)^2 + \cdots$$

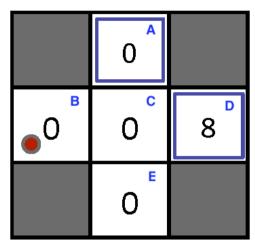
And hence:

$$\overline{x}_n = \frac{x_n + (1 - \alpha)x_{n-1} + (1 - \alpha)^2 x_{n-2} + \cdots}{1 + (1 - \alpha) + (1 - \alpha)^2 + (1 - \alpha)^3 + \cdots}$$

a weighted average that exponentially forgets about the past.

Example: TD Value learning

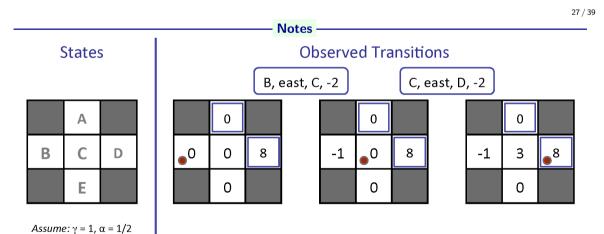
$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



 \triangleright Values represent initial V(s)

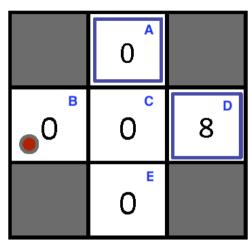
 $V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + \alpha \left[R(s, \pi(s), s') + \gamma V^{\pi}(s')\right]$

Assume: $\gamma = 1, \alpha = 0.5, \pi(s) = \rightarrow$



Example: TD Value learning

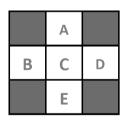
$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



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

27 / 39





Assume: $\gamma = 1$, $\alpha = 1/2$





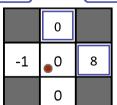
B, east, C, -2

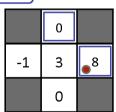
8

0

0

0



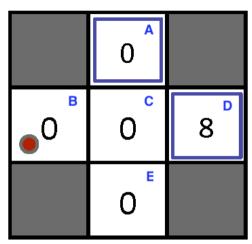


C, east, D, -2

$$V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + \alpha \left[R(s, \pi(s), s') + \gamma V^{\pi}(s')\right]$$

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 $(B, \rightarrow, C), -2, \Rightarrow V(B)$?
- \triangleright $(C, \rightarrow, D), -2, \Rightarrow V(C)$?

27 / 39

States

Α C В D E

Assume: $\gamma = 1$, $\alpha = 1/2$



Observed Transitions

-1

B, east, C, -2

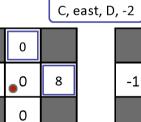
8

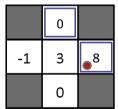
0

0

0

0





$$V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + \alpha \left[R(s, \pi(s), s') + \gamma V^{\pi}(s')\right]$$

Temporal difference value learning: algorithm

Input: the policy π to be evaluated

Algorithm parameter: step size $\alpha \in (0,1]$

Initialize V(s), for all $s \in S^+$, arbitrarily except that V(terminal) = 0

Loop for each episode:

Initialize S

Loop 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.

- Notes -

29 / 39

Learn Q-values, not V-values, and make the action selection model-free too!

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?

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

29 / 39

Notes -

Learn Q-values, not V-values, and make the action selection model-free too!

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?

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

Notes

29 / 39

Learn Q-values, not V-values, and make the action selection model-free too!

Q-learning

30 / 39

- Notes -

So far we walked as prescribed by a $\pi(s)$ because we did not know how to act better.

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

Q values more useful (think about updating π)

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

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

Notes

31/39

Draw the (s)-(s,a)-(s')-(s',a') tree. It will be also handy when discussing exploration vs. exploitation – where to drive next.

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

Notes -

There are alternatives how to compute the sample value. SARSA method takes $Q(S_{t+1}, A_{t+1})$ directly, not the max. More next week.

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)

Notes

32 / 39

There are alternatives how to compute the sample value. SARSA method takes $Q(S_{t+1}, A_{t+1})$ directly, not the max. More next week.

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.

Notes

32 / 39

There are alternatives how to compute the sample value. SARSA method takes $Q(S_{t+1}, A_{t+1})$ directly, not the max. More next week.

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 = $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$

Notes

There are alternatives how to compute the sample value. SARSA method takes $Q(S_{t+1}, A_{t+1})$ directly, not the max. More next week.

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 = $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$
- ▶ α update $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\mathsf{sample} Q(S_t, A_t))$ or (the same) $Q(S_t, A_t) \leftarrow (1 \alpha)Q(S_t, A_t) + \alpha \,\mathsf{sample}$

Notes

There are alternatives how to compute the sample value. SARSA method takes $Q(S_{t+1}, A_{t+1})$ directly, not the max. More next week.

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 = $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$
- ▶ α update $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\mathsf{sample} Q(S_t, A_t))$ or (the same) $Q(S_t, A_t) \leftarrow (1 \alpha)Q(S_t, A_t) + \alpha \,\mathsf{sample}$

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

Notes

There are alternatives how to compute the sample value. SARSA method takes $Q(S_{t+1}, A_{t+1})$ directly, not the max. More next week.

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

```
step size 0 < \alpha \le 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 \mathcal{A}(S) take action A, observe R,S' Q(S,A) \leftarrow Q(S,A) + \alpha \big[R + \gamma \max_a Q(S',a) - Q(S,A)\big] S \leftarrow S' until S is terminal until Time is up, . . .
```

33 / 39

Notes

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

34 / 39

Notes -

Q-function for a discrete, finite problem? But what about continous space or discrete but a very large one? Use the (s)-(s,a)-(s')-(s',a') tree to discuss the next-action selection.

- ▶ 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: $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} Q(S_t, A_t))$

Technicalities for the Q-learning agent

34 / 39

Notes

Q-function for a discrete, finite problem? But what about continous space or discrete but a very large one? Use the (s)-(s,a)-(s')-(s',a') tree to discuss the next-action selection.

- ▶ 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?

Notes

Q-function for a discrete, finite problem? But what about continous space or discrete but a very large one? Use the (s)-(s,a)-(s')-(s',a') tree to discuss the next-action selection.

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

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)

Notes

Q-function for a discrete, finite problem? But what about continous space or discrete but a very large one? Use the (s)-(s,a)-(s')-(s',a') tree to discuss the next-action selection.

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?

Notes

Q-function for a discrete, finite problem? But what about continous space or discrete but a very large one? Use the (s)-(s,a)-(s')-(s',a') tree to discuss the next-action selection.







▶ 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 (ϵ -greedy):

Notes -

- ullet We can think about lowering ϵ as the learning progresses.
- Favor unexplored states be optimistic exploration functions f(u, n) = u + k/n, where u is the value estimated, and n is the visit count, and k is the training/simulation episode.

Random (ϵ -greedy):

Flip a coin every step.

Notes -

- ullet We can think about lowering ϵ as the learning progresses.
- Favor unexplored states be optimistic exploration functions f(u, n) = u + k/n, where u is the value estimated, and n is the visit count, and k is the training/simulation episode.

Random (ϵ -greedy):

- ► Flip a coin every step.
- \blacktriangleright With probability ϵ , act randomly.

Notes

- ullet We can think about lowering ϵ as the learning progresses.
- Favor unexplored states be optimistic exploration functions f(u, n) = u + k/n, where u is the value estimated, and n is the visit count, and k is the training/simulation episode.

Random (ϵ -greedy):

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

Notes

- ullet We can think about lowering ϵ as the learning progresses.
- Favor unexplored states be optimistic exploration functions f(u, n) = u + k/n, where u is the value estimated, and n is the visit count, and k is the training/simulation episode.

Random (ϵ -greedy):

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

Problems with randomness?

Notes

- ullet We can think about lowering ϵ as the learning progresses.
- Favor unexplored states be optimistic exploration functions f(u, n) = u + k/n, where u is the value estimated, and n is the visit count, and k is the training/simulation episode.

Random (ϵ -greedy):

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

Problems with randomness?

► Keeps exploring forever.

Notes

- ullet We can think about lowering ϵ as the learning progresses.
- Favor unexplored states be optimistic exploration functions f(u, n) = u + k/n, where u is the value estimated, and n is the visit count, and k is the training/simulation episode.

Random (ϵ -greedy):

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

Problems with randomness?

- Keeps exploring forever.
- ▶ Should we keep ϵ fixed (over learning)?

Notes

- We can think about lowering ϵ as the learning progresses.
- Favor unexplored states be optimistic exploration functions f(u, n) = u + k/n, where u is the value estimated, and n is the visit count, and k is the training/simulation episode.

Random (ϵ -greedy):

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

Problems with randomness?

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

Notes

- We can think about lowering ϵ as the learning progresses.
- Favor unexplored states be optimistic exploration functions f(u, n) = u + k/n, where u is the value estimated, and n is the visit count, and k is the training/simulation episode.

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

37 / 39

Notes -

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.incomplete ideas.net/book/the-book-2nd.html.