

Reinforcement learning II

Active learning

Tomáš Svoboda

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

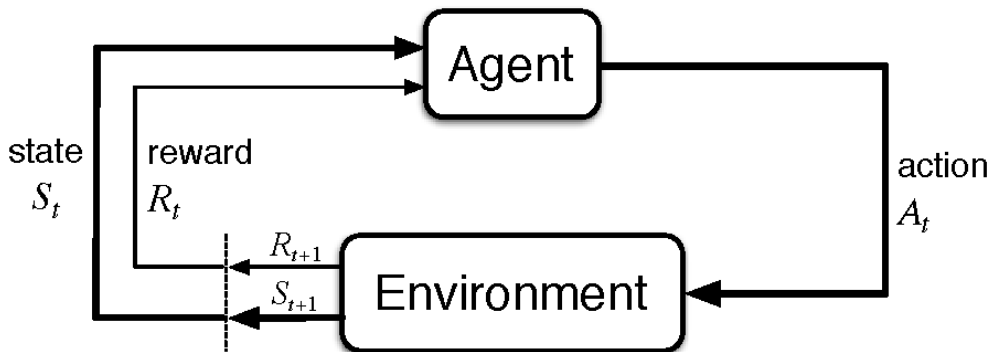
April 9, 2026

1 / 37

Notes

Not all slides with notes. What can be noted about the title page, eh?

Recap: Reinforcement Learning



1

- ▶ Feedback in form of Rewards
- ▶ Learn to act so as to maximize sum of expected rewards.
- ▶ In `kuimaze` package, `env.step(action)` is the method.

¹Scheme from [3]

2 / 37

Notes

Robot/agent action changes environment.

- Environment is everything ...
- battery state
- robot position
- object position (manipulation task)
- ...

Learning to control flippers



- ▶ What are the states?
- ▶ How to design rewards?
- ▶ How to perform training episodes (roll-outs)?
- ▶ Simulator to reality gap.

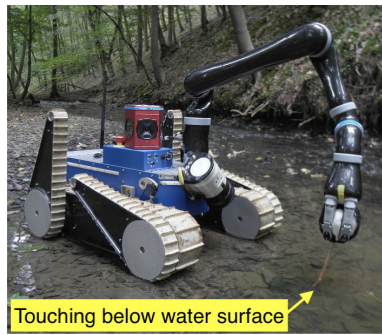
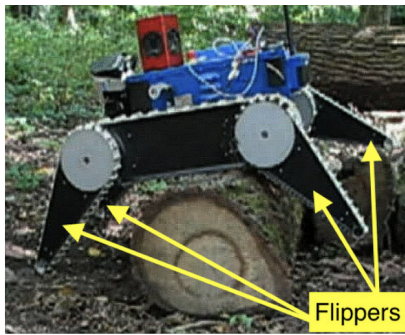
<http://cyber.felk.cvut.cz/vras/>

3 / 37

Notes

- States may contain interoceptive as well as exteroceptive sensing.
- Reward shaping.
- Train in simulator, then go for a real roll-out, back to simulator and so on.
- Physical simulator for robot terrain interactions.
- Sensor models.

Next few slides display a possible parameterization of the flipper control task.



- ◆ **Construction:** 2× main tracks, 4× subtracks (flippers), differential break
great stability and climbing capability
- ◆ **Sensor suite:** SICK LMS-151 range finder, Ladybug omnicam, Xsens MTi-G IMU
3D sensing and localization
- ◆ **Control inputs:** Velocity vector, 4× flipper angle, 4× flipper stiffness,
differential break (0/1)
difficult to control all of them manually!

Notes

This and the next three slides introduce some ideas and approaches published in:

- Karel Zimmermann, Petr Zuzanek, Michal Reinstein, Vaclav Hlavac. Adaptive traversability of unknown complex terrain with obstacles for mobile robots. In 2014 IEEE international conference on robotics and automation (ICRA).

The work has been extended in several directions:

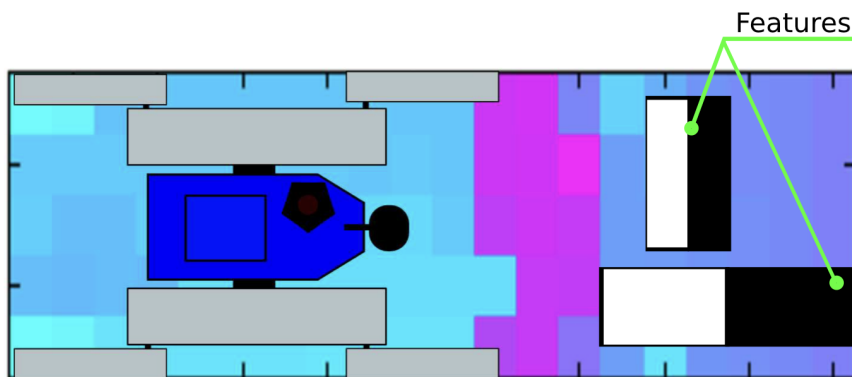
- Martin Pecka, Vojtěch Šalanský, Karel Zimmermann, Tomáš Svoboda. Autonomous flipper control with safety constraints. In 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS).
- M. Pecka, K. Zimmermann, M. Reinstein, and T. Svoboda. Controlling Robot Morphology from Incomplete Measurements. In *IEEE Transactions on Industrial Electronics*, Feb 2017, Vol 64, Issue: 2, pp. 1773-1782
- M. Pecka, K. Zimmermann, T. Svoboda. Fast Simulation of Vehicles with Non-deformable Tracks. In *Intelligent Robots and Systems (IROS)*, 2017.
- Martin Pecka, Karel Zimmermann, Matěj Petrlík, Tomáš Svoboda. Data-driven Policy Transfer with Imprecise Perception Simulation. *IEEE Robotics and Automation Letters*, Vol. 3, Issue 4, Oct 2018

State $s \in \mathcal{S} \subset \mathbb{R}^n$ concatenates:

◆ Proprioceptive measurements: roll, pitch, torques, velocity, acceleration.

◆ Local exteroceptive measurements.

features on digital elevation map
with fixed size.



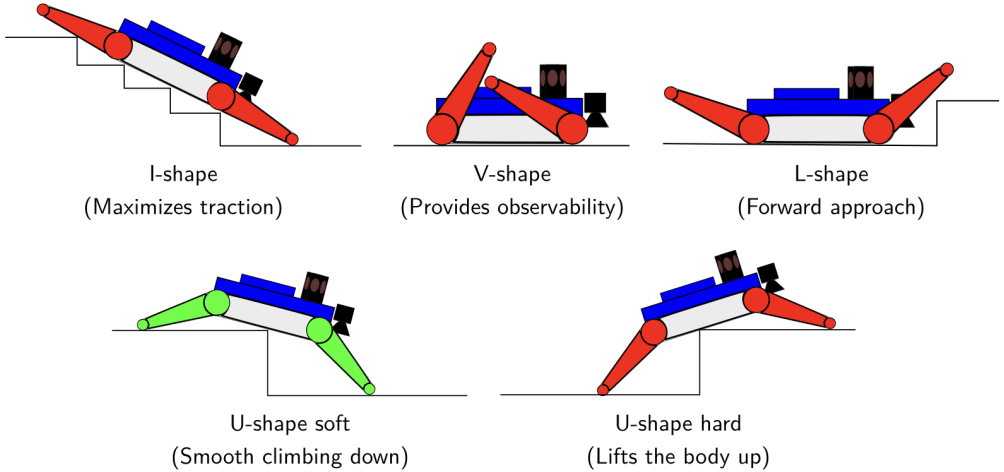
5 / 37

Notes

- Colors encode height of the terrain.
- Haar's features: Think about sum of heights in white area - sum of heights in the black area.

Instead of $\mathbf{a} \in \mathcal{A} \subset \mathbb{R}^8$ we consider only 5 configurations²:

$$\mathcal{A} = \{\text{I-shape}, \text{V-shape}, \text{L-shape}, \text{U-shape soft}, \text{U-shape hard}\}$$

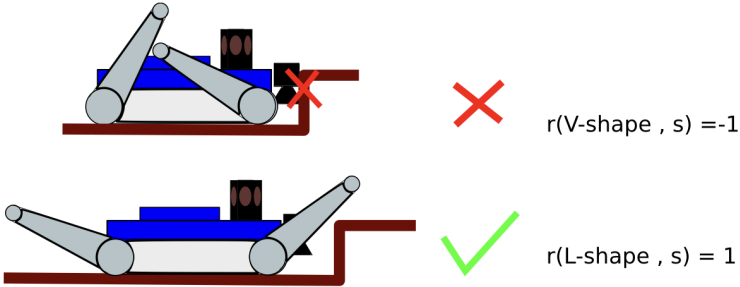


Notes

- Discretization of the Action space.
- Colors of the flippers encode whether they are stiff (red) or soft – terrain compliant (green).

Reward $r(a, s) : \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is a weighted sum of following contributions:

1. Safe pitch and roll reward, avoiding tipping over
2. Smoothness reward, suppresses body hits
3. Speed reward, drives robot forward
4. User denoted reward (penalty) indicating the success (failure) of the particular maneuver indicates failure/possible damages



7 / 37

Notes

- Hand-crafted reward function.

Discuss the difference with AlphaGo Zero, Atari games etc. There you can afford to learn about your performance only at the end of the game, as you can play many many games.

In robotics, you usually can't do that...

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, $a \in \mathcal{A}(s)$
- ▶ A transition model $p(s'|s, a)$ (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 p and reward r functions not known.
- ▶ Agent/robot must act and learn from experience.

Notes

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 were correct.

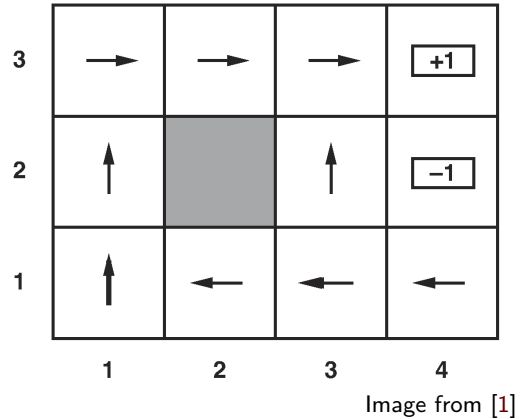
Learning MDP model:

- ▶ Try s, a , observe s' , count s, a, s' .
- ▶ Normalize to get and estimate of $p(s'|s, a)$
- ▶ Discover each $r(s, a, s')$ when experienced.

Solve the learned MDP.

Model-free learning

- ▶ r, p not known.
- ▶ Move around, observe.
- ▶ And learn on the way.
- ▶ **Goal:** Learn the state value $v(s)$, or (better), q-value $q(s, a)$ functions.



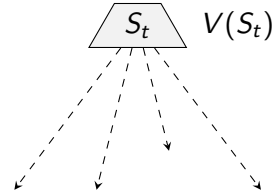
Notes

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

Recap: V -learning

Learn $V(s)$ values as the robot/agent goes (temporal difference).

- ▶ time t , at S_t
 - ▶ select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
 - ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma V(S_{t+1})$
 - ▶ α temporal difference update
 $V(S_t) \leftarrow V(S_t) + \alpha(\text{sample} - V(S_t))$
 - ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)



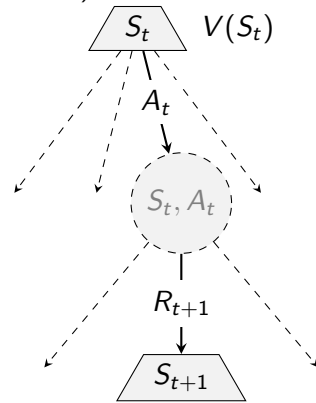
Notes

The S_t, A_t chance node is depicted in gray purposely. It is not directly observable, it is our model of uncertain outcome of action A_t taken in state S_t .

Recap: V - learning

Learn $V(s)$ values as the robot/agent goes (temporal difference).

- ▶ time t , at S_t
- ▶ select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
- ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma V(S_{t+1})$
- ▶ α temporal difference update
 $V(S_t) \leftarrow V(S_t) + \alpha(\text{sample} - V(S_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)



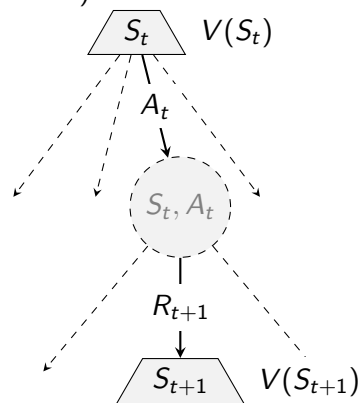
Notes

The S_t, A_t chance node is depicted in gray purposively. It is not directly observable, it is our model of uncertain outcome of action A_t taken in state S_t .

Recap: V - learning

Learn $V(s)$ values as the robot/agent goes (temporal difference).

- ▶ time t , at S_t
- ▶ select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
- ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma V(S_{t+1})$
- ▶ α temporal difference update
 $V(S_t) \leftarrow V(S_t) + \alpha(\text{sample} - V(S_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)



11 / 37

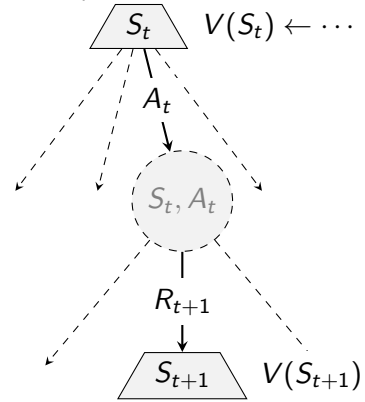
Notes

The S_t, A_t chance node is depicted in gray purposively. It is not directly observable, it is our model of uncertain outcome of action A_t taken in state S_t .

Recap: V - learning

Learn $V(s)$ values as the robot/agent goes (temporal difference).

- ▶ time t , at S_t
- ▶ select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
- ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma V(S_{t+1})$
- ▶ α temporal difference update
 $V(S_t) \leftarrow V(S_t) + \alpha(\text{sample} - V(S_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)



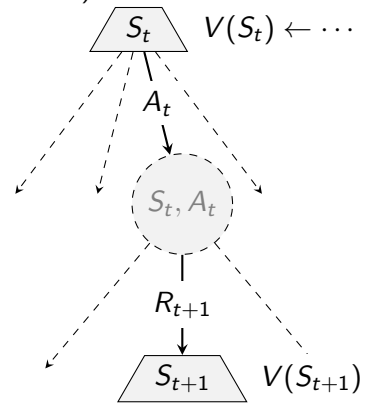
Notes

The S_t, A_t chance node is depicted in gray purposively. It is not directly observable, it is our model of uncertain outcome of action A_t taken in state S_t .

Recap: V - learning

Learn $V(s)$ values as the robot/agent goes (temporal difference).

- ▶ time t , at S_t
- ▶ select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
- ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma V(S_{t+1})$
- ▶ α temporal difference update
 $V(S_t) \leftarrow V(S_t) + \alpha(\text{sample} - V(S_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)



Notes

The S_t, A_t chance node is depicted in gray purposively. It is not directly observable, it is our model of uncertain outcome of action A_t taken in state S_t .

Recap: V - values, converged ...

$\gamma = 1$, rewards $-1, +10, -10$, and deterministic robot

7.00	8.00	9.00	10.00
6.00		8.00	-10.00
5.00	6.00	7.00	6.00

$$V(S_t) = R_{t+1} + V(S_{t+1})$$

12 / 37

Notes

$\gamma = 1$, Rewards $-1, +10, -10$, and no uncertainty on the outcome of actions – deterministic robot/agent. Rewards associated with leaving the state. Q values close next to terminal state includes the actual reward and the transition cost stepping in, or better, leaving the last living state.

How would the Values change if $\gamma = 0.9$?

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?

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

$$\triangleright \pi(s) = \arg \max_a Q(s, a)$$

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?

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

$$\blacktriangleright \pi(s) = \arg \max_a Q(s, a)$$

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?

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

$$\blacktriangleright \pi(s) = \arg \max_a Q(s, a)$$

Notes

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

Model-free TD learning, updating after each transition

- ▶ Observe, experience environment through learning episodes, collecting:

$$S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, \dots$$

- ▶ Update by mimicking Bellman updates after each transition $(S_t, A_t, R_{t+1}, S_{t+1})$

Notes

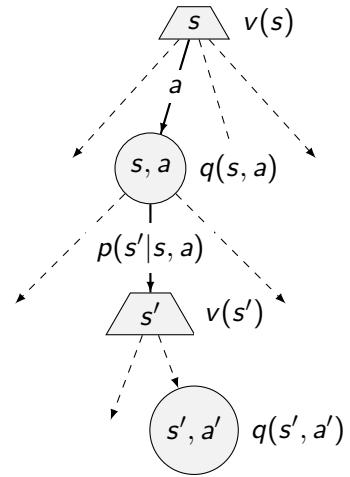
Think about $S_t - A_t - S_{t+1} - A_{t+1} - S_{t+2}$ tree with associated rewards. Episode starts in a start state and ends in a terminal state.

Recap: Bellman optimality equations for $v(s)$ and $q(s, a)$

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

The value of a q -state (s, a) :

$$\begin{aligned}q(s, a) &= \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')] \\ &= \sum_{s'} p(s'|s, a) \left[r(s, a, s') + \gamma \max_{a'} q(s', a') \right]\end{aligned}$$



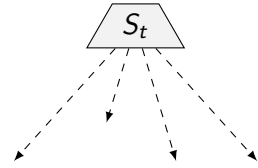
Notes

The tree continues from s' through a' and so on until it terminates.

Q-learning (off-policy TD control)

Learn policy (Q-values) as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ▶ time t , at S_t
 - ▶ select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
 - ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
 - ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
 - ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)



In each step Q directly approximates the optimal q^* function (learns optimal policy).

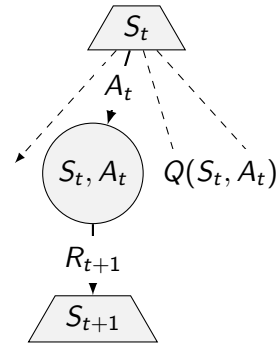
Notes

Why **off-policy** ? It learns Q-values answering “What would be this action worth in this state, assuming abandoning current policy and start according the policy that chooses the best action”

Q-learning (off-policy TD control)

Learn policy (Q-values) as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ▶ time t , at S_t
- ▶ select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
 - ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
 - ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
 - ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)



In each step Q directly approximates the optimal q^* function (learns optimal policy).

Notes

16 / 37

Why **off-policy** ? It learns Q-values answering "What would be this action worth in this state, assuming abandoning current policy and start according to the policy that chooses the best action"

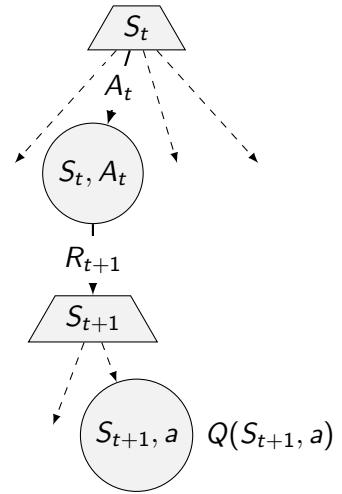
Q-learning (off-policy TD control)

Learn policy (Q-values) as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ▶ time t , at S_t
- ▶ select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
- ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$

- ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)

In each step Q directly approximates the optimal q^* function (learns optimal policy).



Notes

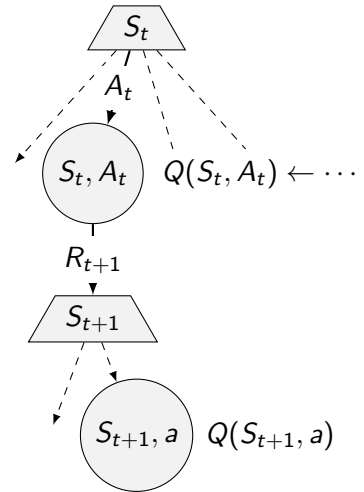
Why **off-policy** ? It learns Q-values answering "What would be this action worth in this state, assuming abandoning current policy and start according to the policy that chooses the best action"

Q-learning (off-policy TD control)

Learn policy (Q-values) as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ▶ time t , at S_t
- ▶ select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
- ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
 - ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)

In each step Q directly approximates the optimal q^* function (learns optimal policy).



Notes

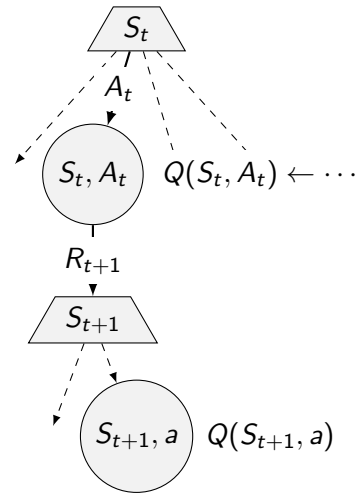
Why **off-policy** ? It learns Q-values answering "What would be this action worth in this state, assuming abandoning current policy and start according to the policy that chooses the best action"

Q-learning (off-policy TD control)

Learn policy (Q-values) as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ▶ time t , at S_t
- ▶ select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
- ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)

In each step Q directly approximates the optimal q^* function (learns optimal policy).



Notes

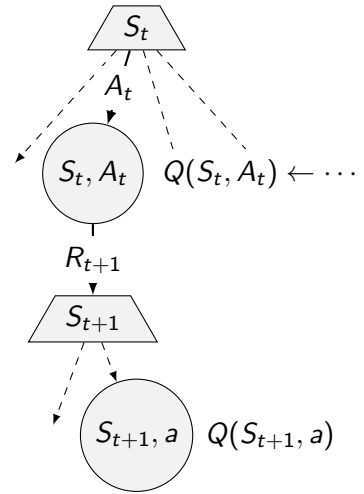
Why **off-policy** ? It learns Q-values answering "What would be this action worth in this state, assuming abandoning current policy and start according to the policy that chooses the best action"

Q-learning (off-policy TD control)

Learn policy (Q-values) as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ▶ time t , at S_t
- ▶ select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
- ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)

In each step Q directly approximates the optimal q^* function (learns optimal policy).



Notes

Why **off-policy** ? It learns Q-values answering "What would be this action worth in this state, assuming abandoning current policy and start according to the policy that chooses the best action"

Recap: V - and Q - values, converged ...

$\gamma = 1$, rewards $-1, +10, -10$, and deterministic robot

7.00	8.00	9.00	10.00	<table border="1"> <tr> <td>6.00</td> <td>7.00</td> <td>8.00</td> <td>0.00</td> </tr> <tr> <td>6.00</td> <td>7.00</td> <td>6.00</td> <td>8.00</td> </tr> <tr> <td>5.00</td> <td>7.00</td> <td>7.00</td> <td>0.00</td> </tr> </table>	6.00	7.00	8.00	0.00	6.00	7.00	6.00	8.00	5.00	7.00	7.00	0.00
6.00	7.00	8.00	0.00													
6.00	7.00	6.00	8.00													
5.00	7.00	7.00	0.00													
6.00		8.00	-10.00	<table border="1"> <tr> <td>6.00</td> <td></td> <td>8.00</td> <td>0.00</td> </tr> <tr> <td>5.00</td> <td>5.00</td> <td>7.00</td> <td>-11.00</td> </tr> <tr> <td>4.00</td> <td></td> <td>6.00</td> <td>0.00</td> </tr> </table>	6.00		8.00	0.00	5.00	5.00	7.00	-11.00	4.00		6.00	0.00
6.00		8.00	0.00													
5.00	5.00	7.00	-11.00													
4.00		6.00	0.00													
5.00	6.00	7.00	6.00	<table border="1"> <tr> <td>5.00</td> <td>5.00</td> <td>7.00</td> <td>-11.00</td> </tr> <tr> <td>4.00</td> <td>5.00</td> <td>4.00</td> <td>6.00</td> </tr> <tr> <td>4.00</td> <td>5.00</td> <td>5.00</td> <td>5.00</td> </tr> </table>	5.00	5.00	7.00	-11.00	4.00	5.00	4.00	6.00	4.00	5.00	5.00	5.00
5.00	5.00	7.00	-11.00													
4.00	5.00	4.00	6.00													
4.00	5.00	5.00	5.00													

$$V(S_t) = R_{t+1} + V(S_{t+1})$$

$$Q(S_t, A_t) = R_{t+1} + \max_a Q(S_{t+1}, a)$$

17 / 37

Notes

$\gamma = 1$, Rewards $-1, +10, -10$, and no uncertainty on the outcome of actions – deterministic robot/agent. Rewards associated with leaving the state. Q values close next to terminal state includes the actual reward and the transition cost stepping in, or better, leaving the last living state.

$Q(s, a)$ – expected sum of rewards having taken the action and acting according to the (optimal) policy.

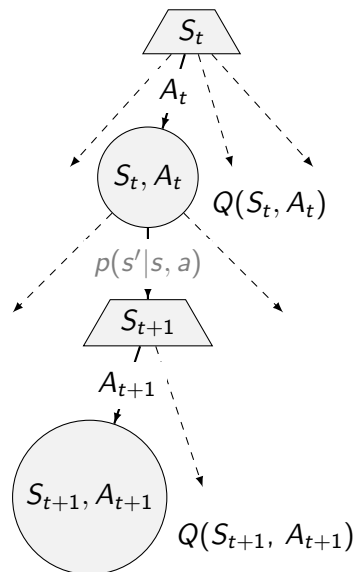
How would the (q)values change if $\gamma = 0.9$?

Sarsa (on-policy TD control)

Learn Q values as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ▶ time t , at S_t , select $A_t \in \mathcal{A}(S_t)$
- ▶ take A_t , observe R_{t+1}, S_{t+1}
 - ▶ select $A_{t+1} \in \mathcal{A}(S_{t+1})$
 - ▶ which gives sample estimate
sample = $R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$
 - ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
 - ▶ $S_t \leftarrow S_{t+1}, A_t \leftarrow A_{t+1}$ and repeat (unless S_t is terminal)

In each step learns Q .



Notes

on-policy : It learns Q answering “What would be this action worth in this state, assuming I stick with my policy?” SARSA – State-Action-Reward-State-Action.

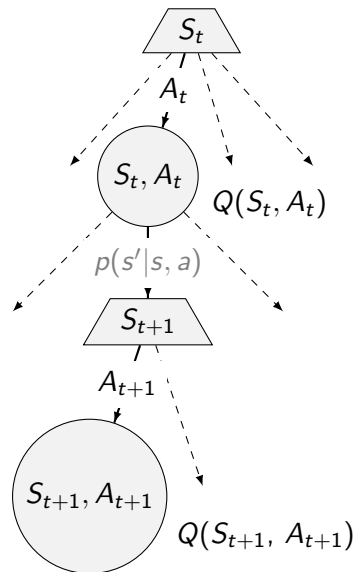
- Q-learning: Learning from 4-tuples $S_t, A_t, R_{t+1}, S_{t+1}$.
- SARSA: Learning from 5-tuples $S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}$.

Sarsa (on-policy TD control)

Learn Q values as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ▶ time t , at S_t , select $A_t \in \mathcal{A}(S_t)$
- ▶ take A_t , observe R_{t+1}, S_{t+1}
- ▶ **select** $A_{t+1} \in \mathcal{A}(S_{t+1})$
 - ▶ which gives sample estimate
sample = $R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$
 - ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
 - ▶ $S_t \leftarrow S_{t+1}, A_t \leftarrow A_{t+1}$ and repeat (unless S_t is terminal)

In each step learns Q .



Notes

on-policy : It learns Q answering “What would be this action worth in this state, assuming I stick with my policy?” SARSA – State-Action-Reward-State-Action.

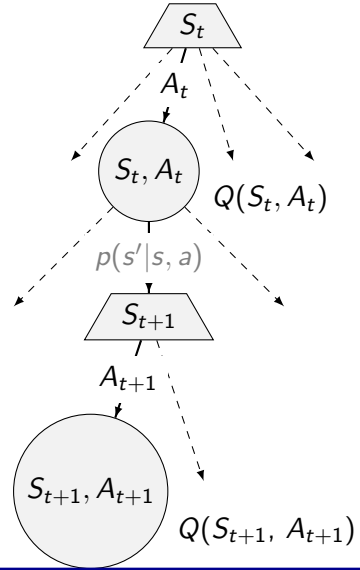
- Q-learning: Learning from 4-tuples $S_t, A_t, R_{t+1}, S_{t+1}$.
- SARSA: Learning from 5-tuples $S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}$.

Sarsa (on-policy TD control)

Learn Q values as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ▶ time t , at S_t , select $A_t \in \mathcal{A}(S_t)$
- ▶ take A_t , observe R_{t+1}, S_{t+1}
- ▶ **select** $A_{t+1} \in \mathcal{A}(S_{t+1})$
- ▶ which gives sample estimate
sample = $R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$
- ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
- ▶ $S_t \leftarrow S_{t+1}, A_t \leftarrow A_{t+1}$ and repeat (unless S_t is terminal)

In each step learns Q .



18 / 37

Notes

on-policy : It learns Q answering "What would be this action worth in this state, assuming I stick with my policy?" SARSA – State-Action-Reward-State-Action.

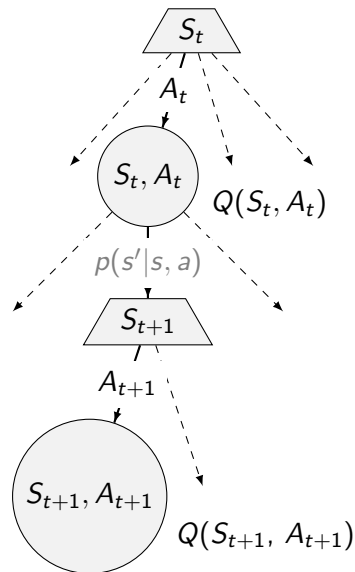
- Q-learning: Learning from 4-tuples $S_t, A_t, R_{t+1}, S_{t+1}$.
- SARSA: Learning from 5-tuples $S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}$.

Sarsa (on-policy TD control)

Learn Q values as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ▶ time t , at S_t , select $A_t \in \mathcal{A}(S_t)$
- ▶ take A_t , observe R_{t+1}, S_{t+1}
- ▶ **select** $A_{t+1} \in \mathcal{A}(S_{t+1})$
- ▶ which gives sample estimate
sample = $R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$
- ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
- ▶ $S_t \leftarrow S_{t+1}, A_t \leftarrow A_{t+1}$ and repeat (unless S_t is terminal)

In each step learns Q .



Notes

on-policy : It learns Q answering “What would be this action worth in this state, assuming I stick with my policy?” SARSA – State-Action-Reward-State-Action.

- Q-learning: Learning from 4-tuples $S_t, A_t, R_{t+1}, S_{t+1}$.
- SARSA: Learning from 5-tuples $S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}$.

Q-learning: algorithm

step size $0 < \alpha \leq 1$

initialize $Q(s, a)$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

repeat episodes:

 initialize S

for each step of episode **do**

 choose A from $\mathcal{A}(S)$

 take action A , observe R, S'

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

 until S is terminal

until Time is up, ...

Sarsa: algorithm

step size $0 < \alpha \leq 1$

initialize $Q(s, a)$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

repeat episodes:

 initialize S

 choose A from $\mathcal{A}(S)$

for each step of episode: **do**

 take action A , observe R, S'

 choose A' from $\mathcal{A}(S')$

$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S', A \leftarrow A'$

 until S is terminal

until Time is up, ...

How to select A_t in S_t ? What policy?

- ▶ time t , at S_t
- ▶ take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
- ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)

How to select A_t in S_t ? What policy?

- ▶ time t , at S_t
- ▶ take A_t derived from Q , observe R_{t+1}, S_{t+1}
- ▶ compute sample estimate at time t
sample = $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
- ▶ α temporal difference update
 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)

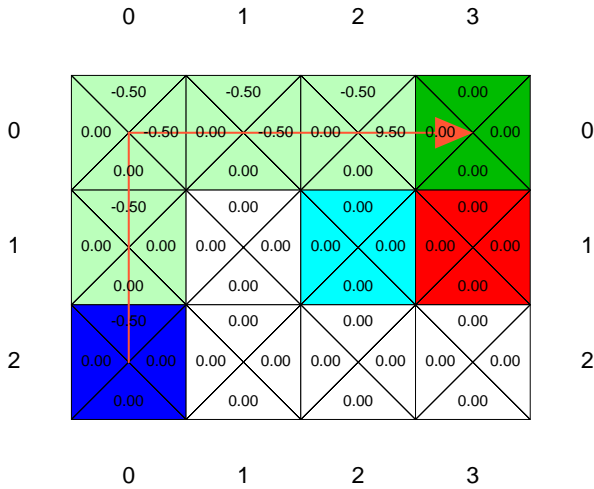
... A_t derived from Q

What about keeping optimality, taking max?

$$A_t = \arg \max_a Q(S_t, a)$$

see the demo run of `rl_agents.py`.

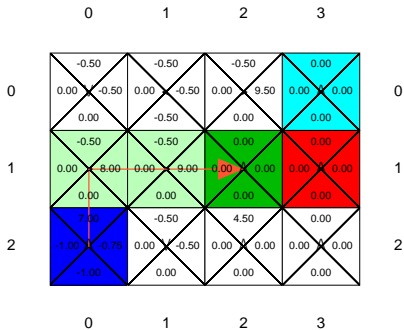
Two good goal states



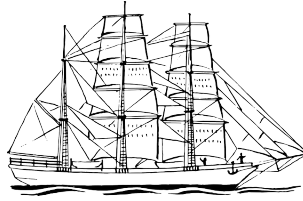
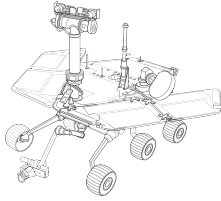
Notes

Discuss the on-line demo with two good goal states. $\gamma = 1, \alpha = 0.5$, Living reward -1 , $R(1, 2) = 10, R(0, 3) = 20, R(1, 1) = -10$. Taking the action, corresponding the max Q . If equal options, than in the 0, 1, 2, 3 action order. 50 training episodes. What happened?

- No exploration step: <https://youtu.be/Y5yLttbkPMM>
- Exploring steps involved (will be talking in a few minutes): https://youtu.be/cAr-lrawF_c



Exploration vs Exploitation



- ▶ 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 new one?
- ▶ Go to bussiness or study a demanding program?
- ▶ ...

How to explore?

Random (ϵ -greedy):

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

Problems with randomness?

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

Notes

We can think about lowering ϵ as the learning progresses.

How to explore?

Random (ϵ -greedy):

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

Problems with randomness?

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

Notes

We can think about lowering ϵ as the learning progresses.

How to explore?

Random (ϵ -greedy):

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

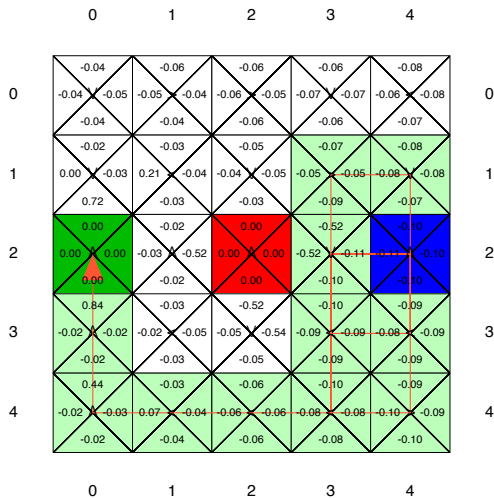
Problems with randomness?

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

Notes

We can think about lowering ϵ as the learning progresses.

How to evaluate the result? When to stop learning?



- ▶ What is the actual result of q-learning?
- ▶ How to evaluate it?
- ▶ When to stop learning?

Notes

- Accurate estimation of q-values in every state, which gives:
 1. Estimate of the sum of expected future rewards from every state.
 2. Policy – which action to take in every state. Simply a max over q-values.
- Note that the policy can be poorly estimated in areas less visited.
- Evaluating learning progress. Recall (Q-learning, but similar for SARSA):
 - $\text{sample} = R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$
 - α temporal difference update
 - $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} - Q(S_t, A_t))$
 - learning progress / convergence: $(\text{sample} - Q(S_t, A_t)) \rightarrow 0$
- When to stop learning
 - The more learning the better - think about visiting all places/states.
 - Never... If you can afford it, keep learning with a small learning rate...
 - Note: the learning rate is embedded in two parameters: α and ϵ

Going beyond tables – generalizing across states

	0	1	2	3	4	
0	0.84	0.88	0.92	0.96	1.00	0
	0	1	2	3	4	

27 / 37

Notes

We were talking about v - and q - functions but what was the representation? (look-up) Tables. Looking at $v(s)$, we need a table for each of the state!

Btw., we always visualize RL on the grid but note that the agent does not know about the topology of the world. It only knows about q -values! Even in model-based RL: One learns also the transition functions, but these are still do not give a map of the state space. With stochastic action outcomes, the agent can at most estimate how the world looks like. This knowledge is, however, irrelevant for the algorithms we study. This world is small, but think bigger!

Going beyond tables – generalizing across states

	0	1	2	3	4	
0	0.84	0.80	0.76	0.72	0.68	0
1	0.88	0.84	0.80	0.76	0.72	1
2	0.92	0.88	0.84	0.80	0.76	2
3	0.96	0.92	0.88	0.84	0.80	3
4	1.00	0.96	0.92	0.88	0.84	4
	0	1	2	3	4	

Notes

Looking a $V(s)$, we need a table for each of the states! This world is small, but think bigger!

$v(s)$ not as a table but as an approximation function $\hat{v}(s, \mathbf{w})$

	0	1	2	3	4	
0	0.84	0.88	0.92	0.96	1.00	0
	0	1	2	3	4	

$$\hat{v}(s, \mathbf{w}) = w_0 + w_1 s$$

What are w_0, w_1 equal to?

Instead of the complete table, only 2 parameters to learn $\mathbf{w} = [w_0, w_1]^T$

29 / 37

Notes

Note: we can approximate $v(s)$ or $q(s, a)$.

Two key benefits as opposed to keeping the table with $v(s) / q(s, a)$ for every state:

- Space complexity – for large worlds, the table simply won't fit in memory.
- **Generalization.**
 - Pacman example (UC Berkeley, Lecture 11 Reinforcement Learning II). Running away from the ghost is what matters, not running away only when pacman is in state (5,7) and ghost is in position (7,8)...
 - Tracked robot – obstacle type not its detailed shape/coordinates determine the action...

What are w_0, w_1 equal to?, we can start from left, target is the true $v(s = 0) = 0.84$, next target is $v(s = 1) = 0.88$, ...

Note about notation. Bold lower case letters are used to denote vectors. Vectors are always considered oriented column-wise unless explicitly stated otherwise.

Linear value functions

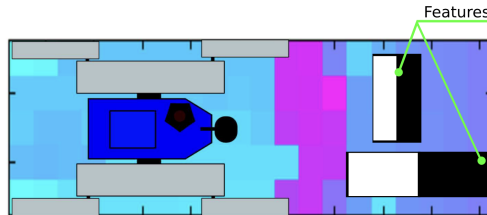
State $s \in \mathcal{S} \subset \mathbb{R}^n$ concatenates:

- ◆ Proprioceptive measurements: roll, pitch, torques, velocity, acceleration.

7.00	8.00	9.00	10.00
6.00		8.00	-10.00
5.00	6.00	7.00	6.00

- ◆ Local exteroceptive measurements.

features on digital elevation map
with fixed size.



$$\hat{v}(s, \mathbf{w}) = w_1 f_1(s) + w_2 f_2(s) + w_3 f_3(s) + \dots + w_n f_n(s)$$
$$\hat{q}(s, a, \mathbf{w}) = w_1 f_1(s, a) + w_2 f_2(s, a) + w_3 f_3(s, a) + \dots + w_n f_n(s, a)$$

30 / 37

Notes

Note: the state description/encoding has to contain useful information about the context.

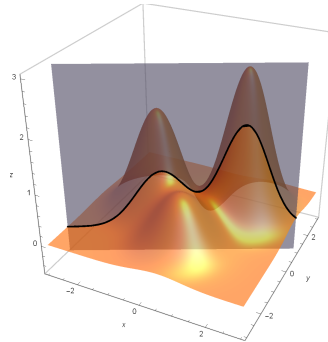
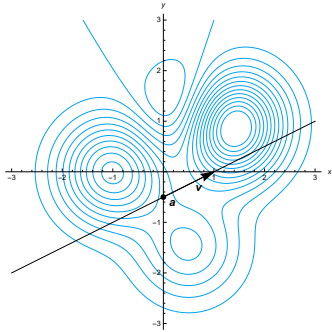
What could be the f functions for the grid world?

- Coordinates: (x,y) ? Probably not.
- Pit is ahead... These kind of things would be useful, but currently not part of the state description...

Obviously, when data are available, we can fit. How to do it on-line?

Směrová a parciální derivace (a stolen slide)

- Ať $f : D \subseteq \mathbb{R}^2 \rightarrow \mathbb{R}$ přiřazuje bodům na mapě D nadmořskou výšku.
- V mapě se vydáme z bodu a rovnoměrně přímočaře rychlostí v . Jaká bude okamžitá změna nadmořské výšky v bodě a ?



Learning \mathbf{w} by Stochastic Gradient Descent (SGD)

- ▶ assume $\hat{v}(s, \mathbf{w})$ differentiable in all states
- ▶ we update \mathbf{w} in discrete time steps t
- ▶ in each step S_t we observe a new example of (true) $v^\pi(S_t)$
- ▶ $\hat{v}(S_t, \mathbf{w})$ is an approximator \rightarrow error = $v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)$
 - ▶ minimize squared error (function, depending on \mathbf{w}) $[v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)]^2$

Shifting \mathbf{w} along the gradient:

$$\begin{aligned}\mathbf{w}_{t+1} &\doteq \mathbf{w}_t - \frac{1}{2} \alpha \nabla [v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)]^2 \\ &= \mathbf{w}_t + \alpha [v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w}_t)\end{aligned}$$

$$\nabla f(\mathbf{w}) \doteq \left[\frac{\partial f(\mathbf{w})}{\partial w_1}, \frac{\partial f(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial f(\mathbf{w})}{\partial w_d} \right]^\top$$

32 / 37

Notes

Gradient descent - all samples are known, Stochastic GD - update after each sample
 α is a kind of damping factor, convergence of SGD requires that α decreases over time.

$\hat{v}(s, \mathbf{w})$ could be quite complex, e.g. a Multi Layer Perceptron (MLP), Deep Network, and \mathbf{w} represents the weights. See, e.g.

- <https://skymind.ai/wiki/deep-reinforcement-learning>
- Vision for robotics course you may take next term. <https://cw.fel.cvut.cz/wiki/courses/b3b33vir/start>

Learning \mathbf{w} by Stochastic Gradient Descent (SGD)

- ▶ assume $\hat{v}(s, \mathbf{w})$ differentiable in all states
- ▶ we update \mathbf{w} in discrete time steps t
- ▶ in each step S_t we observe a new example of (true) $v^\pi(S_t)$
- ▶ $\hat{v}(S_t, \mathbf{w})$ is an approximator \rightarrow error = $v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)$
- ▶ minimize squared error (function, depending on \mathbf{w}) $\left[v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t) \right]^2$

Shifting \mathbf{w} along the gradient:

$$\begin{aligned} \mathbf{w}_{t+1} &\doteq \mathbf{w}_t - \frac{1}{2} \alpha \nabla \left[v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t) \right]^2 \\ &= \mathbf{w}_t + \alpha \left[v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t) \right] \nabla \hat{v}(S_t, \mathbf{w}_t) \end{aligned}$$

$$\nabla f(\mathbf{w}) \doteq \left[\frac{\partial f(\mathbf{w})}{\partial w_1}, \frac{\partial f(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial f(\mathbf{w})}{\partial w_d} \right]^\top$$

32 / 37

Notes

Gradient descent - all samples are known, Stochastic GD - update after each sample
 α is a kind of damping factor, convergence of SGD requires that α decreases over time.

$\hat{v}(s, \mathbf{w})$ could be quite complex, e.g. a Multi Layer Perceptron (MLP), Deep Network, and \mathbf{w} represents the weights. See, e.g.

- <https://skymind.ai/wiki/deep-reinforcement-learning>
- Vision for robotics course you may take next term. <https://cw.fel.cvut.cz/wiki/courses/b3b33vir/start>

Learning \mathbf{w} by Stochastic Gradient Descent (SGD)

- ▶ assume $\hat{v}(s, \mathbf{w})$ differentiable in all states
- ▶ we update \mathbf{w} in discrete time steps t
- ▶ in each step S_t we observe a new example of (true) $v^\pi(S_t)$
- ▶ $\hat{v}(S_t, \mathbf{w})$ is an approximator \rightarrow error = $v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)$
- ▶ minimize squared error (function, depending on \mathbf{w}) $\left[v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t) \right]^2$

Shifting \mathbf{w} along the gradient:

$$\begin{aligned}\mathbf{w}_{t+1} &\doteq \mathbf{w}_t - \frac{1}{2}\alpha \nabla \left[v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t) \right]^2 \\ &= \mathbf{w}_t + \alpha \left[v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t) \right] \nabla \hat{v}(S_t, \mathbf{w}_t)\end{aligned}$$

$$\nabla f(\mathbf{w}) \doteq \left[\frac{\partial f(\mathbf{w})}{\partial w_1}, \frac{\partial f(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial f(\mathbf{w})}{\partial w_d} \right]^\top$$

32 / 37

Notes

Gradient descent - all samples are known, Stochastic GD - update after each sample
 α is a kind of damping factor, convergence of SGD requires that α decreases over time.

$\hat{v}(s, \mathbf{w})$ could be quite complex, e.g. a Multi Layer Perceptron (MLP), Deep Network, and \mathbf{w} represents the weights. See, e.g.

- <https://skymind.ai/wiki/deep-reinforcement-learning>
- Vision for robotics course you may take next term. <https://cw.fel.cvut.cz/wiki/courses/b3b33vir/start>

Approximate Q-learning (of a linear combination)

$$\hat{q}(s, a, \mathbf{w}) = w_1 f_1(s, a) + w_2 f_2(s, a) + w_3 f_3(s, a) + \dots + w_n f_n(s, a)$$

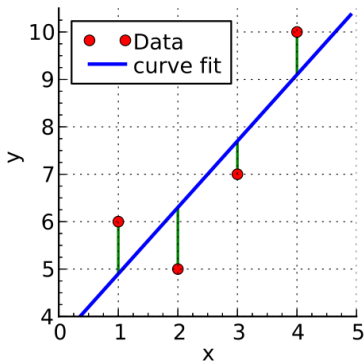
- ▶ transition = $S_t, A_t, R_{t+1}, S_{t+1}$
- ▶ sample $R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a, \mathbf{w}_t)$
- ▶ diff = $\left[R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a, \mathbf{w}_t) \right] - \hat{q}(S_t, A_t, \mathbf{w}_t)$
- ▶ Update: $\mathbf{w} = [w_1, w_2, \dots, w_d]^\top$
from previous slide we know that $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha \left[v^\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t) \right] \nabla \hat{v}(S_t, \mathbf{w}_t)$
and $\hat{q}(s, a, \mathbf{w})$ is linear in \mathbf{w}
 $w_i \leftarrow w_i + \alpha [\text{diff}] f_i(S_t, A_t)$

33 / 37

Notes

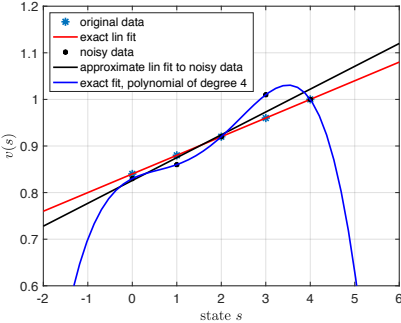
- We are minimizing the error at the point where we measure (sample).
- However, we know we only approximate.
- α is a kind of damping factor; convergence of SGD requires that α decreases over time.

How is it possible at all? On-line least squares!



How to design the q-function? Overfitting ...

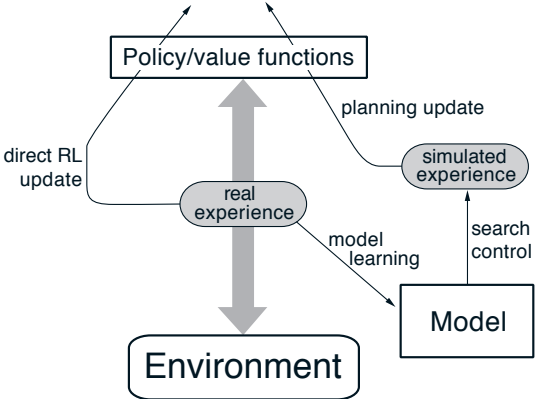
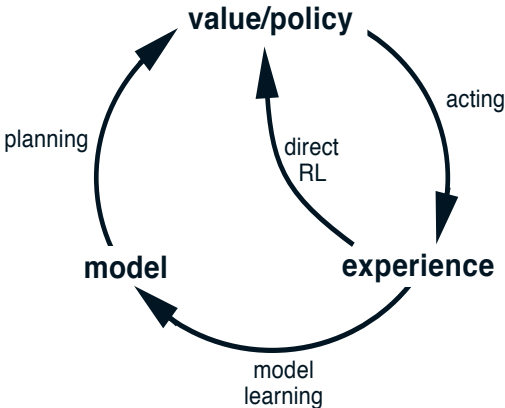
	0	1	2	3	4	
0	0.84	0.88	0.92	0.96	1.00	0
	0	1	2	3	4	



Notes

See the `fitdemo.m` run, higher degree polynomials perfectly fits, but poorly generalizes outside the range

Going beyond - Dyna-Q integration planning, acting, learning



2

²Schemes from [3]

References I

Further reading: Chapter 21 of [1] (chapter 23 of [2]). More detailed discussion in [3] Chapters 6 and 9. You can read about strategies for exploratory moves at various places, [Tensor Flow related](#)³. More RL URLs at the [course pages](#)⁴.

[1] Stuart Russell and Peter Norvig.
Artificial Intelligence: A Modern Approach.
Prentice Hall, 3rd edition, 2010.
<http://aima.cs.berkeley.edu/>.

[2] Stuart Russell and Peter Norvig.
Artificial Intelligence: A Modern Approach.
Prentice Hall, 4th edition, 2021.

References II

- [3] 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>.

³[https://medium.com/emergent-future/
simple-reinforcement-learning-with-tensorflow-part-7-action-selection-strategies-for-exploration-d3a97b7ccef](https://medium.com/emergent-future/simple-reinforcement-learning-with-tensorflow-part-7-action-selection-strategies-for-exploration-d3a97b7ccef)
⁴[https:
//cw.fel.cvut.cz/wiki/courses/b3b33kui/cviceni/program_po_tydnech/tyden_09#reinforcement_learning_plus](https://cw.fel.cvut.cz/wiki/courses/b3b33kui/cviceni/program_po_tydnech/tyden_09#reinforcement_learning_plus)

37 / 37