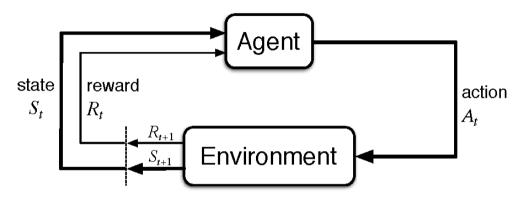
Reinforcement learning II Active learning

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 11, 2025

Recap: Reinforcement Learning



- ► 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.

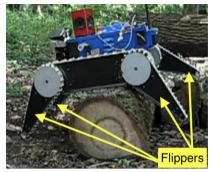
¹Scheme from [3]

Learning to control flippers



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

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





- ◆ 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)

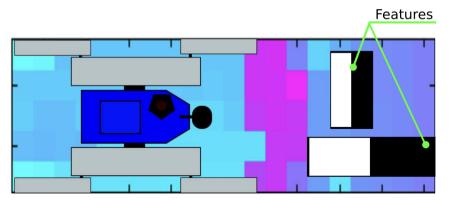
difficult to control all of them manually!

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

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

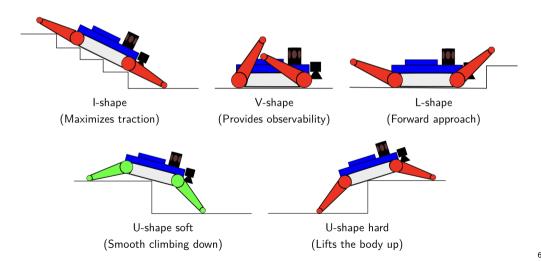
Local exteroceptive measurements.

features on digital elevation map with fixed size.



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

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



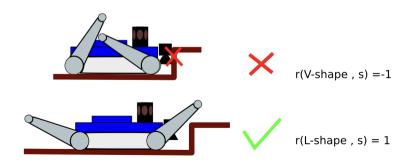
Reward $r(a, \mathbf{s}) : \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ is a weighted sum of following contributions:

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



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, $a \in A(s)$
- ▶ 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.

(Transition) Model-based learning

The main idea: Do something and:

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

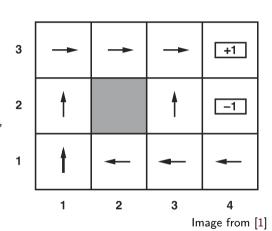
Learning MDP model:

- ▶ 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.

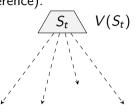
Solve the learned MDP.

Model-free learning

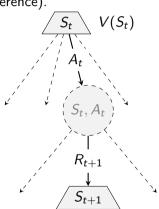
- ightharpoonup r, p not known.
- ► Move around, observe.
- And learn on the way.
- ▶ **Goal:** Learn the state value v(s), or (better), q-value q(s, a) functions.



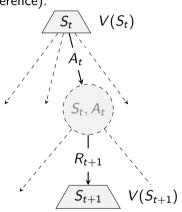
- ightharpoonup time t, at S_t
- lacktriangle select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
- compute sample estimate at time tsample = $R_{t+1} + \gamma V(S_{t+1})$
- ho α temporal difference update $V(S_t) \leftarrow V(S_t) + \alpha(\text{sample} V(S_t))$
- $ightharpoonup S_t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



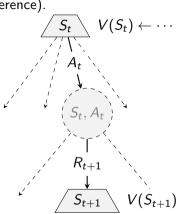
- ightharpoonup time t, at S_t
- ▶ select and take $A_t \in A(S_t)$, observe R_{t+1}, S_{t+1}
- compute sample estimate at time t sample = $R_{t+1} + \gamma V(S_{t+1})$
- $\sim \alpha$ temporal difference update $V(S_t) \leftarrow V(S_t) + \alpha(\text{sample} V(S_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and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



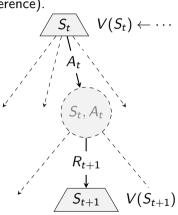
- \triangleright time t, at S_t
- ▶ select and take $A_t \in A(S_t)$, observe R_{t+1}, S_{t+1}
- compute sample estimate at time t sample = $R_{t+1} + \gamma V(S_{t+1})$
- ▶ α temporal difference update $V(S_t) \leftarrow V(S_t) + \alpha(\text{sample} V(S_t))$
- \triangleright $S_t \leftarrow S_{t+1}$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



- \triangleright time t, at S_t
- ▶ select and take $A_t \in A(S_t)$, 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 $V(S_t) \leftarrow V(S_t) + \alpha (\text{sample} V(S_t))$
- $ightharpoonup S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)



- ightharpoonup time t, at S_t
- ▶ select and take $A_t \in A(S_t)$, 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)



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

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

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) = \underset{a}{\operatorname{arg max}} \sum_{s'} p(s' \mid s, a) \left[ r(s, a, s') + \gamma V(s') \right]
\pi(s) = \underset{a}{\operatorname{arg max}} Q(s, a)
```

What is wrong with the temporal difference Value learning?

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

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

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

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

What is wrong with the temporal difference Value learning?

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

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

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

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

Model-free TD learning, updating after each transition

► Observe, experience environment through learning episodes, collecting:

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

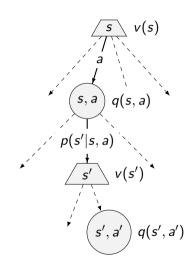
▶ Update by mimicking Bellman updates after each transition $(S_t, A_t, R_{t+1}, S_{t+1})$

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

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

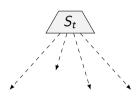
The value of a q-state (s, a):

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



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

- ightharpoonup time t, at S_t
- ightharpoonup 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 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\text{sample} Q(S_t, A_t))$
- $ightharpoonup S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)

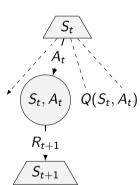


Learn policy (Q-values) as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ightharpoonup time t, at S_t
- lacktriangle select and take $A_t \in \mathcal{A}(S_t)$, observe R_{t+1}, S_{t+1}
- compute sample estimate at time t sample = $R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a)$
- ightharpoonup lpha temporal difference update

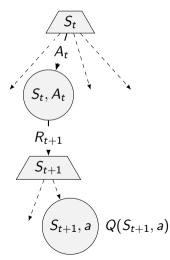
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\mathsf{sample} - Q(S_t, A_t))$$

 \triangleright $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal



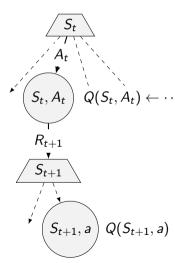
Learn policy (Q-values) as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

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



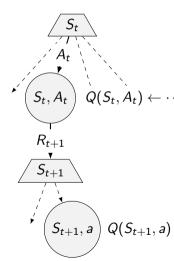
Learn policy (Q-values) as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ightharpoonup time t, at S_t
- ▶ select and take $A_t \in A(S_t)$, 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 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\mathsf{sample} Q(S_t, A_t))$
- \triangleright $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)



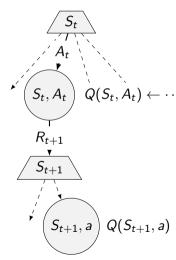
Learn policy (Q-values) as the robot/agent goes (temporal difference). If some Q quantity not known, initialize.

- ightharpoonup time t, at S_t
- ▶ select and take $A_t \in A(S_t)$, 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 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\mathsf{sample} Q(S_t, A_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)



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

- ightharpoonup 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 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\mathsf{sample} Q(S_t, A_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)



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

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

	, .	•		
7.00	8.00	9.00	10.00	6.00 7.00 8.00 0.00 6.00 7.00 6.00 8.00 7.00 9.00 0.00 0.00 5.00 7.00 7.00 0.00
6.00		8.00	-10.00	6.00 5.00 5.00 6.00 6.00 6.00 6.00 6.00
5.00	6.00	7.00	6.00	5.00 5.00 7.00 -11.00 4.00 5.00 6.00 5.00 5.00 6.00 5.00

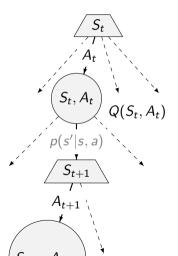
$$V(S_t) = R_{t+1} + V(S_{t+1})$$

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

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

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

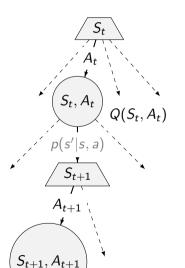
In each step learns Q.



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

- ▶ time t, at S_t , select $A_t \in \mathcal{A}(S_t)$
- ightharpoonup take A_t , observe R_{t+1}, S_{t+1}
- ▶ select $A_{t+1} \in A(S_{t+1})$
- which gives sample estimate sample = $R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$
- α temporal difference update $\alpha(S, A_1) + \alpha(S, A_2) + \alpha(S, A_3)$
- ▶ $S_t \leftarrow S_{t+1}$, $A_t \leftarrow A_{t+1}$ and repeat (unless S_t is terminal)

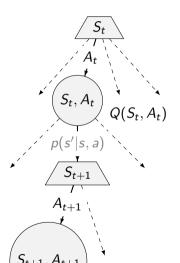
In each step learns Q.



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

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

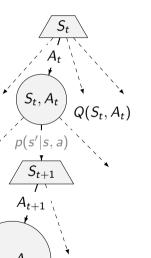
In each step learns Q



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

- \blacktriangleright time t, at S_t , select $A_t \in \mathcal{A}(S_t)$
- ightharpoonup take A_t , observe R_{t+1}, S_{t+1}
- ▶ select $A_{t+1} \in A(S_{t+1})$
- which gives sample estimate
- sample = $R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$ $\triangleright \alpha$ 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.



Q-learning: algorithm

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

Sarsa: algorithm

```
step size 0 < \alpha < 1
initialize Q(s, a) for all s \in \mathcal{S}, a \in \mathcal{A}(s)
repeat episodes:
    initialize S
    choose A from A(S)
    for each step of episode: do
        take action A. observe R, S'
        choose A' from 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?

- \triangleright 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 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\mathsf{sample} Q(S_t, A_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)

How to select A_t in S_t ? What policy?

- \triangleright 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 $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(\mathsf{sample} Q(S_t, A_t))$
- ▶ $S_t \leftarrow S_{t+1}$ and repeat (unless S_t is terminal)

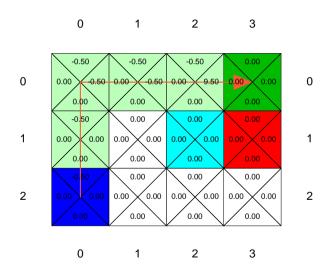
 $...A_t$ derived from Q

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



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.
- \blacktriangleright With probability ϵ , act randomly.
- ▶ With probability 1ϵ , use the policy.

Problems with randomness?

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

How to explore?

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?

How to explore?

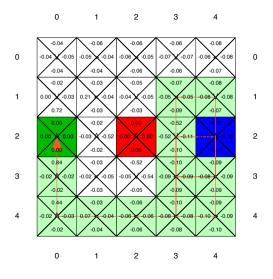
Random (ϵ -greedy):

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

Problems with randomness?

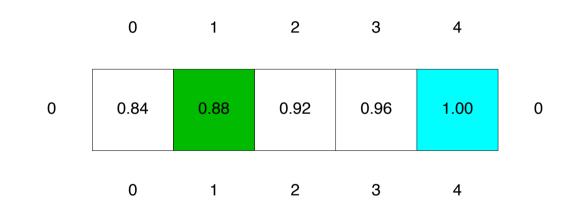
- ► Keeps exploring forever.
- ▶ Should we keep ϵ fixed (over learning)?
- ightharpoonup ϵ same everywhere?

How to evaluate the result? When to stop learning?

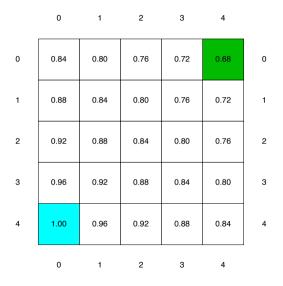


- ► What is the actual result of q-learning?
- How to evaluate it?
- When to stop learning?

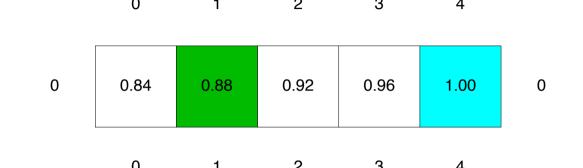
Going beyond tables – generalizing across states



Going beyond tables – generalizing across states



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



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

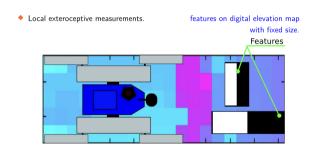
What are w_0, w_1 equal to? Instead of the complete table, only 2 parameters to learn $\mathbf{w} = [w_0, w_1]^{\top}$

Linear value functions

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

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

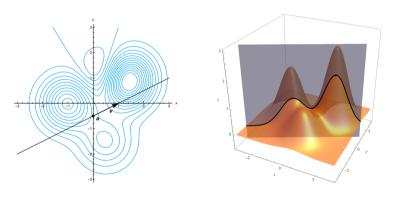




$$\hat{v}(s, \mathbf{w}) = w_1 f_1(s) + w_2 f_2(s) + w_3 f_3(s) + \dots + w_n f_n(s)
\hat{q}(s, a, \mathbf{w}) = w_1 f_1(s, a) + w_2 f_2(s, a) + w_3 f_3(s, a) + \dots + w_n f_n(s, a)$$

Směrová a parciální derivace (a stolen slide)

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



Martin Bohata Matematická analýza 2 Směrová a parciální derivace 2 / 22

Learning w by Stochastic Gradient Descent (SGD)

- **assume** $\hat{v}(s, \mathbf{w})$ differentiable in all states
- we update **w** in discrete time steps t
- ightharpoonup in each step S_t we observe a new example of (true) $v^{\pi}(S_t)$
- $\hat{m{v}}(S_t, {m{w}})$ is an approximator o error $= m{v}^\pi(S_t) \hat{m{v}}(S_t, {m{w}}_t)$

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_{t} - \frac{1}{2} \alpha \nabla \left[v^{\pi}(S_{t}) - \hat{v}(S_{t}, \mathbf{w}_{t}) \right]^{2}$$

$$= \mathbf{w}_{t} + \alpha \left[v^{\pi}(S_{t}) - \hat{v}(S_{t}, \mathbf{w}_{t}) \right] \nabla \hat{v}(S_{t}, \mathbf{w}_{t})$$

$$\nabla f(\mathbf{w}) \doteq \left[\frac{\partial f(\mathbf{w})}{\partial w_{1}}, \frac{\partial f(\mathbf{w})}{\partial w_{2}}, \cdots, \frac{\partial f(\mathbf{w})}{\partial w_{d}} \right]^{\top}$$

Learning w by Stochastic Gradient Descent (SGD)

- **assume** $\hat{v}(s, \mathbf{w})$ differentiable in all states
- ightharpoonup we update $m{f w}$ in discrete time steps t
- ightharpoonup in each step S_t we observe a new example of (true) $v^{\pi}(S_t)$
- $m{\hat{v}}(S_t, m{w})$ is an approximator o error $= v^\pi(S_t) \hat{v}(S_t, m{w}_t)$

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

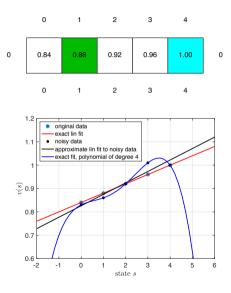
$$\nabla f(\mathbf{w}) \doteq \left[\frac{\partial f(\mathbf{w})}{\partial w_1}, \frac{\partial f(\mathbf{w})}{\partial w_2}, \cdots, \frac{\partial f(\mathbf{w})}{\partial w_d} \right]^{\top}$$

Approximate Q-learning (of a linear combination)

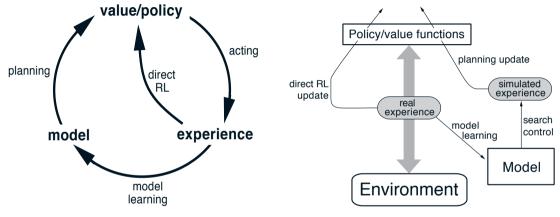
$$\hat{q}(s, a, \mathbf{w}) = w_1 f_1(s, a) + w_2 f_2(s, a) + w_3 f_3(s, a) + \cdots + w_n f_n(s, a)$$

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

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



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 related³. More RL URLs at the course pages⁴.

[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.

³https://medium.com/emergent-future/ simple-reinforcement-learning-with-tensorflow-part-7-action-selection-strategies-for-exploration-d3a97b7cceaf ⁴https:

 $^{//}cw.fel.cvut.cz/wiki/courses/b3b33kui/cviceni/program_po_tydnech/tyden_09\#reinforcement_learning_plus$