

DTU

02465: Introduction to reinforcement learning and control

Model-Free Control with tabular and linear methods

Tue Herlau

DTU Compute, Technical University of Denmark (DTU)

DTU Compute
Department of Applied Mathematics and Computer Science



Lecture Schedule

- | | |
|--|---|
| Dynamical programming <ul style="list-style-type: none"> ① The finite-horizon decision problem 7 February ② Dynamical Programming 14 February ③ DP reformulations and introduction to Control 21 February ④ Discretization and PID control 28 February ⑤ Direct methods and control by optimization 7 March ⑥ Linear-quadratic problems in control 14 March ⑦ Linearization and iterative LQR 21 March | Reinforcement learning <ul style="list-style-type: none"> ⑧ Exploration and Bandits 28 March ⑨ Bellmans equations and exact planning 4 April ⑩ Monte-carlo methods and TD learning 11 April ⑪ Model-Free Control with tabular and linear methods 25 April ⑫ Eligibility traces 2 May ⑬ Deep-Q learning 9 May |
|--|---|
- Syllabus: <https://02465material.pages.compute.dtu.dk/02465public>
Help improve lecture by giving feedback on DTU learn

2 DTU Compute

Lecture 11 25 April, 2025

Reading material:

- [SB18, Chapter 6.4-6.5; 7-7.2; 9-9.3; 10.1]

Learning Objectives

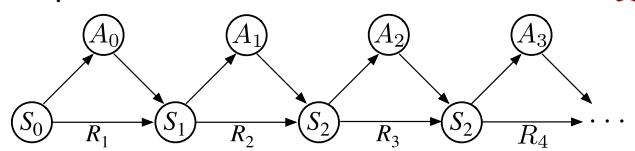
- Sarsa on-policy learning
- Q off-policy learning
- the n-step return
- value-function approximations and linear methods

3 DTU Compute

Lecture 11 25 April, 2025



Recap: First-Visit Monte-Carlo value estimation



We want to calculate the value function $v_\pi(s) = \mathbb{E}[G_t | S_t = s]$.
Simulate an episode of experience $s_0, a_0, r_1, s_1, a_1, r_2, \dots, r_T$ using π

- First step t we visit a state s
- Measure return $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$ for rest of the episode
- Estimate value function as $v_\pi(s_t) = \mathbb{E}[G_t | S_t = s] \approx \frac{1}{n} \sum_{i=1}^n G_t^{(n)}$
- The average can be computed incrementally:

$$V(s) \leftarrow V(s) + \frac{1}{n}(G_t - V(s))$$

- We use a fixed learning rate α

$$V(s) \leftarrow V(s) + \alpha(G_t - V(s))$$

4 DTU Compute

Lecture 11 25 April, 2025

Dynamical Programming

| Bellman equation | Learning algorithm |
|---|--|
| Bellman expectation equation for v_π $v_\pi(s) = \mathbb{E}_\pi [R + \gamma v_\pi(S') s]$ | Iterative policy evaluation to learn v_π $V(s) \leftarrow \mathbb{E}_\pi [R + \gamma V(S') s]$ |
| Bellman expectation equation for q_π $q_\pi(s, a) = \mathbb{E}_\pi [R + \gamma q_\pi(S', A') s, a]$ | Iterative policy evaluation to learn q_π $Q(s, a) \leftarrow \mathbb{E}_\pi [R + \gamma Q(S', A') s, a]$ |
| Policy iteration: Use policy evaluation to estimate v_π or q_π Improve by acting greedily: $\pi'(s) \leftarrow \arg \max_a q_\pi(s, a)$ | |
| Bellman optimality equation for v_* $v_*(s) = \max_a \mathbb{E}[R + \gamma v_*(S') s, a]$ | Value iteration $V(s) \leftarrow \max_a \mathbb{E}[R + \gamma V(S') s, a]$ |
| Bellman optimality equation for q_* $q_*(s, a) = \mathbb{E}[R + \gamma \max_{a'} q_*(S', a') s, a]$ | Q-value iteration $Q(s, a) \leftarrow \mathbb{E}[R + \gamma \max_{a'} Q(S', a') s, a]$ |

5 DTU Compute

Lecture 11 25 April, 2025



Sarsa control TD and MC value estimation

- Recall $v_\pi(s) = \mathbb{E}[G_t | S_t = s]$
- MC learning: G_t estimate of $v_\pi(s)$; update:
$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$$
- Bellman equation:
$$v_\pi(s) = \mathbb{E}[R_{t+1} + \gamma V(S_{t+1}) | S_t = s]$$
- TD learning: $R_{t+1} + \gamma V(S_{t+1})$ is also an estimate of $v_\pi(s)$; update:
$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$
- TD learning has several advantages
 - Lower variance
 - Don't have to wait for episode to finish
- Natural idea: Apply TD to $Q(s, a)$
 - Still ϵ -greedy policy improvement
 - Update Q estimates at each time step

6 DTU Compute

Lecture 11 25 April, 2025

Sarsa estimation of action-value function

- Bellman equation:

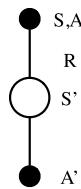
$$q_\pi(s, a) = \mathbb{E}[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$

- Implies $R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1})$ is an estimate of $q_\pi(s, a)$

- Implies the update equation

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

- We use bootstrapping (i.e. biased estimate)

**Sarsa control****Sarsa (on-policy TD control) for estimating $Q \approx q_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

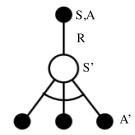
 Initialize S
 Choose A from S using policy derived from Q (e.g., ε -greedy)
 Loop for each step of episode:
 Take action A , observe R, S'
 Choose A' from S' using policy derived from Q (e.g., ε -greedy)
 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$
 $S \leftarrow S'; A \leftarrow A'$;
 until S is terminal

Convergence of Sarsa

Sarsa converge to optimal action-value function $Q \rightarrow q_*$ assuming

- GLIE sequence of policies (decreasing but non-trivial exploration)
- Robbins-Monro sequence of step-sizes α_t

$$\sum_{t=1}^{\infty} \alpha_t = \infty, \quad \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

Using the Bellman optimality equation

- Bellman equation:

$$q_*(s, a) = \mathbb{E} [R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a]$$

- Implies $R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a')$ is a Monte-Carlo estimate of $q_*(s, a)$
- Implied update equation

$$Q(S, A) \leftarrow Q(S, A) + \alpha (R + \gamma \max_{a'} Q(S', a') - Q(S, A))$$

- Note we use bootstrapping (i.e. biased estimate)

Q-learning
Q-learning is off-policy

$$Q(S, A) \leftarrow Q(S, A) + \alpha (R + \gamma \max_{a'} Q(S', a') - Q(S, A))$$

- The **behavior policy** determines which S_t, A_t are visited
- The environment determines what happens next (S')
- The Q -values are updated **without** reference to the **behavior policy**
- Q -learning is therefore **off-policy**

Q-learning
Q*-learning*Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S
 Loop for each step of episode:
 Choose A from S using policy derived from Q (e.g., ε -greedy)
 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

Q-learning Exam question: Q-learning



- a. The first step in training a Q -learning agent is to compute the set of all states the agent can be in
- b. The Q -table $Q(s, a)$ in Q -learning is a measure of the reward the agent will obtain in the very next step multiplied by γ
- c. Q -learning still works if we initialize the Q -table to -1 , i.e. $Q(s, a) = -1$ for all $s \in \mathcal{S}$
- d. When Q -learning is applied to a deterministic environment, the agent will follow a deterministic policy
- e. Don't know.

13 DTU Compute

Lecture 11 25 April, 2025

Convergence of Q -learning

Q -learning converge to optimal action-value function $Q \rightarrow q_*$ assuming

- All s, a pairs visited infinitely often
- Robbins-Monro sequence of step-sizes α_t

$$\sum_{t=1}^{\infty} \alpha_t = \infty, \quad \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

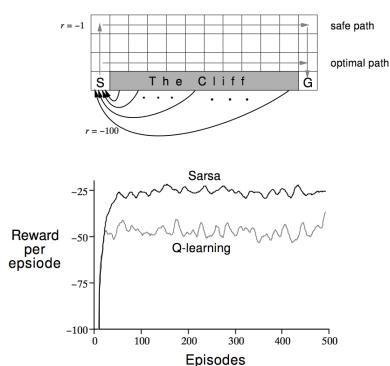
14 DTU Compute

Lecture 11 25 April, 2025

Q-learning Comparing Q -learning and SARSA



- Reward -100 if we fall
- Reward -1 per step
- Both use ϵ -greedy exploration



lecture_11_sarsa_cliff.py, lecture_11_q_cliff.py

15 DTU Compute

Lecture 11 25 April, 2025

Q-learning Algorithms so far

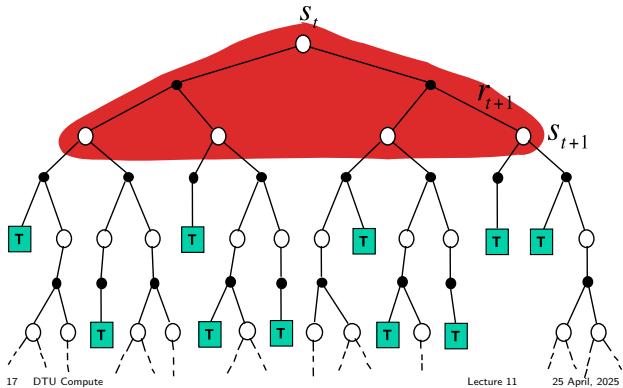
| Bellman equation | Learning algorithm | TD Learning |
|---|---|--|
| $V(s) \leftarrow \mathbb{E}_\pi [R + \gamma V(S') s]$ | Iterative policy evaluation to learn v_π | $V(S) \leftarrow R + \gamma V(S')$ |
| $q_\pi(s, a) = \mathbb{E}_\pi [R + \gamma q_\pi(S', A') s, a]$ | Iterative policy evaluation to learn q_π | Sarsa $Q(S, A) \leftarrow R + \gamma Q(S', A')$ |
| $v_*(s) = \max_a \mathbb{E}[R + \gamma v_*(S') s, a]$ | Policy iteration: Use policy evaluation to estimate v_π or q_π Improve by acting greedily: $\pi'(s) \leftarrow \arg \max_a q_\pi(s, a)$ | Q-value iteration $Q(S, A) \leftarrow \mathbb{E}[R + \gamma \max_{a'} Q(S', a') s, a]$ |
| $q_*(s, a) = \mathbb{E}[R + \gamma \max_{a'} q_*(S', a') s, a]$ | Value iteration | $Q(S, A) \leftarrow \mathbb{E}[R + \gamma \max_{a'} Q(S', a') s, a]$ |

16 DTU Compute

n -step backups From two weeks ago: DP backups



$$V(S_t) \leftarrow \mathbb{E}_\pi [R_{t+1} + \gamma V(S_{t+1})]$$

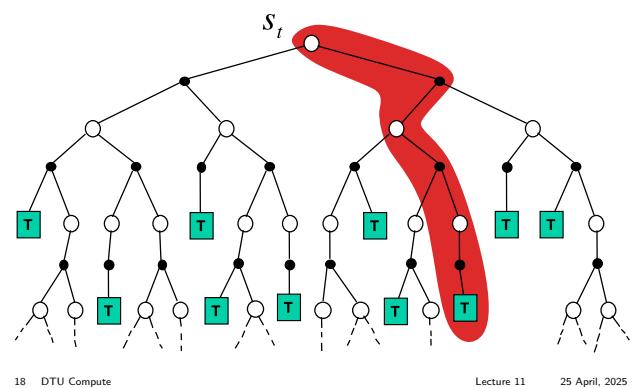


17 DTU Compute

Lecture 11 25 April, 2025

n -step backups Last week: MC backups

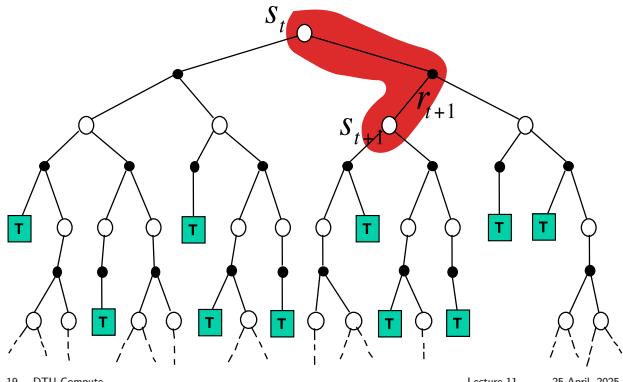
$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



18 DTU Compute

Lecture 11 25 April, 2025

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

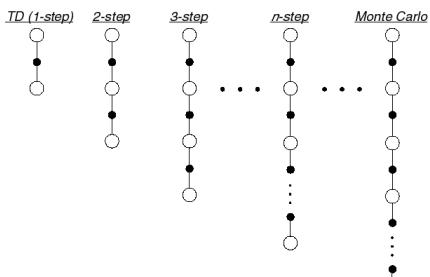


- **Bootstrapping:** Update involves an estimate (e.g. V)
 - TD and DP bootstraps
 - MC does not bootstrap

- **Sampling:** Update involves a sample estimate of an expectation
 - MC and TD sample
 - DP does not sample

Let's combine methods and avoid either/or choices

- Let TD target look n steps into the future



- Recall return is $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots$

$$n=1: (\text{TD}) \quad G_t^{(1)} = R_{t+1} + \gamma G_{t+1}$$

$$n=2: \quad G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 G_{t+2}$$

$$n: \quad G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n G_{t+n}$$

$$n=\infty (\text{MC}): \quad G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

- Using the rules of expectations:

$$\begin{aligned} v_\pi(s) &= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n G_{t+n}|s] \\ &= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \mathbb{E}[\gamma^n G_{t+n}|S_{t+n}]|S_t=s] \\ &= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n v_\pi(S_{t+n})|S_t=s] \end{aligned}$$

Therefore, the n -step return is an estimate of $V(S_t)$

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

- This gives n -step temporal difference update:

$$V(S_t) \leftarrow V(S_t) + \alpha(G_{t:t+n} - V(S_t))$$

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

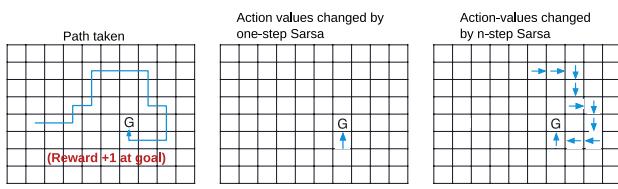
$$V(S_t) \leftarrow V(S_t) + \alpha(G_t^{(n)} - V(S_t))$$

- We cannot compute $G_t^{(n)}$ until we have the n next steps episodes

- Maintain buffer of size n

- At end of episode, we are still missing $n-1$ updates

- Do a for-loop and perform missing updates



Recall the decomposition:

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n G_{t+n}$$

- As before:

$$\begin{aligned} q_\pi(s, a) &= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n G_{t+n}|S_t=s, A_t=a] \\ &= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n q_\pi(S_{t+n}, A_{t+n})|S_t=s, A_t=a] \end{aligned}$$

- Therefore, the following n -step action-value return is an unbiased estimate of q_π

$$q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n q_\pi(S_{t+n}, A_{t+n})$$

- Suggest the following bootstrap update of the action-value function

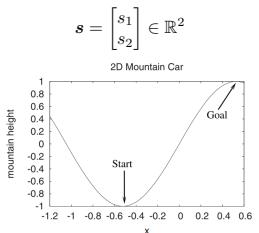
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(q_t^{(n)} - Q(S_t, A_t))$$

Value-function approximations

Scaling up reinforcement learning

We want to apply RL to large problems

- Chess: $> 10^{40}$ states
- Go: $> 10^{170}$ states
- Robot arm: continuous state space
- **Example: Mountain-Car** position, velocity. Discrete actions



lecture_11_mountaincar_nolearn.py

25 DTU Compute

Lecture 11 25 April, 2025



Value-function approximations

Value Function Approximation

- We have used lookup table representation (stored $Q(s, a)$ as a big table)

- Every state s has an entry $V(s)$ or
- Every state-action pair s, a has an entry $Q(s, a)$

- Issues with lookup tables

- There are too many states and/or actions to store in memory
- It is too slow to learn the value of each state individually

- Idea:

- Estimate value function or state-action value with function approximation

$$\hat{v}(s, \mathbf{w}) \approx v_\pi(s)$$

$$\hat{q}(s, a, \mathbf{w}) \approx q_\pi(s, a)$$

- Generalize from seen states to unseen states

26 DTU Compute

Lecture 11 25 April, 2025



Value-function approximations

Feature Vectors and linear representations

- Represent value function by a linear combination of features

$$\hat{v}(s, \mathbf{w}) = \mathbf{x}(s)^\top \mathbf{w}, \quad \mathbf{w} \in \mathbb{R}^d$$

Where **feature vector** is defined as:

$$\mathbf{x}(s) = \begin{bmatrix} x_1(s) \\ \vdots \\ x_d(s) \end{bmatrix}$$

- The gradient is simply:

$$\nabla \hat{v}(s, \mathbf{w}) = \mathbf{x}(s)$$

- For Q -values we only need to change the feature vector:

$$\hat{q}(s, a, \mathbf{w}) = \mathbf{x}(s, a)^\top \mathbf{w}$$

27 DTU Compute

Lecture 11 25 April, 2025



Value-function approximations

Example: Mountain-car

- Mountain-car has two **actions** $a = 0, 1$ (left, right); the **reward** is $R_t = -1$.

- The **state** is two-dimensional $s = (\text{position, velocity}) = (s_1, s_2)$

- **Naive example 1:** $\mathbf{x}(s) = \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$. **Naive example 2:** $\mathbf{x}(s) = \begin{bmatrix} s_1 \\ \sin(s_1) \cos(s_2) \\ \vdots \end{bmatrix}$

- How about $\mathbf{x}(s, a)$? Idea: Re-use $\mathbf{x}(s)$ by zero-padding:

$$\mathbf{x}(s, a = 0)^\top = [\underbrace{0, 0, 0, 0, 0, \dots, 0}_{\text{pad with } |\mathbf{x}(s)| \text{ zeros}}, \mathbf{x}(s)^\top]$$

$$\mathbf{x}(s, a = 1)^\top = [\mathbf{x}(s)^\top, \underbrace{0, 0, 0, 0, 0, \dots, 0}_{\text{pad with } |\mathbf{x}(s)| \text{ zeros}}]$$

28 DTU Compute

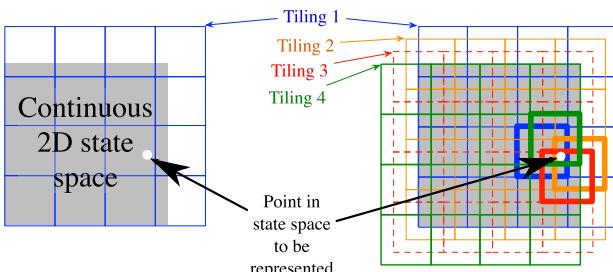
Lecture 11 25 April, 2025



Value-function approximations

Details: Tile coding

- Divide each dimension of s into a number of tiles n_T
- Translate tiles in fraction of tile width to get overlap



- \mathbf{x} has now n_T non-zero elements corresponding to the number of active tiles

lecture_11_mountaincar_random_weights.py

29 DTU Compute

Lecture 11 25 April, 2025



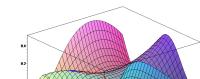
Value-function approximations

Recall from 02450: Gradient Descent

- Let $E(\mathbf{w})$ be a differentiable function of parameter vector \mathbf{w}

- The gradient of $E(\mathbf{w})$ is

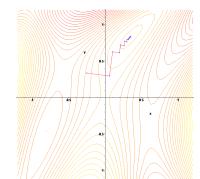
$$\nabla_{\mathbf{w}} E(\mathbf{w}) = \begin{bmatrix} \frac{\partial E(\mathbf{w})}{\partial w_1} \\ \vdots \\ \frac{\partial E(\mathbf{w})}{\partial w_n} \end{bmatrix}$$



- Adjust \mathbf{w} in direction of negative gradient to find a **local minimum** of $E(\mathbf{w})$

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} E(\mathbf{w})$$

with step-size parameter α (**learning rate**)



30 DTU Compute

Lecture 11 25 April, 2025



Value-function approximations Using the approximations



- Consider TD learning which implements Bellman equation:

$$v_\pi(s) = \mathbb{E}[R + \gamma v(S')|s]$$

- Standard TD update

$$V(s) \leftarrow V(s) + \alpha(r + \gamma V(s') - V(s))$$

- Easy to plug in $\hat{v}(s, w)$ instead of $V(s)$ on right-hand side

$$\hat{v}(s, w) \leftarrow \hat{v}(s, w) + \alpha(r + \gamma \hat{v}(s', w) - \hat{v}(s, w))$$

- ..but how do we update w on the left-hand side so $\hat{v}(s, w)$ agrees with r.h.s.?

31 DTU Compute

Lecture 11 25 April, 2025

Value-function approximations Take a step back: What do we want to do?

- No function approximators: $v(s) = \mathbb{E}[R + \gamma v(S')|s]$
- With function approximators: Find w so that:

$$\hat{v}(s, w) = \mathbb{E}[R + \gamma v(S')|s]$$

- Find w so that:

$$w = \arg \min_w \frac{1}{2} (\hat{v}(s, w) - \mathbb{E}[R + \gamma v(S')|s])^2$$

- Find w using gradient descent:

$$\begin{aligned} w &\leftarrow w + \alpha \nabla_w \frac{1}{2} (\hat{v}(s, w) - \mathbb{E}[R + \gamma v(S')|s])^2 \\ &= w + \alpha (\hat{v}(s, w) - \underbrace{\mathbb{E}[R + \gamma v(S')|s]}_{\approx \frac{1}{B} \sum_{n=1}^B r^{(n)} + v(s^{(n)})}) \nabla \hat{v}(s, w) \end{aligned}$$

- Use a sample-size of $B = 1$ to compute the average

$$w \leftarrow w + \alpha (\hat{v}(s, w) - r + \gamma v(s')) \nabla \hat{v}(s, w)$$

32 DTU Compute

Lecture 11 25 April, 2025

Value-function approximations Summary



- Given $f(x) = \mathbb{E}_z[g(x, z)]$ and approximation-function $\hat{f}(x, w)$
 - To find w such that $\hat{f}(x, w) \approx f(x)$ iterate:
- $$w \leftarrow w + \alpha (g(x, z) - \hat{f}(x, w)) \nabla \hat{f}(x, w)$$
- TD learning: $V(s) = \mathbb{E}[R + \gamma V(S')|s]$ and $\hat{v}(s, w) \approx v(s)$
- $$V(s) \leftarrow V(s) + \alpha(r + \gamma V(s') - V(s))$$
- $$w \leftarrow w + \alpha (r + \gamma \hat{v}(s', w) - \hat{v}(s, w)) \nabla \hat{v}(s, w)$$
- Sarsa learning: $q(s, s) = \mathbb{E}[R + \gamma q(S', A')|s, a]$ and $\hat{q}(s, a, w) \approx q(s, a)$
- $$q(s, a) \leftarrow q(s, a) + \alpha (r + \gamma q(s', a') - q(s, a))$$
- $$w \leftarrow w + \alpha (r + \gamma \hat{q}(s', a', w) - \hat{q}(s, a, w)) \nabla \hat{q}(s, a, w)$$
- Q -learning: $q(s, s) = \mathbb{E}[R + \gamma \max_a q(S', a')|s, a]$ and $\hat{q}(s, a, w) \approx q(s, a)$
- $$q(s, a) \leftarrow q(s, a) + \alpha (r + \gamma \max_{a'} q(s', a') - q(s, a))$$
- $$w \leftarrow w + \alpha (r + \gamma \max_{a'} \hat{q}(s', a', w) - \hat{q}(s, a, w)) \nabla \hat{q}(s, a, w)$$
- Remember that $\nabla \hat{q}(s, a, w) = x(s, a)$ and $\nabla v(s, w) = x(s)$

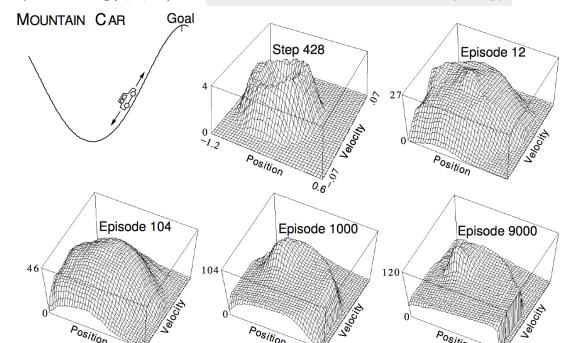
33 DTU Compute

Lecture 11 25 April, 2025

Value-function approximations Linear Sarsa with tile coding in mountain car



- Sarsa with linear tile-encoding ($|x(s)| \approx 2000$)
- We plot $\max_a \hat{q}(s, a; w)$: [lecture_11_mountaincar_feature_space.py](#)



34 DTU Compute

Lecture 11 25 April, 2025

Value-function approximations Quiz: Linear function approximators



Which of the following statements is true about reinforcement learning and linear function approximators?

- Linear function approximators can only be used with continuous state spaces and not with discrete spaces.
- Linear function approximators provide a way to generalize from known states to unknown states, which can be useful in tabular reinforcement learning situations with large state spaces.
- Linear function approximators in SARSA or Q-learning requires that we store all state-action pairs.
- When using linear function approximators the policy will be deterministic
- Don't know.

35 DTU Compute

Lecture 11 25 April, 2025

Value-function approximations Implementing this



```

1 # semi_grad_q.py
2 class LinearSemiGradAgent(Agent):
3     def __init__(self, env, gamma=1.0, alpha=0.5, epsilon=0.1, q_encoder=None):
4         """The Q-values, as implemented using a function approximator, can now be accessed as follows:
5
6         >> self.Q(s, a) # Compute q-value
7         >> self.Q.x(s, a) # Compute gradient of the above expression wrt. w
8         >> self.Q.w # get weight-vector.
9
10        I would recommend inserting a breakpoint and investigating the above expressions yourself;
11        you can of course also check the class LinearQEncoder if you want to see how it is done in practice.
12        """
13        super().__init__(env, gamma, epsilon-epsilon, alpha=alpha)
14        self.Q = LinearQEncoder(env, tilings=8) if q_encoder is None else q_encoder

```

36 DTU Compute

Lecture 11 25 April, 2025

- Suppose f is a real-valued function $f : \mathcal{X} \mapsto \mathbb{R}$ which happens to be defined using an expectation:

$$f(x) = \mathbb{E}_z [g(x, z)] = \int p(z|x)g(x, z)dz$$

- Assume that $\hat{f}(x, w)$ is a neural network we want to use to approximate f with
- Problem:** How do we find w such that $\hat{f}(x, w) \approx f(x)$?
- Idea:** Select w to minimize

$$w^* = \arg \min_w \mathbb{E}_x \left[[\hat{f}(x, w) - f(x)]^2 \right] \quad (1)$$

- Solve this using gradient descent:

$$w \leftarrow w - \alpha \nabla \left(\mathbb{E} [\hat{f}(x, w) - f(x)]^2 \right) \quad (2)$$

Richard S. Sutton and Andrew G. Barto.
Reinforcement Learning: An Introduction.
The MIT Press, second edition, 2018.
(Freely available online).

Appendix Evaluating the gradient

$$\begin{aligned} \nabla \left(\mathbb{E} [\hat{f}(x, w) - f(x)]^2 \right) &= \mathbb{E} \left[\nabla (\hat{f}(x, w) - f(x))^2 \right] \\ &= 2\mathbb{E} \left[(\hat{f}(x, w) - f(x)) \nabla \hat{f}(x, w) \right] \\ &= 2\mathbb{E} \left[(\hat{f}(x, w) - \mathbb{E}_z [g(x, z)]) \nabla \hat{f}(x, w) \right] \end{aligned}$$

Implication: Given samples $x \sim p$ and $z \sim p(z|x)$ then

$$2 (\hat{f}(x, w) - g(x, z)) \nabla \hat{f}(x, w)$$

is an **unbiased estimate** of the gradient

Stochastic gradient descent

Given minimization problem $\arg \min F(w)$ and (technical conditions!) then

$$w_{t+1} \leftarrow w_t - \alpha_t \hat{g}(w_t)$$

converge to w^* provided $\hat{g}(w)$ is an **unbiased estimate** of the gradient $\nabla F(w)$