

Lecture Schedule

Dynamical programming

- 1 The finite-horizon decision problem
31 January
- 2 Dynamical Programming
7 February
- 3 DP reformulations and introduction to Control
14 February

Control

- 4 Discretization and PID control
21 February
- 5 Direct methods and control by optimization
28 February
- 6 Linear-quadratic problems in control
7 March
- 7 Linearization and iterative LQR
14 March

Syllabus: <https://02465material.pages.compute.dtu.dk/02465public>
Help improve lecture by giving feedback on DTU learn

Reinforcement learning

- 8 Exploration and Bandits
21 March
- 9 Policy and value iteration
4 April
- 10 Monte-carlo methods and TD learning
11 April
- 11 Model-Free Control with tabular and linear methods
18 April
- 12 Eligibility traces and value-function approximations
25 April
- 13 Q-learning and deep-Q learning
2 May

Reading material:

- [SB18, Chapter 6.7-6.9; 8-8.4; 16-16.2; 16.5; 16.6]

Learning Objectives

- Double-Q learning
- Dyna-Q and the replay buffer
- Deep-Q learning

Housekeeping

- Unofficial exam Q/A about one week before the exam (the 20th?). Please put wishes on blackboard.
- I have added a survey on the course (what went well/ less well /what can be improved). You can find it in the menu to the right on DTU Learn.
- I have updated the video on preparing for the exam, <https://www2.compute.dtu.dk/courses/02465/exam.html>, and uploaded solutions to the previous exams.
- Exam is planned to be in English as last year (only one language). Please let me know before Tuesday the 7th if this is not acceptable.
- Test exam at <https://eksamen.dtu.dk/studerende/proeve/7482/tilmeld/3a1b13368489ef57c103c1e4642d6ff2> (Hopefully this works!)

- **Bellman optimality condition:**

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

- **Theorem:** q_* satisfies the above recursions if (and only if) it corresponds to the **optimal value function**
- **Value iteration:** Replace q_* arbitrary Q and iterate:

$$Q(s, a) \leftarrow \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a') \mid S_t = s, A_t = a \right]$$

- **Theorem:** Q will converge to q_*
- **Q-learning:** Given $(S_t, A_t, R_{t+1}, S_{t+1}) = (s, a, r, s')$ transition, update

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

Uses that red expression is a **biased** but **consistent** estimate of Q

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 $s \in S^+$, $a \in \mathcal{A}(s)$, 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

Convergence of Q-learning

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

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

- Value iteration uses a **model** of the environment to **plan** a policy

$$Q(s, a) \leftarrow \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a') \mid S_t = s, A_t = a \right]$$

- Q-learning uses **samples** from the environment (s, a, r, s') to **learn** a policy

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

- Both uses value functions and backups
- Can we combine these ideas?

- A **distributional model** is an estimate of the MDP $p(s', r|s, a)$
- A **sample model** is a mechanism to generate samples (s, a, r, s') from the MDP (**weaker assumption**)
- Idea: Learn sample model and use it to improve value function by regular backups
- Allows re-use of data for faster convergence (**sample efficiency**)

Random-sample one-step tabular Q-planning

Loop forever:

1. Select a state, $S \in \mathcal{S}$, and an action, $A \in \mathcal{A}(S)$, at random
2. Send S, A to a sample model, and obtain
a sample next reward, R , and a sample next state, S'
3. Apply one-step tabular Q-learning to S, A, R, S' :

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

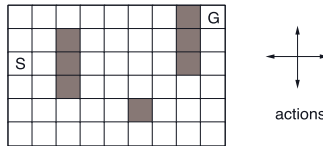
Tabular Dyna-Q

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

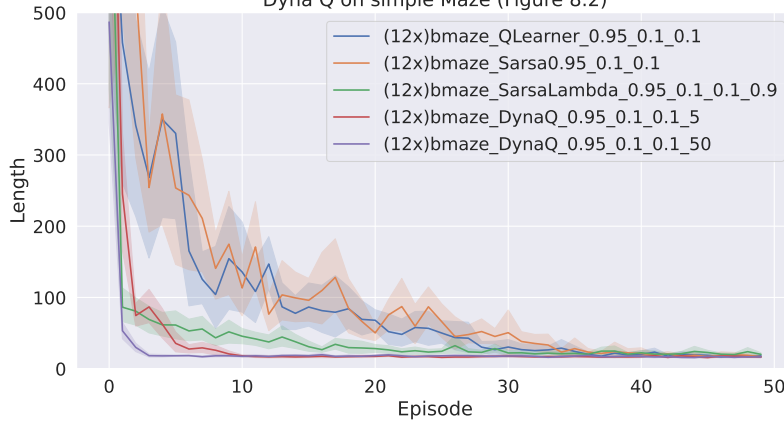
Loop forever:

- (a) $S \leftarrow$ current (nonterminal) state
- (b) $A \leftarrow \varepsilon$ -greedy(S, Q)
- (c) Take action A ; observe resultant reward, R , and state, S'
- (d) $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
- (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)
- (f) Loop repeat n times:
 - $S \leftarrow$ random previously observed state
 - $A \leftarrow$ random action previously taken in S
 - $R, S' \leftarrow Model(S, A)$
 - $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

Dyna-Q on deterministic Maze environment



Dyna Q on simple Maze (Figure 8.2)



🔧 `lecture_13_Q_maze.py` , 🔧 `lecture_13_dyna_q_5_maze.py` ,

🔧 `lecture_13_sarsa_lambda_maze.py`

Tabular Dyna-Q

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

Loop forever:

(a) $S \leftarrow$ current (nonterminal) state

(b) $A \leftarrow \epsilon$ -greedy(S, Q)

(c) Take action A ; observe resultant reward, R , and state, S'

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

(e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)

(f) Loop repeat n times:

$S \leftarrow$ random previously observed state

$A \leftarrow$ random action previously taken in S

$R, S' \leftarrow Model(S, A)$

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

- The **model** is simply a list of experience (a **replay buffer**)
- Deterministic assumption not used

- Target for the Q -values can be considered noisy (random)

$$r + \max_{a'} Q(s', a').$$

- Q -update is

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left(r + \max_{a'} Q(s', a') - Q(s, a) \right)$$

- By chance some of the $Q(s', a')$ values are likely to be unusually large
- This leads to over-estimate $Q(s, a)$:

$$\mathbb{E}[\max(X_1, X_2)] \geq \max(\mathbb{E}[X_1], \mathbb{E}[X_2])$$

- **Conclusion:**

- Q -values systematically over-estimated
- the worse the estimate of a state, the more we will prefer it

Given transition $(S_t, A_t, R_{t+1}, S_{t+1}) = (s, a, r, s')$

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s', \arg \max_a Q(s', a)) \right] Q_2(s', \arg \max_a Q(s', a))$$

- Where Q_2 is another Q -function
- Q_2 is independent of Q which avoids systematic over-estimation

Double Q -learning, for estimating $Q_1 \approx Q_2 \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q_1(s, a)$ and $Q_2(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, such that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

 Choose A from S using the policy ε -greedy in $Q_1 + Q_2$

 Take action A , observe R, S'

 With 0.5 probability:

$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha \left(R + \gamma Q_2(S', \arg \max_a Q_1(S', a)) - Q_1(S, A) \right)$$

 else:

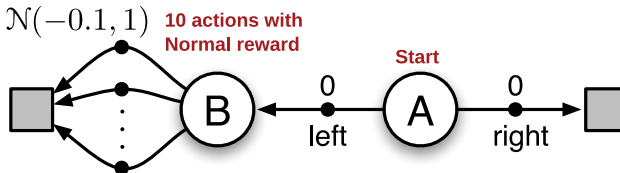
$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha \left(R + \gamma Q_1(S', \arg \max_a Q_2(S', a)) - Q_2(S, A) \right)$$

$S \leftarrow S'$

 until S is terminal

- Twice as slow to learn

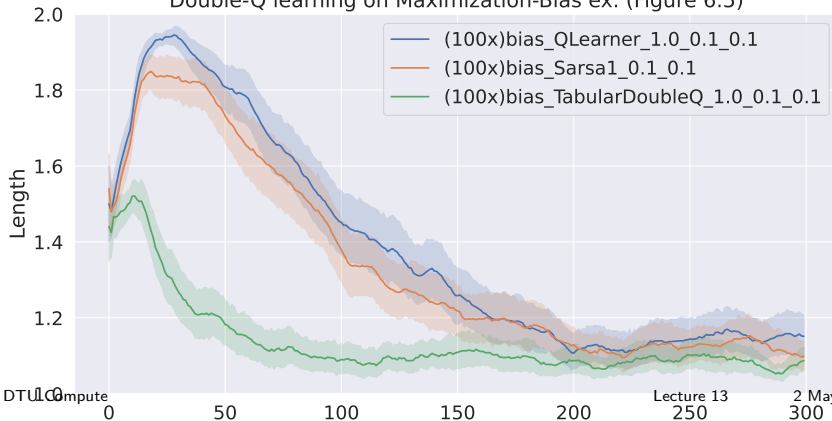
Double-Q learning on bias-example environment

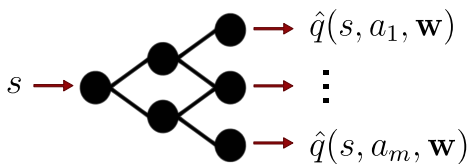


Trajectory length 2
(average reward -0.1)

Trajectory length 1
(average reward 0)

Double-Q learning on Maximization-Bias ex. (Figure 6.5)





- We want an approximation of the Q -values $Q(s, a)$
- Assume $\mathbf{y} = \hat{q}_\phi(s)$ is a vector of dimension $|\mathcal{A}|$ such that

$$y_a \approx Q(s, a)$$

is our approximation of the Q -value

- In practice, $\hat{q}_\phi : \mathbb{R}^d \mapsto \mathbb{R}^{|\mathcal{A}|}$ is a deep network
 - Input-dimension is dimension of each state $s \in \mathcal{S} = \mathbb{R}^d$
 - Output dimension $|\mathcal{A}|$

Regular Q -learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

Regular Q -learning with function approximators

- Given $(S_t, A_t, R_{t+1}, S_{t+1}) = (s, a, r, s')$ update:

$$\phi \leftarrow \phi + \alpha \left(r + \gamma \max_{a'} \hat{q}_\phi(s', a') - \hat{q}_\phi(s, a) \right) \nabla_\phi \hat{q}_\phi(s, a)$$

- Defining $y = r + \gamma \max_{a'} \hat{q}_\phi(s', a')$ this can be written as

$$\phi \leftarrow \phi - \alpha \frac{1}{2} \nabla_\phi (y - \hat{q}_\phi(s, a))^2$$

Fitted Q-iteration algorithm

① At step t observe $(s_t, a_t, r_{t+1}, s_{t+1})$

② $y_t = r_{t+1} + \gamma \max_{a'} \hat{q}_\phi(s_{t+1}, a')$

③ Repeat fit step one or more times:

- $\phi \leftarrow \phi - \alpha \nabla_\phi \left[\frac{1}{2} (y_t - \hat{q}_\phi(s_t, a_t))^2 \right]$

- The use of a **single** sample gives a high variance in the gradient estimate
- The samples are only used once

Initialize a **replay buffer** \mathcal{B}

Q-learning with a replay buffer

① At step t observe $(s_t, a_t, r_{t+1}, s_{t+1})$ and add it to \mathcal{B}

② Repeat K times:

① Sample a **batch** $(s_i, a_i, r_i, s'_i)_{i=1}^B$ from \mathcal{B}

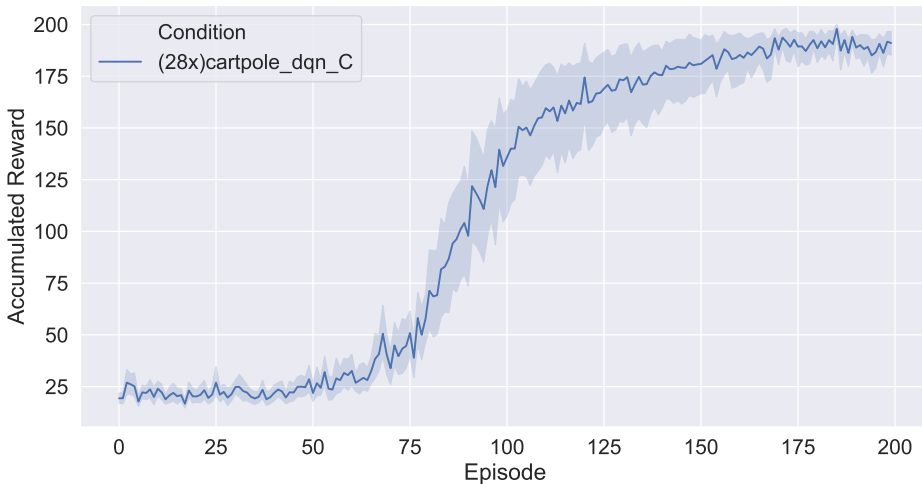
② Set $y_i = r_i + \gamma \max_{a'} \hat{q}_\phi(s'_i, a')$

③ $\phi \leftarrow \phi - \alpha \nabla_\phi \left[\frac{1}{2B} \sum_{i=1}^B (y_i - \hat{q}_\phi(s_i, a_i))^2 \right]$

- Similar to dyna-Q
- Lower gradient variance, quicker convergence
- Replay buffer should be large (thousands to a few millions)
- **You can implement this in the exercises**

Q-learning and function approximators

Basic deep Q learning on Cartpole



- Consider the target
 - ① $y = r_{t+1} + \gamma \max_{a'} \hat{q}_\phi(s_{t+1}, a')$
 - ② $\phi \leftarrow \phi - \alpha \nabla_\phi \left[\frac{1}{2} (y - \hat{q}_\phi(s, a))^2 \right]$
- We don't compute gradients through y
- This is to a great extent why deep-Q sometimes do not converge: We adapt towards y , without taking into account that y changes during the adaption
- **Idea 1:** Use an alternative weight network ϕ'

$$y = r_{t+1} + \gamma \max_{a'} \hat{q}_{\phi'}(s_{t+1}, a')$$

- **Idea 2:** Let ϕ' be an old version of ϕ

Initialize \mathcal{B} and make a copy $\phi' \leftarrow \phi$ of the weights

Deep-Q learning

- 1 At step t observe $(s_t, a_t, r_{t+1}, s_{t+1})$ and add it to \mathcal{B}
 - 2 Repeat K times:
 - 1 Sample a batch $(s_i, a_i, r_i, s'_i)_{i=1}^B$ from \mathcal{B}
 - 2 Set $y_i = r_i + \gamma \max_{a'} \hat{q}_{\phi'}(s'_i, a')$
 - 3 $\phi \leftarrow \phi - \alpha \nabla_{\phi} \left[\frac{1}{2B} \sum_{i=1}^B (y_i - \hat{q}_{\phi}(s_i, a_i))^2 \right]$
 - 3 Update $\phi' \leftarrow \phi' + \tau(\phi - \phi')$ (Slow changes, e.g. $\tau = 0.08$ or less)
- Can we also address the over-estimation problem of the Q-values?

Initialize \mathcal{B} and make a copy $\phi' \leftarrow \phi$ of the weights

Double-Q learning

① At step t observe $(s_t, a_t, r_{t+1}, s_{t+1})$ and add it to \mathcal{B}

② Repeat K times:

① Sample a batch $(s_i, a_i, r_i, s'_i)_{i=1}^B$ from \mathcal{B}

② Set $y_i = r_i + \gamma \hat{q}_{\phi'}(s'_i, \arg \max_{a'} \hat{q}_{\phi}(s', a'))$

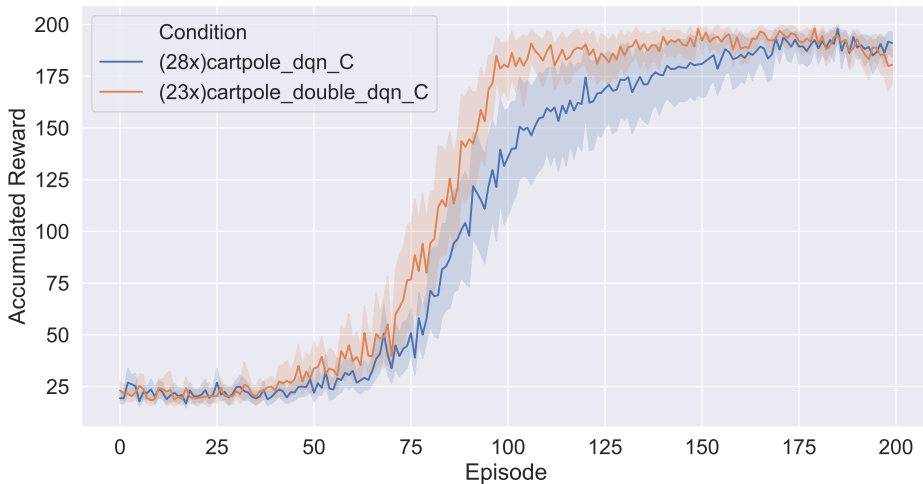
③ $\phi \leftarrow \phi - \alpha \nabla_{\phi} \left[\frac{1}{2B} \sum_{i=1}^B (y_i - \hat{q}_{\phi}(s_i, a_i))^2 \right]$

③ Update $\phi' \leftarrow \phi' + \tau(\phi - \phi')$

- Double-Q: Select actions according to ϕ , but evaluate according to ϕ'
- **We will implement this in the exercises**

Q-learning and function approximators

Double-deep Q learning on Cartpole



The buffer is a list with a sample function

```
1  # deepq_agent.py
2  self.memory = BasicBuffer(replay_buffer_size) if buffer is None else buffer
3  self.memory.push(s, a, r, sp, done) # save current observation
4  """ First we sample from replay buffer. Returns numpy Arrays of dimension
5  > [self.batch_size] x [...]
6  for instance 'a' will be of dimension [self.batch_size x 1].
7  """
8  s,a,r,sp,done = self.memory.sample(self.batch_size)
```

First dimension is batch dimension

$$(\text{batch_size} \times d)$$

Implemented in separate class

```

1  # irlc/ex13/lecture_12_examples.py
2  # Initialize a network class
3  self.Q = Network(env, trainable=True) # initialize the network
4  """ Assuming s has dimension [batch_dim x d] this returns a float numpy Array
5  array of Q-values of [batch_dim x actions], such that qvals[i,a] = Q(s_i,a) """
6  qvals = self.Q(s)
7  actions = env.action_space.n # number of actions
8  """ Assume we initialize target to be of dimension [batch_dim x actions]
9  > target = [batch_dim x actions]
10 The following function will fit the weights in self.Q by minimizing
11 > ||self.Q(s)-target||2
12 (averaged over Batch dimension) using one step of gradient descent
13 """
14 self.Q.fit(s, target)

```

i.e. select `target` appropriately to implement loss

$$\frac{1}{B} \sum_{i=1}^B (\hat{q}_{\phi}(s_i, a_i) - y_i)^2$$



```
1 # irlc/ex13/lecture_12_examples.py
2 self.Q2 = Network(env, trainable=True)
3 """ Update weights in self.Q2 (target, phi') towards those in Q (source, phi)
4 with a factor of tau. tau=0 is no change, tau=1 means overwriting weights
5 (useful for initialization) """
6 self.Q2.update_Phi(Q2, tau=0.1)
```

Updates weights ϕ' in `q2` towards ϕ in `q`

$$\phi' = \phi' + \tau(\phi - \phi')$$

- Parameters: Decrease exploration rate ε_t and learning rate α_t through training
- Networks
 - Clip gradients or use Huber loss
 - Batch normalization
 - Tune parameters; linear \rightarrow shallow \rightarrow deep
- Methods:
 - Double-Q learning always a good idea
 - Replay buffer always a good idea
 - Prioritizing samples (PER) improves convergence speed
 - Check out **Rainbow** for current(ish) state of the art(ish) [HMHV⁺18]
- Lots of training and results highly variable across seeds

FIN!

-  Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver.
Rainbow: Combining improvements in deep reinforcement learning.
In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
-  Richard S. Sutton and Andrew G. Barto.
Reinforcement Learning: An Introduction.
The MIT Press, second edition, 2018.
(Freely available online).