



Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Conclusion

Theoretical aspects of Q-learning

Masters thesis defense

Jacob Harder

Department of Mathematical Sciences

University of Copenhagen

26 June, 2020



Overview

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Conclusion

① Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

② Model-dependent algorithms

Approximation and algorithmic errors

Approximation using ANNs

Approximation with Bernstein polynomials

③ Model-free algorithms

Finite case

History dependent setting



Q-learning as AI

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

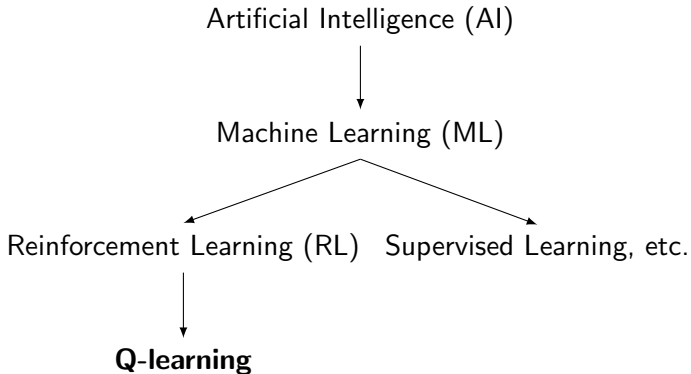
Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion





Machine learning

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Machine Learning is “the study of computer algorithms that improve automatically through *experience*”.

- **Supervised learning**: Tasks are learned from data based on feedback from a *supervisor*. E.g. image classification.
- **Unsupervised learning**: Data is given without evaluatory feedback, general trends about the data are analysed. E.g. principal component analysis, and cluster analysis.
- \rightarrow^1 **Reinforcement learning**: Algorithms which learns through interactions with an *environment*.

¹ “ \rightarrow ”: Our main area of focus in this thesis.



Challenges in RL

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Challenges in Reinforcement Learning include:

- **Exploration-exploitation trade-off.** Training and performing occurs simultaneously so one optimizes the total reward on some time horizon. This is studied in e.g. the multi-armed bandit problem.
- → **Deriving optimal policies.** Training and performing is distinguished and emphasis is put on the expected performance of the final derived policy rather than rewards occurring during training.



The environment

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

The **environment** in RL is often formalized as a **Markov decision process** (MDP), which consists of

- \mathcal{S} a measurable space of states.
- \mathcal{A} a measurable space of actions.
- $P : \mathcal{S} \times \mathcal{A} \rightsquigarrow \mathcal{S}$ a transition kernel².
- $R : \mathcal{S} \times \mathcal{A} \rightsquigarrow \mathbb{R}$ a reward kernel discounted by
- a discount factor $\gamma \in [0, 1)$.
- $\mathfrak{A}(s) \subseteq \mathcal{A}$ a set of admissible actions for each $s \in \mathcal{S}$.

²Here \rightsquigarrow denotes a *stochastic mapping* (to be defined soon).



Examples of MDPs

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Examples of Markov decision processes include

- Board games where one plays against a fixed opponent, e.g. *chess* where the set of states \mathcal{S} is the set of all obtainable chess-positions.
- Time-descretized physics simulations with action inputs and reward outputs, including most single player video games and the classic *cartpole* example (balancing a stick).



The probability kernels

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Probability kernel

A **probability kernel** (also called a *stochastic mapping*, *stochastic kernel* or *Markov kernel*) $\kappa : \mathcal{X} \rightsquigarrow \mathcal{Y}$ is a collection of probability measures $\kappa(\cdot \mid x)$, one for each $x \in \mathcal{X}$ such that for any measurable set $B \subseteq \mathcal{Y}$ the function $x \mapsto \kappa(B \mid x)$ is measurable.

The transition probability measure $P(\cdot \mid s, a)$ of the pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ determines what states are likely to follow after *being* in state s and *choosing* action a . Similarly from the reward kernel R one obtains the measure $R(\cdot \mid s, a)$ determining the reward distribution following the timestep (s, a) .



Policies

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Given a Markov decision process one can define a **policy** π by sequence of probability kernels $\pi = (\pi_1, \pi_2, \dots)$ where $\pi_i : \mathcal{H}_i \rightsquigarrow \mathcal{A}$ and $\mathcal{H}_i = \mathcal{S} \times \mathcal{A} \times \dots \times \mathcal{S}$ is the *history space* at the i th timestep.



Stochastic processes

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

An MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ together with a policy $\pi = (\pi_1, \pi_2, \dots)$ and a distribution μ on \mathcal{S} give rise to a stochastic process $(S_1, A_1, S_2, A_2, \dots) \sim \kappa_\pi \mu$ such that for any $i \in \mathbb{N}$ we have $(S_1, A_1, \dots, S_i) \sim P\pi_{i-1} \dots P\pi_1 \mu$ where $P\pi_{i-1} \dots P\pi_1$ denotes the *kernel-composition* of the probability kernels $P, \pi_1, \dots, \pi_{i-1}$. We denote by \mathbb{E}_s^π expectation over $\kappa_\pi \mu$ where $\mu = \delta_s$, that is, $S_1 = s$ a.s.



Policy evaluation

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

For a policy π we can define the policy evaluation function:

Policy evaluation

Denote by $r(s, a) = \int x \, dR(x \mid s, a)$ the *expected reward function*.

We define the **policy evaluation function** by

$$V_{\pi}(s) = \mathbb{E}_s^{\pi} \sum_{i=1}^{\infty} \gamma^{i-1} r \circ \rho_i$$

where ρ_i is projection onto $(\mathcal{S}_i, \mathcal{A}_i)$.

This an example of a (state-) *value function*, as it assigns a real number to every state $s \in \mathcal{S}$.



Finite policy evaluation

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Similar to the infinite horizon policy evaluation we can also consider a finite horizon version:

Definition: Finite policy evaluation

We define the function $V_{n,\pi} : \mathcal{S} \rightarrow \mathbb{R}$ by

$$V_{n,\pi}(s) = \mathbb{E}_s^\pi \sum_{i=1}^n \gamma^{i-1} r \circ \rho_i$$

called the k th **finite policy evaluation**^a.

^aWhen $n = 0$ we say $V_{0,\pi} = V_0 := 0$ for any π .



Optimal value function

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model-
dependent
algorithms

Model-free
algorithms

Conclusion

Definition: Optimal value functions

$$V_n^*(s) := \sup_{\pi \in R\Pi} V_{n,\pi}(s) = \sup_{\pi \in R\Pi} \mathbb{E}_s^\pi \sum_{i=1}^n r_i$$

$$V^*(s) := \sup_{\pi \in R\Pi} V_\pi(s) = \sup_{\pi \in R\Pi} \mathbb{E}_s^\pi \sum_{i=1}^{\infty} r_i$$

This is called the **optimal value function** (and the n th optimal value function). A policy $\pi^* \in R\Pi$ for which $V_{\pi^*} = V^*$ is called an **optimal policy**. If $V_{n,\pi^*} = V_n^*$ then π^* is called n -optimal.

Provided such an optimal policy π^* exists, obtaining such a policy is the ultimate goal of Reinforcement Learning.



Greediness

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

In order to show existence of optimal policies and talk about algorithms which can determine such policies, we define the concept of *greediness*.



Greedy actions

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

The purpose of (most) value functions $V : \mathcal{S} \rightarrow \mathbb{R}$ is to give an estimate on how *good* a certain state is, in terms of the rewards one may expect after visiting it.

This give rise to the idea of *greedy actions*, that is, actions leading to states high *values* (according to V).

Definition greedy actions

Let $V : \mathcal{S} \rightarrow \mathbb{R}$ be a measurable value-function. We define

$$G_V(s) = \operatorname{argmax}_{a \in \mathcal{A}(s)} T_a V(s) \subseteq \mathcal{A}(s)$$

as the set of **greedy** actions w.r.t. V .



Greedy policies

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Greedy actions leads to *greedy policies*:

Definition: Greedy policy

Let $V : \mathcal{S} \rightarrow \mathbb{R}$ be a measurable value-function and let $\tau : \mathcal{S} \rightsquigarrow \mathcal{A} \in \mathcal{SII}$ be a stationary policy. If there exists a measurable $G_V^\tau(s) \subseteq G_V(s)$ such that

$$\tau(G_V^\tau(s) \mid s) = 1$$

for every $s \in \mathcal{S}$, then τ is called greedy w.r.t. V . We will often denote a V -greedy policy by τ_V .



Existence of greedy policies

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Theorem (Existence of greedy policies)

Suppose $V : \mathcal{S} \rightarrow \mathbb{R}$ is *upper semicontinuous* and that

1. \mathcal{S} and \mathcal{A} are standard Borel.
2. The set of admissible actions $\mathfrak{A}(s) \subseteq \mathcal{A}$ is compact for all $s \in \mathcal{S}$ and $\Gamma = \{(s, a) \in \mathcal{S} \times \mathcal{A} \mid a \in \mathfrak{A}(s)\}$ is a closed subset of $\mathcal{S} \times \mathcal{A}$.
3. The transition kernel P is continuous.
4. The expected reward function $r = \int r' \, dR(r' \mid \cdot)$ is upper semicontinuous and bounded from above.

Then there exists a deterministic policy π_V which is greedy for V .

If assumptions 1.-4. hold we will say that the MDP is *greedy*.



Policy iteration

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

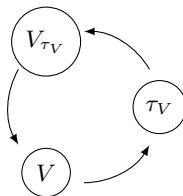
Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Using the concepts we have defined one can get the following idea:
Iteratively generate value functions by picking greedy policies and then
evaluating these policies. This is called *policy iteration*.



Policy iteration is a well studied algorithm and can be shown to converge to optimum for a variety of environments. We will however move on to talk about a related concept called *value iteration*.



Operators on value functions

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Before defining value iteration we introduce some operators

The T -operators

For a stationary policy $\tau \in \mathcal{S}\Pi$ and a value function $V : \mathcal{S} \rightarrow \mathbb{R} \in \mathcal{L}_\infty(\mathcal{S})$ we define the operators

The policy evaluation operator:

$$T_\tau V := s \mapsto \int r(s, a) + \gamma V(s') \, d(P\tau)(a, s' \mid s)$$

The Bellman optimality operator:

$$TV := s \mapsto \sup_{a \in \mathcal{A}(s)} \left(r(s, a) + \gamma \int V(s') \, dP(s' \mid s, a) \right)$$



Properties of the T -operators

Theoretical
aspects of
Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model-
dependent
algorithms

Model-free
algorithms

Conclusion

Proposition (Properties of the T -operators)

$$V_{k,\pi} = T_{\tau_1} V_{k-1,(\tau_2,\dots)} = T_{\tau_1} \dots T_{\tau_k} V_0.$$

$$V_\pi = \lim_{k \rightarrow \infty} T_{\tau_1} \dots T_{\tau_k} V_0$$

For the stationary policy τ we have $T_\tau V_\tau = V_\tau$.

T and T_τ are γ -contractive on $\mathcal{L}_\infty(\mathcal{S})$.

V_τ is the unique bounded fixed point of T_τ in $\mathcal{L}_\infty(\mathcal{S})$.

This way T_τ can be interpreted as a 'one-step policy evaluation, when following the policy τ '. On the other hand T can be interpreted as a 'one-step evaluation, when always choosing greedy actions'.



Value iteration

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model-dependent algorithms

Model-free algorithms

Conclusion

Value iteration is the iterative application of the T -operator. The following theorem show why value iteration is a central idea Reinforcement Learning³.

Theorem (Existence optimal policies & convergence of value iteration)

Given a greedy MDP we have that

$$V_k^* = T^k V_0 = T_{\tau_{k-1}^*} \dots T_{\tau_0^*} V_0 = V_{k, (\tau_{k-1}^*, \dots, \tau_0^*)}$$

The policy $(\tau_{k-1}^*, \dots, \tau_0^*)$ is a deterministic k -optimal policy where $\tau_k^* = \tau_{T^k V_0}$ is any deterministic greedy policy for $T^k V_0$ for any $k \in \mathbb{N}$. Furthermore $V^* = \lim_{k \rightarrow \infty} T^k V_0^*$, the greedy policy $\tau^* = \tau_{V^*}$ exists and an optimal policy.

³Actually value iteration is inherited from dynamic programming.



Convergence rates

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

We can also show that the optimal value function V^* is a fixed point of the Bellman optimality operator T .

$$TV^* = V^*$$

This is often called *Bellman's optimality equation*.

Recalling that T is γ -contractive, by Banach's fixed point theorem we get exponential convergence rates for value iteration:

$$\|T^k V - V^*\| \leq \gamma^k \|V - V^*\|_\infty = \mathcal{O}(\gamma^k)$$



Example: Gridworld

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model-

dependent algorithms

Model-free algorithms

Conclusion

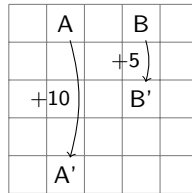
The *gridworld*

MDP consist of 25 states $\mathcal{S} = [5]^2$ and 4 actions $\mathcal{A} = \{U, D, L, R\}$

for *up*, *down*, *left* and *right*

and moves the agent 1 square up, down, left or right. A reward of 0 is given by default, except when

- *hitting the boundary* a reward of -1 is given
- when in $A = (2, 1)$ any action moves to $A' = (2, 5)$ and is rewarded 10.
- when in $B = (4, 1)$ any action moves to $B' = (4, 3)$ and is rewarded 5.



Finally $\gamma = 0.9$ is the standard value of the discount factor in this example.



Example: Gridworld

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model-

dependent algorithms

Model-free algorithms

Conclusion

-0.50	10.00	-0.25	5.00	-0.50
-0.25	0.00	0.00	0.00	-0.25
-0.25	0.00	0.00	0.00	-0.25
-0.25	0.00	0.00	0.00	-0.25
-0.50	-0.25	-0.25	-0.25	-0.50

V_1, τ_r

3.31	8.79	4.43	5.32	1.49
1.52	2.99	2.25	1.91	0.55
0.05	0.74	0.67	0.36	-0.40
-0.97	-0.44	-0.35	-0.59	-1.18
-1.86	-1.35	-1.23	-1.42	-1.98

V_{400}, τ_r

Figure: Policy evaluations of the gridworld environment. Note that $V_{\max} \cdot \gamma^{400} = 100 \cdot (0.9)^{400} \approx 4.97 \cdot 10^{-17}$ so $V_{\tau_r, 400}$ are very close to the true infinite horizon value functions V_{τ_r} (providing numerical errors are insignificant).



Example: Gridworld

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model-dependent algorithms

Model-free algorithms

Conclusion

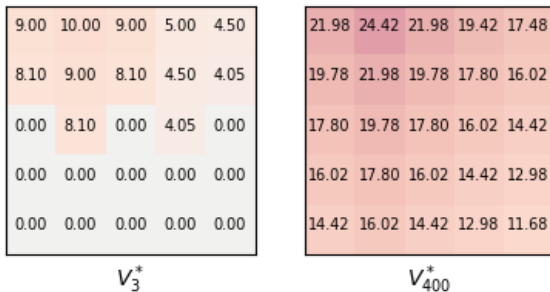


Figure: Optimal value functions of the gridworld environment. By the same upper bound as before we have $\|V^* - V_{400}^*\|_\infty < 4.97 \cdot 10^{-17}$.



Example: Gridworld

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

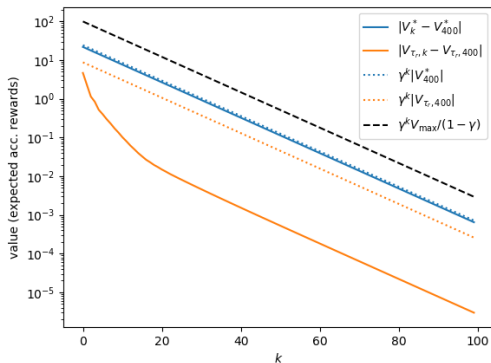
Q-iteration

Model-dependent algorithms

Model-free algorithms

Model-free algorithms

Conclusion



Convergence of gridworld value functions compared with the theoretical bounds. The black dashed line is the general theoretical bound for both T and T_τ by Banachs fixed point theorem and the maximum value $V_{\max} = R_{\max}/(1 - \gamma)$. The dotted blue and orange uses $|V_k^*|$ and $|V_{\tau,k}|$ respectively, which might not be available. ($\gamma = 0.9$).



Q-functions

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

A **Q-function** is simply any function assigning a real number to every state-action pair. They are also called (state-) *action value functions*.

A **Q-learning** algorithm is any algorithm which uses Q-functions to derive a policy for an environment⁴.

⁴Some authors refer to Q-learning as a specific variation of temporal difference learning, but this fails to capture many algorithms which are also referred to as *Q-learning algorithms*.



Motivation for Q-functions

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

A clear advantage of working with Q-function $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ rather than a value function $V : \mathcal{S} \rightarrow \mathbb{R}$, is that finding the optimal action $a^* \in \mathcal{A}(s)$ at state s requires only a maximization over the Q-function itself:

$a^* = \operatorname{argmax}_{a \in \mathcal{A}(s)} Q(s, a)$. This should be compared to finding an optimal action according to a value function V :

$$a^* = \operatorname{argmax}_{a \in \mathcal{A}(s)} r(s, a) + \gamma \mathbb{E}_{P(\cdot|s,a)} V.$$



Greed with Q-functions

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Formally we define greedy actions and policies w.r.t. a Q-function as Let $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ be a measurable Q-function and $\tau : \mathcal{S} \rightsquigarrow \mathcal{A}$ be a (stationary) policy.

Greedy policy

Define the set of *greedy actions* by $G_Q(s) := \operatorname{argmax}_{a \in \mathcal{A}(s)} Q(s, a)$. If there exist a measurable set $G_Q^\tau(s) \subseteq G_Q(s)$ for every $s \in \mathcal{S}$ such that

$$\tau \left(G_Q^\tau(s) \mid s \right) = 1$$

then τ is said to be **greedy** with respect to Q and is denoted τ_Q .



Q-function operators

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model-

dependent algorithms

Model-free algorithms

Conclusion

Moreover we define T -operators similar to ones for value functions

Operators for Q-functions

For any stationary policy $\tau \in \mathcal{S}\Pi$ and integrable Q-function

$Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} \in \mathcal{L}_\infty(\mathcal{S} \times \mathcal{A})$ we define

Next-step operator:

$$P_\tau Q(s, a) = \int Q(s', a') \, d\tau P(s', a' \mid s, a)$$

Policy evaluation operator:

$$T_\tau Q(s, a) = r(s, a) + \gamma \int Q(s', a') \, d\tau P(s', a' \mid s, a)$$

Bellman optimality operator:

$$TQ(s, a) = r(s, a) + \gamma \int \max_{a' \in \mathcal{A}} Q(s', a') \, dP(s' \mid s, a)$$

where $T_a = T_{\delta_a}$.



Relation between value- and Q-functions

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Theorem (Relations between Q- and value functions)

Let $\pi = (\tau_1, \tau_2, \dots) \in M\Pi$ be a Markov policy and $\tau \in S\Pi$ stationary.

Then

- Policy evaluations are related by $\mathbb{E}_{\tau(\cdot|s)} Q_{k,\pi} = V_{k+1,(\tau,\pi)}(s)$.
- T_{τ} -operators are related by $T_{\tau} Q_{k,\pi}(s, a) = r + \gamma \mathbb{E}_{P(\cdot|s,a)} T_{\tau} V_{k,\pi}$.
- τ is greedy for $Q_{k,\pi}$ if and only if τ is greedy for $V_{k,\pi}$ and τ is greedy for Q_{π} if and only if τ is greedy for V_{π} .
- Optimal policies are related by $\max_{a \in \mathfrak{A}(s)} Q^*(s, a) = V^*(s)$ and
$$Q_k^*(s, a) = r(s, a) + \gamma \mathbb{E}_{P(\cdot|s,a)} V_k^*, \quad Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{P(\cdot|s,a)} V^*$$



Properties of Q-functions

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Because of the close relations many properties are inherited from value function to Q-functions:

Proposition (Properties of Q-functions)

Let $\pi = (\tau_1, \tau_2, \dots) \in M\Pi$ be a Markov policy and $\tau \in S\Pi$ stationary. Then

$$Q_{k,\pi} = T_{\tau_1} \dots T_{\tau_k} Q_0 \text{ and } Q_k^* = T^k Q_0^*.$$

$$Q_\pi = \lim_{k \rightarrow \infty} Q_{k,\pi} \text{ and } Q^* = \lim_{k \rightarrow \infty} Q_k^*.$$

T, T_τ are γ -contractive on $\mathcal{L}_\infty(\mathcal{S} \times \mathcal{A})$ and Q^*, Q_τ are their unique fixed points.

$$Q^* = Q_{\tau^*}$$



Q-iteration

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Q-iteration is the analogue of value-iteration for Q-function. It can be stated in the form of an algorithm as follows:

Algorithm (Q-iteration)

Data: MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$, number of iterations K

Initialize expected reward function $r \leftarrow \int x \, dR(x \mid \cdot)$ and $\tilde{Q}_0 \leftarrow r$.

for $k = 0, 1, \dots, K - 1$ **do**

$\tilde{Q}_{k+1} \leftarrow T\tilde{Q}_k$

end

Output: \tilde{Q}_K

In the context of a greedy MDP we immediately have that the output of the Q-iteration algorithm $\tilde{Q}_K = Q_K^*$ is K -optimal.



Value iteration with Q-functions

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

Similar to value-iteration we can use Banach fixed point theorem with the contractive properties of the T -operator for Q-functions to obtain exponential convergence of Q-iteration:

Proposition (Convergence of Q-iteration)

Suppose the Q-iteration algorithm is run with a greedy MDP. Then the output $\tilde{Q}_K = Q_K^*$ satisfy

$$\|Q^* - Q_K^*\|_\infty \leq \gamma^K V_{\max}$$



What have we done so far?

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

For greedy MDPs we have proved

- Existence of optimal policies.
- Exponential convergence of Q-iteration to the optimal policy.



Are we not done?

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and
the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

We have exponential convergence for the broad class of problems expressible as a greedy MDP. This class includes highly difficult environments such as control problems in time-descretized simulation environments such as computer games, including the game of *chess*. Are we then done?



Problems of Q-iteration

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model-dependent algorithms

Model-free algorithms

Conclusion

Model-dependency: It is assumed that we know how to integrate over P and R .

- The distributions of P and R might be impractical to work with in a computer.
- It is common in RL to assume that P and R are unknown, thus including a variety of environments, which we have not yet considered.

Representation infeasibility: It is assumed that we know how to represent Q functions in a feasible way in a computer.

- Representing the functions arising from use of the T -operator may be difficult as these functions are defined by successive integration over potentially complex transition kernels.
- Even in the finite case where Q can be represented as a table, this table may be excessively large.



Example: Chess

Theoretical aspects of Q-learning

Introduction

The environment

Value functions and the goal of RL

Value iteration

Q-functions

Q-iteration

Model- dependent algorithms

Model-free algorithms

Conclusion

The state space of chess is very large (roughly $|\mathcal{S}_{\text{chess}}| \geq 10^{43}$). This means that if we were to use Q-iteration naively (with finite implementation as in the gridworld example) then we would have to store a vector of roughly $N \cdot 10^{43}$ real numbers for each Q-function we define, where N is the average number of admissible actions at each state $\mathcal{A}(s)$, $s \in \mathcal{S}$ which has been estimated to around $N \approx 35$ for chess. This requires roughly $1.4 \cdot 10^{45}$ bytes, if each number is stored as a single precision floating point number (4 bytes). For comparison the entire digital data capacity in the world is estimated less than 10^{23} bytes as of 2020. Needless to say this is beyond any practical relevance.



Model-dependent algorithms

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

To deal with the problem of representation infeasibility we will investigate what happens when using approximations of TQ in the Q-iteration algorithm. The hope is then that we can choose an class of approximation-functions

$$\mathcal{F} \subseteq \mathcal{Q} = \{Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\}$$

that is both *representable* (in a computer) and *dense* (to some degree) in $T\mathcal{F} = \{TQ \mid Q \in \mathcal{Q}\}$.



Approximation and algorithmic errors

Theoretical aspects of Q-learning

Introduction

Model- dependent algorithms

Approximation and algorithmic errors

Approximation using ANNs

Approximation with Bernstein polynomials

Model-free algorithms

Conclusion

Let $\|\cdot\|$ be any norm the space of Q-functions \mathcal{Q} . Suppose we can approximate $T\tilde{Q}_0$ by some $\tilde{Q}_1 \in \mathcal{F}$ to $\varepsilon_1 > 0$ precision and then approximate $T\tilde{Q}_1$ by $\tilde{Q}_2 \in \mathcal{F}$ and so on. This way we get a sequence of Q-functions satisfying

$$\left\| T\tilde{Q}_{k-1} - \tilde{Q}_k \right\| \leq \varepsilon_k, \forall k \in \mathbb{N}$$



Approximation and algorithmic errors

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

First observe that

$$\begin{aligned}\|T^k \tilde{Q}_0 - \tilde{Q}_k\| &\leq \|T^k \tilde{Q}_0 - T \tilde{Q}_{k-1}\| + \|T \tilde{Q}_{k-1} - \tilde{Q}_k\| \\ &\leq \gamma \|T^{k-1} \tilde{Q}_0 - \tilde{Q}_{k-1}\| + \|T \tilde{Q}_{k-1} - \tilde{Q}_k\|\end{aligned}$$

Using this iteratively we get

$$\|T^k \tilde{Q}_0 - \tilde{Q}_k\| \leq \sum_{i=1}^k \gamma^{k-i} \varepsilon_i$$

Using this we can bound ...



Approximation and algorithmic errors

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

$$\begin{aligned}\|Q^* - \tilde{Q}_k\| &\leq \|Q^* - T^k \tilde{Q}_0\| + \|T^k \tilde{Q}_0 - \tilde{Q}_k\| \\ &\leq \gamma^k \|Q^* - \tilde{Q}_0\| + \sum_{i=1}^k \gamma^{k-i} \varepsilon_i\end{aligned}$$

- The first term is called the **algorithmic error** and decreases exponentially with k , so we will be content with this.
- The second term $\sum_{i=1}^k \gamma^{k-i} \varepsilon_i$ is where the trouble from our approximative strategy arises and is therefore called the **approximation error**. We denote it $\varepsilon_{\text{approx}}(k)$. It is this the approximation error which we will now struggle to bound.



Theoretical approximation errors

Theoretical aspects of Q-learning

Introduction

Model-dependent algorithms

Approximation and algorithmic errors

Approximation using ANNs

Approximation with Bernstein polynomials

Model-free algorithms

Conclusion

We here give a few examples of how the approximation error behave in relation to the step-wise errors ε_i .

- If $\varepsilon_i(k) = \varepsilon$ for some $\varepsilon > 0$ we easily get the bound

$$\varepsilon_{\text{approx}}(k) = \varepsilon \frac{1 - \gamma^k}{1 - \gamma} \leq \frac{\varepsilon}{1 - \gamma}$$

- If $\varepsilon_i \leq c\gamma^i$ we get

$$\varepsilon_{\text{approx}}(k) \leq ck\gamma^k \rightarrow 0$$

as $k \rightarrow \infty$.

- Generally we have $\sum_{i=1}^k \gamma^{k-i} \varepsilon_i \rightarrow 0$ whenever $\varepsilon_k \rightarrow 0$ as $k \rightarrow \infty$.



Artificial Neural Networks

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

We will now consider the approach of (model-dependent) Q-iteration using artificial neural networks (ANNs) as function approximators.

Definition: Artificial neural network

An **artificial neural network** (ANN) with $L \in \mathbb{N}_0$ hidden layers, structure $(d_i)_{i=0}^{L+1} \subseteq \mathbb{N}$, activation functions $(\sigma_i)_{i=1}^L$, weights $(W_i)_{i=1}^{L+1} \in M^{d_i \times d_{i-1}}$ and biases $(v_i)_{i=1}^{L+1} \in \mathbb{R}^{d_i}$ is the function $f : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{d_{L+1}}$ defined by

$$f = w_{L+1} \circ \sigma_L \circ w_L \circ \sigma_{L-1} \circ \cdots \circ \sigma_1 \circ w_1$$

where w_i is the affine function $x \mapsto W_i x + v_i$, and $\sigma_i : \mathbb{R}^{d_i} \rightarrow \mathbb{R}^{d_i}$ is coordinate-wise application of components $\sigma_{ij} : \mathbb{R} \rightarrow \mathbb{R}$. We denote the class of these networks (or functions)

$$\mathcal{DN} \left((\sigma_i)_{i=1}^L, (d_i)_{i=0}^{L+1} \right)$$

An ANN is called *deep* if there are two or more hidden layers.



ANNs as graphs

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

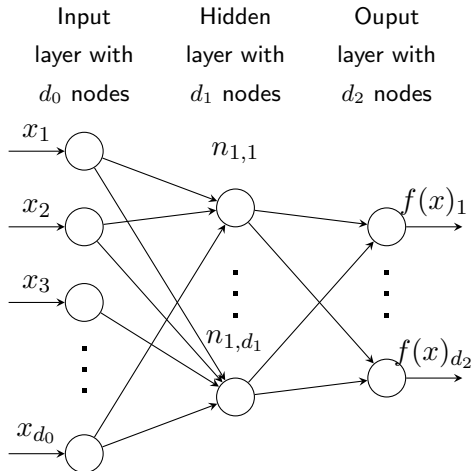


Figure: An ANN with one hidden layer ($L = 1$). Notice that there is no edge from $n_{0,3}$ to $n_{1,1}$ which means that $W_1(1,3) = 0$.



Universal approximation theorem for ANNs

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

Theorem (Universal approximation theorem for ANNs)

Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be non-constant, bounded and continuous activation function. Let $\varepsilon > 0$ and $f \in C([0, 1]^w)$. Then there exists an $N \in \mathbb{N}$ and a network $F \in \mathcal{DN}(\sigma, (w, N, 1))$ with one hidden layer, unbiased final layer (that is $v_2 = 0$) and activation function σ such that

$$\|F - f\|_{\infty} < \varepsilon$$

In other words $\bigcup_{N \in \mathbb{N}} \mathcal{DN}(\sigma, (w, N, 1))$ is dense in $C([0, 1]^w)$.



Notation and setting

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

Rectified linear unit (ReLU) activation function:

$$\sigma_r(x) = \max(0, x)$$

Class of *ReLU networks*:

$$\mathcal{RN} \left((d_i)_{i=0}^{L+1} \right) := \mathcal{DN} \left(\sigma_r, (d_i)_{i=0}^{L+1} \right)$$

Setting: Greedy MDP with state-space $\mathcal{S} = [0, 1]^w \subseteq \mathbb{R}^w$ inside the unit (hyper-) cube, finite action space $|\mathcal{A}| < \infty$ and expected reward function r being (fully) *continuous*.



Q-iteration with ANN-approximation - bound

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

Theorem

Let $\varepsilon > 0$. Assume that one of the following two conditions hold:

1. P is deterministic with $P(\cdot \mid s, a) = \delta_{p(s,a)}$. For some continuous $p : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$.
2. $P(\cdot \mid s, a)$ is absolutely continuous with respect to the same measure ν on \mathcal{S} for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ with density $p(\cdot \mid s, a)$ which is continuous in s .

Then for every $k \in \mathbb{N}$ there exists a $N \in \mathbb{N}$ and a sequence of Q-networks $(\tilde{Q}_i)_{i=1}^k \subseteq \mathcal{RN}(w|\mathcal{A}|, N, 1)$ such that

$$\varepsilon_{\text{approx}}(i) = \left\| T\tilde{Q}_{i-1} - \tilde{Q}_i \right\|_{\infty} < \varepsilon$$

for all $i \in [k]$. In particular

$$\left| Q^* - \tilde{Q}_k \right| < \gamma^k V_{\max} + \varepsilon / (1 - \gamma)$$



Q-iteration with ANN: conclusion

Theoretical aspects of Q-learning

Introduction

Model- dependent algorithms

Approximation and algorithmic errors

Approximation using ANNs

Approximation with Bernstein polynomials

Model-free algorithms

Conclusion

What we have accomplished: Arbitrarily precise approximation of Q^* using a concrete class of approximators, namely ANNs.

What is still missing:

- How large must the network depth N be? This requires going through a constructive proof of the universal approximation theorem, which is available in literature.
- How to obtain the each approximator \tilde{Q}_i ? Optimization using ANNs are extensively studied, so this might be answered in literature.



Bernstein polynomials

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

We now turn to a new approximation technique:

Definition: Bernstein polynomial

The (multivariate) Bernstein polynomial $B_{f,n}$ of degree $n = (n_1, \dots, n_w) \in \mathbb{N}^w$ approximating the function $f : [0, 1]^w \rightarrow \mathbb{R}$ is defined by

$$B_{f,n}(x_1, \dots, x_w) = \sum_{j=1}^w \sum_{k_j=0}^{n_j} f\left(\frac{k_1}{n_1}, \dots, \frac{k_w}{n_w}\right) \prod_{\ell=1}^w \binom{n_\ell}{k_\ell} x_\ell^{k_\ell} (1 - x_\ell)^{n_\ell - k_\ell}$$



Approximation with Bernstein polynomials

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

Theorem: Approximation properties of Bernstein polynomials

Let $f : [0, 1]^w \rightarrow \mathbb{R}$ be Lipschitz w.r.t. the standard euclidean 2-norm induced metrics on $[0, 1]^w$ and \mathbb{R} with constant L . Let $n = (n_1, \dots, n_w) \in \mathbb{N}^w$. The Bernstein polynomial

$B_{f,n} : [0, 1]^w \rightarrow \mathbb{R}$ satisfies

1. $\|f - B_{f,n}\|_{\infty} \leq \frac{L}{2} \sqrt{\sum_{j=1}^w \frac{1}{n_j}}$
2. $\|B_{f,n}\|_{\infty} \leq \|f\|_{\infty}$



Q-iteration with Bernstein polynomials

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

Proposition

Suppose we are given a greedy MDP with $\mathcal{S} = [0, 1]^w$ and finite action space. Assume that there exists a probability measure $\mu \in \mathcal{P}(\mathcal{S})$ such that $P(\cdot \mid s, a)$ has density $p(\cdot \mid s, a) : \mathcal{S} \rightarrow \mathbb{R}$ w.r.t μ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Furthermore assume that $r(\cdot, a)$ and $p(s \mid \cdot, a)$ are $\|\cdot\|_2$ -Lipschitz with constants L_r, L_p for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

Using Bernstein polynomials of degree $n = (m, \dots, m) \in \mathbb{N}^w$ for $m \in \mathbb{N}$ we have the following bound

$$\left\| Q^* - \tilde{Q}_k \right\|_{\infty} \leq \gamma^k V_{\max} + \frac{L_r + \gamma V_{\max} L_p}{2(1 - \gamma)} \sqrt{w} \frac{1}{\sqrt{m}}$$

In particular $\left\| Q^* - \tilde{Q}_k \right\|_{\infty} = \mathcal{O}(\gamma^k + \frac{1}{\sqrt{m}})$ when using k iterations.



Conclusion on Bernstein polynomials

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Approximation and
algorithmic errors

Approximation using
ANNs

Approximation with
Bernstein
polynomials

Model-free
algorithms

Conclusion

What we have achieved: $\mathcal{O}(\gamma^k + 1/\sqrt{m})$ convergence bound for Q-iteration with (m, \dots, m) -degree Bernstein polynomial approximation. This solves two problems with approximation with ANNs of

1. the function class not being concrete
2. not knowing how to obtain each approximator \tilde{Q}_i

We mention two weak points

1. The restriction on P : For example we cannot use the bound on deterministic decision processes, since if P is deterministic then there are no measure $\mu \in \mathcal{P}(\mathcal{S})$ which allows for a density $p(\cdot \mid s, a)$ (i.e. $p \cdot \mu = P(\cdot \mid s, a)$), unless $P(\cdot \mid s, a) = \delta_{s'}$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, which would lead to a quite boring environment. Generally the processes with fast convergence bounds must be very stochastic.
2. We lack understanding of the computational complexity involved.



Model-free algorithms: Overview

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion

We deal with the model-free algorithm of this thesis in two parts:

1. The first part (section 3.1 and 3.2) is a survey on various old and recent convergence bounds. Most bounds are either confined to finite processes (section 3.1) or not establishing concrete bounds (section 3.2).
2. The second section (section 3.3) is a detailed presentation of a preprint [Fan et al. (2020+)] proving convergence of a Q-learning algorithm in a continuous state-space setting.



What is model-free algorithms?

Theoretical aspects of Q-learning

Introduction

Model- dependent algorithms

Model-free algorithms

Finite case

History dependent setting

Linear function approximation

Deep fitted Q-iteration

Conclusion

In model-free algorithms are based on *sampling* from the environment dynamics instead of using these distributions directly (e.g. by taking expectations).

This means that we cannot use any of the approaches so far considered without modification because taking expectations over the kernels P and R is no longer possible (directly).



Motivating model-free algorithms

Theoretical aspects of Q-learning

Introduction

Model- dependent algorithms

Model-free algorithms

Finite case

History dependent setting

Linear function approximation

Deep fitted Q-iteration

Conclusion

Why put the model-free restriction upon our algorithms?

- Big actors in both research and industry (such as *Google DeepMind*) favor model-free approaches. Thus much research goes in this direction.
- Some environments may be hard to model, but easy to sample from, such as the stock-market.
- *Works out of the box*: No need for adaptation to each model as only access to samples are required.
- Philosophical interest: Are humans model-free algorithms?
- Experience that Monte Carlo approaches can lead to faster computation.



Non-example: Model-free Q-iteration

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion

Model-free Q-iteration (example of thought)

Input: MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$, number of iterations K

For all $(s, a) \in \mathcal{S} \times \mathcal{A}$ let $\tilde{Q}_0(s, a)$ be sampled from $R(\cdot \mid s, a)$.

for $k = 0, 1, 2, \dots, K - 1$ **do**

 For all $(s, a) \in \mathcal{S} \times \mathcal{A}$ sample a reward r' from $R(\cdot \mid s, a)$, a next-state s'
 from $P(\cdot \mid s, a)$ and let

$$\tilde{Q}_{k+1}(s, a) \leftarrow r' + \gamma \sup_{a' \in \mathcal{A}} \tilde{Q}_k(s', a')$$

end

Define π_K as the greedy policy w.r.t. \tilde{Q}_K

Output: An estimator \tilde{Q}_K of Q^* and policy π_K

Problems?

- Uncountable \mathcal{S} or \mathcal{A} does not make much sense: How are the estimations \tilde{Q}_i well-defined?
- Even in finite case: Convergence is not possible for non-deterministic processes.



Direct vs. indirect approaches

Theoretical aspects of Q-learning

Introduction

Model- dependent algorithms

Model-free algorithms

Finite case

History dependent setting

Linear function approximation

Deep fitted Q-iteration

Conclusion

There are broadly two ways of dealing with these problems. In the *indirect* approaches one tries to first estimate P and R by sampling. We can then use the model-dependent methods of the last chapter with estimated kernels \tilde{P} and \tilde{R} . We will not go further into indirect methods in this thesis. The *direct* approaches covers *the rest*, and it is mainly these were are going to look at throughout this thesis.



The finite setting

Theoretical aspects of Q-learning

Introduction

Model- dependent algorithms

Model-free algorithms

Finite case

History dependent setting

Linear function approximation

Deep fitted Q-iteration

Conclusion

Finite action-uniform MDP

We consider a (discounted) MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ where $|\mathcal{S}| = s \in \mathbb{N}$ and $|\mathcal{A}| = a \in \mathbb{N}$ are finite. Also we assume unrestricted actions, i.e. that the set of admissible actions $\mathcal{A}(s) = \mathcal{A}$ is the same for all $s \in \mathcal{S}$

This setting is a special case of a greedy MDP, and thus we get existence of optimal value functions and convergence of $Q_k^* \rightarrow Q^*$.



TD-learning

Theoretical aspects of Q-learning

Introduction

Model- dependent algorithms

Model-free algorithms

Finite case

History dependent setting

Linear function approximation

Deep fitted Q-iteration

Conclusion

A popular way of dealing with the problems of naive model-free Q-iteration, is called *temporal difference learning* (TD-learning) and is based on the following *update-step*:

$$\tilde{Q}_{k+1}(s, a) \leftarrow (1 - \alpha_k) \tilde{Q}_k(s, a) + \alpha_k (r' + \gamma \cdot \max_{a' \in \mathcal{A}} \tilde{Q}_k(s', a'))$$

Where

- r' and s' are the reward and next-state drawn from the reward and transition kernels,
- $\alpha_k \in [0, 1]$ is called the **learning rate**,
- $r' + \gamma \max_{a' \in \mathcal{A}} \tilde{Q}_k(s', a')$ is called the **T-value**, due to its similarity to $T\tilde{Q}_k(s, a) = r(s, a) + \gamma \int \max_{a' \in \mathcal{A}} \tilde{Q}_k(s', a') dP(s' | s, a)$,
- the **temporal difference** term
 $\alpha_k (r' + \gamma \cdot \max_{a \in \mathcal{A}} \tilde{Q}_k(s', a') - \tilde{Q}_k(s, a))$.



Finite asynchronous Q-learning

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion

Algorithm: Finite asynchronous Q-learning

Input: Finite MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$, number of iterations K ,
state-action pairs $(s_1, a_1, \dots, s_K, a_K)$, learning rates
 $(\alpha_1, \dots, \alpha_K)$, initial $\tilde{Q}_0 : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$

for $k = 1, 2, \dots, K$ **do**

 Sample $r' \sim R(\cdot \mid s_k, a_k)$, $s' \sim P(\cdot \mid s_k, a_k)$.

 Update action-value function: For all $(s, a) \in \mathcal{S} \times \mathcal{A}$ let

$$\tilde{Q}_k(s, a) \leftarrow \begin{cases} \tilde{Q}_{k-1}(s, a) & (s, a) \neq (s_k, a_k) \\ (1 - \alpha_k)\tilde{Q}_{k-1} + \alpha_k(r' + \gamma \max_{a' \in \mathcal{A}} \tilde{Q}_{k-1}(s', a')) & (s, a) = (s_k, a_k) \end{cases}$$

end

Define $\tilde{\pi}_K$ as the greedy policy w.r.t. \tilde{Q}_K

Output: An estimator \tilde{Q}_K of Q^* and policy $\tilde{\pi}_K$



Convergence result

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion

Theorem (Watkins, Dayan 1992): Convergence of asynchronous QL

Let $s_1, a_1, s_2, a_2, \dots \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{A} \times \dots$ be random variables and $\alpha_1, \alpha_2, \dots \in [0, 1]$ be a sequence of learning rates. The output \tilde{Q}_K of the finite asynchronous Q-learning algorithm converges to Q^* provided that for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ it holds that

- $\mathbb{P} \left(\sum_{i=1}^{\infty} \alpha_i [(s_i, a_i) = (s, a)] = \infty \right) = 1$ and
 $\mathbb{P} \left(\sum_{i=1}^{\infty} \alpha_i^2 [(s_i, a_i) = (s, a)] < \infty \right) = 1.$
- $\mathbb{V}(R(\cdot \mid s, a)) < \infty$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}.$

Here $[(s_i, a_i) = (s, a)]$ is the Bernoulli random variable with $\mathbb{P}([(s_i, a_i) = (s, a)] = 1) = \mathbb{P}(s_i = s \text{ and } a_i = a).$



Bound on convergence

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion

Theorem (Szepesvári 1997)

Using the finite asynchronous algorithm let $K \in \mathbb{N}$ and $(s_1, a_1), (s_2, a_2) \dots, (s_K, a_K)$ be sampled i.i.d. from a distribution $p \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$ with $\text{supp}(p) = \mathcal{S} \times \mathcal{A}$, i.e. all state-action pairs have non-zero probability of occurring. Set the learning rates such that $\alpha_k = |\{i \in [k] \mid (s_i, a_i) = (s_k, a_k)\}|^{-1}$, i.e. they are the reciprocal of the frequency of (s_k, a_k) at step k . Let $\beta = \max_{x \in \mathcal{S} \times \mathcal{A}} p(x) / \min_{x \in \mathcal{S} \times \mathcal{A}} p(x)$. Then for some $B > 0$ the following holds asymptotically almost surely

$$\left| \tilde{Q}_K - Q^* \right| \leq B \frac{1}{K^{\beta(1-\gamma)}} \quad (1)$$

and

$$\left| \tilde{Q}_K - Q^* \right| \leq B \sqrt{\frac{\log \log K}{K}} \quad (2)$$

For the output \tilde{Q}_K of the finite asynchronous Q-learning algorithm with K -steps.

Note that eq. (1) is tightest when $\gamma > 1 - \beta/2$ otherwise eq. (2) is tighter.



History dependent setting: Introduction

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

**History dependent
setting**

Linear function
approximation

Deep fitted
Q-iteration

Conclusion

Staying in setting of decision processes with finite states and actions, we will now turn to the setting of history dependent processes, and present a result found in [Majeed and Hutter (2018)].



History dependent setting

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion

Define $\mathcal{H}^* = \bigcup_{i \in \mathbb{N}} \mathcal{H}_i$ as the space of all *finite histories*. We will write (history-) tuples $h_k = (s_1, a_1, s_2, \dots, a_{k-1}, s_k) \in \mathcal{H}_k$ as strings $h_k = s_1 a_1 s_2 \dots a_{k-1} s_k$ for convenience.

- Transition dynamics are history dependent

$$P : \mathcal{H}^* \times \mathcal{A} \rightsquigarrow \mathcal{S}.$$

- Rewards are history dependent, deterministic

$r : \mathcal{H}^* \times \mathcal{A} \rightarrow \mathbb{R}$, and discounted by $\gamma \in [0, 1)$, that is

$r(h_n, a) \in [-\gamma^{n-2} R_{\max}, \gamma^{n-2} R_{\max}]$ for any $h_n \in \mathcal{H}_n$ and $a \in \mathcal{A}$.

- Actions are unrestricted, so that $\mathfrak{A}(s) = \mathcal{A}$ for all $s \in \mathcal{S}$.



Remark on existence of optimal value functions

Theoretical aspects of Q-learning

Introduction

Model- dependent algorithms

Model-free algorithms

Finite case

History dependent setting

Linear function approximation

Deep fitted Q-iteration

Conclusion

The finite history dependent setting can be analysed with theorems developed by [Schäl, 1974], which considers history dependent processes generalized to arbitrary standard Borel \mathcal{S} and \mathcal{A} . This implies by the existence of an optimal $\pi^* \in R\Pi$ and that $V_n^* \rightarrow V^*$ ⁵.

⁵These remarks was not discussed in the original article by Majeed and Hutter.



History based Q-functions

Theoretical aspects of Q-learning

Introduction

Model- dependent algorithms

Model-free algorithms

Finite case

History dependent setting

Linear function approximation

Deep fitted Q-iteration

Conclusion

Q-functions are generalized so that they are taking values in $\mathcal{H}^* \times \mathcal{A}$. Likewise the T -operator is generalized by defining

$$TQ(h, a) := r(ha) + \gamma \sum_{s \in \mathcal{S}} \max_{a' \in \mathcal{A}} Q(has, a') P(s \mid ha)$$

The optimal Q-function Q^* is defined as the fixed point of the T operator in $\mathcal{L}_\infty(\mathcal{H}^* \times \mathcal{A})$.



Partial observability

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion

A function $\phi : \mathcal{H}^* \rightarrow \mathcal{X}$ is introduced which maps a history to a new finite space \mathcal{X} . The intuition here is that $x_n = \phi(h_{n-1}a_{n-1}s_n)$ is the state s_n as it is perceived by the agent. This is called **partial observability**. Using ϕ we are now considering a class of problems which is wider than a history dependent decision process (HDP). Namely a partially observable HDP or shortened: POHDP. A kernel for the observable process is defined for each $h \in \mathcal{H}^*$

$$p_h : \mathcal{A} \rightsquigarrow \mathcal{X}$$
$$p_h(x' \mid a) = \sum_{s: \phi(has)=x'} P(s \mid ha)$$



Theoretical aspects of Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion

We can now consider yet another kind of Q-functions defined on $\mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$, which we will call *partial* to avoid confusion. The optimal partial Q-function q_h^* , naturally dependent on a history $h \in \mathcal{H}^*$ and is defined by the equation

$$q_h^*(x, a) = r(h, a) + \gamma \sum_{x' \in \mathcal{X}} \max_{a' \in \mathcal{A}} q_h^*(x', a') p_h(x' \mid xa) \quad (3)$$



State-uniformity

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion

We can now introduce a central assumption:

Assumption: State-uniformity

For any $h, h' \in \mathcal{H}^*$ we have

$$\phi(h) = \phi(h') \implies Q^*(h, \cdot) = Q^*(h', \cdot)$$

Under this assumption the process is called (by Majeed and Hutter) a *Q-Value uniform decision process* (QDP).

State-uniformity is important because of the following theorem:

Theorem (Hutter, 2016)

Under the state-uniformity assumption we have

$$q_{h'}^*(\phi(h), a) = Q^*(h, a) \text{ for any } h' \in \mathcal{H}^*.$$



History dependent TD-learning

Theoretical aspects of Q-learning

Introduction

Model- dependent algorithms

Model-free algorithms

Finite case

History dependent setting

Linear function approximation

Deep fitted Q-iteration

Conclusion

Motivated by fact that the partial q -functions attains the (non-partial) optimum, we will use the asynchronous Q-learning algorithm with update step

$$q_{t+1}(x, a) = q_t(x, a) + \alpha_t(x, a) \left(r' + \gamma \max_{a' \in \mathcal{A}} q_t(x', a') - q_t(x, a) \right),$$



Convergence theorem

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion

Theorem

Assume

1. State-uniformity.
2. Any state is reached eventually under any policy.
3. Learning rate satisfies

$$\sum_{t=0}^{\infty} \alpha_t(x, a) = \infty, \quad \sum_{t=0}^{\infty} \alpha_t(x, a)^2 < \infty$$

Then starting with any $q_0 : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ the asynchronous Q-learning algorithm yields a sequence $(q_t)_{t \in \mathbb{N}}$ such that

$$q_t \rightarrow Q^*$$



HDP classes

Theoretical aspects of Q-learning

Introduction

Model- dependent algorithms

Model-free algorithms

Finite case

History dependent setting

Linear function approximation

Deep fitted Q-iteration

Conclusion

It seems relevant to ask how restrictive the state-uniformity assumption is. [Majeed, Hutter (2018)] answers this by an array of examples showing the following relations of the classes of decision processes:

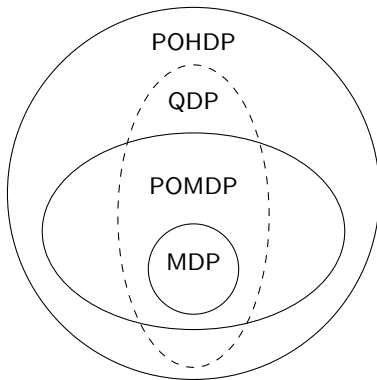


Figure: Classes of finite decision processes considered in [Majeed, Hutter (2018)].



Setting

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion



Linear function classes

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

**Linear function
approximation**

Deep fitted
Q-iteration

Conclusion



TD gradient steps

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion



The algorithm

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion



The theorem

Theoretical aspects of Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

**Linear function
approximation**

Deep fitted
Q-iteration

Conclusion



Conclusion

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion



Introduction

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion



Fitted Q-iteration

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Finite case

History dependent
setting

Linear function
approximation

Deep fitted
Q-iteration

Conclusion



Conclusion

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Conclusion

Errata



Error corrections

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Conclusion

Errata

- p. 31 Continuous MDP should also be greedy (i.e. condition on P is missing).
- p. 33 Top: Discussion of proofs of theorem 2.77: "...to obtain a contradiction to the statement that $\bigcup_{N \in \mathbb{N}} \mathcal{DN}(\sigma, (w, N, 1))$ is **not** dense in $C([0, 1]^w)$ ".
- p. 33 Middle: the bound on $\left| Q^* - \tilde{Q}_k \right|$ should be $\leq \gamma^k V_{\max} + \varepsilon/(1 - \gamma)$ (and not $< \varepsilon/(1 - \gamma)$).
- p. 36 Corollary 2.84: $\left\| Q^* - \tilde{Q}_k \right\|_{\infty} \leq \gamma^k V_{\max} + \dots$
- p. 46 Top: "we may view FQI as a class of **algorithms**, because ...".



Algorithm reference error

Theoretical
aspects of
Q-learning

Introduction

Model-
dependent
algorithms

Model-free
algorithms

Conclusion

Errata

Due to an unexpected behavior of the \LaTeX -package *cleveref* when adding line numbers to algorithms, many references to 'algorithm n ' was accidentally changed to 'line k '.

p. 19 Replace 'line 3' with 'algorithm 1'.

p. 24 Replace 'line 3' with 'algorithm 2'.

p. 25 Top: "and therefore **algorithm 2** could be applied".

Middle: "To use **algorithm 1** and **algorithm 2**".

Bottom: "update step $\tilde{V}_{k+1} \leftarrow T_\tau \tilde{V}_k$ in **algorithm 2** becomes".

p. 26 "Similarly in **algorithm 2** the update step ...".

p. 29 "we thus have convergence of **algorithm 3**".

p. 30 Top: "which are used in **algorithm 3** are not [...] This means that if we were to use **algorithm 3** naively ...".

p. 37 "It is clear that in the model-free setting **algorithm 3** will not work ...".