

# A Theoretical Analysis of Fitted Q-Iteration

Jacob Harder

University of Copenhagen

April 16, 2020

## 0.1 Abstract

## 0.2 Foreword

The main purpose of this master thesis for me, has been to investigate what has been proven about the convergence of Q-learning algorithms. In particular Q-learning algorithms using (deep) ANNs.

## 0.3 Disambiguation

- $\text{id} := x \mapsto x$  the identity function.
- $[\phi] := 1$  when  $\phi$  is true/holds and 0 otherwise, for a logical formula  $\phi$ .
- $[q] := \{1, \dots, q\}$  for  $q \in \mathbb{N}$ .
- $C_{\mathbb{K}}(X) := \{f : X \rightarrow \mathbb{K} \mid f \text{ continuous}\}$ ,  $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$ .  $C(X) = C_{\mathbb{R}}(X)$
- $C_b(X) := \{f \in C(X) \mid \exists K : \|f\|_{\infty} < K\}$ . I.e. the space of bounded continuous functions.
- ANN: abrv. artificial neural network see definition 12.
- $\delta_a$ : Dirac-measure of point  $a$ . I.e.  $\delta_a(A) = [a \in A]$ .
- $(\Omega, \mathcal{F}, \mathbb{P})$ : the underlying measure space of all random variables and processes when not otherwise specified.
- $\mathbb{B}_n$  the  $n$ -dimensional Borel  $\sigma$ -algebra.
- $\lambda^n$  the  $n$ -dimension Lebesgue measure.

## 1 Introduction

### 1.1 Reinforcement learning in general

#### 1.1.1 Reinforcement learning

In Reinforcement Learning (RL) we are concerned with finding an optimal policy for an agent in some environment. This environment is described by a so-called decision process consisting of a sequence of state and action spaces  $\mathcal{S}_1, \mathcal{A}_1, \mathcal{S}_2, \dots$  and rules  $P_1, R_1, P_2, \dots$  specifying which states and rewards and likely to follow after some action is chosen. In the general case  $\mathcal{S}_i, \mathcal{A}_i$  are measurable spaces and  $P_i, R_i$  are probability kernels on  $\mathcal{S}_{i+1}$  and  $\overline{\mathbb{R}}$ , respectively. One then

attempts to find a *policy* (behavior or strategy) that maximizes the rewards returned from the environment. A policy and an distribution over starting states  $\mu \in \mathcal{P}(\mathcal{S}_1)$  give rise to a countable stochastic process,  $(X_i)_{i \in \mathbb{N}} = (S_i, A_i, R_i)_{i \in \mathbb{N}}$  that is a probability measure  $P_\mu^\pi$  on  $\mathcal{S}_1 \times \mathcal{A}_1 \times \dots$ . Intuitively  $S_1$  is drawn from  $\mu$ , then for all  $i \in \mathbb{N}$   $A_i$  is drawn from  $\pi(\cdot | S_i)$ , a reward is then drawn from  $R(\cdot | S_i, A_i)$ , then  $S_{i+1}$  is drawn from  $P(\cdot | S_i, A_i)$  and so on.

In RL there is a categorization of algorithms into *off-policy* and *on-policy* classes. This is simply whether the algorithm learns from data (states, actions and rewards) arising from following its own policy (on-policy) or it can learn from more arbitrary data (off-policy). This *more arbitrary data* could for example be the trajectory of another algorithm when interacting with a decision process, or simply state-action-reward pairs drawn from some distribution. In this paper we exclusively consider off-policy algorithms.

## 1.2 Measure theory

We work with a background probability space  $(\Omega, \Sigma_\Omega, \mathbb{P})$ . For a measurable space  $(\mathcal{X}, \Sigma_\mathcal{X})$  we denote the set of probability measures on this space  $\mathcal{P}(\Sigma_\mathcal{X})$  or simply  $\mathcal{P}(\mathcal{X})$  when the  $\sigma$ -algebra is unambiguous. When taking cartesian products  $\mathcal{X} \times \mathcal{Y}$  of measurable spaces  $(\mathcal{X}, \Sigma_\mathcal{X}), (\mathcal{Y}, \Sigma_\mathcal{Y})$  we always endow such with the product  $\sigma$ -algebra  $\Sigma_\mathcal{X} \otimes \Sigma_\mathcal{Y}$ , unless otherwise specified. A map  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is called  $\Sigma_\mathcal{X}$ - $\Sigma_\mathcal{Y}$  measurable provided  $f^{-1}(\Sigma_\mathcal{Y}) \subseteq \Sigma_\mathcal{X}$  and we denote the set of such functions  $\mathcal{M}(\Sigma_\mathcal{X}, \Sigma_\mathcal{Y})$ . By a random variable  $X$  on  $(\mathcal{X}, \Sigma_\mathcal{X})$  mean a  $\Sigma_\Omega$ - $\Sigma_\mathcal{X}$  measurable map.

### 1.2.1 Kernels

**Definition 1** (Probability kernel). Let  $(\mathcal{X}, \Sigma_\mathcal{X}), (\mathcal{Y}, \Sigma_\mathcal{Y})$  be measurable spaces. A function

$$\kappa(\cdot | \cdot) : \Sigma_\mathcal{Y} \times \mathcal{X} \rightarrow [0, 1]$$

is a  $(\mathcal{X}, \Sigma_\mathcal{X})$ -**probability kernel** on  $(\mathcal{Y}, \Sigma_\mathcal{Y})$  provided

1.  $B \mapsto \kappa(B | x) \in \mathcal{P}(\Sigma_\mathcal{Y})$  that is  $\kappa(\cdot | x)$  is a probability measure for any  $x \in \mathcal{X}$ .
2.  $x \mapsto \kappa(B | x) \in \mathcal{M}(\Sigma_\mathcal{X}, \Sigma_\mathcal{Y})$  that is  $\kappa(B | \cdot)$  is  $(\Sigma_\mathcal{X}$ - $\Sigma_\mathcal{Y})$  measurable for any  $B \in \Sigma_\mathcal{Y}$ .

When the  $\sigma$ -algebras are unambiguous we shall simply say an  $\mathcal{X} \rightsquigarrow \mathcal{Y}$  kernel. For any  $x \in \mathcal{X}$  and  $f \in \mathcal{L}_1(\kappa(\cdot | x))$  we write the integral of  $f$  over  $\kappa(\cdot | x)$  as  $\int f(y) d\kappa(y | x)$ .

We now state some fundamental results on probability kernels

**Theorem 1** (Integration of a kernel). Let  $\mu \in \mathcal{P}(\mathcal{X})$  and  $\kappa : \mathcal{X} \rightsquigarrow \mathcal{Y}$ . Then there exists a uniquely determined probability measure  $\lambda \in \mathcal{P}(\Sigma_\mathcal{X} \otimes \Sigma_\mathcal{Y})$  such that

$$\lambda(A \times B) = \int_A \kappa(B, x) d\mu(x)$$

We denote this measure  $\lambda = \kappa\mu$ .

*Proof.* We refer to [ref to EH markov, thm. 1.2.1]. □

Notice that by theorem 1 besides getting a probability measure on  $\mathcal{X} \times \mathcal{Y}$  we get an induced probability measure on  $\mathcal{Y}$  defined by  $B \mapsto (\kappa\mu)(\mathcal{X} \times B)$ . We will denote this measure by  $\kappa \circ \mu$ . This way  $\kappa$  can also be seen as a mapping from  $\mathcal{P}(\mathcal{X}) \rightarrow \mathcal{P}(\mathcal{Y})$ . Also note that  $\kappa \circ \delta_x = \kappa(\cdot | x)$ .

For an idea how to actually compute integrals over kernel derived measures we here include

**Theorem 2** (Extended Tonelli and Fubini). Let  $\mu \in \mathcal{P}(\mathcal{X})$ ,  $f \in \mathcal{M}(\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{Y}}, \mathbb{B})$  be a measurable function and  $\kappa : \mathcal{X} \rightsquigarrow \mathcal{Y}$  be a probability kernel. Then

$$\int |f| d\kappa \circ \mu = \int \int |f| d\kappa(\cdot | x) d\mu(x)$$

Furthermore if this is finite, i.e.  $f \in \mathcal{L}_1(\kappa(\cdot, \mu))$  then  $A_0 := \{x \in \mathcal{X} \mid \int f d\kappa(\cdot | x) < \infty\} \in \Sigma_{\mathcal{X}}$  with  $\mu(A_0) = 1$ ,

$$x \mapsto \begin{cases} \int f d\kappa(\cdot | x) & x \in A_0 \\ 0 & x \notin A_0 \end{cases}$$

is  $\Sigma_{\mathcal{X}}\text{-}\mathbb{B}$  measurable and

$$\int f d\kappa(\cdot | \mu) = \int_{A_0} \int f d\kappa(\cdot | x) d\mu(x)$$

*Proof.* We refer to [ref to EH markov, thm. 1.3.2 + 1.3.3] □

**Proposition 1** (Composition of kernels). Let  $\kappa : \mathcal{X} \rightsquigarrow \mathcal{Y}, \psi : \mathcal{Y} \rightsquigarrow \mathcal{Z}$  be probability kernels. Then

$$(\psi \circ \kappa)(A | x) := \int \psi(A | y) d\kappa(y | x), \quad \forall A \in \Sigma_{\mathcal{Z}}, x \in \mathcal{X}$$

is a  $\mathcal{X} \rightsquigarrow \mathcal{Z}$  probability kernel called the composition of  $\kappa$  and  $\psi$ . The composition operator  $\circ$  is associative, i.e. if  $\phi : \mathcal{Z} \rightsquigarrow \mathcal{W}$  is a third probability kernel then  $(\phi \circ \psi) \circ \kappa = \phi \circ (\psi \circ \kappa)$ . The associativity also extends to measures, i.e.  $\forall \mu \in \mathcal{X} : (\psi \circ \kappa) \circ \mu = \psi \circ (\kappa \circ \mu)$  and this is uniquely determined by  $\psi, \kappa$  and  $\mu$ .

*Proof.* The first assertion is a trivial verification of the two conditions in definition 1 and left as an exercise. For the associativity we refer to [todo ref to EH markov, lem. 4.5.4]. □

Proposition 1 actually makes the class of measurable spaces into a category [todo ref: see Lawvere, The Category of Probabilistic Mappings], with identity  $\text{id}_{\mathcal{X}}(\cdot | x) = \delta_x$ . Notice that the mapping  $(A, x) \mapsto \delta_x(A) \kappa(A | x)$  defines a probability kernel  $\mathcal{X} \rightsquigarrow \mathcal{X} \times \mathcal{Y}$  which we could denote  $\text{id}_{\mathcal{X}} \times \kappa$ . Now if  $\psi : \mathcal{X} \times \mathcal{Y} \rightsquigarrow \mathcal{Z}$  is a kernel then by proposition 1 the composition  $(\text{id}_{\mathcal{X} \times \mathcal{Y}} \times \psi) \circ (\text{id}_{\mathcal{X}} \times \kappa)$  is a kernel  $\mathcal{X} \rightarrow \mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$  which we will denote  $\psi \kappa$ . It inherits associativity from  $\circ$  and again this associativity extends to application on measures: if  $\mu$  is a measure on  $\mathcal{X}$  then  $\psi(\kappa \mu) = (\psi \kappa) \mu$ .

**Proposition 2.** Let  $\kappa : \mathcal{X} \rightarrow \mathcal{Y}$  be a probability kernel and  $f : \mathcal{Y} \rightarrow \overline{\mathbb{R}}$  be integrabel. Then  $x \mapsto \int f d\kappa(\cdot | x)$  is measurable into  $(\overline{\mathbb{R}}, \mathbb{B})$ .

*Proof.* Simple functions are measurable since  $\kappa$  is a kernel. Now extend by sums and limits. □

### 1.2.2 Kernel derived processes

Let  $(\mathcal{X}_n, \Sigma_{\mathcal{X}_n})_{n \in \mathbb{N}}$  be a sequence of measurable spaces. For each  $n \in \mathbb{N}$  define  $\mathcal{X}^n := \mathcal{X}_1 \times \cdots \times \mathcal{X}_n$ ,  $\Sigma_{\mathcal{X}^n} := \Sigma_{\mathcal{X}_1} \otimes \cdots \otimes \Sigma_{\mathcal{X}_n}$  and let  $\kappa_n : \mathcal{X}^n \rightsquigarrow \mathcal{X}_{n+1}$  be a probability kernel. Then  $\kappa^n := \kappa_n \dots \kappa_1$  is a kernel from  $\mathcal{X}_1$  to  $\mathcal{X}^n$ . So for any probability measure  $\rho_1 \in \mathcal{P}(\mathcal{X}_1)$  there exists a unique probability measure  $\rho_n$  on  $\mathcal{X}^n$  defined by  $\kappa^n \rho_1$ .

Let  $\mathcal{X}^{\infty} := \prod_{n \in \mathbb{N}} \mathcal{X}_n$  and  $\Sigma_{\mathcal{X}^{\infty}} := \bigotimes_{n \in \mathbb{N}} \Sigma_{\mathcal{X}_n}$ . We are not equipped to establish existence of a kernel generated measure on  $(\mathcal{X}^{\infty}, \Sigma_{\mathcal{X}^{\infty}})$  yet which we will need. This problem was solved by Cassius Ionescu-Tulcea in 1949:

**Theorem 3** (Ionescu-Tulcea extension theorem). For every  $\mu \in \mathcal{P}(\mathcal{X}_1)$  there exists a unique probability measure  $\rho \in \mathcal{P}(\mathcal{X}^\infty)$  such that

$$\rho_n(A) = \rho \left( A \times \prod_{k=n+1}^{\infty} \mathcal{X}_k \right), \quad \forall A \in \Sigma_{\mathcal{X}^n}, n \in \mathbb{N}$$

We denote this measure  $\dots \kappa_2 \kappa_1 \mu = \prod_{i=1}^{\infty} \kappa_i \mu := \rho$ .

*Proof.* Todo: what about this. □

**Proposition 3** (Ionescu-Tulcea kernel). Let  $\mu_x$  denote the Ionescu-Tulcea measure of a sequence of probability kernels  $\kappa_i : \mathcal{X}^i \rightarrow \mathcal{X}_{i+1}$  with starting measure  $\delta_x$  on  $\mathcal{X}_1$  for any  $x \in \mathcal{X}_1$ . Then  $\kappa(A \mid x) = \mu_x(A)$  defines a probability kernel  $\kappa : \mathcal{X}_1 \rightarrow \mathcal{X}^\infty$ .

*Proof.* Since we already know that  $\mu_x$  is a probability measure for any  $x \in \mathcal{X}_1$ , we just have to show that  $\kappa(A \mid x) = \mu_x(A)$  is measurable for all  $A \in \bigotimes_i \Sigma_{\mathcal{X}_i}$ . ... todo □

**Lemma 1.** The Ionescu-Tulcea kernel satisfies  $\prod_{i=1}^{\infty} \kappa_i = (\prod_{i=2}^{\infty} \kappa_i) \kappa_1$ .

*Proof.* Let  $x \in \mathcal{X}_1$ . Notice that by associativity of the finitely induced measures  $\kappa_n \dots \kappa_1 \delta_x = (\kappa_n \dots \kappa_2)(\kappa_1 \delta_x)$ . This implies that

$$\prod_{i=1}^{\infty} \kappa_i \delta_x \left( A \times \prod_{k=n+1}^{\infty} \mathcal{X}_k \right) = \prod_{i=2}^{\infty} \kappa_i \kappa_1 \delta_x \left( A \times \prod_{k=n+1}^{\infty} \mathcal{X}_k \right)$$

for all  $n \in \mathbb{N}$  and  $A \in \Sigma_{\mathcal{X}^n}$ . By the uniqueness in theorem 3 we are done. □

**Proposition 4.** Let  $\mathcal{X}, \mathcal{Y}$  be separable and metrizable,  $\kappa : \mathcal{X} \rightarrow \mathcal{Y}$  be a continuous probability kernel and  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \overline{\mathbb{R}}$  be Borel-measurable satisfying one of  $f \leq 0, f \geq 0, |f| < \infty$ . If  $f$  is bounded from above (below) and upper (lower) semicontinuous then

$$x \mapsto \int f d\kappa(\cdot \mid x)$$

is bounded from above (below) and upper (lower) semicontinuous.

*Proof.* We refer to [BS SOC, prop. 7.31]. □

## 2 Decision models and value functions

### 2.1 General decision model

In the quest to have a united framework to talk about results from several different sources we define here a quite general model. One which is quite close to in generality can be found in [ref. to Schal]. In this section recall that  $\underline{\mathbb{R}} = \mathbb{R} \cup \{-\infty\}$ ,  $\overline{\mathbb{R}} = \mathbb{R} \cup \{\infty\}$  and  $\overline{\underline{\mathbb{R}}} = \mathbb{R} \cup \{\pm\infty\}$ .

**Definition 2** (Decision model). A general **decision** model is determined by

1.  $(\mathcal{S}_n, \Sigma_{\mathcal{S}_n})_{n \in \mathbb{N}}$  a measurable space of **states** for each timestep.
2.  $(\mathcal{A}_n, \Sigma_{\mathcal{A}_n})_{n \in \mathbb{N}}$  a measurable space of **actions** for each timestep.

for each  $n \in \mathbb{N}$  we define the so called **history** spaces

$$\mathcal{H}_n = \mathcal{S}_1 \times \mathcal{A}_1 \times \mathcal{S}_2 \times \overline{\underline{\mathbb{R}}} \times \mathcal{A}_2 \times \mathcal{S}_3 \times \overline{\underline{\mathbb{R}}} \dots \times \mathcal{S}_n, \mathcal{H}_\infty = \mathcal{S}_1 \times \mathcal{A}_1 \times \mathcal{S}_2 \times \overline{\underline{\mathbb{R}}} \times \dots$$

with associated product  $\sigma$ -algebras

3.  $(P_n)_{n \in \mathbb{N}}$  a sequence of  $\mathcal{H}_n \times \mathcal{A}_n \rightsquigarrow \mathcal{S}_{n+1}$  kernels called the **transition** kernels.
4.  $(R_n)_{n \in \mathbb{N}}$  a sequence of  $\mathcal{H}_{n+1} \rightsquigarrow \overline{\mathbb{R}}$  kernels called the **reward** kernels.

Some authors refers to this as a *dynamic programming* model. Notice the slight irregularity in the beginning of the history spaces: We are missing a reward state after  $\mathcal{S}_1$ . We could avoid this by introducing some start reward, but we will be careless.

The vast majority of sources considered in this paper actually specialize the DP model with the following:

**Assumption 1** (One state and action space).  $\mathcal{S}_1 = \mathcal{S}_2 = \dots := \mathcal{S}$   $\mathcal{A}_1 = \mathcal{A}_2 = \dots := \mathcal{A}$

However we will do without this for the rest of this section in order to present some results in the generality they deserve. One could ask if it is possible to embed the general DP model into one with assumption 1 by setting  $\mathcal{S} := \bigcup_{i \in \mathbb{N}} \mathcal{S}_i$  and  $\mathcal{A} := \bigcup_{i \in \mathbb{N}} \mathcal{A}_i$  or similar. One attempt at this can be found in [BS SOC, chp. 10], but this will not be covered here.

Other ways to specialize include reducing one or both of the transition and reward kernels to functions defined on  $\mathcal{S} \times \mathcal{A}$ . These models are often called *deterministic*, but the exact definitions vary across sources, and we will instead specify each setting individually.

For a DP model we can define

**Definition 3** (Policy). A (randomized) **policy**  $\pi = (\pi_n)_{n \in \mathbb{N}}$  is a sequence of  $\mathcal{H}_n \rightsquigarrow \mathcal{A}_n$  kernels. The set of all policies we denote  $R\Pi$ . The policy  $\pi$  is called **semi Markov** if each  $\pi_i$  only depends on the first and last state in the history and is called **Markov** if only the last. The sets are denoted  $s\Pi$  and  $\Pi$ . Furthermore  $\pi$  is called **deterministic** if all  $\pi_i$  are degenerate, i.e. are actually measurable functions from  $\mathcal{H}_n$  to  $\mathcal{A}_n$ . Under assumption 1 it makes sense to make a (Markov) policy  $(\pi, \pi, \dots)$ , where  $\pi$  only depends on the last state. Such a policy is called **stationary**, and the set of them denoted  $S\Pi$ .

We have the following inclusions

$$S\Pi \subseteq \Pi \subseteq s\Pi \subseteq R\Pi$$

$$DS\Pi \subseteq D\Pi \subseteq DSs\Pi \subseteq Ds\Pi$$

**Proposition 5.** A dynamic programming model together with a policy  $\pi$  defines a probability kernel  $\kappa_\pi : \mathcal{S}_1 \rightarrow \mathcal{H}_\infty$ .

*Proof.* This is the Ionescu-Tulcea kernel generated by  $\dots R_2 P_2 \pi_2 R_1 P_1 \pi_1$ . □

This kernel yields a probability measure  $\kappa_\pi \mu$  on  $\mathcal{H}_\infty$  for every  $\mu \in \mathcal{S}_1$ . In particular for any  $s \in \mathcal{S}_1$   $\kappa_\pi \delta_s$  yields the measure  $\kappa(\cdot \mid s)$  and we shall occasionally write this  $\kappa_\pi s$  and integration with respect to it  $\mathbb{E}_s^\pi$ .

Across litterature generally any function mapping a state space  $\mathcal{S}$  to  $\overline{\mathbb{R}}$  can be called a (state) **value** function. Similarly any  $\overline{\mathbb{R}}$  valued function on pairs of states and actions can be called (state) **action value** or **Q-** function. The idea behind such functions are commonly to estimate the cumulative rewards associated with a state or state-action pair and the trajectory of states it can lead to. In order to define some of the most standard of value functions, which we call **ideal** to avoid confusion, we will need one of the following conditions:

**Condition  $F^+$ .**  $R_i(\{\infty\} \mid h) = 0$  for all  $h \in \mathcal{H}_{i+1}$  and  $i \in \mathbb{N}$

**Condition  $F^-$ .**  $R_i(\{-\infty\} \mid h) = 0$  for all  $h \in \mathcal{H}_{i+1}$  and  $i \in \mathbb{N}$

When assuming either of  $(F^+)$  or  $(F^-)$  adding rewards cannot lead to a  $\infty - \infty$  situation, and the following definition makes sense

**Definition 4.** Let  $\mathcal{R}_i : \mathcal{H}_\infty \rightarrow \overline{\mathbb{R}}$  be the projection onto the  $i$ th reward. Define

$$V_{n,\pi}(s) = \mathbb{E}_s^\pi \sum_{i=1}^n \mathcal{R}_i$$

called the  $n$ th finite **ideal** value function. When  $n = 0 \forall \pi : V_{0,\pi} = V_0 := 0$ .

These are also called *finite horizon* value functions.

We would like to extend this to an infinite horizon value function, i.e. letting  $n$  tend to  $\infty$ . To ensure that the integral is well-defined we need one of the following conditions

**Condition P.**  $R_i([0, \infty] \mid h) = 1, \forall h \in \mathcal{H}_{i+1}, i \in \mathbb{N}$

**Condition N.**  $R_i([-\infty, 0] \mid h) = 1 \forall h \in \mathcal{H}_{i+1}, i \in \mathbb{N}$

**Condition D.** There exist a bound  $R_{\max} > 0$  and a  $\gamma \in [0, 1)$  called the **discount** factor such that  $R_i([-R_{\max}\gamma^i, R_{\max}\gamma^i] \mid h) = 1 \forall h \in \mathcal{H}_{i+1}, i \in \mathbb{N}$

**Definition 5.** Assume (P), (N) or (D). We define the (infinite horizon) ideal value function by

$$V_\pi(s) = \mathbb{E}_s^\pi \lim_{n \rightarrow \infty} \sum_{i=1}^n \mathcal{R}_i$$

**Proposition 6.** The ideal value functions  $V_{n,\pi}, V_\pi$  are measurable into  $(\overline{\mathbb{R}}, \overline{\mathbb{B}})$ .

*Proof.* Use proposition 2. □

**Proposition 7.** Under P, N or D we have  $\lim_{n \rightarrow \infty} V_{n,\pi} = V_\pi$  for all  $\pi \in RII$ .

*Proof.* By monotone or dominated convergence. □

### 2.1.1 Optimal policies

Let  $(\mathcal{S}_n, \mathcal{A}_n, P_n, R_n)_{n \in \mathbb{N}}$  be a DP model.

**Assumption 2.** (Reward independence)  $P_n, R_n$  and policies are only allowed to depend on the states and actions.

In all sources known to this writer assumption 2 is assumed. This is a bit of a puzzle since it is obvious that one could want to define algorithms (policies) that take into account which rewards they received in the past. We will also do this but stick to the standard and never attempt to evaluate ideal value functions of policies that depend on rewards. Thus we will assume assumption 2 henceforth with including the shrinkage of the set of general policies  $RII$  that it entails.

A neat consequence of assumption 2 when talking about value functions is that we can reduce the reward kernels to functions  $r_i : \mathcal{H}_{i+1} \rightarrow \overline{\mathbb{R}} = h \rightarrow \int r dR_i(r \mid h)$  which are measurable (due to proposition 2).

**Definition 6** (Optimal value functions).

$$V_n^*(s) := \sup_{\pi \in RII} V_n^\pi(s) \qquad V^*(s) := \sup_{\pi \in RII} V^\pi(s)$$

are called the **optimal** value functions. A policy  $\pi^* \in RII$  for which  $V_{\pi^*} = V^*$  is called an **optimal** policy. If  $V_{n,\pi^*} = V_n^*$  it is called  $n$ -optimal.

An interesting fact about the optimal value functions is that they might not be Borel measurable [todo ref to counterexample] even in the finite case. After all we are taking a supremum over sets of policies which have cardinality of at least the continuum. However it is sometimes possible to show that they are universally measurable, thus Lebesgue measurable and therefore standard Lebesgue integration is possible. We will take these discussions as they occur in various settings.

At this point some central questions can be asked.

1. To which extend does an optimal policy  $\pi^*$  exist?
2. Does  $V_n^*$  converge to  $V^*$ ?
3. When can optimal policies be chosen to be Markov, deterministic, etc.?
4. Can an algorithm be designed to efficiently find  $V^*$  and  $\pi^*$ ?

These questions has been answered in a variety of settings. We will try to address these question in order by strength of assumptions they require.

In a quite general setting, questions 1 and 2 was investigated by M. Schäl in 1974 [todo ref. to On Dynamic Programming: Compactness of the space of policies, 1974]. Here some additional structure on our model is imposed:

**Setting 1** (Schäl). 1.  $V_\pi < \infty$  for all policies  $\pi \in R\Pi$ .

2.  $(\mathcal{S}_n, \Sigma_{\mathcal{S}_n})$  is assumed to be standard Borel. I.e.  $\mathcal{S}_n$  is a non-empty Borel subset of a Polish space and  $\Sigma_{\mathcal{S}_n}$  is the Borel subsets of  $\mathcal{S}_n$ .
3.  $(\mathcal{A}_n, \Sigma_{\mathcal{A}_n})$  is similarly assumed to be standard Borel.
4.  $\mathcal{A}_n$  is compact.
5.  $\forall s \in \mathcal{S}_1 : Z_n = \sup_{N \geq n} \sup_{\pi \in R\Pi} \sum_{t=n+1}^N \mathbb{E}_s^\pi r_t \rightarrow 0$  as  $n \rightarrow \infty$ .

In this setting Schäl introduced two set of criteria for the existence of an optimal policy:

**Condition S.** 1. The function

$$(a_1, a_2, \dots, a_n) \mapsto P_n(\cdot \mid s_1, a_1, s_2, a_2, \dots, s_n, a_n)$$

is set-wise continuous (hence the name **S**) for all  $s_1, \dots, s_n \in \mathcal{S}^n$ .

2.  $r_n$  is upper semi-continuous.

**Condition W.** 1. The function

$$(h_n, a_n) \mapsto P_n(\cdot \mid h_n, a_n)$$

is weakly continuous (hence the name **W**).

2.  $r_n$  is continuous.

**Theorem 4** (Existence and convergence of optimal policies in DP). When either S or W hold then

1. There exist an optimal policy  $\pi^* \in R\Pi$ .
2.  $V_n^* \rightarrow V^*$  as  $n \rightarrow \infty$ .

*Proof.* We refer to [todo ref: On Dynamic Programming: Compactness of the space of policies, M. Schäl 1974]. □

## 2.2 Markov decision model

**Setting 2** (Markov decision model). • A decision model (definition 2) under Assumption 1 i.e. there is only one state and action space  $\mathcal{S}, \mathcal{A}$ .

- $P_n$  depends only on  $s_n$  and  $a_n$  and does not differ with  $n$ . I.e. there exists a kernel  $P$  such that  $P_n(\cdot \mid s_1, \dots, s_n, a_n) = P(\cdot \mid s_n, a_n)$  for all  $n \in \mathbb{N}$ . We will write  $P$  instead of  $P_n$  understanding kernel compositions as if using  $P_n$ .
- $r_n$  depends only on  $s_n$  and  $a_n$  and does not differ with  $n$  except for a potential discount. I.e. there exists a measurable function  $r : \mathcal{S} \times \mathcal{A} \rightarrow \overline{\mathbb{R}}$  such that  $r = r_n/\gamma^{n-1}$  for all  $n \in \mathbb{N}$  (in the case where we are not discounting set  $\gamma = 1$ ).

## 2.3 Bertsekas-Shreve framework

The theory described here is largely based on [ref to Bertsekas-Shreve, Stochastic Optimal Control]. Their framework is cost-based as opposed to the this paper reward-based outset. This means that positive and negative, upper and lower, supremum and infimum, ect. are opposite to the source.

**Setting 3** (BS). • We consider an MDM  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$  (see setting 2).

- $\mathcal{S}$  and  $\mathcal{A}$  are Borel spaces.
- $\mathcal{A}$  is compact.
- $P(S \mid \cdot)$  is continuous for any  $S \in \Sigma_{\mathcal{S}}$ .
- $r(s, a) = \gamma^{1-i} \int x dR(x \mid s, a)$  is upper semicontinuous and uniformly bounded from above (least upper bound denoted  $0 < R_{\max} < \infty$ ).
- The policies must consist of universally measurable probability kernels.

The original setup in [ref to Bertsekas-Shreve, Stochastic Optimal Control] is slightly different than the setup here presented. Besides having a state and action space, it also features a non-empty Borel space called the *disturbance space*  $W$ , a *disturbance kernel*  $p : \mathcal{S} \times \mathcal{A} \rightarrow W$ , instead of a transition kernel which on the other hand is a deterministic *system function*  $f : \mathcal{S} \times \mathcal{A} \times W \rightarrow \mathcal{S}$  which should be Borel measurable. Moreover it allows for constrains on the action space for each state. This is made precise by a function  $U : \mathcal{S} \rightarrow \Sigma_{\mathcal{A}}$  and a restriction on RII that all policies  $\pi$  should satisfy  $\pi(U(s) \mid s) = 1$ . Lastly the rewards are interpreted as negative costs, and thus  $g$  is required to be semi *lower*continuous.

By setting  $P(\cdot \mid s, a) = f(s, a, p(\cdot \mid s, a))$  and maximizing rewards of upper semicontinuous instead of minimizing lower semicontinuous ones, we fully capture all aspects of the original model and its results, except the for the action constrains.

Notice that setting 3 implies  $(F^+)$ . Throughout this section (P), (N) or (D) are always assumed. Some results only hold for some of these conditions and we will indicate this by e.g. (D) (P) when the result only holds for the discounted and positive case. At this point it makes sense to define

**Definition 7** (The  $T$ -operators). For a stationary policy  $\pi$  and measurable  $V : \mathcal{S} \rightarrow \overline{\mathbb{R}}$  with  $V \geq 0$ ,  $V \leq 0$  or  $|V| < \infty$  we define the operators

$$P_{\pi}V := s \mapsto \int V(s') dP\pi(s' \mid s)$$



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

$$TV := s \mapsto \sup_{a \in \mathcal{A}} T_a V(s)$$

where  $T_a = T_{\delta_a}$ .

**Proposition 8.** The operators  $P_\pi, T_\pi$  and  $T$  commutes with limits.

*Proof.* By monotone or dominated convergence theorems.  $\square$

**Proposition 9.** Let  $\pi = (\pi_1, \pi_2, \dots)$  be a Markov policy. Then  $V_{k,\pi} = T_{\pi_1} \dots T_{\pi_k} V_0$  and  $V_\pi = \lim_{k \rightarrow \infty} T_{\pi_1} \dots T_{\pi_k} V_0$ .

*Proof.* todo  $\square$

**Proposition 10.** Let  $\pi$  be a stationary policy then  $T_\pi V_\pi = V_\pi$ .

*Proof.* By proposition 9  $T_\pi V_\pi = T_\pi \lim_{k \rightarrow \infty} T_\pi^k V_0 = \lim_{k \rightarrow \infty} T_\pi^{k+1} V_0 = V_\pi$ .  $\square$

**Proposition 11.** Under (D) for any  $\pi \in R\Pi$  we have

$$|V_{n,\pi}|, |V_\pi|, |V_k^*|, |V^*| \leq V_{\max} := R_{\max}/(1 - \gamma)$$

*Proof.* For any  $\pi \in R\Pi$

$$|V_\pi(s)| \leq \mathbb{E}_s^\pi \sum_{i \in \mathbb{N}} |r_i| \leq \sum_{i \in \mathbb{N}} \gamma^{i-1} R_{\max} = R_{\max}/(1 - \gamma)$$

This also covers  $V_{n,\pi}$ .  $\square$

**Proposition 12.** (D)

$T$  and  $T_\pi$  are  $\gamma$ -contractive on  $\mathcal{L}_\infty(\mathcal{S})$ .

*Proof.* Let  $V, V' \in \mathcal{L}_\infty(\mathcal{S})$  and let  $K = \|V - V'\|_\infty$ . Then

$$|T^\pi V - T^\pi V'| = \gamma \left| \int V(s') - V'(s') dP\pi(s' | s) \right| \leq \gamma K$$

For  $T$  use that same argument and the fact that  $\left| \sup_x f(x) - \sup_y g(y) \right| \leq |\sup_x f(x) - g(x)|$  for any  $f, g : X \rightarrow \mathbb{R}$ .  $\square$

**Corollary 1.** (D)

Let  $\pi \in S\Pi$  be a stationary policy.  $V_\pi$  is the unique bounded fixed point of  $T_\pi$  in  $\mathcal{L}_\infty(\mathcal{S})$ .

*Proof.* By proposition 10  $V_\pi$  is a fixed point. By proposition 12 and Banach fixed point theorem we get uniqueness.  $\square$

**Proposition 13** (Prop. 8.6 in BS).  $V_k^* = T^k V_0$  and is upper semicontinuous. Furthermore for any  $k \in \mathbb{N}$  there exists a deterministic, Markov, Borel-measurable  $k$ -optimal policy  $\pi_k^* = (\pi_{k,1}^*, \pi_{k,2}^*, \dots, \pi_{k,k}^*, \dots) \in D\Pi$ . These policies satisfy for any  $i < k$   $\pi_i^* = (\pi_{k,k-i}^*, \dots, \pi_{k,k}^*, \dots)$ .

**Theorem 5** (Cor. 9.17.2 in BS). Under (N) or (D)  $V^* = \lim_{k \rightarrow \infty} V_k^*$  and is upper semicontinuous. Furthermore there exist a deterministic stationary, Borel-measurable policy  $\pi^*$ .

**Proposition 14.**  $V^* = T_{\pi^*} V^* = TV^*$

(D)  $V^*$  is the unique fixed point of  $T$  in  $\mathcal{L}_\infty(\mathcal{S})$ .

*Proof.* Since  $\pi^*$  is optimal  $V^* = V_{\pi^*} = T_{\pi^*} V_{\pi^*}$  by proposition 10. By theorem 5  $TV^* = T \lim_{k \rightarrow \infty} T^k V_0 = \lim_{k \rightarrow \infty} T^{k+1} V_0 = V^*$ . If (D) holds  $V^* \in \mathcal{L}_\infty(\mathcal{S})$  so proposition 12 and Banach fixed point theorem ensures uniqueness.  $\square$

### 2.3.1 Analytic setting

**Setting 4** (BS Analytic). The same as setting 3 except:  $P$  is not necessarily continuous.  $r$  is upper semianalytic.  $\mathcal{A}$  is not necessarily compact, but there exists a  $k \in \mathbb{N}$  such that  $\forall \lambda \in \mathbb{R}, n \geq k, s \in \mathcal{S}$

$$A_n^\lambda(s) = \left\{ a \in \mathcal{A} \mid r(s, a) + \gamma \int V_n^* P(\cdot \mid s, a) \geq \lambda \right\}$$

is a compact subset of  $\mathcal{A}$ .

**Theorem 6** (Prop. 9.17 BS). Under setting 4 we have  $V^* = \lim_{n \rightarrow \infty} V_n^*$  for all  $s \in \mathcal{S}$  and there exists a optimal policy  $\pi^*$  which is stationary and deterministic.

*Proof.* We refer to [todo ref to Bertsekas and Schreve, Stochastic Optimal Control: The Discrete-Time Case, prop. 9.17].  $\square$

## 2.4 Q-functions

The letter Q originates to a PhD thesis by C. Watkins from 1989 [todo ref C. Watkins, 1989]. Upon his definition he noted

“This is much simpler to calculate than  $[V_\pi]$  for to calculate  $[Q_\pi]$  it is only necessary to look one step ahead [...]

A clear advantage of working with Q-function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \overline{\mathbb{R}}$  rather than a value function  $V : \mathcal{S} \rightarrow \overline{\mathbb{R}}$ , is that finding the optimal action in state  $s$  requires only a maximization over the Q-function itself:  $a = \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a)$ . This should be seen in contrast to finding a best action according to  $V$ :  $a = \operatorname{argmax}_{a \in \mathcal{A}} r(s, a) + \gamma \mathbb{E}_{P(\cdot \mid s, a)} V$ . This requires taking an expectation with respect to the transition kernel  $P$ . Later we will study settings where we are not allowed to know the transition kernel when attempting to find the optimal strategy. In these situations the advantage of Q-functions is clear. For now however the transition kernel will remain known and we will in this section see how the results of state-value functions translate to Q-functions. The results in this section is original in the generality here presented, as I was unable to find them elsewhere.

**Definition 8.** Let  $\pi \in R\Pi$ . Define

$$Q_{k,\pi}(s, a) = r(s, a) + \gamma \mathbb{E}_{P(\cdot \mid s, a)} V_{k,\pi}, \quad Q_\pi = r(s, a) + \gamma \mathbb{E}_{P(\cdot \mid s, a)} V_\pi$$

We define  $Q_{0,\pi} = r + \gamma \mathbb{E} V_0 = r := Q_0$ .

**Proposition 15.**  $\lim_{k \rightarrow \infty} Q_{k,\pi} = Q_\pi$

*Proof.* (D) By dominated convergence or (P/N) by monotone convergence theorem.  $\square$

**Definition 9.**

$$Q_k^* = \sup_{\pi \in R\Pi} Q_{k,\pi}, \quad Q^* = \sup_{\pi \in R\Pi} Q_\pi$$

**Definition 10** ( $T$  operators for Q-functions). For any stationary policy  $\pi \in S\Pi$  and measurable  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \overline{\mathbb{R}}$  with  $Q \geq 0, Q \leq 0$  or  $|Q| < \infty$  we define

$$\begin{aligned} P_\pi Q(s, a) &= \int Q(s', a') d\pi P(s', a' \mid s, a) \\ T_\pi Q &= r + \gamma P_\pi Q \\ TQ(s, a) &= r(s, a) + \gamma \int \sup_{a' \in \mathcal{A}} Q(s', a') dP(\cdot \mid s, a) \end{aligned}$$

where  $T_a = T_{\delta_a}$ .

**Proposition 16.** Let  $\pi = (\pi_1, \dots, \pi_k, \dots) \in M\Pi$  then

- $T_{\pi_k} Q_{k-1, \pi} = r + \gamma \mathbb{E} T_{\pi_k} V_{k-1, \pi}$
- $Q_{k, \pi} = T_{\pi_1} \dots T_{\pi_k} r$ .

*Proof.* The first statement is almost by definition. For the second use the first inductively.  $\square$

**Proposition 17.**  $Q_k^* = r + \gamma \mathbb{E} V_k^*$  and  $Q^* = r + \gamma \mathbb{E} V^*$ .

*Proof.*  $Q^* = \sup_{\pi} Q_{\pi} = \sup_{\pi} (r + \gamma \mathbb{E} V_{\pi}) = r + \gamma \sup_{\pi} \mathbb{E} V_{\pi} = r + \gamma \mathbb{E} V^*$  where in the fourth equality we have used that  $V^* \geq V_{\pi}$  uniformly so  $\sup_{\pi} \mathbb{E} V_{\pi} \leq \mathbb{E} V^*$  while trivially  $\mathbb{E} V^* = \mathbb{E} V_{\pi^*} \leq \sup_{\pi} \mathbb{E} V_{\pi}$ . For  $Q_k^*$  the argument is similar.  $\square$

**Proposition 18.**  $\sup_{a \in \mathcal{A}} Q^*(s, a) = V^*(s)$

*Proof.* This is by definition after considering proposition 17.  $\square$

**Proposition 19.**  $TQ_k^* = r + \gamma \mathbb{E} T V_k^*$  and if  $\pi^* = (\pi_1^*, \dots, \pi_k^*, \dots)$  is  $k$ -optimal then  $Q_k^* = T_{\pi_1^*} \dots T_{\pi_k^*} r = T^k r$ .

*Proof.*

$$\begin{aligned} TQ_k^*(s, a) &= T(r + \gamma \mathbb{E} V_k^*)(s, a) \\ &= r(s, a) + \gamma \int \sup_{a' \in \mathcal{A}} (r(s', a') + \gamma \mathbb{E}_{P(\cdot | s', a')} V_k^*) dP(s' | s, a) \\ &= r(s, a) + \gamma \int \sup_{a' \in \mathcal{A}} \left( r(s', a') + \gamma \int V_k^*(s'') dP(s'' | s', a') \right) dP(s' | s, a) \\ &= r(s, a) + \gamma \int T V_k^*(s') dP(s' | s, a) \end{aligned}$$

To get  $Q_k^* = T^k r$  use this inductively  $Q_k^* = r + \gamma \mathbb{E} V_k^* = r + \gamma T V_{k-1}^* = T Q_{k-1}^* = \dots$ . For  $Q_k^* = T_{\pi_1^*} \dots T_{\pi_k^*} r$  use  $Q_k^* = r + \gamma \mathbb{E} V_k^* = r + \gamma T_{\pi_1^*} \dots T_{\pi_k^*} V_0$  and first statement in proposition 16 inductively.  $\square$

The proof of proposition 19 also shows

**Proposition 20.**  $TQ^* = r + \gamma \mathbb{E} T V^*$ .

implying

**Proposition 21.**  $TQ^* = Q^*$ .

**Proposition 22.** For stationary  $\pi \in S\Pi$  we have  $T_{\pi} Q_{\pi} = Q_{\pi}$ .

*Proof.* Using proposition 16 and proposition 10  $T_{\pi} Q_{\pi} = T_{\pi} (r + \gamma \mathbb{E} \lim_{k \rightarrow \infty} T_{\pi}^k V_0) = \lim_{k \rightarrow \infty} T_{\pi} (r + \gamma \mathbb{E} T_{\pi}^k V_0) = \lim_{k \rightarrow \infty} (r + \gamma \mathbb{E} T_{\pi}^{k+1} V_0) = r + \gamma \mathbb{E} \lim_{k \rightarrow \infty} T_{\pi}^{k+1} V_0 = r + \gamma \mathbb{E} V_{\pi} = Q_{\pi}$ .  $\square$

**Proposition 23.**  $Q^* = Q_{\pi^*}$  and  $T_{\pi^*} Q^* = Q^*$ .

*Proof.*  $Q^* = r + \gamma \mathbb{E} V^* = r + \gamma \mathbb{E} V_{\pi^*} = Q_{\pi^*}$  for the second statement use proposition 22.  $\square$

**Proposition 24.**  $Q^* = \lim_{k \rightarrow \infty} Q_k^*$ .

*Proof.* By monotone or dominated convergence and theorem 5.  $\square$

**Proposition 25.**  $T$  and  $T_\pi$  is  $\gamma$ -contractive on  $\mathcal{L}_\infty(\mathcal{S} \times \mathcal{A})$ . If furthermore (D) holds, then  $|Q_{k,\pi}|, |Q_\pi|, |Q_k^*|, |Q^*| \leq V_{\max}$  and  $Q_\pi, Q^*$  are the unique fixed points of  $T_\pi, T$  in  $\mathcal{L}_\infty(\mathcal{S} \times \mathcal{A})$ .

*Proof.* The contrativeness of  $T, T_\pi$  follows from the same argument as in proposition 12. If (D) holds the boundedness of the  $Q$  functions follow from an argument similar to the proof of proposition 11. Then proposition 21, proposition 22 and Banach fixed point theorem implies uniqueness.  $\square$

**Proposition 26.** (D)

For any  $Q \in \mathcal{L}_\infty(\mathcal{S} \times \mathcal{A})$   $T^k Q$  converges to  $Q^*$  with rate  $\gamma^k$ .

That is  $\|T^k Q - Q^*\|_\infty \leq \gamma^k \|Q - Q^*\|_\infty$ .

*Proof.* By proposition 25  $T$   $\gamma$ -contracts so

$$\begin{aligned} \|T^k Q - Q^*\|_\infty &= \|T^k Q - T^k Q^*\| \\ &\leq \gamma^k \|Q - Q^*\| \end{aligned}$$

$\square$

**Proposition 27.** (D)(N)

$Q^*(s, \cdot)$  is upper semicontinuous.

*Proof.* Since  $P$  is continuous and  $V^*$  is upper semicontinuous by theorem 5 the proposition follow by proposition 4.  $\square$

**Definition 11.** Let  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  be a stationary policy. If

$$\pi \left( \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a) \mid s \right) = 1$$

then  $\pi$  is said to be **greedy** with respect to  $Q$  and is denoted  $\pi_Q$ .

**Proposition 28.** (D)(N)

Let  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \overline{\mathbb{R}}$  be measurable and upper semicontinuous in the second entry. Then there exists a deterministic greedy policy for  $Q$ .

*Proof.* Since  $Q$  is upper semicontinuous in the second entry the set  $A_s = \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a)$  is non-empty and measurable for all  $s$ . Pick (by axiom of choice) an  $a_s \in A_s$  for every  $s \in \mathcal{S}$ . Then  $\pi(\cdot \mid s) = \delta_{a_s}$  is greedy with respect to  $Q$ .  $\square$

**Proposition 29.** For any  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \overline{\mathbb{R}}$  if  $\pi_Q$  is greedy with respect to  $Q$  then  $T_{\pi_Q} Q = TQ$ .

*Proof.*

$$\begin{aligned} T_{\pi_Q} Q &= r + \gamma \int Q(s, a) d\pi_Q(s, a \mid \cdot) \\ &= r + \gamma \int \int Q(s, a) d\pi_Q(a \mid s) dP(s \mid \cdot) \\ &= r + \gamma \int \max_{a \in \mathcal{A}} Q(s, a) dP(s \mid \cdot) \\ &= TQ \end{aligned}$$

$\square$

**Proposition 30.** Let  $\pi_i$  be greedy with respect to  $Q_{i-1}^*$  for  $k \in \mathbb{N}$ . Then  $Q_k^* = T_{\pi_1} \dots T_{\pi_k} Q_0$  for all  $k \in \mathbb{N}$ .

**Proposition 31.** (D)(N)

Any greedy policy with respect to  $Q^*$  is optimal and can be chosen to be deterministic.

*Proof.* By proposition 28 we can pick a greedy policy  $\pi$  for  $Q^*$  which can be chosen to be deterministic but let  $\pi$  stay general. Then by proposition 29  $T_\pi Q^* = TQ^*$  and by proposition 23  $Q_\pi = Q^*$  implying that  $\pi$  is optimal.  $\square$

---

**Algorithm 1:** Simple theoretical Q-iteration

---

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

$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : r(s, a) \leftarrow \int x dR(x \mid s, a).$

$\tilde{Q}_0 \leftarrow r$

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

$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \tilde{Q}_{k+1}(s, a) \leftarrow r(s, a) + \gamma \int \sup_{a' \in \mathcal{A}} \tilde{Q}_k(s', a') dP(s' \mid s, a)$

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

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

---

**2.4.1 Finite Q-iteration**

Concluding on the results so far we have showed how if one knows the dynamics of a stationary decision model satisfying rather broad criteria, such as continuity and compactness, the optimal policy and state-value function can be found simply by iteration over the  $T$ -operator and picking a greedy strategy (see proposition 26). Of course this is practical computationally, only if the resulting  $Q$  functions can be represented and computed in finite space and time. This is trivially the case when

**Assumption 3.**  $\mathcal{S} \times \mathcal{A}$  is finite.

Say  $|\mathcal{S}| = k$  and  $|\mathcal{A}| = \ell$ . In this case the transition operator  $P$  can be represented as a matrix of *transition probabilities*

$$P := \begin{pmatrix} P(s_1 \mid s_1, a_1) & \dots & P(s_k \mid s_1, a_1) \\ \vdots & \ddots & \vdots \\ P(s_1 \mid s_k, a_\ell) & \dots & P(s_k \mid s_k, a_\ell) \end{pmatrix}$$

then the algorithm becomes

---

**Algorithm 2:** Simple finite Q-iteration

---

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

Set  $r \leftarrow (\int r dR(\cdot \mid s_1, a_1), \dots, \int r dR(\cdot \mid s_k, a_\ell))^T$

and  $\tilde{Q}_0 = r$ .

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

Set  $m(\tilde{Q}_k) := (\max_{a \in \mathcal{A}} \tilde{Q}(s_1, a), \dots, \max_{a \in \mathcal{A}} \tilde{Q}(s_k, a))^T$

Update action-value function:

$$\tilde{Q}_{k+1} \leftarrow r + \gamma P m(\tilde{Q}_k)$$

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

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

---

**Proposition 32.** The output  $\tilde{Q}_K$  from algorithm 2 is  $K$ -optimal and  $\|\tilde{Q}_K - Q^*\|_\infty \leq \gamma^K \|Q^*\|_\infty$ .

*Proof.* See proposition 26 and proposition 30.  $\square$

## 2.5 Approximation

In this section we will look at what happens if we instead use approximations the  $Q$ -functions and  $T$  operator. We first look at a naive approach using  $Q$ -functions.

Let  $\tilde{Q}_0$  be any bounded  $Q$ -function. Suppose we approximate  $T\tilde{Q}_0$  by a  $Q$ -function  $\tilde{Q}_1$  to  $\varepsilon_1 > 0$  precision and then approximate  $T\tilde{Q}_1$  and so on getting a sequence of  $Q$ -functions satisfying

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

First observe that

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

Using this iteratively we get

$$\left| T^k \tilde{Q}_0 - \tilde{Q}_k \right| \leq \sum_{i=1}^k \gamma^{k-i} \varepsilon_i := \varepsilon_a(k)$$

Then we can bound

$$\begin{aligned} \left| Q^* - \tilde{Q}_k \right| &\leq \left| Q^* - T^k \tilde{Q}_0 \right| + \left| T^k \tilde{Q}_0 - \tilde{Q}_k \right| \\ &\leq \gamma^k \left| Q^* - \tilde{Q}_0 \right| + \varepsilon_a(k) \end{aligned}$$

The first term converges quickly while the other depends on our step-wise approximations. For example  $\varepsilon_i(k) = \varepsilon$  we easily get the bound  $\varepsilon_a(k) = \varepsilon \frac{1-\gamma^k}{1-\gamma} \leq \frac{\varepsilon}{1-\gamma}$ . Or if  $\varepsilon_i \leq c\gamma^i$  we get  $\varepsilon_a(k) \leq ck\gamma^k \rightarrow 0$  as  $k \rightarrow \infty$ . Generally if one can show that  $\varepsilon_i \rightarrow 0$  we have

**Proposition 33.**  $\sum_{i=1}^k \gamma^{k-i} \varepsilon_i \rightarrow 0$  whenever  $\varepsilon_k \rightarrow 0$  as  $k \rightarrow \infty$ .

*Proof.* Let  $\varepsilon > 0$ . Find  $N$  such that  $\varepsilon_n \leq \varepsilon(1-\gamma)/2$  for all  $n > N$  and find  $M > N$  such that  $\gamma^M \leq \varepsilon\gamma^N \left( \sum_{i=1}^N \gamma^{N-i} \varepsilon_i \right)^{-1}$ . Then for all  $m > M$

$$\sum_{i=1}^m \gamma^{m-i} \varepsilon_i \leq \gamma^{m-N} \sum_{i=1}^N \gamma^{N-i} \varepsilon_i + \sum_{i=N+1}^m \gamma^{m-i} \varepsilon(1-\gamma)/2 \leq \varepsilon/2 + \varepsilon/2 \leq \varepsilon$$

□

## 3 Hidden dynamics

In this section we will look at what can be done when the process dynamics are unknown. In this case we cannot calculate directly neither  $r, T_\pi Q$  nor  $TQ$  because the transition and reward kernels  $P, R$  are unknown.

It is clear that algorithm 1 will not work without modification in this case. Simply because  $R$  and  $P$  are not available. To make the scheme work anyway we could simply avoid taking expectations and use the random outcomes of the kernels. Leading to

---

**Algorithm 3:** Random theoretical Q-iteration (example of thought)

---

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

$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \tilde{Q}_0(s, a) \leftarrow X \sim R(\cdot \mid s, a)$ .

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

$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \tilde{Q}_{k+1}(s, a) \leftarrow r' + \gamma \sup_{a' \in \mathcal{A}} \tilde{Q}_k(s', a')$   
    where  $r' \sim R(\cdot \mid s, a), s' \sim P(\cdot \mid s, a)$ .

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

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

---

We immediately run into problems in the uncountable case, because drawing uncountably many times from a distribution is not easily defined in a sensible way. Even in the finite case, even though the  $\tilde{Q}_k$ s are well defined, they cannot converge if  $R$  is not deterministic. Therefore this approach is not attractive in a continuous or stochastic setting.

### 3.0.1 Finite case

A common way to overcome the problem of convergence is called *temporal difference* (TD) learning and is based on the following update scheme

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

Here  $r'$  and  $s'$  are the reward and next-state drawn from the reward and transition kernels, and  $\alpha_k \in [0, 1]$  is the so-called **learning rate** (of the  $k$ th step). The 'temporal difference' is also the name of term  $\alpha_k (r' + \gamma \cdot \max_{a' \in \mathcal{A}} \tilde{Q}_k(s', a') - \tilde{Q}_k(s, a))$  occurring from rearranging eq. (1). Usually the learning rate is fixed before running the algorithm (does not depend on the history) and is set to decay from 1 to 0 in some fashion as  $k \rightarrow \infty$ .

We will now look at a convergence result obtained by [Jaakkola, Jordan, Singh, 1993] of a TD algorithm using Q-functions

---

**Algorithm 4:** Simple Q-learning

---

**Input:** MDP  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$  such that  $|\mathcal{S}||\mathcal{A}| < \infty$ , number of iterations  $K$ , state-action pairs  $(s_1, a_1, \dots, s_K, a_K)$ , learning rates  $(\alpha'_1, \dots, \alpha'_K)$ , initial  $\tilde{Q}_0 : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$

Put  $\alpha_k = \delta_{(s_k, a_k)} \alpha'_k$ .

**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:

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

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

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

---

Note that only the value of the pair  $(s_k, a_k)$  are updated in each step of the algorithm (since  $\alpha_k(s, a) = 0$  for all  $(s, a) \neq (s_k, a_k)$ ).

**Theorem 7.** (Jaakkola, Jordan, Singh) 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  $\alpha_1, \alpha_2, \dots \in [0, 1]$ . The output  $\tilde{Q}_K$  of algorithm 4 converges to  $Q^*$  provided

1.  $\mathbb{P} \left( \sum_{i=1}^{\infty} \alpha_i(s, a) = \infty \right) = 1, \mathbb{P} \left( \sum_{i=1}^{\infty} \alpha_i^2(s, a) < \infty \right) = 1$ .
2.  $\text{Var}(R(\cdot \mid s, a)) < \infty$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .

3. If  $\gamma = 1$  all policies lead to a reward-free terminal state almost surely.

In the original formulation the sums of learning rates were supposed to converge *uniformly*. However this is equivalent to this formulation because of the fact that  $\mathbb{P}(\sup_{(s,a) \in \mathcal{S} \times \mathcal{A}} |f_n(s,a)| \rightarrow 0) = 1 \iff \mathbb{P}(|f_n(s,a)| \rightarrow 0) = 1, \forall (s,a) \in \mathcal{S} \times \mathcal{A}$  whenever  $\mathcal{S}, \mathcal{A}$  is finite. Notice that the first condition implies that all state-action pairs must occur infinitely often almost surely. Also notice that the second condition is automatically fulfilled under (D) since then  $\text{Var}(R(\cdot | s, a)) < \mathbb{E}(2R_{\max})^2 = 4R_{\max}$ .

### 3.0.2 Perspectives

Another approach could be to estimate  $R$  and  $P$  before or while using an algorithm like algorithm 1 using the estimated kernels. I was not able to find sources that did this, however you can argue that this idea is already employed in temporal difference learning and others.

## 4 Fan et al.

### 4.1 Introduction

#### 4.1.1 The decision model

**Setting 5** (Fan et al.). 1. An MDP ?? i.e. one state and action space  $(\mathcal{S}, \mathcal{A})$  and one transition and reward kernel  $P, R$  which only depends on the previous state-action pair.

2. Discounted with a factor strictly in the unit interval  $\gamma \in (0, 1)$ .

#### 4.1.2 Artificial neural networks

**Definition 12.** An **ANN** (Artificial Neural Network) with structure  $\{d_i\}_{i=0}^{L+1} \subseteq \mathbb{N}$ , activation functions  $\sigma_i = (\sigma_{ij} : \mathbb{R} \rightarrow \mathbb{R})_{j=1}^{d_i}$  and weights  $\{W_i \in M^{d_i \times d_{i-1}}, v_i \in \mathbb{R}^{d_i}\}_{i=1}^{L+1}$  is the function  $F : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{d_{L+1}}$

$$F = w_{L+1} \circ \sigma_L \circ w_L \circ \sigma_{L-1} \circ \dots \circ w_1$$

where  $w_i$  is the affine function  $x \mapsto W_i x + v_i$  for all  $i$ .

Here  $\sigma_i(x_1, \dots, x_{d_i}) = (\sigma_{i1}(x_1), \dots, \sigma_{id_i}(x_{d_i}))$ .

$L \in \mathbb{N}_0$  is called the number of hidden layers.

$d_i$  is the number of neurons or nodes in layer  $i$ .

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

**Definition 13** (Sparse ReLU Networks). For  $s, V \in \mathbb{R}$  a **(s,V)-Sparse ReLU Network** is an ANN  $f$  with any structure  $\{d_i\}_{i \in [L+1]}$ , all activation functions being *ReLU* i.e.  $\sigma_{ij} = \max(\cdot, 0)$  and any weights  $(W_\ell, v_\ell)$  satisfying

$$\bullet \max_{\ell \in [L+1]} \|\widetilde{W}_\ell\|_\infty \leq 1 \quad \bullet \sum_{\ell=1}^{L+1} \|\widetilde{W}_\ell\|_0 \leq s \quad \bullet \max_{j \in [d_{L+1}]} \|f_j\|_\infty \leq V$$

Here  $\widetilde{W}_\ell = (W_\ell, v_\ell)$ . The set of them we denote  $\mathcal{F}(L, \{d_i\}_{i=0}^{L+1}, s, V)$ .



### 4.1.3 Fitted Q-Iteration

We here present the algorithm which everything in this paper revolves around:

---

**Algorithm 5:** Fitted Q-Iteration Algorithm

---

**Input:** MDP  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ , function class  $\mathcal{F}$ , sampling distribution  $\nu$ , number of iterations  $K$ , number of samples  $n$ , initial estimator  $\tilde{Q}_0$

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

    Sample i.i.d. observations  $\{(S_i, A_i), i \in [n]\}$  from  $\nu$  obtain  $R_i \sim R(S_i, A_i)$  and

$S'_i \sim P(S_i, A_i)$

    Let  $Y_i = R_i + \gamma \cdot \max_{a \in \mathcal{A}} \tilde{Q}_k(S'_i, a)$

    Update action-value function:

$$\tilde{Q}_{k+1} \leftarrow \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (Y_i - f(S_i, A_i))^2$$

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

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

---

## 4.2 Notational deviations from [TODO ref YangXieWang]

Because  $\sigma$  is used ambiguously in theorem 8 we denote the probability distribution  $\sigma$  from [YangXieWang, thm. 6.2, p. 20] by  $\nu$  instead.

I avoid the shorthand defined in [YangXieWang, p. 26 bottom]:  $\|f\|_n^2 = 1/n \cdot \sum_{i=1}^n f(X_i)^2$ . and use  $p$ -norms instead. The conversion to my notation thus becomes  $\|f\|_n \rightsquigarrow \|f\|/n$ .

## 4.3 Assumptions

### 4.3.1 Hölder Smoothness

**Definition 14** (Hölder smoothness). For  $f : \mathcal{S} \rightarrow \mathbb{R}$  we define

$$\|f\|_{C_r} := \sum_{|\alpha| < \beta} \|\partial^\alpha f\|_\infty + \sum_{\|\alpha\|_1 = \lfloor \beta \rfloor} \sup_{x \neq y} \frac{|\partial^\alpha (f(x) - f(y))|}{\|x - y\|_\infty^{\beta - \lfloor \beta \rfloor}} \quad (2)$$

Where  $\alpha = (\alpha_1, \dots, \alpha_r) \in \mathbb{N}_0^r$ . And  $\partial^k$  is the partial derivative w.r.t. the  $k$ th variable. If  $\|f\|_{C_r} < \infty$  then  $f$  is **Hölder smooth**. Given a compact subset  $\mathcal{D} \subseteq \mathbb{R}^r$  the space of Hölder smooth functions on  $\mathcal{D}$  with norm bounded by  $H > 0$  is denoted

$$C_r(\mathcal{D}, \beta, H) := \left\{ f : \mathcal{D} \rightarrow \mathbb{R} \mid \|f\|_{C_r} \leq H \right\}$$

**Definition 15.** Let  $t_j, p_j \in \mathbb{N}$ ,  $t_j \leq p_j$  and  $H_j, \beta_j > 0$  for  $j \in [q]$ . We say that  $f$  is a **composition of Hölder smooth functions** when

$$f = g_q \circ \dots \circ g_1$$

for some functions  $g_j : [a_j, b_j]^{p_j} \rightarrow [a_{j+1}, b_{j+1}]^{p_{j+1}}$  that only depend on  $t_j$  of their inputs for each of their components  $g_{jk}$ , and satisfies  $g_{jk} \in C_{t_j}([a_j, b_j]_{t_j}^{t_j}, \beta_j, H_j)$ , i.e. they are Hölder smooth. We denote the class of these functions

$$\mathcal{G}(\{p_j, t_j, \beta_j, H_j\}_{j \in [q]})$$

**Definition 16.** Define

$$\mathcal{F}_0 = \{f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} \mid f(\cdot, a) \in \mathcal{F}(s, V) \forall a \in \mathcal{A}\}$$

and

$$\mathcal{G}_0 = \left\{f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} \mid f(\cdot, a) = \mathcal{G}(\{p_j, t_j, \beta_t, H_j\}_{j \in [q]}) \forall a \in \mathcal{A}\right\}$$

**Assumption 4.**  $T\mathcal{F}_0 \subseteq \mathcal{G}_0$ . I.e. it is assumed that  $Tf \in \mathcal{G}_0$  for any  $f \in \mathcal{F}_0$ , so when using the Bellman optimality operator on our sparse ReLU networks, we should stay in the class of compositions of Hölder smooth functions.

If also  $\mathcal{G}_r$  is well approximated by functions in  $\mathcal{F}_0$  then this assumption implies that  $\mathcal{F}_0$  is approximately closed under the Bellman operator  $T$  and thus that  $Q^*$  is close to  $\mathcal{F}_0$ .

We now look at a case where  $\mathcal{D} = [0, 1]^r$ ,  $q = 1$  and both the expected reward function and transition kernel is Hölder smooth.

**Proposition 34.** Assume for all  $a \in \mathcal{A}$  that  $P(s, a)$  and  $R(s, a)$  are absolutely continuous w.r.t.  $\lambda^k$  and for all  $s' \in \mathcal{S}$  that  $s \mapsto P(s' \mid s, a)$  and  $s \mapsto \mathbb{E}R(s, a)$  are both in  $C_r([0, 1]^r, \beta, H)$ . Then  $T\mathcal{F}_0 \subseteq C_r([0, 1]^r, \beta, (1 + \gamma V_{\max})H)$ .

*Proof.* Let  $f \in \mathcal{F}_0$  and  $\alpha \in \mathbb{N}_0^r$ . Observe that

$$\begin{aligned} \partial^\alpha(Tf)(s, a) &= \partial_s^\alpha(\mathbb{E}R(s, a)) + \gamma \int_{\mathcal{S}} \partial_s^\alpha \left[ \max_{a' \in \mathcal{A}} f(s', a') P(s' \mid s, a) \right] ds' \\ &\leq \partial_s^\alpha(\mathbb{E}R(s, a)) + \gamma V_{\max} \sup_{s' \in \mathcal{S}} \partial_s^\alpha P(s' \mid s, a) \end{aligned}$$

similarly

$$\begin{aligned} \partial^\alpha(Tf)(s, a) - \partial^\alpha(Tf)(s', a) &\leq \partial_s^\alpha(\mathbb{E}R(s, a)) - \partial_s^\alpha(\mathbb{E}R(s', a)) \\ &\quad + \gamma V_{\max} \sup_{s'' \in \mathcal{S}} (\partial_s^\alpha P(s'' \mid s, a) - \partial_s^\alpha P(s'' \mid s', a)) \end{aligned}$$

Thus

$$\begin{aligned} \|Tf\|_{C_r} &\leq \sum_{|\alpha| < \beta} \left( \|\partial^\alpha \mathbb{E}R(\cdot, a)\|_\infty + \gamma V_{\max} \sup_{s \in \mathcal{S}} \|\partial^\alpha P(s \mid \cdot, a)\|_\infty \right) \\ &\quad + \sum_{\|\alpha\|_1 = [\beta]} \sup_{x \neq y} \left( \frac{|\partial^\alpha(\mathbb{E}R(x, a) - \mathbb{E}R(y, a))|}{\|x - y\|_\infty^{\beta - [\beta]}} + \gamma V_{\max} \sup_{s \in \mathcal{S}} \frac{|\partial^\alpha(\mathbb{E}P(s \mid x, a) - \mathbb{E}P(s \mid y, a))|}{\|x - y\|_\infty^{\beta - [\beta]}} \right) \\ &\leq H + \gamma V_{\max} H = (1 + \gamma V_{\max})H \end{aligned}$$

□

#### 4.3.2 Concentration coefficients

**Definition 17** (Concentration coefficients). Let  $\nu_1, \nu_2 \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$  be probability measures, absolutely continuous w.r.t.  $m_\lambda$ . Define

$$\kappa(m, \nu_1, \nu_2) = \sup_{\pi_1, \dots, \pi_m} \left[ \mathbb{E}_{v_2} \left( \frac{d(P^{\pi_m} \dots P^{\pi_1} \nu_1)}{d\nu_2} \right)^2 \right]^{1/2}$$

**Assumption 5.** Let  $\nu$  be the sampling distribution from the algorithm, and  $\mu$  the distribution over which we measure the error in the main theorem, then we assume

$$(1 - \gamma)^2 \sum_{m \geq 1} \gamma^{m-1} m \kappa(m, \mu, \nu) = \phi_{\mu, \nu} < \infty$$

#### 4.4 Main theorem

**Theorem 8** (Yang, Xie, Wang). Let  $\mu$  be any distribution over  $\mathcal{S} \times \mathcal{A}$ . Make assumption 4 and assumption 5 with the constants  $\phi_{\mu,\nu} > 0$ ,  $q \in \mathbb{N}$  and  $\{p_j, t_j, \beta_j, H_j\}_{j \in [q]}$ . Furthermore assume that there exists a constant  $\xi > 0$  such that

$$\max \left\{ \sum_{j=1}^q (t_j + \beta_j + 1)^{3+t_k}, \sum_{j=1}^q \log(t_j + \beta_j), \max_{j \in [q]} p_j \right\} \leq (\log n)^\xi$$

Set  $\beta_j^* = \beta_j \prod_{\ell=j+1}^q \min(\beta_\ell, 1)$  for  $j \in [q-1]$ ,  $\beta_q^* = 1$ ,  $\alpha^* = \max_{j \in [q]} t_j / (2\beta_j^* + t_j)$ ,  $\xi^* = 1 + 2\xi$  and  $\kappa^* = \min_{j \in [q]} \beta_j^* / t_j$ . Then there exists a class of ReLU networks

$$\mathcal{F}_0 = \{f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} : f(\cdot, a) \in \mathcal{F}(\tilde{L}, \{\tilde{d}_j\}_{j=0}^{\tilde{L}+1}, \tilde{s}) \mid a \in \mathcal{A}\}$$

with structure satisfying

$$\tilde{L} \lesssim (\log n)^{\xi^*}, \tilde{d}_0 = r, \tilde{d}_j \leq 6n^{\alpha^*} (\log n)^{\xi^*}, d_{L+1} = 1, \tilde{s} \lesssim n^{\alpha^*} \cdot (\log n)^{\xi^*}$$

such that when running algorithm 5 with  $\mathcal{F}_0$  and  $n$  is sufficiently large

$$\|Q^* - Q^{\pi_K}\|_{1,\mu} \leq C_\varepsilon \frac{\phi_{\mu,\nu}\gamma}{(1-\gamma)^2} V_{\max}^2 n^{\max\{-2\alpha^*\kappa^*, (\alpha^*-1)/2\}} \log(n)^{1+2\xi^*} + \frac{4\gamma}{(1-\gamma)^2} R_{\max} \gamma^K$$

where  $C_\varepsilon > 0$  is a constant not depending on  $n$  or  $K$ .

#### 4.5 Proofs

**Theorem 9** (Error Propagation). Let  $\{\tilde{Q}_i\}_{0 \leq i \leq K}$  be the iterates of the fitted Q-iteration algorithm. Then

$$\|Q^* - Q^{\pi_K}\|_{1,\mu} \leq \frac{2\phi_{\mu,\nu}\gamma}{(1-\gamma)^2} \cdot \varepsilon_{\max} + \frac{4\gamma^{K+1}}{(1-\gamma)^2} \cdot R_{\max}$$

Where

$$\varepsilon_{\max} = \max_{k \in [K]} \|T\tilde{Q}_{k-1} - \tilde{Q}_k\|_{2,\nu}$$

**Theorem 10** (One-step Approximation Error). Let

- $\mathcal{F} \subseteq \mathcal{B}(\mathcal{S} \times \mathcal{A}, V_{\max})$  be a class of bounded measurable functions
- $\mathcal{G} = T(\mathcal{F})$  the class of functions obtainable by applying  $T$  to some function in  $\mathcal{F}$ .
- $\nu \in \mathcal{P}(\mathcal{S}, \mathcal{A})$  be a probability measure
- $(S_i, A_i)_{i \in [n]}$  be  $n$  i.i.d. samples following  $\nu$
- $(R_i, S'_i)_{i \in [n]}$  be the rewards and next states corresponding to the samples
- $Q \in \mathcal{F}$  be fixed
- $Y_i = R_i + \gamma \max_{a \in \mathcal{A}} Q(S'_i, a)$
- $\hat{Q} = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(S_i, A_i) - Y_i)^2$
- $\kappa \in (0, 1]$ ,  $\delta > 0$  be fixed
- $\mathcal{N}(\delta, \mathcal{F}, \|\cdot\|_\infty)$  a minimal  $\delta$ -covering of  $\mathcal{F}$  w.r.t.  $\|\cdot\|_\infty$
- $N_\delta = |\mathcal{N}(\delta, \mathcal{F}, \|\cdot\|_\infty)|$  the number of elements in this covering

Then

$$\begin{aligned} \left\| \hat{Q} - TQ \right\|_{\nu}^2 &\leq \frac{(1+\kappa)^2}{\kappa} \frac{1}{n} C_2^2 V_{\max}^2 \log(N_{\delta}) + (1+\kappa) \left( \delta C_2^2 V_{\max}^2 \log(N_{\delta}) + \omega(\mathcal{F}) \right) \\ &\quad + 8\sqrt{2} V_{\max} n^{-1/2} \sqrt{\log N_{\delta}} + 8V_{\max}(n^{-1} + \delta) \end{aligned}$$

Where

$$\omega(\mathcal{F}) = \sup_{g \in \mathcal{G}} \inf_{f \in \mathcal{F}} \frac{1}{n} \mathbb{E} \|f - g\|_{\nu}^2$$

*Proof of main theorem.* Using theorem 9 we get

$$\|Q^* - Q^{\pi_K}\|_{1,\mu} \leq \frac{2\phi_{\mu,\nu}\gamma}{(1-\gamma)^2} \varepsilon_{\max} + \frac{4\gamma^{K+1}}{(1-\gamma)^2} R_{\max} \quad (3)$$

where  $\varepsilon_{\max} = \max_{k \in [K]} \|T\tilde{Q}_{k-1} - \tilde{Q}_k\|_{2,\nu}$ . Using theorem 10 with  $Q = \tilde{Q}_{k-1}$ ,  $\mathcal{F} = \mathcal{F}_0$ ,  $\epsilon = 1$  and  $\delta = 1/n$ , we get

$$\varepsilon_{\max} \leq 6n^{-1} C_2^2 V_{\max}^2 \log(N_0) + 2\omega(\mathcal{F}_0) + 8\sqrt{2} V_{\max} n^{-1/2} \sqrt{\log N_0} + 16V_{\max} n^{-1} \quad (4)$$

where  $N_0 = |\mathcal{N}(1/n, \mathcal{F}_0, \|\cdot\|_{\infty})|$ . The remains only to bound  $\omega(\mathcal{F}_0)$  and  $N_0$ , starting with  $\omega(\mathcal{F}_0)$ .

**Step 1.** We want to employ the following lemma by [Schmidt-Hieber 2019, thm. 5, p. 22]

**Lemma 2** (Approximation of Hölder Smooth Functions). Let  $m, M \in \mathbb{Z}_+$  with  $N \geq \max\{(\beta + 1)^r, (H + 1)e^r\}$ ,  $L = 8 + (m + 5)(1 + \lceil \log_2(r + \beta) \rceil)$ ,  $d_0 = r$ ,  $d_j = 6(r + \lceil \beta \rceil)N$ ,  $d_{L+1} = 1$ . Then for any  $g \in \mathcal{C}_r([0, 1]^r, \beta, H)$  there exists a ReLU network  $f \in \mathcal{F}(L, \{d_j\}_{j=0}^{L+1}, s, \infty)$  with  $s \leq 141(r + \beta + 1)^{3+r} N(m + 6)$  such that

$$\|f - g\|_{\infty} \leq (2H + 1)6^r N(1 + r^2 + \beta^2)2^{-m} + H3^{\beta} N^{-\beta/r}$$

to each Hölder smooth part of  $g$  and then piece it together somehow, using that ReLU networks are easily stitched together into bigger ReLU networks. Therefore the first step is to refit our Hölder Smooth compositions in  $\mathcal{G}_0$  to be defined on a hyper-cube instead. This is a relatively simple procedure:

Let  $f \in \mathcal{G}_0$  then  $f(\cdot, a) \in \mathcal{G}(\{p_j, t_j, \beta_j, H_j\})$  for all  $a \in \mathcal{A}$ . Therefore  $f(\cdot, a) = g_q \circ \dots \circ g_1$  where the (sub-)components  $(g_{jk})_{k=1}^{p_j+1} = g_j$  satisfy

$$g_{jk} \in C_{t_j}([a_j, b_j]^{t_j}, \beta_j, H_j), \quad j \in [q], k \in [p_j + 1] \quad (5)$$

Here  $a_1 = 0, b_1 = 1$  and,  $a_j < b_j \in \mathbb{R}$  are some real numbers for  $2 \leq j \leq q$ . Notice that the Hölder smooth condition implies that  $g_{jk}([a_j, b_j]^{t_j}) \subseteq [-H_j, H_j]$ . Define

$$\begin{aligned} h_1 &= g_1 / (2H_1) + 1/2 \\ h_j(u) &= g_j(2H_{j-1}u - H_{j-1}) / (2H_j) + 1/2, & j \in \{2, \dots, q-1\} \\ h_q(u) &= g_q(2H_{q-1}u - H_{q-1}) \end{aligned} \quad (6)$$

Then  $g_q \circ \dots \circ g_1 = h_q \circ \dots \circ h_1$  and

$$\begin{aligned} h_{1k} &\in C_{t_1}([0, 1]^{t_1}, \beta_1, 1) \\ h_{jk} &\in C_{t_j}([0, 1]^{t_j}, \beta_j, (2H_{j-1})^{\beta_j}), & j \in \{2, \dots, q-1\} \\ h_q &\in C_{t_q}([0, 1]^{t_q}, \beta_q, H_q(2H_{q-1})^{\beta_q}) \end{aligned} \quad (7)$$

Define  $N := \max_{j \in [q]} n^{t_j/(2\beta_j^* + t_j)}$ ,  $\eta := \log \left( (2W + 1)6^{t_j} N / (W3^{\beta_j} N^{-\beta_j/t_j}) \right)$ , and  $m := \eta \lceil \log_2 n \rceil$ , and assume  $n$  is sufficiently large such that  $N \geq \max \{(\beta_j + 1)^{t_j}, (H_j + 1)e^{t_j} \mid j \in [q]\}$ .

$$W := \max \left( \left\{ (2H_{j-1})^{\beta_j} \mid 1 \leq j \leq q-1 \right\} \cup \left\{ H_q(2H_{q-1})^{\beta_q}, 1 \right\} \right) \quad (8)$$

By lemma 2 there exists a ReLU network

$$\hat{h}_{jk} \in \mathcal{F} \left( L_j + 2, \left\{ t_j, \tilde{d}_j p_{j+1}, \dots, \tilde{d}_j p_{j+1}, p_{j+1} \right\}, (\tilde{s}_j + 4) \cdot p_{j+1} \right) \quad (9)$$

where  $\tilde{d}_j = 6(t_j + \lceil \beta_j \rceil)N$  and  $\tilde{s}_j \leq 141(t_j + \beta_j + 1)^{3+t_j} N(m + 6)$  such that

$$\left\| \hat{h}_{jk} - h_{jk} \right\|_{\infty} \leq (2W + 1)6^{t_j} N 2^{-m} + W3^{\beta_j} N^{-\beta_j/t_j} \leq 2W3^{\beta_j} N^{-\beta_j/t_j} \quad (10)$$

Since  $h_{j+1}$  is defined on  $[0, 1]^{t_{j+1}}$  but  $\tilde{h}_j$  takes values in  $\mathbb{R}$  we need to restrict  $\tilde{h}_j$  somehow to stitch the two together (by function composition). This is easily done by

**Lemma 3.** Restriction to  $[0, 1]$  is expressible as a two-layer ReLU network with 4 non-zero weights.

*Proof.* Namely  $\tau(u) = 1 - (1 - u)_+ = \min \{ \max \{ u, 0 \}, 1 \}$ .  $\square$

Now define  $\tilde{h}_{jk} = \tau \circ \hat{h}_{jk}$  (and  $\tilde{h}_j = (\tilde{h}_{jk})_{k \in [p_{j+1}]}$ ). Then

$$\tilde{h}_{jk} \in \mathcal{F} \left( L_j + 2, \left\{ t_j, \tilde{d}_j, \dots, \tilde{d}_j, 1 \right\}, (\tilde{s}_j + 4)p_{j+1} \right) \quad (11)$$

and since  $h_{jk}([0, 1]^{t_j}) \in [0, 1]$  by eq. (10)

$$\left\| \tilde{h}_{jk} - h_{jk} \right\|_{\infty} = \left\| \tau \circ \hat{h}_{jk} - \tau \circ h_{jk} \right\|_{\infty} \quad (12)$$

$$\leq \left\| \hat{h}_{jk} - h_{jk} \right\|_{\infty} \quad (13)$$

$$\leq 2W3^{-\beta_j} N^{-\beta_j/t_j} \quad (14)$$

**Step 2.** Now define  $\tilde{f} : \mathcal{S} \rightarrow \mathbb{R}$  as  $\tilde{f} = \tilde{h}_1 \circ \dots \circ \hat{h}_1$ . If we set  $\tilde{L} := \sum_{j=1}^q (L_j + 2)$ ,  $\tilde{d} := \max_{j \in [q]} \tilde{d}_j p_{j+1}$  and  $\tilde{s} := \sum_{j=1}^q (\tilde{s}_j + 4)p_{j+1}$ . Then  $\tilde{f} \in \mathcal{F} \left( \tilde{L}, \left\{ r, \tilde{d}, \dots, \tilde{d}, 1 \right\}, \tilde{s} \right)$ . We now take a moment to verify the size of the constants involved in the network. Starting with  $\tilde{L}$ .

$$\begin{aligned} \tilde{L} &\leq \sum_{j=1}^q (L_j + 2) \\ &= \sum_{j=1}^q (8 + (\eta \lceil \log_2 n \rceil + 5)(1 + \lceil \log_2(\beta_j + t_j) \rceil)) \\ &\leq \sum_{j=1}^q (8 + (\eta \log_2 n + \eta + 5)(2 + \log_2(\beta_j + t_j))) \\ &\leq 8q + (2\eta + 5) \log_2(n) \sum_{j=1}^q (2 + \log_2(\beta_j + t_j)) \\ &\leq 8q + (2\eta + 5) \log_2(n) (2q + \log(n)^\xi) \\ &\leq (10q + 1)(2\eta + 5) \log_2(e) \log(n)^{1+\xi} \\ &\leq C_{\tilde{L}} \log(n)^{1+2\xi} \end{aligned}$$

where  $C_{\tilde{L}} = (10q + 1)(2\eta + 5) \log_2(e)$ . For  $\tilde{d}$  we have

$$\begin{aligned}
\tilde{d} &= \max_{j \in [q]} \tilde{d}_j p_{j+1} \\
&= \max_{j \in [q]} 6(t_j + \beta_j + 1) N p_{j+1} \\
&\leq 6N (\max_{j \in [q]} p_j) (\max_{j \in [q]} (t_j + \beta_j + 1)) \\
&\leq 6N (\log n)^{2\xi} \\
&\leq 6n^{\alpha^*} (\log n)^{\xi^*}
\end{aligned}$$

and for  $\tilde{s}$

$$\begin{aligned}
\tilde{s} &= \sum_{j=1}^q (\tilde{s}_j + 4) p_{j+1} \\
&\leq \sum_{j=1}^q (141N(m+6)(t_j + \beta_j + 1)^{3+t_j} + 4) p_{j+1} \\
&\leq 142N (\log n)^\xi (2\eta + 6) \log_2(n) \sum_{j=1}^q (t_j + \beta_j + 1)^{3+t_j} \\
&\leq 142N (\log n)^\xi (2\eta + 6) \log_2(e) \log(n) (\log n)^\xi \\
&= 142N (2\eta + 6) \log_2(e) (\log n)^{1+2\xi} \\
&= C_{\tilde{s}} n^{\alpha^*} (\log n)^{\xi^*}
\end{aligned}$$

where  $C_{\tilde{s}} = 142(2\eta + 6) \log_2(e)$ . Now we bound  $\|\tilde{f} - f(\cdot, a)\|_\infty$ . Define  $G_j = h_j \circ \dots \circ h_1$ ,  $\tilde{G}_j = \tilde{h}_j \circ \dots \circ \tilde{h}_1$  for  $j \in [q]$ ,  $\lambda_j = \prod_{\ell=j+1}^q (\beta_\ell \wedge 1)$  for all  $j \in [q-1]$  and  $\lambda_q = 1$ . We have

$$\begin{aligned}
\|G_j - \tilde{G}_j\|_\infty &= \|h_j \circ G_{j-1} - h_j \circ \tilde{G}_{j-1} + h_j \circ \tilde{G}_{j-1} - \tilde{h}_j \circ \tilde{G}_{j-1}\| \\
&\leq \|h_j \circ \tilde{G}_{j-1} - h_j \circ G_{j-1}\|_\infty + \|h_j \circ \tilde{G}_{j-1} - \tilde{h}_j \circ \tilde{G}_{j-1}\|_\infty \\
&\leq W \|G_{j-1} - \tilde{G}_{j-1}\|_\infty^{\beta_j \wedge 1} + \|\tilde{h}_j - h_j\|_\infty^{\lambda_j}
\end{aligned}$$

so by induction and eq. (10)

$$\begin{aligned}
\|f(\cdot, a) - \tilde{f}\|_\infty &= \|G_q - \tilde{G}_q\|_\infty \\
&\leq W^q \sum_{j=1}^q \|\tilde{h}_j - h_j\|_\infty^{\lambda_j} \\
&\leq W^q \sum_{j=1}^q \left( 2W 3^{\beta_j} N^{-\beta_j/t_j} \right)^{\lambda_j} \\
&\leq 2q 3^{\max_{j \in [q]} \beta_j^*} W^{q+1} \max_{j \in [q]} N^{-\beta_j^*/t_j} \\
&\leq c_N^{1/2} \max_{j \in [q]} n^{-\alpha^* \beta_j^*/t_j} \\
&\leq c_N^{1/2} n^{-\alpha^* \min_{j \in [q]} \beta_j^*/t_j}
\end{aligned}$$

and therefore

$$\begin{aligned}
\omega(\mathcal{F}_0) &= \sup_{g \in \mathcal{G}} \inf_{f \in \mathcal{F}} \\
&\leq C_N n^{-2\alpha^* \min_{j \in [q]} \beta_j^*/t_j} \\
&\leq C_N n^{-2\alpha^* \kappa^*}
\end{aligned} \tag{15}$$

where we define  $\kappa^* = \min_{j \in [q]} \beta_j^*/t_j$ .

**Step 3.** Finally what is left is to bound the covering number of  $\mathcal{F}_0$ . Denote by  $\mathcal{N}_\delta$  the  $\delta$ -covering of  $\mathcal{F}(\tilde{L}, \{\tilde{d}_j\}_{j=1}^{\tilde{L}+1}, \tilde{s})$  by

$$\mathcal{N}_\delta := \mathcal{N}\left(\delta, \mathcal{F}\left(\tilde{L}, \{\tilde{d}_j\}_{j=1}^{\tilde{L}+1}, \tilde{s}\right), \|\cdot\|_\infty\right)$$

Since  $\mathcal{N}_\delta$  is a covering, for any  $f \in \mathcal{F}_0$  and  $a \in \mathcal{A}$  you can find a  $g_a \in \mathcal{N}_\delta$  such that  $\|f(\cdot, a) - g_a\|_\infty < \delta$ . Now let  $g : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} = (s, a) \mapsto g_a(s)$ . Then  $\|f - g\|_\infty < \delta$ , so we can bound the covering number of  $\mathcal{F}_0$  by

$$|\mathcal{N}(\delta, \mathcal{F}, \|\cdot\|_\infty)| \leq |\mathcal{N}_\delta|^{|\mathcal{A}|}$$

We now utilize a lemma found in [todo: ref to Anthony Bartlett 2009]

**Lemma 4** (Covering Number of ReLU Networks). Consider the family of ReLU networks

$$\mathcal{F}\left(L, \{d_j\}_{j=0}^{L+1}, s, V_{\max}\right)$$

where  $\mathcal{F}$  is defined in definition 13. Let  $D := \prod_{\ell=1}^{L+1} (d_\ell + 1)$ . Then for any  $\delta > 0$

$$\mathcal{N}\left(\delta, \mathcal{F}\left(L, \{d_j\}_{j=0}^{L+1}, s, V_{\max}\right), \|\cdot\|_\infty\right) \leq (2(L+1)D^2/\delta)^{s+1}$$

*Proof.* We refer to theorem 14.5 in [todo: ref to Anthony Bartlett, thm. 14.5]. □

With lemma 4 and  $n$  sufficiently large we can bound

$$\begin{aligned} \log N_0 &= \log |\mathcal{N}(1/n, \mathcal{F}_0, \|\cdot\|_\infty)| \\ &\leq |\mathcal{A}| \cdot \log |\mathcal{N}_{1/n}| \\ &\leq |\mathcal{A}| (\tilde{s} + 1) \log(2(\tilde{L} + 1)\tilde{D}^2 n) \\ &\leq |\mathcal{A}| (c_{\tilde{s}} n^{\alpha^*} \log(n)^{\xi^*} + 1) 2 \log \left( 2(c_{\tilde{L}} \log(n)^{\xi^*} + 1) \prod_{\ell=1}^{\tilde{L}+1} (\tilde{d}_\ell + 1) \right) \\ &\leq 2|\mathcal{A}| (c_{\tilde{s}} n^{\alpha^*} \log(n)^{\xi^*} + 1) \log \left( 2(c_{\tilde{L}} \log(n)^{\xi^*} + 1) (6n^{\alpha^*} \log(n)^{\xi^*} + 1)^{\tilde{L}+1} \right) \\ &\leq 4|\mathcal{A}| c_{\tilde{s}} n^{\alpha^*} \log(n)^{\xi^*} (\tilde{L} + 1) \log \left( 24c_{\tilde{L}} \log(n)^{\xi^*} n^{\alpha^*} \log(n)^{\xi^*} \right) \\ &\leq 8|\mathcal{A}| c_{\tilde{s}} n^{\alpha^*} \log(n)^{\xi^*} c_{\tilde{L}} \log(n)^{\xi^*} (\alpha^* + 2) \log(n) \\ &= 8c_{\tilde{s}} c_{\tilde{L}} (\alpha^* + 2) n^{\alpha^*} \log(n)^{1+2\xi^*} \end{aligned} \tag{16}$$

Finally, combining eq. (3), eq. (4), eq. (15), eq. (16) and fiddling around with constants one obtains

$$\|Q^* - Q^{\pi_K}\|_{1,\mu} \leq C_\varepsilon \frac{\phi_{\mu,\nu}\gamma}{(1-\gamma)^2} V_{\max}^2 n^{\max\{-2\alpha^* \kappa^*, (\alpha^*-1)/2\}} \log(n)^{1+2\xi^*} + \frac{4\gamma}{(1-\gamma)^2} R_{\max} \gamma^K$$

where

$$C_\varepsilon = 160C_2^2 C_{\tilde{s}} C_{\tilde{L}} (\alpha^* + 2) + 4C_N + 32$$

only depends on the constants in assumption 5 finishing the proof. □

**Lemma 5.**  $TQ \geq T^\pi Q$  for any policy  $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$  and any action value function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ .

*Proof.*

$$\begin{aligned}
(TQ)(s, a) &= \mathbb{E} \left( R(s, a) + \gamma \max_{a'} Q(S', a') \mid S' \sim P(\cdot \mid s, a) \right) \\
&\geq \mathbb{E} (R(s, a) + \gamma Q(S', A') \mid S' \sim P(\cdot \mid s, a), A' \sim \pi(\cdot \mid S')) \\
&= T^\pi Q(s, a)
\end{aligned}$$

□

**Lemma 6.** Let  $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  be an action-value function,  $\tau_1, \dots, \tau_m$  be policies and  $\mu \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$  be a probability measure. Then

$$\mathbb{E}_\mu[(P^{\tau_m} \dots P^{\tau_1})(f)] \leq \kappa(k - i + j; \mu, \nu) \|f\|_{2, \nu}$$

For any measure  $\nu \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$  which is absolutely continuous w.r.t.  $(P^{\tau_m} \dots P^{\tau_1})(\mu)$ . Here  $\kappa$  is the concentration coefficients defined in definition 17.

*Proof.* Recall that

$$\begin{aligned}
\kappa(m; \mu, \nu) &:= \sup_{\pi_1, \dots, \pi_m} \left[ \mathbb{E}_\nu \left| \frac{d(P^{\pi_m} \dots P^{\pi_1} \mu)}{d\nu} \right|^2 \right]^{1/2} \\
&= \sup_{\pi_1, \dots, \pi_m} \left\| \frac{d(P^{\pi_m} \dots P^{\pi_1} \mu)}{d\nu} \right\|_{2, \nu}
\end{aligned}$$

Thus

$$\mathbb{E}_\mu[(P^{\tau_m} \dots P^{\tau_1})(f)] = \int (P^{\tau_m} \dots P^{\tau_1})(f) d\mu \quad (17)$$

$$= \int f d(P^{\tau_m} \dots P^{\tau_1} \mu) \quad (18)$$

$$= \int f \frac{d(P^{\tau_m} \dots P^{\tau_1} \mu)}{d\nu} d\nu \quad (19)$$

$$\leq \left\| \frac{d(P^{\tau_m} \dots P^{\tau_1} \mu)}{d\nu} \right\|_{2, \nu} \cdot \|f\|_{2, \nu} \quad (20)$$

$$\leq \kappa(m, \mu, \nu) \|f\|_{2, \nu} \quad (21)$$

Where eq. (19) is due to the Radon-Nikodym theorem and eq. (20) is Cauchy-Schwarz. □

*Proof of theorem 9.* First some things to keep in mind during the proof. Recall that  $V_{\max} = R_{\max}/(1 - \gamma)$  and that  $\pi_Q$  is the greedy policy w.r.t.  $Q$ . Denote

$$\pi_i = \pi_{\tilde{Q}_i}, \quad Q_{i+1} = T\tilde{Q}_i, \quad \varrho_i = Q_i - \tilde{Q}_i, \quad \text{for } i \in \{0, \dots, K+1\}$$

Note that for any policy  $\pi$ ,  $P^\pi$  is linear and 1-contrative on  $\mathcal{L}^\infty(\mathcal{S} \times \mathcal{A})$ . Also

$$T^\pi Q^\pi = Q^\pi, \quad TQ = T^{\pi_Q} Q, \quad TQ^* = Q^* = Q^{\pi^*}$$

where  $\pi^*$  is greedy w.r.t.  $Q^*$ . If  $f > f'$  for  $f, f' : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  then  $P^\pi f \geq P^\pi f'$ .

The proof consists of four steps.



**Step 1** We start by relating  $Q^* - Q^{\pi_K}$ , the quantity of interest, to  $Q^* - \tilde{Q}_K$ , which is more related to the output of the algorithm. Using lemma 5 we can make the upper bound

$$\begin{aligned}
Q^* - Q^{\pi_K} &= T^{\pi^*} Q^* - T^{\pi_K} Q^{\pi_K} \\
&= T^{\pi^*} Q^* + (T^{\pi^*} \tilde{Q}_K - T^{\pi^*} \tilde{Q}_K) + (T \tilde{Q}_K - T \tilde{Q}_K) - T^{\pi_K} Q^{\pi_K} \\
&= (T^{\pi^*} \tilde{Q}_K - T \tilde{Q}_K) + (T^{\pi^*} Q^* - T^{\pi^*} \tilde{Q}_K) + (T \tilde{Q}_K - T^{\pi_K} Q^{\pi_K}) \\
&\leq (T^{\pi^*} Q^* - T^{\pi^*} \tilde{Q}_K) + (T \tilde{Q}_K - T^{\pi_K} Q^{\pi_K}) \\
&= (T^{\pi^*} Q^* - T^{\pi^*} \tilde{Q}_K) + (T^{\pi_K} \tilde{Q}_K - T^{\pi_K} Q^{\pi_K}) \\
&= \gamma P^{\pi^*} (Q^* - \tilde{Q}_K) + \gamma P^{\pi_K} (\tilde{Q}_K - Q^{\pi_K}) \\
&= \gamma (P^{\pi^*} - P^{\pi_K}) (Q^* - \tilde{Q}_K) + \gamma P^{\pi_K} (Q^* - Q^{\pi_K})
\end{aligned} \tag{22}$$

This implies

$$(I - \gamma P^{\pi_K}) (Q^* - Q^{\pi_K}) \leq \gamma (P^{\pi^*} - P^{\pi_K}) (Q^* - \tilde{Q}_K)$$

Since  $\gamma P^{\pi_K}$  is  $\gamma$ -contractive,  $U = (I - \gamma P^{\pi_K})^{-1}$  exists as a bounded operator on  $\mathcal{L}^\infty(\mathcal{S} \times \mathcal{A})$  and equals

$$U = \sum_{i=0}^{\infty} \gamma^i (P^{\pi_K})^i$$

From this we also see that  $f \geq f' \implies Uf \geq Uf'$  for any  $f, f' : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ . Therefore we can apply  $U$  on both sides of eq. (22) to obtain

$$Q^* - Q^{\pi_K} \leq \gamma U^{-1} (P^{\pi^*} (Q^* - \tilde{Q}_K) - P^{\pi_K} (Q^* - \tilde{Q}_K)) \tag{23}$$

**Step 2** Using lemma 5 for any  $i \in [K]$  we can get an upper bound

$$\begin{aligned}
Q^* - \tilde{Q}_{i+1} &= Q^* + (T \tilde{Q}_i - T \tilde{Q}_i) - \tilde{Q}_{i+1} + (T^{\pi^*} \tilde{Q}_i - T^{\pi^*} \tilde{Q}_i) \\
&= (Q^* - T^{\pi^*} \tilde{Q}_i) + (T \tilde{Q}_i - \tilde{Q}_{i+1}) + (T^{\pi^*} \tilde{Q}_i - T \tilde{Q}_i) \\
&= (T^{\pi^*} Q^* - T^{\pi^*} \tilde{Q}_i) + \varrho_{i+1} + (T^{\pi^*} \tilde{Q}_i - T \tilde{Q}_i) \\
&\leq T^{\pi^*} Q^* - T^{\pi^*} \tilde{Q}_i + \varrho_{i+1} \\
&= \gamma P^{\pi^*} (Q^* - \tilde{Q}_i) + \varrho_{i+1}
\end{aligned} \tag{24}$$

and a lower bound

$$\begin{aligned}
Q^* - \tilde{Q}_{i+1} &= Q^* + (T \tilde{Q}_i - T \tilde{Q}_i) - \tilde{Q}_{i+1} + (T^{\pi_i} Q^* - T^{\pi_i} Q^*) \\
&= (T^{\pi_i} Q^* - T^{\pi_i} \tilde{Q}_i) + \varrho_{i+1} + (T Q^* - T^{\pi_i} Q^*) \\
&\geq T^{\pi_i} Q^* - T^{\pi_i} \tilde{Q}_i + \varrho_{i+1} \\
&= \gamma P^{\pi_i} (Q^* - \tilde{Q}_i) + \varrho_{i+1}
\end{aligned} \tag{25}$$

Applying eq. (24) and eq. (25) iteratively we get

$$Q^* - \tilde{Q}_K \leq \gamma^K (P^{\pi^*})^K (Q^* - \tilde{Q}_0) + \sum_{i=0}^{K-1} \gamma^{K-1-i} (P^{\pi^*})^{K-1-i} \varrho_{i+1} \tag{26}$$

and

$$Q^* - \tilde{Q}_K \geq \gamma^K (P^{\pi_{K-1}} \dots P^{\pi_0}) (Q^* - \tilde{Q}_0) + \sum_{i=0}^{K-1} \gamma^{K-1-i} (P^{\pi_{K-1}} \dots P^{\pi_{i+1}}) \varrho_{i+1} \tag{27}$$

**Step 3** Combining eq. (26) and eq. (27) with eq. (23) we get

$$Q^* - Q^{\pi_K} \leq U^{-1} \left( \gamma^{K+1} ((P^{\pi_K})^{K+1} - P^{\pi_K} \dots P^{\pi_0})(Q^* - \tilde{Q}_0) + \sum_{i=0}^{K-1} \gamma^{K-i} ((P^{\pi_K})^{K-i} - P^{\pi_K} \dots P^{\pi_{i+1}}) \varrho_{i+1} \right) \quad (28)$$

For shorthand define constants

$$\alpha_i = \frac{(1-\gamma)\gamma^{K-i-1}}{1-\gamma^{K+1}} \text{ for } 0 \leq i \leq K-1 \text{ and } \alpha_K = \frac{(1-\gamma)\gamma^K}{1-\gamma^{K+1}} \quad (29)$$

(note that  $\sum_{i=0}^K \alpha_i = 1$ ) and operators

$$O_i = (1-\gamma)/2U^{-1}[(P^{\pi_K})^{K-i} + (P^{\pi_K} \dots P^{\pi_{i+1}})] \quad (30)$$

$$O_K = (1-\gamma)/2U^{-1}[(P^{\pi_K})^{K+1} + (P^{\pi_K} \dots P^{\pi_0})] \quad (31)$$

Then by eq. (28)

$$|Q^* - Q^{\pi_K}| \leq \frac{2\gamma(1-\gamma^{K+1})}{(1-\gamma)^2} \left[ \sum_{i=0}^{K-1} \alpha_i O_i |\varrho_{i+1}| + \alpha_K O_K |Q^* - \tilde{Q}_0| \right] \quad (32)$$

So by linearity of expectation

$$\|Q^* - Q^{\pi_K}\|_{1,\mu} = \mathbb{E}_\mu |Q^* - Q^{\pi_K}| \quad (33)$$

$$\leq \frac{2\gamma(1-\gamma^{K+1})}{(1-\gamma)^2} \left[ \sum_{i=0}^{K-1} \alpha_i \mathbb{E}_\mu (O_i |\varrho_{i+1}|) + \alpha_K \mathbb{E}_\mu (O_K |Q^* - \tilde{Q}_0|) \right] \quad (34)$$

With the bound on rewards we (crudely) estimate

$$\mathbb{E}_\mu O_K |Q^* - \tilde{Q}_0| \leq 2V_{\max} = 2R_{\max}/(1-\gamma) \quad (35)$$

The remaining difficulty lies in  $\mathbb{E}_\mu (O_i |\varrho_{i+1}|)$ .

**Step 4** Using the sum expansion of  $U^{-1}$  we get

$$\mathbb{E}_\mu (O_i |\varrho_{i+1}|) \quad (36)$$

$$= \frac{1-\gamma}{2} \mathbb{E}_\mu \left( U^{-1} [(P^{\pi_K})^{K-i} + P^{\pi_K} \dots P^{\pi_{i+1}}] |\varrho_{i+1}| \right) \quad (37)$$

$$= \frac{1-\gamma}{2} \mathbb{E}_\mu \left( \sum_{j=0}^{\infty} [(P^{\pi_K})^j (P^{\pi_K})^{K-i} + (P^{\pi_K})^{j+1} P^{\pi_{K-1}} \dots P^{\pi_{i+1}}] |\varrho_{i+1}| \right) \quad (38)$$

$$= \frac{1-\gamma}{2} \sum_{j=0}^{\infty} \mathbb{E}_\mu \left( [(P^{\pi_K})^j (P^{\pi_K})^{K-i} + (P^{\pi_K})^{j+1} P^{\pi_{K-1}} \dots P^{\pi_{i+1}}] |\varrho_{i+1}| \right) \quad (39)$$

Notice that there are  $K-i+j$   $P$ -operators on both terms in the sum. Therefore we can employ lemma 6 twice. Moreover define  $\varepsilon_{\max} = \max_{i \in [K]} \|\varrho_i\|_{2,\nu}$ . Then

$$\begin{aligned} \mathbb{E}_\mu (O_i |\varrho_{i+1}|) &\leq (1-\gamma) \sum_{j=0}^{\infty} \gamma^j \kappa(K-i+j; \mu, \nu) \|\varrho_{i+1}\|_{2,\nu} \\ &\leq \varepsilon_{\max} (1-\gamma) \sum_{j=0}^{\infty} \gamma^j \kappa(K-i+j; \mu, \nu) \end{aligned} \quad (40)$$

Using eq. (34), eq. (35) and eq. (40)

$$\begin{aligned} \|Q^* - Q^{\pi_K}\|_{1,\mu} &\leq \frac{2\gamma(1-\gamma^{K+1})}{1-\gamma} \left[ \sum_{i=0}^{K-1} \sum_{j=0}^{\infty} \alpha_i \gamma^j \kappa(K-i+j; \mu, \nu) \right] \varepsilon_{\max} \\ &\quad + \frac{4\gamma(1-\gamma^{K+1})}{(1-\gamma)^3} \alpha_K R_{\max} \end{aligned} \quad (41)$$

Focusing on the first term on RHS of eq. (41), if we then we can take the norm out of the sum as a constant. We are left with

$$\begin{aligned} &\sum_{i=0}^{K-1} \sum_{j=0}^{\infty} \alpha_i \gamma^j \kappa(K-i+j; \mu, \nu) \\ &= \sum_{i=0}^{K-1} \sum_{j=0}^{\infty} \frac{(1-\gamma)\gamma^{K-i+j-1}}{1-\gamma^{K+1}} \kappa(K-i+j; \mu, \nu) \\ &= \frac{1-\gamma}{1-\gamma^{K+1}} \sum_{j=0}^{\infty} \sum_{i=0}^{K-1} \gamma^{K-i+j-1} \kappa(K-i+j; \mu, \nu) \\ &\leq \frac{1-\gamma}{1-\gamma^{K+1}} \sum_{m=0}^{\infty} \gamma^{m-1} \cdot m \cdot \kappa(m; \mu, \nu) \\ &\leq \frac{1}{1-\gamma^{K+1}(1-\gamma)} \phi_{\mu, \nu} \end{aligned} \quad (42)$$

Where the last inequality is due to assumption 5. Combining eq. (41) and eq. (42) we arrive at

$$\|Q^* - Q^{\pi_K}\|_{1,\mu} \leq \frac{2\gamma \cdot \phi_{\mu, \nu}}{(1-\gamma)^2} \cdot \varepsilon_{\max} + \frac{4\gamma^{K+1}}{(1-\gamma)^2} \cdot R_{\max} \quad (43)$$

□

**Lemma 7** (Rotation invariance). Let  $(X_i)_{i=1}^n$  be independent, centered and sub-gaussian. Then  $\sum_{i=1}^n X_i$  is centered and sub-gaussian with

$$\left\| \sum_{i=1}^n X_i \right\|_{\psi_2}^2 \leq C \sum_{i=1}^n \|X_i\|_{\psi_2}^2$$

*Proof.* See [Vershynin 2010, p. 12].

□

**Definition 18** (Sub-exponential norm). For a random variable define

$$\|X\|_{\psi_1} = \sup_{p \geq 1} p^{-1} \|X\|_p$$

called the sub-exponential norm, said to 'exist' if finite. In that case  $X$  is said to be 'sub-exponential'.

**Lemma 8** (Sub-gaussian squared is sub-exponential). A random variable  $X$  is sub-gaussian if and only if  $X^2$  is sub-exponential and

$$\|X\|_{\psi_2}^2 \leq \left\| X^2 \right\|_{\psi_1} \leq 2\|X\|_{\psi_2}^2$$

*Proof.* See [Vershynin 2010, p. 14]

□

**Proposition 35.** Let  $v$  be a random vector in  $\mathbb{R}^n$  then

$$\mathbb{E}\|v\|_1 \leq \sqrt{n}\sqrt{\mathbb{E}\|v\|_2^2}$$

*Proof.* Denote  $v$ 's coordinates  $v = (v_1, \dots, v_n)$ . Cauchy-Schwarz applied to some vector  $w$  and  $(1, \dots, 1)$  yields

$$\|w\|_1 \leq \sqrt{n}\|w\|_2$$

Now let  $w = (\mathbb{E}v_1, \dots, \mathbb{E}v_n)$ . Then by linearity of expectation and Jensens inequality

$$\mathbb{E}\|v\| = \|w\| \leq \sqrt{n}\sqrt{\sum_{i=1}^n (\mathbb{E}v_i)^2} \leq \sqrt{n}\sqrt{\mathbb{E}\sum_{i=1}^n v_i^2} = \sqrt{n}\sqrt{\mathbb{E}\|v\|_2^2}$$

□

**Theorem 11** (Bernstein's inequality). Suppose  $U_1, \dots, U_n$  are independent with  $\mathbb{E}U_i = 0$ ,  $|U_i| \leq M$  for all  $i \in [n]$ . Then for all  $t > 0$

$$\mathbb{P}\left(\left|\sum_{i=1}^n U_i\right| \geq t\right) \leq \exp\left(\frac{-t^2}{2/3Mt + 2\sigma^2}\right)$$

where  $\sigma^2 = \sum_{i=1}^n V(U_i)$ .

*Proof of theorem 10.* First some introductory fixing of notation and variables. Fix a minimal  $\delta$ -covering of  $\mathcal{F}$  with centers  $f_1, \dots, f_{N_\delta}$ . Define

$$\tilde{Q} := \operatorname{argmin}_{f \in \mathcal{F}} \|f - TQ\|_\nu^2$$

$$k^* := \operatorname{argmin}_{k \in [N_\delta]} \|f_k - \hat{Q}\|_\infty$$

and  $X_i := (S_i, A_i)$ . Notice that  $\tilde{Q}$  differs from  $\hat{Q}$  in that  $\tilde{Q}$  approximates  $TQ$  w.r.t.  $\|\cdot\|_\nu^2$  while  $\hat{Q}$  approximates  $Y = (Y_1, \dots, Y_n)$  in mean squared error over  $X = (X_1, \dots, X_n)$ . We shall be loose about applying functions to vectors (of random variables) in the sense that they are applied entry-wise. We use  $\|\cdot\|_p$  to denote the (finite dimensional)  $p$ -norm ( $p$  omitted when  $p = 2$ ). When talking about  $p$ -norms on the random variables we always specify the distribution (e.g.  $\|\cdot\|_\nu$ ). When the sample (e.g.  $X$ ) is clear from context we omit it writing  $\|f\| = \|f(X)\|$ .

**Step 1** By definition (of  $\hat{Q}$ ) for all  $f \in \mathcal{F}$  we have  $\|\hat{Q}(X) - Y\|^2 \leq \|f(X) - Y\|^2$ , leading to

$$\|Y\|^2 + \|\hat{Q}\|^2 - 2Y \cdot \hat{Q} \leq \|Y\|^2 + \|f\|^2 - 2Y \cdot f \quad (44)$$

$$\iff \|\hat{Q}\|^2 + \|TQ\|^2 - 2\hat{Q} \cdot TQ \leq \|f\|^2 + \|TQ\|^2 - 2f \cdot TQ + 2Y \cdot \hat{Q} - 2Y \cdot f - 2\hat{Q} \cdot TQ + 2f \cdot TQ \quad (45)$$

$$\iff \|\hat{Q} - TQ\|^2 \leq \|f - TQ\|^2 + 2(Y - TQ) \cdot (\hat{Q} - f) \quad (46)$$

$$\iff \|\hat{Q} - TQ\|^2 \leq \|f - TQ\|^2 + 2\xi \cdot (\hat{Q} - f) \quad (47)$$

Where  $\xi_i := Y_i - TQ(X_i)$  and  $\xi := (\xi_1, \dots, \xi_n)$ . Let  $\Sigma = (X_1, \dots, X_n)^{-1}(\mathbb{B}_n) \in \mathcal{H}$  be the  $\sigma$ -algebra generated by the samples. Now we proof a minor lemma

**Proposition 36.**  $\mathbb{E}(\xi_i \mid \Sigma) = 0$  and thus  $\mathbb{E}(\xi_i g(X_i)) = 0$  for any function  $g : \mathbb{R} \rightarrow \mathbb{R}$ .

*Proof.* Recall that  $X_i = (S_i, A_i)$ ,

$$Y_i = R_i + \gamma \max_{a \in \mathcal{A}} Q(S_{i+1}, a)$$

where  $S_{i+1} \sim P(X_i)$ ,  $R_i \sim R(X_i)$  and

$$TQ(X_i) = \mathbb{E}_\Sigma R'_i + \gamma \mathbb{E}_\Sigma Q(S', \operatorname{argmax}_{a \in \mathcal{A}} Q(S', a))$$

where  $S' \sim P(X_i)$ ,  $R'_i \sim R(X_i)$ . Since  $S'$  and  $S_{i+1}$  are i.i.d.

$$\begin{aligned} \mathbb{E}_\Sigma \xi_i &= \mathbb{E}_\Sigma (Y_i - TQ(X_i)) \\ &= \mathbb{E}_\Sigma R_i - \mathbb{E}_\Sigma R'_i + \gamma \left( \mathbb{E}_\Sigma \left( \max_{a \in \mathcal{A}} Q(S_{i+1}, a) \right) - \mathbb{E}_\Sigma \operatorname{argmax}_{a \in \mathcal{A}} (Q(S', a)) \right) \\ &= 0 \end{aligned}$$

Therefore  $\mathbb{E}(\xi_i \mid \Sigma) = 0$ . □

By this lemma we can deduce

$$\mathbb{E}(\xi \cdot (\hat{Q} - f)) = \mathbb{E}(\xi \cdot (\hat{Q} - TQ)) \quad (48)$$

To bound this we insert  $f_{k*}$  by the triangle inequality

$$\left| \mathbb{E}(\xi \cdot (\hat{Q} - TQ)) \right| \leq \left| \mathbb{E}(\xi \cdot (\hat{Q} - f_{k*})) \right| + \left| \mathbb{E}(\xi \cdot (f_{k*} - TQ)) \right| \quad (49)$$

We now bound these two terms. The first by Cauchy-Schwarz

$$\left| \mathbb{E} \xi \cdot (\hat{Q} - f_{k*}) \right| \leq \mathbb{E} \left( \|\xi\| \|\hat{Q} - f_{k*}\| \right) \leq \mathbb{E}(\|\xi\|) \sqrt{n} \delta \leq 2n V_{\max} \delta \quad (50)$$

where we have used that  $\|\hat{Q} - f_{k*}\|_\infty \leq \delta$  so

$$\|\hat{Q} - f_{k*}\|^2 = \sum_{i=1}^n (\hat{Q}(X_i) - f_{k*}(X_i))^2 \leq \sum_{i=1}^n \delta^2 = n\delta^2 \quad (51)$$

and that  $|Y_i|, TQ(X_i) \leq V_{\max}$  so

$$\|\xi\|^2 = \sum_{i=1}^n (Y_i - TQ(X_i))^2 \leq \sum_{i=1}^n (2V_{\max})^2 = 4V_{\max}^2 n \quad (52)$$

To bound the second term in eq. (49) define

$$Z_j := \xi \cdot (f_j - TQ) \|f_j - TQ\|^{-1} \quad (53)$$

Note that since  $\xi_i$  are centered  $Z_j$ . For a sub- $\sigma$ -algebra  $\Sigma$  define the *sub-gaussian* norm by

**Definition 19** (Sub-gaussian norm).

$$\|W\|_{\psi_2, \Sigma} := \sup_{p \geq 1} p^{-1/2} (\mathbb{E}_\Sigma |W|^p)^{1/p}$$

Because of proposition 36  $\xi_i(f_j(X_i) - TQ(X_i))$  is centered for any  $i \in [n]$  and

$$\|\xi_i(f_j(X_i) - TQ(X_i))\|_{\psi_2, \Sigma} \leq 2V_{\max} |f_j(X_i) - TQ(X_i)| \quad (54)$$

Therefore by lemma 7

$$\|Z_j\|_{\psi_2, \Sigma}^2 \leq \|f_j - TQ\|^{-2} \left\| \sum_{i=1}^n \xi_i(f_j(X_i) - TQ(X_i)) \right\|_{\psi_2, \Sigma}^2 \quad (55)$$

$$\leq \|f_j - TQ\|^{-2} C_1 \sum_{i=1}^n \|\xi_i(f_j(X_i) - TQ(X_i))\|_{\psi_2, \Sigma}^2 \quad (56)$$

$$\leq \|f_j - TQ\|^{-2} C_1 \sum_{i=1}^n 4V_{\max} |f_j(X_i) - TQ(X_i)|^2 \quad (57)$$

$$= 4V_{\max}^2 C_1 \quad (58)$$

Observe that  $\|X\|_p \leq \sqrt{p} \sup_{p \geq 1} \|X\|_{\psi_2}$ . Thus by [Vershynin 2010, p. 11 and Lemma 5.5]

$$\mathbb{E} \exp \left( cZ_j^2 / \|Z_j\|_{\psi_2}^2 \right) \leq e \quad (59)$$

so

$$\mathbb{E} \max_{j \in N_\delta} Z_j^2 = \frac{\max_{j \in [N_\delta]} \|Z_j\|_{\psi_2}^2}{c} \mathbb{E} \left( \max_{j \in [N_\delta]} \frac{cZ_j^2}{\max_{k \in [N_\delta]} \|Z_k\|_{\psi_2}^2} \right) \quad (60)$$

$$\leq \frac{4V_{\max}^2 C_1}{c} \mathbb{E} \left( \max_{j \in N_\delta} \frac{cZ_j^2}{\|Z_j\|_{\psi_2}^2} \right) \quad (61)$$

$$\leq \frac{4V_{\max}^2 C_1}{c} \log \left( \mathbb{E} \max_{j \in N_\delta} \exp \left( \frac{cZ_j^2}{\|Z_j\|_{\psi_2}^2} \right) \right) \quad (62)$$

$$\leq \frac{4V_{\max}^2 C_1}{c} \log \left( \sum_{j \in [N_\delta]} \mathbb{E} \exp \left( \frac{cZ_j^2}{\|Z_j\|_{\psi_2}^2} \right) \right) \quad (63)$$

$$\leq \frac{4V_{\max}^2 C_1}{c} \log (eN_\delta) \quad (64)$$

$$\leq C_2^2 V_{\max}^2 \log(N_\delta) \quad (65)$$

Where  $C_2 := \sqrt{8C_1/c}$ . Now we can bound

$$\mathbb{E} (\xi \cdot (f_{k*} - TQ)) = \mathbb{E} (\|f_{k*} - TQ\| |Z_{k*}|) \quad (66)$$

$$\leq \mathbb{E} \left( \left( \|\hat{Q} - TQ\| + \|\hat{Q} - f_{k*}\| \right) |Z_{k*}| \right) \quad (67)$$

$$\leq \mathbb{E} \left( \left( \|\hat{Q} - TQ\| + n\delta \right) |Z_{k*}| \right) \quad (68)$$

$$\leq \left( \mathbb{E} \left( \|\hat{Q} - TQ\| + n\delta \right)^2 \right)^{1/2} \left( \mathbb{E} Z_{k*}^2 \right)^{1/2} \quad (69)$$

$$\leq \mathbb{E} \left( \|\hat{Q} - TQ\| + n\delta \right) \left( \mathbb{E} Z_{k*}^2 \right)^{1/2} \quad (70)$$

$$\leq \left( \sqrt{\mathbb{E} \|\hat{Q} - TQ\|_2^2} + n\delta \right) \left( \mathbb{E} Z_{k*}^2 \right)^{1/2} \quad (71)$$

$$\leq \left( \sqrt{\mathbb{E} \|\hat{Q} - TQ\|_2^2} + n\delta \right) C_2^2 V_{\max}^2 \log(N_\delta) \quad (72)$$

Where eq. (66) to eq. (67) is by the triangle inequality and eq. (70) to eq. (71) is proposition 35. Combining eq. (47), eq. (49), eq. (50) and eq. (72)

$$\mathbb{E}\|\hat{Q} - TQ\|^2 \leq \mathbb{E}\|f - TQ\|^2 + 4nV_{\max}\delta + \left(\sqrt{\mathbb{E}\|\hat{Q} - TQ\|^2} + \sqrt{n\delta}\right) C_2 V_{\max} \sqrt{\log(N_\delta)} \quad (73)$$

$$= C_2 V_{\max} \sqrt{n \log(N_\delta)} \sqrt{\mathbb{E}\|\hat{Q} - TQ\|^2} + nC_2^2 \delta V_{\max}^2 \log(N_\delta) + \mathbb{E}\|f - TQ\|^2 \quad (74)$$

**Lemma 9.** Let  $a, b > 0, \kappa \in (0, 1]$  then

$$a^2 \leq 2ab + c \implies a^2 \leq (1 + \kappa)^2 b^2 / \kappa + (1 + \kappa)c$$

*Proof.*  $0 \leq (x - y)^2 = x^2 + y^2 - 2xy \implies 2xy \leq x^2 + y^2$  for any  $x, y \in \mathbb{R}$  so

$$\begin{aligned} 2ab &= 2\sqrt{\frac{\kappa}{1 + \kappa}} a \sqrt{\frac{1 + \kappa}{\kappa}} b \\ &\leq \frac{\kappa}{1 + \kappa} a^2 + \frac{1 + \kappa}{\kappa} b^2 \end{aligned}$$

□

By lemma 9 applied to eq. (74)

$$\frac{1}{n} \mathbb{E}\|\hat{Q} - TQ\|^2 \leq \frac{(1 + \kappa)^2}{\kappa} \frac{1}{n} C_2^2 V_{\max}^2 \log(N_\delta) + (1 + \kappa) \left( \delta C_2^2 V_{\max}^2 \log(N_\delta) + \frac{1}{n} \mathbb{E}\|f - TQ\|^2 \right) \quad (75)$$

We now take a closer look at the last term. Since  $f$  and  $TQ$  doesn't depend on the  $X_i$ 's we have

$$\begin{aligned} \frac{1}{n} \mathbb{E}\|f - TQ\|^2 &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}(f(X_i) - TQ(X_i))^2 \\ &= \mathbb{E}(f(X_i) - TQ(X_i))^2 \\ &= \|f - TQ\|_\nu^2 \end{aligned}$$

Now since eq. (75) holds for any  $f \in \mathcal{F}$  we can further say

$$\begin{aligned} \frac{1}{n} \mathbb{E}\|\hat{Q} - TQ\|^2 &\leq \frac{(1 + \kappa)^2}{\kappa} \delta C_2^2 V_{\max}^2 \log(N_\delta) \\ &\quad + (1 + \kappa) \left( C_2^2 V_{\max}^2 \log(N_\delta) + \sup_{g \in \mathcal{G}} \inf_{f \in \mathcal{F}} \|f - g\|_\nu^2 \right) \\ &\leq \frac{(1 + \kappa)^2}{\kappa} \frac{1}{n} C_2^2 V_{\max}^2 \log(N_\delta) + (1 + \kappa) \left( C_2^2 V_{\max}^2 \log(N_\delta) + \omega(\mathcal{F}) \right) \quad (76) \end{aligned}$$

where we take the supremum over  $\mathcal{G}$  (recall  $TQ \in \mathcal{G}$ ).

**Step 2** Here we link up  $\|\hat{Q} - TQ\|_\sigma^2$  with  $\mathbb{E} \frac{1}{n} \|\hat{Q} - TQ\|^2$ . First note that

$$\left| \left( \hat{Q}(x) - TQ(x) \right)^2 - \left( f_{k*}(x) - TQ(x) \right)^2 \right| = \left| \hat{Q}(x) - f_{k*}(x) \right| \cdot \left| \hat{Q}(x) + f_{k*}(x) - 2TQ(x) \right| \quad (77)$$

$$\leq 4V_{\max}\delta \quad (78)$$

Using this twice we can say

$$(\hat{Q}(\hat{X}_i) - TQ(\hat{X}_i))^2 \quad (79)$$

$$\leq (\hat{Q}(\tilde{X}_i) - TQ(\tilde{X}_i))^2 - (f_{k*}(\tilde{X}_i) - TQ(\tilde{X}_i))^2 + (f_{k*}(\tilde{X}_i) - TQ(\tilde{X}_i))^2 \quad (80)$$

$$\leq (f_{k*}(\tilde{X}_i) - TQ(\tilde{X}_i))^2 + (\hat{Q}(X_i) - TQ(X_i))^2 - (\hat{Q}(X_i) - TQ(X_i))^2 + (f_{k*}(X_i) - TQ(X_i))^2 - (f_{k*}(X_i) - TQ(X_i))^2 + 4V_{\max}\delta \quad (81)$$

$$\leq (\hat{Q}(X_i) - TQ(X_i))^2 + (f_{k*}(\tilde{X}_i) - TQ(\tilde{X}_i))^2 - (f_{k*}(X_i) - TQ(X_i))^2 + 8V_{\max}\delta \quad (82)$$

Thus we get

$$\left\| \hat{Q} - TQ \right\|_{\sigma}^2 \quad (83)$$

$$= \mathbb{E} \frac{1}{n} \sum_{i=1}^n (\hat{Q}(\tilde{X}_i) - TQ(\tilde{X}_i))^2 \quad (84)$$

$$\leq \mathbb{E} \frac{1}{n} \sum_{i=1}^n \left( (\hat{Q}(X_i) - TQ(X_i))^2 + (f_{k*}(\tilde{X}_i) - TQ(\tilde{X}_i))^2 - (f_{k*}(X_i) - TQ(X_i))^2 \right) + 8V_{\max}\delta \quad (85)$$

$$= \frac{1}{n} \left\| \hat{Q} - TQ \right\|^2 + \frac{1}{n} \sum_{i=1}^n h_{k*}(X_i, \tilde{X}_i) + 8V_{\max}\delta \quad (86)$$

Where we define

$$h_j(x, y) := (f_j(y) - TQ(y))^2 - (f_j(x) - TQ(x))^2 \quad (87)$$

For any  $j \in [N_{\delta}]$ . Define  $\Upsilon = 2V_{\max}$  and

$$T := \max_{j \in [N_{\delta}]} \left| \sum_{i=1}^n h_j(X_i, \tilde{X}_i) / \Upsilon \right| \quad (88)$$

Then we can bound the middle term in eq. (86)

$$\mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n h_{k*}(X_i, \tilde{X}_i) \right) \leq \Upsilon / n \mathbb{E} \max_{j \in [N_{\delta}]} \left( \left| \sum_{i=1}^n h_j(X_i, \tilde{X}_i) / \Upsilon \right| \right) \quad (89)$$

$$\leq \Upsilon / n \mathbb{E} T \quad (90)$$

We want to use Bernsteins inequality (theorem 11) with  $U_i = h_j(X_i, \tilde{X}_i)$ . Therefore notice that  $|h_j| \leq \Upsilon^2$  and

$$\text{Var} h_j(X_i, \tilde{X}_i) = 2 \text{Var} (f_j(X_i) - TQ(X_i))^2 \quad (91)$$

$$\leq 2 \mathbb{E} (f_j(X_i) - TQ(X_i))^4 \quad (92)$$

$$\leq 2\Upsilon^4 \quad (93)$$



so by union bounding for any  $u < 6n\Upsilon$  we have

$$\mathbb{E}T = \int_0^\infty \mathbb{P}(T \geq t) \quad (94)$$

$$\leq u + \int_u^\infty \mathbb{P}(T \geq t) dt \quad (95)$$

$$\leq u + \int_u^\infty 2N_\delta \exp\left(\frac{-t^2}{2\Upsilon t/3 + 4n\Upsilon^2}\right) dt \quad (96)$$

$$\leq u + 2N_\delta \int_u^\infty \exp\left(\frac{-t^2}{2\Upsilon^2(t/(3\Upsilon) + 2n)}\right) dt \quad (97)$$

$$\leq u + 2N_\delta \left( \int_u^{6n\Upsilon} \exp\left(\frac{-t^2}{8n\Upsilon^2}\right) dt + \int_{6n\Upsilon}^\infty \exp\left(\frac{-t}{4/3\Upsilon}\right) dt \right) \quad (98)$$

$$\leq u + 2N_\delta \left( \frac{8n\Upsilon}{2u} \exp\left(\frac{-u^2}{8n\Upsilon}\right) + \frac{4\Upsilon}{3} \exp\left(\frac{-24n\Upsilon}{3\Upsilon}\right) \right) \quad (99)$$

where we use lemma 10 from eq. (98) to eq. (99). Now set  $u = \Upsilon\sqrt{8n\log N_\delta}$  continuing from eq. (99) we have

$$\dots = \Upsilon\sqrt{8n\log N_\delta} + \frac{\Upsilon^2 8nN_\delta}{\Upsilon\sqrt{8n\log N_\delta}} \exp(-\log N_\delta) + 8/3N_\delta\Upsilon \exp(-9/2n) \quad (100)$$

$$= \Upsilon 2\sqrt{2n} \left( \log N_\delta + \frac{1}{\log N_\delta} \right) + 8/3N_\delta e^{-9/2n} \quad (101)$$

$$\leq 4\sqrt{2}\Upsilon\sqrt{n\log N_\delta} + 8/3\Upsilon \quad (102)$$

Inserting eq. (102) and eq. (76) into eq. (86)

$$\left\| \hat{Q} - TQ \right\|_\nu^2 \leq \frac{1}{n} \mathbb{E} \left\| \hat{Q} - TQ \right\|^2 + 8\sqrt{2}V_{\max}n^{-1/2}\sqrt{\log N_\delta} + 8V_{\max}(n^{-1} + \delta) \quad (103)$$

$$\begin{aligned} &\leq \frac{(1+\kappa)^2}{\kappa} \frac{1}{n} C_2^2 V_{\max}^2 \log(N_\delta) + (1+\kappa) \left( \delta C_2^2 V_{\max}^2 \log(N_\delta) + \omega(\mathcal{F}) \right) \\ &\quad + 8\sqrt{2}V_{\max}n^{-1/2}\sqrt{\log N_\delta} + 8V_{\max}(n^{-1} + \delta) \end{aligned} \quad (104)$$

□

## 5 Appendices

### 5.1 Lemmas for Fan et al.

**Lemma 10.** For  $x > 0$ .

$$\int_x^\infty e^{-t^2/2} dt \leq \frac{1}{x} e^{-x^2/2}$$

*Proof.* Observe that for  $t \geq x > 0$  we have  $1 \leq t/x$  so

$$\begin{aligned} \int_x^\infty e^{-t^2/2} dt &\leq \int_x^\infty \frac{t}{x} e^{-t^2/2} dt \\ &\leq \frac{1}{x} e^{-x^2/2} \end{aligned}$$

□

## 5.2 Other notes

**Definition 20** (Almost sure uniform convergence of random processes). A sequence of random processes  $X_n : \mathcal{X} \times \Omega \rightarrow \mathbb{R}$  is said to converge **almost surely uniformly** to  $X : \mathcal{X} \times \Omega \rightarrow \mathbb{R}$  if and only if

$$\mathbb{P}(\sup_{x \in \mathcal{X}} |X_n(x) - X(x)| \rightarrow 0) = 1$$

**Definition 21** (Uniform convergence in probability of random processes). A sequence of random processes  $X_n : \mathcal{X} \times \Omega \rightarrow \mathbb{R}$  is said to converge **uniformly in probability** to  $X : \mathcal{X} \times \Omega \rightarrow \mathbb{R}$  if and only if

$$\sup_{x \in \mathcal{X}} |X_n(x) - X(x)| \xrightarrow{P} 0$$