

A Theoretical Analysis of Fitted Q-Iteration

Jacob Harder
University of Copenhagen

March 12, 2020

1 Abstract

2 Foreword

The main purpose of this master thesis for me, has been to uncover what (at present) it is possible to say (mathematically) about the convergence of Q-learning algorithms. In particular Q-learning algorithms using (deep) ANNs.

I came to realize during my reading of [TODO ref to YangXieWang] that it is quite error-prone with some errors not obviously fixable.

3 Disambiguation

- $[\phi] = 1$ when ϕ is true/holds and 0 otherwise, for a logical formula ϕ .
- $[q] = \{1, \dots, q\}$ for $q \in \mathbb{N}$.
- $C_{\mathbb{K}}(X) = \{f : X \rightarrow \mathbb{K} \mid f \text{ continuous}\}$, $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$. $C(X) = C_{\mathbb{R}}(X)$
- ANN abrv. artificial neural network see definition 2.
- δ_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 σ -algebra.

3.1 Notational deviations from [TODO ref YangXieWang]

Because σ is used ambiguously in theorem 1 we denote the probability distribution ' σ ' from [YangXieWang, thm. 6.2, p. 20] by ν instead.

I dislike the shorthand defined in [YangXieWang, p. 26 bottom]: $\|f\|_n^2 = 1/n \cdot \sum_{i=1}^n f(X_i)^2$. This is partially due to inconsistencies and abuse of this notation employed. For example $\|f\|_n$ is used as $1/n \sum_{i=1}^n f(X_i)$ as opposed another likely interpretation $\sqrt{\|f\|_n^2}$, whereas $\|f\|_n^{-1}$ is used to mean $1/(\|f\|_n)$. This is avoided by using finite dimensional p -norms instead. The conversion to my notation thus becomes $\|f\|_n \rightsquigarrow \|f\|_1 / n$, $\|f\|_n^2 \rightsquigarrow \|f\|^2 / n$, $\|f\|_n^{-1} \rightsquigarrow n\|f\|_1^{-1}$.

4 Introduction

4.1 Reinforcement Learning

In Reinforcement Learning (RL) we are concerned with finding an optimal policy for an agent in some environment. Typically (also in the case of Q-learning) this environment is a Markov decision process

Definition 1. A Markov decision process (MDP) $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ consists of

- \mathcal{S} a set of states
- \mathcal{A} a set of actions
- $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$ its Markov transition kernel
- $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathbb{R})$ its immediate reward distribution
- $\gamma \in (0, 1)$ the discount factor

A policy (for an MDP) is a function

$$\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$$

With this we can define the state-value function $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$

$$V^\pi(s) = \mathbb{E} \left(\sum_{t \geq 0} \gamma^t R_t \mid R_t \sim R(S_t, A_t), S_t \sim P(S_{t-1}, A_{t-1}), A_t \sim \pi(S_t), S_0 = s \right)$$

And the state-action-value (Q-) function $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$

$$Q^\pi(s, a) = \mathbb{E}(R(s, a) + \gamma V^\pi(S_0) \mid S_0 \sim P(s, a))$$

The optimal Q-function is defined as

$$Q^*(s, a) = \sup_{\pi} Q^\pi(s, a)$$

One can show that there is a policy π^* such that $Q^* = Q^{\pi^*}$. This is the optimal policy - the goal of RL.

Note that V^π , Q^π and Q^* are usually infeasible to calculate to machine precision, unless $\mathcal{S} \times \mathcal{A}$ is finite and not very big.

4.2 Q-Learning

Let $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ be a policy. We define the operator

$$(P^\pi Q)(s, a) = \mathbb{E}(Q(S', A') \mid S' \sim P(s, a), A' \sim \pi(S'))$$

Intuitively this operator yields the expected state-action-value function when looking *one step ahead* following the policy π and taking expectation of Q .

We define the operator T^π called the Bellman operator by

$$(T^\pi Q)(s, a) = \mathbb{E}R(s, a) + \gamma(P^\pi Q)(s, a)$$

This operator adjust the Q function to look more like Q^π making one "iteration" of "propagation of rewards" discounting with γ . Indeed it is easily seen that Q^π is a fixed point for T^π .

A *greedy* policy π with respect to a state-action value function Q is a policy which deterministically chooses an action with maximal value of Q for each state. That is $\pi(s) = \delta_a$ for some $a \in \operatorname{argmax}_a Q(s, a)$. We then write $\pi = \pi_Q$. With this we can define the operator T :

$$TQ = T^{\pi_Q} Q$$

called the Bellman *optimality* operator.

The Bellman optimality *equation* can then be written $Q^* = TQ^*$.

Proposition 1. Q^π is the unique fixed point of T^π .

Proof. Clearly $T^\pi Q^\pi = Q^\pi$. [TODO: rest of this proof]

□

4.3 Artificial Neural Networks

Definition 2. 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.

4.4 Fitted Q-Iteration

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

Algorithm 1: Fitted Q-Iteration Algorithm

Input: MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$, function class \mathcal{F} , sampling distribution ν , 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 ν 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 π_K as the greedy policy w.r.t. \tilde{Q}_K

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

5 Measure Theory

We are mostly concerned with a random process

$$(Z_i)_{i=1}^K = (S_i, A_i, R_i)_{i=1}^K \in (\mathcal{S} \times \mathcal{A} \times (0, R_{\max}))^K \quad (1)$$

where $\mathcal{S} \subseteq \mathbb{R}^d$ is compact and \mathcal{A} is finite, so we can model this as a discrete (and finite) time random process in a compact subset of \mathbb{R}^{d+1} having the Markov property, namely that

$$\mathbb{P}(Z_j \in A \mid Z_{j-1}, \dots, Z_1) = \mathbb{P}(Z_j \in A \mid Z_{j-1}) \quad (2)$$

These random variables live on some background probability space, denote this $(\Omega, \mathcal{H}, \mathbb{P})$.

6 Assumptions

6.1 Assumption 1: Holder Smoothness

Definition 3. 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_ℓ, v_ℓ) satisfying

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

Here $\widetilde{W}_\ell = (W_\ell, v_\ell)$.

The set of them we denote $\mathcal{F}(s, V)$.

Definition 4. Let $\mathcal{D} \subseteq \mathbb{R}^r$ be compact and $\beta, H > 0$. A function $f : \mathcal{D} \rightarrow \mathbb{R}$ we call Holder smooth if

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

Where $\alpha = (\alpha_1, \dots, \alpha_r) \in \mathbb{N}^r$. We write $f \in C_r(\mathcal{D}, \beta, H)$.

Definition 5. 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 Holder 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 $g_{jk} \in C_{t_j}([a_j, b_j]^{t_j}, \beta_j, H_j)$, i.e. they are Holder smooth. We denote the class of these functions

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

Definition 6. 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

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

Assumption 1. It is assumed that $Tf \in \mathcal{G}_0$ for any $f \in \mathcal{F}_0$.

I.e. when using the Bellman optimality operator on our sparse ReLU networks, we should stay in the class of compositions of Holder smooth functions.

6.2 Assumption 2: Concentration Coefficients

Definition 7 (Concentration coefficients). Let $\nu_1, \nu_2 \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$ be probability measures, absolutely continuous w.r.t. m_λ . Define

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

Assumption 2. Let ν be the sampling distribution from the algorithm, and μ the distribution over which we measure the error in the main theorem, then we assume

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

7 Main theorem

Theorem 1 (Yang, Xie, Wang). For any $K \in \mathbb{N}$ let Q^{π_K} be the action-value function corresponding to policy π_K which is returned by Algorithm 1, when run with a sparse ReLU network on the form

$$\mathcal{F}_0 = \{f(\cdot, a) \in \mathcal{F}(L^*, \{d_j^*\}_{j=0}^{L^*+1}, s^*) \mid a \in \mathcal{A}\}$$

where

$$L^* \lesssim (\log n)^{\xi'}, d_0 = r, d_j^* = 1, \lesssim n^{\xi'}, s^* \asymp n^{\alpha^*} \cdot (\log n)^{\xi'}$$

Let μ be any distribution over $\mathcal{S} \times \mathcal{A}$. Under assumption 1 and assumption 2

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

Here $C, \xi', \xi^*, \phi_{\mu, \nu} \in \mathbb{R}_+$ and $\alpha^* \in (0, 1)$ are constants depending on the assumptions and R_{\max} the maximum possible reward.

8 Proofs

Proof of main theorem. Using theorem 2 we get

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

where $\varepsilon_{\max} = \max_{k \in [K]} \|T\tilde{Q}_{k-1} - \tilde{Q}_k\|_{2,\nu}$. Using theorem 3 with $Q = \tilde{Q}_{k-1}$, $\mathcal{F} = \mathcal{F}_0$, $\epsilon = 1$ and $\delta = 1/n$, we get

$$\varepsilon_{\max} \leq 4n^{-1}V_{\max}^2 + 12\delta V_{\max} + 2\omega(\mathcal{F}_0) + 8\sqrt{2}V_{\max}n^{-1/2}\sqrt{\log N_0} + 16V_{\max}n^{-1} \quad (4)$$

where $N_0 = |\mathcal{N}(1/n, \mathcal{F}_0, \|\cdot\|_{\infty})|$. To bound $\omega(\mathcal{F}_0)$ we want to employ the following lemma by

Lemma 1 (Approximation of Hölder Smooth Functions). Let $m, M \in \mathbb{N}$ with $N \geq \max\{(\beta + 1)^r, (H + 1)\}$, $L = 8 + (m + 5)(1 + \lceil \log_2 r \rceil)$, $d_0 = r, d_j = 12rN, d_{L+1} = 1$. Then for any $g \in \mathcal{C}_r([0, 1]^r, \beta, H)$ there exists a ReLU network $f \in \mathcal{F}(L, \{d_j\}_{j=0}^{L+1}, s, V_{\max})$ with $s \leq 94r^2(\beta + 1)^{2r}N(m + 6)(1 + \lceil \log_2 r \rceil)$ such that

$$\|f - g\|_{\infty} \leq (2H + 1)3^{r+1}N2^{-m} + H2^{\beta}N^{-\beta/r}$$

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 $f(\cdot, a) \in \mathcal{G}(\{p_j, t_j, \beta_j, H_j\})$ so let $f(\cdot, a) = g_q \circ \dots \circ g_1$ where $g_{jk} \in C_{t_j}([a_j, b_j]^{t_j}, \beta_j, H_j)$ for each $j \in [q]$ and $k \in [p_{j+1}]$. Define

$$h_1 = g_1 / (2H_1) + 1/2 \quad (5)$$

$$h_j(u) = g_j(2H_{j-1}u - H_{j-1}) / (2H_j) + 1/2, \quad j \in \{2, \dots, q-1\} \quad (6)$$

$$h_q(u) = g_q(2H_{q-1}u - H_{q-1}) \quad (7)$$

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

$$h_{1k} \in C_{t_1}([0, 1]^{t_1}, \beta_1, 1) \quad (8)$$

$$h_{jk} \in C_{t_j}([0, 1]^{t_j}, \beta_j, (2H_{j-1})^{\beta_j}) \quad (9)$$

$$h_q \in C_{t_q}([0, 1]^{t_q}, \beta_q, H_q(2H_{q-1})^{\beta_q}) \quad (10)$$

□

Theorem 2 (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

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

Lemma 2. $TQ \geq T^{\pi}Q$ 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 3. Let $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ be an action-value function, τ_1, \dots, τ_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 κ is the concentration coefficients defined in definition 7.

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

$$\mathbb{E}_\mu[(P^{\tau_m} \dots P^{\tau_1})(f)] = \int (P^{\tau_m} \dots P^{\tau_1})(f) d\mu \quad (11)$$

$$= \int f d(P^{\tau_m} \dots P^{\tau_1} \mu) \quad (12)$$

$$= \int f \frac{d(P^{\tau_m} \dots P^{\tau_1} \mu)}{d\nu} d\nu \quad (13)$$

$$\leq \left\| \frac{d(P^{\tau_m} \dots P^{\tau_1} \mu)}{d\nu} \right\|_{2, \nu} \cdot \|f\|_{2, \nu} \quad (14)$$

$$\leq \kappa(m, \mu, \nu) \|f\|_{2, \nu} \quad (15)$$

Where eq. (13) is due to the Radon-Nikodym theorem and eq. (14) is Cauchy-Schwarz. \square

Proof of theorem 2. First some things to keep in mind during the proof. Recall that $V_{\max} = R_{\max}/(1 - \gamma)$ and that π_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 π , P^π 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 π^* 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 2 we can make the upper bound

$$\begin{aligned}Q^* - Q^{\pi_K} &= T^{\pi^*} Q^* - T^{\pi_K} Q^{\pi_K} \\ &= T^{\pi^*} Q^* + (T^{\pi^*} \tilde{Q}_K - T^{\pi^*} \tilde{Q}_K) + (T\tilde{Q}_K - T\tilde{Q}_K) - T^{\pi_K} Q^{\pi_K} \\ &= (T^{\pi^*} \tilde{Q}_K - T\tilde{Q}_K) + (T^{\pi^*} Q^* - T^{\pi^*} \tilde{Q}_K) + (T\tilde{Q}_K - T^{\pi_K} Q^{\pi_K}) \\ &\leq (T^{\pi^*} Q^* - T^{\pi^*} \tilde{Q}_K) + (T\tilde{Q}_K - T^{\pi_K} Q^{\pi_K}) \\ &= (T^{\pi^*} Q^* - T^{\pi^*} \tilde{Q}_K) + (T^{\pi_K} \tilde{Q}_K - T^{\pi_K} Q^{\pi_K}) \\ &= \gamma P^{\pi^*} (Q^* - \tilde{Q}_K) + \gamma P^{\pi_K} (\tilde{Q}_K - Q^{\pi_K}) \\ &= \gamma (P^{\pi^*} - P^{\pi_K}) (Q^* - \tilde{Q}_K) + \gamma P^{\pi_K} (Q^* - Q^{\pi_K})\end{aligned} \quad (16)$$

This implies

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

Since γP^{π_K} is γ -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. (16) to obtain

$$Q^* - Q^{\pi_K} \leq \gamma U^{-1} (P^{\pi^*} (Q^* - \tilde{Q}_K) - P^{\pi_K} (Q^* - \tilde{Q}_K)) \quad (17)$$

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

$$\begin{aligned}
Q^* - \tilde{Q}_{i+1} &= Q^* + (T\tilde{Q}_i - T\tilde{Q}_i) - \tilde{Q}_{i+1} + (T^{\pi^*}\tilde{Q}_i - T^{\pi^*}\tilde{Q}_i) \\
&= (Q^* - T^{\pi^*}\tilde{Q}_i) + (T\tilde{Q}_i - \tilde{Q}_{i+1}) + (T^{\pi^*}\tilde{Q}_i - T\tilde{Q}_i) \\
&= (T^{\pi^*}Q^* - T^{\pi^*}\tilde{Q}_i) + \varrho_{i+1} + (T^{\pi^*}\tilde{Q}_i - T\tilde{Q}_i) \\
&\leq T^{\pi^*}Q^* - T^{\pi^*}\tilde{Q}_i + \varrho_{i+1} \\
&= \gamma P^{\pi^*}(Q^* - \tilde{Q}_i) + \varrho_{i+1}
\end{aligned} \tag{18}$$

and a lower bound

$$\begin{aligned}
Q^* - \tilde{Q}_{i+1} &= Q^* + (T\tilde{Q}_i - T\tilde{Q}_i) - \tilde{Q}_{i+1} + (T^{\pi_i}Q^* - T^{\pi_i}Q^*) \\
&= (T^{\pi_i}Q^* - T^{\pi_i}\tilde{Q}_i) + \varrho_{i+1} + (TQ^* - T^{\pi_i}Q^*) \\
&\geq T^{\pi_i}Q^* - T^{\pi_i}\tilde{Q}_i + \varrho_{i+1} \\
&= \gamma P^{\pi_i}(Q^* - \tilde{Q}_i) + \varrho_{i+1}
\end{aligned} \tag{19}$$

Applying eq. (18) and eq. (19) iteratively we get

$$Q^* - \tilde{Q}_K \leq \gamma^K (P^{\pi^*})^K (Q^* - \tilde{Q}_0) + \sum_{i=0}^{K-1} \gamma^{K-1-i} (P^{\pi^*})^{K-1-i} \varrho_{i+1} \tag{20}$$

and

$$Q^* - \tilde{Q}_K \geq \gamma^K (P^{\pi_{K-1}} \dots P^{\pi_0}) (Q^* - \tilde{Q}_0) + \sum_{i=0}^{K-1} \gamma^{K-1-i} (P^{\pi_{K-1}} \dots P^{\pi_{i+1}}) \varrho_{i+1} \tag{21}$$

Step 3 Combining eq. (20) and eq. (21) with eq. (17) we get

$$\begin{aligned}
Q^* - Q^{\pi_K} &\leq U^{-1} \left(\gamma^{K+1} ((P^{\pi^*})^{K+1} - P^{\pi_K} \dots P^{\pi_0}) (Q^* - \tilde{Q}_0) \right. \\
&\quad \left. + \sum_{i=0}^{K-1} \gamma^{K-i} ((P^*)^{K-i} - P^{\pi_K} \dots P^{\pi_{i+1}}) \varrho_{i+1} \right)
\end{aligned} \tag{22}$$

For shorthand define constants

$$\alpha_i = \frac{(1-\gamma)\gamma^{K-i-1}}{1-\gamma^{K+1}} \text{ for } 0 \leq i \leq K-1 \text{ and } \alpha_K = \frac{(1-\gamma)\gamma^K}{1-\gamma^{K+1}} \tag{23}$$

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

$$O_i = (1-\gamma)/2U^{-1}[(P^{\pi^*})^{K-i} + (P^{\pi_K} \dots P^{\pi_{i+1}})] \tag{24}$$

$$O_K = (1-\gamma)/2U^{-1}[(P^{\pi^*})^{K+1} + (P^{\pi_K} \dots P^{\pi_0})] \tag{25}$$

Then by eq. (22)

$$|Q^* - Q^{\pi_K}| \leq \frac{2\gamma(1-\gamma^{K+1})}{(1-\gamma)^2} \left[\sum_{i=0}^{K-1} \alpha_i O_i |\varrho_{i+1}| + \alpha_K O_K |Q^* - \tilde{Q}_0| \right] \tag{26}$$

So by linearity of expectation

$$\|Q^* - Q^{\pi_K}\|_{1,\mu} = \mathbb{E}_\mu |Q^* - Q^{\pi_K}| \tag{27}$$

$$\leq \frac{2\gamma(1-\gamma^{K+1})}{(1-\gamma)^2} \left[\sum_{i=0}^{K-1} \alpha_i \mathbb{E}_\mu (O_i |\varrho_{i+1}|) + \alpha_K \mathbb{E}_\mu (O_K |Q^* - \tilde{Q}_0|) \right] \tag{28}$$

With the bound on rewards we (crudely) estimate

$$\mathbb{E}_\mu O_K |Q^* - \tilde{Q}_0| \leq 2V_{\max} = 2R_{\max}/(1-\gamma) \tag{29}$$

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

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

$$\mathbb{E}_\mu(O_i|\varrho_{i+1}|) \tag{30}$$

$$= \frac{1-\gamma}{2} \mathbb{E}_\mu \left(U^{-1}[(P^{\pi_K})^{K-i} + P^{\pi_K} \dots P^{\pi_{i+1}}]|\varrho_{i+1}| \right) \tag{31}$$

$$= \frac{1-\gamma}{2} \mathbb{E}_\mu \left(\sum_{j=0}^{\infty} [(P^{\pi_K})^j (P^{\pi_K})^{K-i} + (P^{\pi_K})^{j+1} P^{\pi_{K-1}} \dots P^{\pi_{i+1}}]|\varrho_{i+1}| \right) \tag{32}$$

$$= \frac{1-\gamma}{2} \sum_{j=0}^{\infty} \mathbb{E}_\mu \left([(P^{\pi_K})^j (P^{\pi_K})^{K-i} + (P^{\pi_K})^{j+1} P^{\pi_{K-1}} \dots P^{\pi_{i+1}}]|\varrho_{i+1}| \right) \tag{33}$$

Notice that there are $K-i+j$ P -operators on both terms in the sum. Therefore we can employ lemma 3 twice. Moreover define $\varepsilon_{\max} = \max_{i \in [K]} \|\varrho_i\|_{2,\nu}$. Then

$$\begin{aligned} \mathbb{E}_\mu(O_i|\varrho_{i+1}|) &\leq (1-\gamma) \sum_{j=0}^{\infty} \gamma^j \kappa(K-i+j; \mu, \nu) \|\varrho_{i+1}\|_{2,\nu} \\ &\leq \varepsilon_{\max} (1-\gamma) \sum_{j=0}^{\infty} \gamma^j \kappa(K-i+j; \mu, \nu) \end{aligned} \tag{34}$$

Using eq. (28), eq. (29) and eq. (34)

$$\begin{aligned} \|Q^* - Q^{\pi_K}\|_{1,\mu} &\leq \frac{2\gamma(1-\gamma^{K+1})}{1-\gamma} \left[\sum_{i=0}^{K-1} \sum_{j=0}^{\infty} \alpha_i \gamma^j \kappa(K-i+j; \mu, \nu) \right] \varepsilon_{\max} \\ &\quad + \frac{4\gamma(1-\gamma^{K+1})}{(1-\gamma)^3} \alpha_K R_{\max} \end{aligned} \tag{35}$$

Focusing on the first term on RHS of eq. (35), if we then we can take the norm out of the sum as a constant. We are left with

$$\begin{aligned} &\sum_{i=0}^{K-1} \sum_{j=0}^{\infty} \alpha_i \gamma^j \kappa(K-i+j; \mu, \nu) \\ &= \sum_{i=0}^{K-1} \sum_{j=0}^{\infty} \frac{(1-\gamma)\gamma^{K-i+j-1}}{1-\gamma^{K+1}} \kappa(K-i+j; \mu, \nu) \\ &= \frac{1-\gamma}{1-\gamma^{K+1}} \sum_{j=0}^{\infty} \sum_{i=0}^{K-1} \gamma^{K-i+j-1} \kappa(K-i+j; \mu, \nu) \\ &\leq \frac{1-\gamma}{1-\gamma^{K+1}} \sum_{m=0}^{\infty} \gamma^{m-1} \cdot m \cdot \kappa(m; \mu, \nu) \\ &\leq \frac{1}{1-\gamma^{K+1}(1-\gamma)} \phi_{\mu,\nu} \end{aligned} \tag{36}$$

Where the last inequality is due to assumption 2. Combining eq. (35) and eq. (36) we arrive at

$$\|Q^* - Q^{\pi_K}\|_{1,\mu} \leq \frac{2\gamma \cdot \phi_{\mu,\nu}}{(1-\gamma)^2} \cdot \varepsilon_{\max} + \frac{4\gamma^{K+1}}{(1-\gamma)^2} \cdot R_{\max} \tag{37}$$

□

Theorem 3 (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 ν
- $(R_i, S'_i)_{i \in [n]}$ 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 δ -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

$$\begin{aligned} \left\| \hat{Q} - TQ \right\|_\nu^2 &\leq \frac{(1 + \kappa)^2}{\kappa} \frac{1}{n} V_{\max}^2 + (1 + \kappa) (6\delta V_{\max} + \omega(\mathcal{F})) \\ &\quad + 8\sqrt{2} V_{\max} n^{-1/2} \sqrt{\log N_\delta} + 8V_{\max} (n^{-1} + \delta) \end{aligned}$$

Where

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

where

Lemma 4 (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]. □

Proposition 2. 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 4 (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 $t > 0$

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

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

Proof of theorem 3. First some introductory fixing of notation and variables. Fix a minimal δ -covering of \mathcal{F} with centers f_1, \dots, f_{N_δ} . 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

$$\|Y\|^2 + \|\hat{Q}\|^2 - 2Y \cdot \hat{Q} \leq \|Y\|^2 + \|f\|^2 - 2Y \cdot f \quad (38)$$

$$\iff \|\hat{Q}\|^2 + \|TQ\|^2 - 2\hat{Q} \cdot TQ \leq \|f\|^2 + \|TQ\|^2 - 2f \cdot TQ + 2Y \cdot \hat{Q} - 2Y \cdot f - 2\hat{Q} \cdot TQ + 2f \cdot TQ \quad (39)$$

$$\iff \|\hat{Q} - TQ\|^2 \leq \|f - TQ\|^2 + 2(Y - TQ) \cdot (\hat{Q} - f) \quad (40)$$

$$\iff \|\hat{Q} - TQ\|^2 \leq \|f - TQ\|^2 + 2\xi \cdot (\hat{Q} - f) \quad (41)$$

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 σ -algebra generated by the samples. Now we proof a minor lemma

Proposition 3. $\mathbb{E}(\xi_i \mid \Sigma) = 0$ 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

$$\mathbb{E}(\xi \cdot (\hat{Q} - f)) = \mathbb{E}(\xi \cdot (\hat{Q} - TQ)) \quad (42)$$

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

$$\left| \mathbb{E}(\xi \cdot (\hat{Q} - TQ)) \right| \leq \left| \mathbb{E}(\xi \cdot (\hat{Q} - f_{k^*})) \right| + \left| \mathbb{E}(\xi \cdot (f_{k^*} - TQ)) \right| \quad (43)$$

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

$$\left| \mathbb{E}(\xi \cdot (\hat{Q} - f_{k^*})) \right| \leq \mathbb{E}(\|\xi\| \|\hat{Q} - f_{k^*}\|) \leq \mathbb{E}(\|\xi\|) \sqrt{n} \delta \leq 2n V_{\max} \delta \quad (44)$$

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

$$\|\widehat{Q} - f_{k^*}\|^2 = \sum_{i=1}^n (\widehat{Q}(X_i) - f_{k^*}(X_i))^2 \leq \sum_{i=1}^n \delta^2 = n\delta^2 \quad (45)$$

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

$$\|\xi\|^2 = \sum_{i=1}^n (Y_i - TQ(X_i))^2 \leq \sum_{i=1}^n (2V_{\max})^2 = 4V_{\max}^2 n \quad (46)$$

To bound the second term in eq. (43) define

$$Z_j := \xi \cdot (f_j - TQ) \|f_j - TQ\|^{-1} \quad (47)$$

Note that since ξ_i are centered Z_j . For a sub- σ -algebra Σ define the *sub-gaussian* norm by

Definition 8 (Subgaussian norm).

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

Because of proposition 3 $\xi_i(f_j(X_i) - TQ(X_i))$ is centered for any $i \in [n]$ and

$$\|\xi_i(f_j(X_i) - TQ(X_i))\|_{\psi_2, \Sigma} \leq 2V_{\max} |f_j(X_i) - TQ(X_i)| \quad (48)$$

Therefore by lemma 4

$$\|Z_i\|_{\psi_2, \Sigma}^2 \leq \|f_j - TQ\|^{-2} \left\| \sum_{i=1}^n \xi_i(f_j(X_i) - TQ(X_i)) \right\|_{\psi_2, \Sigma}^2 \quad (49)$$

$$\leq \|f_j - TQ\|^{-1} C_1 \sum_{i=1}^n \|\xi_i(f_j(X_i) - TQ(X_i))\|_{\psi_2, \Sigma}^2 \quad (50)$$

$$\leq \|f_j - TQ\|^{-1} C_1 \sum_{i=1}^n 4V_{\max} |f_j(X_i) - TQ(X_i)|^2 \quad (51)$$

$$\mathbb{E} Z_{k^*}^2 = \mathbb{E} \max_{j \in [N_\delta]} Z_j^2 \quad (52)$$

$$\leq 32C_1^2 V_{\max}^4 \quad (53)$$

Then

$$\mathbb{E} (\xi \cdot (f_{k^*} - TQ)) = \mathbb{E} (\|f_{k^*} - TQ\| |Z_{k^*}|) \quad (54)$$

$$\leq \mathbb{E} \left(\left(\|\widehat{Q} - TQ\| + \|\widehat{Q} - f_{k^*}\| \right) |Z_{k^*}| \right) \quad (55)$$

$$\leq \mathbb{E} \left(\left(\|\widehat{Q} - TQ\| + n\delta \right) |Z_{k^*}| \right) \quad (56)$$

$$\leq \left(\mathbb{E} \left(\|\widehat{Q} - TQ\| + n\delta \right)^2 \right)^{1/2} \left(\mathbb{E} Z_{k^*}^2 \right)^{1/2} \quad (57)$$

$$\leq \mathbb{E} \left(\|\widehat{Q} - TQ\| + n\delta \right) \left(\mathbb{E} Z_{k^*}^2 \right)^{1/2} \quad (58)$$

$$\leq \left(\sqrt{\mathbb{E} \|\widehat{Q} - TQ\|_2^2} + n\delta \right) \left(\mathbb{E} Z_{k^*}^2 \right)^{1/2} \quad (59)$$

$$\leq \left(\sqrt{\mathbb{E} \|\widehat{Q} - TQ\|_2^2} + n\delta \right) 2V_{\max} \sqrt{n} \quad (60)$$

Where eq. (54) to eq. (55) is by the triangle inequality, eq. (58) to eq. (59) is proposition 2 and eq. (59) to eq. (60) is due to that Combining eq. (41), eq. (43), eq. (44) and eq. (60)

$$\mathbb{E}\|\widehat{Q} - TQ\|^2 \leq \mathbb{E}\|f - TQ\|^2 + 4nV_{\max}\delta + \left(\sqrt{\mathbb{E}\|\widehat{Q} - TQ\|^2} + \sqrt{n}\delta\right) 2V_{\max} \quad (61)$$

$$= 2V_{\max}\sqrt{n}\sqrt{\mathbb{E}\|\widehat{Q} - TQ\|^2} + 6n\delta V_{\max} + \mathbb{E}\|f - TQ\|^2 \quad (62)$$

Lemma 5. 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 5 applied to eq. (62)

$$\frac{1}{n}\mathbb{E}\|\widehat{Q} - TQ\|^2 \leq \frac{(1 + \kappa)^2}{\kappa} \frac{1}{n}V_{\max}^2 + (1 + \kappa) \left(6\delta V_{\max} + \frac{1}{n}\mathbb{E}\|f - TQ\|^2\right) \quad (63)$$

Since this holds for any $f \in \mathcal{F}$ we can further say

$$\frac{1}{n}\mathbb{E}\|\widehat{Q} - TQ\|^2 \leq \frac{(1 + \kappa)^2}{\kappa} \frac{1}{n}V_{\max}^2 + (1 + \kappa) \left(6\delta V_{\max} + \inf_{f \in \mathcal{F}} \frac{1}{n}\mathbb{E}\|f - TQ\|^2\right) \quad (64)$$

$$\leq \frac{(1 + \kappa)^2}{\kappa} \frac{1}{n}V_{\max}^2 + (1 + \kappa) (6\delta V_{\max} + \omega(\mathcal{F})) \quad (65)$$

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

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

$$\left| \left(\widehat{Q}(x) - TQ(x) \right)^2 - \left(f_{k^*}(x) - TQ(x) \right)^2 \right| = \left| \widehat{Q}(x) - f_{k^*}(x) \right| \cdot \left| \widehat{Q}(x) + f_{k^*}(x) - 2TQ(x) \right| \quad (66)$$

$$\leq 4V_{\max}\delta \quad (67)$$

Using this twice we can say

$$(\widehat{Q}(\widehat{X}_i) - TQ(\widehat{X}_i))^2 \quad (68)$$

$$\leq (\widehat{Q}(\widetilde{X}_i) - TQ(\widetilde{X}_i))^2 - (f_{k^*}(\widetilde{X}_i) - TQ(\widetilde{X}_i))^2 + (f_{k^*}(\widetilde{X}_i) - TQ(\widetilde{X}_i))^2 \quad (69)$$

$$\begin{aligned} &\leq (f_{k^*}(\widetilde{X}_i) - TQ(\widetilde{X}_i))^2 + (\widehat{Q}(X_i) - TQ(X_i))^2 - (\widehat{Q}(X_i) - TQ(X_i))^2 \\ &\quad + (f_{k^*}(X_i) - TQ(X_i))^2 - (f_{k^*}(X_i) - TQ(X_i))^2 + 4V_{\max}\delta \end{aligned} \quad (70)$$

$$\leq (\widehat{Q}(X_i) - TQ(X_i))^2 + (f_{k^*}(\widetilde{X}_i) - TQ(\widetilde{X}_i))^2 - (f_{k^*}(X_i) - TQ(X_i))^2 + 8V_{\max}\delta \quad (71)$$

Thus we get

$$\|\widehat{Q} - TQ\|_{\sigma}^2 \quad (72)$$

$$= \mathbb{E} \frac{1}{n} \sum_{i=1}^n (\widehat{Q}(\widetilde{X}_i) - TQ(\widetilde{X}_i))^2 \quad (73)$$

$$\leq \mathbb{E} \frac{1}{n} \sum_{i=1}^n \left((\widehat{Q}(X_i) - TQ(X_i))^2 + (f_{k^*}(\widetilde{X}_i) - TQ(\widetilde{X}_i))^2 - (f_{k^*}(X_i) - TQ(X_i))^2 \right) + 8V_{\max}\delta \quad (74)$$

$$= \frac{1}{n} \|\widehat{Q} - TQ\|^2 + \frac{1}{n} \sum_{i=1}^n h_{k^*}(X_i, \widetilde{X}_i) + 8V_{\max}\delta \quad (75)$$

Where we define

$$h_j(x, y) := (f_j(y) - TQ(y))^2 - (f_j(x) - TQ(x))^2 \quad (76)$$

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

$$T := \max_{j \in [N_\delta]} \left| \sum_{i=1}^n h_j(X_i, \tilde{X}_i) / \Upsilon \right| \quad (77)$$

Then we can bound the middle term in eq. (75)

$$\mathbb{E} \left(\frac{1}{n} \sum_{i=1}^n h_{k^*}(X_i, \tilde{X}_i) \right) \leq \Upsilon / n \mathbb{E} \max_{j \in [N_\delta]} \left(\left| \sum_{i=1}^n h_j(X_i, \tilde{X}_i) / \Upsilon \right| \right) \quad (78)$$

$$\leq \Upsilon / n \mathbb{E} T \quad (79)$$

We want to use Bernsteins inequality (theorem 4) with $U_i = h_j(X_i, \tilde{X}_i)$. Therefore notice that $|h_j| \leq \Upsilon^2$ and

$$\mathbb{V} h_j(X_i, \tilde{X}_i) = 2\mathbb{V} (f_j(X_i) - TQ(X_i))^2 \quad (80)$$

$$\leq 2\mathbb{E} (f_j(X_i) - TQ(X_i))^4 \quad (81)$$

$$\leq 2\Upsilon^4 \quad (82)$$

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

$$\mathbb{E} T = \int_0^\infty \mathbb{P}(T \geq t) \quad (83)$$

$$\leq u + \int_u^\infty \mathbb{P}(T \geq t) dt \quad (84)$$

$$\leq u + \int_u^\infty 2N_\delta \exp \left(\frac{-t^2}{2\Upsilon t/3 + 4n\Upsilon^2} \right) dt \quad (85)$$

$$\leq u + 2N_\delta \int_u^\infty \exp \left(\frac{-t^2}{2\Upsilon^2(t/(3\Upsilon) + 2n)} \right) dt \quad (86)$$

$$\leq u + 2N_\delta \left(\int_u^{6n\Upsilon} \exp \left(\frac{-t^2}{8n\Upsilon^2} \right) dt + \int_{6n\Upsilon}^\infty \exp \left(\frac{-t}{4/3\Upsilon} \right) dt \right) \quad (87)$$

$$\leq u + 2N_\delta \left(\frac{8n\Upsilon}{2u} \exp \left(\frac{-u^2}{8n\Upsilon} \right) + \frac{4\Upsilon}{3} \exp \left(\frac{-24n\Upsilon}{3\Upsilon} \right) \right) \quad (88)$$

where we use lemma 6 from eq. (87) to eq. (88). Now set $u = \Upsilon \sqrt{8n \log N_\delta}$ continuing from eq. (88) we have

$$\dots = \Upsilon \sqrt{8n \log N_\delta} + \frac{\Upsilon^2 8n N_\delta}{\Upsilon \sqrt{8n \log N_\delta}} \exp(-\log N_\delta) + 8/3 N_\delta \Upsilon \exp(-9/2n) \quad (89)$$

$$= \Upsilon 2\sqrt{2n} \left(\log N_\delta + \frac{1}{\log N_\delta} \right) + 8/3 N_\delta e^{-9/2n} \quad (90)$$

$$\leq 4\sqrt{2}\Upsilon \sqrt{n \log N_\delta} + 8/3\Upsilon \quad (91)$$

Inserting eq. (91) and eq. (65) into eq. (75)

$$\left\| \hat{Q} - TQ \right\|_\nu^2 \leq \frac{1}{n} \mathbb{E} \left\| \hat{Q} - TQ \right\|^2 + 8\sqrt{2}V_{\max} n^{-1/2} \sqrt{\log N_\delta} + 8V_{\max}(n^{-1} + \delta) \quad (92)$$

$$\begin{aligned} &\leq \frac{(1+\kappa)^2}{\kappa} \frac{1}{n} V_{\max}^2 + (1+\kappa) (6\delta V_{\max} + \omega(\mathcal{F})) \\ &\quad + 8\sqrt{2}V_{\max} n^{-1/2} \sqrt{\log N_\delta} + 8V_{\max}(n^{-1} + \delta) \end{aligned} \quad (93)$$

□

9 Appendices

9.1 Various lemmas

Lemma 6. 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

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

□