

# Fitted Q-iteration in continuous action-space MDPs

**Authored by:**

Rémi Munos  
Csaba Szepesvári  
Andrés Antos

## **Abstract**

We consider continuous state, continuous action batch reinforcement learning where the goal is to learn a good policy from a sufficiently rich trajectory generated by another policy. We study a variant of fitted Q-iteration, where the greedy action selection is replaced by searching for a policy in a restricted set of candidate policies by maximizing the average action values. We provide a rigorous theoretical analysis of this algorithm, proving what we believe is the first finite-time bounds for value-function based algorithms for continuous state- and action-space problems.

## **1 Paper Body**

We consider continuous state, continuous action batch reinforcement learning where the goal is to learn a good policy from a sufficiently rich trajectory generated by some policy. We study a variant of fitted Q-iteration, where the greedy action selection is replaced by searching for a policy in a restricted set of candidate policies by maximizing the average action values. We provide a rigorous analysis of this algorithm, proving what we believe is the first finite-time bound for value-function based algorithms for continuous state and action problems.

1

### Preliminaries

We will build on the results from [1, 2, 3] and for this reason we use the same notation as these papers. The unattributed results cited in this section can be found in the book [4]. A discounted MDP is defined by a quintuple  $(X, A, P, S, \gamma)$ , where  $X$  is the (possibly infinite) state space,  $A$  is the set of actions,  $P : X \times A \rightarrow \mathcal{M}(X)$  is the transition probability kernel with  $P(\cdot|x, a)$  defining the next-state distribution upon taking action  $a$  from state  $x$ ,  $S(\cdot|x, a)$  gives the corresponding distribution of immediate rewards, and  $\gamma \in (0, 1)$  is the discount factor. Here  $X$  is a measurable space and  $\mathcal{M}(X)$  denotes the set of all probability measures over  $X$ . The Lebesgue measure shall be denoted

by  $\pi$ . We start with the following mild assumption on the MDP: Assumption A1 (MDP Regularity)  $X$  is a compact subset of the  $d_X$ -dimensional Euclidean space,  $A$  is a compact subset of  $[A^-, A^+]^{d_A}$ . The random immediate rewards are bounded by  $R$  and that the expected immediate reward function,  $r(x, a) = \mathbb{E}[R_t | X_t = x, A_t = a]$ , is uniformly bounded by  $R_{\max}$ . A policy determines the next action given the past observations. Here we shall deal with stationary (Markovian) policies which choose an action in a stochastic way based on the last observation only. The value of a policy  $\pi$  when it is started from a state  $x$  is defined as the expected discounted total reward that is encountered while the policy is executed:  $V^\pi(x) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R_t | X_0 = x]$ . Here  $R_t \in \mathbb{R}$  is the reward received at time step  $t$ , the state,  $X_t$ , evolves according to  $X_{t+1} = S(X_t, A_t)$ . Also with: Computer and Automation Research Inst. of the Hungarian Academy of Sciences Kende u. 13-17, Budapest 1111, Hungary.

1

$P(\cdot | X_t, A_t)$ , where  $A_t$  is sampled from the distribution determined by  $\pi$ . We use  $Q^\pi : X \times A \rightarrow \mathbb{R}$  to denote the action-value function of policy  $\pi$ :  $Q^\pi(x, a) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R_t | X_0 = x, A_0 = a]$ . The goal is to find a policy that attains the best possible values,  $V^*(x) = \sup_{\pi} V^\pi(x)$ , at all states  $x \in X$ . Here  $V^*$  is called the optimal value function and a policy  $\pi^*$  that satisfies  $V^\pi(x) = V^*(x)$  for all  $x \in X$  is called optimal. The optimal action-value function  $Q^*(x, a)$  is  $Q^*(x, a) = \sup_{\pi} Q^\pi(x, a)$ . We say that a (deterministic stationary) policy  $\pi$  is greedy w.r.t. an action-value function  $Q$  if  $\pi(a) = \arg\max_{a \in A} Q(x, a)$ , and we write  $\pi = \arg\max_a Q(x, a)$ , if, for all  $x \in X$ ,  $\pi(x) = \arg\max_a Q(x, a)$ . Under mild technical assumptions, such a greedy policy always exists. Any greedy policy w.r.t.  $Q^*$  is optimal. For  $\pi : X \times A \rightarrow \mathbb{R}$ , by  $R$  define its evaluation operator,  $T : B(X \times A) \rightarrow B(X \times A)$  by  $(TQ)(x, a) = r(x, a) + \gamma \mathbb{E}_{\pi} [Q(y, b) | (x, a)]$ . It is known that  $Q^* = TQ^*$ . Further, if we let the Bellman operator,  $T : B(X \times A) \rightarrow B(X \times A)$ , defined by  $(TQ)(x, a) = r(x, a) + \gamma \sup_{b \in A} Q(x, b)$  then  $Q^* = TQ^*$ . It is known that  $V^*$  and  $Q^*$  are bounded by  $R_{\max}/(1-\gamma)$ , just like  $Q$  and  $V$ . For  $\pi : X \times A \rightarrow \mathbb{R}$ , the operator  $E_\pi : B(X \times A) \rightarrow B(X)$  is defined by  $(E_\pi Q)(x) = Q(x, \pi(x))$ , while  $E : B(X \times A) \rightarrow B(X)$  is defined by  $(EQ)(x) = \sup_{a \in A} Q(x, a)$ . Throughout the paper  $\mathcal{F} = \{f : X \times A \rightarrow \mathbb{R}\}$  will denote a subset of real-valued functions over the state-action space  $X \times A$  and  $\mathcal{A}X$  will be a set of policies. For  $\mu \in \mathcal{M}(X)$  and  $f : X \rightarrow \mathbb{R}$  measurable, we let  $(\int f d\mu)_p = \int |f(x)|^p d\mu(x)$ . We simply write  $\|f\|_p$  for  $(\int f^2 d\mu)^{1/2}$ . Further, we extend  $\|f\|_p$  to  $\mathcal{F}$  by  $\|f\|_p = (\int_X \int_A |f(x, a)|^p d\pi(a) d\mu(x))^{1/p}$ , where  $\pi$  is the uniform distribution over  $A$ . We shall use the shorthand notation  $\int f$  to denote the integral  $\int f(x) d\mu(x)$ . We denote the space of bounded measurable functions with domain  $X$  by  $B(X)$ . Further, the space of measurable functions bounded by  $0 \leq f \leq K$  shall be denoted by  $B(X; K)$ . We let  $\|f\|_\infty$  denote the supremum norm.

2

Fitted Q-iteration with approximate policy maximization

We assume that we are given a finite trajectory,  $\{(X_t, A_t, R_t)\}_{t=0}^{T-1}$ , generated by some stochastic stationary policy  $\pi_b$ , called the behavior policy:

$A_t \sim p(\cdot | X_t)$ ,  $X_{t+1} \sim P(\cdot | X_t, A_t)$ ,  $R_t \sim \text{def } S(\cdot | X_t, A_t)$ , where  $p(\cdot | x)$  is a density with  $\int p(\cdot | x) dx = 1$ . The generic recipe for fitted Q-iteration (FQI) [5] is  $Q_{k+1} = \text{Regress}(D_k(Q_k))$ ,

(1)

where  $\text{Regress}$  is an appropriate regression procedure and  $D_k(Q_k)$  is a dataset defining a regression problem in the form of a list of data-point pairs:  $D_k(Q_k) = \{(X_t, A_t), R_t + \gamma \max_b Q_k(X_{t+1}, b) | b \in \mathcal{A}\}$

1  $\leq t \leq N$

Fitted Q-iteration can be viewed as approximate value iteration applied to action-value functions. To see this note that value iteration would assign the value  $(TQ_k)(x, a) = r(x, a) + \gamma \max_b Q_k(y, b) \int p(dy | x, a)$  to  $Q_{k+1}(x, a)$  [6]. Now, remember that the regression function for the jointly distributed random variables  $(Z, Y)$  is defined by the conditional expectation of  $Y$  given  $Z$ :  $m(Z) = E[Y | Z]$ . Since for any fixed function  $Q$ ,  $E[R_t + \gamma \max_b Q(X_{t+1}, b) | X_t, A_t] = (TQ)(X_t, A_t)$ , the regression function corresponding to the data  $D_k(Q)$  is indeed  $TQ$  and hence if FQI solved the regression problem defined by  $Q_k$  exactly, it would simulate value iteration exactly. However, this argument itself does not directly lead to a rigorous analysis of FQI: Since  $Q_k$  is obtained based on the data, it is itself a random function. Hence, after the first iteration, the "target" function in FQI becomes random. Furthermore, this function depends on the same data that is used to define the regression problem. Will FQI still work despite these issues? To illustrate the potential difficulties consider a dataset where  $X_1, \dots, X_N$  is a sequence of independent random variables, which are all distributed uniformly at random in  $[0, 1]$ . Further, let  $M$  be a random integer greater than  $N$  which is independent of the dataset  $(X_t)_{t=1}^N$ . Let  $U$  be another random variable, uniformly distributed in  $[0, 1]$ . Now define the regression problem by  $Y_t = f_{M,U}(X_t)$ , where  $f_{M,U}(x) = \text{sgn}(\sin(2\pi M U (x + U)))$ . Then it is not hard to see that no matter how big  $N$  is, no procedure can

Since the designer controls  $Q_k$ , we may assume that it is continuous, hence the maximum exists.

2

estimate the regression function  $f_{M,U}$  with a small error (in expectation, or with high probability), even if the procedure could exploit the knowledge of the specific form of  $f_{M,U}$ . On the other hand, if we restricted  $M$  to a finite range then the estimation problem could be solved successfully. The example shows that if the complexity of the random functions defining the regression problem is uncontrolled then successful estimation might be impossible. Amongst the many regression methods in this paper we have chosen to work with least-squares methods. In this case Equation (1) takes the form  $Q_{k+1} = \text{argmin}_Q \sum_{t=1}^N \sum_{(X_t, A_t)} (R_t + \gamma \max_b Q(X_{t+1}, b) - Q(X_t, A_t))^2$ . We call this method the least-squares fitted Q-iteration (LSFQI) method. Here we introduced the weighting  $1/p(A_t | X_t)$  since we do not want to give more weight to those actions that are preferred by the behavior policy. Besides this weighting, the only parameter of the method is the function set  $\mathcal{F}$ . This function set should be chosen carefully, to keep a

balance between the representation power and the number of samples. As a specific example for  $F$  consider neural networks with some fixed architecture. In this case the function set is generated by assigning weights in all possible ways to the neural net. Then the above minimization becomes the problem of tuning the weights. Another example is to use linearly parameterized function approximation methods with appropriately selected basis functions. In this case the weight tuning problem would be less demanding. Yet another possibility is to let  $F$  be an appropriate restriction of a Reproducing Kernel Hilbert Space (e.g., in a ball). In this case the training procedure becomes similar to LS-SVM training [7]. As indicated above, the analysis of this algorithm is complicated by the fact that the new dataset is defined in terms of the previous iterate, which is already a function of the dataset. Another complication is that the samples in a trajectory are in general correlated and that the bias introduced by the imperfections of the approximation architecture may yield to an explosion of the error of the procedure, as documented in a number of cases in, e.g., [8]. Nevertheless, at least for finite action sets, the tools developed in [1, 3, 2] look suitable to show that under appropriate conditions these problems can be overcome if the function set is chosen in a judicious way. However, the results of these works would become essentially useless in the case of an infinite number of actions since these previous bounds grow to infinity with the number of actions. Actually, we believe that this is not an artifact of the proof techniques of these works, as suggested by the counterexample that involved random targets. The following result elaborates this point further: Proposition 2.1. Let  $F \subset B(X \rightarrow A)$ . Then even if the pseudo-dimension of  $F$  is finite, the fatshattering function of  $F_{\max} = \{VQ : VQ(\cdot) = \max_{a \in A} Q(\cdot, a), Q \in F\}$

2

can be infinite over  $(0, 1/2)$ .

Without going into further details, let us just note that the finiteness of the fat-shattering function is a sufficient and necessary condition for learnability and the finiteness of the fat-shattering function is implied by the finiteness of the pseudo-dimension [9]. The above proposition thus shows that without imposing further special conditions on  $F$ , the learning problem may become infeasible. One possibility is of course to discretize the action space, e.g., by using a uniform grid. However, if the action space has a really high dimensionality, this approach becomes unfeasible (even enumerating  $2^d A$  points could be impossible when  $dA$  is large). Therefore we prefer alternate solutions. Another possibility is to make the functions in  $F$ , e.g., uniformly Lipschitz in their state coordinates. Then the same property will hold for functions in  $F_{\max}$  and hence by a classical result we can bound the capacity of this set (cf. pp. 353-357 of [10]). One potential problem with this approach is that this way it might be difficult to get a fine control of the capacity of the resulting set. 2 The proof of this and the other results are given in the appendix, available in the extended version of this paper, downloadable from <http://hal.inria.fr/inria-00185311/en/>.

3

In the approach explored here we modify the fitted Q-iteration algorithm by introducing a policy set  $\Pi$  and a search over this set for an approximately

greedy policy in a sense that will be made precise in a minute. Our algorithm thus has four parameters:  $F, \mathcal{F}, K, Q_0$ . Here  $F$  is as before,  $\mathcal{F}$  is a user-chosen set of policies (mappings from  $X$  to  $A$ ),  $K$  is the number of iterations and  $Q_0$  is an initial value function (a typical choice is  $Q_0 \equiv 0$ ). The algorithm computes a sequence of iterates  $(Q_k, \pi_k)$ ,  $k = 0, \dots, K$ , defined by the following equations:  $\pi_0 \equiv Q_{k+1} \equiv \pi_{k+1}$

$$\begin{aligned} &= \\ &\arg\max_{\pi \in \mathcal{F}} \\ &= \\ &\arg\min_{\pi \in \mathcal{F}} \\ &\arg\max_{\pi \in \mathcal{F}} \\ &Q_0(X_t, \pi(X_t)), \\ &t=1 \\ &Q \leftarrow F \\ &= \\ &\sum_{t=1}^N \sum_{X \in \mathcal{X}} \\ &\sum_{t=1}^N \sum_{X \in \mathcal{X}} Q(X_t, A_t) - R_t + \gamma Q(X_{t+1}, \pi(X_{t+1})) - \gamma b(A_t - X_t) \\ (3) \\ &Q_{k+1}(X_t, \pi(X_t)). \\ (4) \\ &t=1 \end{aligned}$$

Thus, (3) is similar to (2), while (4) defines the policy search problem. The policy search will generally be solved by a gradient procedure or some other appropriate method. The cost of this step will be primarily determined by how well-behaving the iterates  $Q_{k+1}$  are in their action arguments. For example, if they were quadratic and if  $\gamma$  was linear then the problem would be a quadratic optimization problem. However, except for special cases<sup>3</sup> the action value functions will be more complicated, in which case this step can be expensive. Still, this cost could be similar to that of searching for the maximizing actions for each  $t = 1, \dots, N$  if the approximately maximizing actions are similar across similar states. This algorithm, which we could also call a fitted actor-critic algorithm, will be shown to overcome the above mentioned complexity control problem provided that the complexity of  $\mathcal{F}$  is controlled appropriately. Indeed, in this case the set of possible regression problems is determined by the set  $F \equiv \{V : V(\pi) = Q(\pi, \pi(\pi)), Q \in F, \pi \in \mathcal{F}\}$ , and the proof will rely on controlling the complexity of  $F$  by selecting  $F$  and  $\mathcal{F}$  appropriately.

### 3.3.1

The main theoretical result Outline of the analysis

In order to gain some insight into the behavior of the algorithm, we provide a brief summary of its error analysis. The main result will be presented subsequently. For  $f, Q \in F$  and a policy  $\pi$ , we define the  $t$ th TD-error as follows:  $d_t(f; Q, \pi) = R_t + \gamma Q(X_{t+1}, \pi(X_{t+1})) - f(X_t, A_t)$ . Further, we define the empirical loss function by  $\sum_{t=1}^N d_t^2(f; Q, \pi) = \sum_{t=1}^N (f(X_t, A_t) - \gamma Q(X_{t+1}, \pi(X_{t+1})))^2$ .

where the normalization with  $\frac{1}{\sqrt{N}}$  is introduced for mathematical convenience. Then (3) can be written compactly as  $Q_{k+1} = \arg\min_{Q \in \mathcal{F}} \frac{1}{\sqrt{N}} \sum_{i=1}^N L(f_i; Q, \pi)$ . The algorithm can then be motivated by the observation that for any  $f, Q$ , and  $\pi$ ,  $L(f; Q, \pi)$  is an unbiased estimate of  $\mathcal{L}(f; Q, \pi) = \mathbb{E} [T Q_k + L(f; Q, \pi)]$ , where the first term is the error we are interested in and the second term captures the variance of the random samples:  $Z_L(f; Q, \pi) = \mathbb{E} [\text{Var} [R_1 + Q(X_2, \pi(X_2)) - X_1, A_1 = a]]$ . A linear quadratic regulation is such a nice case. It is interesting to note that in this special case the obvious choices for  $\mathcal{F}$  and  $\pi$  yield zero error in the limit, as can be proven based on the main result of this paper.

4

Since the variance term in (5) is independent of  $f$ ,  $\arg\min_{Q \in \mathcal{F}} \frac{1}{\sqrt{N}} \sum_{i=1}^N L(f_i; Q, \pi) = \arg\min_{Q \in \mathcal{F}} \frac{1}{\sqrt{N}} \sum_{i=1}^N L(f_i; Q, \pi)$ . Thus, if  $\pi_k$  were greedy w.r.t.  $Q_k$  then  $\arg\min_{Q \in \mathcal{F}} \frac{1}{\sqrt{N}} \sum_{i=1}^N L(f_i; Q, \pi_k) = 2 \arg\min_{Q \in \mathcal{F}} \frac{1}{\sqrt{N}} \sum_{i=1}^N L(f_i; Q, \pi_k)$ . Hence we can still think of the procedure as approximate value iteration over the space of action-value functions, projecting  $T Q_k$  using empirical risk minimization on the space  $\mathcal{F}$  w.r.t.  $\frac{1}{\sqrt{N}} \sum_{i=1}^N L(f_i; Q, \pi_k)$  distances in an approximate manner. Since  $\pi_k$  is only approximately greedy, we will have to deal with both the error coming from the approximate projection and the error coming from the choice of  $\pi_k$ . To make this clear, we write the iteration in the form  $Q_{k+1} = T \pi_k Q_k + \eta_k = T Q_k + \eta_k + (T \pi_k Q_k - T Q_k) = T Q_k + \eta_k$ , def

where  $\eta_k$  is the error committed while computing  $T \pi_k Q_k$ ,  $\eta_{0k} = T \pi_k Q_k - T Q_k$  is the error committed because the greedy policy is computed approximately and  $\eta_k = \eta_k + \eta_{0k}$  is the total error of step  $k$ . Hence, in order to show that the procedure is well behaved, one needs to show that both errors are controlled and that when the errors are propagated through these equations, the resulting error stays controlled, too. Since we are ultimately interested in the performance of the policy obtained, we will also need to show that small action-value approximation errors yield small performance losses. For these we need a number of assumptions that concern either the training data, the MDP, or the function sets used for learning.

### Assumptions

#### 3.2.1

##### Assumptions on the training data

We shall assume that the data is rich, is in a steady state, and is fast-mixing, where, informally, mixing means that future depends weakly on the past. Assumption A2 (Sample Path Properties) Assume that  $\{(X_t, A_t, R_t)\}_{t=1, \dots, N}$  is the sample path of  $\pi_b$ , a stochastic stationary policy. Further, assume that  $\{X_t\}$  is strictly stationary ( $X_t \sim M(X)$ ) and exponentially  $\beta$ -mixing with the actual rate given by the parameters  $(\beta, b, \gamma)$ .<sup>4</sup> We further assume that the sampling policy  $\pi_b$  satisfies  $\gamma_0 = \inf_{(x,a)} \sum_{i=1}^{\infty} \gamma_i^a(x) > 0$ . The  $\beta$ -mixing property will be used to establish tail inequalities for certain empirical processes.<sup>5</sup> Note that the mixing coefficients do not need to be known. In the case when no mixing condition is satisfied, learning might be impossible. To see this just consider the case when  $X_1 = X_2 = \dots = X_N$ . Thus, in this

case the learner has many copies of the same random variable and successful generalization is thus impossible. We believe that the assumption that the process is in a steady state is not essential for our result, as when the process reaches its steady state quickly then (at the price of a more involved proof) the result would still hold. 3.2.2

### Assumptions on the MDP

In order to prevent the uncontrolled growth of the errors as they are propagated through the updates, we shall need some assumptions on the MDP. A convenient assumption is the following one [11]: Assumption A3 (Uniformly stochastic transitions) For all  $x \in X$  and  $a \in A$ , assume that  $P(\cdot|x, a)$  is absolutely continuous w.r.t.  $\mu$  and the derivative of  $P$  w.r.t.  $\mu$  is bounded.  $\mu$  Radon-Nikodym def  $\frac{dP}{d\mu}(\cdot|x, a)$  uniformly with bound  $C$ :  $C = \sup_{x \in X, a \in A} \int \frac{dP}{d\mu}(\cdot|x, a) d\mu \leq C$ .

Note that by the definition of measure differentiation, Assumption A3 means that  $P(\cdot|x, a) \neq C(\cdot|?)$ . This assumption essentially requires the transitions to be noisy. We will also prove (weaker) results under the following, weaker assumption: 4

For the definition of  $\gamma$ -mixing, see e.g. [2]. We say  $\gamma$ -empirical process? and  $\gamma$ -empirical measure?, but note that in this work these are based on dependent (mixing) samples. 5

5

Assumption A4 (Discounted-average concentrability of future-state distributions) Given  $\mathcal{S}, \mathcal{A}$ ,  $m \geq 1$  and an arbitrary sequence of stationary policies  $\{\pi_m\}_{m \geq 1}$ , assume that the future state distribution  $\mathbb{P}^{\pi_1} \mathbb{P}^{\pi_2} \dots \mathbb{P}^{\pi_m}$  is absolutely continuous w.r.t.  $\mathbb{P}$ . Assume that  $c(m) = \mathbb{P}^{\pi_1} \mathbb{P}^{\pi_2} \dots \mathbb{P}^{\pi_m} \left( \sup_{\mathcal{S}} \frac{d(\mathbb{P}^{\pi_1} \mathbb{P}^{\pi_2} \dots \mathbb{P}^{\pi_m})}{d\mathbb{P}} \right)$  satisfies  $m \geq 1 \implies m \geq 1 \implies c(m) \leq \gamma$ . We shall call  $C^{\gamma, \gamma} = \mathbb{P}^{\pi_1} \mathbb{P}^{\pi_2} \dots \mathbb{P}^{\pi_m} \max_{\mathcal{S}} (1 - \gamma)^2 m \geq 1 \implies m \geq 1 \implies c(m), (1 - \gamma)^2 m \geq 1 \implies m \geq 1 \implies m c(m)$  the discounted-average concentrability coefficient of the future-state distributions. The number  $c(m)$  measures how much  $\mathbb{P}$  can get amplified in  $m$  steps as compared to the reference distribution  $\mathbb{P}$ . Hence, in general we expect  $c(m)$  to grow with  $m$ . In fact, the condition that  $C^{\gamma, \gamma}$  is finite is a growth rate condition on  $c(m)$ . Thanks to discounting,  $C^{\gamma, \gamma}$  is finite for a reasonably large class of systems (see the discussion in [11]). A related assumption is needed in the error analysis of the approximate greedy step of the algorithm: Assumption A5 (The random policy  $\pi$  makes no peak-states?) Consider the distribution  $\mathbb{P} = (\mathbb{P}^{\pi_1} \mathbb{P}^{\pi_2} \dots \mathbb{P}^{\pi_m})^{\pi}$  which is the distribution of a state that results from sampling an initial state according to  $\mathbb{P}$  and then executing an action which is selected uniformly at random. Then  $\mathbb{P} = \mathbb{P}^{\pi_1} \mathbb{P}^{\pi_2} \dots \mathbb{P}^{\pi_m} \mathbb{P}^{\pi}$ . Note that under Assumption A3 we have  $\mathbb{P} \leq C^{\gamma, \gamma} \mathbb{P}$ . This (very mild) assumption means that after one step, starting from  $\mathbb{P}$  and executing this random policy, the probability of the next state being in a set is upper bounded by  $C^{\gamma, \gamma}$ -times the probability of the starting state being in the same set.  $\square$

Besides, we assume that  $A$  has the following regularity property: Let  $\text{Py}(a, h, ?) = ? \times (a_0, v) \times \text{Rd}A + 1 : k_a \times a_0 k_1 \times ?, 0 \times v/h \times 1 \times k_a \times a_0 k_1 / ?$  denote the pyramid with hight  $h$  and base given by the  $1$ -ball  $B(a, ?) = a_0 \times \text{Rd}A : k_a \times a_0 k_1 \times ?$  centered at  $a$ . Assumption A6 (Regularity of

the action space) We assume that there exists  $\epsilon > 0$ , such that for all  $a \in A$ , for all  $\gamma \in [0, 1]$ ,  $\|P_\gamma(a, \cdot, \cdot) - P_\gamma(a, \cdot, \cdot)\| \leq \epsilon$  (where  $P_\gamma(a, \cdot, \cdot) = \int P(\cdot | x, a) d\gamma(x)$ ). For example, if  $A$  is an '1'-ball itself, then this assumption will be satisfied with  $\epsilon = 2\gamma dA$ . Without assuming any smoothness of the MDP, learning in infinite MDPs looks hard (see, e.g., [12, 13]). Here we employ the following extra condition: Assumption A7 (Lipschitzness of the MDP in the actions) Assume that the transition probabilities and rewards are Lipschitz w.r.t. their action variable, i.e., there exists  $L_P, L_r \geq 0$  such that for all  $(x, a, a_0) \in X \times A \times A$  and measurable set  $B$  of  $X$ ,  $\|P(B | x, a) - P(B | x, a_0)\| \leq L_P \|a - a_0\|$ ,  $\|r(x, a) - r(x, a_0)\| \leq L_r \|a - a_0\|$ .

Note that previously Lipschitzness w.r.t. the state variables was used, e.g., in [11] to construct consistent planning algorithms.

Assumptions on the function sets used by the algorithm

These assumptions are less demanding since they are under the control of the user of the algorithm. However, the choice of these function sets will greatly influence the performance of the algorithm, as we shall see it from the bounds. The first assumption concerns the class  $F$ : Assumption A8 (Lipschitzness of candidate action-value functions) Assume  $F \subset B(X \times A)$  and that any elements of  $F$  is uniformly Lipschitz in its action-argument in the sense that  $\|Q(x, a) - Q(x, a_0)\| \leq L_A \|a - a_0\|$  holds for any  $x \in X$ ,  $a, a_0 \in A$ , and  $Q \in F$ .

Remember that  $\gamma_A$  denotes the uniform distribution over the action set  $A$ .

6

We shall also need to control the capacity of our function sets. We assume that the reader is familiar with the concept of VC-dimension.<sup>7</sup> Here we use the pseudo-dimension of function sets that builds upon the concept of VC-dimension: Definition 3.1 (Pseudo-dimension). The pseudo-dimension  $VF$  of  $F$  is defined as the VC-dimension of the subgraphs of functions in  $F$  (hence it is also called the VC-subgraph dimension of  $F$ ). Since  $A$  is multidimensional, we define  $V^+_F$  to be the sum of the pseudo-dimensions of the coordinate projection spaces,  $\sum_k VF_k$  of  $F$ :

$$\begin{aligned} V^+_F &= \sum_{k=1}^d VF_k \\ &= \sum_{k=1}^d \dim \{ \pi_k \circ f : f \in F \} \end{aligned}$$

Now we are ready to state our assumptions on our function sets: Assumption A9 (Capacity of the function and policy sets) Assume that  $F \subset B(X \times A; Q_{\max})$  for  $Q_{\max} \geq 0$  and  $VF + V^+_F < \infty$ . Also,  $A \in \mathcal{A}$  and  $V^+_A < \infty$ . Besides their capacity, one shall also control the approximation power of the function sets involved. Let us first consider the policy set  $\Pi$ . Introduce  $e_\Pi(F, \gamma) = \sup_{Q \in F} \inf_{\pi \in \Pi} \|EQ - E\pi(Q)\|$ .

Note that  $\inf_{\pi \in \Pi} \|EQ - E\pi(Q)\|$  measures the quality of approximating  $EQ$  by  $E\pi(Q)$ . Hence,  $e_\Pi(F, \gamma)$  measures the worst-case approximation error of  $EQ$  as  $Q$  is changed within  $F$ . This can be made small by choosing  $\gamma$  large.



Another related quantity is the one-step Bellman-error of  $F$  w.r.t.  $\pi$ . This is defined as follows: For a fixed policy  $\pi$ , the one-step Bellman-error of  $F$  w.r.t.  $\pi$  is defined as  $E_1(F; \pi) = \sup_{\mathbf{x}} \inf_{\mathbf{a}} Q_0^\pi(\mathbf{x}, \mathbf{a}) - T^\pi Q^\pi F(\mathbf{x}, \mathbf{a})$

Taking again a pessimistic approach, the one-step Bellman-error of  $F$  is defined as  $E_1(F, \pi) = \sup_{\mathbf{x}} E_1(F; \pi)$ .

Typically by increasing  $F$ ,  $E_1(F, \pi)$  can be made smaller (this is discussed at some length in [3]). However, it also holds for both  $\pi$  and  $F$  that making them bigger will increase their capacity (pseudo-dimensions) which leads to an increase of the estimation errors. Hence,  $F$  and  $\pi$  must be selected to balance the approximation and estimation errors, just like in supervised learning.

The main result

**Theorem 3.2.** Let  $\pi_K$  be a greedy policy w.r.t.  $Q_K$ , i.e.  $\pi_K(\mathbf{x}) = \arg\max_{\mathbf{a}} A Q_K(\mathbf{x}, \mathbf{a})$ . Then under Assumptions A1, A2, and A5-A9, for all  $\epsilon > 0$  we have with probability at least  $1 - \epsilon$ : given Assumption A3 (respectively A4),  $kV \pi \pi V \pi K k_1$  (resp.  $kV \pi \pi V \pi K k_1$ ), is bounded by  $\epsilon \leq 1 + 1/\epsilon + 1/\epsilon^2 A \pi \pi 4? (\log N + \log(K/\epsilon)) K \pi + \pi$ ,  $C \pi E_1(F, \pi) + e^\pi(F, \pi) + 1/4 \pi \pi N \pi \pi A$  where  $C$  depends on  $dA$ ,  $VF + \pi$ ,  $(V\pi + \pi)dk=1$ ,  $\pi$ ,  $\pi$ ,  $b$ ,  $\pi$ ,  $C^\pi$  (resp.  $C^\pi, \pi$ ),  $\pi \pi$ ,  $LA$ ,  $LP$ ,  $Lr$ ,  $\pi$ ,  $\pi(A)$ ,  $\pi_0$ ,  $k$

$\pi + 1$

$\pi \max$ , and  $A^\pi$ . In particular,  $C$  scales with  $V^{4(dA+1)}$ , where  $V = 2VF + \pi V\pi + Q_{\max}$ ,  $R_{\max}$ ,  $R$  plays the role of the combined effective dimension of  $F$  and  $\pi$ . Readers not familiar with VC-dimension are suggested to consult a book, such as the one by Anthony and Bartlett [14].

7

4

Discussion

We have presented what we believe is the first finite-time bounds for continuous-state and actionspace RL that uses value functions. Further, this is the first analysis of fitted Q-iteration, an algorithm that has proved to be useful in a number of cases, even when used with non-averagers for which no previous theoretical analysis existed (e.g., [15, 16]). In fact, our main motivation was to show that there is a systematic way of making these algorithms work and to point at possible problem sources the same time. We discussed why it can be difficult to make these algorithms work in practice. We suggested that either the set of action-value candidates has to be carefully controlled (e.g., assuming uniform Lipschitzness w.r.t. the state variables), or a policy search step is needed, just like in actorcritic algorithms. The bound in this paper is similar in many respects to a previous bound of a Bellman-residual minimization algorithm [2]. It looks that the techniques developed here can be used to obtain results for that algorithm when it is applied to continuous action spaces. Finally, although we have not explored them here, consistency results for FQI can be obtained from our results using standard methods, like the methods of sieves. We believe that the methods developed here will eventually lead to algorithms where the function approximation methods are chosen based on the data (similar to adaptive regression methods) so as to optimize performance, which in our opinion is one of the biggest open questions in RL. Currently we are exploring this pos-

sibility. Acknowledgments Andr as Antos would like to acknowledge support for this project from the Hungarian Academy of Sciences (Bolyai Fellowship). Csaba Szepesv ari greatly acknowledges the support received from the Alberta Ingenuity Fund, NSERC, the Computer and Automation Research Institute of the Hungarian Academy of Sciences.

## 2 References

- [1] A. Antos, Cs. Szepesv ari, and R. Munos. Learning near-optimal policies with Bellman-residual minimization based fitted policy iteration and a single sample path. In COLT-19, pages 574?588, 2006.
- [2] A. Antos, Cs. Szepesv ari, and R. Munos. Learning near-optimal policies with Bellman-residual minimization based fitted policy iteration and a single sample path. Machine Learning, 2007. (accepted).
- [3] A. Antos, Cs. Szepesv ari, and R. Munos. Value-iteration based fitted policy iteration: learning with a single trajectory. In IEEE ADPRL, pages 330?337, 2007.
- [4] D. P. Bertsekas and S.E. Shreve. Stochastic Optimal Control (The Discrete Time Case). Academic Press, New York, 1978.
- [5] D. Ernst, P. Geurts, and L. Wehenkel. Tree-based batch mode reinforcement learning. Journal of Machine Learning Research, 6:503?556, 2005.
- [6] R.S. Sutton and A.G. Barto. Reinforcement Learning: An Introduction. Bradford Book. MIT Press, 1998.
- [7] N. Cristianini and J. Shawe-Taylor. An introduction to support vector machines (and other kernel-based learning methods). Cambridge University Press, 2000.
- [8] J.A. Boyan and A.W. Moore. Generalization in reinforcement learning: Safely approximating the value function. In NIPS-7, pages 369?376, 1995.
- [9] P.L. Bartlett, P.M. Long, and R.C. Williamson. Fat-shattering and the learnability of real-valued functions. Journal of Computer and System Sciences, 52:434?452, 1996.
- [10] A.N. Kolmogorov and V.M. Tihomirov.  $\epsilon$ -entropy and  $\epsilon$ -capacity of sets in functional space. American Mathematical Society Translations, 17(2):277?364, 1961.
- [11] R. Munos and Cs. Szepesv ari. Finite time bounds for sampling based fitted value iteration. Technical report, Computer and Automation Research Institute of the Hungarian Academy of Sciences, Kende u. 13-17, Budapest 1111, Hungary, 2006.
- [12] A.Y. Ng and M. Jordan. PEGASUS: A policy search method for large MDPs and POMDPs. In Proceedings of the 16th Conference in Uncertainty in Artificial Intelligence, pages 406?415, 2000.
- [13] P.L. Bartlett and A. Tewari. Sample complexity of policy search with known dynamics. In NIPS-19. MIT Press, 2007.
- [14] M. Anthony and P. L. Bartlett. Neural Network Learning: Theoretical Foundations. Cambridge University Press, 1999.
- [15] M. Riedmiller. Neural fitted Q iteration ? first experiences with a data efficient neural reinforcement learning method. In 16th European Conference on Machine Learning, pages 317?328, 2005.
- [16] S. Kalyanakrishnan and P. Stone. Batch reinforcement learning in a complex domain. In AAMAS-07, 2007.