

Notes on Fitted Q-iteration

Nan Jiang

October 9, 2020

Let $M = (\mathcal{S}, \mathcal{A}, P, R, \gamma, d_0)$ be an MDP, where d_0 is the initial distribution over states. Given a dataset $\{(s, a, r, s')\}$ generated from M and a Q-function class $\mathcal{F} \subset \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$, we want to analyze the guarantee of Fitted Q-Iteration. This note is inspired by and scrutinizes the results in Approximate Value/Policy Iteration literature [e.g., 1, 2, 3] under simplification assumptions.

Setup and Assumptions

1. \mathcal{F} is finite but can be exponentially large.
2. Realizability: $Q^* \in \mathcal{F}$.
3. \mathcal{F} is closed under Bellman update: $\forall f \in \mathcal{F}, \mathcal{T}f \in \mathcal{F}$. (For finite \mathcal{F} , this implies realizability.)
4. The dataset $D = \{(s, a, r, s')\}$ is generated as follows: $(s, a) \sim \mu \times U(\mathcal{A})$ ($U(\mathcal{A})$ is uniform over actions), $r \sim R(s, a)$, $s' \sim P(s, a)$. Define the empirical update $\widehat{\mathcal{T}}_{\mathcal{F}} f'$ as

$$\mathcal{L}_D(f; f') := \frac{1}{|D|} \sum_{(s, a, r, s') \in D} (f(s, a) - r - \gamma V_{f'}(s'))^2.$$
$$\widehat{\mathcal{T}}_{\mathcal{F}} f' := \arg \min_{f \in \mathcal{F}} \mathcal{L}_D(f; f'),$$

where $V_{f'}(s') := \max_{a'} f'(s', a')$. Note that by completeness, $\mathcal{T}f' \in \mathcal{F}$ is the Bayes optimal regressor for the regression problem defined in $\mathcal{L}_D(f; f')$. It will also be useful to define

$$\mathcal{L}_{\mu \times U}(f; f') := \mathbb{E}_D[\mathcal{L}_D(f; f')].$$

5. For any function $g : \mathcal{S} \rightarrow \mathbb{R}$, any distribution $\nu \in \Delta(\mathcal{S})$, and $p \geq 1$, define $\|g\|_{p, \nu} := (\mathbb{E}_{s \sim \nu}[|g(s)|^p])^{1/p}$, and let $\|g\|_{\nu}$ be a shorthand for $\|g\|_{2, \nu}$. Such norms are similarly defined for functions over $\mathcal{S} \times \mathcal{A}$.
6. Let d_h^{π} be the distribution of s_h under π , that is, $d_h^{\pi}(s) := \Pr[s_h = s \mid s_1 \sim d_0, \pi]$.
7. μ is exploratory: for a distribution $\nu \in \Delta(\mathcal{S})$ generated by any (non-stationary) policy at any time step (that is, any distribution ν of the form d_h^{π} where π may be non-stationary),

$$\forall s \in \mathcal{S}, \frac{\nu(s)}{\mu(s)} \leq C.$$

As a consequence, $\|\cdot\|_{\nu} \leq \sqrt{C} \|\cdot\|_{\mu}$. Similarly, when we couple μ with a uniform distribution over \mathcal{A} , we have similar results for state-action distributions: $\|\cdot\|_{\nu \times \pi} \leq \sqrt{|\mathcal{A}|C} \|\cdot\|_{\mu \times U}$. See slides for example scenarios where C is naturally bounded.

8. Algorithm (simplified for analysis): let $f_0 \equiv \mathbf{0}$ (assuming $\mathbf{0} \in \mathcal{F}$), and for $k \geq 1$, $f_k := \widehat{\mathcal{T}}_{\mathcal{F}} f_{k-1}$.
9. Uniform deviation bound (can be obtained by concentration inequalities and union bound):

$$\forall f, f' \in \mathcal{F}, |\mathcal{L}_D(f; f') - \mathcal{L}_{\mu \times U}(f; f')| \leq \epsilon.$$

(Note: at the end we will show how to obtain fast rates.)

Goal Let $\hat{\pi} := \pi_{f_k}$. Derive an upper bound on $J(\pi^*) - J(\hat{\pi})$.

Analysis

$$\begin{aligned} J(\pi^*) - J(\hat{\pi}) &= \sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{E}_{s \sim d_h^{\hat{\pi}}} [V^*(s) - Q^*(s, \hat{\pi})] \\ &\leq \sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{E}_{s \sim d_h^{\hat{\pi}}} [Q^*(s, \pi^*) - f_k(s, \pi^*) + f_k(s, \hat{\pi}) - Q^*(s, \hat{\pi})] \\ &\leq \sum_{h=1}^{\infty} \gamma^{h-1} \left(\|Q^* - f_k\|_{1, d_h^{\hat{\pi}} \times \pi^*} + \|Q^* - f_k\|_{1, d_h^{\hat{\pi}} \times \hat{\pi}} \right) \\ &\leq \sum_{h=1}^{\infty} \gamma^{h-1} \left(\|Q^* - f_k\|_{d_h^{\hat{\pi}} \times \pi^*} + \|Q^* - f_k\|_{d_h^{\hat{\pi}} \times \hat{\pi}} \right). \end{aligned} \quad (1)$$

The last line contains two terms, both in the form of $\|Q^* - f_k\|_{\nu \times \pi}$. So it remains to bound $\|Q^* - f_k\|_{\nu \times \pi}$ for any $\nu \times \pi \in \Delta(\mathcal{S} \times \mathcal{A})$ that combines any $\nu \in \Delta(\mathcal{S})$ that satisfies bullet 4 with any $\pi : \mathcal{S} \rightarrow \mathcal{A}$.

First a helper lemma:

Lemma 1. Define $\pi_{f, f_k}(s) := \arg \max_{a \in \mathcal{A}} \max\{f(s, a), f_k(s, a)\}$. Then we have $\forall \nu \in \Delta(\mathcal{S})$,

$$\|V_f - V_{f_k}\|_{\nu} \leq \|f - f_k\|_{\nu \times \pi_{f, f_k}}.$$

Proof.

$$\begin{aligned} \|V_f - V_{f_k}\|_{\nu}^2 &= \sum_{s \in \mathcal{S}} \nu(s) (\max_{a \in \mathcal{A}} f(s, a) - \max_{a' \in \mathcal{A}} f_k(s, a'))^2 \\ &\leq \sum_{s \in \mathcal{S}} \nu(s) (f(s, \pi_{f, f_k}) - f_k(s, \pi_{f, f_k}))^2 = \|f - f_k\|_{\nu \times \pi_{f, f_k}}^2. \end{aligned} \quad \square$$

Now we can bound $\|Q^* - f_k\|_{\nu \times \pi}$ using Lemma 1. Define $P(\nu \times \pi)$ as a distribution over \mathcal{S} generated as $s' \sim P(\nu \times \pi) \Leftrightarrow (s, a) \sim \nu \times \pi, s' \sim P(s, a)$, and

$$\begin{aligned} \|f_k - Q^*\|_{\nu \times \pi} &= \|f_k - \mathcal{T}f_{k-1} + \mathcal{T}f_{k-1} - Q^*\|_{\nu \times \pi} \\ &\leq \|f_k - \mathcal{T}f_{k-1}\|_{\nu \times \pi} + \|\mathcal{T}f_{k-1} - \mathcal{T}Q^*\|_{\nu \times \pi} \\ &\leq \sqrt{|\mathcal{A}|C} \|f_k - \mathcal{T}f_{k-1}\|_{\mu \times U} + \gamma \|V_{f_{k-1}} - V^*\|_{P(\nu \times \pi)} \quad (*) \\ &\leq \sqrt{|\mathcal{A}|C} \|f_k - \mathcal{T}f_{k-1}\|_{\mu \times U} + \gamma \|f_{k-1} - Q^*\|_{P(\nu \times \pi) \times \pi_{f_{k-1}, Q^*}}. \end{aligned} \quad (\text{Lemma 1})$$

Step (*) holds because:

$$\begin{aligned}
\|\mathcal{T}f_{k-1} - \mathcal{T}Q^*\|_{\nu \times \pi}^2 &= \mathbb{E}_{(s,a) \sim \nu \times \pi} \left[((\mathcal{T}f_{k-1})(s,a) - (\mathcal{T}Q^*)(s,a))^2 \right] \\
&= \mathbb{E}_{(s,a) \sim \nu \times \pi} \left[\left(\gamma \mathbb{E}_{s' \sim P(s,a)} [V_{f_{k-1}}(s') - V^*(s')] \right)^2 \right] \\
&\leq \gamma^2 \mathbb{E}_{(s,a) \sim \nu \times \pi, s' \sim P(s,a)} \left[(V_{f_{k-1}}(s') - V^*(s'))^2 \right] && \text{(Jensen)} \\
&= \gamma^2 \mathbb{E}_{s' \sim P(\nu \times \pi)} \left[(V_{f_{k-1}}(s') - V^*(s'))^2 \right] = \gamma^2 \|V_{f_{k-1}} - V^*\|_{P(\nu \times \pi)}^2.
\end{aligned}$$

Note that we can apply the same analysis on $P(\nu \times \pi) \times \pi_{f_{k-1}, Q^*}$ and expand the inequality k times. It then suffices to upper bound $\|f_k - \mathcal{T}f_{k-1}\|_{\mu \times U}$.

$$\begin{aligned}
\|f_k - \mathcal{T}f_{k-1}\|_{\mu \times U}^2 &= \mathcal{L}_{\mu \times U}(f_k; f_{k-1}) - \mathcal{L}_{\mu \times U}(\mathcal{T}f_{k-1}; f_{k-1}) \quad (\mathcal{L} \text{ squared loss} + \mathcal{T}f_{k-1} \text{ Bayes optimal}) \\
&\leq \mathcal{L}_D(f_k; f_{k-1}) - \mathcal{L}_D(\mathcal{T}f_{k-1}; f_{k-1}) + 2\epsilon && (\mathcal{T}f_{k-1} \in \mathcal{F}) \\
&\leq 2\epsilon. && (f_k \text{ minimizes } \mathcal{L}_D(\cdot; f_{k-1}))
\end{aligned}$$

Note that the RHS does not depend on k , so we conclude that

$$\|f_k - Q^*\|_{\nu \times \pi} \leq \frac{1 - \gamma^k}{1 - \gamma} \sqrt{2|\mathcal{A}|C\epsilon} + \gamma^k \frac{R_{\max}}{1 - \gamma}.$$

Apply this to Equation (1) and we get

$$J(\pi^*) - J(\pi_{f_k}) \leq \frac{2}{1 - \gamma} \left(\frac{1 - \gamma^k}{1 - \gamma} \sqrt{2|\mathcal{A}|C\epsilon} + \gamma^k \frac{R_{\max}}{1 - \gamma} \right).$$

Extension: fast rate The previous bound should have $O(n^{-1/4})$ dependence on sample size $n := |D|$, because ϵ in bullet 6 should be $O(n^{-1/2})$ using Hoeffding's, and the final bound depends on $\sqrt{\epsilon}$. Here we exploit realizability to achieve fast rate so that the final bound is $O(n^{-1/2})$.

Define

$$Y(f; f') := (f(s, a) - r - \gamma V_{f'}(s'))^2 - ((\mathcal{T}f')(s, a) - r - \gamma V_{f'}(s'))^2.$$

Plug each $(s, a, r, s') \in D$ into $Y(f; f')$ and we get i.i.d. variables $Y_1(f; f'), Y_2(f; f'), \dots, Y_n(f; f')$ where $n = |D|$. It is easy to see that

$$\frac{1}{n} \sum_{i=1}^n Y_i(f; f') = \mathcal{L}_D(f; f') - \mathcal{L}_D(\mathcal{T}f'; f'),$$

so we only shift our objective \mathcal{L}_D by a f -independent constant. Our goal is to show that

$$\|\widehat{\mathcal{T}}_{\mathcal{F}} f' - \mathcal{T}f'\|_{\mu \times U}^2 \equiv \mathbb{E}[Y(\widehat{\mathcal{T}}_{\mathcal{F}} f'; f')] = O(1/n).$$

Note that this result can be directly plugged into the previous analysis by letting $f' = f_{k-1}$ (hence $\widehat{\mathcal{T}}_{\mathcal{F}} f' = f_k$), and we immediately obtain a final bound of $O(n^{-1/2})$.

To prove the result, first notice that $\forall f \in \mathcal{F}$,

$$\mathbb{E}[Y(f; f')] = \mathcal{L}_{\mu \times U}(f; f') - \mathcal{L}_{\mu \times U}(\mathcal{T}f'; f') = \|f - \mathcal{T}f'\|_{\mu \times U}^2,$$

thanks to realizability and squared loss. Next we bound variance of Y :

$$\begin{aligned}
\mathbb{V}[Y(f; f')] &\leq \mathbb{E}[Y(f; f')^2] \\
&= \mathbb{E} \left[\left((f(s, a) - r - \gamma V_{f'}(s'))^2 - ((\mathcal{T}f')(s, a) - r - \gamma V_{f'}(s'))^2 \right)^2 \right] \\
&= \mathbb{E} \left[(f(s, a) - (\mathcal{T}f')(s, a))^2 (f(s, a) + (\mathcal{T}f')(s, a) - 2r - 2\gamma V_{f'}(s'))^2 \right] \\
&\leq 4V_{\max}^2 \mathbb{E} \left[(f(s, a) - (\mathcal{T}f')(s, a))^2 \right] \\
&= 4V_{\max}^2 \|f - \mathcal{T}f'\|_{\mu \times U}^2 = 4V_{\max}^2 \mathbb{E}[Y(f; f')],
\end{aligned}$$

where $V_{\max} = R_{\max}/(1 - \gamma)$ is a constant.

Next we apply (one-sided) Bernstein's inequality (see [4]) and union bound over all $f \in \mathcal{F}$. Let $N = |\mathcal{F}|$. For any fixed f' , with probability at least $1 - \delta$, $\forall f \in \mathcal{F}$,

$$\begin{aligned}
\mathbb{E}[Y(f; f')] - \frac{1}{n} \sum_{i=1}^n Y_i(f; f') &\leq \sqrt{\frac{2\mathbb{V}[Y(f; f')] \log \frac{N}{\delta}}{n}} + \frac{4V_{\max}^2 \log \frac{N}{\delta}}{3n} \quad (Y_i \in [-V_{\max}^2, V_{\max}^2]) \\
&= \sqrt{\frac{8V_{\max}^2 \mathbb{E}[Y(f; f')] \log \frac{N}{\delta}}{n}} + \frac{4V_{\max}^2 \log \frac{N}{\delta}}{3n}.
\end{aligned}$$

Since $\widehat{\mathcal{T}}_{\mathcal{F}} f'$ minimizes $\mathcal{L}_D(\cdot; f')$, it also minimizes $\frac{1}{n} \sum_{i=1}^n Y_i(\cdot; f')$ because the two objectives only differ by a constant $\mathcal{L}_D(\mathcal{T}f'; f')$. Hence,

$$\frac{1}{n} \sum_{i=1}^n Y_i(\widehat{\mathcal{T}}_{\mathcal{F}} f'; f') \leq \frac{1}{n} \sum_{i=1}^n Y_i(\mathcal{T}f'; f') = 0.$$

Then,

$$\mathbb{E}[Y(\widehat{\mathcal{T}}_{\mathcal{F}} f'; f')] \leq \sqrt{\frac{8V_{\max}^2 \mathbb{E}[Y(\widehat{\mathcal{T}}_{\mathcal{F}} f'; f')] \log \frac{N}{\delta}}{n}} + \frac{4V_{\max}^2 \log \frac{N}{\delta}}{3n}.$$

Solving for the quadratic formula,

$$\mathbb{E}[Y(\widehat{\mathcal{T}}_{\mathcal{F}} f'; f')] \leq \left(\sqrt{2} + \sqrt{\frac{10}{3}} \right)^2 \frac{V_{\max}^2 \log \frac{N}{\delta}}{n}.$$

References

- [1] Rémi Munos. Error bounds for approximate policy iteration. In *International Conference on Machine Learning*, 2003.
- [2] András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 2008.
- [3] Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research*, 2008.
- [4] Sham Kakade. *Hoeffding, Chernoff, Bennet, and Bernstein Bounds*, 2011. <http://stat.wharton.upenn.edu/~skakade/courses/stat928/lectures/lecture06.pdf>.