

## A. Proofs of Main Theorems

### A.1. Full Algorithm of General FQE

**Algorithm 2** Fitted Q-Evaluation (Le et al., 2019)

**input** Datasets  $\mathcal{D} = \{\mathcal{D}_1, \dots, \mathcal{D}_K\}$ , target policy  $\pi$ , function class  $\mathcal{F}$ , initial state distribution  $\xi_0$ .

1: Initialize  $\widehat{Q}_{H+1}^\pi = 0$ .

2: **for**  $h = H, H-1, \dots, 1$  **do**

3:   Compute regression targets for any  $k \in [K], h' \in [H]$ :

$$y_{h,h'}^k = r_{h'}^k + \int_a \widehat{Q}_{h+1}^\pi(s_{h'+1}^k, a) \pi(a|s_{h'+1}^k) da.$$

4:   Build training set  $\{(s_{h'}^k, a_{h'}^k), y_{h,h'}^k\}_{k \in [K], h' \in [H]}$ .

5:   Solve a supervised learning problem:

$$\widehat{Q}_h^\pi = \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \frac{1}{K} \sum_{k=1}^K \frac{1}{H} \sum_{h'=1}^H (f(s_{h'}^k, a_{h'}^k) - y_{h,h'}^k)^2 + \lambda \rho(f) \right\},$$

where  $\rho(f)$  is a proper regularizer.

6: **end for**

**output**  $\widehat{v}^\pi = \int_s \int_a \widehat{Q}_1^\pi(s, a) \xi_1(s) \pi(a|s) ds da$ .

We restate the full algorithm of FQE in Algorithm 2. Here we simply assume the initial state distribution  $\xi_1$  is known. In practice, we always have the access to sample from  $\xi_1$  and thus we can approximate it by Monte Carlo sampling.

### A.2. Equivalence between FQE and model-based plug-in estimator

We show that the FQE in Algorithm 2 with linear function class  $\mathcal{F}$  is equivalent to a plug-in estimator. This equivalence is helpful to derive the asymptotic normality of FQE and bootstrapping FQE. Define

$$\widehat{M}_\pi = \widehat{\Sigma}^{-1} \sum_{n=1}^N \phi(s_n, a_n) \phi^\pi(s_{n+1})^\top, \widehat{R} = \widehat{\Sigma}^{-1} \sum_{n=1}^N r_n \phi(s_n, a_n), \widehat{\Sigma} = \sum_{n=1}^N \phi(s_n, a_n) \phi(s_n, a_n)^\top + \lambda I_d, \quad (\text{A.1})$$

where  $\phi^\pi(s) = \int_a \phi(s, a) \pi(a|s) da$ ,  $s_{N+1}$  is the terminal state and  $\lambda$  is the regularization parameter. Choosing  $\rho(f) = \lambda I$ , the FQE is equivalent to, for  $h = H, \dots, 1$ ,  $\widehat{Q}_h(s, a) = \phi(s, a)^\top \widehat{w}_h^\pi$  with

$$\begin{aligned} \widehat{w}_h^\pi &= \widehat{\Sigma}^{-1} \sum_{n=1}^N \phi(s_n, a_n) \left( r_n + \int_a \widehat{Q}_{h+1}^\pi(s_{n+1}, a) \pi(a|s_{n+1}) da \right) \\ &= \widehat{\Sigma}^{-1} \sum_{n=1}^N \phi(s_n, a_n) \left( r_n + \int_a \phi(s_{n+1}, a)^\top \widehat{w}_{h+1}^\pi \pi(a|s_{n+1}) da \right) \\ &= \widehat{\Sigma}^{-1} \sum_{n=1}^N \phi(s_n, a_n) r_n + \widehat{\Sigma}^{-1} \sum_{n=1}^N \phi(s_n, a_n) \phi^\pi(s_{n+1})^\top \widehat{w}_{h+1}^\pi \\ &= \widehat{R} + \widehat{M}_\pi \widehat{w}_{h+1}^\pi. \end{aligned}$$

This gives us a recursive form of  $\widehat{w}_h^\pi$ . Denoting  $\widehat{w}_{H+1}^\pi = 0$  and  $\nu_1^\pi = \mathbb{E}_{s \sim \xi_1, a \sim \pi(\cdot|s)}[\phi(s, a)]$ , the FQE estimator can be written into

$$\widehat{v}_\pi = \int_s \int_a \widehat{Q}_1(s, a) \xi_1 \pi(a|s) ds da = (\nu_1^\pi)^\top \widehat{w}_1^\pi = (\nu_1^\pi)^\top \sum_{h=0}^{H-1} (\widehat{M}_\pi)^h \widehat{R}. \quad (\text{A.2})$$

On the other hand, from Condition 4.1, there exists some  $w_r, w_h^\pi \in \mathbb{R}^d$  such that  $Q_h^\pi(\cdot, \cdot) = \phi(\cdot, \cdot)^\top w_h^\pi$  for each  $h \in [H]$  and  $r(\cdot, \cdot) = \phi(\cdot, \cdot)^\top w_r$  and there exists  $M_\pi \in \mathbb{R}^{d \times d}$  such that  $\phi(s, a)^\top M_\pi = \mathbb{E}[\phi^\pi(s')^\top | s, a]$ . From Bellman equation and Condition 4.1,

$$\begin{aligned} Q_h^\pi(s, a) &= r(s, a) + \mathbb{E} \left[ \int_a Q_{h+1}^\pi(s', a) \pi(a|s') da | s, a \right] \\ &= \phi(s, a)^\top w_r + \phi(s, a)^\top \mathbb{E}[\phi^\pi(s')^\top | s, a] w_{h+1}^\pi = \phi(s, a)^\top \left( w_r + M_\pi w_{h+1}^\pi \right) \\ &= \phi(s, a)^\top \sum_{h=0}^{H-h} (M_\pi)^h w_r. \end{aligned} \tag{A.3}$$

Therefore, the true scalar value function can be written as

$$v_\pi = \mathbb{E}_{s \sim \xi_1, a \sim \pi(\cdot|s)} \left[ Q_1^\pi(s, a) \right] = (\nu_1^\pi)^\top \sum_{h=0}^{H-1} (M_\pi)^h w_r,$$

which implies Eq. (A.2) is a plug-in estimator.

### A.3. Proof of Theorem 4.2: Asymptotic normality of FQE

Recall  $\nu_h^\pi = \mathbb{E}^\pi[\phi(x_h, a_h) | x_1 \sim \xi_1]$  and denote  $(\widehat{\nu}_h^\pi)^\top = (\nu_1^\pi)^\top (\widehat{M}_\pi)^{h-1}$ . We follow Lemma B.3 in Duan & Wang (2020) to decompose the error term into following three parts:

$$\sqrt{N}(v_\pi - \widehat{v}_\pi) = E_1 + E_2 + E_3,$$

where

$$\begin{aligned} E_1 &= \frac{1}{\sqrt{N}} \sum_{n=1}^N \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) \left( Q_h^\pi(s_n, a_n) - (r_n + V_{h+1}^\pi(s_{n+1})) \right), \\ E_2 &= \sum_{h=1}^H \left( N(\widehat{\nu}_h^\pi)^\top \widehat{\Sigma}^{-1} - (\nu_h^\pi)^\top \Sigma^{-1} \right) \left( \frac{1}{\sqrt{N}} \sum_{n=1}^N \phi(s_n, a_n) \left( Q_h^\pi(s_n, a_n) - (r_n + V_{h+1}^\pi(s_{n+1})) \right) \right), \\ E_3 &= \lambda \frac{1}{\sqrt{N}} \sum_{h=0}^H (\widehat{\nu}_h^\pi)^\top \widehat{\Sigma}^{-1} w_h^\pi. \end{aligned}$$

To prove the asymptotic normality of  $\sqrt{N}(v_\pi - \widehat{v}_\pi)$ , we will first prove the asymptotic normality of  $E_1$  and then show both  $E_1$  and  $E_2$  are asymptotically negligible.

For  $n = 1, 2, \dots, N$ , we denote

$$e_n = \frac{1}{\sqrt{N}} \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) \left( Q_h^\pi(s_n, a_n) - (r_{n+1} + V_{h+1}^\pi(s_{n+1})) \right).$$

Then  $E_1 = \sum_{n=1}^N e_n$ . Define a filtration  $\{\mathcal{F}_n\}_{n=1, \dots, N}$  with  $\mathcal{F}_n$  generated by  $(s_1, a_1, s_2), \dots, (s_{n-1}, a_{n-1}, s_n)$  and  $(s_n, a_n)$ . From the definition of value function, it is easy to see  $\mathbb{E}[e_n | \mathcal{F}_n] = 0$  that implies that  $\{e_n\}_{n \in [N]}$  is a martingale difference sequence. To show the asymptotic normality, we use the following martingale central limit theorem for triangular arrays.

**Theorem A.1** (Martingale CLT, Corollary 2.8 in (McLeish et al., 1974)). Let  $\{X_{mn}; n = 1, \dots, k_m\}$  be a martingale difference array (row-wise) on the probability triple  $(\Omega, \mathcal{F}, P)$ . Suppose  $X_{mn}$  satisfy the following two conditions:

$$\max_{1 \leq n \leq k_m} |X_{mn}| \xrightarrow{p} 0, \text{ and } \sum_{n=1}^{k_m} X_{mn}^2 \xrightarrow{p} \sigma^2,$$

for  $k_m \rightarrow \infty$ . Then  $\sum_{n=1}^{k_m} X_{mn} \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ .

Recall that the variance  $\sigma^2$  is defined as

$$\sigma^2 = \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} \Omega_{h,h} \Sigma^{-1} \nu_h^\pi + 2 \sum_{h_1 < h_2} (\nu_{h_1}^\pi)^\top \Sigma^{-1} \Omega_{h_1, h_2} \Sigma^{-1} \nu_{h_2}^\pi, \quad (\text{A.4})$$

and for any  $h_1 \in [H], h_2 \in [H]$ ,

$$\Omega_{h_1, h_2} = \mathbb{E} \left[ \frac{1}{H} \sum_{h'=1}^H \phi(s_{h'}, a_{h'}^1) \phi(s_{h'}, a_{h'}^1)^\top \varepsilon_{h_1, h'}^1 \varepsilon_{h_2, h'}^1 \right],$$

where  $\varepsilon_{h_1, h'}^1 = Q_{h_1}^\pi(s_{h'}, a_{h'}^1) - (r_{h'}^1 + V_{h_1+1}^\pi(s_{h'+1}^1))$ . To apply Theorem A.1, we let  $k_m = N$ ,  $X_{mn} = e_n$  and we need to verify the following two conditions:

$$\max_{1 \leq n \leq N} \left| \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} \left( \frac{1}{\sqrt{N}} \phi(s_n, a_n) \left( Q_h^\pi(s_n, a_n) - (r_{n+1} + V_{h+1}^\pi(s_{n+1})) \right) \right) \right| \xrightarrow{P} 0, \text{ as } N \rightarrow \infty, \quad (\text{A.5})$$

and

$$\sum_{n=1}^N \left( \frac{1}{\sqrt{N}} \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) \left( Q_h^\pi(s_n, a_n) - (r_{n+1} + V_{h+1}^\pi(s_{n+1})) \right) \right)^2 \xrightarrow{P} \sigma^2, \text{ as } N \rightarrow \infty. \quad (\text{A.6})$$

**Verify Condition A.5:** Since  $r \in [0, 1]$ , we have  $r_n + V_{h+1}^\pi(s_{n+1}) \in [0, H - h]$ . For any  $n \in [N]$ , we have

$$\begin{aligned} & \left| \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} \left( \frac{1}{\sqrt{N}} \phi(s_n, a_n) \left( Q_h^\pi(s_n, a_n) - (r_{n+1} + V_{h+1}^\pi(s_{n+1})) \right) \right) \right| \\ & \leq \frac{1}{\sqrt{N}} \sum_{h=1}^H \left| (\nu_h^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) \right| \left| Q_h^\pi(s_n, a_n) - (r_n + V_{h+1}^\pi(s_{n+1})) \right| \\ & \leq \frac{1}{\sqrt{N}} \sum_{h=1}^H (H - h + 1) \left| (\nu_h^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) \right|. \end{aligned}$$

Note that  $(\nu_h^\pi)^\top \Sigma^{-1} \phi(s_n, a_n)$  is independent of  $N$ . Then for fixed  $d, H$ , Condition A.5 is satisfied when  $N \rightarrow \infty$ .

**Verify Condition A.6:** Recall the definition of  $\sigma^2$  in Eq. (A.4) and let  $\sigma^2 = \sigma_1^2 + \sigma_2^2$  for

$$\begin{aligned} \sigma_1^2 &= \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} \Omega_{h,h} \Sigma^{-1} \nu_h^\pi, \\ \sigma_2^2 &= 2 \sum_{h_1 < h_2} (\nu_{h_1}^\pi)^\top \Sigma^{-1} \Omega_{h_1, h_2} \Sigma^{-1} \nu_{h_2}^\pi. \end{aligned}$$

Using the following decomposition,

$$\begin{aligned} & \sum_{n=1}^N \left( \frac{1}{\sqrt{N}} \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) \left( Q_h^\pi(s_n, a_n) - (r_{n+1} + V_{h+1}^\pi(s_{n+1})) \right) \right)^2 \\ &= \sum_{n=1}^N \frac{1}{N} \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) \phi(s_n, a_n)^\top \Sigma^{-1} \nu_h^\pi \left( Q_h^\pi(s_n, a_n) - (r_n + V_{h+1}^\pi(s_{n+1})) \right)^2 \\ &+ \sum_{n=1}^N \frac{1}{N} 2 \sum_{h_1 < h_2} (\nu_{h_1}^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) (\nu_{h_2}^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) \\ &\cdot \left( Q_{h_1}^\pi(s_n, a_n) - (r_{n+1} + V_{h_1+1}^\pi(s_{n+1})) \right) \left( Q_{h_2}^\pi(s_n, a_n) - (r_{n+1} + V_{h_2+1}^\pi(s_{n+1})) \right). \end{aligned}$$

We denote the first term as  $I_1$ , the second term as  $I_2$  and separately bound  $I_1 - \sigma_1^2$  and  $I_2 - \sigma_2^2$  as follows:

- We rewrite  $I_1$  in terms of episodes as

$$I_1 = \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1/2} \left( \frac{1}{K} \sum_{k=1}^K \frac{1}{H} \sum_{h'=1}^H \Sigma^{-1/2} \phi(s_{h'}^k, a_{h'}^k) \phi(s_{h'}^k, a_{h'}^k)^\top (\varepsilon_{hh'}^k)^2 \Sigma^{-1/2} \right) \Sigma^{-1/2} \nu_h^\pi.$$

Moreover, denote

$$Z_h = \Sigma^{-1/2} \mathbb{E} \left[ \frac{1}{H} \sum_{h'=1}^H \phi(s_{h'}^1, a_{h'}^1) \phi(s_{h'}^1, a_{h'}^1)^\top (\varepsilon_{hh'}^1)^2 \right] \Sigma^{-1/2} \in \mathbb{R}^{d \times d}.$$

Then we have

$$\begin{aligned} |I_1 - \sigma_1^2| &= \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1/2} \left( \frac{1}{K} \sum_{k=1}^K \frac{1}{H} \sum_{h'=1}^H \Sigma^{-1/2} \phi(s_{h'}^k, a_{h'}^k) \phi(s_{h'}^k, a_{h'}^k)^\top (\varepsilon_{hh'}^k)^2 \Sigma^{-1/2} - Z_h \right) \Sigma^{-1/2} \nu_h^\pi \\ &\leq \sum_{h=1}^H \left\| (\nu_h^\pi)^\top \Sigma^{-1/2} \right\|_2 \left\| \frac{1}{K} \sum_{k=1}^K \left( \frac{1}{H} \sum_{h'=1}^H \Sigma^{-1/2} \phi(s_{h'}^k, a_{h'}^k) \phi(s_{h'}^k, a_{h'}^k)^\top (\varepsilon_{hh'}^k)^2 \Sigma^{-1/2} - Z_h \right) \right\|_2, \end{aligned}$$

where the last inequality is from Cauchy–Schwarz inequality. From Lemma B.7, we reach  $I_1 \xrightarrow{P} \sigma_1^2$  as  $K \rightarrow \infty$ .

- We rewrite  $I_2$  as

$$I_2 = 2 \sum_{h_1 < h_2} (\nu_{h_1}^\pi)^\top \Sigma^{-1/2} \left( \frac{1}{K} \sum_{k=1}^K \frac{1}{H} \sum_{h'=1}^H \Sigma^{-1/2} \phi(s_{h'}^k, a_{h'}^k) \phi(s_{h'}^k, a_{h'}^k)^\top \varepsilon_{h_1 h'}^k \varepsilon_{h_2 h'}^k \Sigma^{-1/2} \right) \Sigma^{-1/2} \nu_{h_2}^\pi,$$

and denote

$$Z_{h_1 h_2} = \Sigma^{-1/2} \mathbb{E} \left[ \frac{1}{H} \sum_{h'=1}^H \phi(s_{h'}^1, a_{h'}^1) \phi(s_{h'}^1, a_{h'}^1)^\top \varepsilon_{h_1 h'}^1 \varepsilon_{h_2 h'}^1 \right] \Sigma^{-1/2} \in \mathbb{R}^{d \times d}.$$

Then we have

$$\begin{aligned} |I_2 - \sigma_2^2| &= 2 \sum_{h_1 < h_2} (\nu_{h_1}^\pi)^\top \Sigma^{-1/2} \left( \frac{1}{K} \sum_{k=1}^K \frac{1}{H} \sum_{h'=1}^H \Sigma^{-1/2} \phi(s_{h'}^k, a_{h'}^k) \phi(s_{h'}^k, a_{h'}^k)^\top \varepsilon_{h_1 h'}^k \varepsilon_{h_2 h'}^k \Sigma^{-1/2} - Z_{h_1 h_2} \right) \Sigma^{-1/2} \nu_{h_2}^\pi \\ &\leq 2 \sum_{h_1 < h_2} \left\| (\nu_{h_1}^\pi)^\top \Sigma^{-1/2} \right\|_2 \left\| (\nu_{h_2}^\pi)^\top \Sigma^{-1/2} \right\|_2 \left\| \frac{1}{K} \sum_{k=1}^K \left( \frac{1}{H} \sum_{h'=1}^H \Sigma^{-1/2} \phi(s_{h'}^k, a_{h'}^k) \phi(s_{h'}^k, a_{h'}^k)^\top \varepsilon_{h_1 h'}^k \varepsilon_{h_2 h'}^k \Sigma^{-1/2} - Z_{h_1 h_2} \right) \right\|_2. \end{aligned}$$

From Lemma B.7, we reach  $I_2 \xrightarrow{P} \sigma_2^2$  as  $K \rightarrow \infty$ .

Putting the above two steps together, we have verified Condition A.6. Then applying Theorem A.1 we obtain that  $E_1 \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ .

On the other hand, according to Lemmas B.6, B.10 in (Duan & Wang, 2020),

$$\begin{aligned} |E_2| &\leq 15 \sqrt{(\nu_0^\pi)^\top (\Sigma^\pi)^{-1} \nu_0^\pi} \cdot \|(\Sigma^\pi)^{1/2} \Sigma^{-1/2}\|_2 \cdot \sqrt{C_1 \kappa_1 (2 + \kappa_2)} \cdot \frac{\ln(8dH/\delta) dH^{3.5}}{\sqrt{N}} \\ |E_3| &\leq \sqrt{(\nu_0^\pi)^\top (\Sigma^\pi)^{-1} \nu_0^\pi} \cdot \|(\Sigma^\pi)^{1/2} \Sigma^{-1/2}\|_2 \cdot \frac{5 \ln(8dH/\delta) C_1 dH^2}{\sqrt{N}}, \end{aligned}$$

with probability at least  $1 - \delta$  and  $\kappa_1, \kappa_2$  are some problem-dependent constants that do not depend on  $N$ . When  $N \rightarrow \infty$ , both  $|E_2|, |E_3|$  converge in probability to 0. By Slutsky's theorem, we have proven the asymptotic normality of  $\sqrt{N}(v_\pi - \hat{v}_\pi)$ . ■

#### A.4. Proof of Theorem 4.5: Efficiency bound

**Influence function.** Recall that our dataset  $\mathcal{D}$  consists of  $K$  *i.i.d.* trajectories, each of which has length  $H$ . Denote

$$\boldsymbol{\tau} := (s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_H, a_H, r_H, s_{H+1}).$$

For simplicity, we assume that the reward  $r_h$  is deterministic given  $(s_h, a_h)$ , *i.e.*  $r_h = r(s_h, a_h)$  for some reward function  $r$ .

The distribution of  $\boldsymbol{\tau}$  is given by

$$\begin{aligned} \mathcal{P}(d\boldsymbol{\tau}) &= \bar{\xi}_1(ds_1, da_1) \mathcal{P}(ds_2 | s_1, a_1) \bar{\pi}(da_2 | s_2) \mathcal{P}(ds_3 | s_2, a_2) \\ &\quad \dots \mathcal{P}(ds_H | s_{H-1}, a_{H-1}) \bar{\pi}(da_H | s_H) \mathcal{P}(ds_{H+1} | s_H, a_H). \end{aligned}$$

Define  $\mathcal{P}_\eta := \mathcal{P} + \eta \Delta \mathcal{P}$  where  $\Delta \mathcal{P}$  satisfies

$$(\Delta \mathcal{P})\mathcal{F} \subseteq \mathcal{F}$$

under condition 4.1. Denote score functions

$$g(\boldsymbol{\tau}) := \frac{\partial}{\partial \eta} \log \mathcal{P}_\eta(d\boldsymbol{\tau}) \quad \text{and} \quad g(s' | s, a) := \frac{\partial}{\partial \eta} \log \mathcal{P}_\eta(ds' | s, a).$$

Note that

$$g(\boldsymbol{\tau}) = \sum_{h=1}^H g(s_{h+1} | s_h, a_h).$$

We consider the pointwise estimation. The objective function  $\psi_{\xi_1}$  is defined as

$$\psi_{\xi_1}(\mathcal{P}_\eta) := \mathbb{E} \left[ \sum_{h=1}^H r_\eta(s_h, a_h) \mid (s_1, a_1) \sim \xi_1, \mathcal{P}_\eta, \pi \right].$$

We calculate the derivative  $\frac{\partial}{\partial \eta} \psi_{\xi_1}(\mathcal{P}_\eta)$  and have

$$\begin{aligned} \frac{\partial}{\partial \eta} \psi_{\xi_1}(\mathcal{P}_\eta) &= \frac{\partial}{\partial \eta} \left[ \sum_{h=1}^H \int_{(\mathcal{S} \times \mathcal{A})^h} r(s_h, a_h) \xi_1(ds_1, da_1) \prod_{j=1}^{h-1} \mathcal{P}_\eta(ds_{j+1} | s_j, a_j) \pi(da_{j+1} | s_{j+1}) \right] \\ &= \sum_{h=1}^H \int_{(\mathcal{S} \times \mathcal{A})^h} r(s_h, a_h) \left( \sum_{j=1}^{h-1} g(s_{j+1} | s_j, a_j) \right) \xi_1(ds_1, da_1) \prod_{j=1}^{h-1} \mathcal{P}_\eta(ds_{j+1} | s_j, a_j) \pi(da_{j+1} | s_{j+1}). \end{aligned}$$

By using Q-functions  $Q_{\eta,j}^\pi(s_j, a_j) := \mathbb{E}[\sum_{h=j}^H r_\eta(s_h, a_h) | (s_j, a_j), \mathcal{P}_\eta, \pi]$  for  $j = 1, 2, \dots, H$ ,  $Q_{\eta,H+1} := 0$ , we find that

$$\begin{aligned} \frac{\partial}{\partial \eta} \psi_{\xi_1}(\mathcal{P}_\eta) &= \int_{(\mathcal{S} \times \mathcal{A})^H} \sum_{j=1}^{H-1} g(s_{j+1} | s_j, a_j) \sum_{h=j+1}^H r_\eta(s_h, a_h) \xi_1(ds_1, da_1) \prod_{i=1}^{H-1} \mathcal{P}_\eta(ds_{i+1} | s_i, a_i) \pi(da_{i+1} | s_{i+1}) \\ &= \int_{(\mathcal{S} \times \mathcal{A})^H} \sum_{j=1}^{H-1} g(s_{j+1} | s_j, a_j) \xi_1(ds_1, da_1) \prod_{i=1}^j \mathcal{P}_\eta(ds_{i+1} | s_i, a_i) \pi(da_{i+1} | s_{i+1}) \\ &\quad \cdot \left( \sum_{h=j+1}^H r_\eta(s_h, a_h) \prod_{i=j+1}^{H-1} \mathcal{P}_\eta(ds_{i+1} | s_i, a_i) \pi(da_{i+1} | s_{i+1}) \right) \\ &= \sum_{j=1}^{H-1} \int_{(\mathcal{S} \times \mathcal{A})^{j+1}} g(s_{j+1} | s_j, a_j) Q_{\eta,j+1}^\pi(s_{j+1}, a_{j+1}) \xi_1(ds_1, da_1) \prod_{i=1}^j \mathcal{P}_\eta(ds_{i+1} | s_i, a_i) \pi(da_{i+1} | s_{i+1}) \\ &= \sum_{h=1}^H \mathbb{E}[g(s_{h+1} | s_h, a_h) V_{\eta,h+1}^\pi(s_{h+1}) | (s_1, a_1) \sim \xi_1, \mathcal{P}_\eta, \pi]. \end{aligned}$$

It follows that

$$\frac{\partial}{\partial \eta} \psi_{\xi_1}(\mathcal{P}_\eta) \Big|_{\eta=0} = \mathbb{E} \left[ \sum_{h=1}^H g(s_{h+1} | s_h, a_h) V_{h+1}^\pi(s_{h+1}) \mid (s_1, a_1) \sim \xi_1, \mathcal{P}, \pi \right].$$

Define  $w_h(s, a) := \phi(s, a)^\top \Sigma^{-1} \nu_h^\pi = \phi(s, a)^\top \Sigma^{-1} \mathbb{E}[\phi(s_h, a_h) \mid (s_1, a_1) \sim \xi_1, \mathcal{P}, \pi]$  for  $h = 1, 2, \dots, H$ . Let  $\mathbb{H}$  For any  $f \in \mathbb{H}$  with  $f(s, a) = \phi(s, a)^\top w_f$ , we have

$$\begin{aligned} \mathbb{E}[f(s_h, a_h) \mid (s_1, a_1) \sim \xi_1, \mathcal{P}, \pi] &= \mathbb{E}[\phi(s_h, a_h)^\top w_f \mid (s_1, a_1) \sim \xi_1, \mathcal{P}, \pi] \\ &= \mathbb{E}[\phi(s_h, a_h)^\top \Sigma^{-1} \mathbb{E}_{(s,a) \sim \bar{\mu}}[\phi(s, a) \phi(s, a)^\top] w_f \mid (s_1, a_1) \sim \xi_1, \mathcal{P}, \pi] \\ &= \mathbb{E}_{(s,a) \sim \bar{\mu}} \left[ \mathbb{E}[\phi(s_h, a_h) \mid (s_1, a_1) \sim \xi_1, \mathcal{P}, \pi]^\top \Sigma^{-1} \phi(s, a) \phi(s, a)^\top w_f \right] \\ &= \mathbb{E}_{(s,a) \sim \bar{\mu}} [w_h(s, a) f(s, a)], \end{aligned}$$

where  $\bar{\mu}$  is the distribution of dataset  $\mathcal{D}$ . Since the mapping  $(s, a) \mapsto \mathbb{E}[g(s' | s, a) V_h^\pi(s') \mid s, a]$  belongs to  $\mathbb{H}$ , therefore,

$$\frac{\partial}{\partial \eta} \psi_{\xi_1}(\mathcal{P}_\eta) \Big|_{\eta=0} = \mathbb{E}_{(s,a) \sim \bar{\mu}} \left[ \sum_{h=1}^H w_h(s, a) g(s' | s, a) V_{h+1}^\pi(s') \right].$$

Note that  $\mathbb{E}[g(s' | s, a) \mid s, a] = 0$ , therefore,

$$\frac{\partial}{\partial \eta} \psi_{\xi_1}(\mathcal{P}_\eta) \Big|_{\eta=0} = \mathbb{E}_{(s,a) \sim \bar{\mu}} \left[ \sum_{h=1}^H w_h(s, a) g(s' | s, a) \left( V_{h+1}^\pi(s') - \mathbb{E}[V_{h+1}^\pi(s') \mid s, a] \right) \right].$$

By definition of  $\mu$ , we have

$$\begin{aligned} &\frac{\partial}{\partial \eta} \psi_{\xi_1}(\mathcal{P}_\eta) \Big|_{\eta=0} \\ &= \frac{1}{H} \sum_{j=1}^H \mathbb{E} \left[ \sum_{h=1}^H w_h(s_j, a_j) g(s_{j+1} | s_j, a_j) \left( V_{h+1}^\pi(s_{j+1}) - \mathbb{E}[V_{h+1}^\pi(s_{j+1}) \mid s_j, a_j] \right) \mid (s_1, a_1) \sim \bar{\xi}_1, \mathcal{P}, \bar{\pi} \right]. \end{aligned}$$

We use the property  $\mathbb{E}[g(s' | s, a) \mid s, a] = 0$  again and derive that

$$\begin{aligned} &\frac{\partial}{\partial \eta} \psi_{\xi_1}(\mathcal{P}_\eta) \Big|_{\eta=0} \\ &= \frac{1}{H} \sum_{j=1}^H \mathbb{E} \left[ \sum_{h=1}^H w_h(s_j, a_j) \left( \sum_{l=1}^H g(s_{l+1} | s_l, a_l) \right) \left( V_{h+1}^\pi(s_{j+1}) - \mathbb{E}[V_{h+1}^\pi(s_{j+1}) \mid s_j, a_j] \right) \mid (s_1, a_1) \sim \bar{\xi}_1, \mathcal{P}, \bar{\pi} \right] \\ &= \frac{1}{H} \mathbb{E} \left[ g(\boldsymbol{\tau}) \sum_{h=1}^H \sum_{j=1}^H w_h(s_j, a_j) \left( V_{h+1}^\pi(s_{j+1}) - \mathbb{E}[V_{h+1}^\pi(s_{j+1}) \mid s_j, a_j] \right) \mid (s_1, a_1) \sim \bar{\xi}_1, \mathcal{P}, \bar{\pi} \right]. \end{aligned}$$

We can conclude that

$$\dot{\psi}_{\mathcal{P}}(\boldsymbol{\tau}) := \frac{1}{H} \sum_{h=1}^H \sum_{h'=1}^H w_{h'}(s_h, a_h) \left( V_{h'+1}^\pi(s_{h+1}) - \mathbb{E}[V_{h'+1}^\pi(s_{h+1}) \mid s_h, a_h] \right),$$

is an influence function.

**Efficiency bound.** For notational convenience, we take shorthands

$$q(s, a, s') := \sum_{h=1}^H w_h(s, a) \left( V_{h+1}^\pi(s') - \mathbb{E}[V_{h+1}^\pi(s') \mid s, a] \right),$$

and rewrite

$$\dot{\psi}_{\mathcal{P}}(\boldsymbol{\tau}) = \frac{1}{H} \sum_{h=1}^H q(s_h, a_h, s_{h+1}).$$

Since  $\mathbb{E}[q(s, a, s') \mid s, a] = 0$ , we find that

$$\mathbb{E}[\dot{\psi}_{\mathcal{P}}^2(\boldsymbol{\tau})] = \frac{1}{H^2} \mathbb{E} \left[ \left( \sum_{h=1}^H q(s_h, a_h, s_{h+1}) \right)^2 \mid \bar{\xi}_1, \mathcal{P}, \bar{\pi} \right] = \frac{1}{H^2} \sum_{h=1}^H \mathbb{E} \left[ q^2(s_h, a_h, s_{h+1}) \mid \bar{\xi}_1, \mathcal{P}, \bar{\pi} \right].$$

It follows that

$$\begin{aligned} \mathbb{E}[\dot{\psi}_{\mathcal{P}}^2(\boldsymbol{\tau})] &= \frac{1}{H^2} \mathbb{E}_{(s,a) \sim \bar{\mu}} \left[ \mathbb{E}[q^2(s, a, s') \mid s, a] \right] = \frac{1}{H^2} \mathbb{E}_{(s,a) \sim \bar{\mu}} \left[ \mathbb{E}[q^2(s, a, s') \mid s, a] \right] \\ &= \frac{1}{H^2} \mathbb{E}_{(s,a) \sim \bar{\mu}} \left[ \left( \phi(s, a)^\top \Sigma^{-1} \sum_{h=1}^H \left( V_{h+1}^\pi(s') - \mathbb{E}[V_{h+1}^\pi(s') \mid s, a] \right) \nu_h^\pi \right)^2 \right], \end{aligned}$$

which coincides with the asymptotic variance of OPE estimator defined in (4.2). ■

### A.5. Proof of Theorem 5.1: Distributional consistency of bootstrapping FQE

In order to simplify the derivation, we assume  $\lambda = 0$  and the empirical covariance matrix  $\sum_{n=1}^N \phi(s_n, a_n) \phi(s_n, a_n)^\top$  is invertible in this section since the effect of  $\lambda$  is asymptotically negligible. For a matrix  $A \in \mathbb{R}^{m \times n}$ , suppose the  $\text{vec}$  operator stacks the column of a matrix such that  $\text{vec}(A) \in \mathbb{R}^{mn \times 1}$ . We use the equivalence form of FQE in Eq. (A.2) such that

$$\hat{v}_\pi = (\nu_1^\pi)^\top \sum_{h=0}^{H-1} (\widehat{M}_\pi)^h \widehat{R}.$$

$\widehat{M}_\pi$  can be viewed as the solution of the following multivariate linear regression:

$$\phi^\pi(s_{n+1})^\top = \phi(s_n, a_n)^\top M_\pi + \eta_n,$$

where  $\eta_n = \phi^\pi(s_{n+1})^\top - \phi(s_n, a_n)^\top M_\pi$ . We first derive the asymptotic distribution of  $\sqrt{N} \text{vec}(\widehat{M}_\pi - M_\pi)$  that follows:

$$\begin{aligned} \sqrt{N} \text{vec}(\widehat{M}_\pi - M_\pi) &= \text{vec} \left( \sqrt{N} \widehat{\Sigma}^{-1} \sum_{n=1}^N \phi(s_n, a_n) \left( \phi^\pi(s'_n)^\top - \phi(s_n, a_n)^\top M_\pi \right) \right) \\ &= (N \widehat{\Sigma}^{-1} \otimes I_d) \frac{1}{\sqrt{K}} \sum_{k=1}^K \text{vec} \left( \frac{1}{\sqrt{H}} \sum_{h=1}^H \phi(s_h^k, a_h^k) \left( \phi^\pi(s_h^{k'})^\top - \phi(s_h^k, a_h^k)^\top M_\pi \right) \right), \end{aligned} \quad (\text{A.7})$$

where  $\otimes$  is kronecker product. Define  $\xi_h^k = \phi^\pi(s_{h+1}^k)^\top - \phi(s_h^k, a_h^k)^\top M_\pi$ . From the definition of  $M_\pi$ , it is easy to see

$$\mathbb{E}[\phi^\pi(s_{h+1}^k)^\top \mid s_h^k, a_h^k] = \int_s \mathbb{P}(s \mid s_h^k, a_h^k) \int_a \pi(a \mid s) \phi(s, a)^\top da ds = \phi(s_h^k, a_h^k)^\top M_\pi.$$

Again with martingale central limit theorem and independence between each episode, we have as  $K \rightarrow \infty$ ,

$$\frac{1}{\sqrt{K}} \sum_{k=1}^K \text{vec} \left( \frac{1}{\sqrt{H}} \sum_{h=1}^H \phi(s_h^k, a_h^k) \xi_h^k \right) \xrightarrow{d} N(0, \Delta), \quad (\text{A.8})$$

where  $\Delta \in \mathbb{R}^{d^2 \times d^2}$  is the covariance matrix defined as: for  $j, k \in [d^2]$

$$\Delta_{jk} = \mathbb{E} \left[ \text{vec} \left( \frac{1}{\sqrt{H}} \sum_{h=1}^H \phi(s_h^k, a_h^k) \xi_h^k \right)_j \text{vec} \left( \frac{1}{\sqrt{H}} \sum_{h=1}^H \phi(s_h^k, a_h^k) \xi_h^k \right)_k \right]. \quad (\text{A.9})$$

Next we start to derive the conditional bootstrap asymptotic distribution. For notation simplicity, denote  $\phi_{hk} = \phi(s_h^k, a_h^k)$  and  $y_{hk} = \phi^\pi(s_{h+1}^k)^\top$ . We rewrite the dataset combined with feature map  $\phi(\cdot, \cdot)$  such that  $\mathcal{D}_k = \{\phi_{hk}, y_{hk}, r_{hk}\}_{h=1}^H$ . Recall that we bootstrap  $\mathcal{D}$  by episodes such that each episode is sampled with replacement to form the starred data  $\mathcal{D}_k^* = \{\phi_{hk}^*, y_{hk}^*, r_{hk}^*\}_{h=1}^H$  for  $k \in [K]$ . More specifically,

$$\phi_{hk}^* = \sum_{k=1}^K W_k^* \phi_{hk}, y_{hk}^* = \sum_{k=1}^K W_k^* y_{hk}, r_{hk}^* = \sum_{k=1}^K W_k^* r_{hk},$$

where  $W^* = (W_1^*, \dots, W_K^*)$  is the bootstrap weight. For example,  $W^*$  could be a multinomial random vector with parameters  $(K; K^{-1}, \dots, K^{-1})$  that forms the standard nonparametric bootstrap. Note that for different  $h \in [H]$ , they have the same bootstrap weight and given the original samples  $\mathcal{D}_1, \dots, \mathcal{D}_K$ , the resampled vectors are independent. Define the corresponding starred quantity  $\widehat{M}_\pi^*, \widehat{R}^*$  as

$$\widehat{M}_\pi^* = \widehat{\Sigma}^{*-1} \sum_{k=1}^K \sum_{h=1}^H \phi_{hk}^* y_{hk}^*, \widehat{R}^* = \widehat{\Sigma}^{*-1} \sum_{k=1}^K \sum_{h=1}^H r_{hk}^* \phi_{hk}^*,$$

where

$$\widehat{\Sigma}^* = \sum_{k=1}^K \sum_{h=1}^H \phi_{hk}^* \phi_{hk}^{*\top}.$$

We will derive the asymptotic distribution of  $\sqrt{N}(\text{vec}(\widehat{M}_\pi^* - \widehat{M}_\pi))$  by using the following decomposition:

$$\begin{aligned} \sqrt{N} \text{vec}(\widehat{M}_\pi^* - \widehat{M}_\pi) &= \sqrt{N} \text{vec}\left(\widehat{\Sigma}^{*-1} \sum_{k=1}^K \sum_{h=1}^H \phi_{hk}^* (y_{hk} - \phi_{hk}^* \widehat{M}_\pi)\right) \\ &= (N \widehat{\Sigma}^{*-1} \otimes I_d) \text{vec}\left(\frac{1}{\sqrt{K}} \sum_{k=1}^K \frac{1}{\sqrt{H}} \sum_{h=1}^H \phi_{hk}^* (y_{hk} - \phi_{hk}^* \widehat{M}_\pi)\right). \end{aligned}$$

We denote

$$Z = \frac{1}{\sqrt{K}} \sum_{k=1}^K \frac{1}{\sqrt{H}} \sum_{h=1}^H \phi_{hk} (y_{hk} - \phi_{hk} M_\pi), Z^* = \frac{1}{\sqrt{K}} \sum_{k=1}^K \frac{1}{\sqrt{H}} \sum_{h=1}^H \phi_{hk}^* (y_{hk}^* - \phi_{hk}^* \widehat{M}_\pi).$$

Both  $Z$  and  $Z^*$  are the sum of independent  $d \times d$  random matrices. We prove the bootstrap consistency using the Mallows metric as a central tool. The Mallows metric, relative to the Euclidean norm  $\|\cdot\|$ , for two probability measures  $\mu, \nu$  in  $\mathbb{R}^d$  is defined as

$$\Lambda_l(\mu, \nu) = \inf_{U \sim \mu, V \sim \nu} \mathbb{E}^{1/l}(\|U - V\|^l),$$

where  $U$  and  $V$  are two random vectors that  $U$  has law  $\mu$  and  $V$  has law  $\nu$ . For random variables  $U, V$ , we sometimes write  $\Lambda_l(U, V)$  as the  $\Lambda_l$ -distance between the laws of  $U$  and  $V$ . We refer [Bickel & Freedman \(1981\)](#); [Freedman et al. \(1981\)](#) for more details about the properties of Mallows metric. Suppose the common distribution of original  $K$  episodes  $\{\mathcal{D}_1, \dots, \mathcal{D}_K$  is  $\mu$  and their empirical distribution is  $\mu_K$ . Both  $\mu$  and  $\mu_K$  are probability in  $\mathbb{R}^{2Hd+H}$ . From [Lemma B.1](#), we know that  $\Lambda_4(\mu_K, \mu) \rightarrow 0$  a.e. as  $K \rightarrow \infty$ .

- **Step 1.** We prove  $\widehat{\Sigma}^*/N$  converges in conditional probability to  $\Sigma$ . From the bootstrap design,  $\frac{1}{H} \sum_{h=1}^H \phi_{kh}^* \phi_{kh}^{*\top}$  is independent of  $\frac{1}{H} \sum_{h=1}^H \phi_{k'h}^* \phi_{k'h}^{*\top}$  for any  $k \neq k'$ . According to [Lemma B.3](#), we have

$$\begin{aligned} \Lambda_1\left(\sum_{k=1}^K \frac{1}{H} \sum_{h=1}^H \phi_{kh}^* \phi_{kh}^{*\top}, \sum_{k=1}^K \frac{1}{H} \sum_{h=1}^H \phi_{kh} \phi_{kh}^\top\right) &\leq \sum_{k=1}^K \Lambda_1\left(\frac{1}{H} \sum_{h=1}^H \phi_{kh}^* \phi_{kh}^{*\top}, \frac{1}{H} \sum_{h=1}^H \phi_{kh} \phi_{kh}^\top\right) \\ &= K \Lambda_1\left(\frac{1}{H} \sum_{h=1}^H \phi_{kh}^* \phi_{kh}^{*\top}, \frac{1}{H} \sum_{h=1}^H \phi_{kh} \phi_{kh}^\top\right). \end{aligned}$$



Both sides of the above inequality are random variables such that the distance is computed between the conditional distribution of the starred quantity and the unconditional distribution of the unstarred quantity. Define a mapping  $f : \mathbb{R}^{Hd} \rightarrow \mathbb{R}^{d \times d}$  such that for any  $x_1, \dots, x_H \in \mathbb{R}^d$ ,

$$f(x_1, \dots, x_H) = \frac{1}{H} \sum_{h=1}^H x_h x_h^\top.$$

From Lemma B.2 with  $f$ , we have as  $K$  goes to infinity

$$\Lambda_1 \left( \frac{1}{H} \sum_{h=1}^H \phi_{kh}^* \phi_{kh}^{*\top}, \frac{1}{H} \sum_{h=1}^H \phi_{kh} \phi_{kh}^\top \right) \rightarrow 0.$$

This implies the conditional law of  $\frac{1}{H} \sum_{h=1}^H \phi_{kh}^* \phi_{kh}^{*\top}$  is close to the unconditional law of  $\frac{1}{H} \sum_{h=1}^H \phi_{kh} \phi_{kh}^\top$ . By the law of large numbers:

$$\frac{1}{K} \sum_{k=1}^K \frac{1}{H} \sum_{h=1}^H \phi_{kh} \phi_{kh}^\top \xrightarrow{p} \Sigma. \quad (\text{A.10})$$

This further implies the conditional on  $\mathcal{D}$ , we have  $\widehat{\Sigma}^*/N \xrightarrow{p} \Sigma$ .

- **Step 2.** We prove  $Z^*$  conditionally converges to a multivariate Gaussian distribution. From Lemma B.4,

$$\Lambda_2(\text{vec}(Z^*), \text{vec}(Z))^2 \leq \Lambda_2 \left( \text{vec} \left( \frac{1}{\sqrt{H}} \sum_{h=1}^H \phi_{hk}^* (y_{hk}^* - \phi_{hk}^* \widehat{M}_\pi) \right), \text{vec} \left( \frac{1}{\sqrt{H}} \sum_{h=1}^H \phi_{hk} (y_{hk} - \phi_{hk} M_\pi) \right) \right)^2.$$

Using Lemma B.5, we have the right side converges to 0, a.e. as  $K \rightarrow \infty$ . This means the conditional law of  $\text{vec}(Z^*)$  is close to the unconditional law of  $\text{vec}(Z)$ , and the latter essentially converges to a multivariate Gaussian distribution with zero mean and covariance matrix  $\Delta$  from Eq. (A.8).

By Slutsky's theorem, we have conditional on  $\mathcal{D}$ ,

$$\sqrt{N} \text{vec}(\widehat{M}_\pi^* - \widehat{M}_\pi) \xrightarrow{d} N \left( 0, (\Sigma^{-1} \otimes I_d) \Delta (\Sigma^{-1} \otimes I_d) \right), \quad (\text{A.11})$$

where  $\Delta$  is defined in Eq. (A.9).

According to the equivalence between FQE and plug-in estimator in Section A.2,

$$\widehat{v}_\pi^* = (\nu_1^\pi)^\top \sum_{h=0}^{H-1} (\widehat{M}_\pi^*)^h \widehat{R}^*, \widehat{v}_\pi = (\nu_1^\pi)^\top \sum_{h=0}^{H-1} (\widehat{M}_\pi)^h \widehat{R}.$$

Define a function  $g : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$  as

$$g(M) := (\nu_1^\pi)^\top \sum_{h=0}^{H-1} (M)^h w_r.$$

By the high-order matrix derivative (Petersen & Pedersen, 2008), we have

$$\frac{\partial}{\partial M} (\nu_1^\pi)^\top (M)^h w_r = \sum_{r=1}^{h-1} (M^r)^\top \nu_1^\pi w_r^\top (M^{h-1-r})^\top \in \mathbb{R}^{d \times d}.$$

This implies the gradient of  $g$  at  $\text{vec}(M_\pi)$

$$\begin{aligned} \nabla g(\text{vec}(M_\pi)) &= \text{vec} \left( \sum_{h=0}^{H-1} \sum_{r=1}^{h-1} (M_\pi^r)^\top \nu_1^\pi w_r^\top (M_\pi^{h-1-r})^\top \right) \\ &= \text{vec} \left( \sum_{h=1}^H \nu_h^\pi w_r^\top \sum_{h'=1}^{H-h} (M_\pi^{h'-1})^\top \right) \in \mathbb{R}^{d^2 \times 1}. \end{aligned}$$

Applying multivariate delta theorem (Theorem B.6) for Eq. (A.11), we have conditional on  $\mathcal{D}$

$$\sqrt{N} \left( g(\widehat{M}_\pi^*) - g(\widehat{M}_\pi) \right) \xrightarrow{d} \mathcal{N} \left( 0, \nabla^\top g(\text{vec}(\widehat{M}_\pi)) (\Sigma^{-1} \otimes I_d) \Delta (\Sigma^{-1} \otimes I_d) \nabla g(\text{vec}(\widehat{M}_\pi)) \right),$$

where  $\Delta$  is defined in Eq. (A.9). From Eq. (A.10), we have  $\widehat{\Sigma}/N \xrightarrow{p} \Sigma$ . Using Slutsky's theorem and Eqs. (A.7)-(A.8), we have

$$\sqrt{N} \left( \text{vec}(\widehat{M}_\pi - M_\pi) \right) \xrightarrow{d} \mathcal{N} \left( 0, (\Sigma^{-1} \otimes I_d) \Delta (\Sigma^{-1} \otimes I_d) \right).$$

This further implies  $\widehat{M}_\pi \xrightarrow{p} M_\pi$ . By continuous mapping theorem,

$$\sqrt{N} \left( g(\widehat{M}_\pi^*) - g(\widehat{M}_\pi) \right) \xrightarrow{d} \mathcal{N} \left( 0, \nabla^\top g(\text{vec}(M_\pi)) (\Sigma^{-1} \otimes I_d) \Delta (\Sigma^{-1} \otimes I_d) \nabla g(\text{vec}(M_\pi)) \right).$$

Now we simplify the variance term as follows:

$$\begin{aligned} & \nabla^\top g(\text{vec}(M_\pi)) (\Sigma^{-1} \otimes I_d) \Delta (\Sigma^{-1} \otimes I_d) \nabla g(\text{vec}(M_\pi)) \\ &= \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} w_r^\top \sum_{h'=1}^{H-h} (M_\pi)^{h'-1} \mathbb{E} \left[ \frac{1}{H} \sum_{h=1}^H \xi_h^\top \phi(s_h^1, a_h^1) \phi(s_h^1, a_h^1)^\top \xi_h \right] \sum_{h=1}^H \sum_{h'=1}^{H-h} (M_\pi)^{h'-1} w_r \Sigma^{-1} (\nu_h^\pi)^\top \\ &= \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} w_r^\top \sum_{h'=1}^{H-h} (M_\pi)^{h'-1} \mathbb{E} \left[ \frac{1}{H} \sum_{h=1}^H \xi_h^\top \phi(s_h^1, a_h^1) \phi(s_h^1, a_h^1)^\top \xi_h \right] \sum_{h'=1}^{H-h} (M_\pi)^{h'-1} w_r \Sigma^{-1} (\nu_h^\pi)^\top \\ &\quad + 2 \sum_{h_1 < h_2} (\nu_{h_1}^\pi)^\top \Sigma^{-1} w_r^\top \sum_{h'=1}^{H-h_1} (M_\pi)^{h'-1} \mathbb{E} \left[ \frac{1}{H} \sum_{h=1}^H \xi_h^\top \phi(s_h^1, a_h^1) \phi(s_h^1, a_h^1)^\top \xi_h \right] \sum_{h'=1}^{H-h_2} (M_\pi)^{h'-1} w_r \Sigma^{-1} (\nu_{h_2}^\pi)^\top \\ &= \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} \mathbb{E} \left[ \frac{1}{H} \sum_{h=1}^H w_r^\top \sum_{h'=1}^{H-h} (M_\pi)^{h'-1} \xi_h^\top \phi(s_h^1, a_h^1) \phi(s_h^1, a_h^1)^\top \xi_h \sum_{h'=1}^{H-h} (M_\pi)^{h'-1} w_r \right] \Sigma^{-1} (\nu_h^\pi)^\top \\ &\quad + 2 \sum_{h_1 < h_2} (\nu_{h_1}^\pi)^\top \Sigma^{-1} \mathbb{E} \left[ \frac{1}{H} \sum_{h=1}^H w_r^\top \sum_{h'=1}^{H-h_1} (M_\pi)^{h'-1} \xi_h^\top \phi(s_h^1, a_h^1) \phi(s_h^1, a_h^1)^\top \xi_h \sum_{h'=1}^{H-h_2} (M_\pi)^{h'-1} w_r \right] \Sigma^{-1} (\nu_{h_2}^\pi)^\top, \end{aligned}$$

where  $\xi_h^1 = \phi(s_h^1, a_h^1)^\top M_\pi - \phi^\pi(s_{h+1}^1)^\top$ . Recall that we define

$$\begin{aligned} \varepsilon_{h,h'}^1 &= Q_h^\pi(s_{h'}^1, a_{h'}^1) - (r_{h'}^1 + V_{h+1}^\pi(s_{h+1}^1)) \\ &= \sum_{h=1}^{H-h_1} \left( \phi(s_h^1, a_h^1)^\top M_\pi - \phi^\pi(s_{h+1}^1)^\top \right) (M_\pi)^{h'-1} w_r = \sum_{h=1}^{H-h} \xi_h^1 (M_\pi)^{h-1} w_r, \end{aligned}$$

where the second equation is from Eq. (A.3). This implies

$$\begin{aligned} & \nabla^\top g(\text{vec}(M_\pi)) (\Sigma^{-1} \otimes I_d) \Delta (\Sigma^{-1} \otimes I_d) \nabla g(\text{vec}(M_\pi)) \\ &= \sum_{h=1}^H (\nu_h^\pi)^\top \Sigma^{-1} \mathbb{E} \left[ \frac{1}{H} \sum_{h'=1}^H \phi(s_{h'}^1, a_{h'}^1) \phi(s_{h'}^1, a_{h'}^1)^\top (\varepsilon_{h,h'}^1)^2 \right] \Sigma^{-1} (\nu_h^\pi)^\top \\ &\quad + 2 \sum_{h_1 < h_2} (\nu_{h_1}^\pi)^\top \Sigma^{-1} \mathbb{E} \left[ \frac{1}{H} \sum_{h'=1}^H \phi(s_{h'}^1, a_{h'}^1) \phi(s_{h'}^1, a_{h'}^1)^\top \varepsilon_{h_1,h'}^1 \varepsilon_{h_2,h'}^1 \right] \Sigma^{-1} (\nu_{h_2}^\pi)^\top = \sigma^2. \end{aligned}$$

Therefore, we have proven that

$$\sqrt{N} \left( g(\widehat{M}_\pi^*) - g(\widehat{M}_\pi) \right) \xrightarrow{d} \mathcal{N} \left( 0, \sigma^2 \right).$$

On the other hand,

$$\widehat{R}^* = (\widehat{\Sigma}^*)^{-1} \sum_{k=1}^K \sum_{h=1}^H r_{hk}^* \phi_{hk}^* = KH (\widehat{\Sigma}^*)^{-1} \frac{1}{K} \sum_{k=1}^K \frac{1}{H} \sum_{h=1}^H r_{hk}^* \phi_{hk}^*.$$

Using Lemma B.3, we have

$$\Lambda_1\left(\frac{1}{K}\sum_{k=1}^K\frac{1}{H}\sum_{h=1}^H r_{hk}^* \phi_{hk}^*, \frac{1}{K}\sum_{k=1}^K\frac{1}{H}\sum_{h=1}^H r_{hk} \phi_{hk}\right) \leq \Lambda_1\left(\frac{1}{H}\sum_{h=1}^H r_{hk}^* \phi_{hk}^*, \frac{1}{H}\sum_{h=1}^H r_{hk} \phi_{hk}\right).$$

The right hand side of the display goes to 0 as  $K \rightarrow \infty$ . From the law of large number,

$$\frac{1}{K}\sum_{k=1}^K\frac{1}{H}\sum_{h=1}^H r_{hk} \phi_{hk} \xrightarrow{p} \mathbb{E}\left[\frac{1}{H}\sum_{h=1}^H \phi(s_h^1, a_h^1) \phi(s_h^1, a_h^1)^\top\right] w_r$$

Combining with the fact that the conditional laws of  $\widehat{\Sigma}^*$  concentrates around  $\Sigma$ , this ends the proof.  $\blacksquare$

## A.6. Proofs of Corollary 5.2 and Corollary 5.4

We prove the consistency of bootstrap confidence interval by using Lemma 23.3 in Van der Vaart (2000). Suppose  $\Psi(t) = \mathbb{P}(\mathcal{N}(0, \sigma^2) \leq t)$ . Combining Theorem 4.2 and Theorem 5.1, we have

$$\mathbb{P}_{\mathcal{D}}\left(\sqrt{N}(\widehat{v}_\pi - v_\pi) \leq t\right) \rightarrow \Psi(t), \quad \mathbb{P}_{W^*|\mathcal{D}}\left(\sqrt{N}(\widehat{v}_\pi^* - \widehat{v}_\pi) \leq t\right) \rightarrow \Psi(t).$$

Using the quantile convergence theorem (Lemma 21.1 in Van der Vaart (2000)), it implies  $q_\delta^\pi \rightarrow \Psi^{-1}(t)$  almost surely. Therefore,

$$\begin{aligned} \mathbb{P}_{\mathcal{D}W^*}\left(v_\pi \leq \widehat{v}_\pi - q_{\delta/2}^\pi\right) &= \mathbb{P}_{\mathcal{D}W^*}\left(\sqrt{N}(\widehat{v}_\pi - v_\pi) \geq q_{\delta/2}^\pi\right) \\ &\rightarrow \mathbb{P}_{\mathcal{D}W^*}\left(\mathcal{N}(0, \sigma^2) \geq \Psi^{-1}(\delta/2)\right) = 1 - \delta/2. \end{aligned}$$

This finishes the proof of Corollary 5.2.

It is well known that the convergence in distribution implies the convergence in moment under the uniform integrability condition. The proof of the consistency of bootstrap moment estimation is straightforward since the condition  $\limsup_{N \rightarrow \infty} \mathbb{E}_{W^*|\mathcal{D}}[(\sqrt{N}(\widehat{v}_\pi^* - \widehat{v}_\pi))^q] < \infty$  for some  $q > 2$  ensures a similar uniform integrability condition. Together with the distributional consistency in Theorem 5.1, we apply Lemma 2.1 in Kato (2011) then we reach the conclusion.  $\blacksquare$

## B. Supporting Results

We present a series of useful lemmas about Mallows metric.

**Lemma B.1** (Lemma 8.4 in Bickel & Freedman (1981)). Let  $\{X_i\}_{i=1}^n$  be independent random variables with common distribution  $\mu$ . Let  $\mu_n$  be the empirical distribution of  $X_1, \dots, X_n$ . Then  $\Lambda_l(\mu_n, \mu) \rightarrow 0$  a.e..

**Lemma B.2** (Lemma 8.5 in Bickel & Freedman (1981)). Suppose  $X_n, X$  are random variables and  $\Lambda_l(X_n, X) \rightarrow 0$ . Let  $f$  be a continuous function. Then  $\Lambda_l(f(X_n), f(X)) \rightarrow 0$ .

**Lemma B.3** (Lemma 8.6 of Bickel & Freedman (1981)). Let  $\{U_i\}_{i=1}^n, \{V_i\}_{i=1}^n$  be independent random vectors. Then we have

$$\Lambda_1\left(\sum_{i=1}^n U_i, \sum_{i=1}^n V_i\right) \leq \sum_{i=1}^n \Lambda_1(U_i, V_i).$$

**Lemma B.4** (Lemma 8.7 of Bickel & Freedman (1981)). Let  $\{U_i\}_{i=1}^n, \{V_i\}_{i=1}^n$  be independent random vectors and  $\mathbb{E}[U_j] = \mathbb{E}[V_j]$ . Then we have

$$\Lambda_2\left(\sum_{i=1}^n U_i, \sum_{i=1}^n V_i\right)^2 \leq \sum_{i=1}^n \Lambda_2(U_i, V_i)^2.$$

Let  $\mu_K$  and  $\mu$  be probabilities on  $\mathbb{R}^{2Hd}$ . A data point in  $\mathbb{R}^{2Hd}$  can be written as  $(x_1, \dots, x_H, y_1, \dots, y_H)$  where  $x_h \in \mathbb{R}^d$  and  $y_h \in \mathbb{R}^d$ . Denote

$$\begin{aligned}\Sigma(\mu) &= \int \frac{1}{H} \sum_{h=1}^H x_h x_h^\top \mu(dx_1, \dots, dx_H, dy_1, \dots, dy_H), \\ M(\mu) &= \Sigma(\mu)^{-1} \int \sum_{h=1}^H x_h y_h^\top \mu(dx_1, \dots, dx_H, dy_1, \dots, dy_H), \\ \varepsilon(\mu, x_1, \dots, x_H, y_1, \dots, y_H) &= \sum_{h=1}^H (y_h - x_h^\top M(\mu)).\end{aligned}$$

**Lemma B.5** (Lemma 7 in [Eck \(2018\)](#)). If  $\Lambda_4(\mu_K, \mu) \rightarrow 0$  as  $K \rightarrow \infty$ , then we have the  $\mu_K$ -law of  $\text{vec}(\sum_{h=1}^H \varepsilon(\mu_K, x_1, \dots, x_H, y_1, \dots, y_H) x_h^\top)$  converges to the  $\mu$ -law of  $\text{vec}(\sum_{h=1}^H \varepsilon(\mu, x_1, \dots, x_H, y_1, \dots, y_H) x_h^\top)$  in  $\Lambda_2$ .

**Theorem B.6** (Multivariate delta theorem). Suppose  $\{T_n\}$  is a sequence of  $k$ -dimensional random vectors such that  $\sqrt{n}(T_n - \theta) \xrightarrow{d} N(0, \Sigma(\theta))$ . Let  $g: \mathbb{R}^k \rightarrow \mathbb{R}$  be once differentiable at  $\theta$  with the gradient matrix  $\nabla g(\theta)$ . Then

$$\sqrt{n}(g(T_n) - g(\theta)) \xrightarrow{d} N(0, \nabla^\top g(\theta) \Sigma(\theta) \nabla g(\theta)).$$

We restate Lemma B.5 in [Duan & Wang \(2020\)](#) in the following that is proven using matrix Bernstein inequality.

**Lemma B.7.** Under the assumption  $\phi(s, a)^\top \Sigma^{-1} \phi(s, a) \leq C_1 d$  for all  $(s, a) \in \mathcal{X}$ , with probability at least  $1 - \delta$ ,

$$\left\| \Sigma^{-1/2} \left( \frac{1}{N} \sum_{n=1}^N \phi(s_n, a_n) \phi(s_n, a_n)^\top \right) \Sigma^{-1/2} - I \right\|_2 \leq \sqrt{\frac{2 \ln(2d/\delta) C_1 d H}{N}} + \frac{2 \ln(2d/\delta) C_1 d H}{3N}. \quad (\text{B.1})$$

## C. Supplement for Experiments

### C.1. Experiment details

The original CliffWalking environment from OpenAI gym has deterministic state transitions. That is, for any state-action pair  $(s, a)$ , there exists a corresponding  $s' \in \mathcal{S}$  such that  $\mathbb{P}(\cdot | s, a) = \delta_{s'}(\cdot)$ . We modify the environment in order to make it stochastic. Specifically, we introduce randomness in state transitions such that given a state-action pair  $(s, a)$ , the transition takes place in the same way as in the deterministic environment with probability  $1 - \epsilon$  and takes place as if the action were a random action  $a'$ , instead of the intended  $a$ , with probability  $\epsilon$ . This is an episodic tabular MDP and the agent stops when falling from the cliff or reaching the terminal point. We also reduce the penalty of falling off the cliff from  $-100$  to  $-50$ .

The original MountainCar environment from OpenAI gym has deterministic state transitions. We modify the environment in order to make it stochastic. Specifically, we introduce randomness in state transitions by adding a Gaussian random force, namely,  $\mathcal{N}(0, \frac{1}{10})$  multiplied by the constant-magnitude force from the original environment. We also increase the gravity parameter from 0.0025 to 0.008, the force parameter from 0.001 to 0.008 and the maximum allowed speed from 0.07 to 0.2.

**Empirical coverage probability.** The preceding discussion leads to the simulation method for estimating the coverage probability of a confidence interval. The simulation method has three steps:

1. Simulate many fresh dataset of episode size  $K$  following the behavior policy.
2. Compute the confidence interval for each dataset.
3. Compute the proportion of dataset for which the true value of target policy is contained in the confidence interval. That proportion is an estimate for the empirical coverage probability for the confidence interval.

## Bootstrapping Fitted Q-Evaluation for Off-Policy Inference

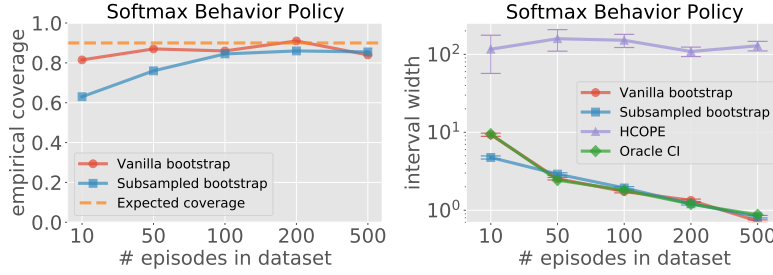


Figure 7. Left: Empirical coverage probability of CI; Right: CI width under different behavior policies.

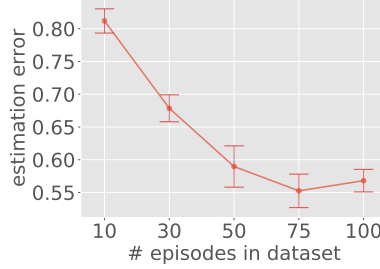


Figure 8. Error of correlation estimates, as data size increases.

The true value of target policy is computed through Monte Carlo rollouts with sufficient number of samples (10000 in our experiments).

With linear function approximation, we use the confidence interval proposed in Section 6 in Duan & Wang (2020) as a baseline since it is only available confidence interval based on FQE. In particular, it shows that with probability at least  $1 - \delta$ ,

$$|\hat{v}^\pi - v^\pi| \leq \sum_{h=1}^H (H - h + 1) \sqrt{(\hat{v}_h^\pi)^\top \hat{\Sigma}^{-1} \hat{v}_h^\pi} \left( \sqrt{2\lambda} + 2\sqrt{2d \log\left(1 + \frac{N}{\lambda d}\right) \log\left(\frac{3N^2 H}{\delta}\right)} + \frac{4}{3} \log\left(\frac{3N^2 H}{\delta}\right) \right),$$

where  $(\hat{v}_h^\pi)^\top = (\nu_1^\pi)^\top (\hat{M}^\pi)^h$  and  $\hat{M}^\pi$  is defined in Eq. (A.1).

### C.2. Additional experiments

In Figure 7, we include the result for soft-max behavior policy in the Cliff Walking environment. In Figure 8, we include the result for correlation estimation in Cliff Walking environment. The behavior policy is 0.1  $\epsilon$ -greedy policy while two target policies are optimal policy and 0.1  $\epsilon$ -greedy policy.

In order to better understand the tradeoff between computational efficiency and accuracy with finite samples, we conduct some empirical demonstrations based on Cliffwalking. We set  $s = K^\gamma$  and the true coverage probability is 0.9. It is relatively safe to set  $\gamma > 0.5$ . Note that  $\gamma = 1$  corresponds to the vanilla bootstrap that has the highest accuracy but heaviest computation.

We argue that bootstrapping sample transitions (which are dependent) would lead to inconsistent estimations of the error distribution and thus output wrong confidence interval and variance estimation. We further run one additional test using the taxi environment and further compute the CI and variance estimations based on different bootstrap distributions in Cliffwalking (CW) and taxi environments. This is already in an asymptotic regime since both the number of episodes and the number of bootstrap samples are  $10e+6$ . It is clear that bootstrapping by sample transition gives an incorrect distribution, thus it is inconsistent.

## Bootstrapping Fitted Q-Evaluation for Off-Policy Inference

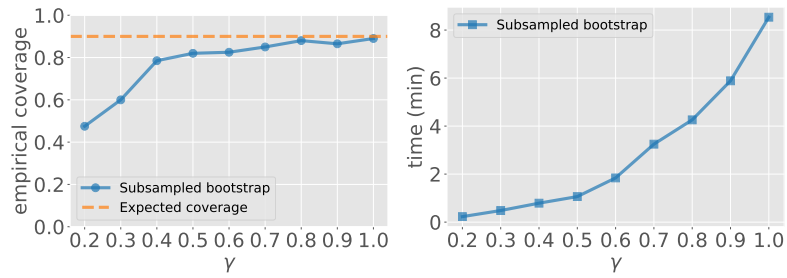


Figure 9. Subsampled bootstrap with  $s = K^\gamma$

	True distribution	By episodes	By sample transition
CI (CW)	(-20.44, -19.74)	(-20.40, -19.74)	(-20.52, -19.58)
Variance (CW)	0.45	0.44	0.082
CI (Taxi)	(2.49, 3.82)	(2.45, 3.73)	(2.01, 4.09)
Variance (Taxi)	0.16	0.16	0.39

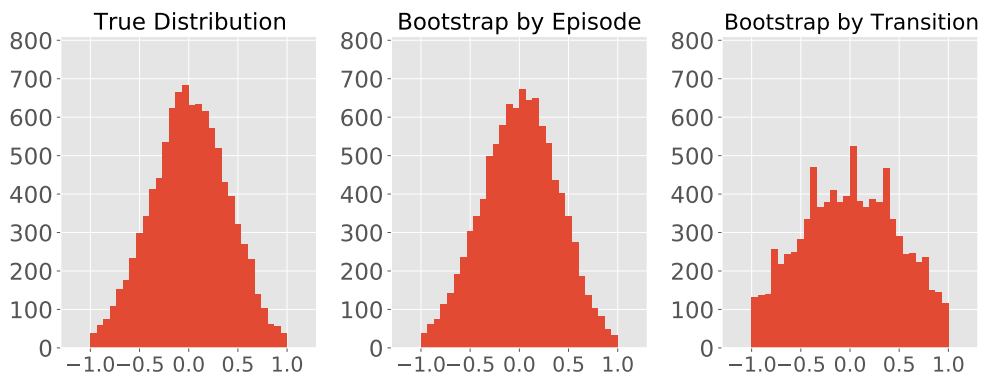


Figure 10. Taxi environment (Dietterich 2000)