
Supplement to “Spatiotemporal Causal Effects Evaluation: A Multi-Agent Reinforcement Learning Framework”

Anonymous Author(s)

Affiliation

Address

email

1 A More on the learning procedure

2 A.1 Estimation of the weight

3 Consider the following optimization problem

$$\hat{\omega}_i = \arg \min_{\omega_i \in \Omega} \sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{T-1} \Delta_{i,t}(\omega_i) f(S_{0,t+1}, S_{i,t+1}, \tilde{S}_{i,t+1}) \right|^2. \quad (1)$$

4 In our implementation, we set \mathcal{F} to a unit ball of a reproducing kernel Hilbert space (RFHS), i.e.,

$$\mathcal{F} = \{f \in \mathcal{H} : \|f\|_{\mathcal{H}} = 1\},$$

5 where

$$\mathcal{H} = \left\{ f(\cdot) = \sum_{t=0}^{T-1} b_t \kappa(S_{0,t+1}, S_{i,t+1}, \tilde{S}_{i,t+1}; \cdot) : \{b_t\}_{t=0}^{T-1} \in \mathbb{R}^T \right\},$$

6 for some positive definite kernel $\kappa(\cdot; \cdot)$. Similar to Theorem 2 of [6], the optimization problem in (1)
7 is then reduced to

$$\hat{\omega}_i = \arg \min_{\omega_i \in \Omega} \sum_{t_1=0}^{T-1} \sum_{t_2=0}^{T-1} \Delta_{i,t_1}(\omega_i) \Delta_{i,t_2}(\omega_i) \kappa(S_{0,t_1+1}, S_{i,t_1+1}, \tilde{S}_{i,t_1+1}; S_{0,t_2+1}, S_{i,t_2+1}, \tilde{S}_{i,t_2+1}).$$

8 We set Ω to the class of neural networks. One could use different parameters to factorize different ω_i
9 such that each $\hat{\omega}_i$ is computed separately. Alternatively, one could allow different ω_i to share some
10 common parameters. We detail our procedure in Algorithm 1.

11 A.2 Estimation of the Q-function and the value

12 We now describe methods to estimate Q_i and $V_i(\pi)$. For two given function classes \mathcal{G} and \mathcal{Q} , define
13 the following penalized estimator

$$\begin{aligned} \hat{g}_i(\cdot, \cdot, \cdot, \cdot; \eta, Q_i) &= \arg \min_{g \in \mathcal{G}} \frac{1}{T} \sum_{t=0}^{T-1} \{R_{i,t} + Q_i(\pi_i, \tilde{A}_i(\pi), S_{0,t+1}, S_{i,t+1}, \tilde{S}_{i,t+1}) \\ &\quad - \eta - Q_i(A_{i,t}, \tilde{A}_{i,t}, S_{0,t}, S_{i,t}, \tilde{S}_{i,t}) - g(A_{i,t}, \tilde{A}_{i,t}, S_{0,t}, S_{i,t}, \tilde{S}_{i,t})\}^2 + \mu J_2^2(g), \\ (\hat{V}_i(\pi), \hat{Q}_i) &= \arg \min_{(\eta, Q_i) \in \mathbb{R} \times \mathcal{Q}} \frac{1}{T} \sum_{t=0}^{T-1} \hat{g}_i^2(A_{i,t}, \tilde{A}_{i,t}, S_{0,t}, S_{i,t}, \tilde{S}_{i,t}; \eta, Q_i) + \lambda J_1^2(Q_i), \end{aligned}$$

Algorithm 1 Estimation of the weight.

Input: The data $\{(S_{0,j}, S_{i,j}, A_{i,j}, R_{i,j}, S_{0,j+1}, S_{i,j+1}) : 1 \leq i \leq N, 0 \leq j < T\}$. A target policy π .

Initial: Initial the density ratio $\omega_i = \omega_{i,\theta}$ for $1 \leq i \leq N$, to be neural networks parameterized by θ .

for iteration = 1, 2, \dots **do**

a Randomly sample a batch \mathcal{M} from $\{0, 1, \dots, T-1\}$.

b **Update** the parameter θ by $\theta \leftarrow \theta - \epsilon N^{-1} \sum_{i=1}^N \nabla_{\theta} D_i(\omega_{i,\theta}/z_{\omega_{i,\theta}})$ where $D_i(\omega_{i,\theta})$ is equal to

$$\frac{1}{|\mathcal{M}|} \sum_{t_1, t_2 \in \mathcal{M}} \Delta_{i,t_1}(\omega_{i,\theta}) \Delta_{i,t_2}(\omega_{i,\theta}) \kappa(S_{0,t_1+1}, S_{i,t_1+1}, \tilde{S}_{i,t_1+1}; S_{0,t_2+1}, S_{i,t_2+1}, \tilde{S}_{i,t_2+1}),$$

and $z_{\omega_{i,\theta}}$ is a normalization constant $z_{\omega_{i,\theta}} = |\mathcal{M}|^{-1} \sum_{t \in \mathcal{M}} \omega_{i,\theta}(S_{0,t+1}, S_{i,t+1}, \tilde{S}_{i,t+1})$.

Output $\omega_{i,\theta}$ for $1 \leq i \leq N$.

14 where J_1 and J_2 denote some penalty functions, μ and λ stand for some tuning parameters. Next we
15 derive the close-form expressions of $(\hat{V}_i(\pi), \hat{Q}_i)$ when RKHS is used to model Q_i and g_i .

16 Define vectors $Z_{i,t} = (A_{i,t}, \tilde{A}_{i,t}, S_{0,t}, S_{i,t}, \tilde{S}_{i,t})^\top$ and $Z_{i,t}^* = (\pi_i, \tilde{A}_i(\pi), S_{0,t+1}, S_{i,t+1}, \tilde{S}_{i,t+1})^\top$.
17 Let K_g and K_Q denote the reproducing kernels used to model g and Q , respectively. In practice, we
18 can use gaussian RBF kernels to model these two functions. For a given Q_i and η , the optimizer of \hat{g}_i
19 can be represented by $\sum_{t=0}^{T-1} \hat{\beta}_{i,t} K_g(Z_{i,t}, \cdot)$. As such, we obtain

$$\begin{aligned} \hat{\beta}_i &= \arg \min_{\beta} \frac{1}{T} \sum_{t=0}^{T-1} \left\{ R_{i,t} + Q_i(Z_{i,t}^*) - \eta - Q_i(Z_{i,t}) - \sum_{j=0}^{T-1} \beta_j K_g(Z_{i,j}, Z_{i,t}) \right\}^2 + \mu \beta^\top K_g \beta \\ &= \frac{1}{T} \beta^\top \{K_g K_g^\top + T\mu K_g\} \beta - \frac{2}{T} \beta^\top K_g (R + Q_i^* - Q_i - \eta \mathbf{1}) + \text{some terms that are independent of } \beta, \end{aligned}$$

20 where $K_g = \{K_g(Z_{i,j_1}, Z_{i,j_2})\}_{j_1, j_2}$ and R, Q_i^* and Q_i the column vectors formed by elements
21 in $R_t, Q_i(Z_{i,t}^*)$ and $Q_i(Z_{i,t})$, respectively. Notice that K_g is symmetric, by some calculations, we
22 obtain

$$\hat{\beta}_i = (K_g K_g^\top + T\mu K_g)^{-1} K_g (R + Q_i^* - Q_i - \eta \mathbf{1}) = (K_g + T\mu I)^{-1} (R + Q_i^* - Q_i - \eta \mathbf{1}).$$

23 As a result, for a given Q_i and η , we have

$$\hat{g}_i(Z_{i,t}; \eta, Q_i) = \hat{\beta}_i^\top K_g e_t,$$

24 where e_t denotes the column vector with the t -th element equals to one and other elements equal to
25 zero. As such,

$$\frac{1}{T} \sum_{t=0}^{T-1} \hat{g}_i^2(A_{i,t}, \tilde{A}_{i,t}, S_{0,t}, S_{i,t}, \tilde{S}_{i,t}; \eta, Q_i) = \frac{1}{T} \hat{\beta}_i^\top K_g K_g^\top \hat{\beta}_i.$$

26 Similarly, we can represent Q_i as $\sum_{t=0}^{2T-1} \hat{\alpha}_{i,t} K_Q(\tilde{Z}_{i,t}, \cdot)$ where $\tilde{Z}_{i,t}$ denotes the t -th element in the
27 vector $(Z_{i,0}^\top, Z_{i,1}^\top, \dots, Z_{i,T-1}^\top, Z_{i,0}^{*\top}, \dots, Z_{i,T-1}^{*\top})^\top$. Let K_Q denotes the corresponding $2T \times 2T$
28 matrix, we have

$$Q_i(Z_{i,t}) = \alpha_i^\top K_Q e_t \quad \text{and} \quad Q_i(Z_{i,t}^*) = \hat{\alpha}_i^\top K_Q e_{t+T+1}.$$

29 It follow that

$$Q_i^* - Q_i = \underbrace{[-I_T, I_T]}_C K_Q \hat{\alpha}_i.$$

30 Note that K_Q is symmetric. Let $E = K_g^\top (K_g + T\mu I)^{-1}$, $\hat{\alpha}_i$ corresponds to the solution of the
31 following optimization problem,

$$\hat{\alpha}_i = \arg \min_{\alpha} (R + C K_Q \alpha - \eta \mathbf{1})^\top E^\top E (R + C K_Q \alpha - \eta \mathbf{1}) + T\lambda \alpha^\top K_Q \alpha.$$

32 Taking derivatives with respect to α and η , we obtain

$$(\hat{\alpha}_i, \hat{V}_i(\pi))^\top = -([C K_Q, -1]^\top E^\top E [C K_Q, -1] + [T\lambda K_Q, 0; 0^\top, 0])^{-1} [C K_Q, -1] E^\top E R.$$

33 A.3 Estimation of the treatment assignment probability

34 Note that $b_i(\pi|S_{0,t}, S_{i,t}, \tilde{S}_{i,t}) = \mathbb{E}\{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\pi))|S_{0,t}, S_{i,t}, \tilde{S}_{i,t}\}$. It can thus be
 35 learned by applying machine learning algorithms to datasets with responses $\{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} =$
 36 $\tilde{A}_i(\pi)) : 0 \leq t < T\}$ and predictors $\{(S_{0,t}, S_{i,t}, \tilde{S}_{i,t}) : 0 \leq t < T\}$.

37 B Additional technical conditions and lemmas

38 B.1 Technical conditions

39 Let $Q_{i,\pi}^*$ denote the function such that $Q_{i,\pi}^*(S_{0,t}, S_{i,t}, \tilde{S}_{i,t}) = Q_i^*(\pi_i, \tilde{A}_i(\pi), S_{0,t}, S_{i,t}, \tilde{S}_{i,t})$ almost
 40 surely for any t and i . Similarly, let $\hat{Q}_{i,\pi}$ denote the function such that $\hat{Q}_{i,\pi}(S_{0,t}, S_{i,t}, \tilde{S}_{i,t}) =$
 41 $\hat{Q}_i(\pi_i, \tilde{A}_i(\pi), S_{0,t}, S_{i,t}, \tilde{S}_{i,t})$ almost surely for any t and i .

42 (A5)(i) $\sum_{i=1}^N |V_i^*(\pi) - \hat{V}_i(\pi)|/N = o_p(1)$; (ii) $\hat{Q}_{i,\pi} \in \mathcal{Q}$, $\hat{\omega}_i \in \mathcal{W}$ almost surely for any i . \mathcal{Q}
 43 and \mathcal{W} satisfy $\sup_Q N(\mathcal{Q}, e_Q, \varepsilon \|F\|_{Q,2}) \leq (A/\varepsilon)^\nu$, $\sup_Q N(\mathcal{W}, e_Q, \varepsilon \|F\|_{Q,2}) \leq (A/\varepsilon)^\nu$ for some
 44 $e \leq A = O(1)$, $\nu = O(NT)$, and their envelope functions are bounded by some constant M . (iii)
 45 $\max_i |\hat{Q}_{i,\pi}(s_0, s_i, \tilde{s}_i) - Q_{i,\pi}^*(s_0, s_i, \tilde{s}_i)|^2 p(b, s_0, s_i, \tilde{s}_i) ds_0 ds_i d\tilde{s}_i = o_p(1)$, $\max_i |\hat{\omega}_i(s_0, s_i, \tilde{s}_i) -$
 46 $\omega_i^*(s_0, s_i, \tilde{s}_i)|^2 p(b, s_0, s_i, \tilde{s}_i) ds_0 ds_i d\tilde{s}_i = o_p(1)$.

47 (A6)(i) $\max_i |V_i^*(\pi) - \hat{V}_i(\pi)|^2 = o_p((NT)^{-1/2})$; (ii) $\max_i |\hat{Q}_{i,\pi}(s_0, s_i, \tilde{s}_i) -$
 48 $Q_{i,\pi}^*(s_0, s_i, \tilde{s}_i)|^2 p(b, s_0, s_i, \tilde{s}_i) ds_0 ds_i d\tilde{s}_i = o_p((NT)^{-1/2})$; (iii) $\max_i |\hat{\omega}_i(s_0, s_i, \tilde{s}_i) -$
 49 $\omega_i^*(s_0, s_i, \tilde{s}_i)|^2 p(b, s_0, s_i, \tilde{s}_i) ds_0 ds_i d\tilde{s}_i = o_p((NT)^{-1/2})$; (iv) $T \gg N\nu^2 \log^4(NT)$.

50 B.2 An auxiliary lemma

51 We briefly introduce our setup before presenting the lemma. Let $\{Z_t : t \geq 0\}$ be a stationary
 52 β -mixing process whose β -mixing coefficients are given by $\{\beta(q) : q \geq 0\}$. Let \mathcal{F} be a pointwise
 53 measurable class of functions that take Z_t as input with a measurable envelope function F . For any
 54 $f \in \mathcal{F}$, suppose $\mathbb{E}f(Z_0) = 0$. Let $\sigma^2 > 0$ be a positive constant such that $\sup_{f \in \mathcal{F}} \mathbb{E}f^2(Z_0) \leq \sigma^2 \leq$
 55 $\mathbb{E}F^2(Z_0)$. In the following, we focus providing an exponential inequality for the empirical process
 56 $\sup_{f \in \mathcal{F}} |\sum_{t=0}^{T-1} f(Z_t)|$.

57 **Lemma B.1** Suppose the envelop function is uniformly bounded by some constant $M > 0$. In
 58 addition, suppose \mathcal{F} belongs to the class of VC-type class such that $\sup_Q N(\mathcal{F}, e_Q, \varepsilon \|F\|_{Q,2}) \leq$
 59 $(A/\varepsilon)^\nu$ [see Definition 2.1 in 3, for details] for some $A \geq e, \nu \geq 1$. Then there exist some constants
 60 $c, C > 0$ such that

$$\begin{aligned} \mathbb{P} \left(\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{T-1} f(Z_t) \right| > c \sqrt{\nu q \sigma^2 T \log \left(\frac{AM}{\sigma} \right)} + c\nu M \log \left(\frac{AM}{\sigma} \right) + cq\tau + Mq \right) \\ \leq Cq \exp \left(-\frac{\tau^2 q}{CT\sigma^2} \right) + Cq \exp \left(-\frac{\tau}{CM} \right) + \frac{T\beta(q)}{q}, \end{aligned}$$

61 for any $\tau > 0, 1 \leq q < T/2$.

62 C Proofs

63 We use c and C to denote some generic constants whose values are allowed to vary from place to
 64 place. For any two positive sequences $\{a_t\}_{t \geq 1}$ and $\{b_t\}_{t \geq 1}$, we write $a_t \leq b_t$ if there exists some
 65 constant $C > 0$ such that $a_t \leq Cb_t$ for any t . The notation $a_t \leq 1$ means $a_t = O(1)$.

66 Lemma 1 can thus be proven in a similar manner as Theorem 1 of [6]. Lemma 2 can be similarly
 67 proven as Lemma 1 of [7]. Theorem 2 can be proven in a similar manner as Theorem 3. In the
 68 following, we focus on proving Theorems 1, 3 and Lemma B.1.

69 C.1 Proof of Theorem 1

70 To prove Theorem 1, we apply the central limit theorem for mixing triangle arrays developed in [5].
 71 Define

$$\widehat{V}_t^{\text{DR}^*}(\boldsymbol{\pi}) = \frac{1}{N} \sum_{i=1}^N \left[V_i^*(\boldsymbol{\pi}) + \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{R_{i,t} + Q_{i,t+1}^*(\boldsymbol{\pi}) - Q_{i,t}^* - V_i^*(\boldsymbol{\pi})\} \right],$$

72 we have $\widehat{V}^{\text{DR}^*}(\boldsymbol{\pi}) = T^{-1} \sum_{t=0}^{T-1} \widehat{V}_t^{\text{DR}^*}(\boldsymbol{\pi})$.

73 Suppose we have shown each $\widehat{V}_t^{\text{DR}^*}(\boldsymbol{\pi})$ is an unbiased estimator for $V(\boldsymbol{\pi})$. For $t \in \{0, 1, \dots, T-1\}$,
 74 let $x_t = (NT)^{-1/2} \{\widehat{V}_t^{\text{DR}^*}(\boldsymbol{\pi}) - V(\boldsymbol{\pi})\}$. It suffices to show the conditions in (1)-(5) of [5] hold for
 75 $\{x_t : 0 \leq t < T\}$. We next verify these conditions.

76 **Condition (1).** Note that $\{R_{i,t}, Q_{i,t}^*, \omega_{i,t}^*, V_i(\boldsymbol{\pi}) : 1 \leq i \leq N, t \geq 0\}$ are uniformly bounded
 77 from infinity, the set of functions $\{b_i : 1 \leq i \leq N\}$ are uniformly bounded from zero. As such,
 78 $\{x_t : 0 \leq t < T\}$ are uniformly bounded. Condition (1) thus holds for any $\nu^* > 0$.

79 **Condition (2).** This condition is automatically implied by the assumption that $NT\text{Var}\{\widehat{V}^{\text{DR}^*}(\boldsymbol{\pi})\} \rightarrow$
 80 $\sigma^2 > 0$.

81 **Condition (3).** This condition holds by setting $\kappa = 0$ and $T_n = 0$ for any n .

82 **Condition (4).** Note that the strong mixing coefficients are upper bounded by the β -mixing coeffi-
 83 cients. Under Condition (A2), we can take the sequence $\alpha(h)$ in Condition (4) by $\kappa_0 \rho^h$.

84 **Condition (5).** Since $\kappa_0 \rho^h$ decays to zero at an exponential rate as h grows to infinity, Condition (5)
 85 is automatically satisfied.

86 It remains to show $\mathbb{E}\widehat{V}_t^{\text{DR}^*}(\boldsymbol{\pi}) = V(\boldsymbol{\pi})$ for any t . Suppose (A4) holds. Under the given conditions,
 87 we have $V_i^*(\boldsymbol{\pi}) = V_i(\boldsymbol{\pi})$. By Lemma 2, we have

$$\mathbb{E}\{R_{i,t} + Q_{i,t+1}^*(\boldsymbol{\pi}) - Q_{i,t}^* - V_i^*(\boldsymbol{\pi}) | \mathbf{A}_t, \mathbf{S}_t\} = 0,$$

88 and hence,

$$\mathbb{E}\omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{R_{i,t} + Q_{i,t+1}^*(\boldsymbol{\pi}) - Q_{i,t}^* - V_i^*(\boldsymbol{\pi})\} = 0.$$

89 Consequently, $\mathbb{E}\widehat{V}_t^{\text{DR}^*}(\boldsymbol{\pi}) = N^{-1} \sum_{i=1}^N V_i(\boldsymbol{\pi}) = V(\boldsymbol{\pi})$.

90 Suppose (A3) holds. Then we have $\omega_{i,t}^* = \omega_{i,t}$ for any i, t where $\omega_{i,t}$ is a shorthand for
 91 $\omega_i(\boldsymbol{\pi}, S_{0,t}, S_{i,t}, \tilde{S}_{i,t})$. As a result, for any i, t , the expectation of the density ratio $\omega_{i,t}^* \mathbb{I}(A_{i,t} =$
 92 $\pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi})) / b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})$ equals one. As such, we have

$$\begin{aligned} & \mathbb{E} \left\{ V_i^*(\boldsymbol{\pi}) - \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} V_i^*(\boldsymbol{\pi}) \right\} \\ &= V_i^*(\boldsymbol{\pi}) \mathbb{E} \left\{ 1 - \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \right\} = 0. \end{aligned} \quad (2)$$

93 In addition, using similar arguments in (2), we have by (A3) that

$$\mathbb{E} \left\{ \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} R_{i,t} \right\} = V_i(\boldsymbol{\pi}). \quad (3)$$

94 Moreover, by some calculations, we have

$$\begin{aligned} \mathbb{E} \left\{ \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} Q_{i,t}^* \right\} &= \mathbb{E} \left\{ \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} Q_{i,t+1}^*(\boldsymbol{\pi}) \right\} \\ &= \int_{s_0, s_i, \tilde{s}_i} Q_i^*(\pi_i, \tilde{A}_i(\boldsymbol{\pi}), s_0, s_i, \tilde{s}_i) p(\boldsymbol{\pi}, s_0, s_i, \tilde{s}_i) ds_0 ds_i d\tilde{s}_i. \end{aligned}$$

Consequently,

$$\mathbb{E} \left[\omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{Q_{i,t+1}^*(\boldsymbol{\pi}) - Q_{i,t}^*\} \right] = 0.$$

This together with (2) and (3) yields

$$\mathbb{E} \left[V_i^*(\boldsymbol{\pi}) + \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{R_{i,t} + Q_{i,t+1}^*(\boldsymbol{\pi}) - Q_{i,t}^* - V_i^*(\boldsymbol{\pi})\} \right] = V_i(\boldsymbol{\pi}).$$

It follows that $\mathbb{E}\hat{V}^{\text{DR}*}(\boldsymbol{\pi}) = V(\boldsymbol{\pi})$.

Thus, $\hat{V}^{\text{DR}*}(\boldsymbol{\pi})$ is unbiased when either (A3) or (A4) holds. The proof is hence completed.

C.2 Proof of Theorem 3

By Theorem 1, it suffices to show $\sqrt{NT}\hat{V}^{\text{DR}}(\boldsymbol{\pi})$ is asymptotically equivalent to $\sqrt{NT}\hat{V}^{\text{DR}*}(\boldsymbol{\pi})$.

Note that $\hat{V}^{\text{DR}}(\boldsymbol{\pi}) - \hat{V}^{\text{DR}*}(\boldsymbol{\pi})$ can be decomposed by $\eta_1 + \eta_2 + \eta_3 + \eta_4 + \eta_5$ where

$$\begin{aligned} \eta_1 &= \frac{1}{NT} \sum_{t=0}^{T-1} \sum_{i=1}^N \left\{ \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} - 1 \right\} \{V_i^*(\boldsymbol{\pi}) - \hat{V}_i(\boldsymbol{\pi})\}, \\ \eta_2 &= \frac{1}{NT} \sum_{t=0}^{T-1} \sum_{i=1}^N \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{\hat{Q}_{i,t+1}(\boldsymbol{\pi}) - \hat{Q}_{i,t} - Q_{i,t+1}^*(\boldsymbol{\pi}) + Q_{i,t}^*\}, \\ \eta_3 &= \frac{1}{NT} \sum_{t=0}^{T-1} \sum_{i=1}^N (\hat{\omega}_{i,t} - \omega_{i,t}^*) \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{R_{i,t} + Q_{i,t+1}^*(\boldsymbol{\pi}) - Q_{i,t}^* - V_i^*(\boldsymbol{\pi})\}, \\ \eta_4 &= \frac{1}{NT} \sum_{t=0}^{T-1} \sum_{i=1}^N (\hat{\omega}_{i,t} - \omega_{i,t}^*) \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{\hat{Q}_{i,t+1}(\boldsymbol{\pi}) - \hat{Q}_{i,t} - Q_{i,t+1}^*(\boldsymbol{\pi}) + Q_{i,t}^*\}, \\ \eta_5 &= \frac{1}{NT} \sum_{t=0}^{T-1} \sum_{i=1}^N (\hat{\omega}_{i,t} - \omega_{i,t}^*) \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{V_i^*(\boldsymbol{\pi}) - \hat{V}_i(\boldsymbol{\pi})\}. \end{aligned}$$

In the following, we show $|\eta_j| = o_p((NT)^{-1/2})$, for $j = 1, 2, \dots, 5$.

Upper bounds on $|\eta_1|$: Note that $\eta_1 = N^{-1} \sum_{i=1}^N \eta_{1,i}$ where

$$\eta_{1,i} = \{V_i^*(\boldsymbol{\pi}) - \hat{V}_i(\boldsymbol{\pi})\} \left[\frac{1}{T} \sum_{t=0}^{T-1} \left\{ \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} - 1 \right\} \right].$$

When (A3) holds, we have $\omega_{i,t}^* = \omega_{i,t}$ for any i, t . The expectation of the density ratio equals one.

As a result, we have

$$\mathbb{E} \left\{ \omega_{i,t} \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} - 1 \right\} = 0,$$

for any i, t . In the following, we apply the Bernstein's inequality for exponential β -mixing processes [2] to bound $|\eta_1|$.

Under Condition (A2), the β -mixing coefficients of the sequence

$$\left\{ \omega_{i,t} \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} - 1 : t \geq 0 \right\}, \quad (4)$$

decays to zero at an exponential rate. In addition, all the terms in (4) are uniformly bounded by some constant $c > 0$. As a result,

$$\max_{t_1, t_2} \mathbb{E} \left| \omega_{i,t_1} \frac{\mathbb{I}(A_{i,t_1} = \pi_i, \tilde{A}_{i,t_1} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t_1}, \tilde{S}_{i,t_1})} - 1 \right| \left| \omega_{i,t_2} \frac{\mathbb{I}(A_{i,t_2} = \pi_i, \tilde{A}_{i,t_2} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t_2}, S_{i,t_2}, \tilde{S}_{i,t_2})} - 1 \right| = O(1).$$

111 It thus follows from Theorem 4.2 of [2] that there exists some constant $C > 0$ such that there exists
 112 some constant $C > 0$ such that for any $\tau \geq 0$ and integer $1 < q < T$,

$$\begin{aligned} & \max_i \mathbb{P} \left(\left| \sum_{t=0}^{T-1} \left\{ \omega_{i,t} \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} - 1 \right\} \right| \geq 6\tau \right) \leq \frac{T}{q} \beta(q) \\ & + \max_i \mathbb{P} \left(\left| \sum_{t \in \mathcal{I}_r} \left\{ \omega_{i,t} \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} - 1 \right\} \right| \geq \tau \right) + 4 \exp \left\{ -\frac{\tau^2}{Cq(T+\tau)} \right\} \end{aligned} \quad (5)$$

113 where $\mathcal{I}_r = \{q\lfloor T/q \rfloor, q\lfloor T/q \rfloor + 1, \dots, T-1\}$. Suppose $\tau \geq qc$. Notice that $|\mathcal{I}_r| \leq q$. It follows
 114 that

$$\max_i \mathbb{P} \left(\left| \sum_{t \in \mathcal{I}_r} \left\{ \omega_{i,t} \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} - 1 \right\} \right| \geq \tau \right) = 0. \quad (6)$$

115 Under (A2), $\beta(q) = O(\rho^q)$. Set $q = -3 \log(NT) / \log \rho$, we obtain $T\beta(q)/q = O(N^{-3}T^{-2})$. Set
 116 $\tau = \max\{2\sqrt{CqT \log(NT)}, 4Cq \log(NT)\}$, we obtain as $T \rightarrow \infty$ that

$$\frac{\tau^2}{2} \geq 2CqT \log(NT) \quad \text{and} \quad \frac{\tau^2}{2} \geq 2Cq\tau \log(nT) \quad \text{and} \quad \tau \geq qc.$$

117 Since $\sqrt{CqT \log(NT)} \gg 2Cq \log(NT)$, it follows from (5) and (6) that

$$\max_i \mathbb{P} \left(\left| \sum_{t=0}^{T-1} \left\{ \omega_{i,t} \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} - 1 \right\} \right| \geq 12\sqrt{CqT \log(NT)} \right) \leq N^{-2}T^{-2}.$$

118 By Bonferroni's inequality, we obtain the following event occurs with probability at least $1 -$
 119 $O(N^{-1}T^{-1})$,

$$\max_i \left| \sum_{t=0}^{T-1} \left\{ \omega_{i,t} \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} - 1 \right\} \right| \leq 12\sqrt{CqT \log(NT)}.$$

120 It follows that

$$|\eta_1| \leq \frac{1}{N} \sum_{i=1}^N |\eta_{1,i}| \leq \frac{\log(NT)}{\sqrt{T}} \left(\frac{1}{N} \sum_{i=1}^N |V_i^*(\boldsymbol{\pi}) - \hat{V}_i(\boldsymbol{\pi})| \right), \quad (7)$$

121 with probability approaching 1. Under (A6) and the condition that $T \gg N \log^4(NT)$, we obtain
 122 $\eta_1 = o_p((NT)^{-1/2})$.

123 **Upper bounds on $|\eta_2|$:** When (A3) holds, we have $\omega_{i,t}^* = \omega_{i,t}$ for any i and t . As discussed in the
 124 proof of Theorem 1, we have $\mathbb{E}\eta_{2,i} = 0$ for any i where

$$\eta_{2,i} = \frac{1}{T} \sum_{t=0}^{T-1} \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{ \hat{Q}_{i,t+1}(\boldsymbol{\pi}) - \hat{Q}_{i,t}(\boldsymbol{\pi}) - Q_{i,t+1}^*(\boldsymbol{\pi}) + Q_{i,t}^*(\boldsymbol{\pi}) \}.$$

125 In addition, notice that $\eta_{2,i}$ can be written as

$$\eta_{2,i} = \frac{1}{T} \sum_{t=0}^{T-1} \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{ \hat{Q}_{i,t+1}(\boldsymbol{\pi}) - \hat{Q}_{i,t}(\boldsymbol{\pi}) - Q_{i,t+1}^*(\boldsymbol{\pi}) + Q_{i,t}^*(\boldsymbol{\pi}) \}.$$

126 We apply Lemma B.1 to bound $\max_i |\eta_{2,i}|$. Define the class of functions $\mathcal{Q}_{i,\varepsilon}$ by

$$\left\{ f \in \mathcal{Q} : \max_i \int_{s_0, s_i, \tilde{s}_i} |f(s_0, s_i, \tilde{s}_i) - Q_{i,\boldsymbol{\pi}}^*(s_0, s_i, \tilde{s}_i)|^2 p(b, s_0, s_i, \tilde{s}_i) ds_0 ds_i d\tilde{s}_i \leq \varepsilon \right\},$$

127 where $\varepsilon = \epsilon N^{-1/2} T^{-1/2}$ for some sufficiently small $\epsilon > 0$. It then follows from (A5)(ii) and (iii)
 128 that $\hat{Q}_{i,\boldsymbol{\pi}} \in \mathcal{Q}_\varepsilon$ for any i with probability tending to 1. As such, we have

$$\begin{aligned} \eta_{2,i} & \leq T^{-1} \sup_{Q_{i,\boldsymbol{\pi}} \in \mathcal{Q}_{i,\varepsilon}} \left| \sum_{t=0}^{T-1} \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\boldsymbol{\pi}))}{b_i(\boldsymbol{\pi}|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{ f(S_{0,t+1}, S_{i,t+1}, \tilde{S}_{i,t+1}) \right. \\ & \quad \left. - f(S_{0,t}, S_{i,t}, \tilde{S}_{i,t}) - Q_{i,t+1}^*(\boldsymbol{\pi}) + Q_{i,t}^*(\boldsymbol{\pi}) \} \right|. \end{aligned}$$

129 Consider the process $\{(S_{0,t}, S_{i,t}, \tilde{S}_{i,t}, A_{i,t}, \tilde{A}_{i,t}, S_{0,t+1}, S_{i,t+1}, \tilde{S}_{i,t+1}) : t \geq 0\}$. Under (A2), such
 130 a process has β -mixing coefficients $\{\beta^*(q) : q \geq 0\}$ that satisfies $\beta^*(q) = O(\rho^q)$ as well. For any f ,
 131 define the function $g = g(f)$ such that

$$g(S_{0,t}, S_{i,t}, \tilde{S}_{i,t}, A_{i,t}, \tilde{A}_{i,t}, S_{0,t+1}, S_{i,t+1}, \tilde{S}_{i,t+1}) = \omega_{i,t}^* \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\pi))}{b_i(\pi|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \\ \times \{f(S_{0,t+1}, S_{i,t+1}, \tilde{S}_{i,t+1}) - f(S_{0,t}, S_{i,t}, \tilde{S}_{i,t}) - Q_{i,t+1}^*(\pi) + Q_{i,t}^*(\pi)\},$$

132 almost surely. Consider the class of functions $\mathcal{G}_{i,\varepsilon} = \{g(f) : f \in \mathcal{Q}_{i,\varepsilon}\}$. Since $\mathcal{Q}_{i,\varepsilon}$ belongs to the
 133 class of VC-type class, so does $\mathcal{G}_{i,\varepsilon}$. Moreover, the VC-index of $\mathcal{G}_{i,\varepsilon}$ is the same as $\mathcal{Q}_{i,\varepsilon}$. Under the
 134 boundedness assumption in Theorem 2, we have

$$\mathbb{E}g^2(S_{0,t}, S_{i,t}, \tilde{S}_{i,t}, A_{i,t}, \tilde{A}_{i,t}, S_{0,t+1}, S_{i,t+1}, \tilde{S}_{i,t+1}) \leq O(1)\varepsilon,$$

135 for some constant $O(1)$. In addition, the envelope function of $\mathcal{G}_{i,\varepsilon}$ is uniformly bounded.

136 Let $Z_{i,t} = (S_{0,t}, S_{i,t}, \tilde{S}_{i,t}, A_{i,t}, \tilde{A}_{i,t}, S_{0,t+1}, S_{i,t+1}, \tilde{S}_{i,t+1})$. Applying Lemma B.1, we obtain

$$\max_i \mathbb{P} \left(\sup_{g \in \mathcal{G}_{i,\varepsilon}} \left| \sum_{t=0}^{T-1} g(Z_{i,t}) \right| > c \sqrt{\nu q \varepsilon T \log \left(\frac{1}{\varepsilon} \right)} + c \nu \log \left(\frac{1}{\varepsilon} \right) + cq\tau + cq \right) \\ \leq cq \exp \left(-\frac{\tau^2 q}{cT\varepsilon} \right) + cq \exp \left(-\frac{\tau}{c} \right) + \frac{T\beta(q)}{q},$$

137 for some constant $c > 0$. Set $q = -2 \log(NT) / \log \rho$, we have $T\beta(q)/q = O(N^{-2}T^{-1})$. Set
 138 $\tau = \max(2c \log(NT), \sqrt{2c\varepsilon T \log(NT)/q})$, the RHS is bounded by $O(N^{-2}T^{-1} \log(NT))$. By
 139 Bonferroni's inequality, we obtain with probability tending to 1 that

$$T|\eta_{2,i}| \leq c \sqrt{\nu q \varepsilon T \log \left(\frac{1}{\varepsilon} \right)} + c \nu \log \left(\frac{1}{\varepsilon} \right) + cq\tau + cq, \quad \forall i \in \{1, \dots, N\},$$

140 or equivalently,

$$\max_i |\eta_{2,i}| \preceq \sqrt{\frac{\varepsilon}{NT}} + o \left(\frac{1}{\sqrt{NT}} \right),$$

141 under the condition that $T \gg N\nu^2 \log^4(NT)$. Since ε can be chosen arbitrarily small, we obtain
 142 $\max_i |\eta_{2,i}| = o_p((NT)^{-1/2})$. This in turn implies $\eta_2 = o_p((NT)^{-1/2})$.

143 **Upper bounds on $|\eta_3|$:** Using similar arguments in proving $\eta_2 = o_p((NT)^{-1/2})$, we can show
 144 $\eta_3 = o_p((NT)^{-1/2})$. We omit the technical details to save space.

145 **Upper bounds on $|\eta_4|$ and $|\eta_5|$:** We show $\eta_4 = o_p((NT)^{-1/2})$ only. Using similar arguments, one
 146 can show $\eta_5 = o_p((NT)^{-1/2})$.

147 Note that

$$\eta_4 = \frac{1}{NT} \sum_{t=0}^{T-1} \sum_{i=1}^N (\hat{\omega}_{i,t} - \omega_{i,t}^*) \frac{\mathbb{I}(A_{i,t} = \pi_i, \tilde{A}_{i,t} = \tilde{A}_i(\pi))}{b_i(\pi|S_{0,t}, S_{i,t}, \tilde{S}_{i,t})} \{\hat{Q}_{i,t+1}(\pi) - \hat{Q}_{i,t}(\pi) - Q_{i,t+1}^*(\pi) + Q_{i,t}^*(\pi)\} \\ \leq O(1) \frac{1}{NT} \sum_{t=0}^{T-1} \sum_{i=1}^N |\hat{\omega}_{i,t} - \omega_{i,t}^*| |\hat{Q}_{i,t+1}(\pi) - \hat{Q}_{i,t}(\pi) - Q_{i,t+1}^*(\pi) + Q_{i,t}^*(\pi)| \\ \leq O(1) \left\{ \frac{1}{NT} \sum_{t=0}^{T-1} \sum_{i=1}^N [(\hat{\omega}_{i,t} - \omega_{i,t}^*)^2 + \{\hat{Q}_{i,t+1}(\pi) - \hat{Q}_{i,t}(\pi) - Q_{i,t+1}^*(\pi) + Q_{i,t}^*(\pi)\}^2] \right\} \\ \leq O(1) \left\{ \frac{1}{NT} \sum_{t=0}^{T-1} \sum_{i=1}^N (\hat{\omega}_{i,t} - \omega_{i,t}^*)^2 \right\} + O(1) \left\{ \frac{1}{NT} \sum_{t=0}^{T-1} \sum_{i=1}^N \{\hat{Q}_{i,t}(\pi) - Q_{i,t}^*(\pi)\}^2 \right\},$$

148 where $O(1)$ denotes some universal constant, and the last two inequalities are due to Cauchy-Schwarz
 149 inequality.

150 To prove $\eta_4 = o_p((NT)^{-1/2})$, it suffices to show

$$\max_i \left[\frac{1}{T} \sum_{t=0}^T \{ \hat{Q}_{i,t}(\boldsymbol{\pi}) - Q_{i,t}^*(\boldsymbol{\pi}) \}^2 \right] = o_p((NT)^{-1/2}), \quad (8)$$

151 and

$$\max_i \left\{ \frac{1}{T} \sum_{t=0}^{T-1} (\hat{\omega}_{i,t} - \omega_{i,t}^*)^2 \right\} = o_p((NT)^{-1/2}). \quad (9)$$

152 The left-hand-side (LHS) of (8) can be upper bounded by

$$\max_i \sup_{f \in \mathcal{Q}_{i,\varepsilon}} \left[\frac{1}{T} \sum_{t=0}^T \{ f(Z_{i,t}) - Q_{i,t}^*(\boldsymbol{\pi}) \}^2 \right],$$

153 with probability tending to 1. Using similar arguments in proving $\eta_2 = o_p((NT)^{-1/2})$, we can show

$$\max_i \sup_{f \in \mathcal{Q}_{i,\varepsilon}} \left| \frac{1}{T} \sum_{t=0}^T \{ f(Z_{i,t}) - Q_{i,t}^*(\boldsymbol{\pi}) \}^2 - \frac{1}{T} \sum_{t=0}^T \mathbb{E} \{ f(Z_{i,t}) - Q_{i,t}^*(\boldsymbol{\pi}) \}^2 \right| \preceq \frac{\epsilon}{\sqrt{NT}} + o\left(\frac{1}{\sqrt{NT}}\right),$$

154 with probability tending to 1. Under (A6), we have

$$\max_i \sup_{f \in \mathcal{Q}_{i,\varepsilon}} \left| \frac{1}{T} \sum_{t=0}^T \mathbb{E} \{ f(Z_{i,t}) - Q_{i,t}^*(\boldsymbol{\pi}) \}^2 \right| \preceq \frac{\epsilon}{\sqrt{NT}}.$$

155 It follows that

$$\max_i \sup_{f \in \mathcal{Q}_{i,\varepsilon}} \left[\frac{1}{T} \sum_{t=0}^T \{ f(Z_{i,t}) - Q_{i,t}^*(\boldsymbol{\pi}) \}^2 \right] \preceq \frac{\epsilon}{\sqrt{NT}} + o\left(\frac{1}{\sqrt{NT}}\right),$$

156 with probability tending to 1. Let $\epsilon \rightarrow 0$, we obtain (8). Similarly, we can show (9) holds. The proof
157 is hence completed.

158 C.3 Proof of Lemma B.1

159 We break the proof into three steps. In the first step, we use Berbee's coupling lemma [see Lemma
160 4.1 in 4] to approximate $\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{T-1} f(Z_t) \right|$ by sum of i.i.d. variables. In the second step, we
161 apply the tail inequality in Lemma 1 of [1] to bound the deviation between the empirical process and
162 its mean. Finally, we apply the maximal inequality in Corollary 5.1 of [3] to bound the expectation of
163 the empirical process.

164 **Step 1.** Following the discussion below Lemma 4.1 of [4], we can construct a sequence of random
165 variables $\{Z_t^0 : t \geq 0\}$ such that

$$\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{T-1} f(Z_t) \right| = \sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{T-1} f(Z_t^0) \right|, \quad (10)$$

166 with probability at least $1 - T\beta(q)/q$, and that the sequences $\{U_{2i}^0 : i \geq 0\}$ and $\{U_{2i+1}^0 : i \geq 0\}$ are
167 i.i.d. where $U_i^0 = (Z_{iq}^0, Z_{iq+1}^0, \dots, Z_{iq+q-1}^0)$.

168 Recall that $\mathcal{I}_r = \{q\lfloor T/q \rfloor, q\lfloor T/q \rfloor + 1, \dots, T-1\}$, we have

$$\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{T-1} f(Z_t^0) \right| \leq \sum_{j=0}^{q-1} \sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/q \rfloor} f(Z_{tq+j}^0) \right| + \sup_{f \in \mathcal{F}} \left| \sum_{t \in \mathcal{I}_r} f(Z_t^0) \right|.$$

169 Under the boundedness assumption on F , the second term on the right-hand-side (RHS) is bounded
170 from above by Mq . Without loss of generality, suppose $\lfloor T/q \rfloor$ is an even number. The first term on

171 the RHS can be bounded from above by $\sum_{j=0}^{2q-1} \sup_{f \in \mathcal{F}} |\sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq+j}^0)|$. To summarize, we
 172 have shown

$$\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{T-1} f(Z_t^0) \right| \leq \sum_{j=0}^{2q-1} \sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq+j}^0) \right| + Mq.$$

173 This together with (10) yields that

$$\mathbb{P} \left(\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{T-1} f(Z_t) \right| > 2\tau q + Mq \right) \leq \mathbb{P} \left(\sum_{j=0}^{2q-1} \sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq+j}^0) \right| > 2\tau q \right) + \frac{T\beta(q)}{q}, \quad (11)$$

174 for any $\tau > 0$. By Bonferroni's inequality, we obtain

$$\mathbb{P} \left(\sum_{j=0}^{2q-1} \sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq+j}^0) \right| > 2\tau q \right) \leq \sum_{j=0}^{2q-1} \mathbb{P} \left(\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq+j}^0) \right| > \tau \right),$$

175 for any $\tau > 0$. Since the process is stationary, we obtain

$$\mathbb{P} \left(\sum_{j=0}^{2q-1} \sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq+j}^0) \right| > 2\tau q \right) \leq 2q \mathbb{P} \left(\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq}^0) \right| > \tau \right).$$

176 Combining this together with (11) yields

$$\mathbb{P} \left(\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{T-1} f(Z_t) \right| > 2\tau q + Mq \right) \leq 2q \mathbb{P} \left(\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq}^0) \right| > \tau \right) + \frac{T\beta(q)}{q}. \quad (12)$$

177 By construction, $\{Z_{2tq}^0 : t \geq 0\}$ are i.i.d. This completes the proof of the first step.

178 **Step 2.** In the second step, we focus on relating the empirical process $\sup_{f \in \mathcal{F}} |\sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq}^0)|$
 179 to its expectation. Without loss of generality, assume $T = kq$ for some integer $k > 0$. Set the
 180 constants η and δ in Lemma 1 of [1] to 1, we obtain

$$\begin{aligned} \mathbb{P} \left(\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq}^0) \right| > 2\mathbb{E} \sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq}^0) \right| + \tau \right) \\ \leq 4 \exp \left(-\frac{\tau^2}{2T\sigma^2/q} \right) + \exp \left(-\frac{\tau}{CM} \right), \end{aligned}$$

181 for some constant $C > 0$. Combining this together with (12), we obtain

$$\begin{aligned} \mathbb{P} \left(\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{T-1} f(Z_t) \right| > 4q \mathbb{E} \sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq}^0) \right| + 2\tau q + Mq \right) \\ \leq 8q \exp \left(-\frac{\tau^2}{2T\sigma^2/q} \right) + 2q \exp \left(-\frac{\tau}{CM} \right) + \frac{T\beta(q)}{q}, \end{aligned} \quad (13)$$

182 for any $\tau > 0$. This completes the proof of the second step.

183 **Step 3.** It remains to bound $\mathbb{E} \sup_{f \in \mathcal{F}} |\sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq}^0)|$. By Corollary 5.1 of [3], we obtain

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{\lfloor T/(2q) \rfloor} f(Z_{2tq}^0) \right| \leq \sqrt{\frac{\nu\sigma^2 T}{q} \log \left(\frac{AM}{\sigma} \right)} + \nu M \log \left(\frac{AM}{\sigma} \right).$$

184 Combining this together with (13), we obtain

$$\begin{aligned} \mathbb{P} \left(\sup_{f \in \mathcal{F}} \left| \sum_{t=0}^{T-1} f(Z_t) \right| > c \sqrt{\nu q \sigma^2 T \log \left(\frac{AM}{\sigma} \right)} + c \nu M \log \left(\frac{AM}{\sigma} \right) + cq\tau + Mq \right) \\ \leq Cq \exp \left(-\frac{\tau^2 q}{CT\sigma^2} \right) + Cq \exp \left(-\frac{\tau}{CM} \right) + \frac{T\beta(q)}{q}, \end{aligned}$$

185 for some constants $c, C > 0$ and any $\tau > 0, 1 \leq q < T/2$. The proof is hence completed.

References

- [1] Radosław Adamczak. A tail inequality for suprema of unbounded empirical processes with applications to Markov chains. *Electron. J. Probab.*, 13:no. 34, 1000–1034, 2008.
- [2] Xiaohong Chen and Timothy M. Christensen. Optimal uniform convergence rates and asymptotic normality for series estimators under weak dependence and weak conditions. *J. Econometrics*, 188(2):447–465, 2015.
- [3] Victor Chernozhukov, Denis Chetverikov, and Kengo Kato. Gaussian approximation of suprema of empirical processes. *Ann. Statist.*, 42(4):1564–1597, 2014.
- [4] Jérôme Dedecker and Sana Louhichi. Maximal inequalities and empirical central limit theorems. In *Empirical process techniques for dependent data*, pages 137–159. Birkhäuser Boston, Boston, MA, 2002.
- [5] Christian Francq and Jean-Michel Zakoïan. A central limit theorem for mixing triangular arrays of variables whose dependence is allowed to grow with the sample size. *Econometric Theory*, 21(6):1165–1171, 2005.
- [6] Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinite-horizon off-policy estimation. In *Advances in Neural Information Processing Systems*, pages 5356–5366, 2018.
- [7] Chengchun Shi, Xiaoyu Wang, Shikai Luo, Rui Song, Hongtu Zhu, and Jieping Ye. A reinforcement learning framework for time-dependent causal effects evaluation in a/b testing. *arXiv preprint arXiv:2002.01711*, 2020.