

A Mathematical preliminaries

Proposition 2 (L^r Convergence Theorem, Loeve (1977)). *Let $0 < r < \infty$, suppose that $\mathbb{E}[|a_n|^r] < \infty$ for all n and that $a_n \xrightarrow{P} a$ as $n \rightarrow \infty$. The following are equivalent:*

- (i) $a_n \rightarrow a$ in L^r as $n \rightarrow \infty$;
- (ii) $\mathbb{E}[|a_n|^r] \rightarrow \mathbb{E}[|a|^r] < \infty$ as $n \rightarrow \infty$;
- (iii) $\{|a_n|^r, n \geq 1\}$ is uniformly integrable.

Proposition 3 (Weak Law of Large Numbers for Martingale, Hall et al. (2014)). *Let $\{S_n = \sum_{i=1}^n X_i, \mathcal{H}_t, t \geq 1\}$ be a martingale and $\{b_n\}$ a sequence of positive constants with $b_n \rightarrow \infty$ as $n \rightarrow \infty$. Then, writing $X_{ni} = X_i \mathbb{1}[|X_i| \leq b_n]$, $1 \leq i \leq n$, we have that $b_n^{-1} S_n \xrightarrow{P} 0$ as $n \rightarrow \infty$ if*

- (i) $\sum_{i=1}^n P(|X_i| > b_n) \rightarrow 0$;
- (ii) $b_n^{-1} \sum_{i=1}^n \mathbb{E}[X_{ni} | \mathcal{H}_{t-1}] \xrightarrow{P} 0$, and;
- (iii) $b_n^{-2} \sum_{i=1}^n \{\mathbb{E}[X_{ni}^2] - \mathbb{E}[\mathbb{E}[X_{ni} | \mathcal{H}_{t-1}]]^2\} \rightarrow 0$.

Remark 6. *The weak law of large numbers for martingale holds when the random variable is bounded by a constant.*

B Proof of Theorem 1

Proof of Theorem 1. We show asymptotic normality of

$$\widehat{R}_T^{\text{ADR}}(\pi^e) = \frac{1}{T} \sum_{t=1}^T \left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) + \phi_2(X_t; \hat{f}_{t-1}) \right\},$$

where

$$\begin{aligned} \phi_1(X_t, A_t, Y_t; g, f) &= \sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] (Y_t - f(a, X_t))}{g(a|X_t)} \\ \phi_2(X_t; f) &= \sum_{a=1}^K \pi^e(a|X_t) f(a, X_t). \end{aligned}$$

Let us define an AIPW estimator with $\hat{f} = f^*$ as

$$\widehat{R}^*(\pi^e) = \frac{1}{T} \sum_{t=1}^T \left\{ \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) + \phi_2(X_t; f^*) \right\}.$$

We decompose $\sqrt{T} (R^{\text{ADRE}}(\pi^e) - R(\pi^e))$ as

$$\sqrt{T} (R^{\text{ADRE}}(\pi^e) - R(\pi^e)) = \sqrt{T} \left(\widehat{R}_T^{\text{ADR}}(\pi^e) - \widehat{R}^*(\pi^e) + \widehat{R}^*(\pi^e) - R(\pi^e) \right).$$

From Proposition 1 of Kato et al. (2020) and Assumption 1 and 3, because $\sqrt{T} (\widehat{R}^*(\pi^e) - R(\pi^e))$ follows asymptotic normal distribution, we want to show

$$\widehat{R}_T^{\text{ADR}}(\pi^e) - \widehat{R}^*(\pi^e) = o_p(1/\sqrt{T}).$$

Here, we have

$$\begin{aligned}
& \widehat{R}_T^{\text{ADR}}(\pi^e) - \widehat{R}^*(\pi^e) \\
&= \frac{1}{T} \sum_{t=1}^T \left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \\
&\quad \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] \right. \\
&\quad \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right] \right\} \\
&\quad + \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) | \Omega_{t-1} \right] + \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) | \Omega_{t-1} \right] \\
&\quad - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_2(X_t; f^*) | \Omega_{t-1} \right].
\end{aligned}$$

In the following parts, we separately show that

$$\begin{aligned}
& \sqrt{T} \frac{1}{T} \sum_{t=1}^T \left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \\
&\quad \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] \right. \\
&\quad \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right] \right\} \\
&= o_p(1);
\end{aligned} \tag{1}$$

and

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) | \Omega_{t-1} \right] + \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) | \Omega_{t-1} \right] \\
&\quad - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_2(X_t; f^*) | \Omega_{t-1} \right] = o_p(1/\sqrt{T}).
\end{aligned} \tag{2}$$

Proof of (1). For any $\varepsilon > 0$, to show that

$$\begin{aligned}
& \mathbb{P} \left(\left| \sqrt{T} \frac{1}{T} \sum_{t=1}^T \left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \right. \\
&\quad \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] \right. \right. \\
&\quad \left. \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right] \right\} \right| > \varepsilon \right) \\
&\rightarrow 0,
\end{aligned}$$

we show that the mean is 0 and the variance of the component converges to 0. Then, from the Chebyshev's inequality, this result yields the statement.

The mean is calculated as

$$\begin{aligned}
& \sqrt{T} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \\
& \quad \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \mid \Omega_{t-1} \right] \right. \right. \\
& \quad \left. \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \mid \Omega_{t-1} \right] \right\} \right] \\
&= \sqrt{T} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\mathbb{E} \left[\left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \right. \\
& \quad \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \mid \Omega_{t-1} \right] \right. \right. \\
& \quad \left. \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \mid \Omega_{t-1} \right] \right\} \mid \Omega_{t-1} \right] \right] \\
&= 0
\end{aligned}$$

Because the mean is 0, the variance is

$$\begin{aligned}
& \text{Var} \left(\sqrt{T} \frac{1}{T} \sum_{t=1}^T \left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \\
& \quad \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \mid \Omega_{t-1} \right] \right. \right. \\
& \quad \left. \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \mid \Omega_{t-1} \right] \right\} \right) \\
&= \mathbb{E} \left[\left(\sqrt{T} \frac{1}{T} \sum_{t=1}^T \left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \right. \\
& \quad \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \mid \Omega_{t-1} \right] \right. \right. \\
& \quad \left. \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \mid \Omega_{t-1} \right] \right\} \right)^2 \right] \\
&= \frac{1}{T} \mathbb{E} \left[\left(\sum_{t=1}^T \left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \right. \\
& \quad \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \mid \Omega_{t-1} \right] \right. \right. \\
& \quad \left. \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \mid \Omega_{t-1} \right] \right\} \right)^2 \right].
\end{aligned}$$

Therefore, we have

$$\begin{aligned}
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \\
&\quad \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \middle| \Omega_{t-1} \right] \right. \right. \\
&\quad \left. \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \middle| \Omega_{t-1} \right] \right)^2 \right] \\
&+ \frac{2}{T} \sum_{t=1}^{T-1} \sum_{s=t+1}^T \mathbb{E} \left[\left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \\
&\quad \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \middle| \Omega_{t-1} \right] \right. \right. \\
&\quad \left. \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \middle| \Omega_{t-1} \right] \right) \right. \\
&\quad \times \left(\phi_1(X_s, A_s, Y_s; \hat{g}_{s-1}, \hat{f}_{s-1}) - \phi_1(X_s, A_s, Y_s; \pi_{s-1}, f^*) \right. \\
&\quad \left. \left. - \mathbb{E} \left[\phi_1(X_s, A_s, Y_s; \hat{g}_{s-1}, \hat{f}_{s-1}) - \phi_1(X_s, A_s, Y_s; \pi_{s-1}, f^*) \middle| \Omega_{s-1} \right] \right. \right. \\
&\quad \left. \left. + \phi_2(X_s; \hat{f}_{s-1}) - \phi_2(X_s; f^*) - \mathbb{E} \left[\phi_2(X_s; \hat{f}_{s-1}) - \phi_2(X_s; f^*) \middle| \Omega_{s-1} \right] \right) \right].
\end{aligned}$$

For $s > t$, we can vanish the covariance terms as

$$\begin{aligned}
&\mathbb{E} \left[\left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \\
&\quad \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \middle| \Omega_{t-1} \right] \right. \right. \\
&\quad \left. \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \middle| \Omega_{t-1} \right] \right) \right. \\
&\quad \times \left(\phi_1(X_s, A_s, Y_s; \hat{g}_{s-1}, \hat{f}_{s-1}) - \phi_1(X_s, A_s, Y_s; \pi_{s-1}, f^*) \right. \\
&\quad \left. \left. - \mathbb{E} \left[\phi_1(X_s, A_s, Y_s; \hat{g}_{s-1}, \hat{f}_{s-1}) - \phi_1(X_s, A_s, Y_s; \pi_{s-1}, f^*) \middle| \Omega_{s-1} \right] \right. \right. \\
&\quad \left. \left. + \phi_2(X_s; \hat{f}_{s-1}) - \phi_2(X_s; f^*) - \mathbb{E} \left[\phi_2(X_s; \hat{f}_{s-1}) - \phi_2(X_s; f^*) \middle| \Omega_{s-1} \right] \right) \right] \\
&= \mathbb{E} \left[U \mathbb{E} \left[\left(\phi_1(X_s, A_s, Y_s; \hat{g}_{s-1}, \hat{f}_{s-1}) - \phi_1(X_s, A_s, Y_s; \pi_{s-1}, f^*) \right. \right. \right. \\
&\quad \left. \left. - \mathbb{E} \left[\phi_1(X_s, A_s, Y_s; \hat{g}_{s-1}, \hat{f}_{s-1}) - \phi_1(X_s, A_s, Y_s; \pi_{s-1}, f^*) \middle| \Omega_{s-1} \right] \right. \right. \\
&\quad \left. \left. + \phi_2(X_s; \hat{f}_{s-1}) - \phi_2(X_s; f^*) - \mathbb{E} \left[\phi_2(X_s; \hat{f}_{s-1}) - \phi_2(X_s; f^*) \middle| \Omega_{s-1} \right] \right) \middle| \Omega_{s-1} \right] \right] \\
&= 0,
\end{aligned}$$

where

$$\begin{aligned}
U = & \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \\
& - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] \\
& + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \\
& \left. - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right] \right).
\end{aligned}$$

Therefore, the variance is calculated as

$$\begin{aligned}
& \text{Var} \left(\sqrt{T} \frac{1}{T} \sum_{t=1}^T \left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \\
& \quad - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] \\
& \quad + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \\
& \quad \left. \left. - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right] \right\} \right) \\
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \\
& \quad - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] \\
& \quad + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \\
& \quad \left. \left. - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right] \right)^2 \right] \\
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\mathbb{E} \left[\left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \right. \\
& \quad - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] \\
& \quad + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \\
& \quad \left. \left. \left. - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right] \right)^2 \middle| \Omega_{t-1} \right] \right] \\
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\text{Var} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right) \right] \\
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\text{Var} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right) \right] \\
& \quad + \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\text{Var} \left(\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right) \right] \\
& \quad + \frac{2}{T} \sum_{t=1}^T \mathbb{E} \left[\text{Cov} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*), \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right) \right].
\end{aligned}$$

Then, we want to show that

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\text{Var} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right) \right] \rightarrow 0, \quad (3)$$

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\text{Var} \left(\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right) \right] \rightarrow 0, \quad (4)$$

$$\frac{2}{T} \sum_{t=1}^T \mathbb{E} \left[\text{Cov} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*), \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right) \right] \rightarrow 0 \quad (5)$$

For showing (3)–(5), we consider showing

$$\text{Var} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right) = o_p(1), \quad (6)$$

$$\text{Var} \left(\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right) = o_p(1) \quad (7)$$

$$\text{Cov} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*), \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right) = o_p(1), \quad (8)$$

The first equation (6) is shown as

$$\begin{aligned} & \text{Var} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right) \\ & \leq \mathbb{E} \left[\left\{ \sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] (Y_t - \hat{f}_{t-1}(a, X_t))}{\hat{g}_{t-1}(a|X_t)} - \sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] (Y_t - f^*(a, X_t))}{\pi_{t-1}(a|X_t)} \right\}^2 | \Omega_{t-1} \right] \\ & = \mathbb{E} \left[\left\{ \sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] (Y_t - \hat{f}_{t-1}(a, X_t))}{\hat{g}_{t-1}(a|X_t)} - \sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] (Y_t - f^*(a, X_t))}{\hat{g}_{t-1}(a|X_t)} \right. \right. \\ & \quad \left. \left. + \sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] (Y_t - f^*(a, X_t))}{\hat{g}_{t-1}(a|X_t)} - \sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] (Y_t - f^*(a, X_t))}{\pi_{t-1}(a|X_t)} \right\}^2 | \Omega_{t-1} \right] \\ & \leq 2\mathbb{E} \left[\left\{ \sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] (Y_t - \hat{f}_{t-1}(a, X_t))}{\hat{g}_{t-1}(a|X_t)} - \sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] (Y_t - f^*(a, X_t))}{\hat{g}_{t-1}(a|X_t)} \right\}^2 w \right] z | \Omega_{t-1} \\ & \quad + 2\mathbb{E} \left[\left\{ \sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] (Y_t - f^*(a, X_t))}{\hat{g}_{t-1}(a|X_t)} - \sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] (Y_t - f^*(a, X_t))}{\pi_{t-1}(a|X_t)} \right\}^2 | \Omega_{t-1} \right] \\ & \leq 2C \|f^* - \hat{f}_{t-1}\|_2^2 + 2 \times 4C^2 \|\hat{g}_{t-1} - \pi_{t-1}\|_2^2 = o_p(1), \end{aligned}$$

where $C > 0$ is a constant. Here, we have used a parallelogram law from the third line to the fourth line. We have used $|\hat{f}_{t-1}| < C$, and $0 < \frac{\pi^e}{\pi_{t-1}} < C$, convergence of π_{t-1} and convergence rate conditions from the third line to the fourth line. Then, from the L^r convergence theorem (Proposition 2) and the boundedness of the random variables, we can show that as $t \rightarrow \infty$,

$$\begin{aligned} & \mathbb{E} \left[\text{Var} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right) \right] \\ & \leq \mathbb{E} \left[\left| \text{Var} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right) \right| \right] \rightarrow 0. \end{aligned}$$

Therefore, for any $\epsilon > 0$, there exists a constant $C > 0$ such that

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\text{Var} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \middle| \Omega_{t-1} \right) \right] \leq C/T + \epsilon.$$

The second equation (7) is derived by Jensen's inequality, and we show (4) as well as (3) by using L^r convergence theorem.

Next, we show the third equation (8) as

$$\begin{aligned} & \text{Cov} \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*), \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \middle| \Omega_{t-1} \right) \\ & \leq \left| \mathbb{E} \left[\left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \right. \\ & \quad \left. \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \middle| \Omega_{t-1} \right] \right) \right. \right. \\ & \quad \left. \left. \times \left(\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \right] \right) \middle| \Omega_{t-1} \right] \right| \\ & \leq \mathbb{E} \left[\left| \left(\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \right. \\ & \quad \left. \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \middle| \Omega_{t-1} \right] \right) \right. \right. \\ & \quad \left. \left. \times \left(\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) \right] \right) \right| \Omega_{t-1} \right] \\ & \leq C \mathbb{E} \left[\left| \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \\ & \quad \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \middle| \Omega_{t-1} \right] \right| \Omega_{t-1} \right] \\ & = o_p(1), \end{aligned}$$

where $C > 0$ is a constant. From the second to third line, we used Jensen's inequality. From the fourth to fifth line, we used consistencies of \hat{f}_{t-1} and \hat{g}_{t-1} , which imply that for all $X_t \in \mathcal{X}$,

$$\begin{aligned} & \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \\ & = \sum_{a=1}^K \left(\frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] \left(Y_t - \hat{f}_{t-1}(a, X_t) \right)}{\hat{g}_{t-1}(a|X_t)} - \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] \left(Y_t - f^*(a, X_t) \right)}{\pi_{t-1}(a|X_t)} \right) \\ & \leq \sum_{a=1}^K \left| \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] \left(Y_t - \hat{f}_{t-1}(a, X_t) \right)}{\hat{g}_{t-1}(a|X_t)} - \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] \left(Y_t - f^*(a, X_t) \right)}{\pi_{t-1}(a|X_t)} \right| \\ & \leq C \sum_{a=1}^K \left| \pi_{t-1}(a|X_t) \left(Y_t - \hat{f}_{t-1}(a, X_t) \right) - \hat{g}_{t-1}(a|X_t) \left(Y_t - f^*(a, X_t) \right) \right| \\ & \leq C \sum_{a=1}^K \left| \pi_{t-1}(a|X_t) - \hat{g}_{t-1}(a|X_t) \right| \\ & \quad - C \sum_{a=1}^K \left| \pi_{t-1}(a|X_t) \hat{f}_{t-1}(a, X_t) - \hat{g}_{t-1}(a|X_t) \hat{f}_{t-1}(a, X_t) + \hat{g}_{t-1}(a|X_t) \hat{f}_{t-1}(a, X_t) - \hat{g}_{t-1}(a|X_t) f^*(a, X_t) \right| \\ & \leq C \sum_{a=1}^K \left| \pi_{t-1}(a|X_t) - \hat{g}_{t-1}(a|X_t) \right| - C \sum_{a=1}^K \left| \hat{f}_{t-1}(a, X_t) - f^*(a, X_t) \right| = o_p(1), \end{aligned}$$

where $C > 0$ is a constant.

Thus, from (3)–(5), the variance of the bias term converges to 0. Then, from Chebyshev's inequality,

$$\begin{aligned}
& \mathbb{P} \left(\left| \sqrt{T} \frac{1}{T} \sum_{t=1}^T \left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \right. \\
& \quad \left. \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] \right. \right. \right. \\
& \quad \left. \left. \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right] \right\} \right| > \varepsilon \right) \\
& \leq \text{Var} \left(\sqrt{T} \frac{1}{T} \sum_{t=1}^T \left\{ \phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) \right. \right. \\
& \quad \left. \left. - \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) - \phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] \right. \right. \\
& \quad \left. \left. \left. + \phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) - \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) - \phi_2(X_t; f^*) | \Omega_{t-1} \right] \right\} \right) / \varepsilon^2 \\
& \rightarrow 0.
\end{aligned}$$

Proof of (2).

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) | \Omega_{t-1} \right] + \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) | \Omega_{t-1} \right] \\
& \quad - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_2(X_t; f^*) | \Omega_{t-1} \right] \\
& = \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] \left(Y_t - \hat{f}_{t-1}(a, X_t) \right)}{\hat{g}_{t-1}(a|X_t)} \Big| \Omega_{t-1} \right] \\
& \quad + \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\sum_{a=1}^K \pi^e(a, X_t) \hat{f}_{t-1}(a, X_t) \Big| \Omega_{t-1} \right] \\
& \quad - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] \left(Y_t - f^*(a, X_t) \right)}{\pi_{t-1}(a|X_t, \Omega_{t-1})} \Big| \Omega_{t-1} \right] \\
& \quad - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\sum_{a=1}^K \pi^e(a, X_t) f^*(a, X_t) \Big| \Omega_{t-1} \right].
\end{aligned} \tag{9}$$

Because (9) is 0,

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \hat{g}_{t-1}, \hat{f}_{t-1}) | \Omega_{t-1} \right] + \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_2(X_t; \hat{f}_{t-1}) | \Omega_{t-1} \right] \\
& - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_1(X_t, A_t, Y_t; \pi_{t-1}, f^*) | \Omega_{t-1} \right] - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\phi_2(X_t; f^*) | \Omega_{t-1} \right] \\
& = \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] \left(Y_t - \hat{f}_{t-1}(a, X_t) \right)}{\hat{g}_{t-1}(a|X_t)} \middle| \Omega_{t-1} \right] \\
& + \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\sum_{a=1}^K \pi^e(a, X_t) \hat{f}_{t-1}(a, X_t) \middle| \Omega_{t-1} \right] - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\sum_{a=1}^K \pi^e(a, X_t) f^*(a, X_t) \middle| \Omega_{t-1} \right] \\
& = \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\sum_{a=1}^K \frac{\pi^e(a|X_t) \mathbb{1}[A_t = a] \left(Y_t - \hat{f}_{t-1}(a, X_t) \right)}{\hat{g}_{t-1}(a|X_t)} \middle| \Omega_{t-1} \right] \\
& - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\sum_{a=1}^K \pi^e(a, X_t) \left(f^*(a, X_t) - \hat{f}_{t-1}(a, X_t) \right) \middle| \Omega_{t-1} \right] \\
& = \frac{1}{T} \sum_{t=1}^T \sum_{a=1}^K \mathbb{E} \left[\mathbb{E} \left[\frac{\pi^e(a|X_t) \pi_{t-1}(a|X_t, \Omega_{t-1}) \left(f^*(a, X_t) - \hat{f}_{t-1}(a, X_t) \right)}{\hat{g}_{t-1}(a|X_t)} \right. \right. \\
& \quad \left. \left. - \pi^e(a, X_t) \left(f^*(a, X_t) - \hat{f}_{t-1}(a, X_t) \right) \middle| X_t, \Omega_{t-1} \right] \middle| \Omega_{t-1} \right] \\
& \leq \frac{1}{T} \sum_{t=1}^T \sum_{a=1}^K \left| \mathbb{E} \left[\frac{\pi^e(a|X_t) \left(\pi_{t-1}(a|X_t) - \hat{g}_{t-1}(a|X_t) \right) \left(f^*(a, X_t) - \hat{f}_{t-1}(a, X_t) \right)}{\hat{g}_{t-1}(a|X_t)} \middle| \Omega_{t-1} \right] \right|.
\end{aligned}$$

By using Hölder's inequality $\|fg\|_1 \leq \|f\|_2 \|g\|_2$, for a constant $C > 0$, we have

$$\begin{aligned}
& \leq \frac{C}{T} \sum_{t=1}^T \left\| \pi_{t-1}(a|X_t, \Omega_{t-1}) - \hat{g}_{t-1}(a|X_t) \right\|_2 \left\| f^*(a, X_t) - \hat{f}_{t-1}(a, X_t) \right\|_2 \\
& = \frac{C}{T} \sum_{t=1}^T o_p(t^{-p}) o_p(t^{-q}) \\
& = \frac{C}{T} \sum_{t=1}^T o_p(t^{-1/2}).
\end{aligned}$$

□

C Adaptive-fitting and batched samples

Section C.1 supplements the description of adaptive-fitting. Next, in Section C.2, we introduce the AIPW estimator when the samples are given in batch form, and the true logging policy is given π_t . This is essentially the same as the generalized method of moment (GMM), which gives an asymptotically normal estimator for martingale difference sequences (MDS), and we do not use adaptive-fitting. Based on this estimator, in Section C.3, we introduce the ADR estimator when the data is given in batch form, but the true logging policy π_t is not given.

C.1 Details of adaptive-fitting

As Section 5, let us define the parameter of interest θ_0 that satisfies $\mathbb{E}[\psi(W_t; \theta_0, \eta_0)] = 0$, where $\{W_t\}_{t=1}^T$ are observations, η_0 is a nuisance parameter, and ψ is a score function. Let us define two

estimators $\hat{\theta}_T$ and $\check{\theta}_T$ as $\frac{1}{T} \sum_{t=1}^T \psi(W_t, \hat{\theta}_T, \eta_{t-1}) = 0$ and $\frac{1}{T} \sum_{t=1}^T \psi(W_t, \check{\theta}_T, \eta_0) = 0$. Suppose that $\check{\theta}_T$ is an asymptotically normal estimator of θ_0 . Then, if $\check{\theta}_T - \hat{\theta}_T$ converges to 0 with convergence rate $o_p(1/\sqrt{T})$, $\hat{\theta}_T$ is also an asymptotically normal estimator. In general, we cannot obtain such a fast convergence rate. However, by using double robustness, we can obtain the convergence rate. In ADR estimator, this conditions appears as $\|\hat{g}_{t-1}(a|X_t) - \pi_{t-1}(a|X_t, \Omega_{t-1})\|_2 = o_p(t^{-p})$, and $\|\hat{f}_{t-1}(a, X_t) - f^*(a, X_t)\|_2 = o_p(t^{-q})$, where $p, q > 0$ such that $p + q = 1/2$, and the expectation of the norm is taken over X_t . This allows us to obtain $o_p(1/\sqrt{T})$ of the asymptotic bias. The image of the vanishing asymptotic bias is shown in Figure 3.

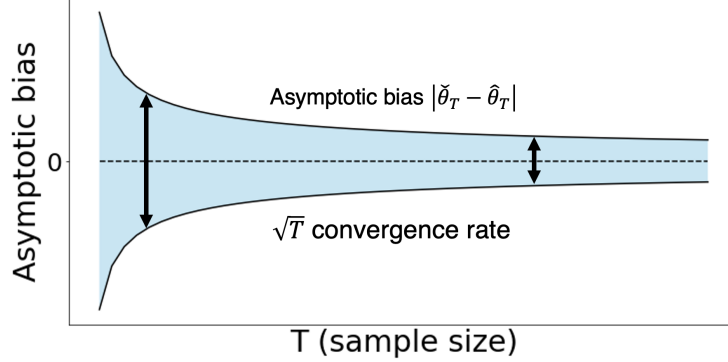


Figure 3: Convergence of asymptotic bias $|\check{\theta}_T - \hat{\theta}_T|$.

C.2 AIPW estimator with batched samples when the true logging policy is known

The proposed adaptive-fitting can be applied when batched samples are given. Let M denote the number of batch updates and $\tau \in I = \{1, 2, \dots, M\}$ denotes the batch index. For $\tau \in I$, the probability is updated at a period t_τ , where $t_\tau - t_{\tau-1} = Tr_\tau$, using samples $\{(X_t, Y_t, A_t)\}_{t=t_{\tau-1}}^{t_\tau}$, where $r_1 + r_2 + \dots + r_M = 1$ and $t_0 = 0$. Thus, in addition to the DGP, we assume that

$$\{(X_t, A_t, Y_t)\}_{t=t_{\tau-1}}^{t_\tau} \stackrel{i.i.d.}{\sim} p(x)\pi_\tau(a|x, \Omega_{t_{\tau-1}})p_a(y|x),$$

where $\pi_\tau(a|x, \Omega_{t_{\tau-1}})$ denotes the probability of choosing an action updated based on samples until the period $t_{\tau-1}$.

We consider asymptotic properties based on the assumption of $t_\tau - t_{\tau-1} \rightarrow \infty$ as $T \rightarrow \infty$ for fixed τ . This strategy is the same as Zhang et al. (2020). Because $\{(X_t, A_t, Y_t)\}_{t=t_{\tau-1}}^{t_\tau}$ is i.i.d., we can use the standard limit theorems for the partial sum of the samples to obtain an asymptotically normal estimator of $\theta_0 = R(\pi^e)$. However, we also have the motivation to use all samples together to increase the efficiency of the estimator. Therefore, by using the GMM, we propose an estimator of θ_0 considering the sample averages of each block as an empirical moment conditions. Although we cannot use standard CLT, we can apply the CLT for the MDS by appropriately constructing an estimator.

For an index of batch $\tau \in I$, a function $f \in \mathcal{F}$ such that $f : \mathcal{A} \times \mathcal{X} \rightarrow \mathbb{R}$ and an evaluation policy $\pi^e \in \Pi$, we define a function h_t as

$$h_t(x, k, y; \tau, R, f, \pi_\tau, \pi^e) = \frac{1}{r_\tau} \xi_t(x, k, y; \tau, R, f, \pi_\tau, \pi^e) \mathbb{1}[t_{\tau-1} < t \leq t_\tau],$$

where $R \in \mathbb{R}$, $\xi_t(x, k, y; \tau, R, f, \pi_\tau, \pi^e) := \phi_t(x, k, y; \tau, f, \pi_\tau, \pi^e) - R$ and

$$\phi_t(x, k, y; \tau, f, \pi_\tau, \pi^e) := \sum_{a=1}^K \pi^e(a|x) \left\{ \frac{\mathbb{1}[k=a] \{y - f(a, x)\}}{\pi_\tau(a|x, \Omega_{t_{\tau-1}})} + f(a, x) \right\}.$$

Let us note that the sequence $\{h_t(X_t, A_t, Y_t; \tau, R(\pi^e), \hat{f}_{t-1}, \pi_\tau, \pi^e)\}_{t=1}^T$ is an MDS: for $h_t(X_t, A_t, Y_t; \tau, R(\pi^e), \hat{f}_{t-1}, \pi_\tau, \pi^e)$, by using $\mathbb{E}[\mathbb{1}[A_t = a] | \mathbb{H}_{t-1}] = \pi_\tau(a | X_t, \Omega_{t-1})$, we have

$$\begin{aligned} & \mathbb{E} \left[h_t(X_t, A_t, Y_t; \tau, R(\pi^e), \hat{f}_{t-1}, \pi_\tau, \pi^e) | \Omega_{t-1} \right] \\ &= \mathbb{E} \left[\frac{\mathbb{1}[t_{\tau-1} < t \leq t_\tau]}{r_\tau} \xi_t(x, k, y; \tau, R(\pi^e), \hat{f}_{t-1}, \pi_\tau, \pi^e) | \Omega_{t-1} \right] \\ &= \frac{\mathbb{1}[t_{\tau-1} < t \leq t_\tau]}{r_\tau} \mathbb{E} \left[\xi_t(x, k, y; \tau, R(\pi^e), \hat{f}_{t-1}, \pi_\tau, \pi^e) | \Omega_{t-1} \right] \\ &= \frac{\mathbb{1}[t_{\tau-1} < t \leq t_\tau]}{r_\tau} \times 0 = 0. \end{aligned}$$

Let us also define

$$\mathbf{h}_t \left(X_t, A_t, Y_t; R, \hat{f}_{t-1}, \pi_\tau, \pi^e \right) := \begin{pmatrix} h_t(X_t, A_t, Y_t; 1, R, \hat{f}_{t-1}, \pi_1, \pi^e) \\ h_t(X_t, A_t, Y_t; 2, R, \hat{f}_{t-1}, \pi_2, \pi^e) \\ \vdots \\ h_t(X_t, A_t, Y_t; M, R, \hat{f}_{t-1}, \pi_M, \pi^e) \end{pmatrix}.$$

Then, the sequence $\{\mathbf{h}_t \left(X_t, A_t, Y_t; R(\pi^e), \hat{f}_{t-1}, \pi_\tau, \pi^e \right)\}_{t=1}^T$ is an MDS with respect to $\{\Omega_t\}_{t=0}^{T-1}$, i.e.,

$$\mathbb{E} \left[\mathbf{h}_t \left(X_t, A_t, Y_t; R(\pi^e), \hat{f}_{t-1}, \pi_\tau, \pi^e \right) | \Omega_{t-1} \right] = \mathbf{0}.$$

Using the sequence $\{\mathbf{h}_t \left(X_t, A_t, Y_t; R, \hat{f}_{t-1}, \pi_\tau, \pi^e \right)\}_{t=1}^T$, we define an estimator of $R(\pi^e)$ as

$$\hat{R}_T^{\text{batch}}(\pi^e) = \arg \min_{R \in \mathbb{R}} (\hat{\mathbf{q}}_T(R))^\top \hat{W}_T(\hat{\mathbf{q}}_T(R)), \quad (10)$$

where $\hat{\mathbf{q}}_T(R) = \frac{1}{T} \sum_{t=1}^T \mathbf{h}_t \left(X_t, A_t, Y_t; R, \hat{f}_{t-1}, \pi_\tau, \pi^e \right)$ and \hat{W}_T is a data-dependent $(M \times M)$ -dimensional positive semi-definite matrix. Let us note that the estimator defined in Eq. (10) is an application of GMM with the moment condition

$$\mathbf{q}(R(\pi^e)) = \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \mathbf{h}_t \left(X_t, A_t, Y_t; R(\pi^e), \hat{f}_{t-1}, \pi_\tau, \pi^e \right) \right] = \mathbf{0}.$$

For the minimization problem defined in Eq. (10), we can analytically calculate the minimizer as

$$\hat{R}_T^{\text{batch}}(\pi^e) = w_T^\top D_T(\pi^e),$$

where $w_T = (w_{T,1} \cdots w_{T,M})^\top$ is an M -dimensional vector such that $\sum_{\tau=1}^M w_{T,\tau} = 1$, and

$$D_T(\pi^e) = \begin{pmatrix} \frac{1}{t_1} \sum_{t=1}^{t_1} \phi_t(X_t, A_t, Y_t; 1, \hat{f}_{t-1}, \pi_1, \pi^e) \\ \frac{1}{t_2 - t_1} \sum_{t=t_1+1}^{t_2} \phi_t(X_t, A_t, Y_t; 2, \hat{f}_{t-1}, \pi_2, \pi^e) \\ \vdots \\ \frac{1}{T - t_{M-1}} \sum_{t=t_{M-1}+1}^T \phi_t(X_t, A_t, Y_t; M, \hat{f}_{t-1}, \pi_M, \pi^e) \end{pmatrix}.$$

Here, we show the asymptotic normality of the proposed estimator $\hat{\theta}_T^{\text{OPE}}$.

Theorem 3 (Asymptotic distribution the proposed estimator). *Suppose that*

- (i) $w_T = (w_{T,1} \cdots w_{T,M})^\top \xrightarrow{P} w = (w_1 \cdots w_M)^\top$;
- (ii) $w_{T,\tau} > 0$ and $\sum_{\tau=1}^M w_{T,\tau} = 1$.

Then, under Assumptions 1, 3, 5, 7,

$$\sqrt{T}(\hat{R}_T^{\text{batch}}(\pi^e) - R(\pi^e)) \xrightarrow{d} \mathcal{N}(0, \sigma^2),$$

where $\sigma^2 = \sum_{\tau=1}^M w_\tau \Psi(\pi^e, \pi_\tau)$.

In the proposed method, we construct a moment condition using martingale difference sequences. On the other hand, for some readers, using martingale difference sequences may look unnecessary because samples are i.i.d in each block between $t_{\tau-1}$ and t_τ . Therefore, such readers also might feel that we can use $f_T(a, x)$, which is an estimator of $\mathbb{E}[Y(a) | x]$ using samples until T -th period, without going through constructing several estimators $\{f_{t_\tau}\}_{\tau=0}^{M-1}$. However, in that case, it is difficult to guarantee the asymptotic normality of the proposed estimator. For example, we can consider Cramér-Wold theorem to consider this problem. This motivations shares with the GMM and the method of Zhang et al. (2020).

The proof of Theorem 3 is shown as follows.

Choice of weight w_T . We discuss the choice of weight w_T . A naive choice is weighting the moment conditions equality; that is, $w_{\tau,T} = \frac{1}{M}$. Next, we consider an efficient weight w_T that minimizes the asymptotic variance of $\hat{R}_T^{\text{batch}}(\pi^e)$. In the GMM, the τ -th element of the efficient weight is given as $w_\tau^* = \frac{1}{\Psi(\pi^e, \pi_\tau)} / \sum_{\tau'=1}^M \frac{1}{\Psi(\pi^e, \pi_{\tau'})}$ (Hamilton, 1994). Here, we use the orthogonality among moment conditions; that is, zero covariance. In this case, the asymptotic variance becomes $1 / \sum_{\tau'=1}^M \frac{1}{\sigma_{\tau'}^2}$. Therefore, for gaining efficiency, we use a weight $\hat{w}_{T,\tau} = \frac{1}{\hat{\Psi}_T(\pi^e, \pi_\tau)} / \sum_{\tau'=1}^M \frac{1}{\hat{\Psi}_T(\pi^e, \pi_{\tau'})}$, where $\hat{\Psi}_T(\pi^e, \pi_\tau)$ is an estimator of $\Psi(\pi^e, \pi_\tau)$.

Proof of Theorem 3. Instead of $\hat{R}_T^{\text{batch}}(\pi^e) = w_T D_T(\pi^e)$, from the original formulation Eq. (10), we consider an estimator $\hat{R}_T^{\text{batch}}(\pi^e) = (I^\top \hat{W}_T I)^{-1} I^\top \hat{W}_T D_T(\pi^e)$, where W_T is a $(M \times M)$ -dimensional positive-definite matrix. Let us note that $w_T = (I^\top \hat{W}_T I)^{-1} I^\top \hat{W}_T$. We prove the following theorem, which is a generalized statement of Theorem 3.

Theorem 4 (Asymptotic distribution the AIPW estimator under batch update). *Suppose that*

- (i) $\hat{W}_T \xrightarrow{p} W$;
- (ii) W is a positive definite;

Then, under Assumptions 1, 3, 5, 7,

$$\sqrt{T}(\hat{R}_T^{\text{batch}}(\pi^e) - R(\pi^e)) \xrightarrow{d} \mathcal{N}(0, \sigma^2),$$

where $\sigma^2 = (I^\top W I)^{-1} I^\top W \Sigma W^\top I (I^\top W I)^{-1}$ and Σ is a $(M \times M)$ diagonal matrix such that the $(\tau \times \tau)$ -element is $\Psi(\pi^e, \pi_\tau)$.

Proof. Let us define $\hat{q}_T(R(\pi^e)) = (\hat{q}_{1,T}(R(\pi^e)) \hat{q}_{2,T}(R(\pi^e)) \cdots \hat{q}_{M,T}(R(\pi^e)))$, where

$$\begin{aligned} \hat{q}_{\tau,T}(R(\pi^e)) = & \frac{1}{T} \frac{1}{r_1} \sum_{t=1}^T \left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - \hat{f}_{t-1}(a, X_t)\}}{\pi_1(a | x, \Omega_0)} + \pi^e(a | X_t) \hat{f}_{t-1}(a, X_t) \right\} - R(\pi^e) \right) \\ & \times \mathbb{1}[t_0 = 0 < t \leq t_1]. \end{aligned}$$

For $\sqrt{T}(\hat{R}_T^{\text{batch}}(\pi^e) - R(\pi^e)) = (I^\top \hat{W}_T I)^{-1} I^\top \hat{W}_T \sqrt{T} \hat{q}_T(R(\pi^e))$, we show that

$$\sqrt{T} \hat{q}_T(R(\pi^e)) \xrightarrow{d} \mathcal{N}(0, \Sigma),$$

where Σ is a diagonal matrix such that the (τ, τ) -element is

$$\frac{1}{r_\tau} \mathbb{E} \left[\sum_{a=1}^K \frac{(\pi^e(a | X))^2 \text{Var}(Y(a) | X)}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \left(\sum_{a=1}^K \pi^e(a | X) \mathbb{E}[Y(a) | X] - R(\pi^e) \right)^2 \right].$$

Then, from Slutsky Theorem, we can show that

$$(I^\top \hat{W}_T I)^{-1} I^\top \hat{W}_T \sqrt{T} \hat{\mathbf{q}}_T(R(\pi^e)) \xrightarrow{d} \mathcal{N}(0, (I^\top W I)^{-1} I^\top W \Sigma W^\top I (I^\top W I)^{-1}).$$

To show this result, we use the CLT for MDS by checking the following conditions:

(a) $(1/T) \sum_{t=1}^T \Sigma_t \rightarrow \Sigma$, where

$$\Sigma_t = \mathbb{E} \left[(\mathbf{h}_t(X_t, A_t, Y_t; R(\pi^e), f_{t-1}, \pi_\tau, \pi^e)) (\mathbf{h}_t(X_t, A_t, Y_t; R(\pi^e), f_{t-1}, \pi_\tau, \pi^e))^\top \right];$$

(b) $\mathbb{E} [\tilde{h}_t(i, R(\pi^e), f_{t-1}, \pi^e) \tilde{h}_t(j, R(\pi^e), f_{t-1}, \pi^e) \tilde{h}_t(a, R(\pi^e), f_{t-1}, \pi^e) \tilde{h}_t(l, R(\pi^e), f_{t-1}, \pi^e)] < \infty$ for $i, j, k, l \in I$, where $\tilde{h}_t(a, R(\pi^e), f_{t-1}, \pi^e) = h_t^{\text{OPE}}(X_t, A_t, Y_t; k, R(\pi^e), f_k, \pi^e)$ for $k \in I$;

(c) $\frac{1}{T} \sum_{t=1}^T (\mathbf{h}_t(X_t, A_t, Y_t; R(\pi^e), f_{t-1}, \pi_\tau, \pi^e)) (\mathbf{h}_t(X_t, A_t, Y_t; R(\pi^e), f_{t-1}, \pi_\tau, \pi^e))^\top \xrightarrow{p} \Sigma$,

Note that the GMM shows the asymptotic normality for more general cases.

Step 1: Condition (a)

From

$$\Sigma_t = \mathbb{E} \left[(\mathbf{h}_t(X_t, A_t, Y_t; R(\pi^e), f_{t-1}, \pi_\tau, \pi^e)) (\mathbf{h}_t(X_t, A_t, Y_t; R(\pi^e), f_{t-1}, \pi_\tau, \pi^e))^\top \right],$$

the matrix $(1/T) \sum_{t=1}^T \Omega_t$ becomes a diagonal matrix such that the (τ, τ) -element is

$$\begin{aligned} & \frac{1}{r_\tau^2 T} \sum_{t=1}^T \mathbb{E} \left[\left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - f_{t-1}(a, X_t)\}}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) f_{t-1}(a, X_t) \right\} - R(\pi^e) \right)^2 \right. \\ & \quad \left. \times \mathbb{1}[t_{\tau-1} < t \leq t_\tau] \right]. \end{aligned}$$

For $\tau \in I$ and t such that $t_{\tau-1} < t \leq t_\tau$,

$$\begin{aligned} & \mathbb{E} \left[\left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - f_{t-1}(a, X_t)\}}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) f_{t-1}(a, X_t) \right\} - R(\pi^e) \right)^2 \right] \\ & - \mathbb{E} \left[\left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - \mathbb{E}[Y_t(a) | X_t]\}}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) \mathbb{E}[Y(a) | X_t] \right\} - R(\pi^e) \right)^2 \right] \\ & \leq \mathbb{E} \left[\left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - f_{t-1}(a, X_t)\}}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) f_{t-1}(a, X_t) \right\} - R(\pi^e) \right)^2 \right. \\ & \quad \left. - \left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - \mathbb{E}[Y_t(a) | X_t]\}}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) \mathbb{E}[Y(a) | X_t] \right\} - R(\pi^e) \right)^2 \right] \end{aligned}$$

Because $\alpha^2 - \beta^2 = (\alpha + \beta)(\alpha - \beta)$, there exists a constant $\gamma_0 > 0$ such that

$$\begin{aligned} & \leq \gamma_0 \mathbb{E} \left[\left[\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - f_{t-1}(a, X_t)\}}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) f_{t-1}(a, X_t) \right. \right. \right. \\ & \quad \left. \left. - \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - \mathbb{E}[Y_t(a) | X_t]\}}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} - \pi^e(a | X_t) \mathbb{E}[Y(a) | X_t] \right\} \right] \right] \end{aligned}$$

Then, there exist constants $\gamma_1 > 0$ such that

$$\leq \gamma_1 \mathbb{E} \left[\sum_{a=1}^K \left| f_{t-1}(a, X_t) - \mathbb{E}[Y(a) | X_t] \right| \right].$$

Here, from the assumption that $f_{t-1}(a, x) - \mathbb{E}[Y(a) | X] \xrightarrow{P} 0$ for $\tau = 2, 3, \dots, M$, and $f_{t_{\tau-1}}(a, x)$ is bounded for $\tau \in I$, we can use L^r convergence theorem. First, to use L^r convergence theorem, we use boundedness of f_{t_m} to derive the uniform integrability of f_{t_m} for $m = 0, 1, \dots, \tau-1$. Then, from L^r convergence theorem, we have $\mathbb{E}[|f_{t_m}(a, X) - \mathbb{E}[Y(a) | X]|] \rightarrow 0$ as $t_m \rightarrow \infty$. Using this results, we can show that, as $t_{\tau-1} \rightarrow \infty$ (this also means $T \rightarrow \infty$),

$$\gamma_1 \sum_{a=1}^K \mathbb{E} \left[\left| f_{t-1}(a, X_t) - \mathbb{E}[Y(a) | X_t] \right| \right] \rightarrow 0.$$

Therefore, as $t_{\tau-1} \rightarrow \infty$ ($T \rightarrow \infty$),

$$\begin{aligned} & \mathbb{E} \left[\left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - f_{t-1}(a, X_t)\}}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) f_{t-1}(a, X_t) \right\} - R(\pi^e) \right)^2 \right] \\ & \rightarrow \mathbb{E} \left[\left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - \mathbb{E}[Y_t(a) | X_t]\}}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) \mathbb{E}[Y(a) | X_t] \right\} - R(\pi^e) \right)^2 \right]. \end{aligned}$$

Then, by using $\mathbb{1}[A_t = a] \mathbb{1}[A_t = l] = 0$, $\mathbb{E} \left[\frac{\mathbb{1}[A_t = a] Y_t^2}{(\pi_\tau(a | X_t, \Omega_{t_{\tau-1}}))^2} \right] = \mathbb{E} \left[\frac{\mathbb{E}[Y_t^2(a) | X_t]}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} \right]$, and

$$\frac{1}{r_\tau T} \sum_{t=1}^T \mathbb{1}[t_{\tau-1} < t \leq t_\tau] = 1,$$

$$\begin{aligned} & \mathbb{E} \left[\left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - \mathbb{E}[Y_t(a) | X_t]\}}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) \mathbb{E}[Y(a) | X_t] \right\} - R(\pi^e) \right)^2 \right] \\ & = \mathbb{E} \left[\sum_{a=1}^K \left\{ \frac{(\pi^e(a | X_t))^2 \text{Var}(Y_t(a) | X_t)}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \left(\pi^e(a | X_t) \mathbb{E}[Y_t(a) | X_t] - R(\pi^e) \right)^2 \right\} \right]. \end{aligned}$$

In addition, the variance does not depend on t . We represent the independence by omitting the subscript t , i.e.,

$$\begin{aligned} & \mathbb{E} \left[\sum_{a=1}^K \left\{ \frac{(\pi^e(a | X_t))^2 \text{Var}(Y_t(a) | X_t)}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \left(\pi^e(a | X_t) \mathbb{E}[Y_t(a) | X_t] - R(\pi^e) \right)^2 \right\} \right] \\ & = \mathbb{E} \left[\sum_{a=1}^K \left\{ \frac{(\pi^e(a | X))^2 \text{Var}(Y(a) | X)}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \left(\pi^e(a | X) \mathbb{E}[Y_t(a) | X] - R(\pi^e) \right)^2 \right\} \right]. \end{aligned}$$

Therefore, we have

$$\begin{aligned} & \frac{1}{r_\tau^2 T} \sum_{t=1}^T \mathbb{E} \left[\left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - f_{t-1}(a, X_t)\}}{\pi_\tau(a | X_t, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) f_{t-1}(a, X_t) \right\} - R(\pi^e) \right)^2 \right. \\ & \quad \left. \times \mathbb{1}[t_{\tau-1} < t \leq t_\tau] \right] \\ & \rightarrow \frac{1}{r_\tau} \mathbb{E} \left[\sum_{a=1}^K \left\{ \frac{(\pi^e(a | X))^2 \text{Var}(Y(a) | X)}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \left(\pi^e(a | X) \mathbb{E}[Y(a) | X] - R(\pi^e) \right)^2 \right\} \right]. \end{aligned}$$

Thus, the matrix $(1/T) \sum_{t=1}^T \Sigma_t$ converges to a diagonal matrix Σ as $T \rightarrow \infty$, where the (τ, τ) -element of Σ is

$$\frac{1}{r_\tau} \mathbb{E} \left[\sum_{a=1}^K \left\{ \frac{(\pi^e(a | X))^2 \text{Var}(Y(a) | X)}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \left(\pi^e(a | X) \mathbb{E}[Y(a) | X] - R(\pi^e) \right)^2 \right\} \right].$$

Step 2: Condition (b)

Because we assume that all variables are bounded, this condition holds.

Step 3: Condition (c)

Here, we check that $(1/T) \sum_{t=1}^T (\mathbf{h}_t(X_t, A_t, Y_t; R(\pi^e), f_{t-1}, \pi_\tau, \pi^e)) (\mathbf{h}_t(X_t, A_t, Y_t; R(\pi^e), f_{t-1}, \pi_\tau, \pi^e))^\top \xrightarrow{p} \Sigma$. The (τ, τ) -element of the matrix is

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \frac{1}{r_\tau^2} \left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - f_{t-1}(a, X_t)\}}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) f_{t-1}(a, X_t) \right\} - \theta \right)^2 \\ & \quad \times \mathbb{1}[t_{\tau-1} < t \leq t_\tau] \\ &= \frac{1}{T} \sum_{t=1}^T \frac{1}{r_\tau^2} \left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - f_{t-1}(a, X_t)\}}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) f_{t-1}(a, X_t) \right\} - \theta \right)^2 \\ & \quad \times \mathbb{1}[t_{\tau-1} < t \leq t_\tau] \\ & - \frac{1}{T} \sum_{t=1}^T \frac{1}{r_\tau^2} \left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - \mathbb{E}[Y(a) | X_t]\}}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) \mathbb{E}[Y(a) | X_t] \right\} - \theta \right)^2 \\ & \quad \times \mathbb{1}[t_{\tau-1} < t \leq t_\tau] \\ & + \frac{1}{T} \sum_{t=1}^T \frac{1}{r_\tau^2} \left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - \mathbb{E}[Y(a) | X_t]\}}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) \mathbb{E}[Y(a) | X_t] \right\} - \theta \right)^2 \\ & \quad \times \mathbb{1}[t_{\tau-1} < t \leq t_\tau]. \end{aligned}$$

The part

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \frac{1}{r_\tau^2} \left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - f_{t-1}(a, X_t)\}}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) f_{t-1}(a, X_t) \right\} - \theta \right)^2 \\ & \quad \times \mathbb{1}[t_{\tau-1} < t \leq t_\tau] \\ & - \frac{1}{T} \sum_{t=1}^T \frac{1}{r_\tau^2} \left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - \mathbb{E}[Y(a) | X_t]\}}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) \mathbb{E}[Y(a) | X_t] \right\} - \theta \right)^2 \\ & \quad \times \mathbb{1}[t_{\tau-1} < t \leq t_\tau] \end{aligned}$$

converges in probability to 0 because $f_{t-1}(a, X_t) \xrightarrow{p} \mathbb{E}[Y(a) | X_t]$. The term

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \frac{1}{r_\tau^2} \left(\sum_{a=1}^K \left\{ \frac{\pi^e(a | X_t) \mathbb{1}[A_t = a] \{Y_t - \mathbb{E}[Y(a) | X_t]\}}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \pi^e(a | X_t) \mathbb{E}[Y(a) | X_t] \right\} - R(\pi^e) \right)^2 \\ & \quad \times \mathbb{1}[t_{\tau-1} < t \leq t_\tau]. \end{aligned}$$

converges in probability to

$$\frac{1}{r_\tau} \mathbb{E} \left[\sum_{a=1}^K \left\{ \frac{(\pi^e(a | X))^2 \text{Var}(Y(a) | X)}{\pi_\tau(a | X, \Omega_{t_{\tau-1}})} + \left(\pi^e(a | X) \mathbb{E}[Y(a) | X] - R(\pi^e) \right)^2 \right\} \right].$$

from the weak law of large numbers for i.i.d. samples as $t_{\tau-1} - t_\tau \rightarrow \infty$ because the samples are i.i.d. between $t_{\tau-1}$ and t_τ . \square

C.3 OPE estimator when the true logging policy is unknown

Then, we consider estimating the policy value without using the true logging policy π_t . We use adaptive-fitting for obtaining an asymptotically normal estimator. As the ADR estimator under Assumption 2, we estimate π_t and f^* only using Ω_{t-1} and denote them as g_{t-1} and \hat{f}_{t-1} , respectively.

Then, we define an estimator of $R(\pi^e)$ as

$$\tilde{R}_T^{\text{batch}}(\pi^e) = \arg \min_{R \in \mathbb{R}} (\tilde{q}_T(R))^\top \hat{W}_T (\tilde{q}_T(R)),$$

where $\tilde{q}_T(R) = \frac{1}{T} \sum_{t=1}^T \mathbf{h}_t(X_t, A_t, Y_t; R, \hat{f}_{t-1}, \hat{g}_{t-1}\pi^e)$ and \hat{W}_T is a data-dependent $(M \times M)$ -dimensional positive semi-definite matrix. This estimator is the ADR estimator under batch update. As well as the proof of Theorem 1, it is hold that

$$\left| \tilde{R}_T^{\text{batch}}(\pi^e) - \hat{R}_T^{\text{batch}}(\pi^e) \right| = o_p(1/\sqrt{T})$$

if $\|\hat{g}_{t-1}(a|X_t) - \pi_{t-1}(a|X_t, \Omega_{t-1})\|_2 = o_p(t^{-p})$, and $\|\hat{f}_{t-1}(a, X_t) - f^*(a, X_t)\|_2 = o_p(t^{-q})$, where $p, q > 0$ such that $p + q = 1/2$. Therefore, we can obtain the following theorem.

Theorem 5 (Asymptotic distribution the ADR estimator under batch update). *Suppose that*

(i) $\hat{W}_T \xrightarrow{p} W$;

(ii) W is a positive definite.

Then, under Assumptions 1, 3, 4–6,

$$\sqrt{T}(\tilde{R}_T^{\text{batch}}(\pi^e) - R(\pi^e)) \xrightarrow{d} \mathcal{N}(0, \sigma^2),$$

where $\sigma^2 = (I^\top W I)^{-1} I^\top W \Sigma W^\top I (I^\top W I)^{-1}$ and Σ is a $(M \times M)$ diagonal matrix such that the $(\tau \times \tau)$ -element is $\Psi(\pi^e, \pi_\tau)$.

D Details of experiments

The description of the dataset is shown in Table 3. We use LinUCB and LinTS policies. We add uniform sampling to make overlap between policies. We can relax this requirement by considering different DGPs, such as batched sampling. However, for brevity, we adopt this setting. Additional results are shown as follows.

For numerical experiments in Section 6.2, we show the result with sample sizes $T = 100, 1,000$ in Table 4. In addition, we show the error distribution with sample size $T = 100$ in Figure 4; we show the error distribution with sample size $T = 500$ in Figure 5; we show the error distribution with sample size $T = 500$ in Figure 6; we show the error distribution with sample size $T = 1,000$ in Figure 7.

For experiments with dependent samples in Section 7, we show the additional results with different settings in Tables 5–10. In Table 5, we show the results using the benchmark datasets with 800 samples generated from the LinUCB algorithm. In Table 6, we show the results using the benchmark datasets with 1,000 samples generated from the LinUCB algorithm. In Table 7, we show the results using the the benchmark datasets with 1,200 samples generated from the LinUCB algorithm. In Table 8, we show the results using the the benchmark datasets with 800 samples generated from the LinTS algorithm. In Table 9, we show the results using the the benchmark datasets with 1,000 samples generated from the LinTS algorithm. In Table 10, we show the results using the the benchmark datasets with 1,200 samples generated from the LinTS algorithm.

Next, we compare the estimators using the benchmark datasets generated from the logistic regression as well as the evaluation weight; that is, the samples are i.i.d. For $\alpha \in \{0.7, 0.4, 0.1\}$ and the sample sizes 800, 1,000, and 1,200, we calculate the RMSEs and the SDs over 10 trials. We show additional results with different settings in Tables 11–13. In Table 11, we show the results using the benchmark datasets with 800 samples. In Table 12, we show the results using the benchmark datasets with 1,000 samples. In Table 13, we show the results using the benchmark datasets with 1,200 samples. In these experiments, AIPW estimators show better performances than the ADR estimator. We conjecture that this is because the logging policy is not unstable unlike the case with dependent samples; that is, the paradox is specific to the case where samples are dependent. Note this case (i.i.d. samples) is not in the scope of the proposed method. In such cases, it is common to use other methods, such as the cross-fitting proposed by Chernozhukov et al. (2018), instead of

the AIPW and ADR estimators discussed in this paper. We show this result to clarify the cause of the paradox. Since the ADR and AIPW estimators have the same asymptotic variance, and since the AIPW estimator uses more information, it is a natural result that the AIPW estimator performs better. However, when the samples are dependent, the ADR estimator paradoxically outperforms the AIPW estimator because of the unstable behavior of the logging policy π_t .

As a surprising discovery, our proposed ADR estimator shows better results than the AIPW estimator, although the AIPW estimator uses more information (true logging policy π_t) than the ADR estimator and their asymptotic properties are the same. As discussed above, we consider that this result is due to the logging policy’s unstable behavior. Even when knowing the true logging policy π_t , we can stabilize estimation by reestimating the logging policy from (A_t, X_t) . This paradox is similar to the well-known property that the IPW estimator using an estimated propensity score shows a smaller asymptotic variance than the IPW using the true propensity score (Hirano et al., 2003; Henmi & Eguchi, 2004; Henmi et al., 2007). However, we consider that they are different phenomena. In previous studies, the paradox is mainly explained by differences in the asymptotic variance between the IPW estimators with the true and an estimated propensity score. On the other hand, for our case, the AIPW and ADR estimators have the same asymptotic variance, unlike IPW-type estimators; therefore, we cannot elucidate the paradox by traditional explanation.

Table 3: Specification of datasets

Dataset	the number of samples	Dimension	the number of classes
mnist	60,000	780	10
satimage	4,435	35	6
sensorless	58,509	48	11
connect-4	67,557	126	3

Table 4: The results of Section 6.1 with sample sizes $T = 100, 1,000$. We show the RMSE, SD, and coverage ratio of the confidence interval (CR). We highlight in red bold two estimators with the lowest RMSE. Estimators with asymptotic normality are marked with †, and estimators that do not require the true logging policy are marked with *.

LinUCB policy																		
T	ADR †*			IPW †			AIPW †			AW-AIPW †			DM *			EIPW *		
	RMSE	SD	CR	RMSE	SD	CR	RMSE	SD	CR	RMSE	SD	CR	RMSE	SD	CR	RMSE	SD	CR
100	0.100	0.010	0.90	0.169	0.056	0.81	0.156	0.046	0.81	0.182	0.081	0.60	0.117	0.015	0.22	0.127	0.022	0.79
1,000	0.026	0.001	0.97	0.070	0.019	0.92	0.062	0.018	0.97	0.067	0.024	0.77	0.034	0.001	0.20	0.136	0.015	0.12

LinTS policy																		
T	ADR †*			IPW †			AIPW †			AW-AIPW †			DM *			EIPW *		
	RMSE	SD	CR	RMSE	SD	CR	RMSE	SD	CR	RMSE	SD	CR	RMSE	SD	CR	RMSE	SD	CR
100	0.092	0.008	0.94	0.183	0.050	0.85	0.153	0.038	0.89	0.165	0.042	0.70	0.104	0.013	0.26	0.141	0.020	0.70
1,000	0.023	0.001	0.97	0.075	0.011	0.92	0.053	0.006	0.93	0.053	0.006	0.72	0.033	0.001	0.17	0.111	0.008	0.17

Table 5: The results of benchmark datasets with the LinUCB policy and $T = 800$. We highlight in red bold the estimator with the lowest RMSE and highlight in under line the estimator with the lowest RMSE among estimators that do not use the true logging policy. Estimators with asymptotic normality are marked with †, and estimators that do not require the true logging policy are marked with *.

mnist		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.047	0.003	0.099	0.011	0.132	0.032	0.279	0.025	0.130	0.019
0.4		0.034	0.002	0.067	0.005	0.115	0.014	0.281	0.015	0.095	0.009
0.1		0.053	0.006	0.079	0.006	0.107	0.013	0.287	0.028	0.119	0.018

satimage		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.015	0.000	0.097	0.013	0.074	0.007	0.035	0.001	0.065	0.003
0.4		0.020	0.000	0.078	0.006	0.088	0.010	0.041	0.001	0.085	0.014
0.1		0.026	0.001	0.080	0.006	0.092	0.011	0.059	0.003	0.097	0.012

sensorless		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.050	0.002	0.253	0.049	0.101	0.016	0.130	0.008	0.121	0.016
0.4		0.065	0.006	0.270	0.039	0.104	0.019	0.150	0.013	0.096	0.012
0.1		0.042	0.001	0.250	0.043	0.075	0.006	0.127	0.008	0.102	0.010

connect-4		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.024	0.001	0.121	0.019	0.042	0.002	0.045	0.001	0.140	0.007
0.4		0.023	0.001	0.125	0.014	0.033	0.001	0.044	0.002	0.090	0.007
0.1		0.028	0.001	0.107	0.012	0.042	0.002	0.068	0.004	0.031	0.001

Table 6: The results of benchmark datasets with the LinUCB policy and $T = 1,000$. We highlight in red bold the estimator with the lowest RMSE and highlight in under line the estimator with the lowest RMSE among estimators that do not use the true logging policy. Estimators with asymptotic normality are marked with †, and estimators that do not require the true logging policy are marked with *.

mnist		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.046	0.002	0.100	0.011	0.162	0.027	0.232	0.013	0.148	0.014
0.4		0.028	0.001	0.068	0.005	0.112	0.009	0.249	0.010	0.080	0.004
0.1		<u>0.086</u>	0.006	0.078	0.006	0.085	0.008	0.299	0.024	0.091	0.008

satimage		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.013	0.000	0.098	0.011	0.060	0.004	0.037	0.001	0.056	0.002
0.4		0.022	0.000	0.078	0.008	0.019	0.000	0.043	0.001	0.060	0.002
0.1		0.029	0.001	0.078	0.005	0.061	0.008	0.041	0.002	0.041	0.002

sensorless		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.043	0.001	0.270	0.049	0.096	0.014	0.109	0.014	0.068	0.007
0.4		0.053	0.004	0.287	0.035	0.070	0.006	0.112	0.008	0.104	0.013
0.1		0.048	0.002	0.260	0.031	0.072	0.006	0.154	0.015	0.088	0.012

connect-4		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.030	0.001	0.135	0.024	0.042	0.002	0.032	0.001	0.173	0.006
0.4		0.015	0.000	0.135	0.017	0.037	0.001	0.038	0.001	0.085	0.002
0.1		0.012	0.000	0.117	0.013	0.030	0.001	0.051	0.002	0.023	0.001

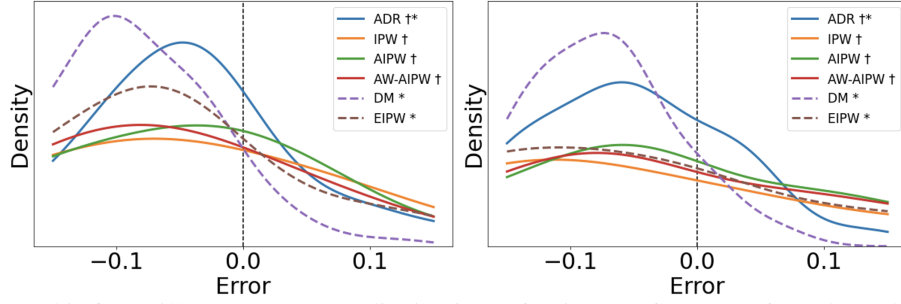


Figure 4: This figure illustrates the error distributions of estimators for OPVE from dependent samples generated with the LinUcB(left) and LinTS(right) with the sample size 100. We smoothed the error distributions using kernel density estimation. Estimators with asymptotic normality are marked with †, and estimators that do not require a true logging policy are marked with *.

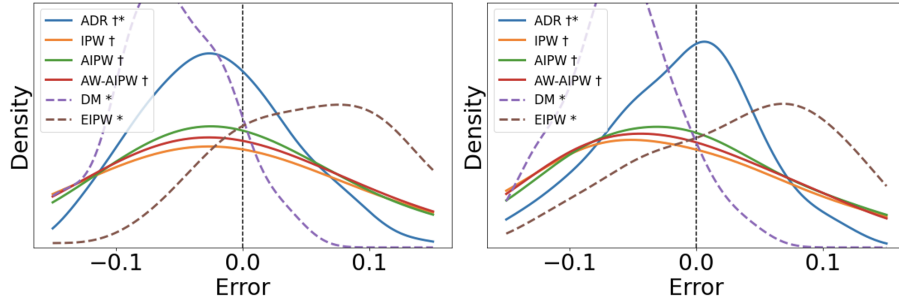


Figure 5: This figure illustrates the error distributions of estimators for OPVE from dependent samples generated with the LinUcB(left) and LinTS(right) with the sample size 250. We smoothed the error distributions using kernel density estimation. Estimators with asymptotic normality are marked with †, and estimators that do not require a true logging policy are marked with *.

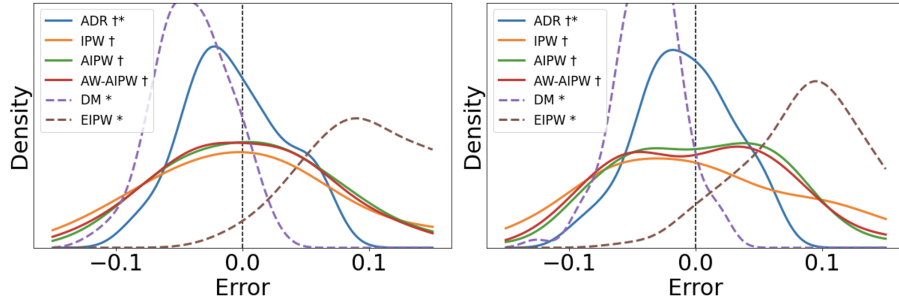


Figure 6: This figure illustrates the error distributions of estimators for OPVE from dependent samples generated with the LinUcB(left) and LinTS(right) with the sample size 500. We smoothed the error distributions using kernel density estimation. Estimators with asymptotic normality are marked with †, and estimators that do not require a true logging policy are marked with *.

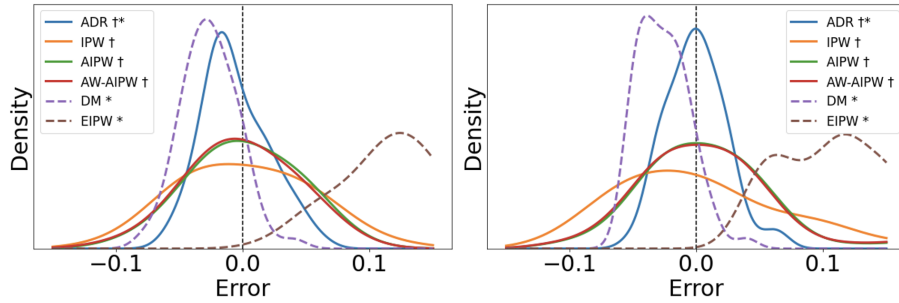


Figure 7: This figure illustrates the error distributions of estimators for OPVE from dependent samples generated with the LinUcB(left) and LinTS(right) with the sample size 1,000. We smoothed the error distributions using kernel density estimation. Estimators with asymptotic normality are marked with †, and estimators that do not require a true logging policy are marked with *.

Table 7: The results of benchmark datasets with the LinUCB policy and $T = 1, 200$. We highlight in red bold the estimator with the lowest RMSE and highlight in under line the estimator with the lowest RMSE among estimators that do not use the true logging policy. Estimators with asymptotic normality are marked with †, and estimators that do not require the true logging policy are marked with *.

mnist	ADR †*		IPW †		AIPW †		DM *		EIPW *	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.063	0.003	0.095	0.010	0.114	0.018	0.204	0.007	0.218	0.021
0.4	0.031	0.002	0.061	0.004	0.100	0.011	0.227	0.010	0.136	0.015
0.1	0.073	0.005	0.076	0.006	0.093	0.013	0.245	0.013	0.059	0.004
satimage	ADR †*		IPW †		AIPW †		DM *		EIPW *	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.015	0.000	0.091	0.010	0.039	0.002	0.035	0.001	0.039	0.001
0.4	0.016	0.000	0.075	0.006	0.041	0.002	0.043	0.001	0.048	0.002
0.1	0.019	0.000	0.075	0.005	0.033	0.001	0.045	0.001	0.073	0.005
sensorless	ADR †*		IPW †		AIPW †		DM *		EIPW *	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.049	0.003	0.272	0.059	0.102	0.018	0.092	0.008	0.078	0.010
0.4	0.047	0.004	0.286	0.037	0.071	0.007	0.096	0.006	0.082	0.009
0.1	0.056	0.004	0.268	0.047	0.061	0.006	0.116	0.008	0.067	0.005
connect-4	ADR †*		IPW †		AIPW †		DM *		EIPW *	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.033	0.001	0.138	0.027	0.046	0.002	0.037	0.002	0.184	0.013
0.4	0.020	0.000	0.130	0.014	0.052	0.003	0.038	0.001	0.114	0.004
0.1	0.020	0.000	0.118	0.013	0.037	0.002	0.034	0.001	0.048	0.003

Table 8: The results of benchmark datasets with the LinTS policy and $T = 800$. We highlight in red bold the estimator with the lowest RMSE and highlight in under line the estimator with the lowest RMSE among estimators that do not use the true logging policy. Estimators with asymptotic normality are marked with †, and estimators that do not require the true logging policy are marked with *.

mnist	ADR †*		IPW †		AIPW †		DM *		EIPW *	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.064	0.006	0.098	0.011	0.150	0.019	0.289	0.026	0.172	0.019
0.4	0.084	0.005	0.060	0.003	0.089	0.007	0.311	0.027	0.116	0.016
0.1	0.117	0.011	0.078	0.006	0.061	0.004	0.331	0.030	0.065	0.004
satimage	ADR †*		IPW †		AIPW †		DM *		EIPW *	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.032	0.001	0.102	0.011	0.094	0.007	0.043	0.002	0.095	0.016
0.4	0.033	0.001	0.081	0.009	0.073	0.005	0.048	0.002	0.098	0.015
0.1	0.028	0.001	0.075	0.004	0.074	0.005	0.043	0.001	0.048	0.003
sensorless	ADR †*		IPW †		AIPW †		DM *		EIPW *	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.055	0.003	0.252	0.053	0.094	0.011	0.158	0.015	0.112	0.014
0.4	0.084	0.006	0.267	0.033	0.078	0.007	0.177	0.024	0.059	0.004
0.1	0.086	0.009	0.247	0.041	0.104	0.017	0.164	0.009	0.079	0.006
connect-4	ADR †*		IPW †		AIPW †		DM *		EIPW *	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.021	0.000	0.129	0.019	0.057	0.003	0.041	0.001	0.115	0.004
0.4	0.025	0.001	0.130	0.018	0.055	0.003	0.050	0.003	0.064	0.004
0.1	0.018	0.000	0.107	0.011	0.054	0.004	0.054	0.002	0.029	0.001

Table 9: The results of benchmark datasets with the LinTS policy and $T = 1,000$. We highlight in red bold the estimator with the lowest RMSE and highlight in under line the estimator with the lowest RMSE among estimators that do not use the true logging policy. Estimators with asymptotic normality are marked with †, and estimators that do not require the true logging policy are marked with *.

mnist		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.053	0.006	0.097	0.010	0.140	0.023	0.248	0.008	0.184	0.020
0.4		<u>0.077</u>	0.008	0.070	0.005	0.098	0.011	0.276	0.028	0.099	0.011
0.1		<u>0.099</u>	0.011	0.076	0.006	0.091	0.008	0.296	0.026	<u>0.088</u>	0.007
satimage		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.016	0.000	0.096	0.011	0.050	0.004	0.044	0.001	0.056	0.003
0.4		<u>0.021</u>	0.001	0.077	0.007	0.041	0.002	0.048	0.001	0.063	0.009
0.1		0.026	0.001	0.079	0.005	0.087	0.011	0.048	0.002	0.091	0.020
sensorless		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.040	0.003	0.273	0.057	0.083	0.009	0.114	0.008	0.087	0.008
0.4		0.034	0.001	0.291	0.035	0.091	0.010	0.099	0.006	0.070	0.007
0.1		0.060	0.002	0.268	0.042	0.075	0.008	0.146	0.011	0.069	0.006
connect-4		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.014	0.000	0.132	0.020	0.059	0.003	0.045	0.001	0.110	0.006
0.4		0.019	0.000	0.136	0.017	0.028	0.001	0.046	0.002	0.066	0.004
0.1		0.016	0.000	0.114	0.012	0.023	0.001	0.050	0.002	0.019	0.000

Table 10: The results of benchmark datasets with the LinTS policy and $T = 1,200$. We highlight in red bold the estimator with the lowest RMSE and highlight in under line the estimator with the lowest RMSE among estimators that do not use the true logging policy. Estimators with asymptotic normality are marked with †, and estimators that do not require the true logging policy are marked with *.

mnist		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.038	0.003	0.095	0.010	0.095	0.010	0.234	0.013	0.203	0.028
0.4		0.048	0.003	0.066	0.004	0.099	0.012	0.245	0.013	0.117	0.011
0.1		0.058	0.002	0.073	0.006	0.052	0.003	0.257	0.021	0.069	0.006
satimage		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.012	0.000	0.096	0.010	0.058	0.004	0.037	0.001	0.059	0.003
0.4		0.020	0.001	0.079	0.007	0.036	0.001	0.043	0.002	0.064	0.010
0.1		0.027	0.001	0.079	0.005	0.045	0.002	0.050	0.001	0.050	0.004
sensorless		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.027	0.001	0.257	0.043	0.063	0.004	0.102	0.005	0.053	0.004
0.4		0.038	0.002	0.277	0.036	0.059	0.003	0.106	0.006	0.087	0.010
0.1		0.046	0.002	0.252	0.035	0.034	0.002	0.138	0.008	<u>0.045</u>	0.003
connect-4		ADR †*		IPW †		AIPW †		DM *		EIPW *	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.014	0.000	0.376	0.244	0.044	0.002	0.029	0.001	0.107	0.007
0.4		0.021	0.000	0.395	0.283	0.025	0.001	0.033	0.001	0.067	0.004
0.1		0.215	0.138	0.385	0.273	0.215	0.138	0.217	0.137	0.215	0.138

Table 11: The results of benchmark datasets with the i.i.d samples and $T = 800$. We highlight in red bold the estimator with the lowest RMSE and highlight in under line the estimator with the lowest RMSE among estimators that do not use the true logging policy. Estimators with asymptotic normality are marked with †, and estimators that do not require the true logging policy are marked with *.

mnist		ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.195	0.090	0.078	0.061	0.078	0.061	0.078	0.061	2.085	9.013
0.4		0.192	0.088	0.078	0.061	0.079	0.061	<u>0.080</u>	0.061	1.187	2.913
0.1		<u>0.105</u>	0.061	0.080	0.061	0.083	0.061	0.108	0.061	0.298	0.189
satimage		ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.025	0.001	0.003	0.00	0.008	0.000	<u>0.006</u>	0.00	0.076	0.014
0.4		0.019	0.001	0.006	0.00	0.017	0.001	<u>0.012</u>	0.00	0.074	0.012
0.1		0.079	0.060	0.079	0.06	0.083	0.060	0.081	0.06	0.098	0.061
sensorless		ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.032	0.003	0.004	0.000	0.005	0.000	0.004	0.000	0.193	0.099
0.4		0.031	0.003	0.007	0.000	0.008	0.000	<u>0.015</u>	0.001	0.154	0.064
0.1		<u>0.052</u>	0.021	0.046	0.021	0.048	0.021	0.057	0.021	0.121	0.050
connect-4		ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.080	0.047	0.069	0.047	0.069	0.047	0.069	0.047	0.234	0.086
0.4		0.074	0.047	0.069	0.047	0.071	0.047	<u>0.071</u>	0.047	0.141	0.052
0.1		0.070	0.047	0.070	0.047	0.073	0.047	0.074	0.047	0.076	0.047

Table 12: The results of benchmark datasets with the i.i.d samples and $T = 1,000$. We highlight in red bold the estimator with the lowest RMSE and highlight in under line the estimator with the lowest RMSE among estimators that do not use the true logging policy. Estimators with asymptotic normality are marked with †, and estimators that do not require the true logging policy are marked with *.

mnist		ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.367	0.062	0.006	0.000	0.011	0.000	0.005	0.000	4.834	2.294
0.4		0.373	0.037	0.016	0.000	0.031	0.001	<u>0.038</u>	0.001	2.672	0.905
0.1		0.165	0.012	0.029	0.001	0.049	0.002	<u>0.155</u>	0.005	0.671	0.098
satimage		ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.043	0.002	0.006	0.000	0.013	0.000	0.010	0.000	0.136	0.010
0.4		0.038	0.001	0.013	0.000	0.035	0.001	<u>0.018</u>	0.000	0.148	0.012
0.1		<u>0.028</u>	0.001	0.026	0.001	0.034	0.001	0.044	0.002	0.113	0.019
sensorless		ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.073	0.005	0.005	0.000	0.012	0.000	<u>0.009</u>	0.000	0.454	0.071
0.4		0.077	0.005	0.009	0.000	0.031	0.001	<u>0.023</u>	0.001	0.390	0.045
0.1		<u>0.066</u>	0.003	0.042	0.003	0.047	0.004	0.083	0.007	0.257	0.066
connect-4		ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α		RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7		0.057	0.002	0.009	0.000	0.014	0.000	<u>0.013</u>	0.000	0.396	0.020
0.4		0.040	0.001	0.008	0.000	0.021	0.001	<u>0.019</u>	0.000	0.218	0.010
0.1		<u>0.027</u>	0.001	0.024	0.001	0.047	0.003	0.042	0.002	0.056	0.003

Table 13: The results of benchmark datasets with the i.i.d samples and $T = 1, 200$. We highlight in red bold the estimator with the lowest RMSE and highlight in under line the estimator with the lowest RMSE among estimators that do not use the true logging policy. Estimators with asymptotic normality are marked with †, and estimators that do not require the true logging policy are marked with *.

mnist	ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.310	0.036	0.006	0.000	0.013	0.000	0.005	0.000	4.709	1.037
0.4	0.346	0.033	0.012	0.000	0.020	0.001	<u>0.040</u>	0.001	2.735	0.461
0.1	0.143	0.013	0.036	0.001	0.076	0.007	0.161	0.006	0.608	0.137

satimage	ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.038	0.001	0.005	0.0	0.015	0.000	<u>0.012</u>	0.000	0.114	0.004
0.4	0.029	0.001	0.011	0.0	0.019	0.000	<u>0.016</u>	0.000	0.159	0.011
0.1	<u>0.021</u>	0.001	0.014	0.0	0.048	0.002	0.044	0.001	0.092	0.007

sensorless	ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.096	0.007	0.007	0.000	0.012	0.000	0.008	0.000	0.481	0.083
0.4	0.097	0.007	0.014	0.000	0.027	0.001	<u>0.025</u>	0.000	0.407	0.058
0.1	<u>0.054</u>	0.002	0.023	0.001	0.038	0.002	0.065	0.002	0.208	0.025

connect-4	ADR †*		IPW †		AIPW †		DM †*		EIPW †*	
α	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD	RMSE	SD
0.7	0.060	0.002	0.005	0.0	0.012	0.000	<u>0.009</u>	0.000	0.378	0.013
0.4	0.043	0.002	0.010	0.0	0.014	0.000	<u>0.023</u>	0.000	0.211	0.012
0.1	0.017	0.000	0.017	0.0	0.037	0.002	0.035	0.001	0.057	0.003