

Supplementary to ‘Doubly Robust Interval Estimation for Optimal Policy Evaluation in Online Learning’

Ye Shen ^{*}, Hengrui Cai [†], Rui Song [‡]

This supplementary article provides sensitivity analyses and all the technical proofs for the established theorems for policy evaluation in online learning under the contextual bandits. Note that the theoretical results in Section 7 can be proven in a similar manner by arguments in Section B. Thus we omit the details here.

A Sensitivity Test for the Choice of p_t

We conduct a sensitivity test for the choice of p_t in this section. We run all the simulations with $p_t = 0.01, 0.05, 0.1$ and find that Algorithm 1 is not sensitive to the choice of p_t .

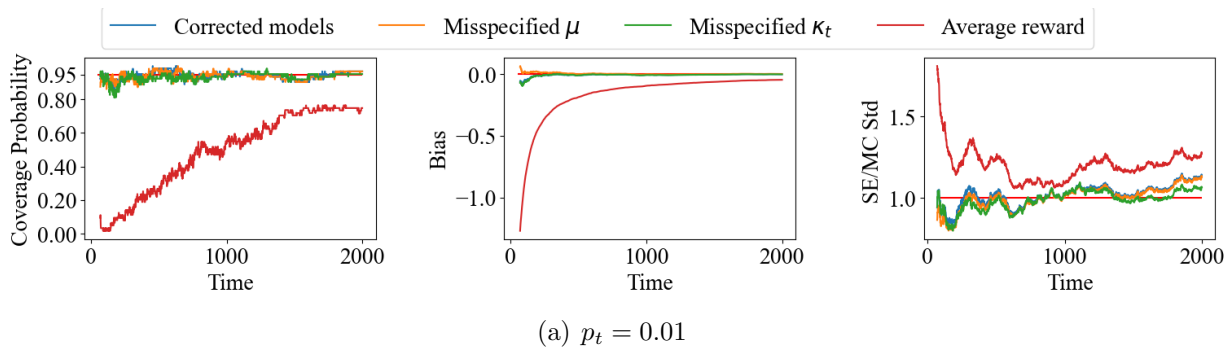


Figure A.1: Results by DREAM under UCB with different model specifications in comparison to the averaged reward. Left panel: the coverage probabilities of the 95% two-sided Wald-type CI, with the red line representing the nominal level at 95%. Middle panel: the bias between the estimated value and the true value. Right panel: the ratio between the standard error and the Monte Carlo standard deviation, with the red line representing the nominal level at 1.

^{*}Equal Contribution, Department of Statistics, North Carolina State University

[†]Equal Contribution, Department of Statistics, University of California Irvine

[‡]Amazon.com, Inc

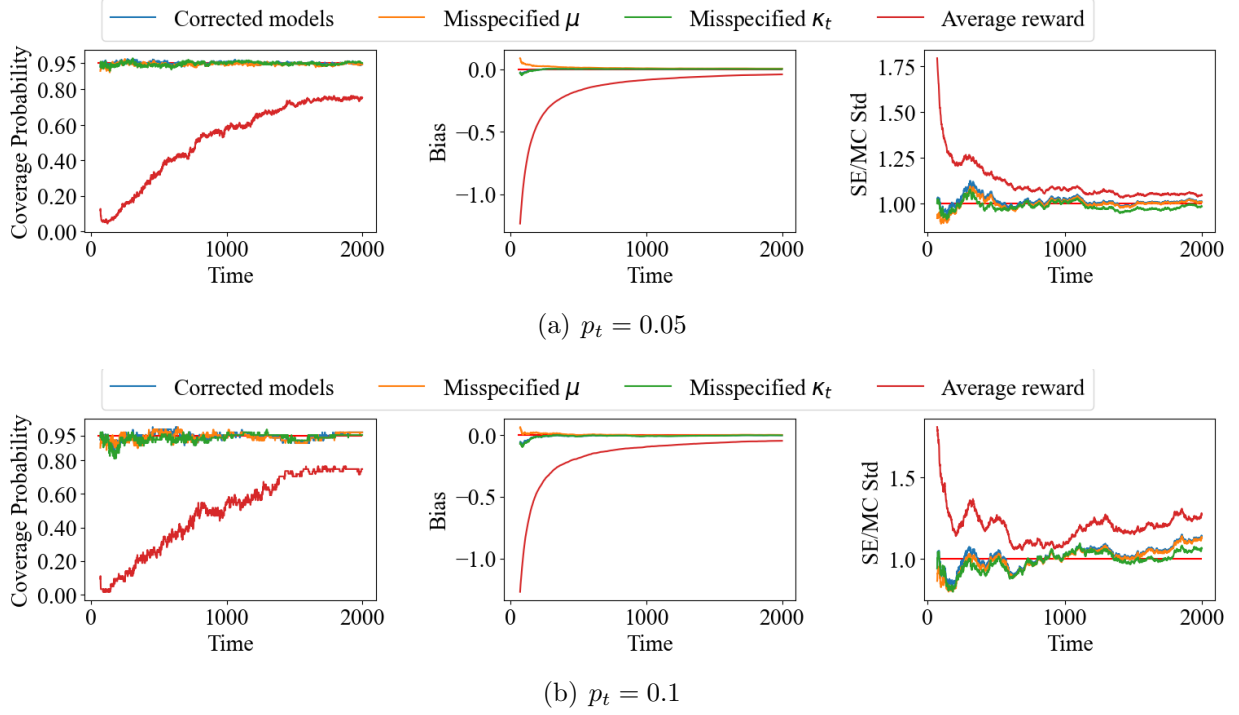
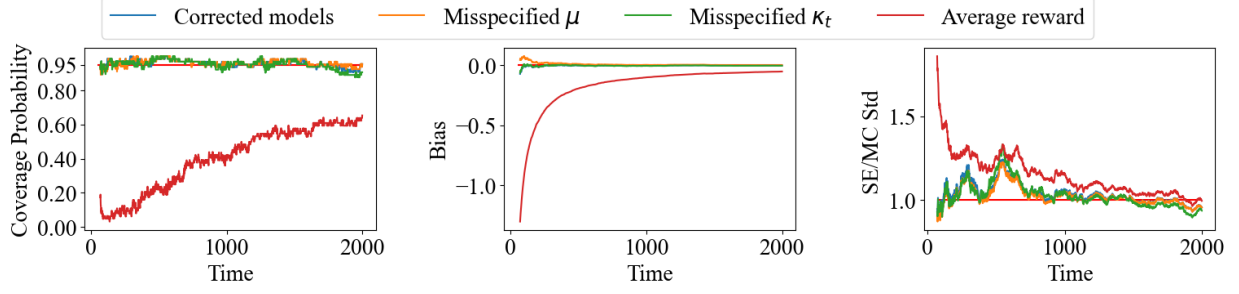
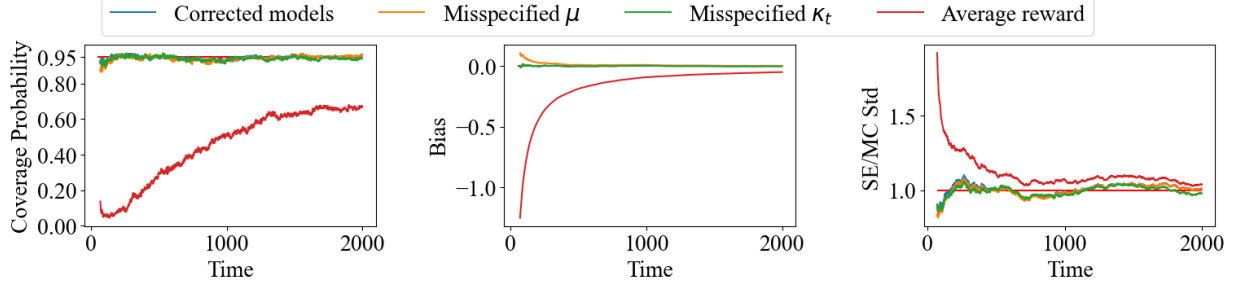


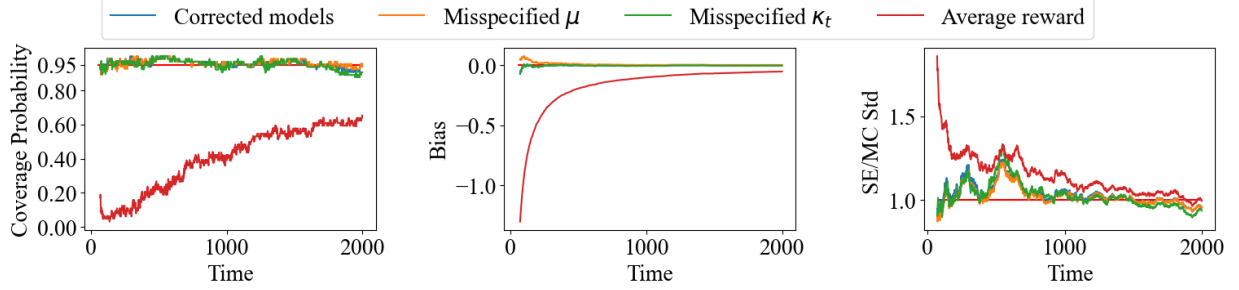
Figure A.1: Results by DREAM under UCB with different model specifications in comparison to the averaged reward. Left panel: the coverage probabilities of the 95% two-sided Wald-type CI, with the red line representing the nominal level at 95%. Middle panel: the bias between the estimated value and the true value. Right panel: the ratio between the standard error and the Monte Carlo standard deviation, with the red line representing the nominal level at 1.



(c) $p_t = 0.01$

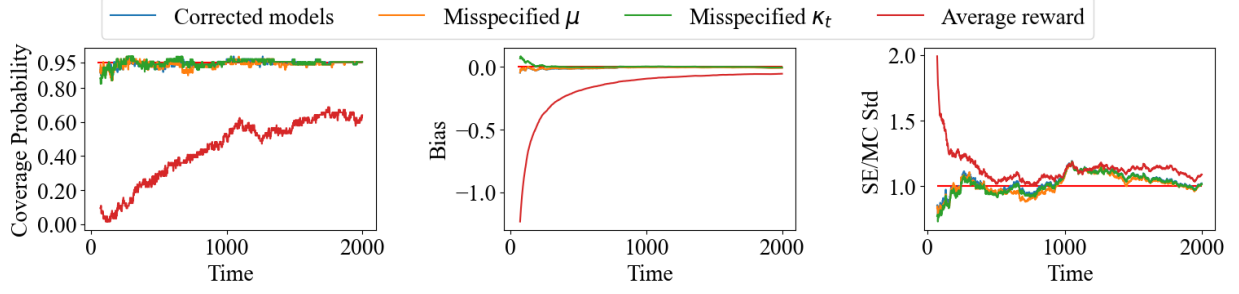


(d) $p_t = 0.05$

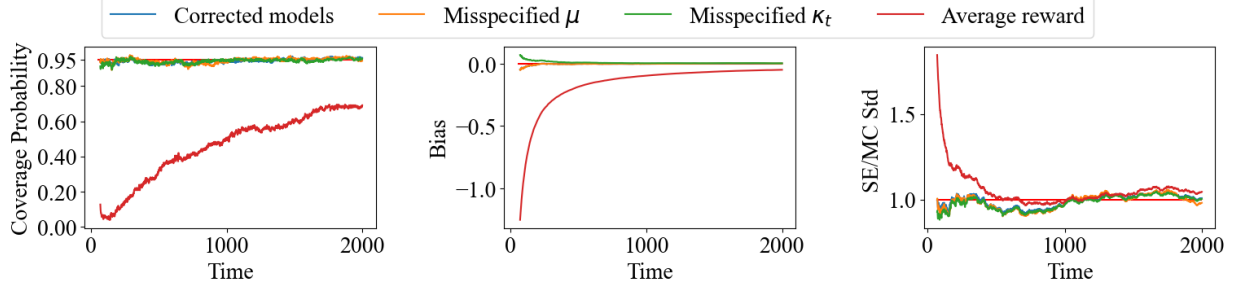


(e) $p_t = 0.1$

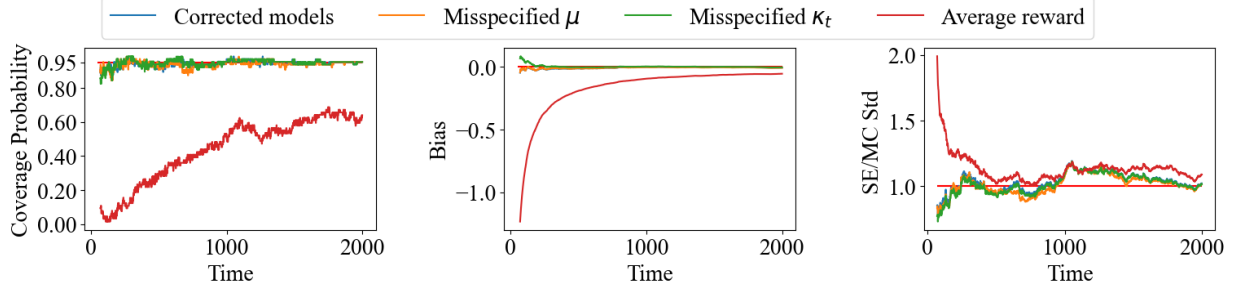
Figure A.2: Results by DREAM under TS with different model specifications in comparison to the averaged reward. Left panel: the coverage probabilities of the 95% two-sided Wald-type CI, with the red line representing the nominal level at 95%. Middle panel: the bias between the estimated value and the true value. Right panel: the ratio between the standard error and the Monte Carlo standard deviation, with the red line representing the nominal level at 1.



(a) $p_t = 0.01$



(b) $p_t = 0.05$



(c) $p_t = 0.1$

Figure A.3: Results by DREAM under EG with different model specifications in comparison to the averaged reward. Left panel: the coverage probabilities of the 95% two-sided Wald-type CI, with the red line representing the nominal level at 95%. Middle panel: the bias between the estimated value and the true value. Right panel: the ratio between the standard error and the Monte Carlo standard deviation, with the red line representing the nominal level at 1.

B Technical Proofs for Main Results

This section provides all the technical proofs for the established theorems for policy evaluation in online learning under the contextual bandits.

B.1 Proof of Lemma 4.1

The proof of Lemma 4.1 consists of three main steps. To be specific, we first reconstruct the target difference $\widehat{\beta}_t(a) - \beta(a)$ and decompose it into two parts. Then, we establish the bound for each part and derive its lower bound $\Pr(\|\widehat{\beta}_t(a) - \beta(a)\|_1 \leq h)$.

Step 1: Recall Equation (1) in the main paper with $\mathbf{D}_{t-1}(a)$ being a $\mathbf{N}_{t-1}(a) \times d$ design matrix at time $t - 1$ with $\mathbf{N}_{t-1}(a)$ as the number of pulls for action a , we have

$$\widehat{\beta}_t(a) = \left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}^{-1} \left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i r_i \right\}.$$

We are interested in the quantity

$$\widehat{\beta}_t(a) - \beta(a) = \left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}^{-1} \left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i r_i \right\} - \beta(a). \quad (\text{B.1})$$

Note that $\beta(a)$ can be written as

$$\left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}^{-1} \left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\} \beta(a),$$

and since $r_i = \mathbf{x}_i^\top \beta(a) + e_i$, we can write (B.1) as

$$\begin{aligned} \widehat{\beta}_t(a) - \beta(a) = & \underbrace{\left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}^{-1} \left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i e_i \right\}}_{\eta_3} \quad (\text{B.2}) \\ & - \underbrace{\left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}^{-1} \frac{\omega}{t} \beta(a)}_{\eta_4}. \end{aligned}$$

Our goal is to find a lower bound of $\Pr(\|\widehat{\boldsymbol{\beta}}_t(a) - \boldsymbol{\beta}(a)\|_1 \leq h)$ for any $h > 0$. Notice that by the triangle inequality we have $\|\widehat{\boldsymbol{\beta}}_t(a) - \boldsymbol{\beta}(a)\|_1 \leq \|\eta_3\|_1 + \|\eta_4\|_1$, thus we can find the lower bound using the inequality as

$$\Pr\left(\left\|\widehat{\boldsymbol{\beta}}_t(a) - \boldsymbol{\beta}(a)\right\|_1 \leq h\right) \geq \Pr\left(\|\eta_3\|_1 + \|\eta_4\|_1 \leq h\right). \quad (\text{B.3})$$

Step 2: We focus on bounding η_4 first. By the relationship between eigenvalues and the L_2 norm of symmetric matrix, we have $\|\mathbf{M}^{-1}\|_2 = \lambda_{\max}(\mathbf{M}^{-1}) = \{\lambda_{\min}(\mathbf{M})\}^{-1}$ for any invertible matrix \mathbf{M} . Thus we can obtain that

$$\begin{aligned} \left\| \left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}^{-1} \right\|_2 &= \left\{ \lambda_{\min} \left(\frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right) \right\}^{-1} \\ &= \frac{1}{\lambda_{\min} \left(\frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top \right) + \frac{1}{t} \omega} \\ &\leq \frac{1}{p_t \lambda_{\min}(\boldsymbol{\Sigma}) + \frac{1}{t} \omega} \leq \frac{1}{p_t \lambda + \frac{1}{t} \omega}, \end{aligned}$$

which leads to

$$\|\eta_4\|_2 \leq \left\| \left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}^{-1} \right\|_2 \left\| \frac{\omega}{t} \boldsymbol{\beta}(a) \right\|_2 \leq \frac{\omega}{t p_t \lambda + \omega} \|\boldsymbol{\beta}(a)\|_2. \quad (\text{B.4})$$

By Cauchy-Schwartz inequality, we further have bound of the L_1 norm of η_4 as

$$\|\eta_4\|_1 \leq \sqrt{d} \|\eta_4\|_2 \leq \frac{\omega \sqrt{d}}{t p_t \lambda + \omega} \|\boldsymbol{\beta}(a)\|_2 \leq \frac{\sqrt{d}}{1 + t p_t \lambda / \omega} \|\boldsymbol{\beta}(a)\|_2 \leq \sqrt{d} \|\boldsymbol{\beta}(a)\|_2. \quad (\text{B.5})$$

Step 3: Lastly, using the results in (B.5), we have

$$\Pr(\|\eta_3\|_1 + \|\eta_4\|_1 \leq h) \geq \Pr\left(\|\eta_3\|_1 \leq h - \sqrt{d} \|\boldsymbol{\beta}(a)\|_2\right). \quad (\text{B.6})$$

By the definition of η_3 and Lemma 2 in Chen et al. (2020), for any constant $c > 0$, we have

$$\Pr(\|\eta_3\|_1 > c) \leq 2d \exp \left\{ -\frac{t \left(\frac{p_t \lambda}{2} \right)^2 c^2}{2d^2 \sigma^2 L_x^2} \right\} = 2d \exp \left\{ -\frac{t p_t^2 \lambda^2 c^2}{8d^2 \sigma^2 L_x^2} \right\}.$$

Therefore, from (B.3) and (B.6), taking $c = h - \sqrt{d} \|\beta(a)\|_2$, we have that under event E_t ,

$$\Pr \left(\left\| \hat{\beta}_t(a) - \beta(a) \right\|_1 > h \right) \leq 2d \exp \left\{ -\frac{tp_t^2 \lambda^2 \left(h - \sqrt{d} \|\beta(a)\|_2 \right)^2}{8d^2 \sigma^2 L_x^2} \right\}. \quad (\text{B.7})$$

Based on the above results, it is immediate that the online ridge estimator $\hat{\beta}_t(a)$ is consistent to $\beta(a)$ if $tp_t^2 \rightarrow \infty$ as $t \rightarrow \infty$. The proof is hence completed.

B.2 Proof of Corollary 1

Since $\hat{\mu}_t(\mathbf{x}_t, a) - \mu(\mathbf{x}_t, a) = \mathbf{x}_t^\top (\hat{\beta}_t(a) - \beta(a))$, by Holder's inequality, we have

$$|\hat{\mu}_t(\mathbf{x}_t, a) - \mu(\mathbf{x}_t, a)| \leq \|\mathbf{x}_t\|_\infty \left\| \hat{\beta}_t(a) - \beta(a) \right\|_1 \leq L_x \left\| \hat{\beta}_t(a) - \beta(a) \right\|_1,$$

which follows

$$\Pr \{ |\hat{\mu}_t(\mathbf{x}_t, a) - \mu(\mathbf{x}_t, a)| > \xi \} \leq \Pr \left\{ L_x \left\| \hat{\beta}_t(a) - \beta(a) \right\|_1 > \xi \right\} = \Pr \left\{ \left\| \hat{\beta}_t(a) - \beta(a) \right\|_1 > \xi / L_x \right\}.$$

By Lemma 4.1, we further have

$$\begin{aligned} \Pr \{ |\hat{\mu}_t(\mathbf{x}_t, a) - \mu(\mathbf{x}_t, a)| > \xi \} &\leq \Pr \left\{ \left\| \hat{\beta}_t(a) - \beta(a) \right\|_1 > \frac{\xi}{L_x} \right\} \\ &\leq 2d \exp \left\{ -\frac{tp_t^2 \lambda^2 \left(\frac{\xi}{L_x} - \sqrt{d} \|\beta(a)\|_2 \right)^2}{8d^2 \sigma^2 L_x^2} \right\} \\ &= 2d \exp \left\{ -\frac{tp_t^2 \lambda^2 \left(\xi - \sqrt{d} L_x \|\beta(a)\|_2 \right)^2}{8d^2 \sigma^2 L_x^4} \right\}. \end{aligned}$$

Note that by the Triangle Inequality,

$$\begin{aligned} |\hat{\mu}_t(\mathbf{x}_t, 1) - \hat{\mu}_t(\mathbf{x}_t, 0) - \Delta_{\mathbf{x}_t}| &= |\{\hat{\mu}_t(\mathbf{x}_t, 1) - \mu(\mathbf{x}_t, 1)\} - \{\hat{\mu}_t(\mathbf{x}_t, 0) - \mu(\mathbf{x}_t, 0)\}| \\ &\leq |\hat{\mu}_t(\mathbf{x}_t, 1) - \mu(\mathbf{x}_t, 1)| + |\hat{\mu}_t(\mathbf{x}_t, 0) - \mu(\mathbf{x}_t, 0)|, \end{aligned}$$

thus for $|\widehat{\mu}_t(\mathbf{x}_t, 1) - \widehat{\mu}_t(\mathbf{x}_t, 0) - \Delta_{\mathbf{x}_t}|$, we have

$$\begin{aligned}
& \Pr \{ |\widehat{\mu}_t(\mathbf{x}_t, 1) - \widehat{\mu}_t(\mathbf{x}_t, 0) - \Delta_{\mathbf{x}_t}| > \xi \} \\
& \leq \Pr \{ |\widehat{\mu}_t(\mathbf{x}_t, 1) - \mu(\mathbf{x}_t, 1)| + |\widehat{\mu}_t(\mathbf{x}_t, 0) - \mu(\mathbf{x}_t, 0)| > \xi \} \\
& \leq \Pr \{ |\widehat{\mu}_t(\mathbf{x}_t, 1) - \mu(\mathbf{x}_t, 1)| > \xi/2 \} + \Pr \{ |\widehat{\mu}_t(\mathbf{x}_t, 0) - \mu(\mathbf{x}_t, 0)| > \xi/2 \} \\
& \leq 2d \exp \left\{ -\frac{tp_t^2 \lambda^2 \left(\xi/2 - \sqrt{d} L_{\mathbf{x}} \|\boldsymbol{\beta}(1)\|_2 \right)^2}{8d^2 \sigma^2 L_{\mathbf{x}}^4} \right\} + 2d \exp \left\{ -\frac{tp_t^2 \lambda^2 \left(\xi/2 - \sqrt{d} L_{\mathbf{x}} \|\boldsymbol{\beta}(0)\|_2 \right)^2}{8d^2 \sigma^2 L_{\mathbf{x}}^4} \right\} \\
& \leq 4d \exp \left\{ -(t-1)p_{t-1}^2 c_{\xi} \right\},
\end{aligned}$$

with

$$c_{\xi} = \frac{\lambda^2 \left[\min \left\{ \left(\xi/2 - \sqrt{d} L_{\mathbf{x}} \|\boldsymbol{\beta}(1)\|_2 \right)^2, \left(\xi/2 - \sqrt{d} L_{\mathbf{x}} \|\boldsymbol{\beta}(0)\|_2 \right)^2 \right\} \right]}{8d^2 \sigma^2 L_{\mathbf{x}}^4},$$

consistent with time t .

B.3 Proof of Theorem 1

The proof of Theorem 1 consists of two main parts to show the probability of exploration under UCB and TS, respectively, by noting the probability of exploration under EG is given by its definition.

B.3.1 Proof for UCB

We first show the probability of exploration under UCB. This proof consists of three main steps stated as following:

1. We first rewrite the target probability by its definition and express it as

$$\Pr \{ |\widehat{\mu}_{t-1}(\mathbf{x}_t, 1) - \widehat{\mu}_{t-1}(\mathbf{x}_t, 0)| < c_t |\widehat{\sigma}_{t-1}(\mathbf{x}_t, 0) - \widehat{\sigma}_{t-1}(\mathbf{x}_t, 1)| \}.$$

2. Then, we establish the bound for the variance estimation such that

$$c_t |\widehat{\sigma}_{t-1}(\mathbf{x}_t, 0) - \widehat{\sigma}_{t-1}(\mathbf{x}_t, 1)| \leq \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}}.$$

3. Lastly, we bound $\Pr \left\{ |\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)| < \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} \right\}$ using the result in Corollary 1.

Step 1: We rewrite the target probability by definition and decompose it into two parts.

Let $\Delta_{\mathbf{x}_t} \equiv \mu(\mathbf{x}_t, 1) - \mu(\mathbf{x}_t, 0)$. Based on the definition of the probability of exploration and the form of the estimated optimal policy $\hat{\pi}_t(\mathbf{x}_t)$, we have

$$\begin{aligned} \kappa_t(\mathbf{x}_t) &= \Pr\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)\} = \mathbb{E}[\mathbb{I}\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)\}] \\ &= \underbrace{\mathbb{E}[\mathbb{I}(a_t = 0)|\hat{\pi}_t(\mathbf{x}_t) = 1]\Pr\{\hat{\pi}_t(\mathbf{x}_t) = 1\}}_{\eta_0} + \underbrace{\mathbb{E}[\mathbb{I}(a_t = 1)|\hat{\pi}_t(\mathbf{x}_t) = 0]\Pr\{\hat{\pi}_t(\mathbf{x}_t) = 0\}}_{\eta_1}, \end{aligned} \quad (\text{B.8})$$

where the expectation is taken with respect to the history \mathcal{H}_{t-1} before time point t .

Next, we rewrite η_0 and η_1 using the estimated mean and variance components $\hat{\mu}_{t-1}(\mathbf{x}_t, a)$ and $\hat{\sigma}_{t-1}(\mathbf{x}_t, a)$, where $a = 0, 1$. We focus on η_0 first.

Given $\hat{\pi}_t(\mathbf{x}_t) = 1$, i.e., $\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0) > 0$, based on the definition of the taken action in Lin-UCB that

$$a_t = \mathbb{I}\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) + c_t \hat{\sigma}_{t-1}(\mathbf{x}_t, 1) > \hat{\mu}_{t-1}(\mathbf{x}_t, 0) + c_t \hat{\sigma}_{t-1}(\mathbf{x}_t, 0)\},$$

the probability of choosing action 0 rather than action 1 is

$$\begin{aligned} &\mathbb{E}[\mathbb{I}(a_t = 0)|\hat{\pi}_t(\mathbf{x}_t) = 1] \\ &= \Pr\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) + c_t \hat{\sigma}_{t-1}(\mathbf{x}_t, 1) < \hat{\mu}_{t-1}(\mathbf{x}_t, 0) + c_t \hat{\sigma}_{t-1}(\mathbf{x}_t, 0)|\hat{\pi}_t(\mathbf{x}_t) = 1\} \\ &= \Pr\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0) < c_t \hat{\sigma}_{t-1}(\mathbf{x}_t, 0) - c_t \hat{\sigma}_{t-1}(\mathbf{x}_t, 1)|\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0) > 0\} \\ &= \Pr[0 < \hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0) < c_t \{\hat{\sigma}_{t-1}(\mathbf{x}_t, 0) - \hat{\sigma}_{t-1}(\mathbf{x}_t, 1)\}] / \Pr\{\hat{\pi}_t(\mathbf{x}_t) = 1\}, \end{aligned}$$

where the second equality is to rearrange the estimated mean and variance components, and the last equality comes from the definition of the conditional probability. Combining this with (B.8), we have

$$\eta_0 = \Pr[0 < \hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0) < c_t \{\hat{\sigma}_{t-1}(\mathbf{x}_t, 0) - \hat{\sigma}_{t-1}(\mathbf{x}_t, 1)\}].$$

Similarly, we have

$$\eta_1 = \Pr[c_t \{\hat{\sigma}_{t-1}(\mathbf{x}_t, 0) - \hat{\sigma}_{t-1}(\mathbf{x}_t, 1)\} < \hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0) < 0].$$

Thus combined with Equation (B.8), we have

$$\begin{aligned}
\kappa_t(\mathbf{x}_t) = \eta_0 + \eta_1 = & \Pr[0 < \hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0) < c_t \{\hat{\sigma}_{t-1}(\mathbf{x}_t, 0) - \hat{\sigma}_{t-1}(\mathbf{x}_t, 1)\}] \\
& + \Pr[c_t \{\hat{\sigma}_{t-1}(\mathbf{x}_t, 0) - \hat{\sigma}_{t-1}(\mathbf{x}_t, 1)\} < \hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0) < 0] \\
= & \Pr\{|\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)| < c_t |\hat{\sigma}_{t-1}(\mathbf{x}_t, 0) - \hat{\sigma}_{t-1}(\mathbf{x}_t, 1)|\}.
\end{aligned} \tag{B.9}$$

The rest of the proof is aims to bound the probability

$$\Pr\{|\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)| < c_t |\hat{\sigma}_{t-1}(\mathbf{x}_t, 0) - \hat{\sigma}_{t-1}(\mathbf{x}_t, 1)|\}.$$

Step 2: Secondly, we bound the variance $\hat{\sigma}_{t-1}(\mathbf{x}_t, 0)$ and $\hat{\sigma}_{t-1}(\mathbf{x}_t, 1)$.

We consider the quantity $\hat{\sigma}_{t-1}(\mathbf{x}_t, 0) = \sqrt{\{\mathbf{x}_t^\top \{\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) + \omega \mathbf{I}_d\}^{-1} \mathbf{x}_t\}}$ first. Let \mathbf{v} be any $d \times 1$ vector, then the sample variance under action 0 is given by

$$\begin{aligned}
\mathbf{x}_t^\top \{\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) + \omega \mathbf{I}_d\}^{-1} \mathbf{x}_t &= \|\mathbf{x}_t\|_2^2 \left(\frac{\mathbf{x}_t}{\|\mathbf{x}_t\|_2} \right)^\top \{\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) + \omega \mathbf{I}_d\}^{-1} \left(\frac{\mathbf{x}_t}{\|\mathbf{x}_t\|_2} \right) \\
&\leq \|\mathbf{x}_t\|_2^2 \max_{\|\mathbf{v}\|_2=1} \mathbf{v}^\top \{\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) + \omega \mathbf{I}_d\}^{-1} \mathbf{v} \\
&\leq \|\mathbf{x}_t\|_2^2 \lambda_{\max}\{(\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) + \omega \mathbf{I}_d)^{-1}\},
\end{aligned} \tag{B.10}$$

where the first inequality is to replace $(\mathbf{x}_t/\|\mathbf{x}_t\|_2)^\top$ with any normalized vector, and the second inequality is due to the definition. According to (B.10), combined with Assumption 4.1, we can further bound $\hat{\sigma}_{t-1}(\mathbf{x}_t, 0)$ by

$$\|\mathbf{x}_t\|_2 \sqrt{\lambda_{\max}\{(\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) + \omega \mathbf{I}_d)^{-1}\}} \leq \frac{L_x}{\sqrt{\lambda_{\min}\{\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) + \omega \mathbf{I}_d\}}}. \tag{B.11}$$

It is immediate from (B.10) and (B.11) that

$$0 < \hat{\sigma}_{t-1}(\mathbf{x}_t, 0) \leq \frac{L_x}{\sqrt{\lambda_{\min}\{\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) + \omega \mathbf{I}_d\}}}. \tag{B.12}$$

Note that

$$\lambda_{\min}\{\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) + \omega \mathbf{I}_d\} = \lambda_{\min}\{\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0)\} + \omega,$$

combined with the fact that $\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) = \sum_{i=1}^{t-1} (1-a_i) \mathbf{x}_i \mathbf{x}_i^\top$, then $\lambda_{\min}\{\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) + \omega \mathbf{I}_d\}$ can be further expressed as

$$\begin{aligned}
\lambda_{\min}\{\mathbf{D}_{t-1}(0)^\top \mathbf{D}_{t-1}(0) + \omega \mathbf{I}_d\} &= (t-1) \lambda_{\min} \left\{ \frac{1}{t-1} \sum_{i=1}^{t-1} (1-a_i) \mathbf{x}_i \mathbf{x}_i^\top \right\} + \omega \\
&> (t-1) p_{t-1} \lambda_{\min}(\mathbf{\Sigma}) + \omega > (t-1) p_{t-1} \lambda + \omega,
\end{aligned}$$

where the first inequality is owing to Assumption 4.2 , and the second inequality is owing to Assumption 4.1. This together with (B.12) gives the lower and upper bounds of $\widehat{\sigma}_{t-1}(\mathbf{x}_t, 0)$ as

$$0 < \widehat{\sigma}_{t-1}(\mathbf{x}_t, 0) \leq \frac{L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda} + \omega} < \frac{L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}}. \quad (\text{B.13})$$

Similarly we have

$$0 < \widehat{\sigma}_{t-1}(\mathbf{x}_t, 1) \leq \frac{L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}}, \quad (\text{B.14})$$

which follows that

$$c_t |\widehat{\sigma}_{t-1}(\mathbf{x}_t, 0) - \widehat{\sigma}_{t-1}(\mathbf{x}_t, 1)| \leq c_t (|\widehat{\sigma}_{t-1}(\mathbf{x}_t, 0)| + |\widehat{\sigma}_{t-1}(\mathbf{x}_t, 1)|) \leq \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}}.$$

Combining (B.9) and the above equation, we get the conclusion that

$$\begin{aligned} \kappa_t(\mathbf{x}_t) &\leq \Pr \{ |\widehat{\mu}_{t-1}(\mathbf{x}_t, 1) - \widehat{\mu}_{t-1}(\mathbf{x}_t, 0)| < c_t |\widehat{\sigma}_{t-1}(\mathbf{x}_t, 0) - \widehat{\sigma}_{t-1}(\mathbf{x}_t, 1)| \} \\ &\leq \Pr \left\{ |\widehat{\mu}_{t-1}(\mathbf{x}_t, 1) - \widehat{\mu}_{t-1}(\mathbf{x}_t, 0)| < \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} \right\}. \end{aligned} \quad (\text{B.15})$$

Step 3: Lastly, we aim to bound $\Pr \left\{ |\widehat{\mu}_{t-1}(\mathbf{x}_t, 1) - \widehat{\mu}_{t-1}(\mathbf{x}_t, 0)| < \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} \right\}$ using the result in Corollary 1.

For any $\xi > 0$, define $E := \{ |\widehat{\mu}_t(\mathbf{x}_t, 1) - \widehat{\mu}_t(\mathbf{x}_t, 0) - \Delta_{\mathbf{x}_t}| \leq \xi \}$, which satisfies $\Pr \{E\} \geq 1 - 4d \exp \{-tp_t^2 c_{\xi}\}$ by Corollary 1. Then on the Event E , we have

$$\begin{aligned} |\widehat{\mu}_{t-1}(\mathbf{x}_t, 1) - \widehat{\mu}_{t-1}(\mathbf{x}_t, 0)| &= |\Delta_{\mathbf{x}_t} + \{\widehat{\mu}_{t-1}(\mathbf{x}_t, 1) - \widehat{\mu}_{t-1}(\mathbf{x}_t, 0) - \Delta_{\mathbf{x}_t}\}| \\ &\geq |\Delta_{\mathbf{x}_t}| - |\widehat{\mu}_t(\mathbf{x}_t, 1) - \widehat{\mu}_t(\mathbf{x}_t, 0) - \Delta_{\mathbf{x}_t}| \geq |\Delta_{\mathbf{x}_t}| - \xi. \end{aligned}$$

Thus for the probability $\Pr \left\{ |\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)| < \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} \right\}$, we have

$$\begin{aligned}
\kappa_t(\mathbf{x}_t) &\leq \Pr \left\{ |\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)| < \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} \right\} \\
&\leq \Pr \left\{ |\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)| < \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} \mid E \right\} + \Pr \{E^c\} \\
&\leq \Pr \left\{ |\Delta_{\mathbf{x}_t}| - \xi < \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} \right\} + 4d \exp \{-(t-1)p_{t-1}^2 c_{\xi}\} \\
&= \Pr \left\{ |\Delta_{\mathbf{x}_t}| < \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} + \xi \right\} + 4d \exp \{-(t-1)p_{t-1}^2 c_{\xi}\}.
\end{aligned} \tag{B.16}$$

Sine $tp_t \rightarrow \infty$ as $t \rightarrow \infty$, for any constant $\delta > \xi$, there exist large enough t satisfying $\frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} \leq \delta - \xi$. Then by Assumption 4.3, there exists some constant γ such that

$$\Pr \left\{ |\Delta_{\mathbf{x}_t}| < \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} + \xi \right\} = \mathcal{O} \left\{ \left(\frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} + \xi \right)^{\gamma} \right\},$$

i.e., there exists some constant C such that

$$\Pr \left\{ |\Delta_{\mathbf{x}_t}| < \frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} + \xi \right\} = C \left(\frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} + \xi \right)^{\gamma}.$$

Therefore, combined with the Equation (B.16), we have

$$\kappa_t(\mathbf{x}_t) \leq C \left(\frac{2c_t L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}} + \xi \right)^{\gamma} + 4d \exp \{-(t-1)p_{t-1}^2 c_{\xi}\}.$$

The proof is hence completed.

B.3.2 Proof for TS

We next show the probability of exploration under TS consisting of three main steps:

1. We firstly define an event $E := \{|\hat{\mu}_t(\mathbf{x}_t, 1) - \hat{\mu}_t(\mathbf{x}_t, 0) - \Delta_{\mathbf{x}_t}| \leq \xi\}$ for any $0 < \xi < |\Delta_{\mathbf{x}_t}|/2$, where the estimated difference between mean functions is close to the true difference. And we have $\Pr\{E\} \geq 1 - 4d \exp\{-tp_t^2 c_\xi\}$ by Corollary 1.
2. Next, we bound the probability of exploration on the event E .
3. Lastly, we combine the results in the previous two steps to get the unconditioned probability of exploration.

Step 1: For any $0 < \xi < |\Delta_{\mathbf{x}_t}|/2$, define $E := \{|\hat{\mu}_t(\mathbf{x}_t, 1) - \hat{\mu}_t(\mathbf{x}_t, 0) - \Delta_{\mathbf{x}_t}| \leq \xi\}$, which satisfies $\Pr\{E\} \geq 1 - 4d \exp\{-(t-1)p_{t-1}^2 c_\xi\}$ by Corollary 1. Then for the probability $\Pr\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)\}$, we have

$$\Pr\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)\} \leq \Pr\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)|E\} + \Pr\{E^c\} \leq \Pr\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)|E\} + 4d \exp\{-(t-1)p_{t-1}^2 c_\xi\}. \quad (\text{B.17})$$

Without loss of generality, we assume $\Delta_{\mathbf{x}_t} > 0$, then $E := \{0 < \Delta_{\mathbf{x}_t} - \xi \leq \hat{\mu}_t(\mathbf{x}_t, 1) - \hat{\mu}_t(\mathbf{x}_t, 0) \leq \Delta_{\mathbf{x}_t} + \xi\}$, which implies $\hat{\pi}_t(\mathbf{x}_t) = 1$.

Using the law of iterated expectations, based on the definition of the probability of exploration and the form of the estimated optimal policy $\hat{\pi}_t(\mathbf{x}_t)$, on the event E , we have

$$\Pr\{a_t \neq 1|E\} = \mathbb{E}[\mathbb{I}\{a_t \neq 1\}|E] = \mathbb{E}(\mathbb{E}[\mathbb{I}\{a_t = 0\}|\hat{\mu}_{t-1}(\mathbf{x}_t, 1), \hat{\mu}_{t-1}(\mathbf{x}_t, 0)]|E). \quad (\text{B.18})$$

Step 2: Next, we focus on deriving the bound of $\mathbb{E}[\mathbb{I}\{a_t = 0\}|\hat{\mu}_{t-1}(\mathbf{x}_t, 1), \hat{\mu}_{t-1}(\mathbf{x}_t, 0)]$ on the Event E .

Recalling the bandit mechanism of TS, we have $a_t = \mathbb{I}\{\mathbf{x}_t^\top \boldsymbol{\beta}_t(1) > \mathbf{x}_t^\top \boldsymbol{\beta}_t(0)\}$, where $\boldsymbol{\beta}_t(a)$ is drawn from the posterior distribution of $\boldsymbol{\beta}(a)$ given by

$$\mathcal{N}_d[\hat{\boldsymbol{\beta}}_{t-1}(a), \rho^2\{\mathbf{D}_{t-1}(a)^\top \mathbf{D}_{t-1}(a) + \omega \mathbf{I}_d\}^{-1}].$$

From the posterior distributions and the definitions of $\hat{\mu}_{t-1}(\mathbf{x}_t, a)$ and $\hat{\sigma}_{t-1}(\mathbf{x}_t, a)$, we have

$$\mathbf{x}_t^\top \boldsymbol{\beta}_t(a) \sim \mathcal{N}[\mathbf{x}_t^\top \hat{\boldsymbol{\beta}}_{t-1}(a), \rho^2 \mathbf{x}_t^\top \{\mathbf{D}_{t-1}(a)^\top \mathbf{D}_{t-1}(a) + \omega \mathbf{I}_d\}^{-1} \mathbf{x}_t],$$

that is,

$$\mathbf{x}_t^\top \boldsymbol{\beta}_t(a) \sim \mathcal{N}[\hat{\mu}_{t-1}(\mathbf{x}_t, a), \rho^2 \hat{\sigma}_{t-1}(\mathbf{x}_t, a)^2].$$

Notice that $\mathbf{x}_t^\top \boldsymbol{\beta}_t(1)$ and $\mathbf{x}_t^\top \boldsymbol{\beta}_t(0)$ are drawn independently, thus,

$$\mathbf{x}_t^\top \boldsymbol{\beta}_t(1) - \mathbf{x}_t^\top \boldsymbol{\beta}_t(0) \sim \mathcal{N}[\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0), \rho^2\{\hat{\sigma}_{t-1}(\mathbf{x}_t, 1)^2 + \hat{\sigma}_{t-1}(\mathbf{x}_t, 0)^2\}]. \quad (\text{B.19})$$

Recall $a_t = \mathbb{I}\{\mathbf{x}_t^\top \boldsymbol{\beta}_t(1) > \mathbf{x}_t^\top \boldsymbol{\beta}_t(0)\}$ in TS, based on the posterior distribution in (B.19). Therefore, on the Event E we have

$$\begin{aligned}
& \mathbb{E}[\mathbb{I}\{a_t = 0\} | \hat{\mu}_{t-1}(\mathbf{x}_t, 1), \hat{\mu}_{t-1}(\mathbf{x}_t, 0)] \\
&= \Pr\{\mathbf{x}_t^\top \boldsymbol{\beta}_t(1) - \mathbf{x}_t^\top \boldsymbol{\beta}_t(0) < 0 | \hat{\mu}_{t-1}(\mathbf{x}_t, 1), \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\} \\
&= \Phi[-\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\} / \sqrt{\rho^2\{\hat{\sigma}_{t-1}(\mathbf{x}_t, 1)^2 + \hat{\sigma}_{t-1}(\mathbf{x}_t, 0)^2\}}] \\
&= 1 - \Phi[\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\} / \sqrt{\rho^2\{\hat{\sigma}_{t-1}(\mathbf{x}_t, 1)^2 + \hat{\sigma}_{t-1}(\mathbf{x}_t, 0)^2\}}],
\end{aligned} \tag{B.20}$$

where $\Phi(\cdot)$ is the cumulative distribution function of the standard normal distribution. Denote $\hat{z}_t \equiv \{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\} / \sqrt{\rho^2\{\hat{\sigma}_{t-1}(\mathbf{x}_t, 1)^2 + \hat{\sigma}_{t-1}(\mathbf{x}_t, 0)^2\}} > 0$, since $\hat{\pi}_t(\mathbf{x}_t) = \mathbb{I}\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) > \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\} = 1$, i.e., $\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0) > 0$. By applying the tail bound established for the normal distribution in Section 7.1 of Feller (2008), we have (B.20) can be bounded as

$$\mathbb{E}[\mathbb{I}\{a_t = 0\} | \hat{\mu}_{t-1}(\mathbf{x}_t, 1), \hat{\mu}_{t-1}(\mathbf{x}_t, 0)] \leq \exp(-\hat{z}_t^2/2).$$

This yields that on the Event E ,

$$\begin{aligned}
& \mathbb{E}[\mathbb{I}\{a_t = 0\} | \hat{\mu}_{t-1}(\mathbf{x}_t, 1), \hat{\mu}_{t-1}(\mathbf{x}_t, 0)] \\
& \leq \mathbb{E}\left\{\exp(-\hat{z}_t^2/2)\right\} = \mathbb{E}\left\{\exp\left(-\frac{\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\}^2}{2\rho^2\{\hat{\sigma}_{t-1}(\mathbf{x}_t, 1)^2 + \hat{\sigma}_{t-1}(\mathbf{x}_t, 0)^2\}}\right)\right\}.
\end{aligned}$$

Using similar arguments in proving (B.13) that $\hat{\sigma}_{t-1}(\mathbf{x}_t, 0) \leq \frac{L_{\mathbf{x}}}{\sqrt{(t-1)p_{t-1}\lambda}}$, we have

$$\hat{\sigma}_{t-1}(\mathbf{x}_t, 1)^2 + \hat{\sigma}_{t-1}(\mathbf{x}_t, 0)^2 \leq \frac{2L_{\mathbf{x}}^2}{(t-1)p_{t-1}\lambda}.$$

Therefore, combining the above two equations leads to

$$\mathbb{E}[\mathbb{I}\{a_t = 0\} | \hat{\mu}_{t-1}(\mathbf{x}_t, 1), \hat{\mu}_{t-1}(\mathbf{x}_t, 0)] \leq \mathbb{E}\left\{\exp\left(-\frac{\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\}^2(t-1)p_{t-1}\lambda}{4\rho^2 L_{\mathbf{x}}^2}\right)\right\},$$

where the expectation is taken with respect to history \mathcal{H}_{t-1} .

Note that on the Event E , we have

$$\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\}^2 \geq (|\Delta_{\mathbf{x}_t}| - \xi)^2,$$

which follows that on the Event E ,

$$\begin{aligned} \mathbb{E}[\mathbb{I}\{a_t = 0\}|\hat{\mu}_{t-1}(\mathbf{x}_t, 1), \hat{\mu}_{t-1}(\mathbf{x}_t, 0)] &\leq \mathbb{E}\left\{\exp\left(-\frac{(|\Delta_{\mathbf{x}_t}| - \xi)^2 (t-1)p_{t-1}\lambda}{4\rho^2 L_{\mathbf{x}}^2}\right)\right\} \\ &\leq \exp\left(-\frac{(|\Delta_{\mathbf{x}_t}| - \xi)^2 (t-1)p_{t-1}\lambda}{4\rho^2 L_{\mathbf{x}}^2}\right). \end{aligned} \quad (\text{B.21})$$

Step 3: Combined with Equation (B.17), Equation (B.18) and Equation (B.21), we have

$$\begin{aligned} \Pr\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)\} &\stackrel{(\text{B.17})}{\leq} \Pr\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)|E\} + 4d \exp\{-(t-1)p_{t-1}^2 c_\xi\} \\ &\stackrel{(\text{B.18})}{\leq} \mathbb{E}(\mathbb{E}[\mathbb{I}\{a_t = 0\}|\hat{\mu}_{t-1}(\mathbf{x}_t, 1), \hat{\mu}_{t-1}(\mathbf{x}_t, 0)]|E) + 4d \exp\{-(t-1)p_{t-1}^2 c_\xi\} \\ &\stackrel{(\text{B.21})}{\leq} \exp\left(-\frac{(|\Delta_{\mathbf{x}_t}| - \xi)^2 (t-1)p_{t-1}\lambda}{4\rho^2 L_{\mathbf{x}}^2}\right) + 4d \exp\{-(t-1)p_{t-1}^2 c_\xi\}. \end{aligned}$$

The proof is hence completed.

B.4 Proof of Theorem 2

We detail the proof of Theorem 2 in this section. Using the similar arguments in (B.2) in the proof of Lemma 4.1, we can rewrite $\sqrt{t}\{\hat{\boldsymbol{\beta}}_t(a) - \boldsymbol{\beta}(a)\}$ as

$$\begin{aligned} \sqrt{t}\{\hat{\boldsymbol{\beta}}_t(a) - \boldsymbol{\beta}(a)\} &= \underbrace{\left\{\frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d\right\}^{-1}}_{\boldsymbol{\xi}} \underbrace{\left\{\frac{1}{\sqrt{t}} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i e_i\right\}}_{\boldsymbol{\eta}_1} \\ &\quad - \underbrace{\left\{\frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d\right\}^{-1}}_{\boldsymbol{\eta}_2} \frac{1}{\sqrt{t}} \boldsymbol{\beta}(a). \end{aligned}$$

Our goal is to prove that $\sqrt{t}\{\hat{\boldsymbol{\beta}}_t(a) - \boldsymbol{\beta}(a)\}$ is asymptotically normal. The proof is to generalize Theorem 3.1 in Chen et al. (2020) by considering commonly used bandit algorithms, including UCB, TS, and EG here. We complete the proof in the following four steps:

- Step 1: Prove that $\boldsymbol{\eta}_1 = (1/\sqrt{t}) \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i e_i \xrightarrow{D} \mathcal{N}_d(\mathbf{0}_d, G_a)$, where G_a is the variance matrix to be spesified shortly.
- Step 2: Prove that $\boldsymbol{\xi} = \left\{ (1/t) \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + (\omega/t) \mathbf{I}_d \right\}^{-1} \xrightarrow{p} \sigma_a^2 G_a^{-1}$, where $\sigma_a^2 = \mathbb{E}(e_t^2 | a_t = a)$ for $a = 0, 1$.
- Step 3: Prove that $\boldsymbol{\eta}_2 = \left\{ (1/t) \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + (\omega/t) \mathbf{I}_d \right\}^{-1} (1/\sqrt{t}) \boldsymbol{\beta}(a) \xrightarrow{p} \mathbf{0}_d$.
- Step 4: Combine above results in steps 1-3 using Slutsky's theorem.

Step 1: We first focus on proving that $\boldsymbol{\eta}_1 = (1/\sqrt{t}) \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i e_i \xrightarrow{D} \mathcal{N}_d(\mathbf{0}_d, G_a)$. Using Cramer-Wold device, it suffices to show that for any $\mathbf{v} \in \mathbb{R}^d$,

$$\boldsymbol{\eta}_1(\mathbf{v}) \equiv \frac{1}{\sqrt{t}} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i e_i \xrightarrow{D} \mathcal{N}_d(0, \mathbf{v}^\top G_a \mathbf{v}).$$

Note that $\mathbb{E}(\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i e_i | \mathcal{H}_{i-1}) = \mathbb{E}(\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i | \mathcal{H}_{i-1}) \mathbb{E}(e_i | \mathcal{H}_{i-1}, a_i = a) = 0$, we have that $\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i e_i$ is a Martingale difference sequence. We next show the asymptotic normality of $\boldsymbol{\eta}_1(\mathbf{v})$ using Martingale central limit theorem, by the following two parts: i) check the conditional Lindeberg condition; ii) derivative the limit of the conditional variance.

Firstly, we check the conditional Lindeberg condition. For any $\delta > 0$, denote

$$\psi = \sum_{i=1}^t \mathbb{E} \left[\frac{1}{t} \mathbb{I}(a_i = a) (\mathbf{v}^\top \mathbf{x}_i)^2 e_i^2 \mathbb{I} \left\{ \left| \frac{1}{\sqrt{t}} \mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i e_i \right| > \delta \right\} | \mathcal{H}_{i-1} \right]. \quad (\text{B.22})$$

Notice that $(\mathbf{v}^\top \mathbf{x}_i)^2 \leq \|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d$, we have

$$\mathbb{I} \left\{ \left| \frac{1}{\sqrt{t}} \mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i e_i \right| > \delta \right\} \leq \mathbb{I} \left\{ \mathbb{I}(a_i = a) e_i^2 \|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d > t \delta^2 \right\} = \mathbb{I} \left\{ \mathbb{I}(a_i = a) e_i^2 > \frac{t \delta^2}{\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d} \right\}.$$

Combining this with (B.22), we obtain that

$$\psi \leq \frac{\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d}{t} \sum_{i=1}^t \mathbb{E} \left(\mathbb{I}(a_i = a) e_i^2 \mathbb{I} \left\{ \mathbb{I}(a_i = a) e_i^2 > \frac{t \delta^2}{\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d} \right\} | \mathcal{H}_{i-1} \right), \quad (\text{B.23})$$

where the right hand side equals

$$\frac{\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d}{t} \sum_{i=1}^t \mathbb{E}(\mathbb{I}(a_i = a) \mid \mathcal{H}_{i-1}) \mathbb{E} \left(e_i^2 \mathbb{I} \left\{ \mathbb{I}(a_i = a) e_i^2 > \frac{\delta^2 t}{\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d} \right\} \mid \mathcal{H}_{i-1} \right).$$

Then, we can further write (B.23) as

$$\psi \leq \frac{\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d}{t} \sum_{i=1}^t \mathbb{E} \left(e_{i(a)}^2 \mathbb{I} \left\{ e_{i(a)}^2 > \frac{\delta^2 t}{\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d} \right\} \mid \mathcal{H}_{i-1} \right). \quad (\text{B.24})$$

where $e_{i(a)} = e_i$ when $a_i = a$ and 0 otherwise. Since e_i conditioned on a_i are *i.i.d.*, $\forall i$, we have the right hand side of the above inequality equals

$$\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d \mathbb{E} \left(e^2 \mathbb{I} \left\{ e^2 > \frac{t \delta^2}{\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d} \right\} \right),$$

where e is the random variable given by $e_i \mid \mathcal{H}_{i-1}$. Note that $e^2 \mathbb{I} \{e^2 > t \delta^2 / (\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d)\}$ is dominated by e^2 with $\mathbb{E} e^2 < \infty$ and converges to 0, as $t \rightarrow \infty$. Then, by Dominated Convergence Theorem, the results in (B.24) can be further bounded by

$$\psi \leq \frac{\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d}{t} \sum_{i=1}^t \mathbb{E} \left(e^2 \mathbb{I} \left\{ e^2 > \frac{t \delta^2}{\|\mathbf{v}\|_2^2 L_{\mathbf{x}}^2 d} \right\} \right) \rightarrow 0, \text{ as } t \rightarrow \infty.$$

Therefore, conditional Lindeberg condition holds.

Secondly, we derive the limit of the conditional variance. Notice that

$$\begin{aligned} \frac{1}{t} \sum_{i=1}^t \mathbb{E} \left[\mathbb{I}(a_i = a) (\mathbf{v}^\top \mathbf{x}_i)^2 e_i^2 \mid \mathcal{H}_{i-1} \right] &= \frac{1}{t} \sum_{i=1}^t \mathbb{E} \left\{ \mathbb{E} \left[\mathbb{I}(a_i = a) (\mathbf{v}^\top \mathbf{x}_i)^2 e_i^2 \mid a_i, \mathbf{x}_i \right] \mid \mathcal{H}_{i-1} \right\} \\ &= \frac{1}{t} \sum_{i=1}^t \mathbb{E} \left\{ \mathbb{I}(a_i = a) (\mathbf{v}^\top \mathbf{x}_i)^2 \mathbb{E} [e_i^2 \mid a_i = a, \mathbf{x}_i] \mid \mathcal{H}_{i-1} \right\}. \end{aligned}$$

Since e_t is independent of \mathcal{H}_{i-1} and \mathbf{x}_i given a_t , and $\mathbb{E} [e_i^2 \mid a_i = a] = \sigma_a^2$, we have

$$\begin{aligned} \frac{1}{t} \sum_{i=1}^t \mathbb{E} \left[\mathbb{I}(a_i = a) (\mathbf{v}^\top \mathbf{x}_i)^2 e_i^2 \mid \mathcal{H}_{i-1} \right] &= \frac{1}{t} \sum_{i=1}^t \mathbb{E} \left[\mathbb{I}(a_i = a) (\mathbf{v}^\top \mathbf{x}_i)^2 \mathbb{E} [e_i^2 \mid a_i = a] \mid \mathcal{H}_{i-1} \right] \\ &= \frac{1}{t} \sum_{i=1}^t \mathbb{E} \left[\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \sigma_a^2 \mid \mathcal{H}_{i-1} \right] \\ &= \frac{1}{t} \sum_{i=1}^t \sigma_a^2 \mathbb{E} \left[\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \mid \mathcal{H}_{i-1} \right], \end{aligned}$$

where

$$\begin{aligned}
\mathbb{E} [\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \mid \mathcal{H}_{i-1}] &= \mathbb{E} \{ \mathbb{E} [\mathbb{I}(a_i \neq \pi^*(\mathbf{x}_i)) \mathbb{I}(\pi^*(\mathbf{x}_i) \neq a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \mid \mathbf{x}_i] \mid \mathcal{H}_{i-1} \} \\
&+ \mathbb{E} \{ \mathbb{E} [\mathbb{I}(a_i = \pi^*(\mathbf{x}_i)) \mathbb{I}(\pi^*(\mathbf{x}_i) = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \mid \mathbf{x}_i] \mid \mathcal{H}_{i-1} \} \\
&= \mathbb{E} \{ \mathbb{E} [\mathbb{I}(a_i \neq \pi^*(\mathbf{x}_i)) \mid \mathbf{x}_i, \mathcal{H}_{i-1}] \mathbb{I}(\pi^*(\mathbf{x}_i) \neq a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \} \\
&+ \mathbb{E} \{ \mathbb{E} [\mathbb{I}(a_i = \pi^*(\mathbf{x}_i)) \mid \mathbf{x}_i, \mathcal{H}_{i-1}] \mathbb{I}(\pi^*(\mathbf{x}_i) = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \}.
\end{aligned}$$

Here, the first equation comes from iteration expectation over \mathbf{x}_i and the fact that $\mathbb{I}(\pi^*(\mathbf{x}_i) \neq a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v}$ is a constant given \mathbf{x}_i and \mathcal{H}_{i-1} .

$$\begin{aligned}
\mathbb{I}(a_i = a) &= \mathbb{I}(a_i = a \neq \pi^*(\mathbf{x}_i)) + \mathbb{I}(a_i = a = \pi^*(\mathbf{x}_i)) \\
&= \mathbb{I}(a_i \neq \pi^*(\mathbf{x}_i)) \mathbb{I}(\pi^*(\mathbf{x}_i) \neq a) + \mathbb{I}(a_i = \pi^*(\mathbf{x}_i)) \mathbb{I}(\pi^*(\mathbf{x}_i) = a),
\end{aligned}$$

and the second equation is owing to the fact that $\mathbb{I}(\pi^*(\mathbf{x}_i) = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v}$ is a constant given \mathbf{x}_i and independent of \mathcal{H}_{i-1} . Define

$$\nu_i(\mathbf{x}_i, \mathcal{H}_{i-1}) \equiv \Pr \{a_i \neq \pi^*(\mathbf{x}_i) \mid \mathbf{x}_i, \mathcal{H}_{i-1}\} = \mathbb{E} [\mathbb{I} \{a_i \neq \pi^*(\mathbf{x}_i)\} \mid \mathbf{x}_i, \mathcal{H}_{i-1}], \quad (\text{B.25})$$

then the conditional variance can be expressed as

$$\begin{aligned}
\frac{1}{t} \sum_{i=1}^t \mathbb{E} [\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \mid \mathcal{H}_{i-1}] &= \frac{1}{t} \sum_{i=1}^t \mathbb{E} \{ \nu_i(\mathbf{x}_i, \mathcal{H}_{i-1}) \mathbb{I}(\pi^*(\mathbf{x}_i) \neq a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \} \\
&+ \frac{1}{t} \sum_{i=1}^t \mathbb{E} [\{1 - \nu_i(\mathbf{x}_i, \mathcal{H}_{i-1})\} \mathbb{I}(\pi^*(\mathbf{x}_i) = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v}],
\end{aligned}$$

which can be expressed as

$$\begin{aligned}
&\frac{1}{t} \sum_{i=1}^t \int \nu_i(\mathbf{x}, \mathcal{H}_{i-1}) \mathbb{I} \{ \mathbf{x}^\top \boldsymbol{\beta}(a) < \mathbf{x}^\top \boldsymbol{\beta}(1-a) \} \mathbf{v}^\top \mathbf{x} \mathbf{x}^\top \mathbf{v} dP_{\mathcal{X}} \\
&+ \frac{1}{t} \sum_{i=1}^t \int \{1 - \nu_i(\mathbf{x}, \mathcal{H}_{i-1})\} \mathbb{I} \{ \mathbf{x}^\top \boldsymbol{\beta}(a) \geq \mathbf{x}^\top \boldsymbol{\beta}(1-a) \} \mathbf{v}^\top \mathbf{x} \mathbf{x}^\top \mathbf{v} dP_{\mathcal{X}} \\
&= \int \frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1}) \mathbb{I} \{ \mathbf{x}^\top \boldsymbol{\beta}(a) < \mathbf{x}^\top \boldsymbol{\beta}(1-a) \} \mathbf{v}^\top \mathbf{x} \mathbf{x}^\top \mathbf{v} dP_{\mathcal{X}} \\
&+ \int \{1 - \frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1})\} \mathbb{I} \{ \mathbf{x}^\top \boldsymbol{\beta}(a) \geq \mathbf{x}^\top \boldsymbol{\beta}(1-a) \} \mathbf{v}^\top \mathbf{x} \mathbf{x}^\top \mathbf{v} dP_{\mathcal{X}}.
\end{aligned}$$

Since $\lim_{i \rightarrow \infty} \Pr\{a_i \neq \pi^*(\mathbf{x})\} = \kappa_\infty(\mathbf{x})$, we have for any $\epsilon > 0$, there exist a constant $t_0 > 0$ such that $|\Pr\{a_i \neq \pi^*(\mathbf{x})\} - \kappa_\infty(\mathbf{x})| < \epsilon$ for all $i \geq t_0$. Therefore, for the expectation of $\frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1})$ over the history, we have

$$\mathbb{E} \left[\frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1}) \right] = \frac{1}{t} \sum_{i=1}^t \mathbb{E} [\nu_i(\mathbf{x}, \mathcal{H}_{i-1})] = \frac{1}{t} \sum_{i=1}^t \Pr\{a_i \neq \pi^*(\mathbf{x})\}.$$

It follows immediately that

$$\begin{aligned} & \mathbb{E} \left[\frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1}) \right] - \kappa_\infty(\mathbf{x}) \\ &= \frac{1}{t} \sum_{i=1}^{t_0} [\Pr\{a_i \neq \pi^*(\mathbf{x})\} - \kappa_\infty(\mathbf{x})] + \frac{1}{t} \sum_{i=t_0}^t [\Pr\{a_i \neq \pi^*(\mathbf{x})\} - \kappa_\infty(\mathbf{x})]. \end{aligned}$$

Therefore, by the triangle inequality, we have

$$\begin{aligned} & \left| \mathbb{E} \left[\frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1}) \right] - \kappa_\infty(\mathbf{x}) \right| \\ &= \frac{1}{t} \sum_{i=1}^{t_0} |\Pr\{a_i \neq \pi^*(\mathbf{x})\} - \kappa_\infty(\mathbf{x})| + \frac{1}{t} \sum_{i=t_0}^t |\Pr\{a_i \neq \pi^*(\mathbf{x})\} - \kappa_\infty(\mathbf{x})| \\ &< \frac{1}{t} \sum_{i=1}^{t_0} [|\Pr\{a_i \neq \pi^*(\mathbf{x})\}| + |\kappa_\infty(\mathbf{x})|] + \frac{1}{t} \sum_{i=t_0}^t \epsilon \\ &\leq \frac{1}{t} \sum_{i=1}^{t_0} 2 + \frac{1}{t} \sum_{i=t_0}^t \epsilon = \frac{2t_0}{t} + \frac{t-t_0}{t} \epsilon. \end{aligned}$$

Since the above equation holds for any $\epsilon > 0$ and $0 < \frac{t-t_0}{t} < 1$, we have

$$\left| \mathbb{E} \left[\frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1}) \right] - \kappa_\infty(\mathbf{x}) \right| \leq \frac{2t_0}{t},$$

which goes to zero as $t \rightarrow \infty$. Thus,

$$\mathbb{E} \left[\frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1}) \right] = \kappa_\infty(\mathbf{x}) + o_p(1). \quad (\text{B.26})$$

Next, we consider the variance of $\frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1})$. Denote $\mathbb{E} \left[\frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1}) \right] = \mu_\nu(\mathbf{x})$, we have $\mu_\nu(\mathbf{x}) = \kappa_\infty(\mathbf{x}) + o_p(1)$. Notice that $\nu_i(\mathbf{x}_i, \mathcal{H}_{i-1}) \equiv \Pr \{a_i \neq \pi^*(\mathbf{x}_i) | \mathbf{x}_i, \mathcal{H}_{i-1}\} \in [0, 1]$, by Lemma B.1, we have

$$\text{Var} \left[\frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1}) \right] \leq \mu_\nu(\mathbf{x}) - \mu_\nu(\mathbf{x})^2 = \kappa_\infty(\mathbf{x}) \{1 - \kappa_\infty(\mathbf{x})\} + o_p(1),$$

which goes to zero as $t \rightarrow \infty$. Combined with Equation (B.26), it follows immediately that as t goes to ∞ , we have

$$\frac{1}{t} \sum_{i=1}^t \nu_i(\mathbf{x}, \mathcal{H}_{i-1}) \rightarrow \kappa_\infty(\mathbf{x}).$$

Therefore, as t goes to ∞ , we have $\frac{1}{t} \sum_{i=1}^t \mathbb{E} [\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} | \mathcal{H}_{i-1}]$ converges to

$$\begin{aligned} & \int \kappa_\infty(\mathbf{x}) \mathbb{I}\{\mathbf{x}^\top \boldsymbol{\beta}(a) < \mathbf{x}^\top \boldsymbol{\beta}(1-a)\} \mathbf{v}^\top \mathbf{x} \mathbf{x}^\top \mathbf{v} dP_{\mathcal{X}} \\ & + \int \{1 - \kappa_\infty(\mathbf{x})\} \mathbb{I}\{\mathbf{x}^\top \boldsymbol{\beta}(a) \geq \mathbf{x}^\top \boldsymbol{\beta}(1-a)\} \mathbf{v}^\top \mathbf{x} \mathbf{x}^\top \mathbf{v} dP_{\mathcal{X}}. \end{aligned} \quad (\text{B.27})$$

Thus, following the similar arguments in S1.2 in Chen et al. (2020), we have

$$\boldsymbol{\eta}_1(\mathbf{v}) = \frac{1}{\sqrt{t}} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i e_i \xrightarrow{D} \mathcal{N}_d(0, \mathbf{v}^\top G_a \mathbf{v}),$$

where

$$\begin{aligned} \mathbf{v}^\top G_a \mathbf{v} = & \sigma_a^2 \left\{ \int \kappa_\infty(\mathbf{x}) \mathbb{I}\{\mathbf{x}^\top \boldsymbol{\beta}(a) < \mathbf{x}^\top \boldsymbol{\beta}(1-a)\} \mathbf{v}^\top \mathbf{x} \mathbf{x}^\top \mathbf{v} dP_{\mathcal{X}} \right. \\ & \left. + \int \{1 - \kappa_\infty(\mathbf{x})\} \mathbb{I}\{\mathbf{x}^\top \boldsymbol{\beta}(a) \geq \mathbf{x}^\top \boldsymbol{\beta}(1-a)\} \mathbf{v}^\top \mathbf{x} \mathbf{x}^\top \mathbf{v} dP_{\mathcal{X}} \right\}. \end{aligned}$$

Finally, by Martingale Central Limit Theorem, we have

$$\boldsymbol{\eta}_1 = \frac{1}{\sqrt{t}} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i e_i \xrightarrow{D} \mathcal{N}_d(0, G_a),$$

where

$$G_a = \sigma_a^2 \left\{ \int \kappa_\infty(\mathbf{x}) \mathbb{I}\{\mathbf{x}^\top \boldsymbol{\beta}(a) < \mathbf{x}^\top \boldsymbol{\beta}(1-a)\} \mathbf{x} \mathbf{x}^\top dP_{\mathcal{X}} + \int \{1 - \kappa_\infty(\mathbf{x})\} \mathbb{I}\{\mathbf{x}^\top \boldsymbol{\beta}(a) \geq \mathbf{x}^\top \boldsymbol{\beta}(1-a)\} \mathbf{x} \mathbf{x}^\top dP_{\mathcal{X}} \right\}. \quad (\text{B.28})$$

The first part is thus completed.

Step 2: We next show that $\boldsymbol{\xi} = \left\{ (1/t) \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}^{-1} \xrightarrow{p} \sigma_a^2 G_a^{-1}$, which is sufficient to find the limit of $\left\{ (1/t) \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}$. By Lemma 6 in Chen et al. (2020), it suffices to show the limit of $\frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v}$ for any $\mathbf{v} \in \mathbb{R}^d$.

Since $\Pr(|\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v}| > h) \leq \Pr(|\mathbf{v}^\top \mathbf{x} \mathbf{x}^\top \mathbf{v}| > h)$ for each $h > 0$ and $i \leq 1$, by Theorem 2.19 in Hall and Heyde (2014), we have

$$\frac{1}{t} \sum_{i=1}^t [\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} - \mathbb{E}\{\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \mid \mathcal{H}_{i-1}\}] \xrightarrow{p} 0, \text{ as } t \rightarrow \infty. \quad (\text{B.29})$$

Recall the results in (B.26) and (B.28), we have $\mathbb{E}\{\mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \mid \mathcal{H}_{i-1}\} = G_a / \sigma_a^2$. Combining this with (B.29), we have

$$\frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{v}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v} \xrightarrow{p} \frac{G_a}{\sigma_a^2}, \text{ as } t \rightarrow \infty.$$

By Lemma 6 in Chen et al. (2020) and Continuous Mapping Theorem, we further have

$$\boldsymbol{\xi} = \left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}^{-1} \xrightarrow{p} \sigma_a^2 G_a^{-1}.$$

Step 3: We focus on proving $\boldsymbol{\eta}_2 = \left\{ (1/t) \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + (\omega/t) \mathbf{I}_d \right\}^{-1} (1/\sqrt{t}) \boldsymbol{\beta}(a) \xrightarrow{p} \mathbf{0}_d$ next. This suffices to show that $\mathbf{b}_i^\top \boldsymbol{\eta}_2 \xrightarrow{p} 0$ holds for any standard basis $\mathbf{b}_i \in \mathbb{R}^d$. Since

$$\begin{aligned} \mathbf{b}_i^\top \boldsymbol{\eta}_2 &= \left\| \mathbf{b}_i^\top \left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}^{-1} \frac{1}{\sqrt{t}} \boldsymbol{\beta}(a) \right\|_2 \\ &\leq \frac{1}{\sqrt{t}} \|\mathbf{b}_i^\top\|_2 \left\| \left\{ \frac{1}{t} \sum_{i=1}^t \mathbb{I}(a_i = a) \mathbf{x}_i \mathbf{x}_i^\top + \frac{1}{t} \omega \mathbf{I}_d \right\}^{-1} \right\|_2 \|\boldsymbol{\beta}(a)\|_2, \end{aligned}$$

and by (B.4), we have

$$\mathbf{b}_i^\top \boldsymbol{\eta}_2 \leq \frac{\|\boldsymbol{\beta}(a)\|_2}{\sqrt{tp_t^2 \lambda} + \frac{1}{\sqrt{t}} \omega}.$$

Thus, we have $\mathbf{b}_i^\top \boldsymbol{\eta}_2 \xrightarrow{p} 0$, as $tp_t^2 \rightarrow \infty$.

Step 4: Finally, we combine the above results using Slutsky's theorem, and conclude that

$$\sqrt{t}\{\widehat{\boldsymbol{\beta}}_t(a) - \boldsymbol{\beta}(a)\} = \boldsymbol{\xi} \boldsymbol{\eta}_1 + \boldsymbol{\eta}_2 \xrightarrow{D} \mathcal{N}_d(\mathbf{0}_d, \sigma_a^4 G_a^{-1}),$$

where G_a is defined in (B.28). Denote the variance term as

$$\begin{aligned} \sigma_{\boldsymbol{\beta}(a)}^2 &= \sigma_a^2 \left[\int \kappa_\infty(\mathbf{x}) \mathbb{I}\{\mathbf{x}^\top \boldsymbol{\beta}(a) < \mathbf{x}^\top \boldsymbol{\beta}(1-a)\} \mathbf{x} \mathbf{x}^\top dP_{\mathcal{X}} \right. \\ &\quad \left. + \int \{1 - \kappa_\infty(\mathbf{x})\} \mathbb{I}\{\mathbf{x}^\top \boldsymbol{\beta}(a) \geq \mathbf{x}^\top \boldsymbol{\beta}(1-a)\} \mathbf{x} \mathbf{x}^\top dP_{\mathcal{X}} \right]^{-1}, \end{aligned}$$

with $\sigma_a^2 = \mathbb{E}(e_t^2 | a_t = a)$ denoting the conditional variance of e_t given $a_t = a$, for $a = 0, 1$, we have

$$\sqrt{t}\{\widehat{\boldsymbol{\beta}}_t(a) - \boldsymbol{\beta}(a)\} \xrightarrow{D} \mathcal{N}_d\{\mathbf{0}_d, \sigma_{\boldsymbol{\beta}(a)}^2\}.$$

The proof is hence completed.

B.5 Proof of Theorem 3

Finally, we prove the asymptotic normality of the proposed value estimator under DREAM in Theorem 3 in this section. The proof consists of four steps. In step 1, we aim to show

$$\widehat{V}_T = \widetilde{V}_T + o_p(T^{-1/2}),$$

where

$$\widetilde{V}_T = \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\}}{1 - \kappa_t(\mathbf{x}_t)} \left[r_t - \mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} \right] + \mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\}.$$

Next, in Step 2, we establish

$$\widetilde{V}_T = \overline{V}_T + o_p(T^{-1/2}),$$

where

$$\overline{V}_T = \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \left[r_t - \mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} \right] + \mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\}.$$

The above two steps yields that

$$\widehat{V}_T = \overline{V}_T + o_p(T^{-1/2}). \quad (\text{B.30})$$

Then, in Step 3, based on (B.30) and Martingale Central Limit Theorem, we show

$$\sqrt{T}(\widehat{V}_T - V^*) \xrightarrow{D} \mathcal{N}(0, \sigma_{DR}^2),$$

with

$$\sigma_{DR}^2 = \int_{\mathbf{x}} \frac{\pi^*(\mathbf{x})\sigma_1^2 + \{1 - \pi^*(\mathbf{x})\}\sigma_0^2}{1 - \kappa_\infty(\mathbf{x})} dP_{\mathcal{X}} + \text{Var}[\mu\{\mathbf{x}, \pi^*(\mathbf{x})\}],$$

where $\sigma_a^2 = \mathbb{E}(e_t^2 | a_t = a)$ for $a = 0, 1$, and $\kappa_t \rightarrow \kappa_\infty$ as $t \rightarrow \infty$.

Lastly, in Step 4, we show the variance estimator in Equation (5) in the main paper is a consistent estimator of σ_{DR}^2 . The proof for Theorem 3 is thus completed.

Step 1: We first show $\widehat{V}_T = \widetilde{V}_T + o_p(T^{-1/2})$. To this end, define a middle term as

$$\widetilde{\phi}_T = \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\}}{1 - \kappa_t(\mathbf{x}_t)} \left[r_t - \widehat{\mu}_{t-1}\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} \right] + \widehat{\mu}_{t-1}\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\}.$$

Thus, it suffices to show $\widehat{V}_T = \widetilde{\phi}_T + o_p(T^{-1/2})$ and $\widetilde{\phi}_T = \widetilde{V}_T + o_p(T^{-1/2})$.

Firstly, we have

$$\begin{aligned} \widehat{V}_T - \widetilde{\phi}_T &= \frac{1}{T} \sum_{t=1}^T \left[\frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\}}{1 - \widehat{\kappa}_t(\mathbf{x}_t)} - \frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\}}{1 - \kappa_t(\mathbf{x}_t)} \right] \left[r_t - \widehat{\mu}_{t-1}\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} \right] \quad (\text{B.31}) \\ &= \frac{1}{T} \sum_{t=1}^T \{\widehat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)\} \left[\frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\} \left[r_t - \widehat{\mu}_{t-1}\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} \right]}{\{1 - \widehat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right]. \end{aligned}$$

We can further decompose (B.31) by

$$\begin{aligned}
\widehat{V}_T - \widetilde{\phi}_T &= \frac{1}{T} \sum_{t=1}^T \{ \widehat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t) \} \\
&\quad \left[\frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\} \left[r_t - \mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} + \mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} - \widehat{\mu}_{t-1}\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} \right]}{\{1 - \widehat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right] \\
&= \frac{1}{T} \sum_{t=1}^T \{ \widehat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t) \} \left[\frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\} \left[r_t - \mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} \right]}{\{1 - \widehat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right] \\
&\quad + \frac{1}{T} \sum_{t=1}^T \{ \widehat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t) \} \left[\frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\} \left[\mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} - \widehat{\mu}_{t-1}\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} \right]}{\{1 - \widehat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right]. \tag{B.32}
\end{aligned}$$

We first show that the first term in (B.32) is $o_p(T^{-1/2})$. Define a class of function

$$\mathcal{F}_\kappa(\mathbf{x}, a, r) = \left\{ \{ \widehat{\kappa}(\mathbf{x}) - \kappa(\mathbf{x}) \} \left[\frac{\mathbb{I}\{a = \pi(\mathbf{x})\} \left[r - \mu\{\mathbf{x}, \pi(\mathbf{x})\} \right]}{\{1 - \widehat{\kappa}(\mathbf{x})\}\{1 - \kappa(\mathbf{x})\}} \right] : \widehat{\kappa}(\cdot), \kappa(\cdot) \in \Lambda, \pi(\cdot) \in \Pi \right\},$$

where Π and Λ are two classes of functions that maps context $\mathbf{x} \in \mathcal{X}$ to a probability. Define the supremum of the empirical process indexed by \mathcal{F}_κ as

$$\|\mathbb{G}_n\|_{\mathcal{F}} \equiv \sup_{\pi \in \Pi} \frac{1}{\sqrt{T}} \sum_{t=1}^T [\mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) - \mathbb{E}\{\mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) | \mathcal{H}_{t-1}\}]. \tag{B.33}$$

Notice that

$$\begin{aligned}
&\mathbb{E}\{\mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) | \mathcal{H}_{t-1}\} \\
&= \mathbb{E} \left(\{ \widehat{\kappa}(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t) \} \left[\frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\} \left[r_t - \mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} \right]}{\{1 - \widehat{\kappa}(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right] \mid \mathcal{H}_{t-1} \right) \\
&= \mathbb{E} \left(\{ \widehat{\kappa}(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t) \} \left[\frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\} e_t}{\{1 - \widehat{\kappa}(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right] \mid \mathcal{H}_{t-1} \right),
\end{aligned}$$

by the definitions and thus, using the iteration of expectation, we have

$$\begin{aligned}
& \mathbb{E}\{\mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) | \mathcal{H}_{t-1}\} \\
&= \mathbb{E} \left\{ \mathbb{E} \left(\{\hat{\kappa}(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)\} \left[\frac{\mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} e_t}{\{1 - \hat{\kappa}(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right] \mid a_t, \mathbf{x}_t \right) \mid \mathcal{H}_{t-1} \right\} \\
&= \mathbb{E} \left\{ \{\hat{\kappa}(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)\} \left[\frac{\mathbb{E}\{\mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} \mid \mathbf{x}_t\}}{\{1 - \hat{\kappa}(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right] \mathbb{E}(e_t \mid a_t, \mathbf{x}_t) \mid \mathcal{H}_{t-1} \right\} = 0,
\end{aligned}$$

where the last equation is due to the definition of the noise e_t . Therefore, Equation (B.33) can be further written as

$$||\mathbb{G}_n||_{\mathcal{F}} \equiv \sup_{\pi \in \Pi} \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t).$$

Next, we show the second moment is bounded by

$$\begin{aligned}
& \mathbb{E} \left\{ \left(\{\hat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)\} \left[\frac{\mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} [r_t - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\}]}{\{1 - \hat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right] \right)^2 \mid \mathcal{H}_{t-1} \right\} \\
&= \mathbb{E} \left\{ \left[\frac{\{\hat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)\}}{\{1 - \hat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right]^2 \mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} [r_t - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\}]^2 \mid \mathcal{H}_{t-1} \right\}, \\
&= \mathbb{E} \left\{ \left[\frac{\{\hat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)\}}{\{1 - \hat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right]^2 \mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} e_t^2 \mid \mathcal{H}_{t-1} \right\}
\end{aligned}$$

by the definitions and thus, using the iteration of expectation, we have

$$\begin{aligned}
& \mathbb{E} \left\{ \left(\{\hat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)\} \left[\frac{\mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} [r_t - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\}]}{\{1 - \hat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right] \right)^2 \mid \mathcal{H}_{t-1} \right\} \\
&= \mathbb{E} \left(\mathbb{E} \left\{ \left[\frac{\{\hat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)\}}{\{1 - \hat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right]^2 \mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} e_t^2 \mid a_t, \mathbf{x}_t \right\} \mid \mathcal{H}_{t-1} \right) \\
&= \mathbb{E} \left(\left[\frac{\{\hat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)\}}{\{1 - \hat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right]^2 \mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} \mathbb{E}\{e_t^2 \mid a_t, \mathbf{x}_t\} \mid \mathcal{H}_{t-1} \right) \\
&= \mathbb{E} \left(\left[\frac{1}{1 - \hat{\kappa}_t(\mathbf{x}_t)} - \frac{1}{1 - \kappa_t(\mathbf{x}_t)} \right]^2 \mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} \sigma_{a_t}^2 \mid \mathcal{H}_{t-1} \right).
\end{aligned}$$

Notice that $\widehat{\kappa}_t(\mathbf{x}_t) \leq C_1 < 1$ and $\kappa_t(\mathbf{x}_t) \leq C_1 < 1$ for sure for some constant $C_1 < 1$ (by definition of a valid bandit algorithm and results of Theorem 1), we have

$$\left[\frac{1}{1 - \widehat{\kappa}_t(\mathbf{x}_t)} - \frac{1}{1 - \kappa_t(\mathbf{x}_t)} \right]^2 \leq \left\{ \left| \frac{1}{1 - \widehat{\kappa}_t(\mathbf{x}_t)} \right| + \left| \frac{1}{1 - \kappa_t(\mathbf{x}_t)} \right| \right\}^2 \leq \left(\frac{2}{1 - C_1} \right)^2 \equiv C_2,$$

where C_2 is a bounded constant. Thus we have

$$\begin{aligned} & \mathbb{E} \left(\left[\frac{1}{1 - \widehat{\kappa}_t(\mathbf{x}_t)} - \frac{1}{1 - \kappa_t(\mathbf{x}_t)} \right]^2 \mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\} \sigma_{a_t}^2 \mid \mathcal{H}_{t-1} \right) \\ & \leq \mathbb{E} (C_2 \mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\} \sigma_{a_t}^2 \mid \mathcal{H}_{t-1}) \\ & \leq C_2 \max\{\sigma_0^2, \sigma_1^2\}. \end{aligned}$$

Therefore, for the second moment of the inner term of the first term, we have

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E} \left\{ \left(\frac{1}{\sqrt{T}} \{ \widehat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t) \} \left[\frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\} [r_t - \mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\}]}{\{1 - \widehat{\kappa}_t(\mathbf{x}_t)\} \{1 - \kappa_t(\mathbf{x}_t)\}} \right] \right)^2 \mid \mathcal{H}_{t-1} \right\} \\ & \leq \sum_{t=1}^T \frac{C_2 \max\{\sigma_0^2, \sigma_1^2\}}{T} = C_2 \max\{\sigma_0^2, \sigma_1^2\} < \infty. \end{aligned}$$

Therefore, we have

$$d_1(f) \equiv \|\mathbb{E}(\mathcal{F}_\kappa(\mathbf{x}_1, a_1, r_1) \mid \mathcal{H}_0)\|_\infty < \infty,$$

and

$$d_2(f) \equiv \|\mathbb{E}((\mathcal{F}_\kappa(\mathbf{x}_1, a_1, r_1))^2 \mid \mathcal{H}_0)\|_\infty^{1/2} < \infty.$$

It follows from the maximal inequality developed in Section 4.2 of Dedecker and Louhichi (2002) that there exist some constant $K \geq 1$ such that

$$\mathbb{E}[\|\mathbb{G}_n\|_{\mathcal{F}}] \lesssim K \left(\sqrt{p} d_2(f) + \frac{1}{\sqrt{T}} \left\| \max_{1 \leq t \leq T} |\mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) - \mathbb{E}(\mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) \mid \mathcal{H}_{t-1})| \right\| \right).$$

The above right-hand-side is upper bounded by

$$O(1) \sqrt{T^{-1/2}},$$

where $O(1)$ denotes some universal constant. Hence, we have

$$\mathbb{E} \left[\|\mathbb{G}_n\|_{\mathcal{F}} \right] = \mathcal{O}_p(T^{-1/2}). \quad (\text{B.34})$$

Combined with Equation (B.33), we have

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T \mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) = \mathcal{O}_p(T^{-1/2}).$$

Therefore, for the first term in (B.32), we have

$$\frac{1}{T} \sum_{t=1}^T \{\widehat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)\} \left[\frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\} \left[r_t - \mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} \right]}{\{1 - \widehat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right] = \mathcal{O}_p(T^{-1}) = o_p(T^{-1/2}).$$

Then we consider the second term in (B.32), where

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \{\widehat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)\} \left[\frac{\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\} \left[\mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} - \widehat{\mu}_{t-1}\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} \right]}{\{1 - \widehat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}} \right] \\ & \leq \frac{1}{T} \sum_{t=1}^T B_\kappa |\widehat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)| |\mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} - \widehat{\mu}_{t-1}\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\}|, \end{aligned}$$

for some B_κ as the bound of $\mathbb{I}\{a_t = \widehat{\pi}_t(\mathbf{x}_t)\}/[\{1 - \widehat{\kappa}_t(\mathbf{x}_t)\}\{1 - \kappa_t(\mathbf{x}_t)\}]$. By Cauchy-Schwartz inequality, we have the above term further bounded by

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T B_\kappa |\widehat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)| |\mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} - \widehat{\mu}_{t-1}\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\}| \\ & \leq B_\kappa \sqrt{\frac{1}{T} \sum_{t=1}^T |\widehat{\kappa}_t(\mathbf{x}_t) - \kappa_t(\mathbf{x}_t)|^2 \frac{1}{T} \sum_{t=1}^T |\mu\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\} - \widehat{\mu}_{t-1}\{\mathbf{x}_t, \widehat{\pi}_t(\mathbf{x}_t)\}|^2}. \end{aligned}$$

Given Assumption 4.4, we have the above bounded by $o_p(T^{-1/2})$, and thus the second term in (B.32) is $o_p(T^{-1/2})$.

Therefore, we have $\widehat{V}_T = \widetilde{\phi}_T + o_p(T^{-1/2})$ hold.

Then, we focus on proving

$$\tilde{\phi}_T - \tilde{V}_T = \frac{1}{T} \sum_{t=1}^T \left[\frac{\mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\}}{1 - \kappa_t(\mathbf{x}_t)} - 1 \right] \left[\mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} - \hat{\mu}_{t-1}\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right], \quad (\text{B.35})$$

is $o_p(T^{-1/2})$. Specifically, since

$$\mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} - \hat{\mu}_{t-1}\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} = \mathcal{O}_p(t^{-1/2}),$$

by Theorem 2, and notice that

$$-1 \leq \frac{\mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\}}{1 - \kappa_t(\mathbf{x}_t)} - 1 \leq \frac{1}{1 - \kappa_t(\mathbf{x}_t)} - -1 \leq \frac{1}{1 - C_1} - 1,$$

which is bounded, we have the inner part of (B.35) satisfies

$$\left[\frac{\mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\}}{1 - \kappa_t(\mathbf{x}_t)} - 1 \right] \left[\mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} - \hat{\mu}_{t-1}\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right] = \mathcal{O}_p(t^{-1/2}) = o_p(t^{-1/2-\alpha}),$$

for $0 < \alpha < 1/2$. Therefore, by Lemma 6 in Luedtke and Van Der Laan (2016), we have

$$\tilde{\phi}_T - \tilde{V}_T = \frac{1}{T} \sum_{t=1}^T o_p(t^{-1/2-\alpha}) = o_p(T^{-1/2-\alpha}) = o_p(T^{-1/2}).$$

Hence, $\tilde{\phi}_T = \tilde{V}_T + o_p(T^{-1/2})$ holds, and thus $\hat{V}_T = \tilde{V}_T + o_p(T^{-1/2})$ holds. The first step is thus completed.

Step 2: We next focus on proving $\tilde{V}_T = \bar{V}_T + o_p(T^{-1/2})$. By definition of \tilde{V}_T and \bar{V}_T , we have

$$\begin{aligned} \sqrt{T}(\tilde{V}_T - \bar{V}_T) = & \underbrace{\frac{1}{\sqrt{T}} \sum_{t=1}^T \left[\frac{\mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\}}{1 - \kappa_t(\mathbf{x}_t)} - 1 \right] \left[\mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right]}_{\eta_5} \quad (\text{B.36}) \\ & + \underbrace{\frac{1}{\sqrt{T}} \sum_{t=1}^T \left[\frac{\mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\}}{1 - \kappa_t(\mathbf{x}_t)} - \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \right] \left[r_t - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right]}_{\eta_6}. \end{aligned}$$

We first show $\eta_5 = o_p(1)$. Since $\kappa_t(\mathbf{x}_t) \leq C_2 < 1$ for sure for some constant $0 < C_2 < 1$ (by definition of a valid bandit algorithm as results for Theorem 1), it suffices to show

$$\psi_5 = T^{-1/2} \sum_{t=1}^T |\mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \pi^*\}| = o_p(1),$$

which is the direct conclusion of Lemma B.3.

Next, we show $\eta_6 = o_p(1)$. Firstly we can express η_6 as

$$\begin{aligned} \eta_6 &= \frac{1}{\sqrt{T}} \sum_{t=1}^T \left[\frac{\mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\}}{1 - \kappa_t(\mathbf{x}_t)} - \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \right] \left[r_t - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right] \\ &= \frac{1}{\sqrt{T}} \sum_{t=1}^T \left[\frac{1}{1 - \kappa_t(\mathbf{x}_t)} - \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \right] \left[r_t - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right] \mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} \\ &\quad - \frac{1}{\sqrt{T}} \sum_{t=1}^T \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \left[r_t - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right] \mathbb{I}\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)\} \\ &= \underbrace{\frac{1}{\sqrt{T}} \sum_{t=1}^T \left[\frac{1}{1 - \kappa_t(\mathbf{x}_t)} - \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \right] e_t \mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\}}_{\psi_6} \\ &\quad - \underbrace{\frac{1}{\sqrt{T}} \sum_{t=1}^T \frac{1}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \left[\mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right] \mathbb{I}\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)\} \mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}_{\psi_7}. \end{aligned}$$

Note that

$$\begin{aligned} \Pr\{a_t = \pi^*(\mathbf{x}_t)\} &\geq \Pr\{a_t = \hat{\pi}_t(\mathbf{x}_t), \hat{\pi}_t(\mathbf{x}_t) = \pi^*(\mathbf{x}_t)\} \\ &\geq \Pr\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} + \Pr\{\hat{\pi}_t(\mathbf{x}_t) = \pi^*(\mathbf{x}_t)\} - 1 \\ &= \Pr\{\hat{\pi}_t(\mathbf{x}_t) = \pi^*(\mathbf{x}_t)\} - \kappa_t(\mathbf{x}_t), \end{aligned}$$

where $\kappa_t(\mathbf{x}_t) = o_p(1)$ by Theorem 1 and $\Pr\{\hat{\pi}_t(\mathbf{x}_t) = \pi^*(\mathbf{x}_t)\} \geq 1 - ct^{-\alpha\gamma}$ by Lemma B.2 as $t \rightarrow \infty$. Thus we have for large enough t ,

$$\Pr\{a_t = \pi^*(\mathbf{x}_t)\} > C_1, \tag{B.37}$$

for some constant $C_1 > 0$.

Then we focus on proving $\psi_6 = o_p(1)$ here. Define a class of function

$$\mathcal{F}_\kappa(\mathbf{x}, a, r) = \left\{ \left[\frac{1}{1 - \kappa(\mathbf{x})} - \frac{\mathbb{I}\{a = \pi^*(\mathbf{x})\}}{\Pr\{a = \pi^*(\mathbf{x})\}} \right] e_t \mathbb{I}\{a = \pi(\mathbf{x})\} : \kappa_t \in \Lambda, \pi(\cdot) \in \Pi \right\},$$

where Π and Λ are two classes of functions that maps context $\mathbf{x} \in \mathcal{X}$ to a probability. Define the supremum of the empirical process indexed by \mathcal{F}_κ as

$$\|\mathbb{G}_n\|_{\mathcal{F}} \equiv \sup_{\pi \in \Pi} \frac{1}{\sqrt{T}} \sum_{t=1}^T [\mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) - \mathbb{E}\{\mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) | \mathcal{H}_{t-1}\}]. \quad (\text{B.38})$$

Firstly we notice that

$$\mathbb{E}\{\mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) | \mathcal{H}_{t-1}\} = \mathbb{E}\left(\left[\frac{1}{1 - \kappa_t(\mathbf{x}_t)} - \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}}\right] e_t \mathbb{I}\{a_t = \pi_t(\mathbf{x}_t)\} \mid \mathcal{H}_{t-1}\right) = 0,$$

since $\mathbb{E}(e_t | \{a_t, \mathcal{H}_{t-1}\}) = 0$.

Secondly since $\kappa_t(\mathbf{x}_t) \leq C_2 < 1$ for sure for some constant $0 < C_2 < 1$ (by definition of a valid bandit algorithm an results for Theorem 1), by the Triangle inequality, we have

$$\left| \frac{1}{1 - \kappa_t(\mathbf{x}_t)} - \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \right| \leq \left| \frac{1}{1 - \kappa_t(\mathbf{x}_t)} \right| + \left| \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \right| \leq \frac{1}{1 - C_2} + \frac{1}{C_1} \triangleq C.$$

Therefore, we have

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E} \left\{ \left(\frac{1}{\sqrt{T}} \left[\frac{1}{1 - \kappa_t(\mathbf{x}_t)} - \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \right] e_t \mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} \right)^2 \mid \mathcal{H}_{t-1} \right\} \\ & \leq \sum_{t=1}^T \mathbb{E} \left\{ \left(\frac{1}{\sqrt{T}} \left[\frac{1}{1 - \kappa_t(\mathbf{x}_t)} - \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \right] e_t \mathbb{I}\{a_t = \hat{\pi}_t(\mathbf{x}_t)\} \right)^2 \mid \mathcal{H}_{t-1} \right\} \\ & \leq \sum_{t=1}^T \frac{C^2}{T} \mathbb{E} \{e_t^2 \mid \mathcal{H}_{t-1}\} \leq \sum_{t=1}^T \frac{C^2}{T} \max\{\sigma_0^2, \sigma_1^2\} = C^2 \max\{\sigma_0^2, \sigma_1^2\} < \infty. \end{aligned}$$

Therefore, we have

$$d_1(f) \equiv \|\mathbb{E}(\|\mathcal{F}_\kappa(\mathbf{x}_1, a_1, r_1)\| | \mathcal{H}_0)\|_\infty < \infty,$$

and

$$d_2(f) \equiv \left\| \mathbb{E}((\mathcal{F}_\kappa(\mathbf{x}_1, a_1, r_1))^2 \mid \mathcal{H}_0) \right\|_\infty^{1/2} < \infty.$$

It follows from the maximal inequality developed in Section 4.2 of Dedecker and Louhichi (2002) that there exist some constant $K \geq 1$ such that

$$\mathbb{E} \left[\|\mathbb{G}_n\|_{\mathcal{F}} \right] \lesssim K \left(\sqrt{p}d_2(f) + \frac{1}{\sqrt{T}} \left\| \max_{1 \leq t \leq T} |\mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) - \mathbb{E}(\mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) \mid \mathcal{H}_{t-1})| \right\| \right)$$

The above right-hand-side is upper bounded by

$$\mathcal{O}(1)\sqrt{T^{-1/2}},$$

where $\mathcal{O}(1)$ denotes some universal constant. Hence, we have

$$\mathbb{E} \left[\|\mathbb{G}_n\|_{\mathcal{F}} \right] = \mathcal{O}_p(T^{-1/2}). \quad (\text{B.39})$$

Combined with Equation (B.38), we have

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T \mathcal{F}_\kappa(\mathbf{x}_t, a_t, r_t) = \mathcal{O}_p(T^{-1/2}).$$

and $\psi_6 = o_p(1)$.

Next, for ψ_7 , by triangle inequality, we have

$$\begin{aligned} |\psi_7| &= \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T \frac{1}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \left[\mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right] \mathbb{I}\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)\} \mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\} \right| \\ &\leq \frac{1}{\sqrt{T}} \sum_{t=1}^T \left| \frac{1}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \left[\mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right] \mathbb{I}\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)\} \mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\} \right| \\ &\leq \frac{1}{\sqrt{T}} \sum_{t=1}^T \left| \frac{1}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \right| \left| \mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right| \mathbb{I}\{a_t \neq \hat{\pi}_t(\mathbf{x}_t)\} \mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\} \\ &\leq \frac{1}{\sqrt{T}} \sum_{t=1}^T \frac{1}{C_1} \left| \mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right|. \end{aligned}$$

Notice that $\psi_5 = \frac{1}{\sqrt{T}} \sum_{t=1}^T \left| \mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} \right| = o_p(1)$, we have

$$|\psi_7| \leq \frac{1}{C_1} \psi_5 = o_p(1).$$

Therefore, $\eta_6 = \psi_6 + \psi_7 = o_p(1)$. Hence, $\tilde{V}_T = \bar{V}_T + o_p(T^{-1/2})$ hold.

Step 3: Then, to show the asymptotic normality of the proposed value estimator under DREAM, based on the above two steps, it is sufficient to show

$$\sqrt{T}(\hat{V}_T - V^*) = \sqrt{T}(\bar{V}_T - V^*) + o_p(1) \xrightarrow{D} \mathcal{N}(0, \sigma_{DR}^2), \quad (\text{B.40})$$

as $T \rightarrow \infty$, using Martingale Central Limit Theorem. By $r_t = \mu\{\mathbf{x}_t, a_t\} + e_t$, we define

$$\begin{aligned} \xi_t &= \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \left[r_t - \mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} \right] + \mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} \\ &= \underbrace{\frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} e_t}_{Z_t} + \underbrace{\mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - V^*}_{W_t}. \end{aligned} \quad (\text{B.41})$$

By (B.41), we have

$$\begin{aligned} \mathbb{E}\{Z_t \mid \mathcal{H}_{t-1}\} &= \mathbb{E}\{\mathbb{E}[Z_t \mid a_t] \mid \mathcal{H}_{t-1}\} = \mathbb{E}\left\{ \mathbb{E}\left[\frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} e_t \mid a_t, \mathbf{x}_t \right] \mid \mathcal{H}_{t-1} \right\} \\ &= \mathbb{E}\left\{ \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \mathbb{E}[e_t \mid a_t, \mathbf{x}_t] \mid \mathcal{H}_{t-1} \right\}. \end{aligned}$$

Since e_t is independent of \mathcal{H}_{i-1} and \mathbf{x}_i given a_t , and $\mathbb{E}\{e_t \mid a_t\} = 0$, we have

$$\mathbb{E}\{Z_t \mid \mathcal{H}_{t-1}\} = \mathbb{E}\left\{ \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \mathbb{E}\{e_t \mid a_t\} \mid \mathcal{H}_{t-1} \right\} = 0.$$

Notice that $\mathbb{E}[\mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\}] = V^*$ by the definition, we have

$$\mathbb{E}\{W_t \mid \mathcal{H}_{t-1}\} = \mathbb{E}\{\mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - V^*\} = 0.$$

Thus, we have $\mathbb{E}\{\xi_t \mid \mathcal{H}_{t-1}\} = 0$, which implies that $\{Z_t\}_{t=1}^T$, $\{W_t\}_{t=1}^T$ and $\{\xi_t\}_{t=1}^T$ are Martingale difference sequences. To prove Equation (B.40), it suffices to prove that $(1/\sqrt{T}) \sum_{t=1}^T \xi_t \xrightarrow{D} \mathcal{N}(0, \sigma_{DR}^2)$, as $T \rightarrow \infty$, using Martingale Central Limit Theorem.

Firstly we calculate the conditional variance of ξ_t given \mathcal{H}_{t-1} . Note that

$$\mathbb{E}(Z_t^2 \mid \mathcal{H}_{t-1}) = \mathbb{E}\left\{ \left(\frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} e_t \right)^2 \mid \mathcal{H}_{t-1} \right\} = \mathbb{E}\left(\frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{[\Pr\{a_t = \pi^*(\mathbf{x}_t)\}]^2} e_t^2 \mid \mathcal{H}_{t-1} \right),$$

and

$$\begin{aligned}\mathbb{E}(Z_t^2 \mid \mathcal{H}_{t-1}) &= \mathbb{E} \left[\mathbb{E} \left(\frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{[\Pr\{a_t = \pi^*(\mathbf{x}_t)\}]^2} e_t^2 \mid a_t, \mathbf{x}_t \right) \mid \mathcal{H}_{t-1} \right] \\ &= \mathbb{E} \left[\frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{[\Pr\{a_t = \pi^*(\mathbf{x}_t)\}]^2} \mathbb{E}(e_t^2 \mid a_t, \mathbf{x}_t) \mid \mathcal{H}_{t-1} \right].\end{aligned}$$

Since e_t is independent of \mathcal{H}_{i-1} and \mathbf{x}_i given a_t , we have

$$\begin{aligned}\frac{1}{T} \sum_{t=1}^T \mathbb{E}(Z_t^2 \mid \mathcal{H}_{t-1}) &= \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{[\Pr\{a_t = \pi^*(\mathbf{x}_t)\}]^2} \mathbb{E}(e_t^2 \mid a_t) \mid \mathcal{H}_{t-1} \right] \\ &= \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left(\frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{[\Pr\{a_t = \pi^*(\mathbf{x}_t)\}]^2} \sigma_{a_t}^2 \mid \mathcal{H}_{t-1} \right). \\ &= \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\frac{1}{[\Pr\{a_t = \pi^*(\mathbf{x}_t)\}]^2} \sigma_{\pi^*(\mathbf{x}_t)}^2 \mathbb{E}(\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\} \mid \mathbf{x}_t, \mathcal{H}_{t-1}) \right].\end{aligned}$$

By the definition of Equation (B.25),

$$\nu_i(\mathbf{x}_i, \mathcal{H}_{i-1}) \equiv \Pr\{a_i \neq \pi^*(\mathbf{x}_i) \mid \mathbf{x}_i, \mathcal{H}_{i-1}\} = \mathbb{E}[\mathbb{I}\{a_i \neq \pi^*(\mathbf{x}_i)\} \mid \mathbf{x}_i, \mathcal{H}_{i-1}],$$

we have

$$\begin{aligned}\frac{1}{T} \sum_{t=1}^T \mathbb{E}(Z_t^2 \mid \mathcal{H}_{t-1}) &= \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\frac{1}{[\Pr\{a_t = \pi^*(\mathbf{x}_t)\}]^2} \sigma_{\pi^*(\mathbf{x}_t)}^2 \{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\} \right] \\ &= \frac{1}{T} \sum_{t=1}^T \int \frac{\{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \sigma_{\pi^*(\mathbf{x})}^2 dP_{\mathcal{X}} \\ &= \int \left[\frac{1}{T} \sum_{t=1}^T \frac{\{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \right] \sigma_{\pi^*(\mathbf{x})}^2 dP_{\mathcal{X}}.\end{aligned}$$

Similar as before, since $\lim_{i \rightarrow \infty} \Pr\{a_i \neq \pi^*(\mathbf{x})\} = \kappa_{\infty}(\mathbf{x})$, we have for any $\epsilon > 0$, there exist a constant $T_0 > 0$ such that $|\Pr\{a_i \neq \pi^*(\mathbf{x})\} - \kappa_{\infty}(\mathbf{x})| < \epsilon$ for all $i \geq T_0$.

We firstly consider the expectation of $(1/T) \sum_{t=1}^T \{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\} / [\Pr\{a_t = \pi^*(\mathbf{x})\}]^2$.

Note that $\Pr\{a_t = \pi^*(\mathbf{x})\}$ is not conditional on \mathcal{H}_{t-1} , thus we have

$$\begin{aligned}\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \frac{\{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \right] &= \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{E} \{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} = \frac{1}{T} \sum_{t=1}^T \frac{\Pr\{a_t = \pi^*(\mathbf{x})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \\ &= \frac{1}{T} \sum_{t=1}^T \frac{1}{\Pr\{a_t = \pi^*(\mathbf{x})\}}.\end{aligned}$$

Therefore by the triangle inequality we have

$$\begin{aligned}& \left| \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \frac{\{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \right] - \frac{1}{1 - \kappa_\infty(\mathbf{x})} \right| = \left| \frac{1}{T} \sum_{t=1}^T \frac{1}{\Pr\{a_t = \pi^*(\mathbf{x})\}} - \frac{1}{1 - \kappa_\infty(\mathbf{x})} \right| \\ &= \left| \frac{1}{T} \sum_{t=1}^T \left[\frac{1}{\Pr\{a_t = \pi^*(\mathbf{x})\}} - \frac{1}{1 - \kappa_\infty(\mathbf{x})} \right] \right| = \left| \frac{1}{T} \sum_{t=1}^T \frac{1 - \kappa_\infty(\mathbf{x}) - \Pr\{a_t = \pi^*(\mathbf{x})\}}{\Pr\{a_t = \pi^*(\mathbf{x})\} \{1 - \kappa_\infty(\mathbf{x})\}} \right| \\ &= \left| \frac{1}{T} \sum_{t=1}^T \frac{\Pr\{a_t \neq \pi^*(\mathbf{x}) - \kappa_\infty(\mathbf{x})\}}{\Pr\{a_t = \pi^*(\mathbf{x})\} \{1 - \kappa_\infty(\mathbf{x})\}} \right| \leq \frac{1}{T} \sum_{t=1}^T \frac{|\Pr\{a_t \neq \pi^*(\mathbf{x}) - \kappa_\infty(\mathbf{x})\}|}{\Pr\{a_t = \pi^*(\mathbf{x})\} \{1 - \kappa_\infty(\mathbf{x})\}} \\ &= \frac{1}{T} \sum_{t=1}^{T_0} \frac{|\Pr\{a_t \neq \pi^*(\mathbf{x}) - \kappa_\infty(\mathbf{x})\}|}{\Pr\{a_t = \pi^*(\mathbf{x})\} \{1 - \kappa_\infty(\mathbf{x})\}} + \frac{1}{T} \sum_{t=T_0}^T \frac{|\Pr\{a_t \neq \pi^*(\mathbf{x}) - \kappa_\infty(\mathbf{x})\}|}{\Pr\{a_t = \pi^*(\mathbf{x})\} \{1 - \kappa_\infty(\mathbf{x})\}} \\ &< \frac{1}{T} \sum_{t=1}^{T_0} \frac{|\Pr\{a_t \neq \pi^*(\mathbf{x})\}| + |\kappa_\infty(\mathbf{x})|}{\Pr\{a_t = \pi^*(\mathbf{x})\} \{1 - \kappa_\infty(\mathbf{x})\}} + \frac{1}{T} \sum_{t=T_0}^T \frac{\epsilon}{\Pr\{a_t = \pi^*(\mathbf{x})\} \{1 - \kappa_\infty(\mathbf{x})\}}.\end{aligned}$$

Since the above equation holds for any $\epsilon > 0$, we have

$$\begin{aligned}& \left| \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \frac{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \right] - \frac{1}{1 - \kappa_\infty(\mathbf{x})} \right| \leq \frac{1}{T} \sum_{t=1}^{T_0} \frac{|\Pr\{a_t = \pi^*(\mathbf{x})\}| + |\kappa_\infty(\mathbf{x})|}{\Pr\{a_t = \pi^*(\mathbf{x})\} \{1 - \kappa_\infty(\mathbf{x})\}} \\ &\leq \frac{1}{T} \sum_{t=1}^{T_0} \frac{2}{\Pr\{a_t = \pi^*(\mathbf{x})\} \{1 - \kappa_\infty(\mathbf{x})\}}.\end{aligned}$$

Since $\kappa_\infty(\mathbf{x}) \leq C_2 < 1$ for sure for some constant $0 < C_2 < 1$ (by definition of a valid bandit algorithm an results for Theorem 1), and by Equation (B.37), we have $\Pr\{a_t = \pi^*(\mathbf{x})\} > C_1$ for some constant $C_1 > 0$, therefore we have

$$\left| \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \frac{\{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \right] - \frac{1}{1 - \kappa_\infty(\mathbf{x})} \right| \leq \frac{1}{T} \sum_{t=1}^{T_0} \frac{2}{C_1 (1 - C_2)} = \frac{2T_0}{C_1 (1 - C_2) T} \rightarrow 0,$$

as $T \rightarrow \infty$.

Then we consider the variance of $(1/T) \sum_{t=1}^T \{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\} / [\Pr\{a_t = \pi^*(\mathbf{x})\}]^2$ over different histories. By Lemma B.1, we have

$$\text{Var} [\nu_t(\mathbf{x}, \mathcal{H}_{t-1})] \leq \Pr\{a_t = \pi^*(\mathbf{x})\} [1 - \Pr\{a_t = \pi^*(\mathbf{x})\}].$$

Therefore, we have

$$\begin{aligned} \text{Var} \left[\frac{1}{T} \sum_{t=1}^T \frac{\{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \right] &= \frac{1}{T^2} \text{Var} \left[\sum_{t=1}^T \frac{\{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \right] \\ &\leq \frac{1}{T} \sum_{t=1}^T \text{Var} \left[\frac{\{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \right] = \frac{1}{T} \sum_{t=1}^T \frac{\text{Var} [\{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}]}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^4} = \frac{1}{T} \sum_{t=1}^T \frac{\text{Var} [\{\nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}]}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^4} \\ &\leq \frac{1}{T} \sum_{t=1}^T \frac{\Pr\{a_t = \pi^*(\mathbf{x})\} [1 - \Pr\{a_t = \pi^*(\mathbf{x})\}]}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^4} = \frac{1}{T} \sum_{t=1}^T \frac{1 - \Pr\{a_t = \pi^*(\mathbf{x})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^3} \\ &\leq \frac{1}{T} \sum_{t=1}^T \frac{1 - \Pr\{a_t = \pi^*(\mathbf{x})\}}{[1 - C_2]^3} = \frac{1}{[1 - C_2]^3} \frac{1}{T} \sum_{t=1}^T \Pr\{a_t \neq \pi^*(\mathbf{x})\}. \end{aligned}$$

Similarly as before, we could proof

$$\frac{1}{T} \sum_{t=1}^T \Pr\{a_t \neq \pi^*(\mathbf{x})\} \rightarrow \kappa_\infty(\mathbf{x}),$$

which follows

$$\text{Var} \left[\frac{1}{T} \sum_{t=1}^T \frac{\{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \right] \rightarrow \frac{1}{[1 - C_2]^3} \kappa_\infty(\mathbf{x}).$$

Therefore, as T goes to infinity, we have

$$\frac{1}{T} \sum_{t=1}^T \frac{\{1 - \nu_t(\mathbf{x}, \mathcal{H}_{t-1})\}}{[\Pr\{a_t = \pi^*(\mathbf{x})\}]^2} \rightarrow \frac{1}{1 - \kappa_\infty(\mathbf{x})},$$

and

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}(Z_t^2 \mid \mathcal{H}_{t-1}) \rightarrow \int \frac{1}{1 - \kappa_\infty(\mathbf{x})} \sigma_{\pi^*(\mathbf{x})}^2 dP_{\mathcal{X}}.$$

Using the same technique of conditioning on a_t and \mathbf{x}_t , we have

$$\begin{aligned}\mathbb{E}(W_t Z_t \mid \mathcal{H}_{t-1}) &= \mathbb{E} \left\{ (\mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - V^*) \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} e_t \mid \mathcal{H}_{t-1} \right\} \\ &= \mathbb{E} \left\{ \mathbb{E} \left[(\mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - V^*) \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} e_t \mid a_t, \mathbf{x}_t \right] \mid \mathcal{H}_{t-1} \right\} \\ &= \mathbb{E} \left\{ (\mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} - V^*) \frac{\mathbb{I}\{a_t = \pi^*(\mathbf{x}_t)\}}{\Pr\{a_t = \pi^*(\mathbf{x}_t)\}} \mathbb{E}(e_t \mid a_t) \mid \mathcal{H}_{t-1} \right\} = 0.\end{aligned}$$

Thus, we further have

$$\mathbb{E}(\xi_t^2 \mid \mathcal{H}_{t-1}) = \mathbb{E}\{(Z_t + W_t)^2 \mid \mathcal{H}_{t-1}\} = \mathbb{E}(Z_t^2 \mid \mathcal{H}_{t-1}) + \mathbb{E}(W_t^2 \mid \mathcal{H}_{t-1}),$$

and

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}(\xi_t^2 \mid \mathcal{H}_{t-1}) = \int \frac{1}{1 - \kappa_\infty(\mathbf{x})} \sigma_{\pi^*(\mathbf{x})}^2 dP_{\mathcal{X}} + \frac{1}{T} \sum_{t=1}^T \text{Var}[\mu\{\mathbf{x}_t, \pi^*(\mathbf{x}_t)\} \mid \mathcal{H}_{t-1}].$$

Therefore as T goes to infinity, we have

$$\sum_{t=1}^T \mathbb{E} \left\{ \left(\frac{1}{\sqrt{T}} \xi_t \right)^2 \mid \mathcal{H}_{t-1} \right\} \longrightarrow \sigma_{DR}^2,$$

where

$$\sigma_{DR}^2 = \int_{\mathbf{x}} \frac{\sigma_1^2 \mathbb{I}\{\mu(\mathbf{x}, 1) > \mu(\mathbf{x}, 0)\} + \sigma_0^2 \mathbb{I}\{\mu(\mathbf{x}, 1) < \mu(\mathbf{x}, 0)\}}{1 - \kappa_\infty(\mathbf{x})} dP_{\mathcal{X}} + \text{Var}[\mu\{\mathbf{x}, \pi^*(\mathbf{x})\}]. \quad (\text{B.42})$$

Then we check the conditional Lindeberg condition. For any $h > 0$, we have

$$\sum_{t=1}^T \mathbb{E} \left\{ \left(\frac{1}{\sqrt{T}} \xi_t \right)^2 \mathbb{I} \left\{ \left| \frac{1}{\sqrt{T}} \xi_t \right| > h \right\} \mid \mathcal{H}_{t-1} \right\} = \frac{1}{T} \sum_{t=1}^T \mathbb{E} \{ \xi_t^2 \mathbb{I} \{ \xi_t^2 > Th^2 \} \mid \mathcal{H}_{t-1} \}.$$

Since $\xi_t^2 \mathbb{I} \{ \xi_t^2 > Th^2 \}$ converges to zero as T goes to infinity and is dominated by ξ_t^2 given \mathcal{H}_{t-1} . Therefore, by Dominated Convergence Theorem, we conclude that

$$\sum_{t=1}^T \mathbb{E} \left\{ \left(\frac{1}{\sqrt{T}} \xi_t \right)^2 \mathbb{I} \left\{ \left| \frac{1}{\sqrt{T}} \xi_t \right| > h \right\} \mid \mathcal{H}_{t-1} \right\} \rightarrow 0, \text{ as } t \rightarrow \infty.$$

Thus the conditional Lindeberg condition is checked.

Next, recall the derived conditional variance in (B.42). By Martingale Central Limit Theorem, we have

$$(1/\sqrt{T}) \sum_{t=1}^T \xi_t \xrightarrow{D} \mathcal{N}(0, \sigma_{DR}^2),$$

as $T \rightarrow \infty$. Hence, we complete the proof of Equation (B.40).

Step 4: Finally, to show the variance estimator in Equation (5) in the main paper is a consistent estimator of σ_{DR}^2 . Recall that the variance estimator is

$$\begin{aligned} \hat{\sigma}_T^2 = & \underbrace{\frac{1}{T} \sum_{t=1}^T \frac{\hat{\sigma}_{1,t-1}^2(\mathbf{x}_t, 1) \mathbb{I}\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) > \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\} + \hat{\sigma}_{0,t-1}^2 \mathbb{I}\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) < \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\}}{1 - \hat{\kappa}_t(\mathbf{x}_t)}}_{\hat{\sigma}_{T,1}^2} \\ & + \underbrace{\frac{1}{T} \sum_{t=1}^T \left[\hat{\mu}_T\{\mathbf{x}_t, \hat{\pi}_T(\mathbf{x}_t)\} - \frac{1}{T} \sum_{t=1}^T \hat{\mu}_T\{\mathbf{x}_t, \hat{\pi}_T(\mathbf{x}_t)\} \right]^2}_{\hat{\sigma}_{T,2}^2}. \end{aligned} \tag{B.43}$$

Firstly we proof the first line of the above Equation (B.43) is a consistent estimator for

$$\int_{\mathbf{x}} \frac{\sigma_1^2 \mathbb{I}\{\mu(\mathbf{x}, 1) > \mu(\mathbf{x}, 0)\} + \sigma_0^2 \mathbb{I}\{\mu(\mathbf{x}, 1) < \mu(\mathbf{x}, 0)\}}{1 - \kappa_\infty(\mathbf{x})} dP_{\mathcal{X}}.$$

Recall that we denote $\hat{\Delta}_{\mathbf{x}_t} = \hat{\mu}_{t-1}(\mathbf{x}_t, 1) - \hat{\mu}_{t-1}(\mathbf{x}_t, 0)$, thus we can rewrite $\hat{\sigma}_T^2$ as

$$\begin{aligned} \hat{\sigma}_{T,1}^2 &= \frac{1}{T} \sum_{t=1}^T \frac{\hat{\sigma}_{1,t-1}^2 \mathbb{I}\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) > \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\} + \hat{\sigma}_{0,t-1}^2 \mathbb{I}\{\hat{\mu}_{t-1}(\mathbf{x}_t, 1) < \hat{\mu}_{t-1}(\mathbf{x}_t, 0)\}}{1 - \hat{\kappa}_t(\mathbf{x}_t)} \\ &= \frac{1}{T} \sum_{t=1}^T \frac{\hat{\sigma}_{1,t-1}^2 \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} + \hat{\sigma}_{0,t-1}^2 \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} < 0\}}{1 - \hat{\kappa}_t(\mathbf{x}_t)}. \end{aligned}$$

We decompose the proposed variance estimator by

$$\begin{aligned}
\hat{\sigma}_{T,1}^2 &= \frac{1}{T} \sum_{t=1}^T \frac{\{\hat{\sigma}_{1,t-1}^2 - \sigma_1^2\} \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} + \{\hat{\sigma}_{0,t-1}^2 - \sigma_0^2\} \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} < 0\}}{1 - \hat{\kappa}_t(\mathbf{x}_t)} \\
&+ \frac{1}{T} \sum_{t=1}^T \frac{\sigma_1^2 \left(\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} \right) + \sigma_0^2 \left(\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} < 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\} \right)}{1 - \hat{\kappa}_t(\mathbf{x}_t)} \\
&+ \frac{1}{T} \sum_{t=1}^T (\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}) \left\{ \frac{1}{1 - \hat{\kappa}_t(\mathbf{x}_t)} - \frac{1}{1 - \kappa_t(\mathbf{x}_t)} \right\} \\
&+ \frac{1}{T} \sum_{t=1}^T \frac{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}}{1 - \kappa_t(\mathbf{x}_t)}.
\end{aligned}$$

Our goal is to prove that the first three lines are all $o_p(1)$.

Firstly, recall that

$$\begin{aligned}
\hat{\sigma}_{a,t}^2 &= \left\{ \sum_{i=1}^t \mathbb{I}(a_i = a) - d \right\}^{-1} \sum_{\substack{1 \leq i \leq t \\ a_i = a}} [\hat{\mu}_i\{\mathbf{x}_i, a_i\} - r_i]^2 \\
&= \left\{ \sum_{i=1}^t \mathbb{I}(a_i = a) - d \right\}^{-1} \sum_{\substack{1 \leq i \leq t \\ a_i = a}} [\mathbf{x}_i^\top \left\{ \hat{\boldsymbol{\beta}}_{i-1}(a) - \boldsymbol{\beta}_{i-1}(a) \right\} - e_i]^2.
\end{aligned}$$

By Lemma 4.1, we have $\|\hat{\boldsymbol{\beta}}_{i-1}(a) - \boldsymbol{\beta}_{i-1}(a)\|_1 = o_p(1)$. Under Assumption 4.1, we have $\mathbf{x}_i^\top \left\{ \hat{\boldsymbol{\beta}}_{i-1}(a) - \boldsymbol{\beta}_{i-1}(a) \right\} = o_p(1)$. Thus by Lemma 6 in Luedtke and Van Der Laan (2016), we have

$$\hat{\sigma}_{a,t}^2 = \left\{ \sum_{i=1}^t \mathbb{I}(a_i = a) - d \right\}^{-1} \sum_{\substack{1 \leq i \leq t \\ a_i = a}} e_i^2 + o_p(1).$$

Since e_i is i.i.d conditional on a_i , and $\mathbb{E}(e_i^2 | a_i = a) = \sigma_a^2$, noting that

$$\lim_{t \rightarrow \infty} \frac{\sum_{i=1}^t \mathbb{I}(a_i = a)}{\sum_{i=1}^t \mathbb{I}(a_i = a) - d} = 1,$$

by Law of Large Numbers we have

$$\hat{\sigma}_{a,t}^2 = \sigma_a^2 + o_p(1).$$

Therefore, the first line is $o_p(1)$.

Secondly, denote the second line as

$$\begin{aligned}\psi_8 &= \frac{1}{T} \sum_{t=1}^T \frac{\sigma_1^2 \left(\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} \right) + \sigma_0^2 \left(\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} < 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\} \right)}{1 - \hat{\kappa}_t(\mathbf{x}_t)} \\ &= \frac{1}{T} \sum_{t=1}^T \frac{\sigma_1^2 \left(\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} \right) + \sigma_0^2 \left(\mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} \right)}{1 - \hat{\kappa}_t(\mathbf{x}_t)} \\ &= \frac{\sigma_1^2 - \sigma_0^2}{T} \sum_{t=1}^T \frac{\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\}}{1 - \hat{\kappa}_t(\mathbf{x}_t)}.\end{aligned}$$

Since $\kappa_t(\mathbf{x}_t) \leq C_2 < 1$ for sure for some constant $0 < C_2 < 1$ (by definition of a valid bandit algorithm and results for Theorem 1), by the triangle inequality, we have

$$\begin{aligned}|\psi_8| &\leq \frac{\sigma_1^2 - \sigma_0^2}{T} \sum_{t=1}^T \frac{\left| \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} \right|}{|1 - \hat{\kappa}_t(\mathbf{x}_t)|} \leq \frac{\sigma_1^2 - \sigma_0^2}{T} \sum_{t=1}^T \frac{\left| \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} \right|}{1 - C_2} \\ &\leq \frac{\sigma_1^2 - \sigma_0^2}{1 - C_2} \frac{1}{T} \sum_{t=1}^T \left| \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} \right|.\end{aligned}$$

Since $\Pr \left(\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} = 0 \right) = 1 - ct^{-\alpha\gamma}$ by Lemma B.2, there exists some constant c such that

$$\Pr \left(\frac{1}{T} \sum_{t=1}^T \left| \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} \right| \neq 0 \right) \leq \sum_{t=1}^T \Pr \left(\left| \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} \right| \neq 0 \right) \leq \sum_{t=1}^T ct^{-\alpha\gamma}.$$

By Lemma 6 in Luedtke and Van Der Laan (2016), we have $\sum_{t=1}^T ct^{-\alpha\gamma} = cT^{-\alpha\gamma}$, thus

$$\Pr(|\psi_8| \neq 0) = \Pr \left(\frac{1}{T} \sum_{t=1}^T \left| \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} \right| \neq 0 \right) \leq cT^{-\alpha\gamma},$$

which follows $\Pr(|\psi_8| \neq 0) = o_p(1)$.

Lastly, under the assumption that $\hat{\kappa}_t(\mathbf{x}_t)$ is a consistent estimator for $\kappa_t(\mathbf{x}_t)$, we have the third line is $o_p(1)$ by the continuous mapping theorem.

Given the above results, we have

$$\hat{\sigma}_{T,1}^2 = \frac{1}{T} \sum_{t=1}^T \frac{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}}{1 - \kappa_t(\mathbf{x}_t)} + o_p(1) = \frac{1}{T} \sum_{t=1}^T \frac{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}}{1 - \Pr\{a_t \neq \hat{\pi}_t(x_t)\}} + o_p(1).$$

Thus, we can further express $\hat{\sigma}_{T,1}^2$ as

$$\begin{aligned} \hat{\sigma}_{T,1}^2 &= \frac{1}{T} \sum_{t=1}^T (\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}) \left(\frac{1}{1 - \Pr\{a_t \neq \hat{\pi}_t(x)\}} - \frac{1}{1 - \Pr\{a_t \neq \pi^*(\mathbf{x})\}} \right) \\ &\quad + \frac{1}{T} \sum_{t=1}^T \frac{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}}{1 - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}} + o_p(1). \end{aligned} \tag{B.44}$$

Notice that

$$\left(\frac{1}{1 - \Pr\{a_t \neq \hat{\pi}_t(x)\}} - \frac{1}{1 - \Pr\{a_t \neq \pi^*(\mathbf{x})\}} \right) = \frac{\Pr\{a_t \neq \hat{\pi}_t(x)\} - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}}{(1 - \Pr\{a_t \neq \hat{\pi}_t(x)\})(1 - \Pr\{a_t \neq \pi^*(\mathbf{x})\})},$$

where

$$\Pr\{a_t \neq \hat{\pi}_t(x)\} - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\} = \mathbb{E}(\mathbb{I}\{a_t \neq \hat{\pi}_t(x)\} - \mathbb{I}\{a_t \neq \pi^*(\mathbf{x}_t)\}) = \mathbb{E}(\mathbb{I}\{\hat{\pi}_t(x) \neq \pi^*(\mathbf{x}_t)\}).$$

We also note that by Lemma B.2, there exists some constant c and $0 < \alpha < \frac{1}{2}$ such that $\alpha\gamma < \frac{1}{2}$ and

$$\mathbb{E}(\mathbb{I}\{\hat{\pi}_t(x) = \pi^*(\mathbf{x}_t)\}) = \Pr(\hat{\pi}_t(x) = \pi^*(\mathbf{x}_t)) \leq ct^{-\alpha\gamma},$$

therefore we have the result that

$$\Pr\{a_t \neq \hat{\pi}_t(x)\} - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\} = \mathbb{E}(\mathbb{I}\{\hat{\pi}_t(x) \neq \pi^*(\mathbf{x}_t)\}) = o_p(1).$$

Thus, Equation (B.44) can be expressed as

$$\begin{aligned}
\hat{\sigma}_{T,1}^2 &= \frac{1}{T} \sum_{t=1}^T (\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}) o_p(1) + \frac{1}{T} \sum_{t=1}^T \frac{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}}{1 - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}} + o_p(1) \\
&= \frac{1}{T} \sum_{t=1}^T \frac{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}}{1 - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}} + o_p(1) \\
&= \underbrace{\frac{1}{T} \sum_{t=1}^T \{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}\} \left[\frac{1}{1 - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}} - \frac{1}{1 - \kappa_\infty(\mathbf{x}_t)} \right]}_{\psi_9} \\
&\quad + \frac{1}{T} \sum_{t=1}^T \frac{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}}{1 - \kappa_\infty(\mathbf{x}_t)} + o_p(1).
\end{aligned}$$

Note that

$$\begin{aligned}
\psi_9 &= \frac{1}{T} \sum_{t=1}^T \{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}\} \left[\frac{1}{1 - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}} - \frac{1}{1 - \kappa_\infty(\mathbf{x}_t)} \right] \\
&= \frac{1}{T} \sum_{t=1}^T \{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}\} \frac{\kappa_\infty(\mathbf{x}_t) - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}}{[1 - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}][1 - \kappa_\infty(\mathbf{x}_t)]},
\end{aligned}$$

which follows

$$|\psi_9| \leq \frac{1}{T} \sum_{t=1}^T \{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}\} \frac{|\kappa_\infty(\mathbf{x}_t) - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}|}{[1 - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}][1 - \kappa_\infty(\mathbf{x}_t)]}.$$

By Equation (B.37), for large enough t , there exists some constant $C_1 > 0$ such that

$$\Pr\{a_t \neq \pi^*(\mathbf{x}_t)\} < C_1.$$

Since $\kappa_t(\mathbf{x}_t) \leq C_2 < 1$ for sure for some constant $0 < C_2 < 1$ (by the definition of a valid bandit algorithm as results shown in Theorem 1), we also have

$$\kappa_\infty(\mathbf{x}_t) < C_2.$$

Therefore,

$$\frac{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}}{[1 - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}][1 - \kappa_\infty(\mathbf{x}_t)]} < \frac{\max\{\sigma_0^2, \sigma_1^2\}}{(1 - C_1)(1 - C_2)} \triangleq C,$$

which follows immediately that

$$|\psi_9| \leq \frac{1}{T} \sum_{t=1}^T C |\kappa_\infty(\mathbf{x}_t) - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}|.$$

Since $\lim_{t \rightarrow \infty} \Pr\{a_t \neq \pi^*(\mathbf{x})\} = \kappa_\infty(\mathbf{x})$ for any \mathbf{x} , therefore for any $\epsilon > 0$, there exists some constant T_0 , such that $|\Pr\{a_t \neq \pi^*(\mathbf{x})\} - \kappa_\infty(\mathbf{x})| < \epsilon$ for any \mathbf{x} with $t \geq T_0$, thus we have

$$\begin{aligned} |\psi_9| &\leq \frac{1}{T} \sum_{t=1}^T C |\kappa_\infty(\mathbf{x}_t) - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}| \\ &= \frac{1}{T} \sum_{t=1}^{T_0} C |\kappa_\infty(\mathbf{x}_t) - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}| + \frac{1}{T} \sum_{t=T_0}^T C |\kappa_\infty(\mathbf{x}_t) - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}| \\ &< \frac{1}{T} \sum_{t=1}^{T_0} C |\kappa_\infty(\mathbf{x}_t) - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}| + \frac{1}{T} \sum_{t=T_0}^T C \epsilon. \end{aligned}$$

Note that by the triangle inequality,

$$|\kappa_\infty(\mathbf{x}_t) - \Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}| \leq |\kappa_\infty(\mathbf{x}_t)| + |\Pr\{a_t \neq \pi^*(\mathbf{x}_t)\}| < C_1 + C_2,$$

thus we have

$$|\psi_9| < \frac{1}{T} \sum_{t=1}^{T_0} C (C_1 + C_2) + \frac{1}{T} \sum_{t=T_0}^T C \epsilon = \frac{T_0 C (C_1 + C_2)}{T} + \frac{T - T_0}{T} C \epsilon.$$

Since the above equation holds for any $\epsilon > 0$, we have

$$|\psi_9| \leq \frac{1}{T} \sum_{t=1}^{T_0} C (C_1 + C_2) + \frac{1}{T} \sum_{t=T_0}^T C \epsilon = \frac{T_0 C (C_1 + C_2)}{T} = \mathcal{O}(\frac{1}{T}),$$

which follows $\psi_9 = 0_p(1)$. Therefore, we have

$$\hat{\sigma}_{T,1}^2 = \frac{1}{T} \sum_{t=1}^T \frac{\sigma_1^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} + \sigma_0^2 \mathbb{I}\{\Delta_{\mathbf{x}_t} < 0\}}{1 - \kappa_\infty(\mathbf{x}_t)} + o_p(1).$$

By Law of Large Numbers, we further have

$$\hat{\sigma}_{T,1}^2 = \int_{\mathbf{x}} \frac{\sigma_1^2 \mathbb{I}\{\mu(\mathbf{x}, 1) > \mu(\mathbf{x}, 0)\} + \sigma_0^2 \mathbb{I}\{\mu(\mathbf{x}, 1) < \mu(\mathbf{x}, 0)\}}{1 - \kappa_{\infty}(\mathbf{x})} dP_{\mathcal{X}} + o_p(1).$$

Next, we proof the second line of Equation (B.43) is a consistent estimator for $\text{Var} [\mu\{\mathbf{x}, \pi^*(\mathbf{x})\}]$. By Central Limit Theorem and Continuous Mapping Theorem, we have

$$\hat{\sigma}_{T,2}^2 = \frac{1}{T} \sum_{t=1}^T \left[\hat{\mu}_T\{\mathbf{x}_t, \hat{\pi}_T(\mathbf{x}_t)\} - \frac{1}{T} \sum_{t=1}^T \hat{\mu}_T\{\mathbf{x}_t, \hat{\pi}_T(\mathbf{x}_t)\} \right]^2.$$

Since \mathbf{x}_t are i.i.d, $\hat{\mu}_T\{\mathbf{x}_t, \hat{\pi}_T(\mathbf{x}_t)\}$ are i.i.d as well. Thus by the Law of Large Numbers, we have

$$\frac{1}{T} \sum_{t=1}^T \hat{\mu}_T\{\mathbf{x}_t, \hat{\pi}_T(\mathbf{x}_t)\} = \mathbb{E} [\hat{\mu}_T\{\mathbf{x}, \hat{\pi}_T(\mathbf{x})\}] + o_p(1),$$

and

$$\hat{\sigma}_{T,2}^2 = \text{Var} [\hat{\mu}_T\{\mathbf{x}, \hat{\pi}_T(\mathbf{x})\}] + o_p(1).$$

Note that

$$\hat{\mu}_T\{\mathbf{x}, \hat{\pi}_T(\mathbf{x})\} = \mathbb{I}\{\hat{\Delta}_{\mathbf{x}} > 0\} \hat{\mu}_T\{\mathbf{x}, 1\} + [1 - \mathbb{I}\{\hat{\Delta}_{\mathbf{x}} > 0\}] \hat{\mu}_T\{\mathbf{x}, 0\},$$

since $\mathbb{I}\{\hat{\Delta}_{\mathbf{x}} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}} > 0\} = o_p(1)$ and the fact that $\hat{\mu}_T\{\mathbf{x}, 0\}$ and $\hat{\mu}_T\{\mathbf{x}, 1\}$ are consistent, we have

$$\hat{\mu}_T\{\mathbf{x}, \hat{\pi}_T(\mathbf{x})\} = \mathbb{I}\{\Delta_{\mathbf{x}} > 0\} \mu\{\mathbf{x}, 1\} + [1 - \mathbb{I}\{\Delta_{\mathbf{x}} > 0\}] \mu\{\mathbf{x}, 0\} = \mu\{\mathbf{x}, \pi(\mathbf{x})\} + o_p(1).$$

Therefore

$$\hat{\sigma}_{T,2}^2 = \text{Var} [\mu\{\mathbf{x}, \pi(\mathbf{x})\} + o_p(1)] = \text{Var} [\mu\{\mathbf{x}, \pi(\mathbf{x})\}] + o_p(1).$$

by Continuous Mapping Theorem. The proof for Theorem 3 is thus completed.

B.6 Results and Proof for Auxiliary Lemmas

Lemma B.1 *Suppose a random variable X is restricted to $[a, b]$ and $\mu = E[X]$, then the variance of X is bounded by $(b - \mu)(\mu - a)$.*

Proof: Firstly consider the case that $a = 0, b = 1$. Notice that we have $E[X^2] \leq E[X]$ since for all $x \in [0, 1], x^2 \leq x$. Therefore,

$$\text{Var}[X] = E[X^2] - (E[X])^2 = E[X^2] - \mu^2 \leq \mu - \mu^2 = \mu(1 - \mu).$$

Then we consider general interval $[a, b]$. Define $Y = \frac{X-a}{b-a}$, which is restricted in $[0, 1]$. Equivalently, $X = (b-a)Y + a$, which follows immediate that

$$\text{Var}[X] = (b-a)^2 \text{Var}[Y] \leq (b-a)^2 \mu_Y (1 - \mu_Y),$$

where the inequality is based on the first result. Now, by substituting $\mu_Y = \frac{\mu_X - a}{b-a}$, the bound equals

$$(b-a)^2 \frac{\mu_X - a}{b-a} \left(1 - \frac{\mu_X - a}{b-a}\right) = (b-a)^2 \frac{\mu_X - a}{b-a} \frac{b - \mu_X}{b-a} = (\mu_X - a)(b - \mu_X),$$

which is the desired result.

Lemma B.2 Suppose the conditions in Theorem 2 hold with Assumption 4.3, then there exists some constant c and $0 < \alpha < \frac{1}{2}$ such that $\alpha\gamma < \frac{1}{2}$ and $\Pr(\hat{\pi}_t(x) \neq \pi^*\{\mathbf{x}_t\}) \geq 1 - ct^{-\alpha\gamma}$.

Proof: By Theorem 2, we have $\hat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t} = \mathcal{O}_p(t^{-\frac{1}{2}})$, thus $\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} = \mathbb{I}\{\Delta_{\mathbf{x}_t} + \mathcal{O}_p(t^{-\frac{1}{2}}) > 0\}$.

By Assumption 4.3, there exists some constant c and $0 < \alpha < \frac{1}{2}$ such that $\alpha\gamma < \frac{1}{2}$ and

$$\Pr\{0 < |\Delta_{\mathbf{x}_t}| < t^{-\alpha}\} \leq ct^{-\alpha\gamma}.$$

Thus, with probability greater than $1 - ct^{-\alpha\gamma}$, we have $|\Delta_{\mathbf{x}_t}| > t^{-\alpha}$, which further implies $\mathbb{I}\{\Delta_{\mathbf{x}_t} + \mathcal{O}_p(t^{-\frac{1}{2}}) > 0\} = \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\}$. In other words, $\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} = 0$ with probability greater than $1 - ct^{-\alpha\gamma}$, which converges to 1 as $t \rightarrow \infty$. Therefore, as $t \rightarrow \infty$, we have

$$\Pr\left(\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} > 0\} - \mathbb{I}\{\Delta_{\mathbf{x}_t} > 0\} = 0\right) \geq 1 - ct^{-\alpha\gamma},$$

i.e.,

$$\Pr(\hat{\pi}_t(x) = \pi^*\{\mathbf{x}_t\}) \geq 1 - ct^{-\alpha\gamma}.$$

Lemma B.3 Suppose conditions in Lemma 4.1 hold. Assuming Assumptions 4.3 with $tp_t^2 \rightarrow \infty$ as $t \rightarrow \infty$, we have

$$T^{-1/2} \sum_{t=1}^T |\mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \pi^*\}| = o_p(1).$$

Proof: Without loss of generality, suppose $\hat{\pi}_t(\mathbf{x}_t) = \mathbb{I}\{\mathbf{x}_t^\top \boldsymbol{\beta}(1) > \mathbf{x}_t^\top \boldsymbol{\beta}(0)\}$. Since $\mu(\mathbf{x}_t, a) = \mathbf{x}_t^\top \boldsymbol{\beta}(a)$, we have

$$\begin{aligned} \mu(\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)) &= \mathbb{I}\{\mathbf{x}_t^\top \hat{\boldsymbol{\beta}}_t(1) > \mathbf{x}_t^\top \hat{\boldsymbol{\beta}}_t(0)\} \mathbf{x}_t^\top \boldsymbol{\beta}(1) + \mathbb{I}\{\mathbf{x}_t^\top \hat{\boldsymbol{\beta}}_t(1) \leq \mathbf{x}_t^\top \hat{\boldsymbol{\beta}}_t(0)\} \mathbf{x}_t^\top \boldsymbol{\beta}(0) \\ &= \mathbb{I}\{\mathbf{x}_t^\top \hat{\boldsymbol{\beta}}_t(1) > \mathbf{x}_t^\top \hat{\boldsymbol{\beta}}_t(0)\} \mathbf{x}_t^\top \{\boldsymbol{\beta}(1) - \boldsymbol{\beta}(0)\} + \mathbf{x}_t^\top \boldsymbol{\beta}(0). \end{aligned} \quad (\text{B.45})$$

Similarly to (B.45), we have

$$\mu(\mathbf{x}_t, \pi^*) = \mathbb{I}\{\mathbf{x}_t^\top \boldsymbol{\beta}(1) > \mathbf{x}_t^\top \boldsymbol{\beta}(0)\} \mathbf{x}_t^\top \{\boldsymbol{\beta}(1) - \boldsymbol{\beta}(0)\} + \mathbf{x}_t^\top \boldsymbol{\beta}(0). \quad (\text{B.46})$$

Combining (B.45) and (B.46), we have

$$\mu(\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)) - \mu(\mathbf{x}_t, \pi^*) = \left[\mathbb{I}\{\mathbf{x}_t^\top \hat{\boldsymbol{\beta}}_t(1) > \mathbf{x}_t^\top \hat{\boldsymbol{\beta}}_t(0)\} - \mathbb{I}\{\mathbf{x}_t^\top \boldsymbol{\beta}(1) > \mathbf{x}_t^\top \boldsymbol{\beta}(0)\} \right] \mathbf{x}_t^\top \{\boldsymbol{\beta}(1) - \boldsymbol{\beta}(0)\}. \quad (\text{B.47})$$

Since $\mathbb{I}\{\mathbf{x}_t^\top \boldsymbol{\beta}(1) > \mathbf{x}_t^\top \boldsymbol{\beta}(0)\} = 1$ by assumption, (B.47) can be simplified as

$$\mu(\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)) - \mu(\mathbf{x}_t, \pi^*) = -\mathbb{I}\{\mathbf{x}_t^\top \hat{\boldsymbol{\beta}}_t(1) - \mathbf{x}_t^\top \hat{\boldsymbol{\beta}}_t(0) \leq 0\} \mathbf{x}_t^\top \{\boldsymbol{\beta}(1) - \boldsymbol{\beta}(0)\} = -\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} \leq 0\} \Delta_{\mathbf{x}_t} \leq 0,$$

where $\Delta_{\mathbf{x}_t} = \mathbf{x}_t^\top \{\boldsymbol{\beta}(1) - \boldsymbol{\beta}(0)\}$ and $\hat{\Delta}_{\mathbf{x}_t} = \mathbf{x}_t^\top \{\hat{\boldsymbol{\beta}}_t(1) - \hat{\boldsymbol{\beta}}_t(0)\}$. Thus, we have

$$T^{-1/2} \sum_{t=1}^T |\mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \pi^*\}| = \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} \leq 0\} \Delta_{\mathbf{x}_t}.$$

To show $T^{-1/2} \sum_{t=1}^T |\mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \pi^*\}| = o_p(1)$, it suffices to show that $(1/\sqrt{T}) \sum_{t=1}^T \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} \leq 0\} \Delta_{\mathbf{x}_t}$ has an upper bound. Since $\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} \leq 0\} \hat{\Delta}_{\mathbf{x}_t} \leq 0$, it suffices to show

$$\zeta = \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} \leq 0\} (\hat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t})$$

has a lower bound. We further notice that for any $\alpha > 0$,

$$\begin{aligned} \zeta = & \underbrace{\Pr\{0 < \Delta_{\mathbf{x}_t} < T^{-\alpha}\} \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbb{I}\{0 < \Delta_{\mathbf{x}_t} < T^{-\alpha}\} \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} \leq 0\} (\hat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t})}_{\zeta_1} \\ & + \underbrace{\Pr\{\Delta_{\mathbf{x}_t} > T^{-\alpha}\} \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbb{I}\{\Delta_{\mathbf{x}_t} > T^{-\alpha}\} \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} \leq 0\} (\hat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t})}_{\zeta_2}. \end{aligned} \quad (\text{B.48})$$

To show ζ has a lower bound, it is sufficient to show $\zeta_1 = o_p(1)$ and $\zeta_2 = o_p(1)$ correspondingly.

Firstly, we are going to show $\zeta_1 = o_p(1)$. By Theorem 2, $\hat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t} = \mathcal{O}_p(t^{-\frac{1}{2}})$, which implies

$$\hat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t} = o_p\{t^{-(\frac{1}{2}-\alpha\gamma)}\}.$$

Thus we have

$$\begin{aligned} & \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbb{I}\{0 < \Delta_{\mathbf{x}_t} < T^{-\alpha}\} \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} \leq 0\} (\hat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t}) \right| \\ & \leq \frac{1}{\sqrt{T}} \sum_{t=1}^T |\mathbb{I}\{0 < \Delta_{\mathbf{x}_t} < T^{-\alpha}\}| |\mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} \leq 0\}| |(\hat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t})| \\ & \leq \frac{1}{\sqrt{T}} \sum_{t=1}^T |(\hat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t})| \leq \frac{\sqrt{T}}{T} \sum_{t=1}^T o_p(t^{-(\frac{1}{2}-\alpha\gamma)}) \stackrel{*}{=} \sqrt{T} o_p(T^{-(\frac{1}{2}-\alpha\gamma)}) = o_p(T^{\alpha\gamma}), \end{aligned}$$

where the equation (*) is derived by Lemma 6 in Luedtke and Van Der Laan (2016). By Assumption 4.3, there exists some constant c and $0 < \alpha < \frac{1}{2}$ such that $\alpha\gamma < \frac{1}{2}$ and

$$\Pr\{0 < \Delta_{\mathbf{x}_t} < T^{-\alpha}\} \leq cT^{-\alpha\gamma}.$$

Therefore we have

$$|\zeta_1| \leq \Pr\{0 < \Delta_{\mathbf{x}_t} < T^{-\alpha}\} o_p(T^{-(\frac{1}{2}-\alpha\gamma)}) = cT^{-\alpha\gamma} o_p(T^{\alpha\gamma}) = o_p(1). \quad (\text{B.49})$$

Notice that

$$\mathbb{I}\{0 < \Delta_{\mathbf{x}_t} < T^{-\alpha}\} \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} \leq 0\} (\hat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t}) \leq \mathbb{I}\{\hat{\Delta}_{\mathbf{x}_t} \leq 0\} \hat{\Delta}_{\mathbf{x}_t} \leq 0, \quad (\text{B.50})$$

where the first inequality holds since $\Delta_{\mathbf{x}_t} \geq 0$. Combining (B.48), (B.49), and (B.50), we have

$$0 \geq \zeta_1 = o_p(1). \quad (\text{B.51})$$

Next, we consider the second part ζ_2 . Note that

$$\mathbb{I}\{\widehat{\Delta}_{\mathbf{x}_t} \leq 0\} = \mathbb{I}\{\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t} \leq -\Delta_{\mathbf{x}_t}\} = \mathbb{I}\{|\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t}| > \Delta_{\mathbf{x}_t}\},$$

we have

$$\left| \mathbb{I}\{\widehat{\Delta}_{\mathbf{x}_t} \leq 0\}(\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t}) \right| = \mathbb{I}\{|\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t}| > \Delta_{\mathbf{x}_t}\} |\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t}| \leq \frac{|\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t}|^2}{\Delta_{\mathbf{x}_t}}, \quad (\text{B.52})$$

where the inequality holds since

$$\mathbb{I}\{|\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t}| > \Delta_{\mathbf{x}_t}\} \leq \frac{|\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t}|}{\Delta_{\mathbf{x}_t}}.$$

Thus, by (B.52), we further have

$$\mathbb{I}\{\widehat{\Delta}_{\mathbf{x}_t} \leq 0\}(\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t}) \geq -\frac{(\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t})^2}{\Delta_{\mathbf{x}_t}}. \quad (\text{B.53})$$

Since $\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t} < 0$, based on (B.53), we have

$$0 \geq \zeta_2 \geq \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbb{I}\{\Delta_{\mathbf{x}_t} > T^{-\alpha}\} \mathbb{I}\{\widehat{\Delta}_{\mathbf{x}_t} \leq 0\}(\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t}) \geq -\frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbb{I}\{\Delta_{\mathbf{x}_t} > T^{-\alpha}\} \frac{(\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t})^2}{\Delta_{\mathbf{x}_t}}. \quad (\text{B.54})$$

Notice that $\mathbb{I}\{\Delta_{\mathbf{x}_t} > T^{-\alpha}\} \leq \Delta_{\mathbf{x}_t} T^\alpha$, combining with (B.54), we further have

$$0 \geq \zeta_2 \geq -\frac{1}{\Delta_{\mathbf{x}_t}} \frac{1}{\sqrt{T}} \sum_{t=1}^T \frac{\Delta_{\mathbf{x}_t}}{T^{-\alpha}} (\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t})^2 = -T^{-\frac{1}{2}+\alpha} \sum_{t=1}^T (\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t})^2 = -T^{\frac{1}{2}+\alpha} \left\{ \frac{1}{T} \sum_{t=1}^T (\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t})^2 \right\}.$$

By Theorem 2, $\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t} = \mathcal{O}_p(t^{-\frac{1}{2}})$, which implies $(\widehat{\Delta}_{\mathbf{x}_t} - \Delta_{\mathbf{x}_t})^2 = o_p(T^{-(\frac{1}{2}+\alpha)})$. And by Lemma 6 in Luedtke and Van Der Laan (2016), $T^{-1} \sum_{t=1}^T o_p\{t^{-(\frac{1}{2}+\alpha)}\} = o_p\{T^{-(\frac{1}{2}+\alpha)}\}$, we have

$$0 \geq \zeta_2 \geq -T^{\frac{1}{2}+\alpha} o_p(T^{-(\frac{1}{2}+\alpha)}) = o_p(1). \quad (\text{B.55})$$

Therefore, combining Equation (B.51) and Equation (B.55), we have

$$0 \geq \zeta = \zeta_1 + \zeta_2 = o_p(1).$$

Thus, we have

$$T^{-1/2} \sum_{t=1}^T |\mu\{\mathbf{x}_t, \hat{\pi}_t(\mathbf{x}_t)\} - \mu\{\mathbf{x}_t, \pi^*\}| = o_p(1).$$

References

- Chen, H., Lu, W. and Song, R. (2020), ‘Statistical inference for online decision making: In a contextual bandit setting’, *Journal of the American Statistical Association* pp. 1–16.
- Dedecker, J. and Louhichi, S. (2002), Maximal inequalities and empirical central limit theorems, *in* ‘Empirical process techniques for dependent data’, Springer, pp. 137–159.
- Feller, W. (2008), *An introduction to probability theory and its applications, vol 2*, John Wiley & Sons.
- Hall, P. and Heyde, C. C. (2014), *Martingale limit theory and its application*, Academic press.
- Luedtke, A. R. and Van Der Laan, M. J. (2016), ‘Statistical inference for the mean outcome under a possibly non-unique optimal treatment strategy’, *Annals of statistics* **44**(2), 713.