

Supplementary Materials for Practical Counterfactual Policy Learning for Top- K Recommendations

Yaxu Liu^{*†}
National Taiwan University
d08944012@ntu.edu.tw

Jui-Nan Yen[†]
National Taiwan University
juinanyen@gmail.com

Bowen Yuan^{†‡}
Amazon
bwuyuan@amazon.com

Rundong Shi
Meituan
shirundong@meituan.com

Peng Yan
Meituan
yanpeng04@meituan.com

Chih-Jen Lin
National Taiwan University
cjlin@csie.ntu.edu.tw

1 DETAILED DERIVATIONS OF THE UNBIASEDNESS OF V_{IPS}^π

From (11), (2), (3), and Assumption 1, we have

$$\begin{aligned} V_{\text{IPS}}^\pi &= \mathbb{E}_{\text{Pr}(\mathbf{u})} \mathbb{E}_{\beta(\mathbf{A}|\mathbf{u};\mathbf{V})} \mathbb{E}_{\text{Pr}(\mathbf{c}|\mathbf{u};\mathbf{V})} \left[\frac{\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \boldsymbol{\theta})}{\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})} r(\mathbf{c}, \mathbf{A}) \right] \\ &= \mathbb{E}_{\text{Pr}(\mathbf{u})} \mathbb{E}_{\text{Pr}(\mathbf{c}|\mathbf{u};\mathbf{V})} \mathbb{E}_{\beta(\mathbf{A}|\mathbf{u};\mathbf{V})} \left[\frac{\pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \boldsymbol{\theta})}{\beta(\mathbf{A} | \mathbf{u}; \mathbf{V})} r(\mathbf{c}, \mathbf{A}) \right] \\ &= \mathbb{E}_{\text{Pr}(\mathbf{u})} \mathbb{E}_{\text{Pr}(\mathbf{c}|\mathbf{u};\mathbf{V})} \left[\sum_{\substack{\mathbf{A} \in G(\mathbb{A}, K) \\ \beta(\mathbf{A}|\mathbf{u};\mathbf{V}) \neq 0}} \pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \boldsymbol{\theta}) r(\mathbf{c}, \mathbf{A}) \right] \\ &= \mathbb{E}_{\text{Pr}(\mathbf{u})} \mathbb{E}_{\text{Pr}(\mathbf{c}|\mathbf{u};\mathbf{V})} \left[\sum_{\mathbf{A} \in G(\mathbb{A}, K)} \pi(\mathbf{A} | \mathbf{u}; \mathbf{V}, \boldsymbol{\theta}) r(\mathbf{c}, \mathbf{A}) \right] \\ &= \mathbb{E}_{\text{Pr}(\mathbf{u})} \mathbb{E}_{\pi(\mathbf{A}|\mathbf{u};\mathbf{V},\boldsymbol{\theta})} \mathbb{E}_{\text{Pr}(\mathbf{c}|\mathbf{u};\mathbf{V})} \left[r(\mathbf{c}, \mathbf{A}) \right] = V^\pi. \end{aligned}$$

2 DETAILED DERIVATIONS OF THEOREM 1

To prove Theorem 1, we need the following Lemma.

LEMMA 1 (HOEFFDING'S INEQUALITY δ -VERSION). *Assume X_1, \dots, X_m to be i.i.d. with 0 mean and $|X_i| \leq M$ almost surely. Then with probability at least $1 - \delta$, we have*

$$\left| \frac{1}{m} \sum_{i=1}^m X_i \right| \leq M \sqrt{\frac{2}{m} \log \frac{2}{\delta}}.$$

Proof of Lemma 1.

Hoeffding's inequality states that for every positive t , $P(|\frac{1}{m} \sum_{i=1}^m X_i| \geq t) \leq 2 \exp(-\frac{mt^2}{2M^2})$. Now let $\delta = 2 \exp(-\frac{mt^2}{2M^2})$. Solving for t , and we get $t = M \sqrt{\frac{2}{m} \log \frac{2}{\delta}}$, which completes our proof.

Proof of Theorem 1.

Recalling the definition,

$$\hat{V}_{\text{IPS}}^\pi(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m w_{\mathbf{A}}^i r(\mathbf{c}_i, \mathbf{A}_i).$$

Let $X_i = w_{\mathbf{A}}^i r(\mathbf{c}_i, \mathbf{A}_i) - V^\pi$, which are i.i.d. From (2), (4), Assumptions 2 and 3, we know that $|X_i| \leq K w_{\max}$. As we have $\mathbb{E}[X_i] = 0$,

^{*}Work done at Meituan as an intern.

[†]Contribute equally for this work.

[‡]Work done at National Taiwan University.

By taking δ' instead, Lemma 1 tells us that for any single π ,

$$P\left(|\hat{V}_{\text{IPS}}^\pi - V^\pi| \leq K w_{\max} \sqrt{\frac{2}{m} \log \frac{2}{\delta'}}\right) \geq 1 - \delta'. \quad (\text{A.1})$$

According to the union bound, for a countable set of events S_1, S_2, S_3, \dots , we have

$$P\left(\bigcup_{i=1}^{\infty} S_i\right) \leq \sum_{i=1}^{\infty} P(S_i). \quad (\text{A.2})$$

With the De Morgan's Law, we have

$$\begin{aligned} P\left(\bigcup_{i=1}^{\infty} S_i^c\right) &\leq \sum_{i=1}^{\infty} P(S_i^c), \\ P\left(\left(\bigcap_{i=1}^{\infty} S_i\right)^c\right) &\leq \sum_{i=1}^{\infty} (1 - P(S_i)), \\ 1 - P\left(\bigcap_{i=1}^{\infty} S_i\right) &\leq \sum_{i=1}^{\infty} (1 - P(S_i)), \\ P\left(\bigcap_{i=1}^{\infty} S_i\right) &\geq 1 - \sum_{i=1}^{\infty} (1 - P(S_i)), \end{aligned} \quad (\text{A.3})$$

where S_i^c means the complement of S_i . Then, by applying (A.1) to (A.3) over the finite policy class \mathcal{H} , where $|\mathcal{H}| = N$, we get

$$\begin{aligned} P\left(\bigcap_{\pi \in \mathcal{H}} \left\{|\hat{V}_{\text{IPS}}^\pi - V^\pi| \leq K w_{\max} \sqrt{\frac{2}{m} \log \frac{2}{\delta'}}\right\}\right) &\geq 1 - N\delta', \\ P\left(\sup_{\pi \in \mathcal{H}} \left\{|\hat{V}_{\text{IPS}}^\pi - V^\pi| \leq K w_{\max} \sqrt{\frac{2}{m} \log \frac{2}{\delta'}}\right\}\right) &\geq 1 - N\delta', \quad (\text{A.4}) \\ P\left(\sup_{\pi \in \mathcal{H}} \left\{|\hat{V}_{\text{IPS}}^\pi - V^\pi| \leq K w_{\max} \sqrt{\frac{2}{m} \log \frac{2N}{\delta}}\right\}\right) &\geq 1 - \delta, \end{aligned}$$

where $\delta \equiv N\delta'$. This completes our proof.

3 DETAILED DERIVATIONS OF THEOREM 2

To prove Theorem 2, we need the following Lemma.

LEMMA 2. *Given $\mathbf{u}, \mathbf{A}_{1:k}$, and c_{A_k} , we have*

$$\mathbb{E}_{P_{\beta}(\mathbf{A}_{k+1:k} | \mathbf{u}, \mathbf{A}_{1:k}, c_{A_k})} \left[\left(\prod_{j=k+1}^K w_{A_j} \right) \right] = 1. \quad (\text{A.5})$$

Proof of Lemma 2.

From our assumption $P_\beta(c_{A_k} | \mathbf{u}, \mathbf{A}_{1:k}, \mathbf{A}_{k+1:K}) = P_\beta(c_{A_k} | \mathbf{u}, \mathbf{A}_{1:k})$, we have

$$\begin{aligned}
& P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}, c_{A_k}) \\
&= \frac{P_\beta(\mathbf{u}, \mathbf{A}_{1:k}, \mathbf{A}_{k+1:K}, c_{A_k})}{P_\beta(\mathbf{u}, \mathbf{A}_{1:k}, c_{A_k})} \\
&= \frac{P_\beta(\mathbf{u}, \mathbf{A}_{1:k}) P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}) P_\beta(c_{A_k} | \mathbf{u}, \mathbf{A}_{1:k}, \mathbf{A}_{k+1:K})}{P_\beta(\mathbf{u}, \mathbf{A}_{1:k}) P_\beta(c_{A_k} | \mathbf{u}, \mathbf{A}_{1:k})} \quad (\text{A.6}) \\
&= \frac{P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}) P_\beta(c_{A_k} | \mathbf{u}, \mathbf{A}_{1:k}, \mathbf{A}_{k+1:K})}{P_\beta(c_{A_k} | \mathbf{u}, \mathbf{A}_{1:k})} \\
&= P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}).
\end{aligned}$$

According to our definition in Section 2.1, $P_\beta(\cdot)$ is the probability distribution decided by the policy β . Thus, from the structure of β in (19), we have

$$P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}) = \beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V}).$$

Then, from (11), we have

$$\mathbb{E}_{P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}, c_{A_k})} \left[\left(\prod_{j=k+1}^K w_{A_j} \right) \right] \quad (\text{A.7})$$

$$= \mathbb{E}_{P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}, c_{A_k})} \left[\frac{\pi(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V}, \boldsymbol{\theta})}{\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V})} \right] \quad (\text{A.8})$$

$$\begin{aligned}
&= \sum_{\substack{\mathbf{A}_{k+1:K} \in G(\mathbb{A} \setminus \mathbf{A}_{1:k}, K-k) \\ P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}, c_{A_k}) \neq 0}} P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}, c_{A_k}) \\
&\quad \times \frac{\pi(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V}, \boldsymbol{\theta})}{\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V})} \quad (\text{A.9})
\end{aligned}$$

$$\begin{aligned}
&= \sum_{\substack{\mathbf{A}_{k+1:K} \in G(\mathbb{A} \setminus \mathbf{A}_{1:k}, K-k) \\ \beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V}) \neq 0}} \beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V}) \frac{\pi(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V}, \boldsymbol{\theta})}{\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V})} \\
&\quad (\text{A.10})
\end{aligned}$$

$$\begin{aligned}
&= \sum_{\substack{\mathbf{A}_{k+1:K} \in G(\mathbb{A} \setminus \mathbf{A}_{1:k}, K-k) \\ \beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V}) \neq 0}} \pi(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V}, \boldsymbol{\theta}) \\
&\quad (\text{A.11})
\end{aligned}$$

$$\begin{aligned}
&= \sum_{\mathbf{A}_{k+1:K} \in G(\mathbb{A} \setminus \mathbf{A}_{1:k}, K-k)} \pi(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V}, \boldsymbol{\theta}), \quad (\text{A.12})
\end{aligned}$$

where (A.10) comes from (A.6), and the equation (A.12) comes from Assumption 1. Finally, we have $\sum_{\mathbf{A}_{k+1:K} \in G(\mathbb{A} \setminus \mathbf{A}_{1:k}, K-k)} \pi(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}; \mathbf{V}, \boldsymbol{\theta}) = 1$, which follows from the fact that π is a probability distribution. This completes our proof.

Proof of Theorem 2.

From (11), (20), and Lemma 2, we have

$$V_{\text{IPS}}^\pi = \mathbb{E}_{P_\beta} \left[\sum_{k=1}^K w_{A_k} c_{A_k} \right] \quad (\text{A.13})$$

$$= \sum_{k=1}^K \mathbb{E}_{P_\beta} \left[w_{A_k} c_{A_k} \right] \quad (\text{A.14})$$

$$\begin{aligned}
&= \sum_{k=1}^K \mathbb{E}_{P_\beta(\mathbf{u}, \mathbf{A}_{1:k}, c_{A_k})} \mathbb{E}_{P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}, c_{A_k})} \left[\left(\prod_{j=1}^k w_{A_j} \right) \right. \\
&\quad \left. \times \left(\prod_{j=k+1}^K w_{A_j} \right) c_{A_k} \right] \quad (\text{A.15})
\end{aligned}$$

$$\begin{aligned}
&= \sum_{k=1}^K \mathbb{E}_{P_\beta(\mathbf{u}, \mathbf{A}_{1:k}, c_{A_k})} \left[\left(\prod_{j=1}^k w_{A_j} \right) c_{A_k} \right. \\
&\quad \left. \times \mathbb{E}_{P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}, c_{A_k})} \left[\left(\prod_{j=k+1}^K w_{A_j} \right) \right] \right] \quad (\text{A.16})
\end{aligned}$$

$$\begin{aligned}
&= \sum_{k=1}^K \mathbb{E}_{P_\beta(\mathbf{u}, \mathbf{A}_{1:k}, c_{A_k})} \left[\left(\prod_{j=1}^k w_{A_j} \right) c_{A_k} \right] \quad (\text{A.17})
\end{aligned}$$

$$\begin{aligned}
&= \sum_{k=1}^K \mathbb{E}_{P_\beta(\mathbf{u}, \mathbf{A}_{1:k}, c_{A_k})} \mathbb{E}_{P_\beta(\mathbf{A}_{k+1:K} | \mathbf{u}, \mathbf{A}_{1:k}, c_{A_k})} \left[\left(\prod_{j=1}^k w_{A_j} \right) c_{A_k} \right] \\
&\quad (\text{A.18})
\end{aligned}$$

$$\begin{aligned}
&= \sum_{k=1}^K \mathbb{E}_{P_\beta} \left[\left(\prod_{j=1}^k w_{A_j} \right) c_{A_k} \right], \quad (\text{A.19})
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_{P_\beta} \left[\sum_{k=1}^K \left(\prod_{j=1}^k w_{A_j} \right) c_{A_k} \right], \quad (\text{A.20})
\end{aligned}$$

where (A.14), (A.16), and (A.20) rely on the linearity of expectation $\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ and $\mathbb{E}[aX] = a\mathbb{E}[X]$.

4 DETAILED DERIVATIONS OF THE SECOND-ORDER APPROXIMATION AS SQUARED LOSS

Let $T(\hat{Y}_{ij}, \tilde{y}_j)$ be the second-order approximation of $\ell_{\log}^-(\hat{Y}_{ij})$ at \tilde{y}_j . By denoting $\ell_{\log}^-(\tilde{y}_j)$ as E_{j0} , $\nabla \ell_{\log}^-(\tilde{y}_j)$ as E_{j1} and $\nabla^2 \ell_{\log}^-(\tilde{y}_j)$ as E_{j2} , we have

$$\begin{aligned}
T(\hat{Y}_{ij}, \tilde{y}_j) &= E_{j0} + E_{j1}(\hat{Y}_{ij} - \tilde{y}_j) + \frac{E_{j2}}{2}(\hat{Y}_{ij} - \tilde{y}_j)^2 \\
&= \frac{1}{2} E_{j2}(\hat{Y}_{ij})^2 + \frac{E_{j1} - E_{j2}\tilde{y}_j}{E_{j2}} \hat{Y}_{ij} \\
&\quad + (E_{j0} - E_{j1}\tilde{y}_j + \frac{1}{2} E_{j2}\tilde{y}_j^2 - \frac{1}{2} \frac{(E_{j1} - E_{j2}\tilde{y}_j)^2}{E_{j2}}), \quad (\text{A.21})
\end{aligned}$$

where the first part is a squared loss and the second part can be omitted as a constant for \hat{Y}_{ij} .