



# Uncertainty-Aware Instance Reweighting for Off-Policy Learning

Xiaoying Zhang, Junpu Chen, Hongning Wang, Hong Xie, Yang Liu, John C.S. Lui, Hang Li



## Background : Off-Policy Learning

**Input:** A logged dataset  $D = \{(x_n, a_n, r_{x_n, a_n})\}_{n=1}^N$  generated by logging policy  $\beta^*(a|x)$ .

**Goal:** Learning a policy  $\pi(a|x)$  that maximize  $V(\pi) = E_{\pi}[r_{x,a}] = E_{\beta^*}[\frac{\pi(a|x)}{\beta^*(a|x)} r_{x,a}]$

## Inverse Propensity Score (IPS)

$$\hat{V}_{IPS}(\pi) = \frac{1}{N} \sum_{n=1}^N \frac{\pi(a_n|x_n)}{\beta^*(a_n|x_n)} r_{x_n, a_n}$$

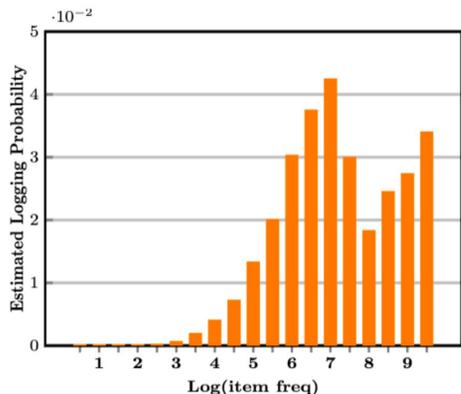
## Approximation when $\beta^*$ is unknown $\rightarrow$ BIPS

In practice,  $\beta^*$  is usually unknown and **approximated** by its estimate  $\hat{\beta}$ .

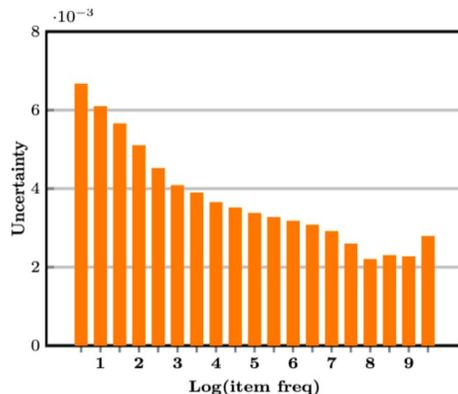
$$\hat{V}_{BIPS}(\pi) = \frac{1}{N} \sum_{n=1}^N \frac{\pi(a_n|x_n)}{\hat{\beta}(a_n|x_n)} r_{x_n, a_n}$$

## BIPS suffers high bias and variance with inaccurate $\hat{\beta}$

- Theoretically, inaccurate and small  $\hat{\beta}(a|x) \rightarrow$  high bias and variance of  $\hat{V}_{BIPS}$ .
- However, **smaller**  $\hat{\beta}(a|x) \rightarrow$  more likely to be **inaccurate** (i.e., high uncertainty in estimation).



(a) Estimated Logging Probability



(b) Uncertainty of Estimation

## Uncertainty-Aware Off-Policy Learning (UIPS)

$$\hat{V}_{UIPS}(\pi_{\vartheta}) = \frac{1}{N} \sum_{n=1}^N \frac{\pi_{\vartheta}(a_n|x_n)}{\hat{\beta}(a_n|x_n)} \cdot \phi_{x_n, a_n} \cdot r_{x_n, a_n}$$

- Per-sample weight  $\phi_{x,a}$ : small value when  $\hat{\beta}(a|x)$  is small and far from  $\beta^*(a|x)$ .

Iterate between two steps

Step 1: Deriving Optimal  $\phi_{x,a}^*$ :  $\hat{V}_{UIPS}(\pi_{\vartheta}) \rightarrow V(\pi_{\vartheta})$

Minimize the upper bound of Mean Square Error (MSE) of  $\hat{V}_{UIPS}(\pi_{\vartheta})$  to its ground-truth  $V(\pi_{\vartheta}) \rightarrow$  Per-sample optimization:

$$\min_{\phi_{x,a}} \lambda \left( \frac{\beta^*(a|x)}{\hat{\beta}(a|x)} \phi_{x,a} - 1 \right)^2 + \frac{\pi_{\vartheta}(a|x)^2}{\hat{\beta}(a|x)^2} \phi_{x,a}^2$$

Direct optimization is infeasible due to the unknown  $\beta^*(a|x)$ , but uncertainty estimation provides its confidence interval  $\beta^*(a|x) \in B_{x,a}$ :

$$\min_{\phi_{x,a}} \max_{\beta_{x,a}} \lambda \left( \frac{\beta_{x,a}}{\hat{\beta}(a|x)} \phi_{x,a} - 1 \right)^2 + \frac{\pi_{\vartheta}(a|x)^2}{\hat{\beta}(a|x)^2} \phi_{x,a}^2$$

$\rightarrow$  Closed-form solution of  $\phi_{x,a}^*$  !

Step 2: Policy Improvement:

$$\nabla_{\vartheta} \hat{V}_{UIPS}(\pi_{\vartheta}) = \frac{1}{N} \sum_{n=1}^N \frac{\pi_{\vartheta}(a_n|x_n)}{\beta^*(a_n|x_n)} \cdot \phi_{x_n, a_n}^* \cdot r_{x_n, a_n} \nabla_{\vartheta} \log(\pi_{\vartheta}(a_n|x_n))$$

## Theoretical Convergence of UIPS

UIPS converges to a stationary point where the true policy gradient is zero, while convergence of policy learning under BIPS is not guaranteed!

## Empirical Results

Improved performance on both synthetic datasets and three unbiased recommendation datasets.

Algorithm	Yahoo			Coat			KuaiRec		
	P@5	R@5	NDCG@5	P@5	R@5	NDCG@5	P@50	R@50	NDCG@50
CE	0.2819±2e <sup>-3</sup>	0.7594±6e <sup>-3</sup>	0.6073±7e <sup>-3</sup>	0.2799±5e <sup>-3</sup>	0.4618±1e <sup>-2</sup>	0.4529±7e <sup>-3</sup>	0.8802±2e <sup>-3</sup>	0.0240±8e <sup>-5</sup>	0.8810±6e <sup>-3</sup>
BIPS-Cap	0.2808±2e <sup>-3</sup>	0.7576±5e <sup>-3</sup>	0.6099±8e <sup>-3</sup>	0.2758±6e <sup>-3</sup>	0.4582±7e <sup>-3</sup>	0.4399±9e <sup>-3</sup>	0.8750±3e <sup>-3</sup>	0.0238±7e <sup>-5</sup>	0.8788±5e <sup>-3</sup>
MinVar	0.2843±4e <sup>-3</sup>	0.7685±1e <sup>-2</sup>	0.6168±1e <sup>-2</sup>	0.2813±3e <sup>-3</sup>	0.4668±9e <sup>-3</sup>	0.4414±8e <sup>-3</sup>	0.8827±1e <sup>-3</sup>	0.0240±5e <sup>-5</sup>	0.8886±2e <sup>-3</sup>
stableVar	0.2787±2e <sup>-3</sup>	0.7499±7e <sup>-3</sup>	0.5919±7e <sup>-3</sup>	0.2840±3e <sup>-3</sup>	0.4662±5e <sup>-3</sup>	0.4393±7e <sup>-3</sup>	0.8524±7e <sup>-3</sup>	0.0231±2e <sup>-4</sup>	0.8570±4e <sup>-3</sup>
Shrinkage	0.2843±3e <sup>-3</sup>	0.7654±8e <sup>-3</sup>	0.6204±7e <sup>-3</sup>	0.2790±5e <sup>-3</sup>	0.4636±4e <sup>-3</sup>	0.4464±1e <sup>-2</sup>	0.8744±3e <sup>-3</sup>	0.0238±9e <sup>-5</sup>	0.8771±6e <sup>-3</sup>
SNIPS	0.2222±4e <sup>-3</sup>	0.5828±1e <sup>-2</sup>	0.4357±1e <sup>-2</sup>	0.2643±7e <sup>-3</sup>	0.4287±1e <sup>-2</sup>	0.4009±9e <sup>-3</sup>	0.8411±6e <sup>-3</sup>	0.0228±2e <sup>-4</sup>	0.8431±6e <sup>-3</sup>
BanditNet	0.2413±8e <sup>-3</sup>	0.6442±2e <sup>-2</sup>	0.4988±2e <sup>-2</sup>	0.2781±8e <sup>-3</sup>	0.4527±1e <sup>-2</sup>	0.4251±1e <sup>-2</sup>	0.8758±5e <sup>-3</sup>	0.0239±2e <sup>-4</sup>	0.8810±4e <sup>-3</sup>
POEM	0.2732±3e <sup>-3</sup>	0.7357±1e <sup>-2</sup>	0.5880±1e <sup>-2</sup>	0.2791±4e <sup>-3</sup>	0.4566±6e <sup>-3</sup>	0.4375±6e <sup>-3</sup>	0.7785±1e <sup>-2</sup>	0.0210±2e <sup>-4</sup>	0.7779±6e <sup>-3</sup>
POXM	0.2250±5e <sup>-3</sup>	0.5940±1e <sup>-2</sup>	0.4542±2e <sup>-2</sup>	0.2663±6e <sup>-3</sup>	0.4308±9e <sup>-3</sup>	0.4006±1e <sup>-2</sup>	0.8962±1e <sup>-3</sup>	0.0245±4e <sup>-4</sup>	0.9041±1e <sup>-2</sup>
Adaptive	0.2762±3e <sup>-3</sup>	0.7451±9e <sup>-3</sup>	0.5919±8e <sup>-3</sup>	0.2830±3e <sup>-3</sup>	0.4634±5e <sup>-3</sup>	0.4217±5e <sup>-3</sup>	0.8375±1e <sup>-2</sup>	0.0227±4e <sup>-4</sup>	0.8460±1e <sup>-2</sup>
ApproxKNN	0.2697±2e <sup>-3</sup>	0.7225±5e <sup>-3</sup>	0.5760±6e <sup>-3</sup>	0.2755±2e <sup>-3</sup>	0.4594±5e <sup>-3</sup>	0.4490±4e <sup>-3</sup>	0.8839±2e <sup>-6</sup>	0.0240±5e <sup>-5</sup>	0.8895±2e <sup>-3</sup>
IPS-C-TS	0.2816±2e <sup>-3</sup>	0.7582±5e <sup>-3</sup>	0.6114±5e <sup>-3</sup>	0.2799±3e <sup>-3</sup>	0.4625±7e <sup>-3</sup>	0.4462±6e <sup>-3</sup>	0.8781±3e <sup>-3</sup>	0.0239±1e <sup>-4</sup>	0.8749±3e <sup>-3</sup>
UIPS-P	0.1829±8e <sup>-3</sup>	0.4560±3e <sup>-2</sup>	0.3300±1e <sup>-2</sup>	0.2685±7e <sup>-3</sup>	0.4364±9e <sup>-3</sup>	0.4087±7e <sup>-3</sup>	0.8638±8e <sup>-3</sup>	0.0235±3e <sup>-4</sup>	0.8685±7e <sup>-3</sup>
UIPS-O	0.1947±3e <sup>-3</sup>	0.4959±1e <sup>-2</sup>	0.3600±8e <sup>-3</sup>	0.2657±5e <sup>-3</sup>	0.4306±9e <sup>-3</sup>	0.4146±9e <sup>-3</sup>	0.8651±8e <sup>-3</sup>	0.0235±2e <sup>-4</sup>	0.8697±7e <sup>-3</sup>
UIPS	<b>0.2868±2e<sup>-3</sup></b>	<b>0.7742±5e<sup>-3</sup></b>	<b>0.6274±5e<sup>-3</sup></b>	<b>0.2877±3e<sup>-3</sup></b>	<b>0.4757±5e<sup>-3</sup></b>	<b>0.4576±8e<sup>-3</sup></b>	<b>0.9120±1e<sup>-3</sup></b>	<b>0.0250±5e<sup>-5</sup></b>	<b>0.9174±7e<sup>-4</sup></b>
p-value	4e <sup>-2</sup>	1e <sup>-2</sup>	3e <sup>-2</sup>	2e <sup>-2</sup>	6e <sup>-4</sup>	5e <sup>-5</sup>	6e <sup>-4</sup>	6e <sup>-4</sup>	1e <sup>-3</sup>

More accurate Off-policy Evaluation

Table 3: MSE of different off-policy estimators. A lower MSE indicates a more accurate estimator.

Algorithm	IPS-GT	BIPS	minVar	stableVar	Shrinkage	UIPS
$\tau = 0.5$	0.0875±4e <sup>-4</sup>	15.786±1.51	0.9021±7e <sup>-13</sup>	0.8612±5e <sup>-8</sup>	0.0718±5e <sup>-6</sup>	<b>0.0210±2e<sup>-6</sup></b>
$\tau = 1.0$	0.0209±8e <sup>-5</sup>	0.5510±0.388	0.9019±8e <sup>-12</sup>	0.8578±2e <sup>-7</sup>	0.1978±2e <sup>-5</sup>	<b>0.0093±1e<sup>-6</sup></b>
$\tau = 2.0$	0.0020±6e <sup>-6</sup>	0.5669±0.013	0.9015±5e <sup>-15</sup>	0.8342±5e <sup>-7</sup>	0.2952±3e <sup>-5</sup>	<b>0.0043±4e<sup>-7</sup></b>

Performance under different uncertainties

- the only off-policy algorithm that outperforms CE on test samples with high uncertainty.

Algorithm	Actions on Samples with High Uncertainty			Actions on Samples with Low Uncertainty		
	P@5(RI)	R@5(RI)	NDCG@5(RI)	P@5(RI)	R@5(RI)	NDCG@5(RI)
CE	0.5190	0.1231	0.5526	0.5913	0.1915	0.6549
BIPS-Cap	0.5117 (-1.41%)	0.1202 (-2.33%)	0.5488 (-0.68%)	0.5913 (+0.00%)	0.1903 (-0.64%)	0.6574 (+0.39%)
Shrinkage	0.5158 (-0.62%)	0.1217 (-1.11%)	0.5505 (-0.37%)	0.5892 (-0.35%)	0.1905 (-0.55%)	0.6546 (-0.05%)
UIPS	0.5222 (+0.61%)	0.1237 (+0.50%)	0.5568 (+0.77%)	0.5994 (+1.38%)	0.1940 (+1.28%)	0.6658 (+1.66%)