

---

# Improved Estimator Selection for Off-Policy Evaluation

---

George Tucker<sup>1</sup> Jonathan Lee<sup>2\*</sup>

## Abstract

Off-policy policy evaluation is a fundamental problem in reinforcement learning. As a result, many estimators with different tradeoffs have been developed; however, selecting the best estimator is challenging with limited data and without additional interactive data collection. Recently, Su et al. (2020b) developed a data-dependent selection procedure that competes with the oracle selection up to a constant and demonstrate its practicality. We refine the analysis to remove an extraneous assumption and improve the procedure. The improved procedure results in a tighter oracle bound and stronger empirical results on a contextual bandit task.

## 1. Introduction

Off-policy policy evaluation is a fundamental problem in reinforcement learning. Moreover, in applications of reinforcement learning, the ability to evaluate potential policies without deploying them is paramount to reducing risk. As a result, developing off-policy estimators has been a subject of intense research activity (e.g., (Horvitz & Thompson, 1952; Precup, 2000; Dudík et al., 2014; Jiang & Li, 2016; Thomas & Brunskill, 2016; Wang et al., 2017; Farajtabar et al., 2018; Liu et al., 2018; Kallus & Zhou, 2018; Voloshin et al., 2019)).

However, selecting the best estimator for a particular application is challenging. There is no straightforward analogue to cross-validation in supervised learning. Recently, Su et al. (2020b) introduced SLOPE, a generic, data-driven selection procedure, which is provable competitive with the oracle selection (up to a constant) and they demonstrate that it empirically outperforms other selection approaches based on MSE surrogates (Thomas & Brunskill, 2016; Wang et al., 2017; Su et al., 2020a).

---

\*Work done while at Google. <sup>1</sup>Google Research, Brain Team <sup>2</sup>Stanford University. Correspondence to: George Tucker <gjt@google.com>.

By refining the analysis underlying their approach, we can remove an extraneous assumption in their procedure and tighten the oracle bound on performance. Furthermore, we show that the improved procedure outperforms SLOPE on the task introduced in Su et al. (2020b).

## 2. Background

We work in the same abstract setup as Su et al. (2020b). Briefly, we would like to estimate a real-valued statistic  $\theta^* := \theta(D)$  defined over an unknown data distribution  $D$ . Given a finite set of estimators  $\{\theta_i(\cdot)\}$  defined over empirical data distributions and a (uniformly-weighted) empirical data distribution  $\hat{D}$  formed from an i.i.d. sample of data  $\{x_i\}$  from  $D$ , our goal is to choose the estimator that minimizes the absolute error to the statistic of interest  $\theta^*$ . In other words, we would like a procedure that selects an estimator index  $\hat{i}$  such that

$$|\hat{\theta}_i - \theta^*| \leq \text{CONST} \times \min_i |\hat{\theta}_i - \theta^*|,$$

where  $\hat{\theta}_i := \theta_i(\hat{D})$  and CONST is a universal constant. To start, we can bound the distance as

$$|\hat{\theta}_i - \theta^*| \leq |\bar{\theta}_i - \theta^*| + |\hat{\theta}_i - \bar{\theta}_i| := \text{BIAS}(i) + \text{DEV}(i),$$

where  $\bar{\theta}_i = \mathbb{E}_{\hat{D}}[\hat{\theta}_i]$ . DEV( $i$ ) characterizes statistical fluctuations, and we can generally bound it with high confidence, whereas BIAS( $i$ ) is typically unknown.

### 2.1. SLOPE

Su et al. (2020b) introduce the SLOPE procedure based on Lepski’s principle (Lepskii, 1991; 1992; 1993; Lepski & Spokoiny, 1997; Mathé, 2006). The procedure assumes a bound on bias  $B(i) \geq \text{BIAS}(i)$ <sup>1</sup> and access to a high confidence bound on DEV( $i$ ) (i.e., with probability at least  $1 - \delta$ ,  $\text{DEV}(i) \leq \text{CNF}(i)$  for all  $i$ ). The procedure also requires monotonicity assumptions:

1.  $B(i) \leq B(i + 1)$
2.  $\exists \kappa > 0$  such that  $\kappa \text{CNF}(i) \leq \text{CNF}(i + 1) \leq \text{CNF}(i)$

Then defining  $I(i) = [\hat{\theta}_i - 2 \text{CNF}(i), \hat{\theta}_i + 2 \text{CNF}(i)]$  and

$$\hat{i} = \max\{i : \cap_{j=1}^i I(j) \neq \emptyset\},$$

---

<sup>1</sup>B( $i$ ) is not used in the procedure but appears in the bound.

with probability at least  $1 - \delta$ ,

$$|\hat{\theta}_i - \theta^*| \leq 6(1 + \kappa^{-1}) \min_i (B(i) + \text{CNF}(i)).$$

So the SLOPE procedure provides an oracle inequality on the competitiveness of the procedure. It is natural for many families of estimators to exhibit a similar trade-off of complexity, captured by the confidence bounds, with bias.

### 3. Methods

We improve the procedure by refining the analysis.

**Theorem 1.** *Given  $\delta > 0$ , high confidence bounds  $\text{CNF}(i)$  on the deviations (i.e., with probability at least  $1 - \delta$ ,  $\text{DEV}(i) \leq \text{CNF}(i)$  for all  $i$ ), and that we have ordered<sup>2</sup> the estimators such that  $\text{CNF}(i) \geq \text{CNF}(i + 1)$ . Selecting our estimator as*

$$\hat{i} = \max\{i : |\hat{\theta}_i - \hat{\theta}_j| \leq \text{CNF}(i) + (\sqrt{6} - 1) \text{CNF}(j), j < i\},$$

ensures that with probability at least  $1 - \delta$ ,

$$|\hat{\theta}_i - \theta^*| \leq (\sqrt{6} + 3) \min_i \left( \max_{j \leq i} \text{BIAS}(j) + \text{CNF}(i) \right).$$

In the common case that the estimator family has increasing bias and tighter bounds on deviation with increasing  $i$  (a family of estimators with a bias/variance trade-off e.g., importance weight clipping), then the bound simplifies to

$$|\hat{\theta}_i - \theta^*| \leq (\sqrt{6} + 3) \min_i (\text{BIAS}(i) + \text{CNF}(i)).$$

In other words, we are able to compete with the oracle estimator (up to a problem-independent constant). Compared to [Su et al. \(2020b\)](#), we do not assume  $\text{CNF}(i)$  decreases slowly and thus our bound does not depend on the problem-dependent parameter  $\kappa$  that controls this decrease. As  $\frac{1}{\kappa} \geq 1$ , we have  $\sqrt{6} + 3 < 6 < 6(1 + 1) \leq 6(1 + \kappa^{-1})$  for the constant on the oracle bound.

*Proof.* We will prove a generalized bound. Let  $B(i) \geq \text{BIAS}(i)$  be a bound on bias satisfying  $B(i) \leq B(i + 1)$ . Then, define  $\hat{i} = \max\{i : |\hat{\theta}_i - \hat{\theta}_j| \leq a \text{CNF}(i) + b \text{CNF}(j), j < i\}$ , where  $a$  and  $b$  are constants that we will optimize.

Lepski's principle leverages the fact that comparing pairs of estimators simplifies the problem. In our case, the difficulty arises from the unknown bias of our estimators, however, we do know that for a pair of estimators

$$\begin{aligned} B(j) - B(i) &\leq |\bar{\theta}_i - \bar{\theta}_j| \\ &\leq |\hat{\theta}_i - \hat{\theta}_j| + |\hat{\theta}_i - \hat{\theta}_j - (\bar{\theta}_i - \bar{\theta}_j)| \\ &\leq |\hat{\theta}_i - \hat{\theta}_j| + \text{CNF}(i) + \text{CNF}(j), \quad (1) \end{aligned}$$

<sup>2</sup>We can increase  $\text{CNF}(i)$  to accommodate any ordering at the expense of a worse bound.

which allows us to control the increase in bias in terms of known quantities.

Let  $i^* = \min_i B(i) + \text{CNF}(i)$ . When  $i^* = \hat{i}$ , we are done, so we consider the other two cases: 1)  $\hat{i} < i^*$  and 2)  $\hat{i} > i^*$ .

1. We know that  $B(\hat{i}) \leq B(i^*)$ , so we need to ensure that  $\text{CNF}$  has not decreased too much. By definition of  $\hat{i}$ , there exists  $j \leq \hat{i}$  such that

$$\begin{aligned} b \text{CNF}(j) + a \text{CNF}(\hat{i} + 1) &< |\hat{\theta}_j - \hat{\theta}_{i+1}| \\ &\leq B(j) + B(\hat{i} + 1) + \text{CNF}(j) + \text{CNF}(\hat{i} + 1). \end{aligned}$$

Implying that

$$\begin{aligned} (b - 1) \text{CNF}(\hat{i}) &\leq (b - 1) \text{CNF}(j) \\ &< B(j) + B(\hat{i} + 1) + (1 - a) \text{CNF}(\hat{i} + 1) \\ &\leq 2B(i^*) + (1 - a) \text{CNF}(\hat{i} + 1) \\ &\leq 2B(i^*) \end{aligned}$$

as long as  $a \geq 1$ . So, with  $b > 1$ , we conclude that

$$\begin{aligned} B(\hat{i}) + \text{CNF}(\hat{i}) &< \left( \frac{2}{b - 1} + 1 \right) B(i^*) \\ &\leq \left( \frac{2}{b - 1} + 1 \right) (B(i^*) + \text{CNF}(i^*)) \end{aligned}$$

2. In the second case,  $\text{CNF}(\hat{i}) \leq \text{CNF}(i^*)$ , so we need to show that the bias is not too much larger. By definition of  $\hat{i}$  and Eq. 1,

$$\begin{aligned} B(\hat{i}) - B(i^*) &\leq |\hat{\theta}_i - \hat{\theta}_{i^*}| + \text{CNF}(\hat{i}) + \text{CNF}(i^*) \\ &\leq (a + 1) \text{CNF}(\hat{i}) + (b + 1) \text{CNF}(i^*) \\ &\leq (a + b + 2) \text{CNF}(i^*) \end{aligned}$$

implying that

$$\begin{aligned} B(\hat{i}) + \text{CNF}(\hat{i}) &\leq B(i^*) + (a + b + 3) \text{CNF}(i^*) \\ &\leq (a + b + 3)(B(i^*) + \text{CNF}(i^*)). \end{aligned}$$

Together, we conclude that

$$B(\hat{i}) + \text{CNF}(\hat{i}) \leq \max \left( \left( \frac{2}{b - 1} + 1 \right), (a + b + 3) \right) (B(i^*) + \text{CNF}(i^*))$$

Minimizing over  $a$  and  $b$  results in  $a = 1$  and  $b = \sqrt{6} - 1$ . Then choosing  $B(i) = \max_{j \leq i} \text{BIAS}(j)$  finishes the proof.  $\square$

Previous works have focused on MSE as the evaluation metric, so as a corollary, we have

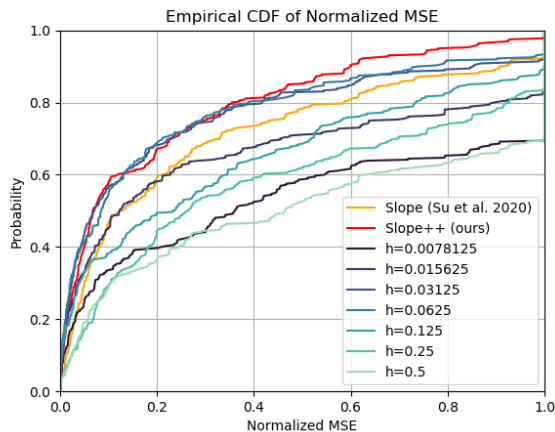


Figure 1. Experimental results for bandwidth selection of OPE estimators for contextual bandits with continuous actions. The plot shows the empirical CDF of normalized MSE (normalized by the worst MSE for that condition) across all conditions.

**Corollary 1.1.** *With the same assumptions as Thm 1 and additionally that  $\theta^*, \hat{\theta}_i \in [0, R]$  a.s. for all  $i$  and that CNF is deterministic. Then, for any  $\delta > 0$ ,*

$$\mathbb{E} \left[ \left( \hat{\theta}_i - \theta^* \right)^2 \right] \leq C \min_i \left( \max_{j \leq i} \text{BIAS}(j)^2 + \text{CNF}(i; \delta)^2 \right) + R^2 \delta,$$

where  $C$  is a universal constant.

The proof follows immediately from Thm 1 and the proof of (Corollary 4; Su et al., 2020b).

## 4. Experiments

Following Su et al. (2020b), we empirically evaluate the procedure on bandwidth selection in a synthetic environment for continuous action contextual bandits. We follow the same experimental protocol and modify the code provided with Su et al. (2020b) at <https://github.com/VowpalWabbit/slope-experiments> to implement our procedure.

Briefly, the experiment evaluates the procedure on environments with varying target and logging policies (NN, tree), “softening” approaches for randomization following Farajtabar et al. (2018) (friendly, adversarial), Lipschitz constant (0.1, 1, 10), and samples (10, 100, 1000)<sup>3</sup>.

We use 7 different choices of geometrically spaced bandwidths  $\{2^{-i} : i \in [7]\}$ . We evaluate the performance of each

<sup>3</sup>From the provided configuration file [https://github.com/VowpalWabbit/slope-experiments/blob/master/scripts/ak\\_commands.py](https://github.com/VowpalWabbit/slope-experiments/blob/master/scripts/ak_commands.py)

of these fixed choices, SLOPE, and the improved procedure (labeled SLOPE++ in the plot). As in Su et al. (2020b), we use the empirical standard deviation to construct the high confidence bound, which is valid asymptotically and usually yields better practical performance than a concentration inequality.

To measure performance, we estimate the ground truth statistic with  $100k$  samples from the target policy. To aggregate results across conditions, we normalize by the worst MSE in each condition, then compute the empirical CDF of the normalized MSE (Figure 1). Increasingly performant methods lie in the top-left quadrant. We see that SLOPE++ is the top performer compared with the fixed bandwidths and SLOPE.

## References

- Dudík, M., Erhan, D., Langford, J., Li, L., et al. Doubly robust policy evaluation and optimization. *Statistical Science*, 29(4):485–511, 2014.
- Farajtabar, M., Chow, Y., and Ghavamzadeh, M. More robust doubly robust off-policy evaluation. In *International Conference on Machine Learning*, pp. 1447–1456. PMLR, 2018.
- Horvitz, D. G. and Thompson, D. J. A generalization of sampling without replacement from a finite universe. *Journal of the American statistical Association*, 47(260):663–685, 1952.
- Jiang, N. and Li, L. Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning*, pp. 652–661. PMLR, 2016.
- Kallus, N. and Zhou, A. Policy evaluation and optimization with continuous treatments. In *International Conference on Artificial Intelligence and Statistics*, pp. 1243–1251. PMLR, 2018.
- Lepski, O. V. and Spokoiny, V. G. Optimal pointwise adaptive methods in nonparametric estimation. *The Annals of Statistics*, pp. 2512–2546, 1997.
- Lepskii, O. On a problem of adaptive estimation in gaussian white noise. *Theory of Probability & Its Applications*, 35(3):454–466, 1991.
- Lepskii, O. Asymptotically minimax adaptive estimation. i: Upper bounds. optimally adaptive estimates. *Theory of Probability & Its Applications*, 36(4):682–697, 1992.
- Lepskii, O. Asymptotically minimax adaptive estimation. ii. schemes without optimal adaptation: Adaptive estimators. *Theory of Probability & Its Applications*, 37(3):433–448, 1993.

- Liu, Q., Li, L., Tang, Z., and Zhou, D. Breaking the curse of horizon: Infinite-horizon off-policy estimation. *arXiv preprint arXiv:1810.12429*, 2018.
- Mathé, P. The lepskii principle revisited. *Inverse problems*, 22(3):L11, 2006.
- Precup, D. Eligibility traces for off-policy policy evaluation. *Computer Science Department Faculty Publication Series*, pp. 80, 2000.
- Su, Y., Dimakopoulou, M., Krishnamurthy, A., and Dudík, M. Doubly robust off-policy evaluation with shrinkage. In *International Conference on Machine Learning*, pp. 9167–9176. PMLR, 2020a.
- Su, Y., Srinath, P., and Krishnamurthy, A. Adaptive estimator selection for off-policy evaluation. In *International Conference on Machine Learning*, pp. 9196–9205. PMLR, 2020b.
- Thomas, P. and Brunskill, E. Data-efficient off-policy policy evaluation for reinforcement learning. In *International Conference on Machine Learning*, pp. 2139–2148. PMLR, 2016.
- Voloshin, C., Le, H. M., Jiang, N., and Yue, Y. Empirical study of off-policy policy evaluation for reinforcement learning. *arXiv preprint arXiv:1911.06854*, 2019.
- Wang, Y.-X., Agarwal, A., and Dudík, M. Optimal and adaptive off-policy evaluation in contextual bandits. In *International Conference on Machine Learning*, pp. 3589–3597. PMLR, 2017.