

A. Proofs

A.1. Proofs for Section 3

Proof of Theorem 3. The proof is similar to that of Corollary 1 in (Mathé, 2006). Define $\tilde{i} = \max\{i : B(i) \leq \text{CNF}(i)\}$. The proof is composed of two steps: first we show that we are competitive with \tilde{i} and then we show that \tilde{i} is competitive with the best index.

Competing with \tilde{i} . Observe that since B is monotonically increasing, and CNF is monotonically decreasing, for $i \leq \tilde{i}$ we have

$$B(i) \leq B(\tilde{i}) \leq \text{CNF}(\tilde{i}) \leq \text{CNF}(i).$$

Therefore, for $i \leq \tilde{i}$

$$\left| \hat{\theta}_i - \theta^* \right| \leq B(i) + \text{CNF}(i) \leq 2\text{CNF}(i).$$

This implies that $\theta^* \in I_i$ for all $i \leq \tilde{i}$.

As a consequence, the definition of our chosen index \hat{i} implies that $\hat{i} \geq \tilde{i}$, which in turn implies that $I_{\tilde{i}} \cap I_{\hat{i}} \neq \emptyset$. So, there exists $x \in I_{\tilde{i}} \cap I_{\hat{i}}$ such that $|x - \hat{\theta}_{\tilde{i}}| \leq 2\text{CNF}(\tilde{i})$ and $|x - \hat{\theta}_{\hat{i}}| \leq 2\text{CNF}(\hat{i})$. As we know that $\theta^* \in I_{\tilde{i}}$, we get

$$|\hat{\theta}_{\tilde{i}} - \theta^*| \leq |\hat{\theta}_{\tilde{i}} - x| + |x - \hat{\theta}_{\hat{i}}| + |\hat{\theta}_{\hat{i}} - \theta^*| \leq 2\text{CNF}(\hat{i}) + 2\text{CNF}(\tilde{i}) + 2\text{CNF}(\tilde{i}) \leq 6\text{CNF}(\tilde{i}). \quad (5)$$

Comparing \tilde{i} to i^* . Define $i^* := \operatorname{argmin}_i \{B(i) + \text{CNF}(i)\}$ which is the index we actually want to compete with in our guarantee. If we compare \tilde{i} with i^* , then by the above argument we can translate to \hat{i} . For this, we consider two cases:

If $i^* \leq \tilde{i}$, then by definition of \tilde{i} , we have

$$B(\tilde{i}) + \text{CNF}(\tilde{i}) \leq 2\text{CNF}(\tilde{i}) \leq 2\text{CNF}(i^*) \leq 2(\text{CNF}(i^*) + B(i^*)),$$

so we are a factor of 2 worse.

On the other hand, if $i^* > \tilde{i}$ then by [Assumption 2](#) and the optimality condition for \tilde{i}

$$\kappa \text{CNF}(\tilde{i}) \leq \text{CNF}(\tilde{i} + 1) \leq B(\tilde{i} + 1) \leq B(i^*).$$

This implies

$$B(\tilde{i}) + \text{CNF}(\tilde{i}) \leq (1 + 1/\kappa)B(i^*).$$

As $\kappa \leq 1$, this bound dominates the previous case, and together, with (5) we have

$$\left| \hat{\theta}_{\tilde{i}} - \theta^* \right| \leq 6 \times \text{CNF}(\tilde{i}) \leq 6(1 + 1/\kappa) \min_{i \in [M]} \{B(i) + \text{CNF}(i)\}. \quad \square$$

Proof of Corollary 4. For the MSE calculation, we simply need to translate from the high probability guarantee to the MSE, which is not difficult under [Assumption 1](#). In particular, fix δ and let \mathcal{E} be the event that all confidence bounds are valid, which holds with probability $1 - \delta$, then we have

$$\begin{aligned} \mathbb{E}(\hat{\theta}_{\tilde{i}} - \theta^*)^2 &= \mathbb{E}(\hat{\theta}_{\tilde{i}} - \theta^*)^2 \mathbf{1}\{\mathcal{E}\} + \mathbb{E}(\hat{\theta}_{\tilde{i}} - \theta^*)^2 \mathbf{1}\{\bar{\mathcal{E}}\} \\ &\leq \mathbb{E}(\hat{\theta}_{\tilde{i}} - \theta^*)^2 \mathbf{1}\{\mathcal{E}\} + R^2 \delta \\ &\leq \mathbb{E} \mathbf{1}\{\mathcal{E}\} \left(6(1 + 1/\kappa) \mathbb{E} \min_{i \in [M]} \{B(i) + \text{CNF}(i; \delta)\} \right)^2 + R^2 \delta \\ &\leq 72(1 + 1/\kappa)^2 \min_{i \in [M]} \{B(i)^2 + \text{CNF}(i; \delta)^2\} + R^2 \delta. \end{aligned}$$

Here in the first line we are introducing the event \mathcal{E} and its complement. In the second, we use that $\hat{\theta}_i, \theta^* \in [0, R]$ almost surely and that $\mathbb{P}[\bar{\mathcal{E}}] \leq \delta$ according to [Assumption 1](#). In the third line, we apply [Theorem 3](#), which holds under event \mathcal{E} . The final step uses the simplification that $(a + b)^2 \leq 2a^2 + 2b^2$. \square

A.2. Proofs for Section 4

Proof of Theorem 5. Let us first verify that the confidence function specified in (2) satisfy [Assumption 1](#). We will apply Bernstein's inequality, which requires variance and range bounds. For the variance, a single sample satisfies

$$\text{Var} \left(\frac{K(|\pi_T(x) - a|/h)r(a)}{h \times \pi_L(a | x)} \right) \leq \frac{1}{h} \mathbb{E} \left[\frac{K(|\pi_T(x) - a|/h)^2}{h \times \pi_L^2(a | x)} \mid a \sim \pi_L(\cdot | x) \right] \leq \frac{1}{2h},$$

where we first use that the variance is upper bounded by the second moment, and then we use that π_L is uniform and the boxcar kernel is at most 1/2. Finally, we use that by a change of variables $K(\cdot/h)/h$ integrates to 1. Note that we are using that $\pi_T \in [\gamma_0, 1 - \gamma_0]$, as we are integrating over the support of π_L .

For the range, we have

$$\sup \frac{K(|\pi_T(x) - a|/h)r(a)}{h \times \pi_L(a | x)} \leq \frac{1}{2h}.$$

Therefore, Bernstein's inequality gives that with probability $1 - \delta$ we have

$$\left| \hat{V}_h(\pi_T) - \mathbb{E} \hat{V}_h(\pi_T) \right| \leq \sqrt{\frac{\log(2/\delta)}{nh}} + \frac{\log(2/\delta)}{3nh},$$

and the first claim follows by a union bound.

Monotonicity is also easy to verify with this definition of CNF. In particular, since $h_i = \gamma h_{i+1}$ and $\gamma < 1$, we immediately have that

$$\gamma^{\text{CNF}(i)} = \gamma \sqrt{\frac{\log(2M/\delta)}{nh_i}} + \gamma \frac{\log(2M/\delta)}{3nh_i} = \gamma \sqrt{\frac{\log(2M/\delta)}{n\gamma h_{i+1}}} + \frac{\log(2M/\delta)}{3nh_{i+1}} \leq \text{CNF}(i+1)$$

Clearly $\text{CNF}(i+1) \leq \text{CNF}(i)$, and so [Assumption 2](#) holds. This verifies that we may apply [Theorem 3](#).

For the last claim, if the rewards are L -Lipschitz, then we claim we can set $B(i) = Lh_i$. To see why, observe that

$$\begin{aligned} \left| \mathbb{E} V_h(\pi_T) - V(\pi_T) \right| &= \left| \mathbb{E}_{(x,r), a \sim \pi_L(x)} \frac{K(|\pi_T(x) - a|/h)r(a)}{h} - r(\pi_T(x)) \right| \\ &= \left| \mathbb{E}_{(x,r)} \int_{a'} \frac{\mathbf{1}\{|\pi_T(x) - a| \leq h\} r(a)}{2h} - r(\pi_T(x)) \right| \\ &= \left| \mathbb{E}_{(x,r)} \int_{a'} \frac{\mathbf{1}\{|\pi_T(x) - a| \leq h\} (r(a) - r(\pi_T(x)))}{2h} \right| \leq Lh. \end{aligned}$$

Clearly this bias bound is monotonic. To apply [Theorem 3](#), it is better to first simplify the confidence function. Observe that as $\theta^*, \hat{\theta}_i \in [0, 1]$, it is always better to clip the estimates $\hat{\theta}_i$ to lie in $[0, 1]$. This has no bearing on the bias and only improves the deviation term, and in particular allows us to replace $\text{CNF}(i)$ with $\min\{\text{CNF}(i), 1\}$. This leads to a further simplification:

$$\sqrt{\frac{\log(2M/\delta)}{nh_i}} + \frac{\log(2M/\delta)}{3nh_i} \leq 1 \Rightarrow \sqrt{\frac{\log(2M/\delta)}{nh_i}} + \frac{\log(2M/\delta)}{3nh_i} \leq \frac{4}{3} \sqrt{\frac{\log(2M/\delta)}{nh_i}}$$

Therefore we may replace CNF with this latter function and by [Theorem 3](#) we guarantee that with probability at least $1 - \delta$

$$\left| \hat{\theta}_i - \theta^* \right| \leq 6(1 + \gamma^{-1}) \min_i \left\{ Lh_i + \frac{4}{3} \sqrt{\frac{\log(2M/\delta)}{nh_i}} \right\}.$$

The optimal choice for h is $h^* := \left(\frac{4}{3L} \sqrt{\frac{\log(2M/\delta)}{n}} \right)^{2/3}$, which will in general not be in our set \mathcal{H} . However, if we use this choice for h , the error rate is $O((L/n)^{1/3})$, and since we know that $\theta^* \in [0, 1]$, if $L > n$ then this error guarantee is trivial.

In other words, the maximum value of L that we are interested in adapting to is $L_{\max} = n$. This will be useful in setting the number of models to search over M .

To set M , we want to ensure that there exists some h_i such that $h_i \leq h^* \leq h_{i+1}$. We first verify the first inequality, which requires that

$$\gamma_0 \gamma^M \leq h^* := \left(\frac{4}{3L} \sqrt{\frac{\log(2M/\delta)}{n}} \right)^{2/3}$$

We will always take $M \geq 2$, which implies that $\log(2M/\delta) \geq 1$. Then, since we are only interested in $L \leq n$, a sufficient condition here is

$$\gamma_0 \gamma^M \leq \frac{4^{2/3}}{3} n^{-1} \Rightarrow M \geq \frac{C_{\gamma_0} \log(n)}{\log(1/\gamma_0)},$$

where C_{γ_0} is a constant that only depends on γ_0 . The upper bound $h^* \leq h_{i+1}$ is satisfied as soon as n is large enough, provided that $L \geq \omega(\sqrt{\log(\log(n))/n})$, which we are assuming. Thus we know that there is i^* such that $h_{i^*} \leq h^* \leq h_{i^*}/\gamma$, and using this choice, we have

$$\begin{aligned} |\hat{\theta}_i - \theta^*| &\leq 6(1 + \gamma^{-1}) \left\{ Lh_{i^*} + \frac{4}{3} \sqrt{\frac{\log(2M/\delta)}{nh_{i^*}}} \right\} \leq 6(1 + \gamma^{-1}) \left\{ Lh^* + \frac{4}{3} \sqrt{\frac{\log(2M/\delta)}{\gamma nh^*}} \right\} \\ &\leq 6(1 + \gamma^{-1}) \cdot \frac{c_1}{\sqrt{\gamma}} (L \log(2M/\delta)/n)^{1/3} \leq C_{\gamma, \gamma_0} (L \log(\log(n)/\delta)/n)^{1/3}, \end{aligned}$$

where C_{γ, γ_0} is a constant that depends only on γ, γ_0 . □

Note that if π_L is non-uniform, but satisfies $\inf_{x,a} \pi_L(a | x) \geq p_{\min}$, then very similar arguments apply. In particular, we have that both variance and range are bounded by $\frac{1}{2hp_{\min}}$, and some Bernstein's inequality in this case yields

$$\left| \hat{V}_h(\pi_T) - \mathbb{E} \hat{V}_h(\pi_T) \right| \leq \sqrt{\frac{\log(2/\delta)}{nhp_{\min}}} + \frac{\log(2/\delta)}{3nhp_{\min}},$$

Monotonicity follows from the same calculation as before and the clipping trick yields a more interpretable final bound, which holds with probability at least $1 - \delta$, of

$$\left| \hat{\theta}_i - \theta^* \right| \leq 6(1 + \gamma^{-1}) \min_i \left\{ Lh_i + \frac{4}{3} \sqrt{\frac{\log(2M/\delta)}{nh_i p_{\min}}} \right\}.$$

The remaining calculation for M is analogous, since this bound is identical to the previous one with n replaced by np_{\min} . Thus, we obtain a final bound of $C_{\gamma, \gamma_0} (L \log(\log(n/\delta))/(np_{\min}))^{1/3}$.

A.3. Proofs for Section 5

We first develop and state the more precise version of [Theorem 6](#). We introduce the doubly robust version of the partial importance weighting estimator. As it is the empirical average over n trajectories, here we will focus on a single trajectory $(x_1, a_1, r_1, \dots, x_H, a_H, r_H)$ sampled by following the logging policy π_L .

Define $\hat{V}_{\text{DR}}^0 := 0$ and

$$\hat{V}_{\text{DR}}^{H+1-h} := \hat{V}(x_h) + p_h(r_h + \gamma \hat{V}_{\text{DR}}^{H-t} - \hat{Q}(x_h, a_h)), \quad p_h := \frac{\pi_T(a_h | x_h)}{\pi_L(a_h | x_h)}.$$

where \hat{Q} is the direct model, trained via supervised learning, and $\hat{V}(x) = \mathbb{E}_{a \sim \pi_T(x)} \hat{Q}(x, a)$. The full horizon doubly-robust estimator is $\hat{V}_{\text{DR}} := \hat{V}_{\text{DR}}^H$. To define the η -step partial estimator, let $\hat{V}_{\text{DM}}^\eta := \rho_\eta \hat{Q}(x_{\eta+1}, \pi_T(x_{\eta+1}))$, which is an estimate

of $\mathbb{E}_{\pi_T} [V(x_{\eta+1})]$. Set $\hat{V}_{\text{DM}}^H := 0$. Then for a false horizon η , we define a similar recursion

$$\hat{V}_\eta^{H+1-h} := \begin{cases} \hat{V}(x_h) + p_h(r_h + \gamma V_\eta^{H-h} - \hat{Q}(x_h, a_h)) & \text{if } 1 \leq h < \eta \\ \hat{V}(x_h) + p_h(r_h + \gamma \hat{V}_{\text{DM}}^\eta - \hat{Q}(x_h, a_h)) & \text{if } h = \eta. \end{cases}$$

The doubly robust variant of the η -step partial importance weighted estimator is \hat{V}_η^H . We also define $\hat{V}_0^H = \hat{V}_{\text{DM}}^0$ which estimates $\mathbb{E}[V(x_1)]$. Observe that if in the definition of \hat{V}_{DR} , we take $\hat{V}, \hat{Q} \equiv 0$ then we obtain the estimator in (4).

Define $\Delta := \log(2(H+1)/\delta)$, $V_{\max} := (1-\gamma)^{-1}$ and recall that $p_{\max} := \max_{x,a} \frac{\pi_T(a|x)}{\pi_L(a|x)}$. Then define

$$\begin{aligned} \mathbf{B}(i) &:= \frac{\gamma^{H-i+1} - 1}{1 - \gamma} \\ \text{CNF}(i) &:= \sqrt{\frac{6V_{\max}^2 \left(1 + \sum_{h=1}^{H-i+1} \gamma^{2(h-1)} p_{\max}^h\right) \Delta}{n}} + \frac{6V_{\max} \left(1 + \sum_{h=1}^{H-i+1} \gamma^{h-1} p_{\max}^h\right) \Delta}{3n} \end{aligned}$$

With these definitions, we now state the theorem

Theorem 7 (Formal version of [Theorem 6](#)). *In the episodic reinforcement learning setting with discount factor γ , consider the doubly robust partial importance weighting estimators $\hat{\theta}_i := \hat{V}_{H-i+1}^H(\pi_T)$ for $i \in \{1, \dots, H+1\}$. Then \mathbf{B} and CNF are valid and monotone, with $\kappa := (1 + \gamma p_{\max})^{-1}$.*

Proof of Theorem 7. We now turn to the proof.

Bias analysis. By repeatedly applying the tower property, the expectation for \hat{V}_η^H is

$$\begin{aligned} \mathbb{E} \left[\hat{V}_\eta^H \right] &= \mathbb{E}_{\pi_L} \left[\hat{V}(x_1) + p_1(r_1 + \gamma \hat{V}_\eta^{H-1} - \hat{Q}(x_1, a_1)) \right] \\ &= \mathbb{E}_{x_1} \left[\hat{V}(x_1) + \mathbb{E}_{a_1 \sim \pi_L(x_1), a_2:H \sim \pi_L} \left[p_1(r_1 + \gamma \hat{V}_\eta^{H-1} - \hat{Q}(x_1, a_1)) \mid x_1 \right] \right] \\ &= \mathbb{E}_{x_1} \left[\hat{V}(x_1) + \mathbb{E}_{a_1 \sim \pi_T(x_1), a_2:H \sim \pi_L} \left[r_1 + \gamma \hat{V}_\eta^{H-1} - \hat{Q}(x_1, a_1) \mid x_1 \right] \right] \\ &= \mathbb{E}_{x_1, a_1 \sim \pi_T(x_1)} [r] + \gamma \mathbb{E}_{x_2 \sim \pi_T, a_2:H \sim \pi_L} \left[\hat{V}_\eta^{H-1} \right] \\ &= \dots \\ &= \mathbb{E}_{\pi_T} \left[\sum_{h=1}^{\eta} \gamma^{h-1} r \right] + \gamma^\eta \mathbb{E}_{x_{\eta+1} \sim \pi_T} \left[\hat{V}_{\text{DM}}^\eta \right]. \end{aligned}$$

Here, we use that p_1 is the one-step importance weight, so it changes the action distribution from π_L to π_T . We also use the relationship between the direct models \hat{Q} and \hat{V} . Therefore, the bias is

$$\left| \mathbb{E} \left[\hat{V}_\eta^H \right] - V(\pi_T) \right| = \gamma^\eta \left| \hat{V}_{\text{DM}}^\eta - \mathbb{E}_{\pi_T} [V(x_{\eta+1})] \right| \leq \frac{\gamma^\eta - \gamma^H}{1 - \gamma} =: \mathbf{B}(H - \eta + 1) \quad (6)$$

The first identity justifies our choice of \hat{V}_{DM}^η which attempts to minimize this bias using the direct model. The inequality here follows from the fact that rewards are in $[0, 1]$, which implies that values at time $\eta+1$ are in $\left[0, \frac{1-\gamma^{H-\eta}}{1-\gamma}\right]$. As $\gamma \in (0, 1)$, clearly we have that $\mathbf{B}(i)$ is monotonically increasing with i increasing. Thus this bias bound is valid.

Variance analysis. For the variance calculation, let $\mathbb{E}_h[\cdot], \text{Var}_h(\cdot)$ denote expectation and variance conditional on all randomness *before* time step h . Adapting Theorem 1 of [Jiang & Li \(2016\)](#) the variance for $1 \leq h < \eta$ is given by the recursive formula:

$$\text{Var}_h(\hat{V}_\eta^{H+1-h}) = \text{Var}_h(\mathbb{E} \left[\hat{V}_\eta^{H+1-h} \mid x_h \right]) + \mathbb{E}_h \left[\text{Var}(p_h \Delta(x_h, a_h) \mid x_h) \right] + \mathbb{E}_h \left[p_h^2 \text{Var}(r_h) \right] + \mathbb{E}_h \left[\gamma^2 p_h^2 \text{Var}_{h+1}(\hat{V}_\eta^{H-h}) \right],$$

where $\Delta(x_h, a_h) := \hat{Q}(x_h, a_h) - Q(x_h, a_h)$. For $h = \eta$ it is identical, except that in the last term we use \hat{V}_{DM}^η instead of $\hat{V}_\eta^{H-\eta}$ (which is not defined).

Unrolling the recursion, the full expression for the variance is

$$\begin{aligned} \text{Var}(\hat{V}_\eta^H) &= \sum_{h=1}^{\eta} \mathbb{E} \left[\gamma^{2(h-1)} \rho_{h-1}^2 \text{Var}_h(\mathbb{E}[\hat{V}_\eta^{H+1-h} | x_h]) \right] \\ &\quad + \sum_{h=1}^{\eta} \mathbb{E} \left[\gamma^{2(h-1)} \rho_{h-1}^2 \mathbb{E}_h [\text{Var}(p_h \Delta(x_h, a_h) | x_h)] \right] \\ &\quad + \sum_{h=1}^{\eta} \mathbb{E} \left[\gamma^{2(h-1)} \rho_{h-1}^2 \mathbb{E}_h [p_h^2 \text{Var}(r_h)] \right] \\ &\quad + \mathbb{E} \left[\gamma^{2\eta} \rho_\eta^2 \text{Var}(\hat{V}_{\text{DM}}^\eta) \right]. \end{aligned}$$

For the variance bound, we do not attempt to obtain the sharpest bound possible. Instead, we use the following facts: (1) rewards are in $[0, 1]$, (2) all values, value estimates, and Q are at most $(1 - \gamma)^{-1} =: V_{\max}$, and (3) for a random variable X that is bounded by B almost surely, we have $\text{Var}(X) \leq B^2$. Using these facts in each term gives

$$\begin{aligned} \text{Var}(\hat{V}_\eta^H) &\leq \sum_{h=1}^{\eta} \mathbb{E} \left[\gamma^{2(h-1)} \rho_{h-1}^2 V_{\max}^2 \right] + \sum_{h=1}^{\eta} \mathbb{E} \left[\gamma^{2(h-1)} \rho_h^2 V_{\max}^2 \right] + \sum_{h=1}^{\eta} \mathbb{E} \left[\gamma^{2(h-1)} \rho_h^2 \right] + \mathbb{E} \left[\gamma^{2\eta} \rho_\eta^2 V_{\max}^2 \right] \\ &= \sum_{h=1}^{\eta+1} \mathbb{E} \left[\gamma^{2(h-1)} \rho_{h-1}^2 V_{\max}^2 \right] + \sum_{h=1}^{\eta} \mathbb{E} \left[\gamma^{2(h-1)} \rho_h^2 (V_{\max}^2 + 1) \right] \\ &\leq 3V_{\max}^2 \sum_{h=1}^{\eta} \mathbb{E} \left[\gamma^{2(h-1)} \rho_h^2 \right] + V_{\max}^2 \end{aligned}$$

Here in the first line we use the three facts we stated above. In the second line we collect the terms. In the third line we note that $\gamma^{2(h-1)} \rho_{h-1}^2 \leq \gamma^{2(h-2)} \rho_{h-1}^2$ since $\gamma \in (0, 1)$, so we can re-index the first summation and group terms again.

To simplify further, let $p_{\max} := \sup_{x,a} \frac{\pi_{\text{T}}(a|x)}{\pi_{\text{L}}(a|x)}$ denote the largest importance weight and note that as $\mathbb{E}_h[p_h] = 1$, we have

$$\sum_{h=1}^{\eta} \mathbb{E} \left[\gamma^{2(h-1)} \rho_h^2 \right] \leq \sum_{h=1}^{\eta} \mathbb{E} \left[\gamma^{2(h-1)} p_{\max} \rho_{h-1}^2 \mathbb{E}_h[w_h] \right] \leq \dots \leq \sum_{h=1}^{\eta} \gamma^{2(h-1)} p_{\max}^h.$$

Therefore, our variance bound will be

$$\text{Var}(\hat{V}_\eta^H) \leq 3V_{\max}^2 \left(1 + \sum_{h=1}^{\eta} \gamma^{2(h-1)} p_{\max}^h \right).$$

For the range, we obtain the recursion (for $1 \leq h < \eta$):

$$\left| \hat{V}_\eta^{H+1-h} \right| \leq V_{\max} + p_{\max}(1 + V_{\max}) + p_{\max}\gamma \left| \hat{V}_\eta^{H-h} \right|,$$

with the terminal condition $\left| \hat{V}_{\text{DM}}^\eta \right| \leq V_{\max}$. A somewhat crude upper bound is

$$\left| \hat{V}_\eta^H \right| \leq 3V_{\max} \left(1 + \sum_{h=1}^{\eta} \gamma^{h-1} p_{\max}^h \right),$$

which has a similar form to the variance expression.

Therefore, Bernstein’s inequality reveals that with probability $1 - \delta$, we have that the n -trajectory empirical averages satisfy

$$\left| \hat{V}_\eta^H - \mathbb{E} \hat{V}_\eta^H \right| \leq \sqrt{\frac{6V_{\max}^2 \left(1 + \sum_{h=1}^\eta \gamma^{2(h-1)} p_{\max}^h\right) \log(2/\delta)}{n}} + \frac{6V_{\max} \left(1 + \sum_{h=1}^\eta \gamma^{h-1} p_{\max}^h\right) \log(2/\delta)}{3n}. \quad (7)$$

This bound is clearly seen to be monotonically increasing in η , which is monotonically decreasing with i as required. The reason is that when we increase η we add one additional non-negative term to both the variance and range expressions.

Finally, we must verify that the bound does not decrease too quickly. For this, we first verify the following elementary fact

Fact 8. *Let $z \geq 0$ and $t \geq 0$ then*

$$\frac{1 + \sum_{\tau=1}^t z^\tau}{1 + \sum_{\tau=1}^{t-1} z^\tau} \leq 1 + z.$$

Proof. Using the geometric series formula, we can rewrite

$$\frac{1 + \sum_{\tau=1}^t z^\tau}{1 + \sum_{\tau=1}^{t-1} z^\tau} = 1 + \frac{z^t}{1 + \sum_{\tau=1}^{t-1} z^\tau} \leq 1 + z. \quad \square$$

Using the above fact, we can see that the variance bound decreases at rate $(1 + \gamma^2 p_{\max})$ and the range bound decreases at rate $(1 + \gamma p_{\max})$. The range bound dominates here, since

$$\sqrt{1 + \gamma^2 p_{\max}} \leq \sqrt{1 + \gamma^2 p_{\max}^2} \leq 1 + \gamma p_{\max}$$

Therefore, we may take the decay constant to be $1/(1 + \gamma p_{\max})$ to verify [Assumption 2](#). □

B. Details for continuous contextual bandits experiments

B.1. The simulation environment.

Here, we explain some of the important details of the simulation environment. The simulator is initialized with a d_x dimensional context space and action space $[0, 1]^d$ for some parameter d . For our experiments we simply take $d = 1$. There is also a hidden parameter matrix $\beta^* \sim \mathcal{N}(0, I)$ with $\beta^* \in \mathbb{R}^{d \times d_x}$. In each round, contexts are sampled iid from $\mathcal{N}(0, I)$, then the optimal action $a^*(x) := \sigma(\beta^* x)$, where $\sigma(z) = \frac{e^z}{e^z + 1}$ is the standard sigmoid, and the function is applied component-wise. This optimal action a^* is used in the design of the reward functions.

We consider two different reward functions called “absolute value” and “quadratic.” The first is simply $\ell(a) := 1 - \min(L \|a - a^*(x)\|_1, 1)$, while the latter is $\ell(a) := 1 - \min(L/4 \sum_{j=1}^d (a_j - a_j^*(x))^2, 1)$. Here L is the Lipschitz constant, which is also a configurable.

For policies, the uniform logging policy simply chooses $a \sim \text{Unif}([0, 1]^d)$ on each round. Other logging and target policies are trained via regression on 10 vector-valued regression samples $(x, a^*(x) + \mathcal{N}(0, 0.5 \cdot I))$ where $x \sim \mathcal{N}(0, I)$. We use two different regression models: linear + sigmoid implemented in PyTorch, and a decision tree implemented in scikit-learn. Both regression procedures yield deterministic policies, and in our experiments we take this policies to be π_T .

For π_L we implement two softening techniques following [Farajtabar et al. \(2018\)](#), called “friendly” and “adversarial,” and both techniques take two parameters α, β . Both methods are defined for discrete action spaces, and to adapt to the continuous setting we partition the continuous action space into m bins (for one-dimensional spaces). We round the deterministic action chosen by the regression model to its associated bin, run the softening procedure to choose a (potentially different) bin, and then sample an action uniformly from this bin. For higher dimensional action spaces, we discretize each dimension individually, so the softening results in a product measure.

Friendly softening with discrete actions is implemented as follows. We sample $U \sim \text{Unif}([-0.5, 0.5])$ and then the updated action is $\pi_{\text{det, disc}}(x)$ with probability $\alpha + \beta U$ and it is uniform over the remaining discrete actions with the remaining probability. Here $\pi_{\text{det, disc}}$ is the deterministic policy obtained by the regression model, discretized to one of the m bins. Adversarial softening instead is uniform over all discrete actions with probability $1 - (\alpha + \beta U)$ and it is uniform over all

but $\pi_{\text{det,disc}}(x)$ with the remaining probability. In both cases, once we have a discrete action, we sample a continuous action from the corresponding bin.

The simulator also supports two different kernel functions: Epanechnikov and boxcar. The boxcar kernel is given by $K(u) = \frac{1}{2}\mathbf{1}\{|u| \leq 1\}$, while Epanechnikov is $K(u) = 0.75 \cdot (1 - u^2)\mathbf{1}\{|u| \leq 1\}$. We address boundary bias by normalizing the kernel appropriately, as opposed to forcing the target policy to choose actions in the interior. This issue is also discussed in Kallus & Zhou (2018).

Finally, we also vary the number of logged samples and the Lipschitz constant of the loss functions.

B.2. Reproducibility Checklist

Data collection process. All data are synthetically generated as described above.

Dataset and Simulation Environment. We will make the simulation environment publicly available.

Excluded Data. No excluded data.

Training/Validation/Testing allocation. There is no training/validation/testing setup in off policy evaluation. Instead all logged data are used for evaluation.

Hyper-parameters. Hyperparameters used in the experimental conditions are: $n \in 10^{1:5}$, $h \in \{2^{-(1:7)}\}$, $L \in \{0.1, 0.3, 1, 3, 10\}$, in addition to the other configurable parameters (e.g., softening technique, kernel, logging policy, target policy).

Evaluation runs. There are 1000 conditions, each with 30 replicates with different random seeds.

Description of experiments. For each condition, determined by logging policy, softening technique, target policy, sample size, lipschitz constant, reward function, and kernel type, we generate n logged samples following π_L , and 100k samples from π_T to estimate the ground truth $V(\pi_T)$. All fixed-bandwidth estimators and SLOPE are calculated based on the same logged data. The MSE is estimated by averaging across the 30 replicates, each with different random seed.

For the learning curve in the right panel of Figure 3 the specific condition shown is: uniform logging policy, linear+sigmoid target policy, $L = 3$, absolute value reward, boxcar kernel. MSE estimates are measured at $n = \{1, 3, 7\} \times 10^{1:3} \cup \{10, 000\}$. We perform 100 replicates for this experiment.

Measure and Statistics. Results are shown in Figure 3. Statistics are based on empirical CDF calculated by aggregating the 1000 conditions. Typically there are no error bars for such plots. Pairwise comparison is based on paired t -test over all pair of methods and conditions, with significance level 0.05. The learning curve is based on 100 replicates, with error bar corresponding to ± 2 standard errors shown in the plots.

Computing infrastructure. Experiments were run on Microsoft Azure.

C. Details for reinforcement learning experiments

C.1. Experiment Details

Environment Description. We provide brief environment description below. More details can be found in Thomas & Brunskill (2016); Voloshin et al. (2019); Brockman et al. (2016).

- Mountain car is a classical benchmark from OpenAI Gym. We make the same modification as Voloshin et al. (2019). The domain has 2-dimensional state space (position and velocity) and one-dimensional action {left,nothing,right}. The reward is $r = -1$ for each timestep before reaching the goal. The initial state has position uniformly distributed in the discrete set $\{-0.4, -0.5, -0.6\}$ with velocity 0. The horizon is set to be $H = 250$ and there is an absorbing state at $(0.5, 0)$. The domain has deterministic dynamics, as well as deterministic, dense reward.
- Graph and Graph-POMDP are adopted from Voloshin et al. (2019). The Graph domain has horizon 16, state space $\{0, 1, 2, \dots, 31\}$ and action space $\{0, 1\}$. The initial state is $x_0 = 0$, and we have the state-independent stochastic transition model with $\mathbb{P}(x_{t+1} = 2t + 1|a = 0) = 0.75$, $\mathbb{P}(x_{t+1} = 2t + 2|a = 0) = 0.25$, $\mathbb{P}(x_{t+1} = 2t + 1|a = 1) = 0.25$, $\mathbb{P}(x_{t+1} = 2t + 2|a = 1) = 0.75$. In the dense reward configuration, we have $r(x_t, a_t, x_{t+1}) = 2(x_{t+1} \bmod 2) - 1 \forall t \leq T$. The sparse reward setting has $r(x_t, a_t, x_{t+1}) = 0 \forall t < T - 1$ with reward only at the last time

step, according to the dense reward function. We also consider a stochastic reward setting, where we change the reward to be $r(x_t, a_t, x_{t+1}) \sim \mathcal{N}(2(x_{t+1} \bmod 2) - 1, 1)$. Graph-POMDP is a modified version of Graph where states are grouped into 6 groups. Only the group information is observed, so the states are aliased.

- Gridworld is also from Voloshin et al. (2019). The state space is an 8×8 grid with four actions [up, down, left, right]. The initial state distribution is uniform over the left column and top row, while the goal is in the bottom right corner. The horizon length is 25. The states belongs to four categories: Field, Hole, Goal, Others. The reward at Field is -0.005, Hole is -0.5, Goal is 1 and Others is -0.01. The exact map can be found in Voloshin et al. (2019).
- Hybrid Domain is from Thomas & Brunskill (2016). It is a composition of two other domains from the same study, called ModelWin and ModelFail. The ModelFail domain has horizon 2, four states $\{s_0, s_1, s_2, s_a\}$ and two actions $\{0, 1\}$. The agent starts at s_0 , goes to s_1 with reward 1 if $a = 0$, and goes to s_2 with reward if $a = 1$. Then it transitions to the absorbing state s_a . This environment has partial observability so that $\{s_0, s_1, s_2\}$ are aliased together. In the hybrid domain the absorbing state s_a is replaced with a new state s_1 in the ModelWin domain. This domain has four states $\{s_1, s_2, s_3, s_a\}$. The action space is $\{0, 1\}$. The agent starts from s_1 . Upon taking action $a = 0$, it goes to s_2 with probability 0.6 and receives reward 1, and goes to s_3 with probability 0.4 and reward -1. If $a = 1$, it does the opposite. From s_2 and s_3 the agent deterministically transitions to s_1 with 0 reward. do a deterministic transition back to s_1 with 0 reward. The horizon here is 20 and $x_{20} = s_a$. The states are fully observable.

Models. Instead of experiment with all possible approaches for direct modeling, which is quite burdensome, we follow the high-level guidelines provided in Table 3 of Voloshin et al. (2019)’s paper: for Graph, PO-Graph, and Mountain Car we use FQE because these environments are stochastic and have severe mismatch between logging and target policy. In contrast, Gridworld has moderate policy mismatch, so we use $Q^\pi(\lambda)$. For the Hybrid domain, we use a simple maximum-likelihood approximate model to predict the full transition operator and rewards, and plan in the model to estimate the value function.

Policy. For Gridworld and Mountain Car, we use ϵ -Greedy polices as logging and target policies. To derive these, we first train a base policy using value iteration and then we take $\pi(a^*|x) = 1 - \epsilon$ and $\pi(a | x) = \epsilon/(|\mathcal{A}| - 1)$ for $a \neq a^*(x)$, where $a^*(x) = \operatorname{argmax} \hat{Q}(x, a)$ for the learned \hat{Q} function. In Gridworld, we take the following policy pairs: $[(1, 0.1), (0.6, 0.1), (0.2, 0.1), (0.1, 0.2), (0.1, 0.6)]$, where the first argument is the ϵ parameter for . For Mountain Car domain, we take the following policy pairs: $[(0.1, 0), (1, 0), (1, 0.1), (0.1, 1)]$ where the first argument is the parameter for π_L and the second is for π_T . For the Graph and Graph-POMDP domain, both logging and target policies are static polices with probability p going left (marked as $a = 0$) and probability $1 - p$ going right (marked as $a = 1$), i.e., $\pi(a = 0|x) = p$ and $\pi(a = 1|x) = 1 - p \forall x$. In both environments, we vary p of the logging policy to be 0.2 and 0.6, while setting p for target policy to be 0.8. For the Hybrid domain, we use the same policy as Thomas & Brunskill (2016). For the first ModelFail part, $\pi_L(a = 0) = 0.88$ and $\pi_L(a = 1) = 0.12$, while the target policy does the opposite. For the second ModelWin part, $\pi_L(a = 0|s_1) = 0.73$ and $\pi_L(a = 1|s_1) = 0.27$, and the target policy does the opposite. For both policies, they select actions uniformly when $s \in \{s_2, s_3\}$.

Other parameters. For both the Graph and Graph-POMDP, we use $\gamma = 0.98$ and $N \in 2^{7:10}$. For Gridworld, $\gamma = 0.99$ and $N \in 2^{7:9}$. For Mountain Car, $\gamma = 0.96$ and $N \in 2^{8:10}$. For Hybrid, $\gamma = 0.99$ and $N \in \{10, 20, 50, \dots, 10000, 20000, 50000\}$. Each condition is averaged over 100 replicates.

C.2. Reproducibility Checklist

Data collection process. All data are synthetically generated as described above.

Dataset and Simulation Environment. The Mountain Car environment is downloadable from OpenAI (Brockman et al., 2016). Graph, Graph-POMDP, Gridworld, and the Hybrid domain are available at <https://github.com/clvoloshin/OPE-tools>, which is the supporting code for Voloshin et al. (2019).

Excluded Data. No excluded data.

Training/Validation/Testing allocation. There is no training/validation/testing setup in off policy evaluation. Instead all logged data are used for evaluation.

Hyper-parameters. Hyperparameters (apart from those optimized by SLOPE) are optimized following the guidelines of Voloshin et al. (2019). For MountainCar, the direct model is trained using a 2-layer fully connected neural network with hidden units 64 and 32. The batch size is 32 and convergence is set to be $1e - 4$, network weights are initialized with

truncated Normal(0, 0.1). For tabular models, convergence of Graph and Graph-POMDP is $1e - 5$ and Gridworld is $4e - 4$.

Evaluation runs. All conditions have 100 replicates with different random seeds.

Description of experiments. For each condition, determined by the choice of environment, stochastic/deterministic reward, sparse/dense reward, stochastic/deterministic transition model, logging policy π_L , target policy π_T and sample size N . We generate N logged trajectories following π_L , and 10000 samples from π_T to compute the ground truth $V(\pi_T)$. All baselines and SLOPE are calculated based on the same logged data. The MSE is estimated by averaging across the 100 replicates, each with different random seed.

Measure and Statistics. Results are shown in [Figure 4](#). Statistics are based on empirical CDF calculated by aggregating the 106 conditions. Typically there are no error bars in ECDF plots. Pairwise comparison is based on paired t -test over all pair of methods over all conditions. Each test has significance level 0.05. Learning curve is based on Hybrid domain with 128 replicates, with error bar corresponding to ± 2 standard errors shown in the plots.

Computing infrastructure. RL experiments were conducted in a Linux compute cluster.