



Towards a Sharp Analysis of Offline Policy Learning for f -Divergence-Regularized Contextual Bandits

Qingyue Zhao^{*1}, Kaixuan Ji^{*1}, Heyang Zhao^{*1}, Tong Zhang² and Quanquan Gu¹

¹University of California, Los Angeles

²University of Illinois Urbana-Champaign



Data Model: Offline Contextual Bandits

- ▶ $(s, a) \in \mathcal{S} \times \mathcal{A}$: input-output (e.g., prompt-response)
- ▶ $\rho \in \Delta(\mathcal{S})$: input distribution
- ▶ $\pi^{\text{ref}} \in \Delta(\mathcal{A}|\mathcal{S})$: reference model
- ▶ $\mathcal{G} \subset (\mathcal{S} \times \mathcal{A} \rightarrow [0, 1])$: known function class
 - ▷ $\forall \epsilon_c > 0, \mathcal{N}_{\mathcal{G}}(\epsilon_c) = \text{poly}(\epsilon_c^{-1})$: low complexity (i.e., nice covering number)
- ▶ $r \in \mathcal{G}$: **realizable** mean reward function
- ▶ $\mathcal{D} := \{(s_i, a_i, r_i)\}_{i=1}^n$: i.i.d. prompt-response-reward dataset
 - ▷ $s_i \sim \rho, a_i \sim \pi^{\text{ref}}(\cdot|s_i)$
 - ▷ $r_i = r(s_i, a_i) + \varepsilon_i$: reward with subgaussian noise
 - ▶ See also the pairwise-comparison (dueling) counterpart in our paper!

Targets of Interest

- ▶ **Objective** $J(\pi) := \mathbb{E}_{s \sim \rho} [\langle r(s, \cdot), \pi(\cdot|s) \rangle - \eta^{-1} D_f(\pi(\cdot|s) \| \pi^{\text{ref}}(\cdot|s))]$
 - ▷ D_f : f -divergence
 - ▶ Different f 's induce different realizations, including TV, KL, ...
 - ▷ η^{-1} : regularization intensity (also known as the “temperature”)
- ▶ **Optimal Policy** $\pi^* \leftarrow \operatorname{argmax}_{\pi} J(\pi)$
 - ▷ e.g., $\pi^*(\cdot|s) \propto \pi^{\text{ref}}(\cdot|s) \exp(\eta \cdot r(s, a))$ for $D_f = \text{KL}$
- ▶ **Learners** We consider learners with general function approximation. Formally, the agent knows \mathcal{D}, \mathcal{G} , and π^{ref}
 - ▷ The agent may optionally know other hyperparameters like η
 - ▷ The agent is not allowed to further sample from ρ, π^{ref} , and $\rho \times \pi^{\text{ref}}$
- ▶ **Performance Metric** For any learner $\hat{\pi}$, we are interested in the sample complexity n for it to achieve $J(\pi^*) - J(\hat{\pi}) \leq \epsilon$
 - ▷ Results regarding $\operatorname{argmax}_{a \in \mathcal{A}} r(s, a) - \langle r(s, \cdot), \hat{\pi}(\cdot|s) \rangle \leq \epsilon$ are well-established, all of which are $\asymp \epsilon^{-2}$, intuitively because the statistical rate for reward estimation is already needs at least $\gtrsim n^{-1/2}$
 - ▷ The convergence against regularized objectives themselves are less-explored even in the i.i.d. setting

Motivation

What is the weakest data **coverage** condition for post-training to be **statistically optimal** against f -divergence-regularized objectives?

- ▶ The answer depends on the “curvature” of the regularizing divergence
- ▶ Scope: offline setting with absolute or dueling feedback

Key Results

Regularizer		2312.11456	2411.04625	Ours
KL	Upper	$d\epsilon^{-2}$	$D_{\text{sup}}^2 \eta \epsilon^{-1} \log \mathcal{N}$	$D_{\pi^*}^2 \eta \epsilon^{-1} \log \mathcal{N}$
	Lower	-	$\eta \epsilon^{-1} \log \mathcal{N}$	$C^{\pi^*} \eta \epsilon^{-1} \log \mathcal{N}$
f -divergence w/ α -strongly convex f	Upper	-	-	$\alpha^{-1} \eta \epsilon^{-1} \log \mathcal{N}$
	Lower	-	-	$\alpha^{-1} \eta \epsilon^{-1} \log \mathcal{N}$

- ▶ $D_{\mathcal{G}}^2((s, a); \pi) = \sup_{g, h \in \mathcal{G}} \frac{(g(s, a) - h(s, a))^2}{\mathbb{E}_{(s', a') \sim \rho \times \pi} [(g(s', a') - h(s', a'))^2]}$
 - ▷ $D_{\mathcal{G}}^2$ -related notions are offline counterparts of **Eluder** dimension
 - ▷ $D_{\text{sup}}^2 = \sup_{(s, a) \in \mathcal{S} \times \mathcal{A}} D_{\mathcal{G}}^2((s, a); \pi^{\text{ref}})$: D^2 -based **all-policy** concentrability
 - ▷ $D_{\pi^*}^2 = \mathbb{E}_{(s, a) \sim \rho \times \pi^*} [D_{\mathcal{G}}^2((s, a); \pi^{\text{ref}})]$: D^2 -based **single-policy** concentrability
- ▶ $D_{\pi^*}^2 \leq D_{\text{sup}}^2$ for sure
 - ▷ D_{sup}^2 can be **much larger** than $D_{\pi^*}^2$ in non-trivial cases
- ▶ $C^{\pi^*} = \sup_{(s, a) \in \mathcal{S} \times \mathcal{A}} \frac{\pi^*(a|s)}{\pi^{\text{ref}}(a|s)}$: density-ratio-based **single-policy** concentrability
 - ▷ **Limitation**: the worst-case notion C^{π^*} may not match the instance-dependent notion $D_{\pi^*}^2$ in general (sometimes $D_{\pi^*}^2 / C^{\pi^*} \approx SA$)
 - ▷ Resolution to this limitation? Check our follow-up refinement on online MABs at <https://arxiv.org/pdf/2603.02155>
- ▶ Subtlety for low temperature: as $\eta \rightarrow \infty$, regularization **vanishes**
 - ▷ In the rigorous hardness results for KL, $\eta \epsilon^{-1}$ is replaced by $\eta \epsilon^{-1} \wedge \epsilon^{-2}$

Reminder: Definition and Realizations of f -divergence

- ▶ $D_f(P \| Q) := \sum_x Q(x) f(P(x)/Q(x))$
 - ▷ Requirement: convex f with $f(1) = 0$
 - ▷ $f(0) := f(0+)$, $0f(0/0) := 0$, and $0f(1/0) = f'(+\infty)$

Realizations	TV($P \ Q$)	KL($P \ Q$)	$\chi^2(P \ Q)$
f	$ x - 1 /2$	$x \log x$	$(x - 1)^2/2$
Convexity of f	convex	strictly convex	strongly convex

- ▶ However, $\pi \mapsto \text{KL}(\pi \| \pi^{\text{ref}})$ is $\Theta(1)$ -**strongly** convex w.r.t. the TV distance

Algorithm Template for $D_f = \text{KL}$: No tricky algorithmic tweaks

Require: regularization η , reference policy π^{ref} , function class \mathcal{G} , dataset \mathcal{D}

- 1: Reward learning via ordinary least square

$$\bar{g} \in \operatorname{argmin}_{g \in \mathcal{G}} \sum_{(s_i, a_i, r_i) \in \mathcal{D}} (g(s_i, a_i) - r_i)^2$$

- 2: Let $\hat{g} \leftarrow \bar{g} - \Gamma_n$, where $\Gamma_n : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_+$ is the pessimism term

Ensure: $\hat{\pi}(a|s) \propto \pi^{\text{ref}}(a|s) \exp(\eta \cdot \hat{g}(s, a))$

Algorithm Template for D_f w/ Strongly Convex f : No pessimism!

- 1: Reward learning via ordinary least square

$$\bar{g} \in \operatorname{argmin}_{g \in \mathcal{G}} \sum_{(s_i, a_i, r_i) \in \mathcal{D}} (g(s_i, a_i) - r_i)^2$$

- 2: Compute the optimal policy under \bar{g} for $s \in \mathcal{S}$ as

$$\hat{\pi}(\cdot|s) \leftarrow \operatorname{argmax}_{\pi(\cdot|s) \in \Delta(\mathcal{A})} \langle \pi(\cdot|s), \bar{g}(s, \cdot) \rangle + \eta^{-1} D_f(\pi(\cdot|s) \| \pi^{\text{ref}}(\cdot|s))$$

Ensure: $\hat{\pi}$

D_f w/ s.c. f provides stronger regularization than KL, $\hat{\pi}$ won't deviate from π^{ref} too much \implies no pess. needed