
Subgaussian and Differentiable Importance Sampling for Off-Policy Evaluation and Learning

Alberto Maria Metelli
DEIB, Politecnico di Milano
albertomaria.metelli@polimi.it

Alessio Russo
DEIB, Politecnico di Milano
alessio.russo@polimi.it

Marcello Restelli
DEIB, Politecnico di Milano
marcello.restelli@polimi.it

Abstract

Importance Sampling (IS) is a widely used building block for a large variety of off-policy estimation and learning algorithms. However, empirical and theoretical studies have progressively shown that vanilla IS leads to poor estimations whenever the behavioral and target policies are too dissimilar. In this paper, we analyze the theoretical properties of the IS estimator by deriving a novel anticoncentration bound that formalizes the intuition behind its undesired behavior. Then, we propose a new class of IS transformations, based on the notion of power mean. To the best of our knowledge, the resulting estimator is the first to achieve, under certain conditions, two key properties: (i) it displays a subgaussian concentration rate; (ii) it preserves the differentiability in the target distribution. Finally, we provide numerical simulations on both synthetic examples and contextual bandits, in comparison with off-policy evaluation and learning baselines.

1 Introduction

The availability of historically collected data is a common scenario in many real-world decision-making problems, including medical treatments [e.g., 17, 67], recommendation systems [e.g., 33, 16], personalized advertising [e.g., 3, 60], finance [e.g., 43], and industrial robot control [e.g., 27, 26]. Historical data can be leveraged to address two classes of problems. First, given data collected with a *behavioral* policy, we want to estimate the performance of a different *target* policy. This problem is known as *off-policy evaluation* [Off-PE, 21]. Second, we want to employ the available data to improve the performance of a baseline policy. This latter problem is named *off-policy learning* [Off-PL 14]. Off-policy methods are studied by both the *reinforcement learning* [RL, 58] and *contextual multi-armed bandit* [CMAB, 30] communities. Given its intrinsic simplicity compared to RL, off-policy methods are nowadays well understood in the CMAB framework [e.g., 44, 1, 14, 64]. Among them, the *doubly robust* estimator [DR, 14] is one of the most promising off-policy methods for CMABs. DR combines a *direct method* (DM), in which the reward is estimated from historical data via regression, with an *importance sampling* [IS, 46] control variate.

IS plays a crucial role in the off-policy methods and counterfactual reasoning. However, IS tends to exhibit problematic behavior for general distributions. This is formalized by its *heavy-tailed* properties [40], which prevent the application of exponential concentration bounds [4]. To cope with this issue, typically, corrections are performed on the importance weight including *truncation* [23] and *self-normalization* [SN, 46], among the most popular. Significant results have recently been derived for both techniques [47, 29, 39]. Nevertheless, we believe that the widespread use of IS calls for a better theoretical understanding of its properties and for the design of principled weight corrections.

Defining the desirable properties of an off-policy estimator is a non-trivial task. Some works employed the *mean squared error* (MSE) as an index of the estimator quality [34, 64]. However, controlling the MSE, while effectively capturing the bias-variance trade-off, does not provide any guarantee on the concentration properties of the estimator, which might still display a heavy-tailed behavior [37]. For this reason, we believe that a more suitable approach is to require that the estimator deviations concentrate at a *subgaussian* rate [12]. Subgaussianity implicitly controls the tail behavior and leads to tight exponential concentration inequalities. Unlike MSE, the probabilistic requirements are non-asymptotic (finite-sample), from which guarantees on the MSE can be derived. While subgaussianity can be considered a satisfactory requirement for Off-PE, additional properties are advisable when switching to Off-PL. In particular, the *differentiability* w.r.t. the target policy parameters is desirable whenever Off-PL is carried out via gradient optimization. For instance, weight truncation, as presented in [47], allows achieving subgaussianity but leads to a non-differentiable objective. Consequently, the optimization phase requires additional care, which sometimes leads to computationally heavy discretizations [47]. On the contrary, the SN estimator is differentiable in the target policy, but fails to achieve subgaussian concentration for general distributions.

In this paper, we take a step towards a better understanding of IS. After having introduced the necessary background (Section 2), we derive an anticoncentration bound for the mean estimation with vanilla IS. We show that polynomial concentration (Chebychev’s inequality) is tight in this setting (Section 3). This result formalizes the intuition behind the undesired behavior of these estimators for general distributions. Hence, we propose a class of importance weight corrections, based on the notion of power mean (Section 4). The rationale behind these corrections is to “smoothly shrink” the weights towards the mean, with different intensities. In this way, we mitigate the heavy-tailed behavior and, in the meantime, we exert control over the induced bias. Then, we derive bounds on the bias and variance that allow obtaining an exponential concentration inequality and, under certain conditions, subgaussian concentration (Section 5). Furthermore, the smooth transformation allows preserving the differentiability in the target policy, unlike some existing transformations, like weight truncation. To the best of our knowledge, this is the first IS correction that preserves the differentiability and is proved to be subgaussian. This correction, however, requires knowledge of a distributional divergence between the target and behavioral policies, which may be unknown or difficult to compute. To this end, we introduce an approach to empirically estimate the correction parameter, preserving the desirable concentration properties (Section 6). After providing a comparative review of the literature (Section 7), we present an experimental study comparing our approach with traditional and modern off-policy baselines on synthetic domains and in the CMAB framework (Section 8). The proofs of the results presented in the main paper can be found in Appendix A. A preliminary version of this work was presented at the “Workshop on Reinforcement Learning Theory” of ICML 2021 [42].¹

2 Preliminaries

We start introducing the background about probability, importance sampling and contextual bandits.

Probability We denote with $\mathcal{P}(\mathcal{Y})$ the set of probability measures over a $(\mathcal{Y}, \mathfrak{F}_{\mathcal{Y}})$. Let $P \in \mathcal{P}(\mathcal{Y})$, $f: \mathcal{Y} \rightarrow \mathbb{R}$ be a function, and $\bar{\mu}_n$ be an estimator for the mean $\mu = \mathbb{E}_{y \sim P}[f(y)]$ obtained with n i.i.d. samples. Suppose that with probability $1 - \delta$ it holds that $|\bar{\mu}_n - \mu| \leq \sqrt{g(n, \delta)}$. For $\beta > 0$, we say that $\bar{\mu}_n$ admits: (i) *polynomial* concentration if $g(n, \delta) = \mathcal{O}(1/(n\delta)^\beta)$; (ii) *exponential* concentration if $g(\delta) = \mathcal{O}((\log(1/\delta)/n)^\beta)$; (iii) *subgaussian* concentration if (ii) holds with $\beta = 1$ [37]. These cases correspond to Chebyshev’s, Bernstein’s, and Höeffding’s inequalities respectively [4].

Importance Sampling Let $P, Q \in \mathcal{P}(\mathcal{Y})$ admitting p and q as density functions, if $P \ll Q$, i.e., P is absolutely continuous w.r.t. Q , for any $\alpha \in (1, 2]$, we introduce the integral: $I_\alpha(P\|Q) = \int_{\mathcal{Y}} p(y)^\alpha q(y)^{1-\alpha} dy$. If $P = Q$ a.s. (almost surely) then $I_\alpha(P\|Q) = 1$. $I_\alpha(P\|Q)$ allows defining several divergences, like Rényi [51]: $(\alpha - 1)^{-1} \log I_\alpha(P\|Q)$. Let $f: \mathcal{Y} \rightarrow \mathbb{R}$ be a function, (*vanilla*) *importance sampling* [IS, 46] allows estimating the expectation of f under the *target* distribution P , i.e., $\mu = \mathbb{E}_{y \sim P}[f(y)]$, using i.i.d. samples $\{y_i\}_{i \in [n]}$ collected with the *behavioral* distribution Q :

$$\hat{\mu}_n = \frac{1}{n} \sum_{i \in [n]} \omega(y_i) f(y_i), \quad \text{where} \quad \omega(y) = \frac{p(y)}{q(y)}, \quad \forall y \in \mathcal{Y}.$$

¹https://lyang36.github.io/icml2021_rltheory/camera_ready/7.pdf.

It is well-known that $\hat{\mu}_n$ is unbiased, i.e., $\mathbb{E}_{y_i \sim Q}[\hat{\mu}_n] = \mu$ [46]. If f is bounded, the variance of the estimator can be upper-bounded as $\mathbb{V}\text{ar}_{y_i \sim Q}[\hat{\mu}_n] \leq \frac{1}{n} \|f\|_\infty I_2(P\|Q)$ [40]. More in general, the integral $I_\alpha(P\|Q)$ represents the α -moment of the importance weight $\omega(y)$ under Q .

Contextual Bandits A *contextual multi-armed bandit* [CMAB, 30] is represented by the tuple $\mathcal{C} = (\mathcal{X}, \mathcal{A}, \rho, p)$, where \mathcal{X} is the set of *contexts*, \mathcal{A} is the finite set of actions (or arms), $\rho \in \mathcal{P}(\mathcal{X})$ is the context distribution, and $p: \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{P}(\mathbb{R})$ is the reward distribution. The agent’s behavior is encoded by a *policy* $\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})$. At each round $t \in \mathbb{N}$, the agent observes a context $x_t \sim \rho$, plays an action $a_t \sim \pi(\cdot|x_t)$, gets the reward $r_t \sim p(\cdot|x_t, a_t)$ and the system moves on to the next round. The *value* of a policy π is given by $v(\pi) = \int_{\mathcal{X}} \rho(x) \sum_{a \in \mathcal{A}} \pi(a|x) \int_{\mathbb{R}} p(r|x, a) r dr dx$. We denote with $r(x, a) = \int_{\mathbb{R}} p(r|x, a) r dr$ the *reward function*. A policy π^* is optimal if it maximizes the value function, i.e., $\pi^* \in \arg \max_{\pi \in \Pi} v(\pi)$, where $\Pi = \{\pi: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})\}$ is the set of all policies.

Let $\mathcal{D} = \{(x_t, a_t, r_t)\}_{t \in [n]}$ be a dataset of samples collected in a CMAB with a behavioral policy $\pi_b \in \Pi$. The *off-policy evaluation* [Off-PE, 21] problem consists in estimating the value function $v(\pi_e)$ of a target policy $\pi_e \in \Pi$ using the samples in \mathcal{D} . The *off-policy learning* [Off-PL, 14] problem consists in estimating an optimal policy $\pi^* \in \Pi$ using the samples in \mathcal{D} . The simplest approach to address the Off-PE/Off-PL problem is to learn from \mathcal{D} an estimate $\hat{r}(x, a)$ of the reward function $r(x, a)$ via regression. This approach is known as *direct method* (DM) and its properties heavily depend on the quality of the estimate \hat{r} . Another approach is to simply apply IS to reweight the samples of \mathcal{D} , leading to the *inverse propensity scoring* [IS, 17] estimator. The two approaches are combined in the *doubly-robust* [DR, 14] estimator, in which the DM estimate is corrected with an IS control variate to reduce the variance using the estimated reward \hat{r} (see also Table 12 in Appendix D).

3 Anticoncentration of Vanilla Importance Sampling

In this section, we analyze the intrinsic limitations of the vanilla IS. It is well-known that under the assumption that for some $\alpha \in (1, 2]$ the divergence $I_\alpha(P\|Q)$ is finite and f is bounded, the vanilla IS estimator $\hat{\mu}_n$ admits *polynomial concentration*, i.e., with probability at least $1 - \delta$:²

$$|\hat{\mu}_n - \mu| \leq \|f\|_\infty \left(\frac{2^{2-\alpha} I_\alpha(P\|Q)}{\delta n^{\alpha-1}} \right)^{\frac{1}{\alpha}}. \quad (1)$$

We now show that the concentration in Equation (1) is tight, by deriving an anticoncentration bound for $|\hat{\mu}_n - \mu|$; then, we discuss its implications and compare it with previous works.

Theorem 3.1. *There exist two distributions $P, Q \in \mathcal{P}(\mathcal{Y})$ with $P \ll Q$ and a bounded measurable function $f: \mathcal{Y} \rightarrow \mathbb{R}$ such that for every $\alpha \in (1, 2]$ and $\delta \in (0, e^{-1})$ if $n \geq \delta e \max\{1, (I_\alpha(P\|Q) - 1)^{\frac{1}{\alpha-1}}\}$, with probability at least δ it holds that:*

$$|\hat{\mu}_n - \mu| \geq \|f\|_\infty \left(\frac{I_\alpha(P\|Q) - 1}{\delta n^{\alpha-1}} \right)^{\frac{1}{\alpha}} \left(1 - \frac{e\delta}{n} \right)^{\frac{n-1}{\alpha}}.$$

First of all, we note the polynomial dependence on the confidence level δ . The bound is vacuous when $I_\alpha(P\|Q) = 1$, i.e., when $P = Q$ a.s., since in an on-policy setting and, being the function f bounded, subgaussian concentration bounds (like Höeffding’s inequality) hold. In particular, for $\alpha = 2$, n and $I_2(P\|Q)$ sufficiently large, the bound has order $\mathcal{O}(\sqrt{I_2(P\|Q)/(\delta n)})$, matching Chebyshev’s and the existing concentration inequalities for vanilla importance sampling [40, 41].

Our result is of independent interest and applies for general distributions. Previous works considered the MAB [34] and CMAB [64] settings proving *minimax* lower bounds in *mean squared error* (MSE) $\mathbb{E}_{y \sim Q}[(\hat{\mu}_n - \mu)^2]$. These results differ from ours in several respects. First, we focus on a *specific estimator*, the vanilla one, while those result are *minimax*. Second, they provide lower bounds to the MSE, while we focus on the deviations in probability.³ From our probabilistic result, it is immediate to derive an MSE guarantee (Corollary A.1 of Appendix A.1). Finally, they assume that the second moment of the importance weight $I_2(P\|Q)$ is finite, whereas our result allows considering scenarios in which only moments of order $\alpha < 2$ are finite.

²The original result [41, Theorem 2] was limited to $\alpha = 2$ and based on Cantelli’s inequality which approaches Chebyshev’s when $\delta \rightarrow 0$. See Theorem A.1 in Appendix A.1, for a proof of Equation (1).

³As noted in [36], when the estimator is not well-concentrated around its mean (e.g., in presence of heavy tails), the MSE is not adequate to capture the error and high-probability bounds are more advisable.

s	$\omega_{s,\lambda}(y)$
$-\infty$ (minimum)	$\min\{\omega(y), 1\}$
-1 (harmonic)	$\frac{\omega(y)}{1-\lambda+\lambda\omega(y)}$
0 (geometric)	$\omega(y)^{1-\lambda}$
1 (arithmetic)	$(1-\lambda)\omega(y) + \lambda$

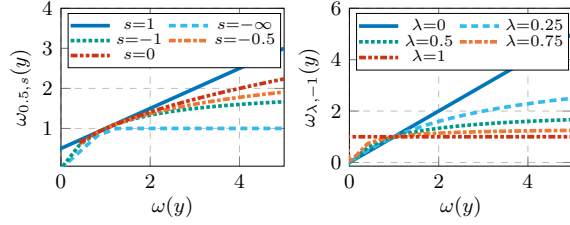


Table 1: Choices of s for the (λ, s) -corrected importance weight of Definition 4.1.

Figure 1: Examples of importance weight corrections of Definition 4.1 for fixed λ (left) and fixed s (right).

4 Power-Mean Correction of Importance Sampling

In this section, motivated by the negative result of Theorem 3.1, we look for a weight correction able to achieve subgaussian concentration. Specifically, we introduce a class of corrections based on the notion of *power mean* [6] and we study its properties. Let us start with the following definition.

Definition 4.1. Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions such that $P \ll Q$, for every $s \in [-\infty, \infty]$ and $\lambda \in [0, 1]$, let $\omega(y) = p(y)/q(y)$, the (λ, s) -corrected importance weight is defined as:

$$\omega_{\lambda,s}(y) = \left((1-\lambda)\omega(y)^s + \lambda \right)^{\frac{1}{s}}, \quad \forall y \in \mathcal{Y}.$$

The correction can be seen as the weighted *power mean* with exponent s between the vanilla importance weight $\omega(y)$ and 1 with weights $1-\lambda$ and λ respectively.⁴ We immediately notice that, regardless of the value of s , for $\lambda=0$, we reduce to the vanilla importance weight $\omega_{0,s}(y) = \omega(y)$ and for $\lambda=1$, we have identically $\omega_{1,s}(y) = 1$. Furthermore, the correction is unbiased when $P=Q$ a.s. regardless s and λ . Thus, the correction “smoothly interpolates” between the vanilla weight $\omega(y)$ and its mean under Q , i.e., 1. s and λ govern the “intensity” of the correction in a continuous way. Differently from the truncation [23], this transformation leads to a differentiable weight. Some specific choices of s and λ are reported in Table 1 and in Figure 1. The following result provides a preliminary characterization of the correction, independent of the properties of the two distributions.

Lemma 4.1. Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions with $P \ll Q$, then for every $\lambda \in [0, 1]$ and $y \in \mathcal{Y}$ it holds that:

- (i) if $s \leq s'$ then $\omega_{\lambda,s}(y) \leq \omega_{\lambda,s'}(y)$;
- (ii) if $s < 0$ then $\omega_{\lambda,s}(y) \leq \lambda^{\frac{1}{s}}$, otherwise if $s > 0$ then $\omega_{\lambda,s}(y) \geq \lambda^{\frac{1}{s}}$;
- (iii) if $s < 1$ then $\mathbb{E}_{y \sim Q}[\omega_{\lambda,s}(y)] \leq 1$, otherwise if $s > 1$ then $\mathbb{E}_{y \sim Q}[\omega_{\lambda,s}(y)] \geq 1$.

From point (ii) we observe that the corrected weight is bounded from below when $s > 0$ and bounded from above when $s < 0$. It is well-known that the inconvenient behavior of IS derives from the heavy-tailed properties [40]. Thus, the *arithmetic* correction ($s=1$) performs just a convex combination between the vanilla weight and 1, having no effect on the tail properties. Any correction with $s > 1$ increases the weight value, making the tail even heavier. So, if we are looking for subgaussianity, we should restrict our attention to $s < 0$, which leads to lighter tails or even bounded weights.

5 Subgaussian and Differentiable Importance Sampling

In this section, we focus on the *harmonic* correction ($s=-1$), which leads to a weight of the form:⁵ $\omega_{\lambda,-1}(y) = \frac{\omega(y)}{1-\lambda+\lambda\omega(y)}$. We start analyzing the bias and variance of this class of estimators. Then, we provide an exponential concentration inequality that, under certain circumstances, results to be subgaussian. Finally, we show that the resulting estimator is differentiable in the target distribution. To lighten the notation we neglect the -1 subscript, abbreviating $\hat{\mu}_\lambda = \hat{\mu}_{\lambda,-1}$.

Bias and Variance We now derive bounds for the bias and the variance induced by the $(\lambda, -1)$ -corrected importance weight.

⁴For $s \in \{-\infty, 0, \infty\}$ the power mean is defined as a limit.

⁵The choice of $s=-1$ is mainly for analytical convenience and, as we shall see, it already allows enforcing the desired properties. We leave investigating the other values of s for future work.

Lemma 5.1. *Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions with $P \ll Q$. For every $\lambda \in [0, 1]$, the bias and variance of the $(\lambda, -1)$ -corrected importance weight can be bounded for every $\alpha \in (1, 2]$ as:*

$$\left| \mathbb{E}_{y \sim Q} [\hat{\mu}_{n,\lambda}] - \mu \right| \leq \|f\|_\infty \lambda^{\alpha-1} I_\alpha(P\|Q), \quad \text{Var}_{y_i \sim Q} [\hat{\mu}_{n,\lambda}] \leq \frac{\|f\|_\infty^2}{n\lambda^{2-\alpha}} I_\alpha(P\|Q).$$

As expected, the bias is zero for $\lambda=0$; it increases with λ and with the divergence term $I_\alpha(P\|Q)$. Indeed, we already observed that the bias is null when $P=Q$ a.s.. In particular, for $\alpha=2$, the bound becomes $\|f\|_\infty \lambda I_2(P\|Q)$. Instead, the variance bound decreases in λ and increases with the divergence $I_\alpha(P\|Q)$. For $\alpha=2$, we obtain the bound $\frac{1}{n} \|f\|_\infty^2 I_2(P\|Q)$. Note that when $P=Q$ a.s., we recover $\frac{1}{n} \|f\|_\infty^2$, which is the Popoviciu's inequality for the variance [50]. Thus, our weight correction allows controlling bias and variance even when $I_2(P\|Q) = \infty$, i.e., when the vanilla IS estimator might have infinite variance. Indeed, our transformed estimator has finite variance provided that there exists $\alpha \in (1, 2)$ so that $I_\alpha(P\|Q) < \infty$. Tighter (but less intelligible) bounds on bias and variance are reported in Appendix A.3.

Concentration Inequality We are now ready to derive the core theoretical result. We prove an exponential concentration inequality for the $(\lambda, -1)$ -corrected IS estimator and we show that, if $I_2(P\|Q)$ is finite, we are able to achieve subgaussian concentration.⁶

Theorem 5.1. *Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions such that $P \ll Q$. For every $\alpha \in (1, 2]$ and $\delta \in (0, 1)$, if we select $\lambda = \lambda_\alpha^*$ then, with probability at least $1 - \delta$ it holds that:*

$$\hat{\mu}_{n,\lambda_\alpha^*} - \mu \leq \|f\|_\infty (2 + \sqrt{3}) \left(\frac{2I_\alpha(P\|Q)^{\frac{1}{\alpha-1}} \log \frac{1}{\delta}}{3(\alpha-1)^2 n} \right)^{1-\frac{1}{\alpha}}, \quad \text{with } \lambda_\alpha^* = \left(\frac{2 \log \frac{1}{\delta}}{3(\alpha-1)^2 I_\alpha(P\|Q)n} \right)^{\frac{1}{\alpha}}.$$

We immediately notice that the dependence on the confidence level δ is the one typical of exponential concentration for every $\alpha \in (1, 2]$. In particular, we observe that the bound is subgaussian when $\alpha=2$, requiring that $I_2(P\|Q) < \infty$. Recalling that $I_2(P\|Q)$ governs the variance of the estimator, this result is in line with the general theory of estimators for which the existence of the variance is an unavoidable requirement to achieve subgaussian concentration [12]. Specifically, for $\alpha=2$ the optimal value of the parameter is $\lambda_2^* = \sqrt{(2 \log(1/\delta)) / (3I_2(P\|Q)n)}$, leading to the bound:

$$\hat{\mu}_{n,\lambda_2^*} - \mu \leq \|f\|_\infty (2 + \sqrt{3}) \sqrt{\frac{2I_2(P\|Q) \log \frac{1}{\delta}}{3n}}. \quad (2)$$

A tighter bound, based on a different choice of the correction parameter λ_α^{**} is derived in Appendix A.3 and it is omitted here for clarity of presentation.

Differentiability As we have already observed, our weight correction, differently from truncation, is smooth and, thus, differentiable in the target policy. We now focus on the properties of the *gradient* of the $(\lambda, -1)$ -corrected estimator and, to this purpose, we constrain the target distribution to belong to a parametric space differentiable distributions $\mathcal{P}_\Theta = \{P_\theta \in \mathcal{P}(\mathcal{Y}) : \theta \in \Theta\}$, where $\Theta \subseteq \mathbb{R}^d$. The gradient of the corrected weight ω_λ w.r.t. the target policy parameters θ is given by:

$$\nabla_\theta \omega_\lambda(y) = \nabla_\theta \frac{p_\theta(y)}{q(y)} = \frac{(1-\lambda)\omega(y)}{(1-\lambda + \lambda\omega(y))^2} \nabla_\theta \log p_\theta(y), \quad \forall y \in \mathcal{Y}.$$

In particular, it can be proved that $\|\nabla_\theta \omega_\lambda(y)\|_\infty \leq \frac{1}{4\lambda} \|\nabla_\theta \log p_\theta(y)\|_\infty$ (Proposition A.1 of Appendix A.3). Thus, if the score $\nabla_\theta \log p_\theta$ is bounded, the gradient will be bounded whenever $\lambda > 0$. This property is advisable for gradient optimization and it is not guaranteed, for example, for vanilla IS ($\lambda=0$). Thus, we can also interpret λ as a regularization parameter for the gradient magnitude.

6 Data-driven Tuning of λ

The computation of the parameter λ_2^* requires the knowledge of the divergence $I_2(P\|Q)$. Even when P and Q are known, computing the $I_2(P\|Q)$ can be challenging, especially for continuous

⁶We introduce our concentration inequalities as a one-sided bounds just for simplicity but they actually hold in both directions. Indeed, by replacing function f with function $-f$, we obtain the reversed one-sided bound.

distributions, since it involves the evaluation of a complex integral.⁷ In principle, we could estimate the divergence $I_2(P\|Q)$ from samples as the empirical second moment of the vanilla importance weights $\frac{1}{n} \sum_{i \in [n]} \omega(y_i)^2$. However, although possibly well-performing in practice [40], this approach would prevent any subgaussian concentration, as the behavior of the non-corrected $\omega(y)^2$ will be surely heavy-tailed whenever $\omega(y)$ is. A general-purpose approach to avoid the divergence estimation is the *Lepski's adaptation method* [32], which only requires knowing an upper and a lower bound on $I_2(P\|Q)$. Unfortunately, this method is known to be computationally intensive.

In this section, we follow a different path inspired by the recent work [66]. If a choice of the parameter λ corrects the weight ω_λ leading to an ideal estimator $\hat{\mu}_{n,\lambda}$, for the mean μ , we may expect that the empirical second moment of the corrected weights ω_λ will provide a reasonable estimation of $I_2(P\|Q)$. Based on this observation, we propose to choose λ by solving the following equation:

$$\lambda^2 \underbrace{\frac{1}{n} \sum_{i \in [n]} \omega_{\lambda n^{1/4}}(y_i)^2}_{\text{empirical second moment}} = \frac{2 \log \frac{1}{\delta}}{3n}. \quad (3)$$

The intuition behind this approach can be stated as follows. If the empirical second moment is close to the divergence, i.e., $\frac{1}{n} \sum_{i \in [n]} \omega_{\lambda n^{1/4}}(y_i)^2 \simeq I_2(P\|Q)$, the solution $\hat{\lambda}$ of Equation (3) approaches the optimal parameter, i.e., $\hat{\lambda} \simeq \sqrt{(2 \log(1/\delta)) / (3I_2(P\|Q)n)} = \lambda_2^*$. We formalize this reasoning in Appendix A.4, proving that Equation (3) admits a unique root $\hat{\lambda} \in [0, 1]$ (Lemma A.4) and that when the number of samples n grows to infinity, $\hat{\lambda}$ converges indeed to λ_2^* (Lemma A.8). The following result provides the concentration properties of the estimator $\hat{\mu}_{n,\lambda}$ when using $\hat{\lambda}$ instead of λ_2^* , under slightly more demanding requirements on the moments of the importance weights.

Theorem 6.1. *Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions such that $P \ll Q$. Let $\hat{\lambda}$ be the solution of Equation (3), then, if $I_3(P\|Q)$ is finite, for sufficiently large n , for every $\delta \in (0, 1)$, with probability at least $1 - 2\delta$ it holds that:*

$$\hat{\mu}_{n,\hat{\lambda}} - \mu \leq \|f\|_\infty \frac{5 + 2\sqrt{3}}{2} \sqrt{\frac{2I_2(P\|Q) \log \frac{1}{\delta}}{3n}}.$$

Compared to Theorem 5.1, this result is weakened in two aspects. First, the inequality holds with a smaller probability $1 - 2\delta$ since two estimation processes with the same samples are needed, i.e., the computation of $\hat{\lambda}$ and the corrected estimator $\hat{\mu}_{n,\hat{\lambda}}$. Second, and most important, the result holds for large n , whose minimum value is reported in the proof and depends on $I_3(P\|Q)$, which must be finite. We think this is a not too strong requirement considering that even the variance of an empirical estimate of $I_2(P\|Q)$ would depend on the fourth moment of the importance weight, i.e., $I_4(P\|Q)$.⁸

7 Related Works

Importance Sampling has a long history in Monte Carlo simulation as an effective technique for variance reduction in presence of rare events and for what-if analysis [25, 55, 20, 9, 52]. Apart from sparse exceptions [e.g., 8, 18], in the machine learning community, IS is primarily employed for off-policy estimation and learning [e.g., 11, 38, 61]. In this setting, it is well-known that IS might display an inconvenient behavior, depending on the behavioral Q and target P distributions [65, 40]. In particular, IS tends to enlarge the range of the estimator up to $\text{esssup}_{y \sim Q} p(y)/q(y)$. Although this term is finite for discrete distributions (if $P \ll Q$), it is likely unbounded for continuous ones [11]. Furthermore, in the latter case, the vanilla IS estimator might have infinite variance and tends to exhibit a heavy-tailed behavior [40, 41]. These properties suggest that a way of addressing this phenomenon is to resort to *robust* statistics, typically employed for mean estimation under heavy-tailed distributions [37]. Methods in this class include the *trimmed mean* [62, 22], the *median of means* [45, 24], and the *Catoni's estimator* [7]. For all of them, subgaussian guarantees were

⁷For some common distributions, including Gaussians, the integral can be computed in closed form [15].

⁸It is possible to circumvent the computation of $I_2(P\|Q)$ by choosing λ independently from $I_2(P\|Q)$ at the price of downgrading the concentration from subgaussian to exponential (Corollary A.2 of Appendix A.4).

Estimator	Maximum Variance	Bias	Correction (order \mathcal{O})	Concentration (order \mathcal{O})	Is subgaussian?	Is unbiased when $P=Q$?	Is differentiable?	
IS [46, 40]	$\text{ess sup } \frac{p}{q}$	$\frac{I_2(P\ Q)}{n}$	0	-	$\sqrt{\frac{I_2(P\ Q)}{\delta n}}$	✗ (poly)	✓	✓
SN-IS [46, 29]	1	V^{SN}	B^{SN}	-	$B^{\text{SN}} + \sqrt{V^{\text{ES}} \log \frac{1}{\delta}}$	✗ (exp)	✓	✓
IS-TR(M) [23, 47]	M	$\frac{I_2(P\ Q)}{n}$	$\frac{I_2(P\ Q)}{M}$	$\sqrt{\frac{nI_2(P\ Q)}{\log \frac{1}{\delta}}}$	$\sqrt{\frac{I_2(P\ Q) \log \frac{1}{\delta}}{n}}$	✓	✗	✗
IS-OS(τ) [56]	$\frac{\sqrt{\tau}}{2}$	$\frac{I_2(P\ Q)}{n}$	$\frac{I_3(P\ Q)}{\tau}$	$\sqrt[3]{\left(\frac{nI_3(P\ Q)}{\log \frac{1}{\delta}}\right)^2}$	$\max_{\beta \in \{2,3\}} \sqrt{\frac{I_\beta(P\ Q) (\log \frac{1}{\delta})^{\beta-1}}{n^{\beta-1}}}$	✗ (exp)	✗	✓
IS- λ (ours)	$\frac{1}{\lambda}$	$\frac{I_2(P\ Q)}{n}$	$\lambda I_2(P\ Q)$	$\sqrt{\frac{\log \frac{1}{\delta}}{I_2(P\ Q)n}}$	$\sqrt{\frac{I_2(P\ Q) \log \frac{1}{\delta}}{n}}$	✓	✓	✓

Table 2: Comparison between several IS estimators assuming $\|f\|_\infty = 1$ and for $\alpha = 2$ w.r.t. several indexes. For the SN-IS estimator V^{SN} is the Efron-Stein estimate of the variance and B^{SN} is the bias. V^{SN} and B^{SN} converge to 0 as $n \rightarrow \infty$, but no convergence rate is provided in [29].

provided [37]. These techniques have been also successfully employed for regret minimization in finite [5] and continuous arm spaces [36]. In principle, these methods could be employed *as-is* in combination with IS, but, being general-purpose, they might disregard the peculiarities of the setting.

Several *ad-hoc* methods to cope with the problematic IS behavior have been progressively developed. An example, devised by the statistical community, is *self-normalization* [SN-IS, 46]: $\tilde{\omega}(y_i) = \omega(y_i) / \sum_{j \in [n]} \omega(y_j)$. This approach has the advantage of controlling the range of the estimator at the price of making all samples interdependent and generating a bias. Although the asymptotic consistency is guaranteed [19, 59], its finite-sample analysis is more challenging. In [40], a polynomial concentration inequality was provided and, more recently, exponential bounds based on Efron-Stein inequalities have been proposed [28, 29]. Nevertheless, the resulting inequality is not guaranteed to decrease with $\mathcal{O}(1/\sqrt{n})$ for general distributions and it is difficult to formally relate its concentration rate to the tail properties of the involved distributions [29]. Another popular technique is the weight *truncation* (or *clipping*) [IS-TR, 23, 3]: $\omega_M^{\text{TR}}(y) = \min\{\omega(y), M\}$, where $M > 0$ is the threshold. Some works rely on empirical selections of the truncation threshold [31, 10], while others focus on more theoretically principled approaches [2, 64, 47]. In particular, in [47] a subgaussian deviation bound is derived by suitably adapting the truncation threshold as a function of the number of samples n and the confidence parameter δ . Another interesting approach, designed for CMABs is the *switch* estimator [DR-SW, 64] that selects between DM and IS (or DR), based on the importance weight value, with also guarantees in MSE. Finally, a not so large part of the literature focuses on less crude transformations than truncation, called *smoothing* [63]. They typically take into explicit consideration the estimator tails [48], also providing asymptotic guarantees. Very recently, shrinkage transformations of the weight were proposed, based on the minimization of different bounds on the MSE, in the CMABs [56] setting. Specifically, the *optimistic shrinkage* [IS-OS, 56] leads to a transformation similar to ours $\omega_\tau^{\text{OS}}(y) = \tau\omega(y)/(\omega(y)^2 + \tau)$. Unfortunately, even when knowing P and Q and setting τ adaptively, IS-OS is unable to achieve subgaussian concentration and requires $I_3(P\|Q)$ to be finite (Appendix E for details). Refer to Table 2 for a comparison of the estimators.

8 Numerical Simulations

In this section, we provide numerical simulations for off-policy evaluation (Section 8.1) and learning (Section 8.2), with the goal of showing that our estimators, while enjoying desirable theoretical properties, are competitive with traditional (e.g., vanilla IS and self-normalization) and modern baselines (e.g., truncation, optimistic shrinkage). The complete results can be found in Appendix B.⁹

8.1 Off-Policy Evaluation

We present two off-policy evaluation experiments. We start with a synthetic example with Gaussian distributions and, then, we move to the CMAB setting.

⁹The code is provided at <https://github.com/albertometelli/subgaussian-is>.

Estimator / n	10	20	50	100	200	500	1000
IS	27.43 ± 13.33	15.70 ± 4.83	10.89 ± 1.81	9.26 ± 0.92	12.41 ± 1.88	9.42 ± 0.68	5.84 ± 0.27
SN-IS	23.89 ± 5.77	15.62 ± 2.62	10.96 ± 1.18	9.53 ± 0.74	8.82 ± 0.62	7.48 ± 0.37	5.14 ± 0.20
IS-TR	23.47 ± 7.52	14.03 ± 2.75	10.32 ± 1.47	8.89 ± 0.79	7.68 ± 0.46	6.21 ± 0.28	4.22 ± 0.15
IS-OS	19.25 ± 8.68	10.93 ± 3.29	8.37 ± 1.35	7.06 ± 0.61	8.69 ± 1.44	6.65 ± 0.47	3.97 ± 0.16
IS- λ^*	21.75 ± 6.36	13.17 ± 2.45	9.26 ± 1.19	7.76 ± 0.62	6.53 ± 0.38	5.29 ± 0.23	3.52 ± 0.12
IS- λ^{**}	20.66 ± 4.08	12.62 ± 2.19	8.86 ± 1.08	7.39 ± 0.57	5.94 ± 0.32	4.74 ± 0.20	3.19 ± 0.10
IS- $\hat{\lambda}$	18.19 ± 3.93	10.27 ± 1.64	7.03 ± 0.75	5.79 ± 0.38	3.85 ± 0.21	2.90 ± 0.10	2.06 ± 0.05

Table 3: Absolute error in the illustrative example varying the number of samples n for the different estimators (mean \pm std, 60 runs). For each column, the estimator with smallest absolute error and the ones not statistically significantly different from that one (Welch’s t-test with $p < 0.02$) are in bold.

s / λ	0	0.1	0.2	0.5
$-\infty$	3.12 ± 0.29	3.12 ± 0.29	3.12 ± 0.29	3.12 ± 0.29
-5	6.73 ± 1.21	2.70 ± 0.30	2.77 ± 0.30	2.57 ± 0.31
-2	6.73 ± 1.21	2.45 ± 0.34	2.42 ± 0.32	2.28 ± 0.32
-1	6.73 ± 1.21	2.72 ± 0.47	2.47 ± 0.37	2.18 ± 0.32
-0.5	6.73 ± 1.21	3.44 ± 0.64	2.71 ± 0.47	2.20 ± 0.34
0	6.73 ± 1.21	4.83 ± 0.89	3.66 ± 0.68	2.38 ± 0.38
0.5	6.73 ± 1.21	5.69 ± 1.05	4.85 ± 0.89	3.03 ± 0.52

Table 4: Absolute error in the illustrative example varying the parameter s of the corrected weight when $n = 500$ (mean \pm std, 60 runs). The estimator with smallest absolute error and the ones not statistically significantly different from that one (Welch’s t-test with $p < 0.02$) are in bold.

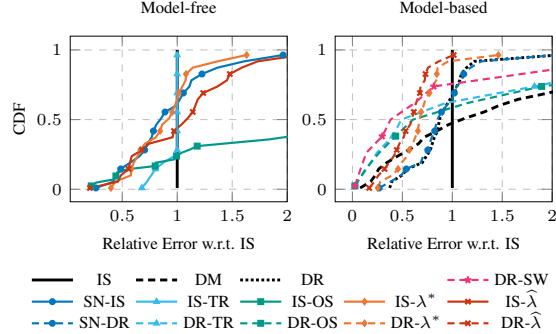


Figure 2: CDF of the absolute error normalized by IS error for stochastic rewards with noise 0.25, across 110 conditions.

8.1.1 Synthetic Experiment

In this experiment, we compare our corrected estimators with IS baselines in a *continuous-distribution* off-policy estimation problem. Specifically, we consider a Gaussian behavioral policy $Q = \mathcal{N}(\mu_Q, \sigma_Q^2)$ and a Gaussian target policy $P = \mathcal{N}(\mu_P, \sigma_P^2)$. We generate n i.i.d. samples from Q and we estimate the expectation of function $f(y) = 100 \cos(2\pi y)$ under P . We select $\mu_Q = 0$, $\mu_P = 0.5$, $\sigma_Q^2 = 1$ and $\sigma_P^2 = 1.9$, leading to a divergence $I_2(P\|Q) \simeq 27.9$. The results with different choices of the σ_P^2 are reported in Appendix B.1.1.

Estimators Comparison In Table 3, we report the absolute error between the estimated and the true mean for the different importance sampling estimators. For our correction, we report the results obtained with optimal value of λ according to Theorem 5.1 with $\alpha = 2$ (IS- λ^*), a value of λ that optimizes a tighter bound reported in Appendix A.3 (IS- λ^{**}), and the value estimated from samples as in Section 6 (IS- $\hat{\lambda}$). We compare these estimators with vanilla IS (IS), self-normalized IS (SN-IS), weight truncation (IS-TR) with optimal threshold selected as in [47], and IS with optimistic shrinkage (IS-OS), where τ is computed by minimizing an MSE bound as in [56]. We notice that our estimators consistently outperform the traditional ones (IS and SN-IS) and overall suffer smaller errors than IS-TR and IS-OS. Interestingly, the minimum error is often obtained by IS- $\hat{\lambda}$, which uses an estimated value $\hat{\lambda}$ that tends to get a higher value than both λ^* and λ^{**} . In this way, the correction is more intense, which, in this specific example, turns out to be beneficial.

Comparison of Different Values of s We empirically test different values of the parameter s employed in Definition 4.1, in the same setting of Table 3 with $n = 500$ for the estimator IS- λ . Since for general value of s , we do not have a principled way to select the correction parameter λ , we consider different values of λ . The results are reported in Table 4. We can see that the best results are obtained with $s \in \{-1, -2\}$.

8.1.2 Contextual Bandits

In this section, we report the experiments about off-policy evaluation in CMABs.

Estimator / n	100	200	500	1000	2000	5000	10000	20000
IS	17.38 ± 1.27	22.26 ± 2.00	15.98 ± 0.82	8.36 ± 0.21	4.67 ± 0.07	2.68 ± 0.03	2.15 ± 0.02	1.10 ± 0.00
SN-IS	23.95 ± 1.68	19.39 ± 1.30	17.94 ± 0.50	11.43 ± 0.22	7.10 ± 0.13	2.54 ± 0.03	1.61 ± 0.01	1.10 ± 0.00
IS-TR	17.38 ± 1.27	18.92 ± 1.36	15.88 ± 0.82	8.36 ± 0.21	4.67 ± 0.07	2.68 ± 0.03	2.15 ± 0.02	1.10 ± 0.00
IS-OS	24.91 ± 1.45	31.93 ± 1.15	15.38 ± 0.56	17.25 ± 0.45	16.41 ± 0.37	30.63 ± 0.15	33.95 ± 0.02	33.61 ± 0.01
IS- λ^*	17.22 ± 1.35	17.10 ± 1.03	11.57 ± 0.45	5.66 ± 0.17	4.86 ± 0.06	2.73 ± 0.02	2.47 ± 0.02	1.27 ± 0.01
IS- λ^{**}	17.21 ± 1.39	16.20 ± 0.93	10.93 ± 0.37	5.55 ± 0.15	5.05 ± 0.06	2.82 ± 0.02	2.68 ± 0.02	1.43 ± 0.01
IS- $\hat{\lambda}$	18.16 ± 1.49	16.52 ± 0.85	11.23 ± 0.29	6.48 ± 0.15	5.85 ± 0.07	3.01 ± 0.03	2.89 ± 0.02	1.50 ± 0.01
DM	20.52 ± 1.18	25.28 ± 0.97	36.19 ± 0.31	36.04 ± 0.08	36.95 ± 0.06	41.99 ± 0.01	42.70 ± 0.01	42.71 ± 0.00
DR	23.00 ± 1.88	25.79 ± 2.38	20.02 ± 0.92	8.30 ± 0.17	4.37 ± 0.08	2.16 ± 0.02	1.38 ± 0.01	0.64 ± 0.00
SN-DR	20.89 ± 1.45	23.38 ± 1.91	20.79 ± 0.74	10.99 ± 0.17	6.48 ± 0.11	2.54 ± 0.02	1.52 ± 0.01	0.99 ± 0.00
DR-TR	18.48 ± 1.13	15.96 ± 0.72	18.58 ± 0.23	15.52 ± 0.09	15.45 ± 0.07	20.33 ± 0.01	21.05 ± 0.01	20.78 ± 0.00
DR-OS	18.47 ± 1.17	18.84 ± 0.60	17.10 ± 0.39	12.19 ± 0.22	8.86 ± 0.11	17.52 ± 0.06	18.40 ± 0.02	19.04 ± 0.02
DR-SW	22.83 ± 1.25	16.81 ± 1.14	4.59 ± 0.18	4.70 ± 0.09	4.86 ± 0.06	0.77 ± 0.01	1.38 ± 0.01	0.78 ± 0.00
DR- λ^*	20.03 ± 1.25	18.70 ± 1.33	13.04 ± 0.61	6.22 ± 0.13	3.82 ± 0.07	1.79 ± 0.02	1.37 ± 0.01	0.61 ± 0.00
DR- λ^{**}	19.40 ± 1.22	17.29 ± 1.17	11.39 ± 0.53	5.60 ± 0.12	3.65 ± 0.06	1.67 ± 0.02	1.38 ± 0.01	0.63 ± 0.00
DR- $\hat{\lambda}$	18.53 ± 1.21	14.92 ± 0.98	9.18 ± 0.44	4.91 ± 0.10	3.39 ± 0.06	1.61 ± 0.02	1.40 ± 0.01	0.65 ± 0.00

Table 5: Absolute error (multiplied by 100) in the *letter* dataset varying the number of samples n for the different estimators, when $\alpha_b = 0.5$ and $\alpha_e = 0.9$ (mean \pm std, 10 runs). For each column, the estimator with smallest absolute error and the ones not statistically significantly different from that one (Welch’s t-test with $p < 0.05$) are in bold.

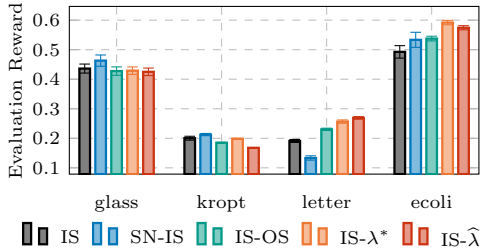


Figure 3: Evaluation reward for four different datasets after 1000 iterations (4000 iterations for *letter*) of gradient ascent with a Boltzmann policy for the model-free estimators (mean \pm std, 10 runs).

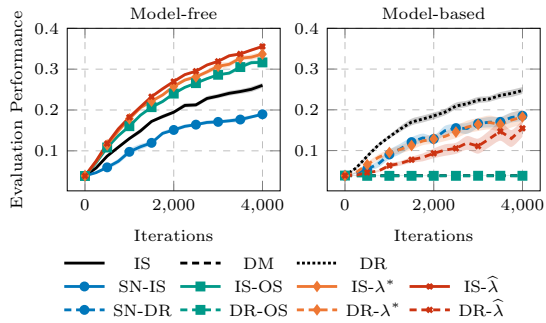


Figure 4: Evaluation reward for the *letter* dataset comparing the learning curve of different estimators (mean \pm std, 10 runs).

Setting We follow the well-established setting of [14, 64, 57, 56]. We consider 11 UCI [13] multi-class classification datasets (see Table 9 in Appendix B.1.2). Each dataset $\mathcal{D}^* = \{(x_i, a_i^*)\}_{i \in [n^*]}$ is mapped to a CMAB problem with action set $\mathcal{A} = [K]$. Every sample (x_i, a) leads to a reward given by $\mathbb{1}\{a = a_i^*\}$. To model noise, the reward is switched with probability $\nu \in [0, 1]$. Each dataset is split into a training set $\mathcal{D}_{\text{train}}$ and an evaluation $\mathcal{D}_{\text{eval}}$ with proportions 30% and 70%. A multi-class classifier \mathcal{C} is trained on $\mathcal{D}_{\text{train}}$. The behavioral policy is obtained as: $\pi_b(a|x) = \alpha_b + \frac{1-\alpha_b}{K}$ if $a = \mathcal{C}(x)$ and $\pi_b(a|x) = \frac{1-\alpha_b}{K}$ otherwise, where $\alpha_b \in [0, 1]$. The target policy π_e is obtained as the behavioral one by training another classifier on $\mathcal{D}_{\text{train}}$ and using $\alpha_e \in [0, 1]$. We employ π_b to generate a dataset $\mathcal{D} = \{(x_i, a_i, r_i)\}_{i \in [n]}$ sampling x_i from $\mathcal{D}_{\text{eval}}$ where $a_i \sim \pi_b(\cdot|x_i)$ and r_i is computed as described before. The ground truth value function is computed as $v(\pi_e) = \frac{1}{n} \sum_{x \in \mathcal{D}_{\text{eval}}} \sum_{a \in \mathcal{A}} \pi_e(a|x) r(x, a)$. For DM and DR, we employ a regressor to learn the reward with a cross-fitting procedure on the full \mathcal{D} .

Estimators Comparison We consider several settings that vary the values of α_b and α_e across all the 11 datasets, generating 110 combinations and a reward noise of $\nu = 0.25$. Details and results for the noiseless case are reported in Appendix B.1.2. To summarize the results, following the approach of [64], we plot in Figure 2 the cumulative distribution function (CDF) of the absolute error normalized by the error of IS. A lower error corresponds to a CDF curve towards the upper-left corner. We distinguish between the approaches that do not make use of the reward estimate \hat{r} (model-free, left) and the ones that do (model-based, right). As for the model-free ones, we note that the performance of our estimator IS- λ^* is very close to that of SN-IS. This is likely because we are dealing with discrete distributions (actions are finite), which implicitly mitigate the degeneracy of the importance weight. Differently, the advantage w.r.t. the optimistic shrinkage (IS-OS) is quite significant. Instead, for the model-based estimators, we observe that our weight correction combined with the DR estimator (DR- λ^* and DR- $\hat{\lambda}$) outperforms the standard DR and its combinations with SN (SN-DR), truncation (DR-TR), and optimistic shrinkage (DR-OS). Instead, the switch estimator (DR-SW) displays a performance comparable to ours.

Estimator / ξ	10	20	50	100	200	500	1000
IS	0.6742	0.5414	0.1754	0.0326	0.0326	0.0198	0.0014
IS- λ^*	0.686	0.5416	0.176	0.0228	0.056	0	0
DR	0.6522	0.4094	0.117	0.0484	0.0378	0.0218	0.0022
DR- λ^*	0.65	0.4046	0.1088	0.03262	0.009	0	0

Table 6: Complementary cumulative distribution of the absolute error (multiplied by 100) $\mathbb{P}(E > \xi)$ in the *glass* dataset varying the number of samples n for the different estimators, when $\alpha_b = 0.9$ and $\alpha_e = 0.9999$ (5000 runs).

For the specific case of the *letter* dataset, we report in Table 5 the results obtained by setting $\alpha_b = 0.5$ and $\alpha_e = 0.9$ for different number of samples n . We notice essentially two behaviors. When the number of samples is very low (e.g., 100, 200) all estimators perform similarly, with poor performance. As n increases, the benefits of the DR-like estimators becomes more visible. In particular, the DR-SW and our corrected estimators (DR- λ^* , DR- λ^{**} , and DR- $\hat{\lambda}$) overall dominate the other baselines.

Tail Behavior Experiment We run 5000 estimation processes using the *glass* dataset, $n = 30$, $\alpha_e = 0.9999$, and $\alpha_b = 0.9$. To compare the tail behavior between vanilla weights and our correction (for both model-free and model-based estimators), we consider the absolute error random (multiplied by 100) variable E (as in Table 5) and we estimate the *complementary cumulative distribution* $\mathbb{P}(E > \xi)$. Thus, for large values of ξ , the larger $\mathbb{P}(E > \xi)$, the heavier the tail, since a larger amount of probability mass accumulates on the right of ξ . Table 6 reports the results for both model-based and model-free estimators. We observe that our corrected estimators consistently display a significantly lighter tail compared to the vanilla ones.

8.2 Off-Policy Learning

Finally, we provide an experiment in which we employ the off-policy methods to improve a baseline policy in the CMAB framework. We refer to the same setting of Section 8.1 with a uniform behavioral policy ($\alpha_b = 0$). For the target policy, we consider a Boltzmann policy in some featurization of the context $\pi_{\theta}(a|x) \propto \exp(\theta_a^T \phi(x))$. We optimize the estimated value function in the parameters θ via gradient ascent. Further details and experiments with regularized objectives are reported in Appendix B.2. We perform Off-PL on four datasets and the results for the model-free estimators are reported in Figure 3. We observe that our weight corrections (DR- λ^* and DR- $\hat{\lambda}$) outperform the considered baselines (IS, SN-IS, and IS-OS) on *ecoli* and *letter* datasets, whereas SN-IS emerges in the *glass* and *kropt* datasets. For the *letter* dataset, we report in Figure 4 the learning curve, distinguishing between model-free (left) and model-based (right) estimators.¹⁰ For the model-free ones, we observe the dominance of our estimators over SN estimator (SN-IS), while the optimistic shrinkage estimator (IS-OS) behaving similarly to ours. Interestingly, for the model-based estimators, plain DR beats the other estimators, including self-normalization that performs almost identically with our DR- λ^* , and IS-OS that fails completely to learn the task.

9 Discussion and Conclusions

In this paper, we have deepened the study of the importance sampling technique for off-policy evaluation and learning. We derived an anticoncentration bound for the vanilla IS estimator, proving polynomial concentration is tight for this setting. Then, we introduced and analyzed a class of importance weight corrections based on the intuition of smoothly shrinking the weight towards one. Assuming that the second moment of the importance weight exists, we have introduced the first transformation that achieves subgaussian concentration and maintains the differentiability of the estimator in the target policy parameters. The experimental evaluation has shown that our theoretically-grounded transformation is competitive with the traditional and modern IS baselines (including self-normalization, truncation, and optimistic shrinkage) in the CMAB framework for both evaluation and learning. The advantages of our correction are more visible in the case of continuous distributions, where the degeneracy of importance sampling is amplified. Future works include the extension of these corrections to the more challenging RL setting with continuous actions.

¹⁰Clearly, the truncated (IS-TR) and the switch (DR-SW) estimators cannot be directly employed in this setting, being non-differentiable.

References

- [1] Heejung Bang and James M Robins. Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–973, 2005.
- [2] Oliver Bembom and Mark J van der Laan. Data-adaptive selection of the truncation level for inverse-probability-of-treatment-weighted estimators. 2008.
- [3] Léon Bottou, Jonas Peters, Joaquin Quiñero Candela, Denis Xavier Charles, Max Chickering, Elon Portugaly, Dipankar Ray, Patrice Y. Simard, and Ed Snelson. Counterfactual reasoning and learning systems: the example of computational advertising. *J. Mach. Learn. Res.*, 14(1):3207–3260, 2013.
- [4] Stéphane Boucheron, Gábor Lugosi, Pascal Massart, et al. On concentration of self-bounding functions. *Electronic Journal of Probability*, 14:1884–1899, 2009.
- [5] Sébastien Bubeck, Nicolò Cesa-Bianchi, and Gábor Lugosi. Bandits with heavy tail. *IEEE Trans. Inf. Theory*, 59(11):7711–7717, 2013.
- [6] Peter S Bullen. *Handbook of means and their inequalities*, volume 560. Springer Science & Business Media, 2013.
- [7] Olivier Catoni. Challenging the empirical mean and empirical variance: a deviation study. In *Annales de l’IHP Probabilités et statistiques*, volume 48, pages 1148–1185, 2012.
- [8] Kamil Andrzej Ciosek and Shimon Whiteson. OFFER: off-environment reinforcement learning. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, pages 1819–1825. AAAI Press, 2017.
- [9] William G. Cochran. *Sampling Techniques, 3rd Edition*. John Wiley, 1977.
- [10] Stephen R Cole and Miguel A Hernán. Constructing inverse probability weights for marginal structural models. *American journal of epidemiology*, 168(6):656–664, 2008.
- [11] Corinna Cortes, Yishay Mansour, and Mehryar Mohri. Learning bounds for importance weighting. In *Advances in Neural Information Processing Systems 23 (NIPS)*, pages 442–450. Curran Associates, Inc., 2010.
- [12] Luc Devroye, Matthieu Lerasle, Gabor Lugosi, Roberto I Oliveira, et al. Sub-gaussian mean estimators. *The Annals of Statistics*, 44(6):2695–2725, 2016.
- [13] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- [14] Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning. In *Proceedings of the 28th International Conference on Machine Learning (ICML)*, pages 1097–1104. Omnipress, 2011.
- [15] M. Gil, Fady Alajaji, and Tamás Linder. Rényi divergence measures for commonly used univariate continuous distributions. *Inf. Sci.*, 249:124–131, 2013.
- [16] Alexandre Gilotte, Clément Calauzènes, Thomas Nedelec, Alexandre Abraham, and Simon Dollé. Offline A/B testing for recommender systems. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining (WSDM)*, pages 198–206. ACM, 2018.
- [17] Jinyong Hahn. On the role of the propensity score in efficient semiparametric estimation of average treatment effects. *Econometrica*, pages 315–331, 1998.
- [18] Josiah P. Hanna, Philip S. Thomas, Peter Stone, and Scott Niekum. Data-efficient policy evaluation through behavior policy search. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, volume 70 of *Proceedings of Machine Learning Research*, pages 1394–1403. PMLR, 2017.
- [19] Tim Hesterberg. Weighted average importance sampling and defensive mixture distributions. *Technometrics*, 37(2):185–194, 1995.

- [20] Timothy Classen Hesterberg. *Advances in importance sampling*. PhD thesis, Citeseer, 1988.
- [21] Daniel G Horvitz and Donovan J Thompson. A generalization of sampling without replacement from a finite universe. *Journal of the American statistical Association*, 47(260):663–685, 1952.
- [22] Peter J Huber. Robust estimation of a location parameter. In *Breakthroughs in statistics*, pages 492–518. Springer, 1992.
- [23] Edward L Ionides. Truncated importance sampling. *Journal of Computational and Graphical Statistics*, 17(2):295–311, 2008.
- [24] Mark Jerrum, Leslie G. Valiant, and Vijay V. Vazirani. Random generation of combinatorial structures from a uniform distribution. *Theor. Comput. Sci.*, 43:169–188, 1986.
- [25] Herman Kahn and Andy W Marshall. Methods of reducing sample size in monte carlo computations. *Journal of the Operations Research Society of America*, 1(5):263–278, 1953.
- [26] Oszel Kilinc, Yang Hu, and Giovanni Montana. Reinforcement learning for robotic manipulation using simulated locomotion demonstrations. *CoRR*, abs/1910.07294, 2019.
- [27] Jens Kober and Jan Peters. *Learning Motor Skills - From Algorithms to Robot Experiments*, volume 97 of *Springer Tracts in Advanced Robotics*. Springer, 2014.
- [28] Ilja Kuzborskij and Csaba Szepesvári. Efron-stein pac-bayesian inequalities. *CoRR*, abs/1909.01931, 2019.
- [29] Ilja Kuzborskij, Claire Vernade, András György, and Csaba Szepesvári. Confident off-policy evaluation and selection through self-normalized importance weighting. In *The 24th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 130 of *Proceedings of Machine Learning Research*, pages 640–648. PMLR, 2021.
- [30] John Langford and Tong Zhang. The epoch-greedy algorithm for multi-armed bandits with side information. In *Advances in Neural Information Processing Systems 20 (NIPS)*, pages 817–824. Curran Associates, Inc., 2007.
- [31] Brian K Lee, Justin Lessler, and Elizabeth A Stuart. Weight trimming and propensity score weighting. *PloS one*, 6(3):e18174, 2011.
- [32] Oleg V Lepski and Vladimir G Spokoiny. Optimal pointwise adaptive methods in nonparametric estimation. *The Annals of Statistics*, pages 2512–2546, 1997.
- [33] Lihong Li, Wei Chu, John Langford, and Xuanhui Wang. Unbiased offline evaluation of contextual-bandit-based news article recommendation algorithms. In *Proceedings of the Forth International Conference on Web Search and Web Data Mining (WSDM)*, pages 297–306. ACM, 2011.
- [34] Lihong Li, Rémi Munos, and Csaba Szepesvári. Toward minimax off-policy value estimation. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 38 of *JMLR Workshop and Conference Proceedings*. JMLR.org, 2015.
- [35] Friedrich Liese and Igor Vajda. *Convex statistical distances*, volume 95. Teubner, 1987.
- [36] Shiyin Lu, Guanghui Wang, Yao Hu, and Lijun Zhang. Optimal algorithms for lipschitz bandits with heavy-tailed rewards. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, volume 97 of *Proceedings of Machine Learning Research*, pages 4154–4163. PMLR, 2019.
- [37] Gabor Lugosi and Shahar Mendelson. Mean estimation and regression under heavy-tailed distributions: A survey. *Found. Comput. Math.*, 19(5):1145–1190, 2019.
- [38] Ashique Rupam Mahmood, Hado van Hasselt, and Richard S. Sutton. Weighted importance sampling for off-policy learning with linear function approximation. In *Advances in Neural Information Processing Systems 27 (NIPS)*, pages 3014–3022, 2014.

- [39] Alberto Maria Metelli, Matteo Papini, Pierluca D’Oro, and Marcello Restelli. Policy optimization as online learning with mediator feedback. In *Thirty-Fifth AAAI Conference on Artificial Intelligence*, pages 8958–8966. AAAI Press, 2021.
- [40] Alberto Maria Metelli, Matteo Papini, Francesco Faccio, and Marcello Restelli. Policy optimization via importance sampling. In *Advances in Neural Information Processing Systems 31 (NeurIPS)*, pages 5447–5459, 2018.
- [41] Alberto Maria Metelli, Matteo Papini, Nico Montali, and Marcello Restelli. Importance sampling techniques for policy optimization. *J. Mach. Learn. Res.*, 21:141:1–141:75, 2020.
- [42] Alberto Maria Metelli, Alessio Russo, and Marcello Restelli. Subgaussian importance sampling for off-policy evaluation and learning. *ICML-21 Workshop on Reinforcement Learning Theory*, 2021.
- [43] John E. Moody and Matthew Saffell. Learning to trade via direct reinforcement. *IEEE Trans. Neural Networks*, 12(4):875–889, 2001.
- [44] Susan A Murphy, Mark J van der Laan, James M Robins, and Conduct Problems Prevention Research Group. Marginal mean models for dynamic regimes. *Journal of the American Statistical Association*, 96(456):1410–1423, 2001.
- [45] Arkadij Semenovič Nemirovskij and David Borisovich Yudin. Problem complexity and method efficiency in optimization. 1983.
- [46] Art B. Owen. *Monte Carlo theory, methods and examples*. 2013.
- [47] Matteo Papini, Alberto Maria Metelli, Lorenzo Lupo, and Marcello Restelli. Optimistic policy optimization via multiple importance sampling. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, volume 97 of *Proceedings of Machine Learning Research*, pages 4989–4999. PMLR, 2019.
- [48] James Pickands III et al. Statistical inference using extreme order statistics. *Annals of statistics*, 3(1):119–131, 1975.
- [49] Iosif Pinelis et al. Best possible bounds of the von bahr–esseen type. *Annals of Functional Analysis*, 6(4):1–29, 2015.
- [50] Tiberiu Popoviciu. Sur les équations algébriques ayant toutes leurs racines réelles. *Mathematica*, 9:129–145, 1935.
- [51] Alfréd Rényi. On measures of entropy and information. Technical report, Hungarian Academy of Sciences Budapest Hungary, 1961.
- [52] Brian D. Ripley. *Stochastic simulation*. Wiley series in probability and mathematical statistics : applied probability and statistics. Wiley, 1987.
- [53] Yuta Saito, Aihara Shunsuke, Matsutani Megumi, and Narita Yusuke. Open bandit dataset and pipeline: Towards realistic and reproducible off-policy evaluation. *CoRR*, abs/2008.07146, 2020.
- [54] Igal Sason. On f -divergences: Integral representations, local behavior, and inequalities. *Entropy*, 20(5):383, 2018.
- [55] David Siegmund. Importance sampling in the monte carlo study of sequential tests. *The Annals of Statistics*, pages 673–684, 1976.
- [56] Yi Su, Maria Dimakopoulou, Akshay Krishnamurthy, and Miroslav Dudík. Doubly robust off-policy evaluation with shrinkage. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119 of *Proceedings of Machine Learning Research*, pages 9167–9176. PMLR, 2020.

- [57] Yi Su, Lequn Wang, Michele Santacatterina, and Thorsten Joachims. CAB: continuous adaptive blending for policy evaluation and learning. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, volume 97 of *Proceedings of Machine Learning Research*, pages 6005–6014. PMLR, 2019.
- [58] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [59] Adith Swaminathan and Thorsten Joachims. The self-normalized estimator for counterfactual learning. In *Advances in Neural Information Processing Systems 28 (NIPS)*, pages 3231–3239, 2015.
- [60] Liang Tang, Rómer Rosales, Ajit Singh, and Deepak Agarwal. Automatic ad format selection via contextual bandits. In *22nd ACM International Conference on Information and Knowledge Management (CIKM'13)*, pages 1587–1594. ACM, 2013.
- [61] Philip S. Thomas, Georgios Theodorou, and Mohammad Ghavamzadeh. High-confidence off-policy evaluation. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pages 3000–3006. AAAI Press, 2015.
- [62] John W Tukey and Donald H McLaughlin. Less vulnerable confidence and significance procedures for location based on a single sample: Trimming/winsorization 1. *Sankhyā: The Indian Journal of Statistics, Series A*, pages 331–352, 1963.
- [63] Aki Vehtari, Daniel Simpson, Andrew Gelman, Yuling Yao, and Jonah Gabry. Pareto smoothed importance sampling. *CoRR*, abs/1507.02646, 2015.
- [64] Yu-Xiang Wang, Alekh Agarwal, and Miroslav Dudík. Optimal and adaptive off-policy evaluation in contextual bandits. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, volume 70 of *Proceedings of Machine Learning Research*, pages 3589–3597. PMLR, 2017.
- [65] Changhe Yuan and Marek J. Druzdzel. How heavy should the tails be? In *Proceedings of the Eighteenth International Florida Artificial Intelligence Research Society Conference (FLAIRS)*, pages 799–805. AAAI Press, 2005.
- [66] Chao Zheng. A new principle for tuning-free huber regression. *Statistica Sinica*, 2020.
- [67] Xin Zhou, Nicole Mayer-Hamblett, Umer Khan, and Michael R Kosorok. Residual weighted learning for estimating individualized treatment rules. *Journal of the American Statistical Association*, 112(517):169–187, 2017.

Appendix

A Proofs and Derivations	16
A.1 Proofs of Section 3	16
A.2 Proofs of Section 4	18
A.3 Proofs of Section 5	18
A.4 Proofs of Section 6	21
B Experiments	26
B.1 Off-Policy Evaluation	26
B.2 Off-Policy Learning	31
C Bound Comparison and Optimization	34
C.1 Numerical Optimization of the Bound of Lemma A.3	34
D Comparison of Estimators for CMABs	36
E Analysis of the IS-OS estimator	36
F Checklist	38

A Proofs and Derivations

In this section, we report the proofs of the results that are reported in the main paper.

A.1 Proofs of Section 3

Theorem A.1. *Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions such that $P \ll Q$. For every $\alpha \in (1, 2]$ and $\delta \in (0, 1)$, with probability at least $1 - \delta$ it holds that:*

$$|\hat{\mu}_n - \mu| \leq \|f\|_\infty \left(\frac{2^{2-\alpha} I_\alpha(P\|Q)}{\delta n^{\alpha-1}} \right)^{\frac{1}{\alpha}}.$$

Proof. We first derive the following inequality concerning the α -absolute central moment of the estimator $\hat{\mu}_n$:

$$\mathbb{E}_{y_i \sim Q} [|\hat{\mu}_n - \mu|^\alpha] \leq \mathbb{E}_{y_i \sim Q} [|\hat{\mu}_n|^\alpha] \quad (\text{P.1})$$

$$\begin{aligned} &= \frac{1}{n^\alpha} \mathbb{E}_{y_i \sim Q} \left[\left| \sum_{i \in [n]} \omega(y_i) f(y_i) \right|^\alpha \right] \\ &\leq \frac{1}{n^{\alpha-1}} 2^{2-\alpha} \mathbb{E}_{y \sim Q} [|\omega(y) f(y)|^\alpha] \quad (\text{P.2}) \end{aligned}$$

$$\leq \frac{1}{n^{\alpha-1}} 2^{2-\alpha} \|f\|_\infty I_\alpha(P\|Q), \quad (\text{P.3})$$

where line (P.1) derives from observing that μ is the expected value of $\hat{\mu}_n$, line (P.2) follows from Equation (1.11) of [49] using as constant $W_\alpha = 2^{2-\alpha}$ of Proposition 1.8 of [49], and line (P.3) derives from the definition of I_α . Now we can derive the concentration inequality:

$$\begin{aligned} \mathbb{P}_{y_i \sim Q} (|\hat{\mu}_n - \mu| \geq \epsilon) &= \mathbb{P}_{y_i \sim Q} (|\hat{\mu}_n - \mu|^\alpha \geq \epsilon^\alpha) \\ &\leq \frac{\mathbb{E}_{y_i \sim Q} [|\hat{\mu}_n - \mu|^\alpha]}{\epsilon^\alpha} \quad (\text{P.4}) \\ &\leq \frac{1}{n^{\alpha-1} \epsilon^\alpha} 2^{2-\alpha} \|f\|_\infty I_\alpha(P\|Q), \end{aligned}$$

where line (P.4) derives from Markov's inequality. By setting the right hand side of the last equation equal to δ , we get the result. \square

Theorem 3.1. *There exist two distributions $P, Q \in \mathcal{P}(\mathcal{Y})$ with $P \ll Q$ and a bounded measurable function $f: \mathcal{Y} \rightarrow \mathbb{R}$ such that for every $\alpha \in (1, 2]$ and $\delta \in (0, e^{-1})$ if $n \geq \delta e \max\{1, (I_\alpha(P\|Q) - 1)^{\frac{1}{\alpha-1}}\}$, with probability at least δ it holds that:*

$$|\hat{\mu}_n - \mu| \geq \|f\|_\infty \left(\frac{I_\alpha(P\|Q) - 1}{\delta n^{\alpha-1}} \right)^{\frac{1}{\alpha}} \left(1 - \frac{e\delta}{n} \right)^{\frac{n-1}{\alpha}}.$$

Proof. The proof is inspired to that of Proposition 6.2 of [7]. We construct a function f and two probability measures P and Q that fulfill the inequality. Let $a > 0$, we consider $\mathcal{Y} = \{-a, 0, a\}$ and $f(y) = y$. First of all, we observe that $a = \|f\|_\infty$. We now define the probability distributions as follows, for $p, q \in [0, 1]$:

$$P(\{-a\}) = P(\{a\}) = \frac{p}{2} \text{ and } P(\{0\}) = 1 - p,$$

$$Q(\{-a\}) = Q(\{a\}) = \frac{q}{2} \text{ and } Q(\{0\}) = 1 - q.$$

We immediately observe that $\mathbb{E}_{y \sim P}[f(y)] = \mathbb{E}_{y \sim Q}[f(y)] = 0$. We select the values p and q as follows, for any $\alpha \in (1, 2]$:

$$\begin{aligned} q &= \left(\frac{a}{n\epsilon} \right)^\alpha \xi, \\ p &= \left(\frac{a}{n\epsilon} \right)^{\alpha-1} \xi, \end{aligned}$$

where $\xi > 0$ will be specified later. First of all, we note that to make these probability valid, we need to enforce:

$$p \leq 1 \implies n \geq \frac{a}{\epsilon} \xi^{\frac{1}{\alpha}}, \quad (\text{P.5})$$

$$q \leq 1 \implies n \geq \frac{a}{\epsilon} \xi^{\frac{1}{\alpha-1}}. \quad (\text{P.6})$$

This choice of p and q ensures that $a \frac{p}{q} = n\epsilon$. Let us now compute the divergence:

$$\begin{aligned} I_\alpha(P\|Q) &= 2 \left(\frac{p}{2}\right)^\alpha \left(\frac{q}{2}\right)^{1-\alpha} + (1-p)^\alpha (1-q)^{1-\alpha} \\ &= p^\alpha q^{1-\alpha} + (1-p)^\alpha (1-q)^{1-\alpha} \\ &= \xi + \left(1 - \xi \left(\frac{a}{n\epsilon}\right)^{\alpha-1}\right)^\alpha \left(1 - \xi \left(\frac{a}{n\epsilon}\right)^\alpha\right)^{1-\alpha} \leq \xi + 1, \end{aligned}$$

where the last inequality is obtained by upper bounding the second addendum under the assumption that $n \geq \frac{a}{\epsilon} \xi^{\frac{1}{\alpha-1}}$:

$$\left(1 - \xi \left(\frac{a}{n\epsilon}\right)^{\alpha-1}\right)^\alpha \left(1 - \xi \left(\frac{a}{n\epsilon}\right)^\alpha\right)^{1-\alpha} \leq \left(1 - \xi \left(\frac{a}{n\epsilon}\right)^{\alpha-1}\right)^\alpha \left(1 - \xi \left(\frac{a}{n\epsilon}\right)^{\alpha-1}\right)^{1-\alpha} = 1 - \xi \left(\frac{a}{n\epsilon}\right)^{\alpha-1} \leq 1.$$

Thus, we select $\xi = I_\alpha(P\|Q) - 1$. Let us now consider the vanilla IS estimator $\hat{\mu}_n$, whose expectation is $\mu = 0$, and the following derivation:

$$\begin{aligned} \mathbb{P}_{y_i \sim Q} (|\hat{\mu}_n - \mu| > \epsilon) &= \mathbb{P}_{y_i \sim Q} (\{\hat{\mu}_n - \mu < -\epsilon\} \cup \{\hat{\mu}_n - \mu > \epsilon\}) \\ &= \mathbb{P}_{y_i \sim Q} (\hat{\mu}_n - \mu < -\epsilon) + \mathbb{P}_{y_i \sim Q} (\hat{\mu}_n - \mu > \epsilon) \end{aligned} \quad (\text{P.7})$$

$$= 2 \mathbb{P}_{y_i \sim Q} (\hat{\mu}_n - \mu > \epsilon), \quad (\text{P.8})$$

where line (P.7) is obtained by observing that the two events are disjoint and line (P.8) comes from the symmetry of the events. We now lower bound the probability:

$$\begin{aligned} \mathbb{P}_{y_i \sim Q} (\hat{\mu}_n - \mu > \epsilon) &\geq \mathbb{P}_{y_i \sim Q} (\text{among the } n \text{ samples, one is } a \text{ and the remaining are } 0) \\ &= n \frac{q}{2} (1-q)^{n-1} \\ &= \frac{1}{2} \left(\frac{a}{\epsilon}\right)^\alpha n^{1-\alpha} \xi \left(1 - \left(\frac{a}{n\epsilon}\right)^\alpha \xi\right)^{n-1}. \end{aligned}$$

Now, we derive a value of $\epsilon > 0$ such that the inequality holds with probability at least δ . We enforce the condition:

$$\frac{1}{2} \left(\frac{a}{\epsilon}\right)^\alpha n^{1-\alpha} \xi \left(1 - \left(\frac{a}{n\epsilon}\right)^\alpha \xi\right)^{n-1} \leq \delta \implies \epsilon \geq a \left(\frac{\xi}{\delta n^{\alpha-1}}\right)^{\frac{1}{\alpha}} \left(1 - \left(\frac{a}{n\epsilon}\right)^\alpha \xi\right)^{\frac{n-1}{\alpha}}. \quad (\text{P.9})$$

We claim that, for $\delta \in (0, e^{-1})$, any value of ϵ fulfilling condition (P.9) must be $\epsilon \leq \epsilon^*$:

$$\epsilon^* = a \left(\frac{\xi}{\delta n^{\alpha-1}}\right)^{\frac{1}{\alpha}} \left(1 - \frac{e\delta}{n}\right)^{\frac{n-1}{\alpha}}$$

Indeed, we have:

$$\begin{aligned} a \left(\frac{\xi}{\delta n^{\alpha-1}}\right)^{\frac{1}{\alpha}} \left(1 - \left(\frac{a}{n\epsilon^*}\right)^\alpha \xi\right)^{\frac{n-1}{\alpha}} &= a \left(\frac{\xi}{\delta n^{\alpha-1}}\right)^{\frac{1}{\alpha}} \left(1 - \left(\frac{a}{n}\right)^\alpha \left(a \left(\frac{\xi}{\delta n^{\alpha-1}}\right)^{\frac{1}{\alpha}} \left(1 - \frac{e\delta}{n}\right)^{\frac{n-1}{\alpha}}\right)^\alpha\right)^{\frac{n-1}{\alpha}} \\ &= a \left(\frac{\xi}{\delta n^{\alpha-1}}\right)^{\frac{1}{\alpha}} \left(1 - \frac{\delta}{n} \left(1 - \frac{e\delta}{n}\right)^{-(n-1)}\right)^{\frac{n-1}{\alpha}} \\ &\geq a \left(\frac{\xi}{\delta n^{\alpha-1}}\right)^{\frac{1}{\alpha}} \left(1 - \frac{\delta e}{n}\right)^{\frac{n-1}{\alpha}} = \epsilon^*, \end{aligned}$$

where the last inequality derives from observing that $\left(1 - \frac{e\delta}{n}\right)^{-(n-1)} \leq e$ if $\delta \in (0, e^{-1})$. Finally, we rephrase conditions (P.5) and (P.6):

$$n \geq \frac{a}{\epsilon^*} \xi^{\frac{1}{\alpha}} \implies n \geq n^{1-\frac{1}{\alpha}} \delta^{\frac{1}{\alpha}} \left(1 - \frac{e\delta}{n}\right)^{-\frac{n-1}{\alpha}} \implies n \geq \delta e,$$

$$n \geq \frac{a}{\epsilon^*} \xi^{\frac{1}{\alpha-1}} \implies n \geq n^{1-\frac{1}{\alpha}} \delta^{\frac{1}{\alpha}} \xi^{\frac{1}{\alpha(\alpha-1)}} \left(1 - \frac{e\delta}{n}\right)^{-\frac{n-1}{\alpha}} \implies n \geq \delta e \xi^{\frac{1}{\alpha-1}},$$

having observed, again, that $\left(1 - \frac{e\delta}{n}\right)^{-\frac{n-1}{\alpha}} \leq e^{\frac{1}{\alpha}}$. Thus, we enforce the condition $n \geq \delta e \max\left\{1, \xi^{\frac{1}{\alpha-1}}\right\}$. \square

Corollary A.1. *There exist two distributions $P, Q \in \mathcal{P}(\mathcal{Y})$ with $P \ll Q$ and a bounded measurable function $f: \mathcal{Y} \rightarrow \mathbb{R}$ such that for every $\alpha \in (1, 2]$ it holds that:*

$$\mathbb{E}_{y_i \sim Q} [|\hat{\mu}_n - \mu|^\alpha] \geq \|f\|_\infty^\alpha \frac{I_\alpha(P\|Q) - 1}{n^{\alpha-1}}.$$

Proof. Let us denote the bad event:

$$\mathcal{E} = \left\{ |\hat{\mu}_n - \mu| \geq \|f\|_\infty \left(\frac{I_\alpha(P\|Q) - 1}{\delta n^{\alpha-1}} \right)^{\frac{1}{\alpha}} \left(1 - \frac{\epsilon \delta}{n} \right)^{\frac{n-1}{\alpha}} \right\}$$

From Theorem 3.1, we know that $\mathbb{P}_{y_i \sim Q}(\mathcal{E}) \geq \delta$. Let us consider the expected absolute error with exponent $\alpha \in (1, 2]$ and apply the law of total expectation:

$$\begin{aligned} \mathbb{E}_{y_i \sim Q} [|\hat{\mu}_n - \mu|^\alpha] &= \mathbb{E}_{y_i \sim Q} [(\hat{\mu}_n - \mu)^\alpha | \mathcal{E}] \mathbb{P}_{y_i \sim Q}(\mathcal{E}) + \mathbb{E}_{y_i \sim Q} [(\hat{\mu}_n - \mu)^\alpha | \mathcal{E}^c] \mathbb{P}_{y_i \sim Q}(\mathcal{E}^c) \\ &\geq \|f\|_\infty^\alpha \frac{I_\alpha(P\|Q) - 1}{\delta n^{\alpha-1}} \left(1 - \frac{\epsilon \delta}{n} \right)^{n-1} \delta + 0. \end{aligned}$$

The result is obtained by setting $\delta \rightarrow 0$. □

A.2 Proofs of Section 4

Lemma 4.1. *Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions with $P \ll Q$, then for every $\lambda \in [0, 1]$ and $y \in \mathcal{Y}$ it holds that:*

- (i) *if $s \leq s'$ then $\omega_{\lambda, s}(y) \leq \omega_{\lambda, s'}(y)$;*
- (ii) *if $s < 0$ then $\omega_{\lambda, s}(y) \leq \lambda^{\frac{1}{s}}$, otherwise if $s > 0$ then $\omega_{\lambda, s}(y) \geq \lambda^{\frac{1}{s}}$;*
- (iii) *if $s < 1$ then $\mathbb{E}_{y \sim Q}[\omega_{\lambda, s}(y)] \leq 1$, otherwise if $s > 1$ then $\mathbb{E}_{y \sim Q}[\omega_{\lambda, s}(y)] \geq 1$.*

Proof. Recall that $\omega_{s, \lambda}(y)$ is the power mean of exponent s between $\omega(y)$ and 1 and weights $(1 - \lambda, \lambda)$. Consequently, (i) follows from the generalized mean inequality [6]. Let us move to (ii), if $s < 0$, we have:

$$\omega_{\lambda, s}(y) = \left((1 - \lambda)\omega(y)^s + \lambda \right)^{\frac{1}{s}} = \frac{1}{\left(\frac{1 - \lambda}{\omega(y)^{-s}} + \lambda \right)^{\frac{1}{s}}} \leq \lambda^{\frac{1}{s}}.$$

Instead for $s > 0$, we have:

$$\omega_{\lambda, s}(y) = \left((1 - \lambda)\omega(y)^s + \lambda \right)^{\frac{1}{s}} \geq \lambda^{\frac{1}{s}}.$$

Concerning (iii), let us first observe that for every $\lambda \in [0, 1]$ and $s = 1$, it holds that $\mathbb{E}_{y \sim Q}[\omega_{\lambda, 1}(y)] = 1$. Following from (i), and from the monotonicity of the expectation, we have that for $s < 1$:

$$\omega_{\lambda, s}(y) \leq \omega_{\lambda, 1}(y) \implies \mathbb{E}_{y \sim Q}[\omega_{\lambda, s}(y)] \leq \mathbb{E}_{y \sim Q}[\omega_{\lambda, 1}(y)] = 1.$$

Symmetrically, for $s > 1$ we have:

$$\omega_{\lambda, s}(y) \geq \omega_{\lambda, 1}(y) \implies \mathbb{E}_{y \sim Q}[\omega_{\lambda, s}(y)] \geq \mathbb{E}_{y \sim Q}[\omega_{\lambda, 1}(y)] = 1. \quad \square$$

A.3 Proofs of Section 5

Before going to the proofs, we introduce the following integral:

$$J_\alpha(P\|Q) = \int_{\mathcal{Y}} q(y) \left| \frac{p(y)}{q(y)} - 1 \right|^\alpha dy.$$

For $\alpha = 1$, $J_1(P\|Q)$ reduces to the total variation divergence. For general values of α , $J_\alpha(P\|Q)$ represents the χ^α -divergence [35, 54]. $J_\alpha(P\|Q)$ can be also seen as the α -absolute central moment of the importance weight $\omega(y) = p(y)/q(y)$. Consequently, we immediacy conclude that $J_\alpha(P\|Q) \leq I_\alpha(P\|Q)$. In particular, for $\alpha = 2$, we have $J_2(P\|Q) = I_2(P\|Q) - 1$.

Lemma A.1. *Let $P, Q \in \mathcal{P}(\mathcal{Y})$ two probability distributions with $P \ll Q$. For every $\lambda \in [0, 1]$, the $(\lambda, -1)$ -corrected importance weight induces a bias that can be bounded for every $\alpha \in (1, 2]$ as:*

$$\left| \mathbb{E}_{y_i \sim Q} [\hat{\mu}_{n, \lambda}] - \mu \right| \leq \|f\|_\infty \lambda^{\alpha-1} J_\alpha(P\|Q)^{\frac{1}{\alpha}} [(1 - \lambda)I_\alpha(P\|Q) + \lambda]^{1 - \frac{1}{\alpha}}.$$

Proof. Let us consider the following derivation:

$$\left| \mathbb{E}_{y_i \sim Q} [\widehat{\mu}_{n,\lambda}] - \mu \right| = \left| \mathbb{E}_{y_i \sim Q} [\widehat{\mu}_{n,\lambda} - \widehat{\mu}_n] \right| \leq \mathbb{E}_{y_i \sim Q} [|\widehat{\mu}_{n,\lambda} - \widehat{\mu}_n|] \leq \|f\|_\infty \mathbb{E}_{y \sim Q} [|\omega_\lambda(y) - \omega(y)|].$$

Thus, we have for $\alpha \in (1, 2]$:

$$\begin{aligned} \mathbb{E}_{y \sim Q} [|\omega_\lambda(y) - \omega(y)|] &= \mathbb{E}_{y \sim Q} \left[\left| \frac{\omega(y)}{1-\lambda + \lambda\omega(y)} - \omega(y) \right| \right] \\ &= \lambda \mathbb{E}_{y \sim Q} \left[\frac{|\omega(y) - 1|}{\frac{1-\lambda}{\omega(y)} + \lambda} \right] \\ &= \lambda \mathbb{E}_{y \sim Q} \left[|\omega(y) - 1| \left(\frac{1}{\frac{1-\lambda}{\omega(y)} + \lambda} \right)^{\alpha-1} \left(\frac{1}{\frac{1-\lambda}{\omega(y)} + \lambda} \right)^{2-\alpha} \right] \\ &\leq \lambda \sup_{v \geq 0} \left(\frac{1}{\frac{1-\lambda}{v} + \lambda} \right)^{2-\alpha} \mathbb{E}_{y \sim Q} \left[|\omega(y) - 1| \left(\frac{1}{\frac{1-\lambda}{\omega(y)} + \lambda} \right)^{\alpha-1} \right]. \end{aligned}$$

Concerning the first term, we observe that the function $\frac{1}{\frac{1-\lambda}{v} + \lambda}$ is monotonically increasing in v and, consequently:

$$\sup_{v \geq 0} \left(\frac{1}{\frac{1-\lambda}{v} + \lambda} \right)^{2-\alpha} = \lim_{v \rightarrow \infty} \left(\frac{1}{\frac{1-\lambda}{v} + \lambda} \right)^{2-\alpha} = \frac{1}{\lambda^{2-\alpha}}.$$

Concerning the second term, we proceed as follows:

$$\mathbb{E}_{y \sim Q} \left[|\omega(y) - 1| \left(\frac{1}{\frac{1-\lambda}{\omega(y)} + \lambda} \right)^{\alpha-1} \right] \leq \mathbb{E}_{y \sim Q} [|\omega(y) - 1|^\alpha]^{\frac{1}{\alpha}} \mathbb{E}_{y \sim Q} \left[\left(\frac{1}{\frac{1-\lambda}{\omega(y)} + \lambda} \right)^\alpha \right]^{1-\frac{1}{\alpha}} \quad (\text{P.10})$$

$$\leq \mathbb{E}_{y \sim Q} [|\omega(y) - 1|^\alpha]^{\frac{1}{\alpha}} \mathbb{E}_{y \sim Q} [((1-\lambda)\omega(y) + \lambda)^\alpha]^{1-\frac{1}{\alpha}} \quad (\text{P.11})$$

$$\leq \mathbb{E}_{y \sim Q} [|\omega(y) - 1|^\alpha]^{\frac{1}{\alpha}} \mathbb{E}_{y \sim Q} [(1-\lambda)\omega(y)^\alpha + \lambda]^{1-\frac{1}{\alpha}} \quad (\text{P.12})$$

$$= J_\alpha(P\|Q)^{\frac{1}{\alpha}} [(1-\lambda)I_\alpha(P\|Q) + \lambda]^{1-\frac{1}{\alpha}},$$

where line (P.10) derives from Hölder's inequality with exponents α and $\frac{\alpha}{\alpha-1}$, line (P.11) is obtained from the power mean inequality [6] by bounding the harmonic mean with the arithmetic mean, line (P.12) follows from Jensen's inequality having observed that the function \cdot^α is a convex function. \square

Lemma A.2. Let $P, Q \in \mathcal{P}(\mathcal{Y})$ two probability distributions with $P \ll Q$. For every $\lambda \in [0, 1]$, the $(\lambda, -1)$ -corrected importance weight induces a variance that can be bounded for every $\alpha \in (1, 2]$ as:

$$\mathbb{V}_{y_i \sim Q} [\widehat{\mu}_{n,\lambda}] \leq \frac{\|f\|_\infty^2}{n\lambda^{2-\alpha}} [(1-\lambda)I_\alpha(P\|Q) + \lambda].$$

Proof. Let us consider the following derivation:

$$\mathbb{V}_{y_i \sim Q} [\widehat{\mu}_{n,\lambda}] = \frac{1}{n} \mathbb{V}_{y \sim Q} [\omega_\lambda(y)f(y)] \leq \frac{1}{n} \mathbb{E}_{y \sim Q} [\omega_\lambda(y)^2 f(y)^2] \leq \frac{1}{n} \|f\|_\infty^2 \mathbb{E}_{y \sim Q} [\omega_\lambda(y)^2].$$

Thus, we have for $\alpha \in (1, 2]$:

$$\begin{aligned} \mathbb{E}_{y \sim Q} [\omega_\lambda(y)^2] &= \mathbb{E}_{y \sim Q} \left[\left(\frac{1}{\frac{1-\lambda}{\omega(y)} + \lambda} \right)^2 \right] \\ &= \mathbb{E}_{y \sim Q} \left[\left(\frac{1}{\frac{1-\lambda}{\omega(y)} + \lambda} \right)^\alpha \left(\frac{1}{\frac{1-\lambda}{\omega(y)} + \lambda} \right)^{2-\alpha} \right] \\ &\leq \sup_{v \geq 0} \left(\frac{1}{\frac{1-\lambda}{v} + \lambda} \right)^{2-\alpha} \mathbb{E}_{y \sim Q} \left[\left(\frac{1}{\frac{1-\lambda}{\omega(y)} + \lambda} \right)^\alpha \right] \\ &\leq \frac{1}{\lambda^{2-\alpha}} [(1-\lambda)I_\alpha(P\|Q) + \lambda], \end{aligned}$$

where the last line is obtained by employing analogous derivations as in Lemma A.1. \square

Lemma 5.1. Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions with $P \ll Q$. For every $\lambda \in [0, 1]$, the bias and variance of the $(\lambda, -1)$ -corrected importance weight can be bounded for every $\alpha \in (1, 2]$ as:

$$\left| \mathbb{E}_{y \sim Q} [\hat{\mu}_{n,\lambda}] - \mu \right| \leq \|f\|_\infty \lambda^{\alpha-1} I_\alpha(P\|Q), \quad \text{Var}_{y_i \sim Q} [\hat{\mu}_{n,\lambda}] \leq \frac{\|f\|_\infty^2}{n\lambda^{2-\alpha}} I_\alpha(P\|Q).$$

Proof. The bias result follows immediately from Lemma A.1 by recalling that $J_\alpha(P\|Q) \leq I_\alpha(P\|Q)$ and observing that $(1-\lambda)I_\alpha(P\|Q) + \lambda \leq I_\alpha(P\|Q)$ as $I_\alpha(P\|Q) \geq 1$. The variance result is obtained from Lemma A.2 by observing that $(1-\lambda)I_\alpha(P\|Q) + \lambda \leq I_\alpha(P\|Q)$ as $I_\alpha(P\|Q) \geq 1$. \square

Lemma A.3. Let $P, Q \in \mathcal{P}(\mathcal{Y})$ two probability distributions such that $P \ll Q$. Let $\{y_i\}_{i \in [n]}$ sampled independently from Q . For every $\alpha \in (1, 2]$ and $\delta \in (0, 1)$ then, for every $\lambda \in [0, 1]$, with probability at least $1 - \delta$ it holds that:

$$\begin{aligned} \hat{\mu}_{n,\lambda} - \mu &\leq \|f\|_\infty \sqrt{\frac{2 \log \frac{1}{\delta}}{n\lambda^{2-\alpha}} [(1-\lambda)I_\alpha(P\|Q) + \lambda]} + \frac{2\|f\|_\infty \log \frac{1}{\delta}}{3\lambda n} \\ &\quad + \|f\|_\infty \lambda^{\alpha-1} J_\alpha(P\|Q)^{\frac{1}{\alpha}} [(1-\lambda)I_\alpha(P\|Q) + \lambda]^{1-\frac{1}{\alpha}}. \end{aligned}$$

Proof. The proof is a straightforward application of Bernstein's inequality together with Lemma A.1 and Lemma A.2. First of all, we highlight the bias in the following decomposition:

$$\hat{\mu}_{n,\lambda} - \mu = \underbrace{\hat{\mu}_{n,\lambda} - \mathbb{E}_{y_i \sim Q} [\hat{\mu}_{n,\lambda}]}_{\text{concentration}} + \underbrace{\mathbb{E}_{y_i \sim Q} [\hat{\mu}_{n,\lambda}] - \mu}_{\text{bias}}.$$

The bias term is bounded by using Lemma A.1, while for the concentration term we apply Bernstein's inequality. Let $\delta \in (0, 1)$, with probability at least $1 - \delta$ it holds that:

$$\begin{aligned} \hat{\mu}_{n,\lambda} - \mathbb{E}_{y_i \sim Q} [\hat{\mu}_{n,\lambda}] &\leq \sqrt{2 \text{Var}_{y_i \sim Q} [\hat{\mu}_{n,\lambda}] \log \frac{1}{\delta}} + \frac{2\|\mu\|_\infty \log \frac{1}{\delta}}{3n} \\ &\leq \|f\|_\infty \sqrt{\frac{2 \log \frac{1}{\delta}}{n\lambda^{2-\alpha}} [(1-\lambda)I_\alpha(P\|Q) + \lambda]} + \frac{2\|f\|_\infty \log \frac{1}{\delta}}{3\lambda n}, \end{aligned}$$

where the last line is obtained by bounding the variance with Lemma A.2 and recalling that $\|\mu\|_\infty \leq \frac{\|f\|_\infty}{\lambda}$. \square

We discuss how to optimize this bound in λ in Appendix C. We now move to a simplified version of the bound.

Theorem 5.1. Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions such that $P \ll Q$. For every $\alpha \in (1, 2]$ and $\delta \in (0, 1)$, if we select $\lambda = \lambda_\alpha^*$ then, with probability at least $1 - \delta$ it holds that:

$$\hat{\mu}_{n,\lambda_\alpha^*} - \mu \leq \|f\|_\infty (2 + \sqrt{3}) \left(\frac{2I_\alpha(P\|Q)^{\frac{1}{\alpha-1}} \log \frac{1}{\delta}}{3(\alpha-1)^2 n} \right)^{1-\frac{1}{\alpha}}, \quad \text{with } \lambda_\alpha^* = \left(\frac{2 \log \frac{1}{\delta}}{3(\alpha-1)^2 I_\alpha(P\|Q)n} \right)^{\frac{1}{\alpha}}.$$

Proof. The derivation is analogous to that of Lemma A.3 using Bernstein's inequality and Lemma 5.1, leading to the inequality:

$$\hat{\mu}_{n,\lambda} - \mu \leq \|f\|_\infty \sqrt{\frac{2 \log \frac{1}{\delta}}{n\lambda^{2-\alpha}} I_\alpha(P\|Q)} + \frac{2 \log \frac{1}{\delta}}{3\lambda n} \|f\|_\infty + \|f\|_\infty \lambda^{\alpha-1} I_\alpha(P\|Q) \quad (\text{P.13})$$

This is a convex function of λ that can be minimized by vanishing the derivative. The derivative is actually a quadratic function in $\lambda^{\frac{2}{\alpha}}$ and its positive solution has a quite complex expression:

$$\lambda_\alpha^\# := \left(\frac{-3\alpha + \sqrt{3}\sqrt{(\alpha+2)(3\alpha-2)} + 6}{6\sqrt{2}(\alpha-1)} \right)^{\frac{2}{\alpha}} \left(\frac{\log \frac{1}{\delta}}{I_\alpha(P\|Q)n} \right)^{\frac{2}{\alpha}} \leq \left(\frac{2}{3(\alpha-1)^2} \right)^{\frac{1}{\alpha}} \left(\frac{\log \frac{1}{\delta}}{I_\alpha(P\|Q)n} \right)^{\frac{1}{\alpha}} =: \lambda_\alpha^*,$$

where the inequality holds with equality when $\alpha = 2$. By substituting this value of λ_α^* we obtain the bound:

$$\hat{\mu}_{n,\lambda_\alpha^*} - \mu \leq \|f\|_\infty \left(2 - \sqrt{3} + \alpha(-2 + \sqrt{3} + \alpha) \right) \left(\frac{2I_\alpha(P\|Q)^{\frac{1}{\alpha-1}} \log \frac{1}{\delta}}{3(\alpha-1)^2 n} \right)^{1-\frac{1}{\alpha}}$$

$$\leq \|f\|_\infty (2 + \sqrt{3}) \left(\frac{2I_\alpha(P\|Q)^{\frac{1}{\alpha-1}} \log \frac{1}{\delta}}{3(\alpha-1)^2 n} \right)^{1-\frac{1}{\alpha}},$$

having observed that $(2 - \sqrt{3} + \alpha(-2 + \sqrt{3} + \alpha))$ is a monotonically increasing function of α . \square

Remark A.1. *In the proof of Theorem 5.1, we did not consider the possibility that $\lambda_\alpha^* > 1$, that would lead to a non-valid correction parameter. We claim that this circumstance occurs for very small values of n and δ only. Indeed:*

$$\lambda_\alpha^* \leq 1 \implies n \geq \frac{2 \log \frac{1}{\delta}}{3(\alpha-1)^2 I_\alpha(P\|Q)}.$$

In any case, if it occurs that $\lambda_\alpha^* > 1$, we conventionally clip it to 1.

Proposition A.1. *Let $\lambda \in [0, 1]$. For every $y \in \mathcal{Y}$, let $\omega(y) = \frac{p_\theta(y)}{q(y)}$, for a target distribution p_θ differentiable in θ . Then, it holds that:*

$$\|\nabla_\theta \omega_\lambda(y)\|_\infty \leq \frac{1}{4\lambda} \|\nabla_\theta \log p_\theta(y)\|_\infty.$$

Proof. Let us first compute the gradient explicitly:

$$\nabla_\theta \omega_\lambda(y) = \frac{\partial \omega_\lambda}{\partial \omega}(y) \nabla_\theta \omega(y) = \frac{1-\lambda}{(1-\lambda + \lambda \omega(y))^2} \omega(y) \nabla_\theta \log p_\theta(y)$$

To get the result, we maximize the value of the following function:

$$g(v) = \frac{(1-\lambda)v}{(1-\lambda + \lambda v)^2}.$$

First of all, we observe that for $v=0$ and $v \rightarrow \infty$, the function has value 0. Thus, the maximum must lie in between. We vanish the derivative to find it:

$$\frac{\partial g(v)}{\partial v} = \frac{(1-\lambda)(1-\lambda-\lambda v)}{(1-\lambda + \lambda v)^3} = 0 \implies v^* = \frac{1}{\lambda} - 1.$$

By substituting the found value, we obtain:

$$g(v^*) = \frac{1}{4\lambda}.$$

The result is obtained by applying the L_∞ -norm. \square

A.4 Proofs of Section 6

For the sake of simplicity, we will denote with $\hat{\eta} = \lambda n^{1/4}$. We introduce the following equation:

$$h(\eta) = \eta^2 \mathbb{E}_{y \sim Q} [\omega_\eta(y)^2] = \frac{2 \log \frac{1}{\delta}}{3\sqrt{n}},$$

and we denote with η^\dagger a solution of this equation. We introduce the corresponding empirical version, that is equivalent to Equation (3):

$$\hat{h}(\eta) = \frac{\eta^2}{n} \sum_{i \in [n]} \omega_\eta(y_i)^2 = \frac{2 \log \frac{1}{\delta}}{3\sqrt{n}},$$

having $\hat{\eta}$ as solution. Clearly, we have $\mathbb{E}_{y_i \sim Q} [\hat{h}(\eta)] = h(\eta)$.

Lemma A.4. *Let $h(\eta) = \eta^2 \mathbb{E}_{y \sim Q} [\omega_\eta(y)^2]$. The following properties hold:*

- (i) *for every $\eta \in [0, 1]$ we have $h(\eta) \in [0, 1]$;*
- (ii) *for every $c \in (0, 1]$, the equation $h(\eta) = c$ admits at most one solution.*

Proof. For (i) we immediately observe that $h(\eta) \geq 0$. Moreover, we have $\omega_\eta(y) \leq \eta^{-1}$, from which the result follows. For (ii) we show that $h(\eta)$ is monotonically increasing in η :

$$\frac{\partial h}{\partial \eta}(\eta) = 2\eta \mathbb{E}_{y \sim Q} \left[\frac{\omega(y)^2}{(1-\eta + \eta \omega(y))^3} \right] > 0.$$

\square

Remark A.2. It might be the case that the equation $\hat{h}(\eta) = \frac{2\log \frac{1}{\delta}}{3\sqrt{n}}$ admits no solution, for instance when $\frac{2\log \frac{1}{\delta}}{3\sqrt{n}} > 1$ or when $\sup_{\eta \in [0,1]} \hat{h}(\eta) < 1$. In these cases, we conventionally set the solution $\eta^\dagger = 1$. We stress that this circumstance occurs only for small values of n , as in Remark A.1. Indeed, the right hand side $\frac{2\log \frac{1}{\delta}}{3\sqrt{n}} \rightarrow 0$ when $n \rightarrow \infty$.

Lemma A.5. Let $h(\eta) = \eta^2 \mathbb{E}_{y \sim Q} [\omega_\eta(y)^2]$. Let $\eta^\dagger \in [0, 1]$ such that:

$$h(\eta^\dagger) = \frac{2\log \frac{1}{\delta}}{3\sqrt{n}} \quad \text{and} \quad \lambda^\dagger = \eta^\dagger n^{-1/4}$$

then it holds that:

$$\lambda_2^* \leq \lambda^\dagger \leq \sqrt{2}\lambda_2^*,$$

where the second inequality holds if $n \geq \frac{4096(I_3(P\|Q) - I_2(P\|Q))^4 (\log \frac{1}{\delta})^2}{9I_2(P\|Q)^6}$, whenever $I_3(P\|Q)$ is finite.

Proof. Let us first observe that:

$$\mathbb{E}_{y \sim Q} [\omega_\eta(y)^2] = \mathbb{E}_{y \sim Q} \left[\frac{1}{\left(\frac{1-\eta}{\omega(y)} + \eta\right)^2} \right] \leq \mathbb{E}_{y \sim Q} [((1-\eta)\omega(y) + \eta)^2] = (I_2(P\|Q) - 1)\eta^2 + 1 \leq I_2(P\|Q),$$

where the first inequality derives from the inequality between the harmonic and arithmetic mean. From the last inequality, we have:

$$h(\eta) \leq \eta^2 I_2(P\|Q) \implies \eta^\dagger \geq \sqrt{\frac{2\log \frac{1}{\delta}}{3I_2(P\|Q)\sqrt{n}}} \implies \lambda^\dagger = \sqrt{\frac{2\log \frac{1}{\delta}}{3I_2(P\|Q)n}} = \lambda_2^*.$$

Concerning the lower bound, we proceed with a second order Taylor expansion centered in $\eta = 0$:

$$\frac{1}{\left(\frac{1-\eta}{\omega(y)} + \eta\right)^2} = \omega(y)^2 - 2\omega(y)^2(\omega(y) - 1)\eta + 3(\omega(y) - 1)^2\omega(y)^2\eta^2 \geq \omega(y)^2 - 2\omega(y)^2(\omega(y) - 1)\eta,$$

for some $\bar{\eta} \in [0, \eta]$. From which, we obtain:

$$\mathbb{E}_{y \sim Q} \left[\frac{1}{\left(\frac{1-\eta}{\omega(y)} + \eta\right)^2} \right] \geq \mathbb{E}_{y \sim Q} [\omega(y)^2 - 2\omega(y)^2(\omega(y) - 1)\eta] = I_2(P\|Q) - 2\eta(I_3(P\|Q) - I_2(P\|Q)).$$

By moving to function $h(\eta)$, and recalling the equation $h(\eta) = \frac{2\log \frac{1}{\delta}}{3\sqrt{n}}$, we have:

$$\begin{aligned} h(\eta) &= \eta^2 \mathbb{E}_{y \sim Q} [\omega_\eta(y)^2] \geq \eta^2 I_2(P\|Q) - 2\eta^3 (I_3(P\|Q) - I_2(P\|Q)) \\ &\implies \eta^2 I_2(P\|Q) - 2\eta^3 (I_3(P\|Q) - I_2(P\|Q)) \leq \frac{2\log \frac{1}{\delta}}{3\sqrt{n}}. \end{aligned}$$

We prove that for sufficiently large n , all solutions η^\dagger of the previous inequality satisfy $\eta \leq \sqrt{\frac{4\log \frac{1}{\delta}}{3I_2(P\|Q)\sqrt{n}}}$:

$$\begin{aligned} \frac{4\log \frac{1}{\delta}}{3I_2(P\|Q)\sqrt{n}} I_2(P\|Q) - 2 \left(\frac{4\log \frac{1}{\delta}}{3I_2(P\|Q)\sqrt{n}} \right)^{\frac{3}{2}} (I_3(P\|Q) - I_2(P\|Q)) &> \frac{2\log \frac{1}{\delta}}{3\sqrt{n}} \\ \implies n &\geq \frac{4096(I_3(P\|Q) - I_2(P\|Q))^4 (\log \frac{1}{\delta})^2}{9I_2(P\|Q)^6}. \end{aligned}$$

This, implies that $\lambda^\dagger \leq \sqrt{\frac{4\log \frac{1}{\delta}}{3I_2(P\|Q)n}} = \sqrt{2}\lambda_2^*$. □

Lemma A.6. Let $h(\eta) = \eta^2 \mathbb{E}_{y \sim Q} [\omega_\eta(y)^2]$, then it holds that:

$$\frac{\partial h(\eta)}{\partial \eta^2} \geq I_2(P\|Q)^{-2}.$$

Proof. Let us first observe that:

$$\frac{\partial h(\eta)}{\partial \eta^2} = \frac{\partial h(\eta)}{\partial \eta} \frac{\partial \eta}{\partial \eta^2} = \frac{\partial h(\eta)}{\partial \eta} \frac{1}{2\eta}.$$

The first factor was already computed in the proof of Lemma A.4. We now lower bound it. Let us first prove the following auxiliary inequality:

$$\begin{aligned} 1 = \mathbb{E}_{y \sim Q} [\omega(y)]^2 &= \mathbb{E}_{y \sim Q} \left[\frac{\omega(y)}{1 - \lambda + \lambda\omega(y)} (1 - \lambda + \lambda\omega(y)) \right]^2 \leq \mathbb{E}_{y \sim Q} \left[\frac{\omega(y)^2}{(1 - \lambda + \lambda\omega(y))^2} \right] \mathbb{E}_{y \sim Q} [(1 - \lambda + \lambda\omega(y))^2] \\ &\leq \mathbb{E}_{y \sim Q} \left[\frac{\omega(y)^2}{(1 - \lambda + \lambda\omega(y))^2} \right] I_2(P\|Q), \end{aligned} \quad (\text{P.14})$$

where the first inequality follows from Cauchy-Schwarz's and the second one by recalling that $\mathbb{E}_{y \sim Q} [(1 - \lambda + \lambda\omega(y))^2] \leq I_2(P\|Q)$. Now, we proceed with Hölder's inequality with $p = \frac{3}{2}$ and $q = 3$:

$$\begin{aligned} \mathbb{E}_{y \sim Q} \left[\frac{\omega(y)^2}{(1 - \lambda + \lambda\omega(y))^2} \right] &= \mathbb{E}_{y \sim Q} \left[\frac{\omega(y)^{\frac{4}{3}}}{(1 - \lambda + \lambda\omega(y))^2} \omega(y)^{\frac{2}{3}} \right] \leq \mathbb{E}_{y \sim Q} \left[\frac{\omega(y)^2}{(1 - \lambda + \lambda\omega(y))^3} \right]^{\frac{2}{3}} \mathbb{E}_{y \sim Q} [\omega(y)^2]^{\frac{1}{3}} \\ &= \mathbb{E}_{y \sim Q} \left[\frac{\omega(y)^2}{(1 - \lambda + \lambda\omega(y))^3} \right]^{\frac{2}{3}} I_2(P\|Q)^{\frac{1}{3}}. \end{aligned} \quad (\text{P.15})$$

Putting together Equation (P.14) and Equation (P.15), we have:

$$\mathbb{E}_{y \sim Q} \left[\frac{\omega(y)^2}{(1 - \lambda + \lambda\omega(y))^3} \right] \geq \mathbb{E}_{y \sim Q} \left[\frac{\omega(y)^2}{(1 - \lambda + \lambda\omega(y))^2} \right]^{\frac{3}{2}} I_2(P\|Q)^{-\frac{1}{2}} \geq I_2(P\|Q)^{-2}.$$

□

Lemma A.7. Let $h(\eta) = \eta^2 \mathbb{E}_{y \sim Q} [\omega_\eta(y)^2]$ and $\hat{h}(\eta) = \frac{\eta^2}{n} \sum_{i \in [n]} \omega_\eta(y_i)^2$. Then, $n\hat{h}(\eta)$ is a self-bounding function. Therefore, for every $\eta \in [0, 1]$ it holds that:

$$\Pr_{y_i \sim Q} \left(\hat{h}(\eta) - h(\eta) \geq \epsilon \right) \leq \exp \left(\frac{-\epsilon^2 n}{2(h(\eta) + \frac{\epsilon}{3})} \right) \quad \text{with } \epsilon > 0, \quad (4)$$

$$\Pr_{y_i \sim Q} \left(h(\eta) - \hat{h}(\eta) \geq \epsilon \right) \leq \exp \left(\frac{-\epsilon^2 n}{2h(\eta)} \right) \quad \text{with } 0 < \epsilon < h(\eta). \quad (5)$$

Proof. We consider the definition of self-bounding function provided in [4, Definition 1]. We denote with $n\hat{h}^{k,z}(\eta)$ the function obtained from $n\hat{h}(\eta)$ by replacing $\omega(y_k)$ with $z \geq 0$. We show that $n\hat{h}(\eta)$ satisfies both conditions:

$$\begin{aligned} n\hat{h}(\eta) - n\hat{h}^{k,z}(\eta) &= \eta^2 (\omega_\eta(y_k)^2 - z^2) \leq \eta^2 \omega_\eta(y_k)^2 \leq 1, \\ \sum_{k \in [n]} \left(n\hat{h}(\eta) - n\hat{h}^{k,z}(\eta) \right)^2 &= \sum_{k \in [n]} (\omega_\eta(y_k)^2 - z^2)^2 \leq \sum_{k \in [n]} (\eta^2 \omega_\eta(y_k)^2)^2 \leq \sum_{k \in [n]} \eta^2 \omega_\eta(y_k)^2 = n\hat{h}(\eta). \end{aligned}$$

having observed that $\eta \omega_\eta(y_k) \leq 1$. By applying the concentration inequalities for the self-bounding functions [4], we obtain that for every $\eta \in [0, 1]$ and $\epsilon > 0$ it holds that:

$$\Pr_{y_i \sim Q} \left(\hat{h}(\eta) - h(\eta) \geq \epsilon \right) \leq \exp \left(\frac{-\epsilon^2 n}{2(h(\eta) + \frac{\epsilon}{3})} \right).$$

Similarly, for every $\eta \in [0, 1]$ and $0 < \epsilon < h(\eta)$ it holds that:

$$\Pr_{y_i \sim Q} \left(h(\eta) - \hat{h}(\eta) \geq \epsilon \right) \leq \exp \left(\frac{-\epsilon^2 n}{2h(\eta)} \right).$$

□

Lemma A.8. Let η^\dagger be the solution of $h(\eta^\dagger) = \frac{2 \log \frac{1}{\delta}}{3\sqrt{n}}$ and $\hat{\eta}$ be the solution of $\hat{h}(\hat{\eta}) = \frac{2 \log \frac{1}{\delta}}{3\sqrt{n}}$. Then, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ it holds that:

$$\frac{1}{2} \leq \frac{\hat{\eta}}{\eta^\dagger} \leq \sqrt{2} \quad \text{and} \quad \frac{1}{2} \leq \frac{\hat{\lambda}}{\lambda^\dagger} \leq \sqrt{2},$$

$$\text{for } n \geq \max \left\{ 544 I_2(P\|Q)^{12} \left(\frac{\log \frac{2}{\delta}}{\log \frac{1}{\delta}} \right)^2, \frac{4096 (I_3(P\|Q) - I_2(P\|Q))^4 (\log \frac{1}{\delta})^2}{9 I_2(P\|Q)^6} \right\}.$$

Proof. Let $\epsilon \in [0, 1]$, consider the event $\left\{ \left| \frac{\hat{\eta}}{\eta^\dagger} - 1 \right| > \epsilon \right\}$. Under the sub-event $\{\hat{\eta} > (1 + \epsilon)\eta^\dagger\}$ recalling that function h and \hat{h} are increasing in η we have:

$$\begin{aligned} \hat{h}(\hat{\eta}) - \hat{h}(\eta^\dagger) &\geq \hat{h}((1 + \epsilon)\eta^\dagger) - \hat{h}(\eta^\dagger) \\ &= \hat{h}((1 + \epsilon)\eta^\dagger) - \hat{h}(\eta^\dagger) \pm h(\eta^\dagger) \pm h((1 + \epsilon)\eta^\dagger) \\ &= \hat{h}((1 + \epsilon)\eta^\dagger) - h((1 + \epsilon)\eta^\dagger) + h(\eta^\dagger) - \hat{h}(\eta^\dagger) + h((1 + \epsilon)\eta^\dagger) - h(\eta^\dagger) \\ &\geq \hat{h}((1 + \epsilon)\eta^\dagger) - h((1 + \epsilon)\eta^\dagger) + h(\eta^\dagger) - \hat{h}(\eta^\dagger) + 2I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2, \end{aligned}$$

where the last inequality follows from Lemma A.6 having applied:

$$h((1 + \epsilon)\eta^\dagger) - h(\eta^\dagger) \geq I_2(P\|Q)^{-2}((1 + \epsilon)^2 - 1)(\eta^\dagger)^2 = I_2(P\|Q)^{-2}(2 + \epsilon)\epsilon(\eta^\dagger)^2 \geq 2I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2.$$

Recalling that $\hat{h}(\hat{\eta}) = h(\eta^\dagger)$, the condition can be further simplified into $h((1 + \epsilon)\eta^\dagger) - \hat{h}((1 + \epsilon)\eta^\dagger) \geq 2I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2$. Symmetrically, under the sub-event $\{\hat{\eta} < (1 - \epsilon)\eta^\dagger\}$ we have:

$$\begin{aligned} \hat{h}(\hat{\eta}) - \hat{h}(\eta^\dagger) &\leq \hat{h}((1 - \epsilon)\eta^\dagger) - \hat{h}(\eta^\dagger) \\ &= \hat{h}((1 - \epsilon)\eta^\dagger) - \hat{h}(\eta^\dagger) \pm h(\eta^\dagger) \pm h((1 - \epsilon)\eta^\dagger) \\ &= \hat{h}((1 - \epsilon)\eta^\dagger) - h((1 - \epsilon)\eta^\dagger) + h(\eta^\dagger) - \hat{h}(\eta^\dagger) + h((1 - \epsilon)\eta^\dagger) - h(\eta^\dagger) \\ &\leq \hat{h}((1 - \epsilon)\eta^\dagger) - h((1 - \epsilon)\eta^\dagger) + h(\eta^\dagger) - \hat{h}(\eta^\dagger) - I_2(P\|Q)^{-2}(1 - (1 - \epsilon)^2)(\eta^\dagger)^2, \end{aligned}$$

that can be simplified, as before, into the condition $\hat{h}((1 - \epsilon)\eta^\dagger) - h((1 - \epsilon)\eta^\dagger) \geq I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2$ since $1 - (1 - \epsilon)^2 = \epsilon(2 - \epsilon) \geq \epsilon$ being $\epsilon < 1$. Thus, we have:

$$\begin{aligned} \Pr_{y \sim Q} \left(\left| \frac{\hat{\eta}}{\eta^\dagger} - 1 \right| > \epsilon \right) &= \Pr_{y \sim Q} \left(\hat{\eta} > (1 + \epsilon)\eta^\dagger \right) + \Pr_{y \sim Q} \left(\hat{\eta} < (1 - \epsilon)\eta^\dagger \right) \\ &\leq \Pr_{y \sim Q} \left(h((1 + \epsilon)\eta^\dagger) - \hat{h}((1 + \epsilon)\eta^\dagger) \geq 2I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2 \right) \\ &\quad + \Pr_{y \sim Q} \left(\hat{h}((1 - \epsilon)\eta^\dagger) - h((1 - \epsilon)\eta^\dagger) \geq I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2 \right). \end{aligned}$$

First of all, we observe that $h((1 + \epsilon)\eta^\dagger) = (1 + \epsilon)^2(\eta^\dagger)^2 \mathbb{E}_{y \sim Q}[\omega_{(1 + \epsilon)\eta^\dagger}(y)^2] \leq 4(\eta^\dagger)^2 I_2(P\|Q)$. Now, recalling that function h is self-bounding as proved in Lemma A.7, we have by Equation (5):

$$\begin{aligned} \Pr \left(h((1 + \epsilon)\eta^\dagger) - \hat{h}((1 + \epsilon)\eta^\dagger) \geq 2I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2 \right) &\leq \exp \left(\frac{-4I_2(P\|Q)^{-4}\epsilon^2(\eta^\dagger)^4 n}{2h((1 + \epsilon)\eta^\dagger)} \right) \\ &\leq \exp \left(\frac{-4I_2(P\|Q)^{-4}\epsilon^2(\eta^\dagger)^4 n}{8(\eta^\dagger)^2 I_2(P\|Q)} \right) \\ &= \exp \left(\frac{-\epsilon^2(\eta^\dagger)^2 n}{2I_2(P\|Q)^5} \right), \end{aligned}$$

provided that $2I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2 \leq h((1 + \epsilon)\eta^\dagger)$, that is fulfilled for every $\epsilon \in [0, 1]$. Indeed, recalling that $h((1 + \epsilon)\eta^\dagger) = (1 + \epsilon)^2(\eta^\dagger)^2 \mathbb{E}_{y \sim Q}[\omega_{(1 + \epsilon)\eta^\dagger}(y)^2] \geq (1 + \epsilon)^2(\eta^\dagger)^2 I_2(P\|Q)^{-2}$ (from Equation (P.14)), we have that $2I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2 \leq (1 + \epsilon)^2(\eta^\dagger)^2 I_2(P\|Q)^{-2}$ is fulfilled for every $\epsilon \in [0, 1]$. Similarly, by Equation (4) and recalling that $h((1 - \epsilon)\eta^\dagger) \leq h(\eta^\dagger) \leq (\eta^\dagger)^2 I_2(P\|Q)$, we have:

$$\begin{aligned} \Pr \left(\hat{h}((1 - \epsilon)\eta^\dagger) - h((1 - \epsilon)\eta^\dagger) \geq I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2 \right) &\leq \exp \left(\frac{-I_2(P\|Q)^{-4}\epsilon^2(\eta^\dagger)^4 n}{2(h((1 - \epsilon)\eta^\dagger) + \frac{1}{3}I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2)} \right) \\ &\leq \exp \left(\frac{-I_2(P\|Q)^{-4}\epsilon^2(\eta^\dagger)^4 n}{2(\eta^\dagger)^2 I_2(P\|Q) + \frac{2}{3}I_2(P\|Q)^{-2}\epsilon(\eta^\dagger)^2} \right) \\ &\quad \exp \left(\frac{-3\epsilon^2(\eta^\dagger)^2 n}{8I_2(P\|Q)^5} \right), \end{aligned}$$

having crudely bounded $I_2(P\|Q)^{-2}\epsilon \leq I_2(P\|Q)$. Putting these inequalities together, we obtain:

$$\Pr \left(\left| \frac{\hat{\eta}}{\eta^\dagger} - 1 \right| > \epsilon \right) \leq \exp \left(\frac{-\epsilon^2(\eta^\dagger)^2 n}{2I_2(P\|Q)^5} \right) + \exp \left(\frac{-3\epsilon^2(\eta^\dagger)^2 n}{2I_2(P\|Q)^5} \right) \leq 2 \exp \left(\frac{-3\epsilon^2(\eta^\dagger)^2 n}{8I_2(P\|Q)^5} \right),$$

leading to the inequality holding with probability at least $1 - \delta$:

$$\left| \frac{\hat{\eta}}{\eta^\dagger} - 1 \right| \leq \sqrt{\frac{8I_2(P\|Q)^5 \log \frac{2}{\delta}}{3n(\eta^\dagger)^2}}.$$

Under Lemma A.5, we know that $\eta^\dagger \geq \sqrt{\frac{2 \log \frac{1}{\delta}}{3I_2(P\|Q)\sqrt{n}}}$. From which we have:

$$\left| \frac{\hat{\eta}}{\eta^\dagger} - 1 \right| \leq \sqrt{\frac{4I_2(P\|Q)^6 \log \frac{2}{\delta}}{\sqrt{n} \log \frac{1}{\delta}}}.$$

Simple calculations allow to conclude that $\frac{1}{2} \leq \frac{\hat{\eta}}{\eta^\dagger} \leq \sqrt{2}$ for $n \geq 544I_2(P\|Q)^{12} \left(\frac{\log \frac{2}{\delta}}{\log \frac{1}{\delta}} \right)^2$. \square

Theorem 6.1. *Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions such that $P \ll Q$. Let $\hat{\lambda}$ be the solution of Equation (3), then, if $I_3(P\|Q)$ is finite, for sufficiently large n , for every $\delta \in (0, 1)$, with probability at least $1 - 2\delta$ it holds that:*

$$\hat{\mu}_{n, \hat{\lambda}} - \mu \leq \|f\|_\infty \frac{5 + 2\sqrt{3}}{2} \sqrt{\frac{2I_2(P\|Q) \log \frac{1}{\delta}}{3n}}.$$

Proof. Let us start observing that if we substitute a value of λ that is proportional to λ_2^* into Equation (P.13), we are able to provide the following bound for $\beta > 0$:

$$\hat{\mu}_{n, \beta \lambda_2^*} - \mu \leq \frac{1 + \sqrt{3}\beta + \beta^2}{\beta} \sqrt{\frac{2I_2(P\|Q) \log \frac{1}{\delta}}{3n}}.$$

Now, we provide sufficient conditions so that $\frac{1}{2}\lambda_2^* \leq \hat{\lambda} \leq 2\lambda_2^*$. First of all, we know from Lemma A.5 that for sufficiently large n we have $1 \leq \frac{\lambda^\dagger}{\lambda_2^*} \leq \sqrt{2}$. Second, from Lemma A.7, we know that for sufficiently large n and with probability at least $1 - \delta$, we have $\frac{1}{2} \leq \frac{\hat{\lambda}}{\lambda^\dagger} \leq \sqrt{2}$. Thus, putting together these results we enforce $\frac{1}{2}\lambda_2^* \leq \hat{\lambda} \leq 2\lambda_2^*$. Therefore, it holds with probability at least $1 - 2\delta$ and sufficiently large n that:

$$\hat{\mu}_{n, \hat{\lambda}} - \mu \leq \frac{\|f\|_\infty}{2} (5 + 2\sqrt{3}) \sqrt{\frac{2I_2(P\|Q) \log \frac{1}{\delta}}{3n}}.$$

\square

Corollary A.2. *Let $P, Q \in \mathcal{P}(\mathcal{Y})$ two probability distributions such that $P \ll Q$. Let $\{y_i\}_{i \in [n]}$ sampled independently from Q . For every $\delta \in (0, 1)$, let*

$$\lambda^\ddagger = \sqrt{\frac{\log \frac{1}{\delta}}{n}}$$

then, with probability at least $1 - \delta$ it holds that:

$$\hat{\mu}_{n, \lambda^\ddagger} - \mu \leq \|f\|_\infty \sqrt{\frac{\log \frac{1}{\delta}}{n}} \left(\frac{2}{3} + \sqrt{2I_2(P\|Q)} + I_2(P\|Q) \right).$$

Proof. The result is simply obtained by substituting λ^\ddagger into Equation (P.13). \square

B Experiments

In this appendix, we report the experimental details and additional experimental results.

Infrastructure The experiments have been run on a machine with two CPUs Intel(R) Xeon(R) CPU E7-8880 v4 @ 2.20GHz (22 cores, 44 thread, 55 MB cache) and 128 GB RAM.

Code The code is built on top of the *Open Bandit Pipeline* [53, <https://github.com/st-tech/zr-obp>], that is licensed under the Apache 2.0 License. In the attached code, the source files that have been modified are marked with an appropriate comment at the beginning.

B.1 Off-Policy Evaluation

B.1.1 Synthetic Example

Experimental Details To accurately estimate the expectation of function f under P , we generate at the beginning 10M from P and we estimate the expectation μ with the sample mean. For all estimators with optimal parameter (truncation threshold or λ), we employ the significance value $\delta = 0.1$.

For the optimistic shrinkage transformation (IS-OS), we compute the correction parameter τ^* , by minimizing an upper bound on the MSE, derived from the one presented in the paper [56], accounting for the fact that we do not have a reward estimate (we are not considering here a DR estimator):

$$\tau^* \in \arg \min_{\tau \geq 0} \underbrace{\widehat{\text{Var}}_{y_i \sim Q} [\omega_\tau^{\text{OS}}(y_i) f(y_i)]}_{\text{estimated variance}} + \underbrace{\frac{\|f\|_\infty^2}{n} \sum_{i \in [n]} (\omega_\tau^{\text{OS}}(y_i) - \omega(y_i))^2}_{\text{estimated bias}},$$

where:

$$\widehat{\text{Var}}_{y_i \sim Q} [\omega_\tau^{\text{OS}}(y_i) r(y_i)] = \frac{1}{n} \sum_{i \in [n]} (\omega_\tau^{\text{OS}}(y_i) f(y_i) - \hat{\mu}_\tau^{\text{OS}})^2, \quad \hat{\mu}_\tau^{\text{OS}} = \frac{1}{n} \sum_{i \in [n]} \omega_\tau^{\text{OS}}(y_i) f(y_i).$$

Complete Results In all experiments, we employ $\mu_P = 0.5$ and $\mu_Q = 0$. The values of σ_P and σ_Q for the different experiments are reported in Table 7. In Table 8 and Figure 5, we report the complete results for the different settings.

σ_Q^2	σ_P^2	$I_2(P\ Q)$
1	1.5	1.904
1	1.9	27.949
1	1.99	5.104e + 11
1	1.999	8.379e + 109

Table 7: Variance values σ_Q^2 and σ_P^2 and divergence $I_2(P\|Q)$ for the different experiments.

$$\sigma_Q^2 = 1, \sigma_P^2 = 1.5$$

Estimator / n	10	20	50	100	200	500	1000
IS	23.52 ± 5.39	15.39 ± 3.26	10.06 ± 1.93	8.35 ± 0.73	6.29 ± 0.32	3.93 ± 0.12	2.54 ± 0.06
SN-IS	23.09 ± 4.62	14.37 ± 2.55	9.15 ± 1.32	8.23 ± 0.63	6.32 ± 0.31	3.96 ± 0.12	2.56 ± 0.06
IS-TR	20.34 ± 4.66	13.48 ± 2.59	8.33 ± 1.08	7.38 ± 0.47	5.88 ± 0.27	3.60 ± 0.11	2.45 ± 0.06
IS-OS	16.55 ± 4.13	11.87 ± 2.79	7.98 ± 1.21	6.53 ± 0.52	5.06 ± 0.26	3.21 ± 0.10	2.17 ± 0.05
IS- λ^*	18.86 ± 4.01	12.20 ± 2.30	7.44 ± 0.92	6.53 ± 0.43	5.14 ± 0.25	3.25 ± 0.10	2.20 ± 0.05
IS- λ^{**}	17.85 ± 3.81	11.32 ± 2.07	6.89 ± 0.79	6.00 ± 0.41	4.81 ± 0.24	3.07 ± 0.09	2.11 ± 0.05
IS- $\hat{\lambda}$	17.98 ± 3.83	11.30 ± 2.07	6.82 ± 0.77	5.89 ± 0.40	4.72 ± 0.23	3.03 ± 0.09	2.09 ± 0.05

$$\sigma_Q^2 = 1, \sigma_P^2 = 1.9$$

Estimator / n	10	20	50	100	200	500	1000
IS	27.43 ± 13.33	15.70 ± 4.83	10.89 ± 1.81	9.26 ± 0.92	12.41 ± 1.88	9.42 ± 0.68	5.84 ± 0.27
SN-IS	23.89 ± 5.77	15.62 ± 2.62	10.96 ± 1.18	9.53 ± 0.74	8.82 ± 0.62	7.48 ± 0.37	5.14 ± 0.20
IS-TR	23.47 ± 7.52	14.03 ± 2.75	10.32 ± 1.47	8.89 ± 0.79	7.68 ± 0.46	6.21 ± 0.28	4.22 ± 0.15
IS-OS	19.25 ± 8.68	10.93 ± 3.29	8.37 ± 1.35	7.06 ± 0.61	8.69 ± 1.44	6.65 ± 0.47	3.97 ± 0.16
IS- λ^*	21.75 ± 6.36	13.17 ± 2.45	9.26 ± 1.19	7.76 ± 0.62	6.53 ± 0.38	5.29 ± 0.23	3.52 ± 0.12
IS- λ^{**}	20.66 ± 4.08	12.62 ± 2.19	8.86 ± 1.08	7.39 ± 0.57	5.94 ± 0.32	4.74 ± 0.20	3.19 ± 0.10
IS- $\hat{\lambda}$	18.19 ± 3.93	10.27 ± 1.64	7.03 ± 0.75	5.79 ± 0.38	3.85 ± 0.21	2.90 ± 0.10	2.06 ± 0.05

$$\sigma_Q^2 = 1, \sigma_P^2 = 1.99$$

Estimator / n	10	20	50	100	200	500	1000
IS	24.42 ± 6.54	25.03 ± 11.38	15.72 ± 3.31	11.10 ± 1.89	8.96 ± 0.74	6.23 ± 0.32	4.77 ± 0.19
SN-IS	25.50 ± 5.84	20.36 ± 3.36	13.99 ± 1.56	9.58 ± 1.08	8.73 ± 0.56	6.08 ± 0.27	4.64 ± 0.16
IS-TR	24.42 ± 6.54	25.03 ± 11.38	15.72 ± 3.31	11.10 ± 1.89	8.96 ± 0.74	6.23 ± 0.32	4.77 ± 0.19
IS-OS	16.39 ± 4.48	16.89 ± 6.36	11.20 ± 1.96	7.66 ± 1.08	6.80 ± 0.48	4.67 ± 0.21	3.62 ± 0.14
IS- λ^*	24.42 ± 6.54	25.03 ± 11.38	15.72 ± 3.31	11.10 ± 1.89	8.96 ± 0.74	6.23 ± 0.32	4.77 ± 0.19
IS- λ^{**}	18.37 ± 4.65	12.95 ± 2.18	15.72 ± 3.31	11.10 ± 1.89	8.96 ± 0.74	6.23 ± 0.32	4.77 ± 0.19
IS- $\hat{\lambda}$	16.12 ± 4.19	12.50 ± 2.04	7.81 ± 0.77	5.19 ± 0.41	4.64 ± 0.24	2.92 ± 0.11	2.25 ± 0.05

$$\sigma_Q^2 = 1, \sigma_P^2 = 1.999$$

Estimator / n	10	20	50	100	200	500	1000
IS	32.44 ± 30.89	22.29 ± 11.21	19.03 ± 5.26	19.39 ± 4.36	15.83 ± 2.03	9.21 ± 0.50	6.96 ± 0.26
SN-IS	21.06 ± 5.75	18.00 ± 3.18	14.78 ± 2.10	11.81 ± 1.39	10.66 ± 0.89	7.94 ± 0.35	6.32 ± 0.20
IS-TR	32.44 ± 30.89	22.29 ± 11.21	19.03 ± 5.26	19.39 ± 4.36	15.83 ± 2.03	9.21 ± 0.50	6.96 ± 0.26
IS-OS	21.32 ± 18.62	15.42 ± 6.75	12.18 ± 3.00	13.50 ± 3.06	10.62 ± 1.25	6.15 ± 0.33	4.68 ± 0.16
IS- λ^*	32.44 ± 30.89	22.29 ± 11.21	19.03 ± 5.26	19.39 ± 4.36	15.83 ± 2.03	9.21 ± 0.50	6.96 ± 0.26
IS- λ^{**}	14.87 ± 3.73	12.81 ± 2.12	8.26 ± 0.91	5.07 ± 0.41	3.56 ± 0.22	2.54 ± 0.07	1.43 ± 0.03
IS- $\hat{\lambda}$	13.36 ± 3.47	11.25 ± 1.67	7.52 ± 0.85	5.27 ± 0.43	3.68 ± 0.20	2.47 ± 0.10	2.20 ± 0.05

Table 8: Absolute error in the illustrative examples varying the number of samples n for the different estimators and the different settings of Table 7 (mean ± std, 60 runs). The estimator with smallest absolute error and the ones not statistically significantly different from that one (Welch's t-test with $p < 0.02$) are in bold.

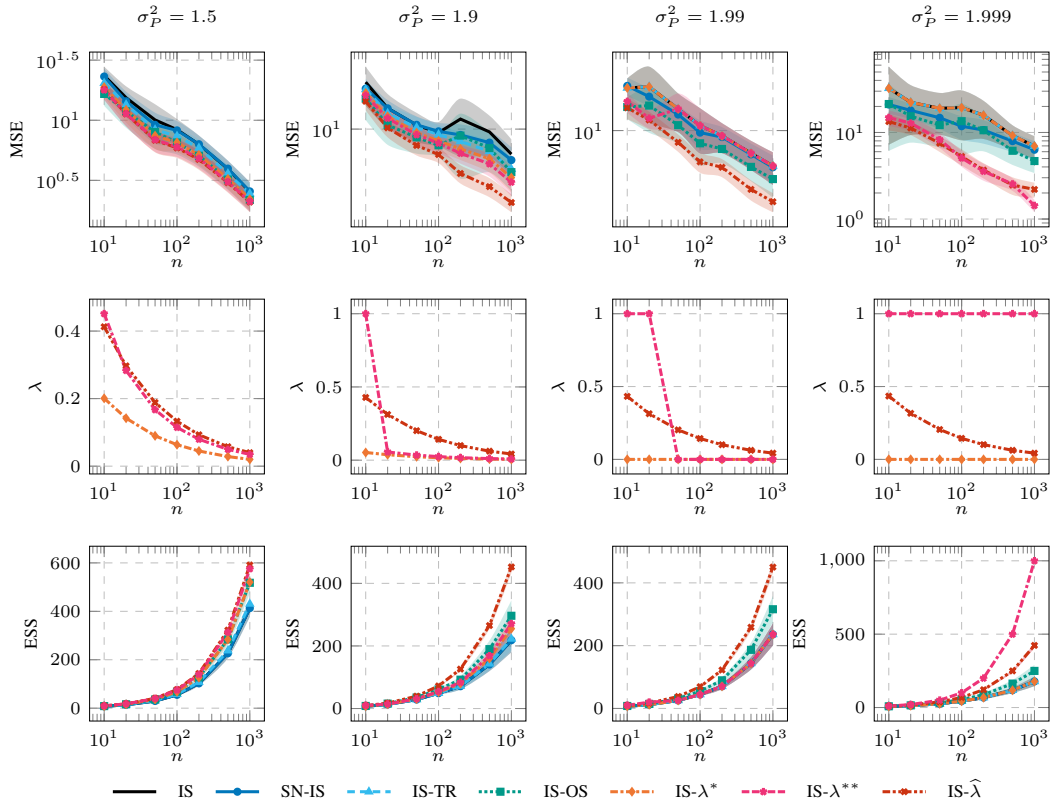


Figure 5: Mean Squared Error (MSE), correction parameter λ , and Effective Sample Size (ESS), computed as $\frac{(\sum_{i \in [n]} \omega(y_i))^2}{\sum_{i \in [n]} \omega(y_i)^2}$, as a function of the number of samples n for the different settings of Table 7 (mean \pm 95% c.i., 60 runs).

B.1.2 Contextual Bandits

Experimental Setting The experimental evaluation is carried out over 11 UCI Machine Learning Repository datasets [13, <https://archive.ics.uci.edu/ml/index.php>] as reported in Table 9. For the estimators requiring the value of the significance, we select $\delta = 0.1$.

Dataset	ecoli	glass	isolet	kropt	letter	optdigits	page-blocks	pendigits	satimage	vehicle	yeast
Dataset size (n^*)	336	214	7797	28056	20000	5620	5473	10992	6435	846	1484
Context dimension	7	9	617	6	16	64	10	16	36	18	8
Classes (K)	8	6	26	18	26	10	5	10	6	4	10

Table 9: The 11 UCI dataset considered in the experiments. For each dataset, we report the number of examples n^* , dimensionality of the context, and number of classes K .

Complete Results In the comprehensive experiment, we consider 110 combinations obtained with a single run over the 11 datasets and 10 values of the pair (α_b, α_e) with $\alpha_b \in \{0.8, 0.9\}$ and $\alpha_e \in \{0.8, 0.85, 0.9, 0.95, 0.99\}$. The experiment with reward noise $\nu = 0.25$ is reported in the main paper (Figure 2), whereas the noiseless $\nu = 0$ (deterministic rewards) is provided in Figure 6. The results are in line with the stochastic case.

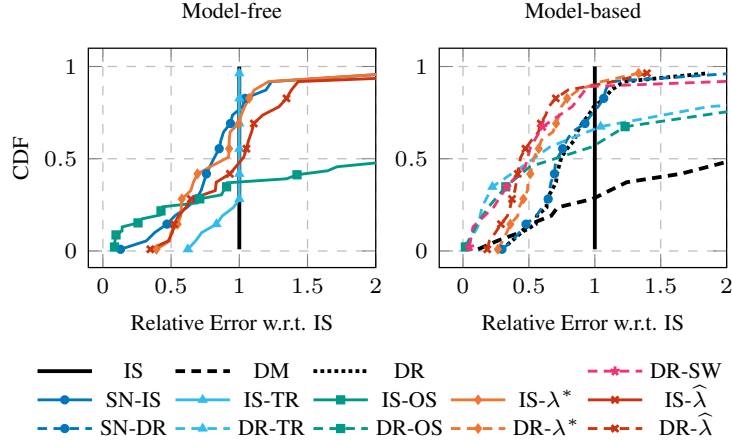


Figure 6: CDF of the absolute error normalized by IS error for deterministic rewards, across 110 conditions for model-free estimators (left) and model-based ones (right).

For the case of the *letter* dataset, we report the experiments with additional choices of α_e (10, and 11).

Estimator / n	100	200	500	1000	2000	5000	10000	20000
IS	20.04 ± 1.24	21.77 ± 2.46	14.03 ± 0.57	8.40 ± 0.20	6.13 ± 0.09	2.77 ± 0.03	1.83 ± 0.01	1.10 ± 0.01
SN-IS	27.34 ± 1.67	23.16 ± 1.40	16.86 ± 0.46	11.94 ± 0.25	7.37 ± 0.13	2.59 ± 0.03	1.74 ± 0.01	1.17 ± 0.01
IS-TR	20.04 ± 1.24	18.17 ± 1.60	13.96 ± 0.57	8.40 ± 0.20	6.13 ± 0.09	2.77 ± 0.03	1.83 ± 0.01	1.10 ± 0.01
IS-OS	24.47 ± 1.50	32.30 ± 1.17	15.37 ± 0.56	17.35 ± 0.46	16.46 ± 0.37	30.70 ± 0.15	34.03 ± 0.02	33.67 ± 0.01
IS- λ^*	20.48 ± 1.33	16.77 ± 1.14	10.06 ± 0.34	6.61 ± 0.16	5.30 ± 0.07	2.88 ± 0.03	2.08 ± 0.01	1.16 ± 0.01
IS- λ^{**}	20.80 ± 1.38	16.11 ± 0.93	9.56 ± 0.30	6.62 ± 0.15	5.19 ± 0.07	3.06 ± 0.03	2.31 ± 0.01	1.29 ± 0.01
IS- $\hat{\lambda}$	22.60 ± 1.52	17.06 ± 0.75	10.22 ± 0.28	7.77 ± 0.16	5.61 ± 0.08	3.32 ± 0.03	2.50 ± 0.02	1.37 ± 0.01
DM	28.86 ± 1.92	27.56 ± 0.95	41.04 ± 0.26	41.94 ± 0.11	42.87 ± 0.05	47.06 ± 0.01	47.58 ± 0.01	47.51 ± 0.00
DR	26.54 ± 4.51	25.56 ± 2.43	16.69 ± 0.72	9.12 ± 0.20	5.62 ± 0.09	2.14 ± 0.02	1.25 ± 0.01	0.83 ± 0.00
SN-DR	25.62 ± 3.21	24.87 ± 1.79	18.94 ± 0.62	12.36 ± 0.23	7.19 ± 0.12	2.46 ± 0.02	1.57 ± 0.01	1.07 ± 0.01
DR-TR	18.97 ± 1.12	16.54 ± 0.70	20.95 ± 0.23	17.93 ± 0.09	17.90 ± 0.06	22.73 ± 0.01	23.45 ± 0.01	23.18 ± 0.00
DR-OS	18.87 ± 1.18	19.21 ± 0.55	17.15 ± 0.38	12.01 ± 0.23	8.67 ± 0.11	17.04 ± 0.06	17.88 ± 0.02	18.49 ± 0.02
DR-SW	23.97 ± 1.28	16.66 ± 1.13	4.58 ± 0.18	4.64 ± 0.09	4.76 ± 0.05	0.75 ± 0.01	1.31 ± 0.01	0.77 ± 0.00
DR- λ^*	21.84 ± 2.30	18.16 ± 1.35	11.26 ± 0.47	6.53 ± 0.14	4.59 ± 0.07	1.78 ± 0.02	1.23 ± 0.01	0.72 ± 0.00
DR- λ^{**}	21.00 ± 2.01	16.70 ± 1.18	10.00 ± 0.41	5.82 ± 0.13	4.27 ± 0.06	1.67 ± 0.02	1.24 ± 0.01	0.69 ± 0.00
DR- $\hat{\lambda}$	19.45 ± 1.62	14.35 ± 0.95	7.89 ± 0.34	4.88 ± 0.11	3.88 ± 0.06	1.60 ± 0.02	1.26 ± 0.01	0.68 ± 0.00

Table 10: Absolute error (multiplied by 100) in the *letter* dataset varying the number of samples n for the different estimators, when $\alpha_b = 0.5$ and $\alpha_e = 0.99$ (mean \pm std, 10 runs). For each column, the estimator with smallest absolute error and the ones not statistically significantly different from that one (Welch’s t-test with $p < 0.05$) are in bold.

Estimator / n	100	200	500	1000	2000	5000	10000	20000
IS	10.08 ± 0.91	20.07 ± 5.66	20.23 ± 1.60	13.52 ± 0.42	12.23 ± 0.24	6.49 ± 0.05	3.62 ± 0.03	2.74 ± 0.01
SN-IS	15.85 ± 1.60	18.18 ± 1.71	26.34 ± 0.75	23.84 ± 0.35	12.91 ± 0.20	5.96 ± 0.05	4.15 ± 0.03	2.14 ± 0.01
IS-TR	10.08 ± 0.91	12.02 ± 1.66	13.94 ± 0.65	11.34 ± 0.22	11.38 ± 0.20	6.40 ± 0.05	3.62 ± 0.03	2.74 ± 0.01
IS-OS	26.61 ± 0.75	53.26 ± 5.63	38.57 ± 1.10	32.35 ± 0.09	30.73 ± 0.06	25.41 ± 0.02	24.48 ± 0.01	23.76 ± 0.00
IS- λ^*	10.17 ± 0.91	12.03 ± 1.62	13.11 ± 0.52	10.33 ± 0.13	8.89 ± 0.09	3.88 ± 0.03	2.74 ± 0.02	2.35 ± 0.01
IS- λ^{**}	10.21 ± 0.91	11.21 ± 1.21	12.15 ± 0.37	9.80 ± 0.09	8.23 ± 0.07	3.54 ± 0.03	2.61 ± 0.02	2.32 ± 0.01
IS- $\hat{\lambda}$	11.00 ± 0.93	9.73 ± 0.37	9.72 ± 0.18	8.95 ± 0.09	7.82 ± 0.07	3.67 ± 0.03	2.98 ± 0.02	2.56 ± 0.01
DM	22.00 ± 1.92	9.78 ± 0.47	7.27 ± 0.19	3.49 ± 0.08	2.83 ± 0.05	9.16 ± 0.02	9.97 ± 0.01	9.89 ± 0.00
DR	35.50 ± 15.08	33.18 ± 7.27	30.82 ± 1.83	24.87 ± 0.82	14.16 ± 0.27	6.19 ± 0.06	3.46 ± 0.03	1.48 ± 0.01
SN-DR	19.47 ± 3.41	20.94 ± 2.62	24.69 ± 1.07	21.41 ± 0.43	13.66 ± 0.16	6.41 ± 0.06	3.57 ± 0.03	1.57 ± 0.01
DR-TR	12.26 ± 1.14	10.90 ± 0.61	6.65 ± 0.23	9.98 ± 0.09	10.29 ± 0.04	2.34 ± 0.01	0.95 ± 0.00	0.60 ± 0.00
DR-OS	12.57 ± 1.23	8.59 ± 0.49	6.73 ± 0.27	5.25 ± 0.13	9.05 ± 0.09	2.60 ± 0.01	1.97 ± 0.01	1.21 ± 0.00
DR-SW	12.46 ± 1.09	11.73 ± 0.64	7.52 ± 0.25	11.31 ± 0.10	11.59 ± 0.04	3.49 ± 0.01	2.09 ± 0.00	1.69 ± 0.00
DR- λ^*	18.78 ± 3.41	16.07 ± 2.09	15.26 ± 0.66	13.55 ± 0.31	8.96 ± 0.15	3.97 ± 0.04	2.44 ± 0.02	1.24 ± 0.01
DR- λ^{**}	17.51 ± 2.59	14.39 ± 1.59	13.17 ± 0.52	11.89 ± 0.24	8.22 ± 0.14	3.55 ± 0.03	2.22 ± 0.02	1.21 ± 0.01
DR- $\hat{\lambda}$	14.58 ± 1.18	10.32 ± 0.58	8.42 ± 0.23	11.33 ± 0.19	11.02 ± 0.23	2.83 ± 0.02	2.00 ± 0.01	1.26 ± 0.00

Table 11: Absolute error (multiplied by 100) in the *letter* dataset varying the number of samples n for the different estimators, when $\alpha_b = 0.9$ and $\alpha_e = 0.99$ (mean \pm std, 10 runs). For each column, the estimator with smallest absolute error and the ones not statistically significantly different from that one (Welch’s t-test with $p < 0.05$) are in bold.

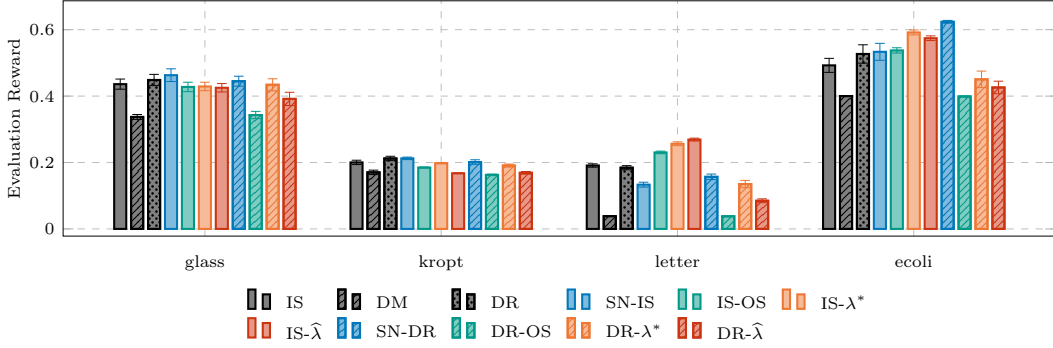


Figure 7: Evaluation reward for four different datasets after 1000 iterations (4000 iterations for *letter*) of gradient ascent with a Boltzmann policy and the non-regularized objective ($\zeta=0$) (mean \pm std, 10 runs).

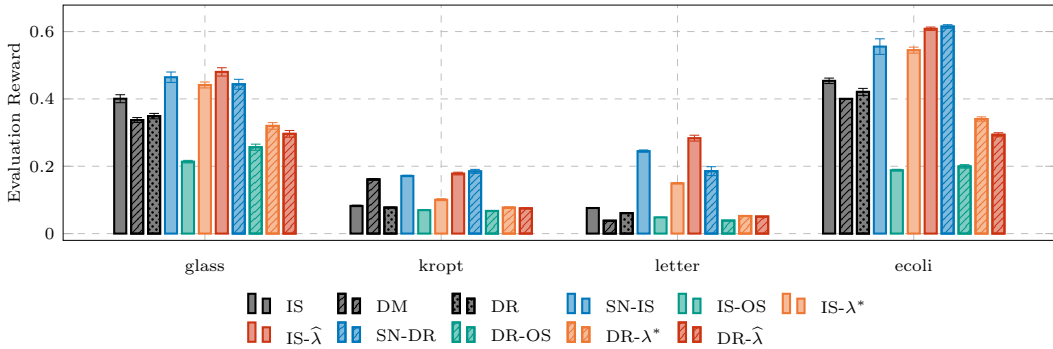


Figure 8: Evaluation reward for four different datasets after 1000 iterations (4000 iterations for *letter*) of gradient ascent with a Boltzmann policy and the regularized objective ($\zeta=0.1$) (mean \pm std, 10 runs).

B.2 Off-Policy Learning

Experimental Setting The optimization is performed by gradient ascent on the objective function:

$$\mathcal{L}(\theta) = \hat{v}(\pi_\theta) - \frac{\zeta}{n} \sum_{i \in [n]} I_2(\pi_\theta(\cdot|x_i) \parallel \pi_b(\cdot|x_i)),$$

where $\hat{v}(\pi_\theta)$ is the estimated value function using the different estimators, that is a function of the target policy π_θ . The second term is the empirical average of the divergence between the target π_θ and the behavioral policy π_b . The regularizer is controlled by the regularization parameter $\zeta \geq 0$. The gradient optimization is performed in mini-batch made of 32 samples and the learning rate is selected with RMSprop, with 0.05 as base learning rate.

Complete Results In Figure 7 and in Figure 8 we report the complete results, in the setting presented in the main paper, for the non-regularized ($\zeta=0$) and the regularized objective ($\zeta=0.1$) respectively. The experiments with the regularized objective are limited to *glass* and *ecoli* datasets. We report the corresponding learning curves for the non-regularized (Figure 9) and the regularized objectives (Figure 10).

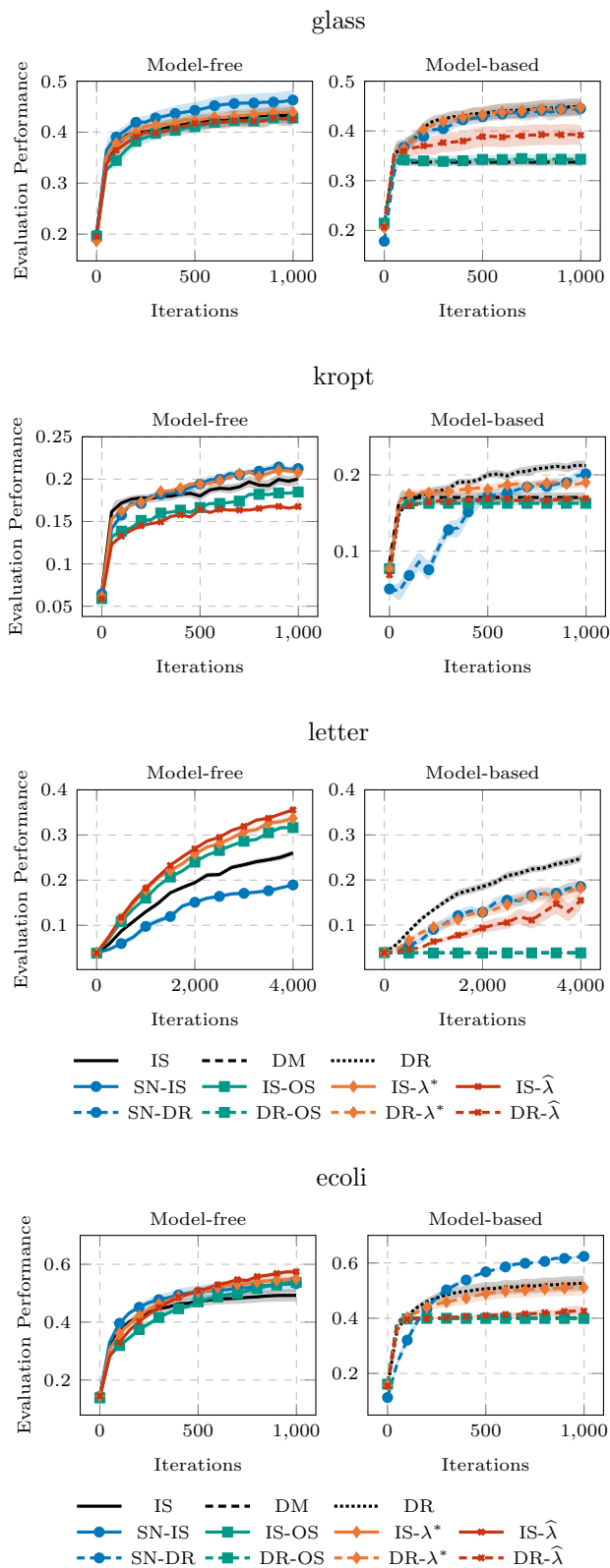


Figure 9: Evaluation reward for the four datasets comparing the learning curve of different estimators with the non-regularized objective ($\zeta = 0$) (mean \pm std, 10 runs).

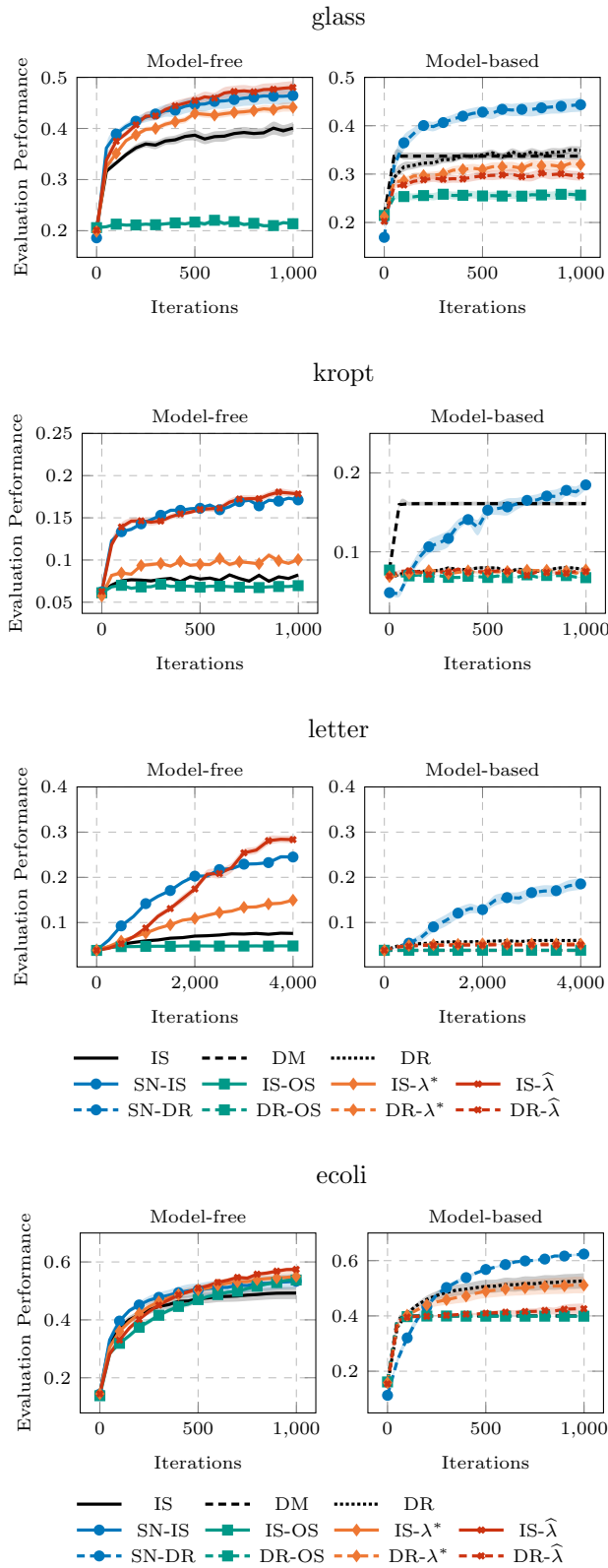


Figure 10: Evaluation reward for the four datasets comparing the learning curve of different estimators with the regularized objective ($\zeta = 0.1$) (mean \pm std, 10 runs).

C Bound Comparison and Optimization

In this appendix, we provide a comparison between the bounds of Lemma A.3 and Theorem 5.1 and show how to numerically optimize the former. For the sake of simplicity, we restrict our attention to $\alpha=2$ and we denote with $B^{**}(\lambda)$ the bound of Lemma A.3, with λ^{**} its global minimum, with $B^*(\lambda)$ the bound of Theorem 5.1, and with λ^* its global minimum.

$B^{**}(\lambda)$ displays a pretty intricate dependence on λ that is not easy to optimize. As we can notice from Figure 11, the bound based on the values of its terms admits either one or two local minima. In any case $\lambda=1$ is a value of interest, leading to a bound of the form:

$$\hat{\mu}_{n,1} - \mu \leq \|f\|_\infty \sqrt{\frac{2 \log \frac{1}{\delta}}{n}} + \frac{2 \|f\|_\infty \log \frac{1}{\delta}}{3n} + \|f\|_\infty \sqrt{J_2(P\|Q)}.$$

In such a case, we are replacing the importance weight with the value of 1 and we are estimating the mean under the target distribution with the mean of the behavioral distribution, paying the whole bias $\sqrt{J_2(P\|Q)} = \sqrt{I_2(P\|Q)} - 1$. Clearly, this circumstance is convenient only when n is sufficiently small.

The bound of Theorem 5.1 B^* is looser compared with that of Lemma A.3 B^{**} . We can see in Figure 12 that bound of B^* is convex and yeilds an optimal value of λ^* that is smaller compared to the optimal value λ^{**} of B^{**} .

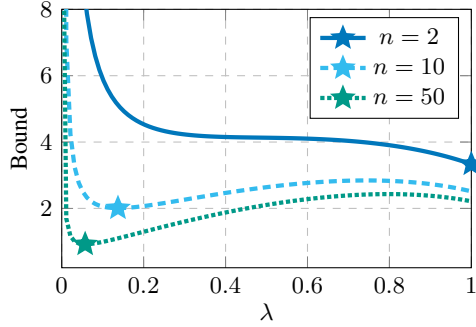


Figure 11: The bound of Lemma A.3 for $\alpha=2$, $I_2(P\|Q)=5$, $\delta=e^{-1}$, and $n \in \{2, 10, 50\}$. The minima are highlighted with the star.

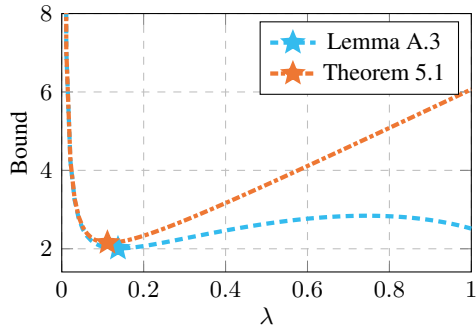


Figure 12: Comparison between the bounds of Lemma A.3 and Theorem 5.1 for $\alpha=2$, $I_2(P\|Q)=5$, $\delta=e^{-1}$, and $n=10$. The minima are highlighted with the star.

C.1 Numerical Optimization of the Bound of Lemma A.3

We now discuss how to find the global minimum of the bound presented in Lemma A.3 $B^{**}(\lambda)$. First of all, we observe that $B^{**}(\lambda)$ is continuously differentiable in λ :

$$\frac{\partial B^{**}(\lambda)}{\partial \lambda} = \frac{2 \log \frac{1}{\delta}}{3n\lambda^2} - \sqrt{(I_2(P\|Q) - 1)((1 - \lambda)I_2(P\|Q) + \lambda)}$$

Algorithm 1 Root finding for bound B^{**} of Lemma A.3

Compute the bound derivative $\frac{\partial B^{**}(\lambda)}{\partial \lambda}$
 Apply Newton's method with λ^* as initial guess obtaining λ_1 as numerical root (if exists)
if Newton's method failed to converge **or** $B^{**}(\lambda_1) < B(1)$ **then**
 return 1
else
 return λ_1
end if

$$-\frac{(I_2(P\|Q) - 1) \left(\sqrt{2 \log \frac{1}{\delta}} + \lambda \sqrt{(I_2(P\|Q) - 1)n} \right)}{2\sqrt{n((1-\lambda)I_2(P\|Q) + \lambda)}}.$$

We start proving that $\frac{\partial B^{**}(\lambda)}{\partial \lambda}$ is a strictly concave function of λ :

$$\begin{aligned} \frac{\partial^2}{\partial \lambda^2} \left(\frac{\partial B^{**}(\lambda)}{\partial \lambda} \right) &= \frac{\partial^3 B^{**}(\lambda)}{\partial \lambda^3} = -\frac{4 \log \frac{1}{\delta}}{n \lambda^4} - \frac{3(I_2(P\|Q) - 1)^{7/2} \lambda}{8((1-\lambda)I_2(P\|Q) + \lambda)^{5/2}} \\ &\quad - \frac{3(I_2(P\|Q) - 1)^{5/2}}{4((1-\lambda)I_2(P\|Q) + \lambda)^{3/2}} - \frac{3(I_2(P\|Q) - 1)^3 \sqrt{\log \frac{1}{\delta}}}{4\sqrt{2n}((1-\lambda)I_2(P\|Q) + \lambda)^{5/2}} < 0. \end{aligned}$$

We now prove that $\frac{\partial B^{**}(\lambda)}{\partial \lambda}$ admits at most two roots. By contradiction, suppose $\frac{\partial B^{**}(\lambda)}{\partial \lambda}$ admits three roots $\lambda_1 < \lambda_2 < \lambda_3$. By Rolle's theorem, there must exist $\lambda_{12} < \lambda_2$ and $\lambda_{23} < \lambda_3$ such that $\frac{\partial^2 B^{**}(\lambda)}{\partial \lambda^2}(\lambda_{12}) = \frac{\partial^2 B^{**}(\lambda)}{\partial \lambda^2}(\lambda_{23}) = 0$. Again, by Rolle's theorem, there must exist $\lambda_{1223} < \lambda_{23}$ such that $\frac{\partial^3 B^{**}(\lambda)}{\partial \lambda^3}(\lambda_{1223}) = 0$, which is a contradiction being $\frac{\partial B^{**}(\lambda)}{\partial \lambda}$ concave. Thus we consider three cases:

- $\frac{\partial B^{**}(\lambda)}{\partial \lambda}$ admits no roots. It follows that the global minimum of B^{**} is on the border $\{0, 1\}$. Since $\lim_{\lambda \rightarrow 0^+} B^{**}(\lambda) = \infty$, the minimum is in $\lambda^{**} = 1$.
- $\frac{\partial B^{**}(\lambda)}{\partial \lambda}$ admits one root. It is simple to prove that for sufficiently large λ (possibly larger than 1, but this does not matter of the sake for the function study) we have $\frac{\partial B^{**}(\lambda)}{\partial \lambda} < 0$. Being also $\lim_{\lambda \rightarrow 0^+} \frac{\partial B^{**}(\lambda)}{\partial \lambda} = -\infty$, we conclude that the root must be a saddle point and, consequently, $\lambda^{**} = 1$.
- $\frac{\partial B^{**}(\lambda)}{\partial \lambda}$ admits two roots $\lambda_1 < \lambda_2$. Thus, there must exist $\lambda_{12} < \lambda_2$ such that $\frac{\partial^2 B^{**}(\lambda)}{\partial \lambda^2}(\lambda_{12}) = 0$. Since $\frac{\partial^2 B^{**}(\lambda)}{\partial \lambda^2}$ is non-increasing, being $\frac{\partial B^{**}(\lambda)}{\partial \lambda}$ concave, it must be that $\frac{\partial^2 B^{**}(\lambda)}{\partial \lambda^2}(\lambda_1) > 0$ and $\frac{\partial^2 B^{**}(\lambda)}{\partial \lambda^2}(\lambda_2) < 0$. Thus, λ_1 is a local minimum and λ_2 a local maximum. It follows that $\lambda^{**} \in \arg \min_{\lambda \in \{\lambda_1, 1\}} B^{**}(\lambda)$.

Thus, based on the function study, it suffices to find numerically the smallest root λ_1 (whenever it exists) of $\frac{\partial B^{**}(\lambda)}{\partial \lambda}$ and compare its bound value $B^{**}(\lambda_1)$ with $B^{**}(1)$. This task can be carried out using numerical root finding, e.g., Newton's method, using as initial guess 0 or λ^* , having observed that in the optimal correction parameter λ^* of the simplified bound B^* the derivative $\frac{\partial B^{**}(\lambda)}{\partial \lambda}$ is negative. The procedure is summarized in Algorithm 1

D Comparison of Estimators for CMABs

Estimator	Formula
Direct method (DM)	$\frac{1}{n} \sum_{i \in [n]} \sum_{a \in \mathcal{A}} \pi_e(a x_i) \hat{r}(x_i, a)$
Inverse propensity scoring (IPS)	$\frac{1}{n} \sum_{i \in [n]} \frac{\pi_e(a_i x_i)}{\pi_b(a_i x_i)} r_i$
Doubly robust (DR)	$\frac{1}{n} \sum_{i \in [n]} \sum_{a \in \mathcal{A}} \pi_e(a x_i) \hat{r}(x_i, a) + \frac{1}{n} \sum_{i \in [n]} \frac{\pi_e(a_i x_i)}{\pi_b(a_i x_i)} (r_i - \hat{r}(x_i, a_i))$

Table 12: Overview of the classical off-policy estimators for CMABs. π_b and π_e denote the behavioral and target policies respectively and \hat{r} the estimated reward function.

E Analysis of the IS-OS estimator

The IS-OS (optimistic shrinkage) [56] is based on the weight transformation:

$$\omega_\tau^{\text{OS}}(y) = \frac{\tau \omega(y)}{\omega(y)^2 + \tau}.$$

First of all, we notice that when $P=Q$ a.s. the weight becomes $\omega_\tau^{\text{OS}}(y) = \frac{\tau}{\tau+1}$, so the estimator is biased. We start by observing that the corrected weight $\omega_\tau^{\text{OS}}(y)$ converges to zero when the non-corrected weight is either zero or infinity. Thus, the maximum value of the weight must be in between. We compute it by vanishing the derivative:

$$\frac{\partial}{\partial \omega} \frac{\tau \omega}{\omega^2 + \tau} = 0 \implies \omega = \sqrt{\tau}.$$

From which, we obtain the maximum value of the weight equal to $\frac{\sqrt{\tau}}{2}$. We now focus on the following result concerning the bias and the variance of the IS-OS estimator.

Lemma E.1. *Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions with $P \ll Q$. For every $\tau \geq 0$, the bias and variance of the IS-OS estimator can be bounded as:*

$$\left| \mathbb{E}_{y \sim Q} [\hat{\mu}_{n,\tau}^{\text{OS}}] - \mu \right| \leq \frac{\|f\|_\infty}{\tau} I_3(P\|Q), \quad \mathbb{V}\text{ar}_{y_i \sim Q} [\hat{\mu}_{n,\tau}^{\text{OS}}] \leq \frac{\|f\|_\infty^2}{n} I_2(P\|Q).$$

Proof. Let us start with the bias. Based also on Lemma A.1, we consider the following inequality:

$$\mathbb{E}_{y \sim Q} \left[\left| \omega_\tau^{\text{OS}}(y) - \omega(y) \right| \right] = \mathbb{E}_{y \sim Q} \left[\frac{\omega(y)^3}{\omega(y)^2 + \tau} \right] \leq \mathbb{E}_{y \sim Q} \left[\frac{\omega(y)^3}{\tau} \right] = \frac{I_3(P\|Q)}{\tau}.$$

We consider now the variance term and derive a bound on the second moment of the OS weight:

$$\mathbb{E}_{y \sim Q} \left[\omega_\tau^{\text{OS}}(y)^2 \right] = \mathbb{E}_{y \sim Q} \left[\left(\frac{\omega(y)\tau}{\omega(y)^2 + \tau} \right)^2 \right] \leq \mathbb{E}_{y \sim Q} \left[\omega(y)^2 \right] = I_2(P\|Q).$$

□

We now move to the concentration result.

Theorem E.1. *Let $P, Q \in \mathcal{P}(\mathcal{Y})$ be two probability distributions with $P \ll Q$. Then, having selected $\tau^* = \left(\frac{6n I_3(P\|Q)}{\log \frac{1}{\delta}} \right)^{\frac{2}{3}}$ for the IS-OS estimator, for every $\delta \in [0, 1]$, with probability at least $1 - \delta$ it holds that:*

$$\hat{\mu}_{n,\tau^*}^{\text{OS}} - \mu \leq \|f\|_\infty \sqrt{\frac{2I_2(P\|Q) \log \frac{1}{\delta}}{n}} + \|f\|_\infty^3 \sqrt{\frac{3I_3(P\|Q) (\log \frac{1}{\delta})^2}{4n^2}}.$$

Proof. We apply Bernstein's inequality to the estimator, starting for the bias and variance bounds of Lemma E.1:

$$\begin{aligned}\hat{\mu}_{n,\tau}^{\text{OS}} - \mu &= \hat{\mu}_{n,\tau}^{\text{OS}} - \mathbb{E}_{y_i \sim Q}[\hat{\mu}_{n,\tau}^{\text{OS}}] + \mathbb{E}_{y_i \sim Q}[\hat{\mu}_{n,\tau}^{\text{OS}}] - \mu \\ &\leq \|f\|_\infty \sqrt{\frac{2I_2(P\|Q) \log \frac{1}{\delta}}{n}} + \frac{\|f\|_\infty \sqrt{\tau} \log \frac{1}{\delta}}{3n} + \frac{\|f\|_\infty}{\tau} I_3(P\|Q).\end{aligned}$$

We now minimize the bound as a function of τ by vanishing the derivative to obtain:

$$\tau^* = \left(\frac{6nI_3(P\|Q)}{\log \frac{1}{\delta}} \right)^{\frac{2}{3}}.$$

By substituting τ^* we obtain the result. □

F Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] See Section 9.
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A] The paper is mostly a theoretical contribution and we think it will not have any potential negative societal impact.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Sections 3, 4, and 5.
 - (b) Did you include complete proofs of all theoretical results? [Yes] They are reported in Appendix A.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] The code is included in the supplementary material, with a README reporting the commands to reproduce the results. The datasets employed in the CMAB experiments are taken from the UCI repository (<https://archive.ics.uci.edu/ml/index.php>).
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Section 8 and Appendix B.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] Their meaning is reported in the corresponding figure captions.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix B
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes] See Appendix B
 - (b) Did you mention the license of the assets? [Yes] See Appendix B
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] The code is attached to the supplementary material.
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]