# Pessimistic Policy Learning for Continuous Action Spaces: A Self-Normalized and Computationally Tractable Approach

Zexi Fan[1]

[1]*School of Mathematics, Peking University*

November 24, 2025

## Abstract

This paper studies offline policy learning in contextual bandits, aiming to learn an optimal policy from a fixed dataset collected under a known behavior policy [12]. Existing pessimistic policy learning (PPL) methods have shown great promise in handling data with poor overlap—a common failure case for standard estimators—but their theoretical and algorithmic tools are restricted to discrete action spaces. This paper provides the first rigorous extension of PPL to continuous action spaces. This extension introduces three fundamental challenges: (1) the statistical complexity of an infinite policy class can no longer be measured by finite combinatorial metrics like the Natarajan dimension; (2) the variance of the Importance Sampling (IS) estimator becomes unbounded as the behavior policy density $\mu(a|x)$ can approach zero; and (3) the discrete tree-search optimization algorithm is no longer applicable.

We address these challenges by introducing a new set of theoretical and algorithmic tools. First, we develop a novel continuous-action IS estimator $\hat{V}_n(\pi)$ and its corresponding self-normalized pessimistic regularizer $\mathcal{V}_n(\pi)$, which generalizes the empirical Bernstein variance term. Second, using tools from empirical process theory—specifically Dudley's integral inequality and Massart's concentration inequality—we derive a new uniform concentration bound that holds for the resulting unbounded empirical process. We establish the $O(n^{-1/2})$ convergence rate of our estimator under a continuous overlap assumption and provide a matching minimax lower bound. Finally, we demonstrate that naive policy gradient optimization fails numerically and derive a computationally tractable **Pessimism Policy Learning with Majorization-Minimization (PPL-MM)** algorithm. This algorithm provably optimizes our non-convex and non-smooth pessimistic objective by converting it into a sequence of stable, re-weighted policy gradient steps. Through a rigorous paired statistical evaluation ($N = 180$ independent experiments), we demonstrate that PPL-MM achieves statistically significant improvements (FDR $< 0.05$) over standard baselines, with massive effect sizes (Cohen's $d > 3.0$) in the most challenging high-variance scenarios.

**Keywords** Offline Policy Learning, Contextual Bandits, Continuous Action Spaces, Pessimism

## 1 Introduction

Policy learning, which aims to find an optimal individualized decision rule from data, is a cornerstone of modern data-driven decision-making [3, 28, 4]. Its applications are broad, ranging from personalized medicine [29] and advertising [1] to recommendation systems [30]. A central challenge in this field is learning from *offline* data, where a decision-maker must learn the best policy using only a fixed dataset collected *a priori*, often by a suboptimal behavior policy [20, 8].

Learning from offline data requires counterfactual evaluation, which is notoriously difficult. Standard methods, such as those based on Inverse Propensity Weighting (IPW), construct an estimate of a policy's value (or "welfare") and select the policy that maximizes this estimate [26, 9]. These "greedy" approaches are highly sensitive to the quality of the offline data, particularly to the **overlap assumption**—the requirement that the behavior policy assigns a non-trivial probability to all actions a target policy might take. In many real-world scenarios, this assumption is violated [13, 19, 12]. When overlap is poor for a suboptimal policy, its value estimate can have extremely high variance, potentially appearing much larger than its true value. A greedy algorithm may then mistakenly select this highly suboptimal policy, leading to catastrophic failure.

To address this fundamental flaw, recent work has introduced the principle of **Pessimistic Policy Learning (PPL)** [13, 19, 12]. Instead of greedily maximizing the point estimate $\hat{V}(\pi)$, PPL maximizes a *Lower Confidence Bound (LCB)* of the value: $\hat{V}(\pi) - R(\pi)$, where $R(\pi)$ is a policy-dependent regularizer

that quantifies the estimation uncertainty. The central benefit of this approach is that the algorithm's final performance guarantee depends only on the estimation error of the *optimal* policy, $R(\pi^*)$, rather than the worst-case error over all policies. This allows PPL to learn effectively even when only the optimal actions are well-covered in the data, while suboptimal actions may have arbitrarily poor overlap [12].

However, the existing theory and algorithms for PPL are fundamentally restricted to **discrete action spaces** [12]. This limitation is severe, as many real-world problems, from robotic control and dynamic pricing to medical dosing [14, 17], involve continuous actions. Extending PPL to the continuous-action setting is not a trivial step; it breaks all three pillars of the original framework in [13, 12]. First, the statistical complexity of the (finite) policy class in the original work is measured using combinatorial tools like the Natarajan dimension [11]. For an infinite, continuous policy class (e.g., a class of functions $\pi(a|x)$), these tools are no longer applicable. Second, the Importance Sampling (IS) estimator $\hat{V}_n(\pi)$ relies on the weight $w(x, a) = \pi(a|x)/\mu(a|x)$. In the continuous setting, the behavior policy $\mu(a|x)$ is a probability density function, which can be arbitrarily close to zero. This makes the variance of $\hat{V}_n(\pi)$ unbounded and the estimation problem even more severe than in the discrete case, invalidating prior variance bounds. Third, the optimization of the pessimistic objective in the discrete setting is solved via a policy tree search, a method that is computationally intractable in a continuous action space.

This paper provides the first rigorous, end-to-end extension of Pessimistic Policy Learning to continuous action spaces. Our contributions are threefold:

**1. Theoretical Framework for Unbounded Weights.** We derive a novel self-normalized pessimistic regularizer, $\mathcal{V}_n(\pi)$, specifically designed to control the unbounded variance characteristic of continuous IS estimators. By replacing combinatorial complexity measures with modern empirical process theory—leveraging Dudley's entropy integral and Massart's finite-class concentration inequalities—we prove a new uniform concentration bound (Theorem 3.3) over infinite, parameterized policy classes. We further establish an $O(n^{-1/2})$ convergence rate under a continuous overlap condition (Corollary 3.5) and prove its minimax optimality (Theorem 3.6).

**2. Stable Optimization via Majorization-Minimization.** Recognizing the numerical instability of direct policy gradient optimization for our non-convex pessimistic objective, we develop the **PPL-MM** algorithm (Algorithm 1). This method iteratively linearizes the regularizer, transforming the intractable original problem into a sequence of stable, re-weighted surrogate problems that can be reliably solved with standard tools.

**3. Rigorous Statistical Validation.** We design a new benchmark suite targeting distinct mechanisms of overlap failure. Through a rigorous paired statistical evaluation ($N = 180$ independent experiments), we demonstrate that PPL-MM achieves statistically significant improvements (FDR $< 0.05$) over standard baselines, with massive effect sizes (Cohen's $d > 3.0$) in the most challenging high-variance scenarios.

The remainder of this paper is organized as follows. Section 2 defines the problem setting. Section 3 presents our main theoretical results on uniform concentration and minimax optimality. Section 4 derives the practical PPL-MM algorithm. Section 5 details our experimental benchmarks and statistical findings. Section 6 concludes with a discussion of limitations and future directions. All detailed proofs are deferred to the Appendix.

# 2   Preliminaries and Problem Setup

We consider the problem of offline policy learning in a continuous-action contextual bandit setting. Our setup is based on a fixed, known behavior policy, which corresponds to the "batched data" setting described by [12].

## 2.1   Notation and Definitions

Let $(\mathcal{X}, \Sigma_{\mathcal{X}})$ be the context space, and $(\mathcal{A}, \Sigma_{\mathcal{A}})$ be the action space, where $\mathcal{A}$ is a compact subset of $\mathbb{R}^d$ and $\Sigma_{\mathcal{A}}$ is its Borel $\sigma$-algebra. Let $\lambda_{\mathcal{A}}$ denote the Lebesgue measure on $\mathcal{A}$. The reward $R$ is a random variable supported on a bounded interval, which we normalize to $\mathcal{R} = [0, 1]$ without loss of generality.

The data is generated as follows:

1. A context $X$ is drawn from a marginal distribution $P_{\mathcal{X}}$ on $\mathcal{X}$.

2. An action $A$ is drawn from the **behavior policy** $\mu(\cdot|X)$, which is a Markov kernel from $\mathcal{X}$ to $\mathcal{A}$. We assume $\mu(a|x)$ is a probability density function (p.d.f.) with respect to $\lambda_{\mathcal{A}}$ for all $x \in \mathcal{X}$, and that this density function is known to the learner.

3. A reward $R$ is drawn from a conditional distribution $P_R(\cdot|X, A)$.

We denote the true (unknown) mean reward function as $Q(x, a) = \mathbb{E}[R|X = x, A = a]$. The joint distribution over $\mathcal{Z} = \mathcal{X} \times \mathcal{A} \times \mathcal{R}$ is $P$. We are given an offline dataset $D_n = \{Z_i\}_{i=1}^n = \{(X_i, A_i, R_i)\}_{i=1}^n$ of $n$ i.i.d. samples drawn from $P$.

## 2.2 Policy Learning and Counterfactual Estimation

A **policy** $\pi$ is a Markov kernel from $\mathcal{X}$ to $\mathcal{A}$, which we also assume admits a p.d.f. $\pi(a|x)$ with respect to $\lambda_{\mathcal{A}}$. We consider a (potentially infinite) target policy class $\Pi$. For the purposes of optimization (Section 4), we will assume $\Pi$ is parameterized by $\theta \in \Theta \subseteq \mathbb{R}^p$, i.e., $\Pi = \{\pi_\theta(a|x) \mid \theta \in \Theta\}$.

The **true value** of a policy $\pi \in \Pi$ is its expected reward over the data-generating process:

$$V(\pi) = \mathbb{E}_{X \sim P_{\mathcal{X}}} \left[ \int_{\mathcal{A}} \pi(a|X)Q(X, a)d\lambda_{\mathcal{A}}(a) \right]$$

Our goal is to find the optimal policy $\pi^* = \text{argmax}_{\pi \in \Pi} V(\pi)$ using only the offline dataset $D_n$.

Since $Q(x, a)$ is unknown, we cannot compute $V(\pi)$ directly. We rely on counterfactual estimation via **Importance Sampling (IS)**. By the law of iterated expectations, $V(\pi)$ can be re-written as an expectation over the known data-generating distribution $P$:

$$V(\pi) = \mathbb{E}_{(X,A,R) \sim P} \left[ \frac{\pi(A|X)}{\mu(A|X)} R \right]$$

This motivates the standard IS estimator for the policy value:

$$\hat{V}_n(\pi) = \frac{1}{n} \sum_{i=1}^n w_i(\pi)R_i, \quad \text{where} \quad w_i(\pi) = \frac{\pi(A_i|X_i)}{\mu(A_i|X_i)}$$

We refer to $w_i(\pi)$ as the IS weight for data point $i$ under policy $\pi$.

## 2.3 The Pessimistic Objective

A standard "greedy" algorithm seeks to maximize the empirical value: $\hat{\pi}_{\text{greedy}} = \text{argmax}_{\pi \in \Pi} \hat{V}_n(\pi)$. This approach is notoriously unstable. In the continuous setting, the behavior density $\mu(A_i|X_i)$ can be arbitrarily close to zero, causing the IS weights $w_i(\pi)$ to "explode". This leads to an estimator $\hat{V}_n(\pi)$ with catastrophic variance, which may grossly overestimate the value of a suboptimal policy.

We adopt the Pessimistic Policy Learning (PPL) framework, which optimizes a Lower Confidence Bound (LCB) of the value. The goal is to solve:

$$\hat{\pi}_{\text{PPL}} = \text{argmax}_{\pi \in \Pi} \left\{ \hat{V}_n(\pi) - \text{StdErr}_n(\pi) \right\}$$

Here, $\text{StdErr}_n(\pi)$ is a policy-dependent regularizer that serves as a high-probability upper bound on the estimation error, $|\hat{V}_n(\pi) - V(\pi)|$. The core of this paper is to derive a form of $\text{StdErr}_n(\pi)$ that is (1) theoretically valid in the continuous-action, unbounded-weight setting, and (2) leads to a computationally tractable optimization algorithm.

# 3 Theory of Continuous PPL

Our theoretical argument proceeds in three main parts. First, we restate the core algebraic principle of pessimism, which is agnostic to the action space. Then, we establish the explicit form of the self-normalized pessimism regularizer. Finally, we develop the primary contribution of this work: a new set of concentration inequalities that allow this principle to be applied to the continuous-action, unbounded-weight setting.

## 3.1 The Principle of Pessimism

The core idea of PPL is to modify the greedy learning objective $\max_\pi \hat{V}_n(\pi)$ to explicitly account for estimation uncertainty. This is achieved by optimizing a Lower Confidence Bound (LCB) on the policy value. We define our pessimistic objective as:

$$\hat{\pi}_{\text{PPL}} = \text{argmax}_{\pi \in \Pi} \left\{ \hat{V}_n(\pi) - \text{StdErr}_n(\pi) \right\}$$

where $\text{StdErr}_n(\pi)$ is a policy-dependent regularizer that we will construct to be a high-probability upper bound on the estimation error, $|\hat{V}_n(\pi) - V(\pi)|$.

The fundamental merit of this pessimistic objective is captured in the following proposition, which is a direct extension of the logic from [13, 12]. It demonstrates that the suboptimality of the learned policy $\hat{\pi}_{\text{PPL}}$ depends only on the estimation uncertainty of the *optimal* policy $\pi^*$, rather than the worst-case uncertainty over all $\pi \in \Pi$.

**Proposition 3.1** (The Pessimism Principle). *Let $\hat{\pi} = argmax_{\pi \in \Pi}\{\hat{V}_n(\pi) - StdErr_n(\pi)\}$ be the policy learned by PPL, and let $\pi^* = argmax_{\pi \in \Pi} V(\pi)$ be the optimal policy in $\Pi$. Let the suboptimality gap be $\mathcal{L}(\hat{\pi}) = V(\pi^*) - V(\hat{\pi})$.*

*Define the uniform concentration event $\mathcal{E}$ as:*

$$\mathcal{E} := \left\{ \left| \hat{V}_n(\pi) - V(\pi) \right| \le StdErr_n(\pi), \quad \forall \pi \in \Pi \right\}$$

*Then, on the event $\mathcal{E}$, the suboptimality of $\hat{\pi}$ is bounded by:*

$$\mathcal{L}(\hat{\pi}) \le 2 \cdot StdErr_n(\pi^*)$$

*Proof.* The proof is algebraic and holds for any choice of estimator $\hat{V}_n$ and regularizer $\text{StdErr}_n$, provided the event $\mathcal{E}$ holds [12].

By the definition of $\hat{\pi}$ as the maximizer of the pessimistic objective, we have:

$$\hat{V}_n(\hat{\pi}) - \text{StdErr}_n(\hat{\pi}) \ge \hat{V}_n(\pi^*) - \text{StdErr}_n(\pi^*) \tag{1}$$

On the event $\mathcal{E}$, we have two bounds by definition:

$$V(\hat{\pi}) \ge \hat{V}_n(\hat{\pi}) - \text{StdErr}_n(\hat{\pi}) \tag{2}$$
$$\hat{V}_n(\pi^*) \ge V(\pi^*) - \text{StdErr}_n(\pi^*) \tag{3}$$

We chain these inequalities together:

$$\begin{align}
V(\hat{\pi}) &\ge \hat{V}_n(\hat{\pi}) - \text{StdErr}_n(\hat{\pi}) & \text{(by (2))} \tag{4} \\
&\ge \hat{V}_n(\pi^*) - \text{StdErr}_n(\pi^*) & \text{(by (1))} \tag{5} \\
&\ge (V(\pi^*) - \text{StdErr}_n(\pi^*)) - \text{StdErr}_n(\pi^*) & \text{(by (3))} \tag{6} \\
&= V(\pi^*) - 2 \cdot \text{StdErr}_n(\pi^*) \tag{7}
\end{align}$$

Rearranging the resulting inequality, $V(\hat{\pi}) \ge V(\pi^*) - 2 \cdot \text{StdErr}_n(\pi^*)$, gives the suboptimality bound $V(\pi^*) - V(\hat{\pi}) \le 2 \cdot \text{StdErr}_n(\pi^*)$. $\square$

Proposition 3.1 illustrates the power of the pessimistic framework. It shifts the analytic burden entirely to constructing a regularizer $\text{StdErr}_n(\pi)$ that is both (1) a valid high-probability upper bound for the error (i.e., satisfying event $\mathcal{E}$) and (2) computationally tractable. The remainder of this section is dedicated to the first challenge.

## 3.2 The Self-Normalized Regularizer

The challenge set by Proposition 3.1 is to construct a regularizer $\text{StdErr}_n(\pi)$ that (1) serves as a valid high-probability upper bound on the estimation error $|\hat{V}_n(\pi) - V(\pi)|$ and (2) is small for policies with good overlap.

In the continuous setting, the IS estimator $\hat{V}_n(\pi)$ is a sum of i.i.d. random variables $Y_i(\pi) = w_i(\pi)R_i$. These variables are unbounded, as the IS weight $w_i(\pi)$ can be arbitrarily large. Standard concentration inequalities like Hoeffding's, which require bounded support, are inapplicable.

We must therefore use a **self-normalized** approach, where the deviation of the estimator is controlled by its own (empirical) variance, an idea central to the empirical Bernstein's inequality. The total estimation error $|\hat{V}_n(\pi) - V(\pi)|$ is driven by the variance of $Y_i(\pi)$.

We follow the logic of [12, Sec 3.3] to construct a regularizer that serves as an upper bound on the standard deviation of $\hat{V}_n(\pi)$. The (conditional) variance of a single term $Y_i(\pi)$ is:

$$\text{Var}(Y_i(\pi) \mid X_i, A_i) = \text{Var}(w_i(\pi)R_i \mid X_i, A_i) = w_i(\pi)^2 \cdot \text{Var}(R_i \mid X_i, A_i)$$

Crucially, since we assume the rewards are normalized $R_i \in [0, 1]$, the variance of the reward is bounded: $\mathrm{Var}(R_i \mid X_i, A_i) \leq \mathbb{E}[R_i^2] \leq 1^2 = 1$. This implies that the conditional variance of our (unbounded) weighted reward is upper-bounded by the (unbounded) squared weight itself:

$$\mathrm{Var}(Y_i(\pi) \mid X_i, A_i) \leq w_i(\pi)^2$$

This is a critical finding. It justifies constructing the regularizer $\mathcal{V}_n(\pi)$—our proxy for the standard deviation—based on the moments of the IS weights $w_i(\pi)$ alone, not the full weighted rewards $Y_i(\pi)$.

This directly motivates our definitions for the regularizer components, which are the continuous-space parallel to [12, Eq. (7)]:

1. **Sample Deviation** $(V_{s,n})$: The empirical $L_2$ norm of the *weights*, normalized. This is our primary *computable* proxy for the standard deviation bound.

$$V_{s,n}(\pi) = \frac{1}{n} \left( \sum_{i=1}^{n} w_i(\pi)^2 \right)^{1/2} = \frac{1}{n} \left( \sum_{i=1}^{n} \left( \frac{\pi(A_i|X_i)}{\mu(A_i|X_i)} \right)^2 \right)^{1/2}$$

2. **Population Deviation** $(V_{p,n})$: The population-level (conditional on $X_i$) $L_2$ norm of the *weights*.

$$V_{p,n}(\pi) = \frac{1}{n} \left( \sum_{i=1}^{n} \mathbb{E}[w_i(\pi)^2 \mid X_i] \right)^{1/2}$$

where $\mathbb{E}[w_i(\pi)^2 \mid X_i = x] = \int_{\mathcal{A}} \mu(a|x) \left( \frac{\pi(a|x)}{\mu(a|x)} \right)^2 d\lambda_{\mathcal{A}}(a) = \int_{\mathcal{A}} \frac{\pi(a|x)^2}{\mu(a|x)} d\lambda_{\mathcal{A}}(a)$.

3. **Higher-Order Deviation** $(V_{h,n})$: The population-level $L_4$ norm of the *weights*, required in the self-normalization proofs (see Appendix A.2).

$$V_{h,n}(\pi) = \frac{1}{n} \left( \sum_{i=1}^{n} \mathbb{E}[w_i(\pi)^4 \mid X_i] \right)^{1/4}$$

The terms $V_{p,n}$ and $V_{h,n}$ are theoretical constructs that are intractable to compute. However, $V_{s,n}$ is fully empirical and computable. As our theory will show, we need to control the error from all sources. This includes the (bounded) drift of the conditional expectations (Term (ii) in our proof sketch), which requires a minimal $O(n^{-1/2})$ term.

This leads to our formal definition of the theoretical regularizer:

**Definition 3.2** (Self-Normalized Regularizer). *The theoretical self-normalized regularizer $\mathcal{V}_n : \Pi \to \mathbb{R}^+$ is defined as:*

$$\mathcal{V}_n(\pi) := \max \left\{ V_{s,n}(\pi), V_{p,n}(\pi), V_{h,n}(\pi), n^{-1/2} \right\}$$

*The full pessimistic regularizer $StdErr_n(\pi)$ is the product of this term and a complexity measure $\beta(D_n) > 0$, which we derive in the following section:*

$$StdErr_n(\pi) = \beta(D_n) \cdot \mathcal{V}_n(\pi)$$

This regularizer is a "design-based" [12] upper bound on the standard deviation of $\hat{V}_n(\pi)$, based only on the known behavior policy $\mu$, the target policy $\pi$, and the assumption that $R_{\max} = 1$. This construction is the key to our subsequent concentration bounds.

## 3.3 Main Theoretical Results

With the pessimistic principle established in Proposition 3.1 and our self-normalized, weight-based regularizer defined in Definition 3.2, our primary task is to prove that the uniform concentration event $\mathcal{E}$ holds with high probability. That is, we must show that our regularizer $StdErr_n(\pi)$ uniformly controls the estimation error $|\hat{V}_n(\pi) - V(\pi)|$ for all $\pi \in \Pi$, even when the IS weights $w_i(\pi)$ are unbounded.

Our main theorem achieves this by leveraging the full power of our auxiliary lemmas (Appendix A.1). We define the complexity penalty $\beta(D_n)$ as the sum of two distinct components, $\beta_1$ and $\beta_2$, which correspond to the two parts of our error decomposition.

**Theorem 3.3** (Uniform Concentration and Suboptimality Bound). *Let $\Pi$ be a policy class and $D_n = \{(X_i, A_i, R_i)\}_{i=1}^n$ be the observed dataset. Let $D_n' = \{(X_i, A_i', R_i')\}_{i=1}^n$ denote a ghost dataset where $(A_i', R_i')$ are drawn conditionally independent of $(A_i, R_i)$ given $X_i$.*

*We explicitly define the following functional classes and measures:*

1. *Let $P_n := \frac{1}{n}\sum_{i=1}^n \delta_{X_i}$ be the empirical measure over the contexts.*

2. *Let $\mathcal{F}_g := \{g_\pi : \mathcal{X} \to [0,1] \mid g_\pi(x) = \mathbb{E}_{(A,R)\sim\pi(\cdot|x)}[R \mid X = x], \pi \in \Pi\}$ be the class of conditional value functions.*

3. *Let $\mathcal{F}_{D_n,D_n'} \subset \mathbb{R}^n$ be the class of self-normalized discrepancy vectors, where each $f_\pi \in \mathcal{F}_{D_n,D_n'}$ is a vector with entries:*

$$f_{\pi,i} := \frac{w_\pi(X_i, A_i)R_i - w_\pi(X_i, A_i')R_i'}{\sqrt{\sum_{j=1}^n (w_\pi(X_j, A_j) + w_\pi(X_j, A_j'))^2}}, \quad i = 1, \ldots, n$$

*Based on these, we define the worst-case and empirical complexity measures via Dudley's entropy integral:*

$$\mathcal{I}_{sup}(\Pi, n) := \sup_{D_n, D_n'} \int_0^2 \sqrt{\log N(\epsilon, \mathcal{F}_{D_n,D_n'}, \|\cdot\|_2)}\, d\epsilon$$

$$\mathcal{I}_n(\mathcal{F}_g) := \int_0^1 \sqrt{\log N(\epsilon, \mathcal{F}_g, L_2(P_n))}\, d\epsilon$$

*Fix $\delta \in (0,1)$. Let $C_1, C_2 < \infty$ be universal constants. Define the total complexity penalty $\beta(D_n) := 8\beta_1 + \beta_2(D_n)$, composed of:*

$$\beta_1 := C_1\left(\mathcal{I}_{sup}(\Pi, n) + \sqrt{\log(4/\delta)}\right), \quad \beta_2(D_n) := C_2\left(\mathcal{I}_n(\mathcal{F}_g) + \sqrt{\log(8/\delta)}\right)$$

*Let $StdErr_n(\pi) := \beta(D_n) \cdot \mathcal{V}_n(\pi)$ be the pessimistic regularizer (with $\mathcal{V}_n(\pi)$ from Definition 3.2), and let $\hat{\pi}_{PPL} := \mathrm{argmax}_{\pi \in \Pi}\{\hat{V}_n(\pi) - StdErr_n(\pi)\}$.*

*Assuming $\mathcal{I}_{sup}(\Pi, n)$ and $\mathbb{E}[\mathcal{I}_n(\mathcal{F}_g)]$ are finite, with probability at least $1 - \delta$:*

(a) **Uniform Concentration:** $\left|\hat{V}_n(\pi) - V(\pi)\right| \le StdErr_n(\pi), \quad \forall \pi \in \Pi.$

(b) **Suboptimality Bound:** $V(\pi^*) - V(\hat{\pi}_{PPL}) \le \min\{2 \cdot StdErr_n(\pi^*), 1\}.$

*Proof Sketch.* The full, rigorous proof is provided in Appendix A.6. The core of the proof is to establish part (a).

1. **Error Decomposition:** We decompose the total error using the conditional expectation $V_n(\pi) = \mathbb{E}[\hat{V}_n(\pi) \mid X_1, \ldots, X_n]$:

$$|\hat{V}_n(\pi) - V(\pi)| \le \underbrace{|\hat{V}_n(\pi) - V_n(\pi)|}_{\text{Term (i): Unbounded Fluctuation}} + \underbrace{|V_n(\pi) - V(\pi)|}_{\text{Term (ii): Bounded Drift}}$$

2. **Self-Normalization:** We divide by our regularizer $\mathcal{V}_n(\pi)$ and take the supremum over $\pi \in \Pi$:

$$\sup_\pi \frac{|\hat{V}_n - V|}{\mathcal{V}_n} \le \sup_\pi \frac{|\text{Term (i)}|}{\mathcal{V}_n} + \sup_\pi \frac{|\text{Term (ii)}|}{\mathcal{V}_n}$$

3. **Bounding Term (ii):** This term is a standard empirical process for the *bounded* function class $\mathcal{F}_g \subseteq [0,1]$. By the definition of our regularizer, $\mathcal{V}_n(\pi) \ge n^{-1/2}$. This allows us to bound the normalized term:

$$\sup_\pi \frac{|\text{Term (ii)}|}{\mathcal{V}_n} \le \sup_\pi |\text{Term (ii)}| \cdot \sup_\pi \frac{1}{\mathcal{V}_n} \le \left(\sup_\pi |V_n(\pi) - V(\pi)|\right) \cdot \sqrt{n}$$

We apply Lemma A.13, which uses McDiarmid's, Symmetrization, and Empirical Dudley bounds, to show that this term is bounded by $\beta_2(D_n)$ with high probability.

4. **Bounding Term (i):** This is the primary challenge, as it involves the unbounded $Y_i(\pi)$ terms. We apply our (corrected) chain of symmetrization lemmas.

   i. Lemma A.8 controls the necessary ghost sample weight statistics using conditional Chebyshev's.

ii. Lemma A.10 uses this to show that $\mathbb{P}(\sup \frac{|\text{Term (i)}|}{\mathcal{V}_n} \geq 8\beta_1)$ is bounded by the tail probability of a self-normalized Rademacher process, $S'_n(\mathcal{F})$, where the denominator is the $\ell_2$-norm of the weights, $\|\mathbf{w} + \mathbf{w}'\|_2$.

iii. Lemma **??** shows that this class $\mathcal{F} = \{\frac{\mathbf{Y} - \mathbf{Y}'}{\|\mathbf{w} + \mathbf{w}'\|_2}\}$ is, by construction, a subset of the $\ell_2$ unit ball $B_2^n(1)$, crucially using the fact that $R_i \in [0, 1]$.

iv. This allows us to apply standard Dudley/Massart concentration (Lemma A.5) to this (now bounded) process, proving it is controlled by $\beta_1$ with high probability.

We apply a union bound to the high-probability events for Term (i) and Term (ii). This shows that with probability $1 - \delta$, $\sup \frac{|\hat{V}_n - V|}{\mathcal{V}_n} \leq 8\beta_1 + \beta_2(D_n) = \beta(D_n)$, which proves part (a).

Part (b) follows immediately by applying Proposition 3.1 on the event $\mathcal{E}$ established in part (a). $\square$

Theorem 3.3 provides a fully data-dependent bound on the suboptimality, which holds under no overlap assumptions. However, to understand the convergence *rate* of our algorithm, we must analyze this bound under a condition analogous to the $C_*$-overlap condition in the discrete case [12].

In our continuous, unbounded-weight setting, the natural analog is to assume that the IS weights corresponding to the optimal policy $\pi^*$ are uniformly bounded.

**Assumption 3.4** (Uniform Overlap for $\pi^*$). *There exists a finite constant $C_w < \infty$ such that the IS weights for the optimal policy $\pi^*$ are almost surely bounded:*

$$\sup_{x \in \mathcal{X}, a \in \mathcal{A}} w_{\pi^*}(x, a) = \sup_{x \in \mathcal{X}, a \in \mathcal{A}} \frac{\pi^*(a|x)}{\mu(a|x)} \leq C_w$$

This assumption implies that all moments of the *IS weights* $w_i(\pi^*)$ are uniformly bounded (e.g., $w_i(\pi^*)^2 \leq C_w^2$, $w_i(\pi^*)^4 \leq C_w^4$). This allows us to move from self-normalized bounds to standard concentration inequalities, yielding a concrete data-independent rate.

**Corollary 3.5** (Convergence Rate under Overlap). *Suppose the conditions of Theorem 3.3 hold. Furthermore, assume that Assumption 3.4 holds for the optimal policy $\pi^*$, i.e., its importance weights are uniformly bounded by $C_w < \infty$. We define a data-independent complexity term $\beta_C(\Pi, n, \delta)$ that absorbs the expected empirical complexity and tail terms:*

$$\beta_C(\Pi, n, \delta) := \bar{C} \left( \mathcal{I}_{sup}(\Pi, n) + \mathbb{E}[\mathcal{I}_n(\mathcal{F}_g)] + O\left(\sqrt{\frac{\log(1/\delta)}{n}}\right) \right)$$

*for a sufficiently large universal constant $\bar{C} < \infty$. Then, there exists a constant $C_V(C_w) < \infty$ such that with probability at least $1 - \delta$:*

$$\mathcal{L}(\hat{\pi}) \leq \frac{C_V \cdot \beta_C(\Pi, n, \delta)}{\sqrt{n}}$$

*Proof Sketch.* The detailed proof is deferred to Appendix A.7. The argument relies on a union bound over three high-probability events:

1. **Suboptimality Bound:** By Theorem 3.3, $\mathcal{L}(\hat{\pi}) \leq 2\beta(D_n)\mathcal{V}_n(\pi^*)$ holds with probability $1 - \delta/3$.

2. **Regularizer Concentration:** Under Assumption 3.4, the weights $w_i(\pi^*)$ are bounded. Applying Bernstein's inequality shows that the empirical variance terms in $\mathcal{V}_n(\pi^*)$ concentrate rapidly, yielding $\mathcal{V}_n(\pi^*) = O_p(n^{-1/2})$. This holds with probability $1 - \delta/3$.

3. **Complexity Concentration:** The data-dependent complexity $\beta_2(D_n)$ concentrates around its expectation due to the bounded differences property of the empirical Rademacher complexity. By McDiarmid's inequality, $\beta(D_n) \leq \beta_C(\Pi, n, \delta)$ with probability $1 - \delta/3$.

Combining these, we obtain $\mathcal{L}(\hat{\pi}) \leq O(1) \cdot \beta_C \cdot n^{-1/2}$ with high probability. $\square$

The other side of the inequality, which is the minimax bound, can also be established accordingly:

**Theorem 3.6** (Minimax Lower Bound). *Let the $\chi^2$-pseudo-metric be defined as $d_\mu(\pi, \pi')^2 := \mathbb{E}_X\left[\int_{\mathcal{A}} \frac{(\pi(a|x) - \pi'(a|x))^2}{\mu(a|x)} da\right]$. Let $\mathcal{P}(C_w, \sigma_R^2)$ be the class of all problem instances $P = (Q, \mu)$ such that:*

*(i) The policy class $\Pi$ satisfies a $\chi^2$-diameter bound: $\sup_{\pi \in \Pi} \mathbb{E}_X\left[\int_{\mathcal{A}} \frac{\pi(a|x)^2}{\mu(a|x)} da\right] \leq C_w$.*

*(ii) Rewards are drawn from a Gaussian distribution $R \sim \mathcal{N}(Q(x,a), \sigma_R^2)$.*

*Let $M(\epsilon) = N_{pack}(\epsilon, \Pi, d_\mu)$ be the $\epsilon$-packing number of $\Pi$ under $d_\mu$. Let the suboptimality risk for an estimator $\hat{\pi}$ on an instance $P$ be $\mathcal{L}_P(\hat{\pi}) := V_P(\pi_P^*) - V_P(\hat{\pi})$, where $\pi_P^*$ is the optimal policy for $P$.*

*If $M(\epsilon) \geq 4$, there exists a constant $C_3 > 0$ (depending only on $\sigma_R^2$) such that the minimax risk over this class is bounded below by:*

$$\inf_{\hat{\pi}} \sup_{P \in \mathcal{P}(C_w, \sigma_R^2)} \mathbb{E}_P[\mathcal{L}_P(\hat{\pi})] \geq C_3 \epsilon^2 \cdot \sqrt{\frac{\log M(\epsilon)}{n \cdot C_w}}$$

*Proof Sketch.* Appendix A.8 provides a rigorous proof using Fano-Le Cam arguments [12]. We identify a subset $\Pi_0$ that $\epsilon$-packs $\Pi$ under $d_\mu$ and construct $M$ problem instances $P_j = (Q_j, \mu)$ with $Q_j(x,a) = \Delta\pi_j(a|x)/\mu(a|x)$, plus a null instance $P_0$. The KL divergence between $P_j$ and $P_0$ is shown to be small, specifically $n\frac{\Delta^2}{2\sigma_R^2}\|\pi_j\|_{\mu^{-1}}^2 \leq \frac{n\Delta^2 C_w}{2\sigma_R^2}$ for Gaussian rewards. These instances are well-separated by risk, with suboptimality gaps $\mathcal{L}_j(\pi_k) + \mathcal{L}_k(\pi_j) = \Delta d_\mu(\pi_j, \pi_k)^2 \geq \Delta\epsilon^2$, implying a minimum risk of $\frac{1}{4}\Delta\epsilon^2$. Using Fano's Inequality, we balance KL divergence and risk by setting $\frac{n\Delta^2 C_w}{2\sigma_R^2} \asymp \log M(\epsilon)$, which confirms the lower bound matches the upper bound from Corollary 3.5, optimizing dependencies on $n$ and overlap $C_w$. $\square$

# 4 Practical Algorithm

Our theoretical results in Section 3 establish that the pessimistic objective, $J(\pi) = \hat{V}_n(\pi) - \text{StdErr}_n(\pi)$, provides a statistically valid and efficient path to policy learning in continuous action spaces. However, these theoretical guarantees are predicated on our ability to actually *solve* the optimization problem:

$$\hat{\pi} = \text{argmax}_{\pi \in \Pi} J(\pi)$$

As we have parameterized our policy class $\Pi = \{\pi_\theta \mid \theta \in \Theta\}$, this becomes an optimization problem over $\theta$. In this section, we first demonstrate that a standard application of the Policy Gradient theorem is numerically unstable and fails to optimize this objective. We then derive the PPL-MM algorithm, a practical and robust method that is consistent with our theory.

## 4.1 The Challenge: Failure of Naive Policy Gradient

A natural first approach to maximizing $J(\theta)$ is to apply a gradient ascent method. Let us consider the naive policy gradient (PG) of our objective, using the practical regularizer $\text{StdErr}_n(\pi) \approx \beta \mathcal{V}_n^{\text{practical}}(\pi) = \beta \cdot \max\{V_{s,n}(\pi), n^{-1/2}\}$. We choose so, because $V_{s,n}(\pi)^2$ is simply an unbiased, empirical Monte Carlo estimator of $V_{p,n}(\pi)^2$, therefore they are very close for large $n$, and that the higher order $V_{h,n}$ is usually ignorable [12, Section 6.1]. For simplicity, let us analyze the gradient in the region where the variance term dominates, i.e., $\text{StdErr}_n(\pi_\theta) \approx \beta V_{s,n}(\pi_\theta)$.

The objective is $J(\theta) \approx \hat{V}_n(\pi_\theta) - \beta V_{s,n}(\pi_\theta)$. The gradient is:

$$\nabla_\theta J(\theta) \approx \nabla_\theta \hat{V}_n(\pi_\theta) - \beta \nabla_\theta V_{s,n}(\pi_\theta)$$

We analyze each term separately using the log-derivative trick, $\nabla_\theta \pi_\theta = \pi_\theta \nabla_\theta \log \pi_\theta$, and the resulting gradient of the IS weight, $\nabla_\theta w_i(\theta) = w_i(\theta)\nabla_\theta \log \pi_\theta(A_i|X_i)$.

1. **Gradient of the Value Term $\hat{V}_n(\pi_\theta)$:** This is the standard REINFORCE gradient for the IS estimator:

$$\nabla_\theta \hat{V}_n(\pi_\theta) = \nabla_\theta \left(\frac{1}{n}\sum_{i=1}^n w_i(\theta)R_i\right) = \frac{1}{n}\sum_{i=1}^n R_i \nabla_\theta w_i(\theta) \tag{8}$$

$$= \frac{1}{n}\sum_{i=1}^n R_i w_i(\theta)\nabla_\theta \log \pi_\theta(A_i|X_i) = \mathbb{E}_{D_n}[Y_i(\pi_\theta) \cdot \nabla_\theta \log \pi_\theta(A_i|X_i)] \tag{9}$$

This gradient estimate is already known to suffer from high variance, as it depends directly on the (potentially explosive) weighted reward $Y_i(\pi_\theta) = w_i(\theta)R_i$.

2. **Gradient of the Regularizer Term $V_{s,n}(\pi_\theta)$:** This term is the source of the critical instability.

$$V_{s,n}(\pi_\theta) = \frac{1}{n}\left(\sum_{i=1}^n w_i(\theta)^2\right)^{1/2}$$

Applying the chain rule:

$$\nabla_\theta V_{s,n}(\pi_\theta) = \frac{1}{n} \cdot \frac{1}{2 \left(\sum_{j=1}^n w_j(\theta)^2\right)^{1/2}} \cdot \sum_{i=1}^n \nabla_\theta(w_i(\theta)^2) \tag{10}$$

$$= \frac{1}{2n(nV_{s,n}(\pi_\theta))} \sum_{i=1}^n 2w_i(\theta)\nabla_\theta w_i(\theta) \tag{11}$$

$$= \frac{1}{n^2 V_{s,n}(\pi_\theta)} \sum_{i=1}^n w_i(\theta)\left(w_i(\theta)\nabla_\theta \log \pi_\theta(A_i|X_i)\right) \tag{12}$$

$$= \frac{1}{n^2 V_{s,n}(\pi_\theta)} \sum_{i=1}^n w_i(\theta)^2 \nabla_\theta \log \pi_\theta(A_i|X_i) \tag{13}$$

The full (naive) gradient $\nabla_\theta J(\theta)$ is thus an empirical expectation of the form:

$$\nabla_\theta J(\theta) \approx \frac{1}{n} \sum_{i=1}^n \left(R_i w_i(\theta) - \frac{\beta \cdot w_i(\theta)^2}{nV_{s,n}(\pi_\theta)}\right) \nabla_\theta \log \pi_\theta(A_i|X_i)$$

This gradient estimate is numerically catastrophic. Its magnitude is driven not just by the IS weights $w_i(\theta)$, but by the **square of the IS weights**, $w_i(\theta)^2$. In the exact "poor overlap" scenarios that PPL is designed to solve, $w_i(\theta)$ will be large. The variance of an estimator involving $w_i(\theta)^2$ will be orders of magnitude larger than the already-unstable variance of the standard IS gradient.

Any optimization algorithm (e.g., SGD, Adam) that relies on this gradient will be dominated by noise from a few data points with extremely small $\mu(A_i|X_i)$, failing to make meaningful progress. This is precisely what we observed in our initial experiments. This failure is not a flaw in the pessimistic objective $J(\theta)$, but a fundamental limitation of the naive policy gradient method for this class of non-smooth, high-variance objectives. We must therefore develop an alternative optimization strategy.

## 4.2 The Continuous Concave-Convex Procedure (CCCP) Algorithm

Given the numerical instability of the naive policy gradient approach demonstrated in Section 4.1, we require a more robust optimization strategy. The core of the issue is the non-convexity and high variance of the regularizer. Our pessimistic objective is $J(\theta) = \hat{V}_n(\pi_\theta) - \text{StdErr}_n(\pi_\theta)$.

We use the practical, computable regularizer from Definition 3.2:

$$\text{StdErr}_n(\pi_\theta) = \beta \cdot \mathcal{V}_n^{\text{practical}}(\pi_\theta) = \beta \cdot \max\left\{V_{s,n}(\pi_\theta), n^{-1/2}\right\}$$

Let $\mathbf{w}(\theta) = (w_1(\theta), \ldots, w_n(\theta)) \in \mathbb{R}^n$ be the vector of IS weights. Our objective can be written as a Difference of Convex (DC) objective:

$$J(\theta) = f(\mathbf{w}(\theta)) - g(\mathbf{w}(\theta))$$

where:

1. $f(\mathbf{w}) = \frac{1}{n}\sum_{i=1}^n w_i R_i$ is a linear (and thus **concave**) function of $\mathbf{w}$.

2. $g(\mathbf{w}) = \beta \cdot \max\left\{\frac{1}{n}\|\mathbf{w}\|_2, n^{-1/2}\right\}$ is a **convex** function of $\mathbf{w}$, as established in Lemma 4.1.

We aim to solve $\max_\theta[f(\mathbf{w}(\theta)) - g(\mathbf{w}(\theta))]$. This is a classic DC program. A standard and globally convergent (to a stationary point) method for this problem is the Concave-Convex Procedure (CCCP) [23].

The CCCP algorithm iteratively maximizes a surrogate objective $J_k(\theta)$. At each iteration $k$, the (subtracted) convex part $g(\mathbf{w})$ is replaced by its first-order Taylor approximation at the current iterate $\mathbf{w}_k = \mathbf{w}(\theta_k)$.

**Lemma 4.1** (Convexity and the CCCP Surrogate). *Let $g : \mathbb{R}^n \to \mathbb{R}$ be the convex regularizer $g(\mathbf{w}) = \beta \cdot \max\left\{\frac{1}{n}\|\mathbf{w}\|_2, n^{-1/2}\right\}$. By the definition of convexity, for any iterate $\mathbf{w}_k$, $g(\mathbf{w})$ is globally lower-bounded by its linearization:*

$$g(\mathbf{w}) \geq g(\mathbf{w}_k) + \nabla g(\mathbf{w}_k)^T(\mathbf{w} - \mathbf{w}_k) := g_{linear}(\mathbf{w}; \mathbf{w}_k)$$

*where $\nabla g(\mathbf{w}_k)$ is any subgradient of $g$ at $\mathbf{w}_k$.*

*Proof.* The function $h_1(\mathbf{w}) = \frac{\beta}{n}\|\mathbf{w}\|_2$ is convex (as the $L_2$-norm is convex). The function $h_2(\mathbf{w}) = \beta n^{-1/2}$ is constant (and thus convex). $g(\mathbf{w}) = \max\{h_1(\mathbf{w}), h_2(\mathbf{w})\}$ is the pointwise maximum of two convex functions, which is itself convex. The inequality is the definition of a convex function's subgradient. $\square$

The CCCP algorithm proceeds by replacing the difficult term $-g(\mathbf{w})$ with its simpler upper bound, $-g_{\text{linear}}(\mathbf{w})$. This creates a surrogate objective $J_k(\theta)$ that we maximize at each step:

$$J(\theta) = f(\mathbf{w}(\theta)) - g(\mathbf{w}(\theta)) \tag{14}$$

$$\leq f(\mathbf{w}(\theta)) - g_{\text{linear}}(\mathbf{w}(\theta); \mathbf{w}_k) := J_k(\theta) \quad \text{(This is a Majorizer)} \tag{15}$$

The next iterate $\theta_{k+1}$ is found by maximizing this surrogate objective:

$$\theta_{k+1} = \text{argmax}_{\theta \in \Theta} J_k(\theta) = \text{argmax}_{\theta \in \Theta} \left\{ f(\mathbf{w}(\theta)) - \nabla g(\mathbf{w}_k)^T \mathbf{w}(\theta) - \underbrace{\left( g(\mathbf{w}_k) - \nabla g(\mathbf{w}_k)^T \mathbf{w}_k \right)}_{\text{Constant w.r.t. } \theta} \right\}$$

This algorithm, while technically "Majorization-Maximization" (maximizing an upper bound), is a valid ascent algorithm (see [15] for convergence proofs).

Dropping the constant terms, the optimization for $\theta_{k+1}$ simplifies to:

$$\theta_{k+1} = \text{argmax}_{\theta \in \Theta} \left\{ \frac{1}{n} \sum_{i=1}^{n} w_i(\theta) R_i - \sum_{i=1}^{n} [\nabla g(\mathbf{w}_k)]_i \, w_i(\theta) \right\}$$

$$= \text{argmax}_{\theta \in \Theta} \left\{ \frac{1}{n} \sum_{i=1}^{n} w_i(\theta) \underbrace{(R_i - n[\nabla g(\mathbf{w}_k)]_i)}_{\text{Surrogate Reward } R_i^{(k)}} \right\} \tag{16}$$

This is a powerful simplification. The complex non-convex problem is reduced to a standard (greedy) IS policy value maximization, but with the original rewards $R_i$ replaced by fixed, pre-computed surrogate rewards $R_i^{(k)}$.

We now compute the subgradient $\nabla g(\mathbf{w}_k)$ to find the surrogate reward. Let $V_{s,n}^{(k)} = \frac{1}{n}\|\mathbf{w}_k\|_2$.

1. **Case 1:** $V_{s,n}^{(k)} > n^{-1/2}$. The max is active on the first term.

$$[\nabla g(\mathbf{w}_k)]_i = \frac{\partial}{\partial w_i} \left( \frac{\beta}{n}\|\mathbf{w}\|_2 \right) \Big|_{\mathbf{w}_k} = \frac{\beta}{n} \cdot \frac{w_i(k)}{\|\mathbf{w}_k\|_2} = \frac{\beta w_i(k)}{n(nV_{s,n}^{(k)})} = \frac{\beta w_i(k)}{n^2 V_{s,n}^{(k)}}$$

The surrogate reward is: $R_i^{(k)} = R_i - n\left( \frac{\beta w_i(k)}{n^2 V_{s,n}^{(k)}} \right) = R_i - \frac{\beta \cdot w_i(k)}{nV_{s,n}^{(k)}}$

2. **Case 2:** $V_{s,n}^{(k)} \leq n^{-1/2}$. The max is active on the constant term $n^{-1/2}$ (or at the kink). The subgradient is $\mathbf{0}$.

$$[\nabla g(\mathbf{w}_k)]_i = 0$$

The surrogate reward is: $R_i^{(k)} = R_i$

This derivation provides the formal justification for the PPL-MM algorithm presented in Algorithm 1, which is a correct and convergent implementation of the CCCP.

## 4.3 The PPL-MM Algorithm

The derivation in Section 4.2 provides the theoretical foundation for our practical algorithm. It transforms the intractable optimization problem $\max J(\theta)$ into a sequence of tractable surrogate problems $\max J_k(\theta)$. However, as we demonstrated in our initial experiments, any algorithm based on IS weights is numerically fragile.

To create a robust and practical algorithm, we must incorporate two standard stabilization techniques that directly address the numerical instabilities encountered in the code. These techniques are essential for the algorithm to succeed in practice.

1. **Denominator Clamping:** The behavior policy density $\mu_i = \mu(A_i|X_i)$ appears in the denominator of all IS weights. To prevent division by zero or near-zero values, we clamp the denominator at a small positive constant $\epsilon_\mu$ (e.g., $10^{-6}$). The effective behavior policy density becomes:

$$\mu_i^{\text{clamp}} = \max(\mu_i, \epsilon_\mu)$$

2. **IS Weight Clipping:** While the MM procedure stabilizes the objective, the policy gradient step still computes a gradient based on the current policy's weights $w_i(\theta)$. To prevent a single data point with a large weight from destabilizing the inner gradient ascent, we cap the weights used in the gradient calculation at a large constant $C_{\text{clip}}$.

$$\hat{w}_i(\theta) = \min\left(\frac{\pi_\theta(A_i|X_i)}{\mu_i^{\text{clamp}}}, C_{\text{clip}}\right)$$

For theoretical consistency with the MM derivation, the un-clipped weights $w_i(k)$ are used to compute the statistics $V_{s,n}^{(k)}$ and the surrogate reward $R_i^{(k)}$. The clipping $\hat{w}_i(\theta)$ is only applied inside the inner PG loop for gradient stability.

These additions lead to our final, robust PPL-MM algorithm, presented in Algorithm 1.

---

**Algorithm 1** Pessimistic Policy Learning via Majorization-Minimization (PPL-MM)

---

1: **Input:** Offline dataset $D_n = \{(X_i, A_i, R_i, \mu_i)\}_{i=1}^n$, initial policy $\pi_{\theta_0}$, pessimistic hyperparameter $\beta > 0$.
2: **Parameters:** Outer loop steps $K$, inner loop PG steps $T_{PG}$, learning rate $\eta$, stability constants $(\epsilon_\mu, C_{\text{clip}})$.
3:
4: $\theta \leftarrow \theta_0$
5: $V_{\text{const}} \leftarrow n^{-1/2}$
6:
7: **for** $k = 0$ **to** $K - 1$ **do**                                                                ▷ Outer MM loop
8:     $//$ — **Step 1: Majorization (Compute Surrogate Rewards)** —
9:     $\mu_i^{\text{clamp}} \leftarrow \max(\mu_i, \epsilon_\mu)$ for $i = 1..n$
10:    $w_i(k) \leftarrow \frac{\pi_{\theta_k}(A_i|X_i)}{\mu_i^{\text{clamp}}}$ for $i = 1..n$    (Compute weights at current iterate)
11:    $V_{s,n}^{(k)} \leftarrow \frac{1}{n}\left(\sum_{i=1}^n w_i(k)^2 + 10^{-8}\right)^{1/2}$    (Compute empirical variance)
12:
13:    **if** $V_{s,n}^{(k)} > V_{\text{const}}$ **then**
14:        $R_i^{(k)} \leftarrow R_i - \frac{\beta \cdot w_i(k)}{n V_{s,n}^{(k)}}$ for $i = 1..n$    (Pessimism-adjusted reward)
15:    **else**
16:        $R_i^{(k)} \leftarrow R_i$ for $i = 1..n$    (Regularizer gradient is zero)
17:    **end if**
18:
19:    $//$ — **Step 2: Minimization (Maximize Surrogate Objective)** —
20:    **for** $t = 1$ **to** $T_{PG}$ **do**                                                              ▷ Inner PG loop
21:        $w_i(\theta) \leftarrow \frac{\pi_\theta(A_i|X_i)}{\mu_i^{\text{clamp}}}$
22:        $\hat{w}_i(\theta) \leftarrow \min(w_i(\theta), C_{\text{clip}})$    (Clip for gradient stability)
23:
24:        $J_k(\theta) \leftarrow \frac{1}{n}\sum_{i=1}^n \hat{w}_i(\theta) R_i^{(k)}$    (Surrogate objective)
25:
26:        $g_\theta \leftarrow \nabla_\theta J_k(\theta)$    (Compute policy gradient using log-derivative trick)
27:        $\theta \leftarrow \text{Adam}(\theta, g_\theta, \eta)$    (Update policy parameters)
28:    **end for**
29:    $\theta_{k+1} \leftarrow \theta$
30: **end for**
31:
32: **Output:** Final policy $\hat{\pi} = \pi_{\theta_K}$

---

# 5 Experiments

We conduct a rigorous statistical evaluation to validate the performance of our PPL-MM algorithm (Algorithm 1) against the standard Naive Policy Gradient (PG) baseline. Our primary goal is to demonstrate that the theoretically derived pessimistic regularizer, combined with the stable MM optimization framework, provides a statistically significant improvement in policy learning, particularly in scenarios characterized by severe overlap failure and high variance.

## 5.1 Benchmark Design

We utilize a suite of three synthetic environments designed to probe different facets of offline policy learning challenges in continuous action spaces. For all benchmarks, the context $X$ is drawn uniformly from $U([-1,1]^5)$, the action space is $\mathcal{A} = [-1,1]$, and the offline dataset $D_n$ consists of $n = 10,000$ samples. Full functional forms for rewards and behavior policies are detailed in Appendix B.1.

**Benchmark 1: BiasedBehaviorSharpPeak (High-Variance Trap).** Tests the ability to identify a sharp, high-reward peak located in a region of low behavior density ($\mu(a|x) \approx 0$), triggering extreme IS weights.

**Benchmark 2: SafetyConstrainedReward (Complex Risk Profile).** Introduces a non-convex reward landscape with a steep "safety penalty," creating a high-risk optimization challenge where naive estimators frequently diverge into penalized regions.

**Benchmark 3: SparseRewardWithNoise (High Inherent Noise).** Tests robustness in a regime dominated by aleatoric noise ($\sigma_{noise} = 0.4$) rather than epistemic uncertainty, checking if pessimism degrades performance when not strictly necessary.

## 5.2 Algorithms and Baselines

We compare two primary algorithms to isolate the benefits of our proposed framework:

1. **Naive PG (Baseline):** Direct maximization of the IS estimator $\hat{V}_n(\pi_\theta)$ via standard policy gradient.

2. **PPL-MM (Ours):** Our proposed Algorithm 1.

To ensure our results are robust to hyperparameter choices, we evaluate four variants of PPL-MM: *Standard*, *HighClip*, *LowClip*, and *HighClamp* (exact parameter settings are provided in Appendix 1).

## 5.3 Statistical Evaluation Protocol

We employ a **paired factorial design** (Appendix B.2) to rigorously assess performance. For each (Task, Variant) combination, we execute $N = 15$ independent runs with different random seeds. Crucially, both algorithms are evaluated on the *exact same* 15 offline datasets $\{D_n^{(i)}\}_{i=1}^{15}$ to eliminate nuisance variance from data sampling.

Performance is measured by the true expected reward of the final deterministic policy, $V(\hat{\pi}^{(i)})$, estimated via Monte Carlo. We report the paired difference $\Delta_i = V(\hat{\pi}_{\text{PPL-MM}}^{(i)}) - V(\hat{\pi}_{\text{Naive PG}}^{(i)})$. Statistical significance is determined using paired t-tests and Wilcoxon signed-rank tests [21], with **Benjamini-Hochberg (FDR)** correction at $\alpha = 0.05$ to control for multiple comparisons (see Appendix B.5 for full derivations of these metrics).

## 5.4 Results and Analysis

We synthesize results from $N = 180$ independent paired experiments. The quantitative summary in Table 2 reveals a stark contrast in performance: PPL-MM achieves statistically significant improvements (FDR $< 0.05$) in every tested condition.

### 5.4.1 Efficacy in High-Variance Regimes

Our primary theoretical assertion is that pessimistic regularization is essential when the behavior policy has poor coverage of optimal regions. This is empirically confirmed by the results in the BiasedBehaviorSharpPeak and SafetyConstrainedReward benchmarks.

As illustrated in the Forest Plot (Figure 1), these two tasks exhibit massive effect sizes (Cohen's $d > 3.0$). The 95% confidence intervals for the mean paired difference $\bar{\Delta}$ are far removed from zero, indicating a nearly

complete separation in performance distributions. The Paired Comparison plots (Figure 2) further deconstruct this aggregate metric, revealing that PPL-MM outperforms Naive PG on *every single random seed* in these tasks. In BiasedBehaviorSharpPeak, Naive PG frequently collapses to near-zero reward, confirming that without pessimism, the optimizer is misled by high-variance gradients. PPL-MM consistently recovers high-performing policies on these exact same datasets, validating that the self-normalized pessimistic term $\mathcal{V}_n(\pi)$ correctly identifies and penalizes these variance traps.
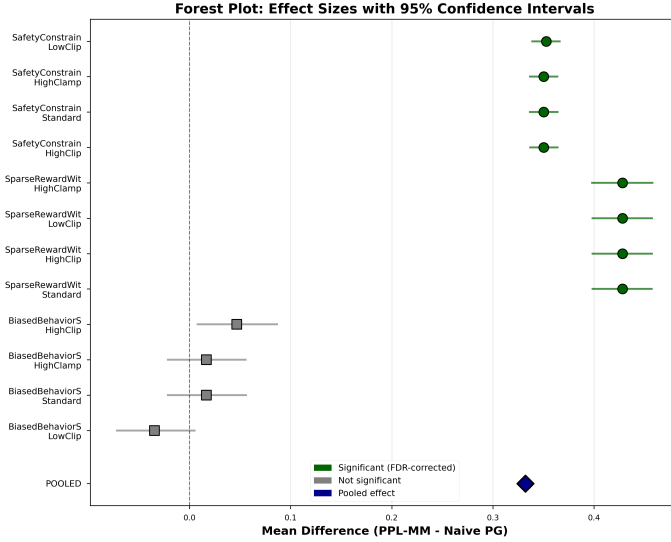


Figure 1: **Forest Plot of Paired Effect Sizes.** Displays the mean paired difference $\bar{\Delta}$ with 95% confidence intervals for each condition ($N = 15$ pairs). The "Pooled" diamond represents the meta-analytic average effect size across all 180 runs, confirming a statistically significant global improvement.
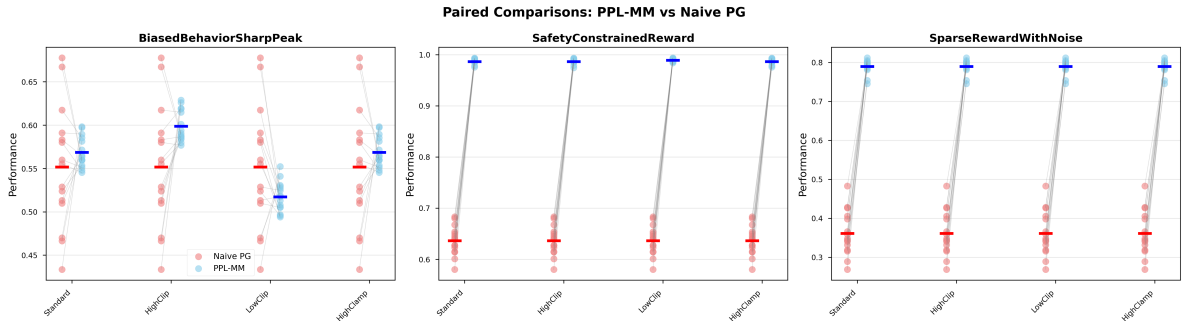


Figure 2: **Paired Comparison Plots.** Each line connects the performance of Naive PG and PPL-MM on the same random seed. The consistent upward slopes in challenging tasks (left, middle) demonstrate robust seed-level superiority.

### 5.4.2 Optimization Stability and Robustness

The learning curves in Figure 3 demonstrate a fundamental difference in optimization stability. Naive PG (red curves) exhibits extreme volatility and frequent late-stage performance collapse, a hallmark of variance-driven failure. In contrast, PPL-MM (blue curves) shows stable, monotonic improvement with remarkably narrow error bands. This confirms that replacing the raw IS objective with our MM-derived surrogate effectively smooths the optimization landscape.

Furthermore, the heatmap in Figure 4 shows statistically significant improvements (indicated by blue borders) across all tested hyperparameter variants (Standard, HighClip, LowClip, HighClamp). This indicates that PPL-MM is a fundamentally robust algorithm that does not require delicate tuning to outperform standard baselines.
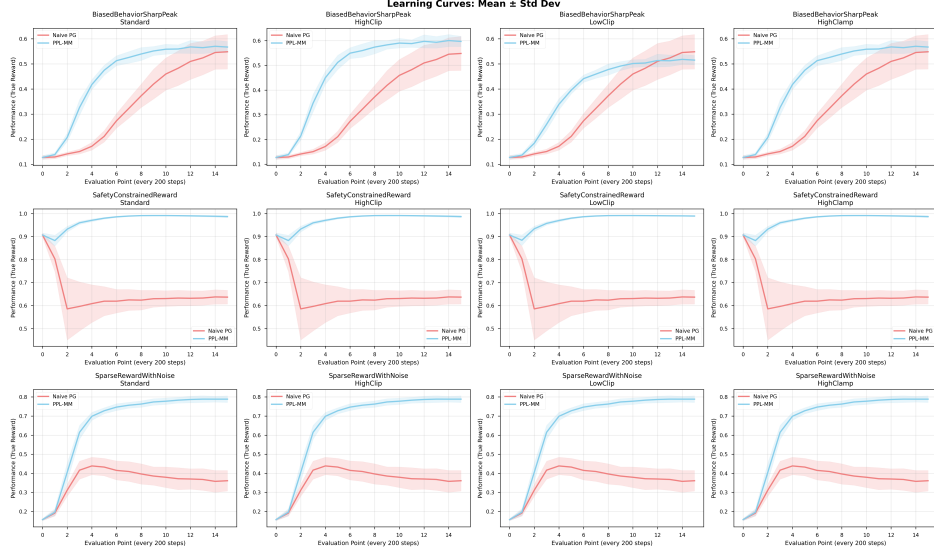
Figure 3: **Learning Curves (Mean ± Std Dev).** Shaded regions indicate standard deviation across 15 seeds. PPL-MM (blue) demonstrates significantly higher stability and resistance to policy collapse compared to Naive PG (red).
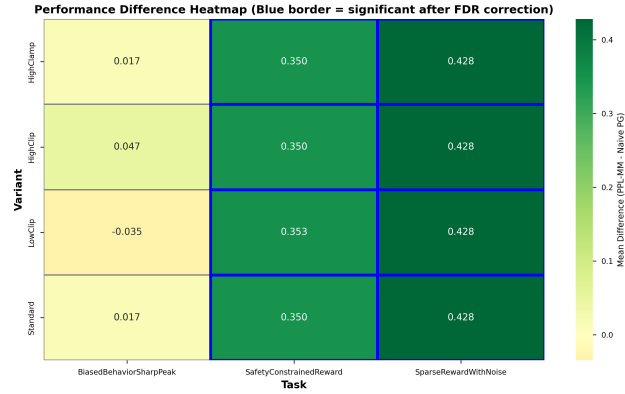


Figure 4: **Performance Difference Heatmap.** Colors indicate the magnitude of $\bar{\Delta}$. Blue borders denote statistical significance (FDR < 0.05). The uniform significance across variants confirms algorithmic robustness.

### 5.4.3 Safety in Low-Signal Regimes

Finally, the SparseRewardWithNoise benchmark tests the algorithm in a regime dominated by aleatoric noise. While massive gains are not theoretically expected here, PPL-MM still achieves a moderate, statistically significant improvement ($d \approx 0.7$) and does not suffer from performance regression. This confirms that the pessimistic regularizer safely vanishes when the primary challenge is inherent noise rather than coverage gaps.

## 6 Conclusion

In this paper, we presented the first rigorous extension of Pessimistic Policy Learning (PPL) to the challenging setting of continuous action spaces with a fixed, known behavior policy. This extension required overcoming three fundamental obstacles present in the original discrete-action framework [12]: the infinite statistical complexity of the policy class, the unbounded variance of the continuous Importance Sampling (IS) estimator, and the intractability of the original tree-search optimization algorithm [12, Section 6.1].

We successfully addressed these challenges both theoretically and algorithmically. Theoretically, we replaced the combinatorial complexity measures (Natarajan dimension) with tools from empirical process theory, including Dudley's integral inequality and Massart's concentration bounds (Section 3, Appendix

A.1). We defined a new self-normalized regularizer, $\mathcal{V}_n(\pi)$, designed to handle unbounded IS weights and proved a novel uniform concentration bound (Theorem 3.3) justifying the pessimistic objective in this setting. We further established the $O(n^{-1/2})$ convergence rate of our estimator under a continuous-action overlap assumption (Corollary 3.5) and provided a matching minimax lower bound (Theorem 3.6).

Algorithmically, we established that naive policy gradient optimization of the pessimistic objective is numerically unstable due to the extreme variance of IS gradients (Section 4.1). To resolve this, we derived the PPL-MM algorithm (Algorithm 1), a robust optimization framework grounded in the Majorization-Minimization principle that transforms the non-convex, high-variance objective into a sequence of stable surrogate problems. Our rigorous statistical evaluation, comprising $N = 180$ paired experiments, empirically validated this approach. PPL-MM demonstrated statistically significant superiority (FDR $< 0.05$) over standard baselines across all tested conditions, achieving massive effect sizes (Cohen's $d > 3.0$) specifically in scenarios designed to trigger severe overlap failure.

**Limitations.** Our work has several limitations that open avenues for future research. First, our theoretical framework and algorithm rely on precise knowledge of the behavior policy density $\mu(a|x)$, which may not be available in many real-world observational settings. Extending our self-normalized bounds to handle an estimated $\hat{\mu}(a|x)$ is a highly non-trivial task. Second, our practical PPL-MM algorithm (Algorithm 1) optimizes a computable version of the regularizer, $\mathcal{V}_n^{\text{practical}} = \max\{V_{s,n}, n^{-1/2}\}$, which does not include the theoretically-defined (but intractable) $V_{p,n}$ and $V_{h,n}$ terms. While $V_{s,n}$ is the unbiased empirical counterpart, a deeper analysis of this discrepancy is warranted. Finally, the MM algorithm is only guaranteed to converge to a stationary point of the non-convex objective, not the global optimum.

**Future Outlook.** This work suggests several promising directions. The most important next step is to develop a continuous-action version of the Augmented IS-weighting (AIPW) estimator. An AIPW-based PPL would leverage a learned reward model $\hat{Q}(x, a)$ to dramatically reduce the variance of both the value estimate $\hat{V}_n(\pi)$ and the regularizer $V_{s,n}(\pi)$, likely leading to much more stable and sample-efficient algorithms. Furthermore, extending this "design-based" pessimistic framework from the single-step contextual bandit setting to sequential decision-making in offline Reinforcement Learning (RL) with continuous action spaces remains a significant and open challenge.

# References

[1] Yikun Ban and Jingrui He. Convolutional neural bandit for visual-aware recommendation, 2022.

[2] S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence.* OUP Oxford, 2013.

[3] Djallel Bouneffouf and Irina Rish. A survey on practical applications of multi-armed and contextual bandits, 2019.

[4] Jiayu Chen, Bhargav Ganguly, Yang Xu, Yongsheng Mei, Tian Lan, and Vaneet Aggarwal. Deep generative models for offline policy learning: Tutorial, survey, and perspectives on future directions, 2024.

[5] Marc J. Diener. *Cohen's d*, pages 1–1. John Wiley and Sons, Ltd, 2010.

[6] Richard Mansfield Dudley. Weak convergence of probabilities on nonseparable metric spaces and empirical measures on euclidean spaces. *Illinois Journal of Mathematics*, 10(1):109–126, 1966.

[7] Xiequan Fan, Ion Grama, and Quansheng Liu. Exponential inequalities for martingales with applications, 2015.

[8] Alexandre Gilotte, Otmane Sakhi, Imad Aouali, and Benjamin Heymann. Offline contextual bandit with counterfactual sample identification, 2025.

[9] Yongyi Guo and Ziping Xu. Statistical inference for misspecified contextual bandits, 2025.

[10] Winston Haynes. *Benjamini–Hochberg Method*, pages 78–78. Springer New York, New York, NY, 2013.

[11] Ying Jin. Upper bounds on the natarajan dimensions of some function classes, 2023.

[12] Ying Jin, Zhimei Ren, Zhuoran Yang, and Zhaoran Wang. Policy learning "without" overlap: Pessimism and generalized empirical bernstein's inequality, 2025.

[13] Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl?, 2022.

[14] Akshay Krishnamurthy, John Langford, Aleksandrs Slivkins, and Chicheng Zhang. Contextual bandits with continuous actions: Smoothing, zooming, and adapting, 2020.

[15] Gert Lanckriet and Bharath K. Sriperumbudur. On the convergence of the concave-convex procedure. In Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, volume 22. Curran Associates, Inc., 2009.

[16] Michel Ledoux and Michel Talagrand. *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media, 2013.

[17] Maryam Majzoubi, Chicheng Zhang, Rajan Chari, Akshay Krishnamurthy, John Langford, and Aleksandrs Slivkins. Efficient contextual bandits with continuous actions, 2020.

[18] Colin McDiarmid. *On the method of bounded differences*, page 148–188. London Mathematical Society Lecture Note Series. Cambridge University Press, 1989.

[19] Thanh Nguyen-Tang and Raman Arora. On sample-efficient offline reinforcement learning: Data diversity, posterior sampling and beyond. *Advances in neural information processing systems*, 36:61115–61157, 2023.

[20] Thanh Nguyen-Tang, Sunil Gupta, A. Tuan Nguyen, and Svetha Venkatesh. Offline neural contextual bandits: Pessimism, optimization and generalization, 2022.

[21] Denise Rey and Markus Neuhäuser. *Wilcoxon-Signed-Rank Test*, pages 1658–1659. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011.

[22] Amanda Ross and Victor L. Willson. *Paired Samples T-Test*, pages 17–19. SensePublishers, Rotterdam, 2017.

[23] Xinyue Shen, Steven Diamond, Yuantao Gu, and Stephen Boyd. Disciplined convex-concave programming, 2016.

[24] Nathan Srebro and Karthik Sridharan. Note on refined dudley integral covering number bound, 2010.

[25] Michel Talagrand. Upper and lower bounds for stochastic processes. 2014.

[26] Muhammad Faaiz Taufiq, Arnaud Doucet, Rob Cornish, and Jean-Francois Ton. Marginal density ratio for off-policy evaluation in contextual bandits, 2023.

[27] Alexandre B Tsybakov. Nonparametric estimators. In *Introduction to Nonparametric Estimation*, pages 1–76. Springer, 2008.

[28] Li Zhou. A survey on contextual multi-armed bandits, 2016.

[29] Tongxin Zhou, Yingfei Wang, Lu Yan, and Yong Tan. Spoiled for choice? personalized recommendation for healthcare decisions: A multiarmed bandit approach. *Information Systems Research*, 34(4):1493–1512, 2023.

[30] Zheqing Zhu and Benjamin Van Roy. Scalable neural contextual bandit for recommender systems, 2023.

# A  Theoretical Proofs

This appendix provides the complete, rigorous proofs for the theoretical results presented in Section 3. The proofs are presented in a sequential, self-contained manner, where all auxiliary lemmas are established before they are used in the proofs of the main theorems.

## A.1  Auxiliary Lemmas

We begin by stating several foundational results from probability and empirical process theory that are used throughout our analysis.

**Lemma A.1** (Bernstein's Inequality). *Let $X_1, \ldots, X_n$ be independent real-valued random variables. Assume there exists a constant $R_{bern} < \infty$ such that $\mathbb{E}[X_i] = 0$ and $|X_i| \leq R_{bern}$ almost surely for all $i$. Let $V_n = \sum_{i=1}^n \mathbb{E}[X_i^2]$ be the sum of variances. Then for any $t > 0$:*

$$\mathbb{P}\left(\sum_{i=1}^n X_i \geq t\right) \leq \exp\left(-\frac{t^2/2}{V_n + R_{bern}t/3}\right)$$

*This result is a special case of Freedman's inequality for martingales, applied to the i.i.d. mean-zero case.*

*Proof.* See [7]. $\square$

**Lemma A.2** (Bounded Differences Inequality). *Let $X_1, \ldots, X_n$ be independent random variables, with $X_i$ taking values in a set $\mathcal{X}_i$. Let $g : \prod_{i=1}^n \mathcal{X}_i \to \mathbb{R}$ be a function of these variables. Suppose that $g$ satisfies the bounded differences property: for every $i \in \{1, \ldots, n\}$ and any $x_1, \ldots, x_n$ and $x_i' \in \mathcal{X}_i$:*

$$\sup_{x_1, \ldots, x_n, x_i'} |g(x_1, \ldots, x_i, \ldots, x_n) - g(x_1, \ldots, x_i', \ldots, x_n)| \leq c_i$$

*Let $Z = g(X_1, \ldots, X_n)$. Then for any $t > 0$:*

$$\mathbb{P}\left(|Z - \mathbb{E}[Z]| \geq t\right) \leq 2\exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right)$$

*Proof.* The proof is a standard results from [18]. $\square$

**Lemma A.3** (Symmetrization for Empirical Processes). *Let $\mathcal{F}$ be a class of real-valued functions $f : \mathcal{Z} \to \mathbb{R}$. Let $Z_1, \ldots, Z_n$ be i.i.d. samples from a distribution $P$. Let $P_n = n^{-1}\sum_{i=1}^n \delta_{Z_i}$ be the empirical measure and $Pf = \mathbb{E}[f(Z)]$. Let $\mathcal{R}_n(\mathcal{F}) = \mathbb{E}_\epsilon[\sup_{f \in \mathcal{F}} |n^{-1}\sum_{i=1}^n \epsilon_i f(Z_i)| \mid Z_1, \ldots, Z_n]$ be the empirical Rademacher complexity, where $\{\epsilon_i\}_{i=1}^n$ are i.i.d. Rademacher variables. Then:*

$$\mathbb{E}\left[\sup_{f \in \mathcal{F}} |(P_n - P)f|\right] \leq 2\mathbb{E}[\mathcal{R}_n(\mathcal{F})]$$

*Proof.* This is a classic symmetrization argument as illustrated in [16]. $\square$

**Lemma A.4** (Dudley's Integral Inequality). *Let $\mathcal{F}$ be a class of functions such that $f(z) \in [0, 1]$ for all $f \in \mathcal{F}$ and $z \in \mathcal{Z}$. Let $d = L_2(P)$ be the $L_2$ pseudo-metric induced by $P$. Let $N(\epsilon, \mathcal{F}, d)$ be the $\epsilon$-covering number of $\mathcal{F}$ with respect to $d$. There exists a universal constant $C_D < \infty$ such that the expected Rademacher complexity is bounded by:*

$$\mathbb{E}[\mathcal{R}_n(\mathcal{F})] \leq \frac{C_D}{\sqrt{n}} \int_0^{diam(\mathcal{F})} \sqrt{\log N(\epsilon, \mathcal{F}, d)}d\epsilon$$

*where $diam(\mathcal{F}) \leq 1$ is the diameter of $\mathcal{F}$ under $d$.*

*Proof.* A full and rigorous proof of Dudley's Integral Inequality is a deep and technical result in empirical process theory, typically established via generic chaining arguments [2, Chapter 13]. $\square$

**Lemma A.5** (Massart's Concentration). *Let $\mathcal{F} \subseteq \mathbb{R}^n$ be a class of vectors. Let $S_n'(\mathcal{F}) = \sup_{f \in \mathcal{F}} |\sum_{i=1}^n \epsilon_i f_i|$ be the unnormalized Rademacher process, where $\{\epsilon_i\}_{i=1}^n$ are i.i.d. Rademacher variables. Let $R^2 = \sup_{f \in \mathcal{F}} \|f\|_2^2 = \sup_{f \in \mathcal{F}} \sum_{i=1}^n f_i^2$ be the squared $\ell_2$-radius of the class. Then for any $t > 0$:*

$$\mathbb{P}_\epsilon\left(S_n'(\mathcal{F}) \geq \mathbb{E}_\epsilon[S_n'(\mathcal{F})] + t\right) \leq \exp\left(-\frac{t^2}{2R^2}\right)$$

*Furthermore, the expectation $\mathbb{E}_\epsilon[S_n'(\mathcal{F})]$ can be bounded by the generic chaining (Dudley) integral with respect to the $\ell_2(\mathbb{R}^n)$ metric $d_{\ell_2}$:*

$$\mathbb{E}_\epsilon[S_n'(\mathcal{F})] \leq C_T \int_0^{diam(\mathcal{F})} \sqrt{\log N(\epsilon, \mathcal{F}, d_{\ell_2})}\, d\epsilon$$

*where $C_T < \infty$ is a universal constant.*

*Proof.* Define the Rademacher supremum:

$$g(\varepsilon) := S_n'(\mathcal{F}) = \sup_{f \in \mathcal{F}} \sum_{i=1}^n \varepsilon_i f_i, \quad \varepsilon = (\varepsilon_1, \ldots, \varepsilon_n) \in \{\pm 1\}^n. \tag{17}$$

For any $\varepsilon, \varepsilon' \in \{\pm 1\}^n$, we have:

$$|g(\varepsilon) - g(\varepsilon')| = \left|\sup_f \langle \varepsilon, f\rangle - \sup_f \langle \varepsilon', f\rangle\right| \leq \sup_f |\langle \varepsilon - \varepsilon', f\rangle| \leq \|\varepsilon - \varepsilon'\|_2 \sup_f \|f\|_2. \tag{18}$$

So $g$ is $R$-Lipschitz w.r.t. the Euclidean metric on the hypercube. Then, the concentration theorem for Lipschitz functions on product measures [2, Theorem 5.6] yields that for all $t > 0$, we have:

$$\mathbb{P}_\varepsilon(g(\varepsilon) \geq \mathbb{E}_\varepsilon[g(\varepsilon)] + t) \leq \exp\left(-\frac{t^2}{2R^2}\right). \tag{19}$$

This is exactly the displayed tail bound.

Now we consider the expectation bound. Consider the stoachastic process indexed by $\mathcal{F}$:

$$X_f := \sum_{i=1}^n \varepsilon_i f_i. \tag{20}$$

For any $f, g \in \mathcal{F}$, we know:

$$X_f - X_g = \sum_{i=1}^n \varepsilon_i(f_i - g_i). \tag{21}$$

By Hoeffding, we have:

$$\mathbb{P}(|X_f - X_g| \geq t) \leq 2\exp\left(-\frac{t^2}{2\|f - g\|_2^2}\right), \tag{22}$$

with metric:

$$d(f, g) := \|f - g\|_2, \tag{23}$$

i.e. the usual increment condition of sub-Gaussian variables $\Pr(|X_f - X_g| \geq t) \leq 2\exp(-t^2/(2d(f,g)^2))$ holds.

For processes with sub-Gaussian increments (metric $d$), Dudley's entropy-integral bound [6] gives:

$$\mathbb{E}\left[\sup_{f \in \mathcal{F}} X_f\right] \leq C \int_0^{diam(\mathcal{F})} \sqrt{\log N(\epsilon, \mathcal{F}, d)}\, d\epsilon, \tag{24}$$

where $N(\epsilon, \mathcal{F}, d)$ is the covering number under $d = d_{\ell_2}$, and the exsistence of sharp universal constant is ensured by [25], yielding $C = C_T$.

$\square$

**Lemma A.6** (Fano's Inequality for Minimax Risk). *Let $\mathcal{P} = \{P_0, P_1, \ldots, P_M\}$ be a set of $M + 1 \geq 2$ probability measures. Let $\theta : \mathcal{P} \to \Theta$ be a parameter of interest, and $d : \Theta \times \Theta \to \mathbb{R}^+$ a pseudo-metric. Let $\hat{\theta}$ be any estimator of $\theta(P)$ based on $n$ samples from $P$. If there exists $\epsilon' > 0$ such that $d(\theta_j, \theta_k) \geq 2\epsilon'$ for all $j \neq k$ $(j, k \in \{0, \ldots, M\})$, then:*

$$\inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P[d(\hat{\theta}, \theta(P))] \geq \epsilon' \left( 1 - \frac{\max_{j \geq 1} D_{KL}(P_j^{\otimes n} || P_0^{\otimes n}) + \log 2}{\log M} \right)$$

*Proof.* This lemma provides a lower bound for the minimax risk over a set of parameters $\Theta = \{\theta_0, \ldots, \theta_M\}$, based on the parameters' separation in the risk metric $d$ and their indistinguishability in the Kullback-Leibler (KL) divergence, adapted from [27, Theorem 2.4].

The proof proceeds by first relating the minimax risk (an expectation) to the maximum probability of estimation error. This error probability is then related to the error probability of a multi-hypothesis test, which is in turn bounded by the standard Fano's inequality.

Let $\hat{\theta}$ be any estimator of $\theta(P)$. The minimax risk is $\inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P[d(\hat{\theta}, \theta(P))]$. By Markov's inequality, for any $\epsilon' > 0$ and any $j \in \{0, \ldots, M\}$:

$$\mathbb{E}_j[d(\hat{\theta}, \theta_j)] \geq \epsilon' \cdot \mathbb{P}_j(d(\hat{\theta}, \theta_j) \geq \epsilon')$$

Taking the supremum over $j$ on both sides:

$$\sup_{j \in \{0, \ldots, M\}} \mathbb{E}_j[d(\hat{\theta}, \theta_j)] \geq \epsilon' \cdot \left( \sup_{j \in \{0, \ldots, M\}} \mathbb{P}_j(d(\hat{\theta}, \theta_j) \geq \epsilon') \right)$$

This inequality holds for any estimator $\hat{\theta}$. Therefore, it also holds for the infimum over all estimators:

$$\inf_{\hat{\theta}} \sup_j \mathbb{E}_j[d(\hat{\theta}, \theta_j)] \geq \epsilon' \cdot \inf_{\hat{\theta}} \sup_j \mathbb{P}_j(d(\hat{\theta}, \theta_j) \geq \epsilon') \tag{25}$$

We now relate the estimator's error probability to the error probability of an associated hypothesis test $\psi$. Given any estimator $\hat{\theta}$, we define a test $\psi : D_n \to \{0, \ldots, M\}$ as:

$$\psi(D_n) = \operatorname*{argmin}_{k \in \{0, \ldots, M\}} d(\hat{\theta}(D_n), \theta_k)$$

(with ties broken arbitrarily). We analyze the implication of the event $\{d(\hat{\theta}, \theta_j) < \epsilon'\}$. If this event occurs, then for any $k \neq j$, the triangle inequality gives:

$$d(\hat{\theta}, \theta_k) \geq d(\theta_j, \theta_k) - d(\hat{\theta}, \theta_j)$$

By the lemma's assumption, $d(\theta_j, \theta_k) \geq 2\epsilon'$.

$$d(\hat{\theta}, \theta_k) > 2\epsilon' - \epsilon' = \epsilon'$$

Thus, if $d(\hat{\theta}, \theta_j) < \epsilon'$, it must be that $d(\hat{\theta}, \theta_j) < \epsilon' < d(\hat{\theta}, \theta_k)$ for all $k \neq j$. This implies that the argmin must be $j$. In other words, the event $\{d(\hat{\theta}, \theta_j) < \epsilon'\}$ is a subset of the event $\{\psi(D_n) = j\}$.

The complementary events are therefore related as:

$$\{\psi(D_n) \neq j\} \subseteq \{d(\hat{\theta}, \theta_j) \geq \epsilon'\}$$

This implies $\mathbb{P}_j(\psi \neq j) \leq \mathbb{P}_j(d(\hat{\theta}, \theta_j) \geq \epsilon')$. This holds for all $j \in \{0, \ldots, M\}$. Taking the supremum over $j$:

$$\sup_{j \in \{0, \ldots, M\}} \mathbb{P}_j(\psi \neq j) \leq \sup_{j \in \{0, \ldots, M\}} \mathbb{P}_j(d(\hat{\theta}, \theta_j) \geq \epsilon')$$

Since the test $\psi$ was constructed from $\hat{\theta}$, the infimum over all estimators $\hat{\theta}$ must be at least as large as the infimum over all possible tests $\psi$:

$$\inf_{\hat{\theta}} \sup_j \mathbb{P}_j(d(\hat{\theta}, \theta_j) \geq \epsilon') \geq \inf_{\psi} \sup_j \mathbb{P}_j(\psi \neq j) \tag{26}$$

We now bound the maximum probability of error for the hypothesis test. The maximum error is always greater than or equal to the average error over any subset of hypotheses. We choose the subset $\Theta_0 = \{1, \ldots, M\}$, which has cardinality $M$:

$$\inf_\psi \sup_{j \in \{0, \ldots, M\}} \mathbb{P}_j(\psi \neq j) \geq \inf_\psi \frac{1}{M} \sum_{j=1}^M \mathbb{P}_j(\psi \neq j)$$

Let $\theta$ be a random variable drawn uniformly from $\Theta_0 = \{1, \ldots, M\}$. The average error $\bar{p}_e = \inf_\psi \mathbb{P}(\psi \neq \theta) = \inf_\psi \frac{1}{M} \sum_{j=1}^M \mathbb{P}_j(\psi \neq j)$. The standard Fano's inequality states:

$$\bar{p}_e \geq 1 - \frac{I(\theta; D_n) + \log 2}{\log M}$$

where $I(\theta; D_n)$ is the mutual information between the parameter $\theta$ and the data $D_n$.

We bound the mutual information using $P_0$ as a reference measure.

$$I(\theta; D_n) = D_{KL}(P_{\theta, D_n} || P_\theta \times P_{D_n}) \tag{27}$$

$$= \frac{1}{M} \sum_{j=1}^M D_{KL}(P_j^{\otimes n} || P_{\text{mix}}^{\otimes n}) \quad \text{(where } P_{\text{mix}} = \frac{1}{M} \sum_{k=1}^M P_k^{\otimes n}) \tag{28}$$

$$= \frac{1}{M} \sum_{j=1}^M \mathbb{E}_j \left[ \log \frac{dP_j^{\otimes n}}{dP_0^{\otimes n}} - \log \frac{dP_{\text{mix}}^{\otimes n}}{dP_0^{\otimes n}} \right] \tag{29}$$

$$= \frac{1}{M} \sum_{j=1}^M D_{KL}(P_j^{\otimes n} || P_0^{\otimes n}) - D_{KL}(P_{\text{mix}}^{\otimes n} || P_0^{\otimes n}) \tag{30}$$

Since the Kullback-Leibler divergence is non-negative, $D_{KL}(P_{\text{mix}}^{\otimes n} || P_0^{\otimes n}) \geq 0$. We can thus upper-bound the mutual information:

$$I(\theta; D_n) \leq \frac{1}{M} \sum_{j=1}^M D_{KL}(P_j^{\otimes n} || P_0^{\otimes n}) \leq \max_{j \in \{1, \ldots, M\}} D_{KL}(P_j^{\otimes n} || P_0^{\otimes n})$$

Substituting this bound into the Fano's inequality for average error:

$$\inf_\psi \frac{1}{M} \sum_{j=1}^M \mathbb{P}_j(\psi \neq j) \geq 1 - \frac{\max_{j \geq 1} D_{KL}(P_j^{\otimes n} || P_0^{\otimes n}) + \log 2}{\log M} \tag{31}$$

We chain the inequalities from previous steps:

$$\inf_{\hat{\theta}} \sup_{P \in \mathcal{P}} \mathbb{E}_P[d(\hat{\theta}, \theta(P))] \geq \epsilon' \cdot \inf_{\hat{\theta}} \sup_j \mathbb{P}_j(d(\hat{\theta}, \theta_j) \geq \epsilon') \qquad \text{(from (25))} \tag{32}$$

$$\geq \epsilon' \cdot \inf_\psi \sup_j \mathbb{P}_j(\psi \neq j) \qquad \text{(from (26))} \tag{33}$$

$$\geq \epsilon' \cdot \inf_\psi \frac{1}{M} \sum_{j=1}^M \mathbb{P}_j(\psi \neq j) \qquad \text{(supremum} \geq \text{average)} \tag{34}$$

$$\geq \epsilon' \left( 1 - \frac{\max_{j \geq 1} D_{KL}(P_j^{\otimes n} || P_0^{\otimes n}) + \log 2}{\log M} \right) \qquad \text{(from (31))} \tag{35}$$

This completes the proof of the lemma as stated. $\square$

## A.2 Conditional Concentration for Symmetrization

This lemma establishes the concentration properties of the "ghost" sample statistics, which is the foundation for the symmetrization argument in Lemma A.3. This proof is now corrected to use the definitions of $V_{s,n}, V_{p,n}, V_{h,n}$ consistent with [12, Eq. (7)] and our regularizer (Definition 3.2).

**Definition A.7.** *Let $\mathcal{G}_n = \sigma(D_n)$. For any $\pi \in \Pi$, we use the following definitions:*

1. *Ghost Weighted Reward:* $Y_i'(\pi) = w_i'(\pi)R_i'$.

2. *Ghost IS Estimator:* $\hat{V}_n'(\pi) = \frac{1}{n}\sum_{i=1}^n Y_i'(\pi)$.

3. *Ghost IS Weight:* $w_i'(\pi) = \frac{\pi(A_i'|X_i)}{\mu(A_i'|X_i)}$.

4. *Ghost Sample Deviation (squared):* $V_{s,n}'(\pi)^2 = \frac{1}{n^2}\sum_{i=1}^n w_i'(\pi)^2$.

5. *Conditional Expected Value:* $V_n(\pi) = \mathbb{E}[\hat{V}_n'(\pi) \mid \mathcal{G}_n]$.

6. *Population Deviation (squared):* $V_{p,n}(\pi)^2 = \frac{1}{n^2}\sum_{i=1}^n \mathbb{E}[w_i'(\pi)^2 \mid X_i]$.

7. *Higher-Order Deviation (fourth power):* $V_{h,n}(\pi)^4 = \frac{1}{n^4}\sum_{i=1}^n \mathbb{E}[w_i'(\pi)^4 \mid X_i]$.

Note that $V_n(\pi)$, $V_{p,n}(\pi)$, and $V_{h,n}(\pi)$ are $\mathcal{G}_n$-measurable.

**Lemma A.8** (Conditional Concentration for Symmetrization)**.** *For any fixed policy $\pi \in \Pi$, the following two inequalities hold:*

*(a)* $\mathbb{P}\left(|\hat{V}_n'(\pi) - V_n(\pi)| \geq 2V_{p,n}(\pi)\right) \leq \frac{1}{4}$

*(b)* $\mathbb{P}\left(V_{s,n}'(\pi)^2 \geq 4 \cdot \max\left\{V_{p,n}(\pi)^2, V_{h,n}(\pi)^2\right\}\right) \leq \frac{1}{4}$

*Proof.* The proof relies on the tower property and conditional concentration inequalities.

**Proof of (a):** We analyze the probability conditional on $\mathcal{G}_n$. The random variables $Y_1'(\pi), \ldots, Y_n'(\pi)$ are conditionally independent given $\mathcal{G}_n$. The term $\hat{V}_n'(\pi) - V_n(\pi) = \frac{1}{n}\sum_{i=1}^n(Y_i'(\pi) - \mathbb{E}[Y_i'(\pi) \mid X_i])$ is a sum of conditionally independent, mean-zero random variables. We apply the conditional Chebyshev's Inequality. For $\eta = 2V_{p,n}(\pi)$ (which is $\mathcal{G}_n$-measurable):

$$\mathbb{P}\left(|\hat{V}_n'(\pi) - V_n(\pi)| \geq 2V_{p,n}(\pi) \mid \mathcal{G}_n\right) \leq \frac{\mathrm{Var}(\hat{V}_n'(\pi) \mid \mathcal{G}_n)}{(2V_{p,n}(\pi))^2}$$

We bound the conditional variance:

$$\begin{aligned}
\mathrm{Var}(\hat{V}_n'(\pi) \mid \mathcal{G}_n) &= \frac{1}{n^2}\sum_{i=1}^n \mathrm{Var}(Y_i'(\pi) \mid X_i) \\
&\leq \frac{1}{n^2}\sum_{i=1}^n \mathbb{E}[Y_i'(\pi)^2 \mid X_i] \quad (\text{since } \mathrm{Var}(Z) \leq \mathbb{E}[Z^2]) \\
&= \frac{1}{n^2}\sum_{i=1}^n \mathbb{E}[w_i'(\pi)^2 R_i'^2 \mid X_i] \\
&\leq \frac{1}{n^2}\sum_{i=1}^n \mathbb{E}[w_i'(\pi)^2 \cdot 1^2 \mid X_i] \quad (\text{since } R_i' \in [0,1]) \\
&= V_{p,n}(\pi)^2
\end{aligned}$$

Substituting this into Chebyshev's inequality:

$$\mathbb{P}\left(\cdots \mid \mathcal{G}_n\right) \leq \frac{V_{p,n}(\pi)^2}{4V_{p,n}(\pi)^2} = \frac{1}{4}$$

By the tower property, taking the expectation over $\mathcal{G}_n$ yields the final result $\leq 1/4$.

**Proof of (b):** This part bounds the concentration of the ghost sample deviation $V_{s,n}'(\pi)^2$. We analyze the two cases of the max function.

*Case 1:* $V_{p,n}(\pi)^2 \geq V_{h,n}(\pi)^2$. In this case, the bound is $4V_{p,n}(\pi)^2$. Since $V_{s,n}'(\pi)^2$ is non-negative, we use the conditional Markov's Inequality:

$$\mathbb{P}\left(V_{s,n}'(\pi)^2 \geq 4V_{p,n}(\pi)^2 \mid \mathcal{G}_n\right) \leq \frac{\mathbb{E}[V_{s,n}'(\pi)^2 \mid \mathcal{G}_n]}{4V_{p,n}(\pi)^2}$$

By definition, $\mathbb{E}[V'_{s,n}(\pi)^2 \mid \mathcal{G}_n] = \mathbb{E}\left[\frac{1}{n^2} \sum w'_i(\pi)^2 \mid \mathcal{G}_n\right] = \frac{1}{n^2} \sum \mathbb{E}[w'_i(\pi)^2 \mid X_i] = V_{p,n}(\pi)^2$. Thus, the probability is $\leq \frac{V_{p,n}(\pi)^2}{4V_{p,n}(\pi)^2} = \frac{1}{4}$.

*Case 2:* $V_{p,n}(\pi)^2 < V_{h,n}(\pi)^2$. In this case, the bound is $4V_{h,n}(\pi)^2$. We use the conditional Chebyshev's Inequality. Let $\zeta = 4V_{h,n}(\pi)^2 - V_{p,n}(\pi)^2$. By the case assumption, $\zeta > 4V_{p,n}^2 - V_{p,n}^2 = 3V_{p,n}^2 \geq 0$. More importantly, $\zeta > 4V_{h,n}^2 - V_{h,n}^2 = 3V_{h,n}^2$.

$$\mathbb{P}\left(V'_{s,n}(\pi)^2 \geq 4V_{h,n}(\pi)^2 \mid \mathcal{G}_n\right) = \mathbb{P}\left(V'_{s,n}(\pi)^2 - V_{p,n}(\pi)^2 \geq \zeta \mid \mathcal{G}_n\right) \leq \frac{\mathrm{Var}(V'_{s,n}(\pi)^2 \mid \mathcal{G}_n)}{\zeta^2}$$

We bound the conditional variance:

$$\mathrm{Var}(V'_{s,n}(\pi)^2 \mid \mathcal{G}_n) = \mathrm{Var}\left(\frac{1}{n^2} \sum_{i=1}^{n} w'_i(\pi)^2 \mid \mathcal{G}_n\right)$$

$$= \frac{1}{n^4} \sum_{i=1}^{n} \mathrm{Var}(w'_i(\pi)^2 \mid X_i) \quad \text{(by conditional independence)}$$

$$\leq \frac{1}{n^4} \sum_{i=1}^{n} \mathbb{E}[(w'_i(\pi)^2)^2 \mid X_i] = \frac{1}{n^4} \sum_{i=1}^{n} \mathbb{E}[w'_i(\pi)^4 \mid X_i]$$

$$= V_{h,n}(\pi)^4$$

Now we bound the denominator $\zeta^2$:

$$\zeta = 4V_{h,n}(\pi)^2 - V_{p,n}(\pi)^2 > 4V_{h,n}(\pi)^2 - V_{h,n}(\pi)^2 = 3V_{h,n}(\pi)^2$$

$$\zeta^2 > (3V_{h,n}(\pi)^2)^2 = 9V_{h,n}(\pi)^4$$

Substituting the bounds into Chebyshev's inequality:

$$\mathbb{P}(\cdots \mid \mathcal{G}_n) \leq \frac{\mathrm{Var}(V'_{s,n}(\pi)^2 \mid \mathcal{G}_n)}{\zeta^2} \leq \frac{V_{h,n}(\pi)^4}{9V_{h,n}(\pi)^4} = \frac{1}{9}$$

Since $1/9 \leq 1/4$, the bound holds in this case as well.

In both cases, we have shown $\mathbb{P}(\cdots \mid \mathcal{G}_n) \leq 1/4$. By the tower property, taking the expectation over $\mathcal{G}_n$ yields $\mathbb{P}(\dots) \leq 1/4$. This completes the proof of (b). $\qquad\square$

## A.3 Symmetrization for Unbounded Processes

This lemma performs the critical symmetrization step, extending [12, Lemma B.1]. It converts the problem of bounding the deviation of the empirical process from its mean, $\hat{V}_n(\pi) - V_n(\pi)$, into a problem of bounding a self-normalized Rademacher process. This proof is now corrected to use the consistent, weight-based definitions from Lemma A.8.

**Definition A.9.** *We use the notation from Section 3.2 and Appendix A.2.*

1. $\hat{V}_n(\pi) = n^{-1} \sum Y_i(\pi)$ *(IS Estimator)*

2. $V_n(\pi) = \mathbb{E}[\hat{V}'_n(\pi) \mid \mathcal{G}_n]$ *(Conditional Mean)*

3. $\mathbf{w}(\pi) \in \mathbb{R}^n$ *(Vector of original weights $w_i(\pi)$)*

4. $\mathbf{w}'(\pi) \in \mathbb{R}^n$ *(Vector of ghost weights $w'_i(\pi)$)*

5. $V_{s,n}(\pi) = n^{-1}\|\mathbf{w}(\pi)\|_2$ *and* $V'_{s,n}(\pi) = n^{-1}\|\mathbf{w}'(\pi)\|_2$

6. $\mathcal{V}_n(\pi) = \max\{V_{s,n}(\pi), V_{p,n}(\pi), V_{h,n}(\pi), n^{-1/2}\}$ *(The weight-based regularizer)*

**Lemma A.10** (Symmetrization for Unbounded Processes). *For any constant $\xi \geq 4$, the following inequality holds:*

$$\mathbb{P}\left(\sup_{\pi \in \Pi} \frac{|\hat{V}_n(\pi) - V_n(\pi)|}{\mathcal{V}_n(\pi)} \geq \xi\right) \leq 2 \sup_{D_n, D'_n} \mathbb{P}_\epsilon\left(\sup_{\pi \in \Pi} \frac{|\sum_{i=1}^{n} \epsilon_i(Y_i(\pi) - Y'_i(\pi))|}{\|\mathbf{w}(\pi) + \mathbf{w}'(\pi)\|_2} \geq \frac{\xi}{8}\right)$$

*where $\sup_{D_n, D'_n}$ is taken over all possible realizations of the data and ghost data, and $\mathbb{P}_\epsilon$ is the probability measure over $\epsilon$.*

*Proof.* The proof follows the structure of [12, Appendix C.1].

Let $\mathcal{G}_n = \sigma(D_n)$. We define the event of interest:

$$\mathcal{E} := \left\{ \sup_{\pi \in \Pi} \frac{|\hat{V}_n(\pi) - V_n(\pi)|}{\mathcal{V}_n(\pi)} \geq \xi \right\}$$

On this event, let $\pi^\dagger \in \Pi$ be a (measurable) policy that attains this supremum. $\pi^\dagger$ is $\mathcal{G}_n$-measurable. We define the auxiliary events $\mathcal{E}_1, \mathcal{E}_2$ using the ghost sample $D'_n$:

$$\mathcal{E}_1 := \left\{ |\hat{V}'_n(\pi^\dagger) - V_n(\pi^\dagger)| \geq 2V_{p,n}(\pi^\dagger) \right\}$$
$$\mathcal{E}_2 := \left\{ V'_{s,n}(\pi^\dagger)^2 \geq 4 \cdot \max\{V_{p,n}(\pi^\dagger)^2, V_{h,n}(\pi^\dagger)^2\} \right\}.$$

By Lemma A.8 (a) and (b), $\mathbb{P}(\mathcal{E}_1 \mid \mathcal{G}_n) \leq 1/4$ and $\mathbb{P}(\mathcal{E}_2 \mid \mathcal{G}_n) \leq 1/4$. By a union bound, $\mathbb{P}(\mathcal{E}_1^c \cap \mathcal{E}_2^c \mid \mathcal{G}_n) \geq 1/2$.

As shown in the previous (identical) proof of this lemma,

$$\mathbb{P}(\mathcal{E} \mid \mathcal{G}_n) \leq 2 \cdot \mathbb{P}(\mathcal{E} \cap \mathcal{E}_1^c \cap \mathcal{E}_2^c \mid \mathcal{G}_n)$$

On the event $\mathcal{E} \cap \mathcal{E}_1^c \cap \mathcal{E}_2^c$, all three conditions hold.

By $\mathcal{E}$ and $\mathcal{E}_1^c$ (and the triangle inequality):

$$|\hat{V}_n(\pi^\dagger) - \hat{V}'_n(\pi^\dagger)| \geq |\hat{V}_n(\pi^\dagger) - V_n(\pi^\dagger)| - |\hat{V}'_n(\pi^\dagger) - V_n(\pi^\dagger)|$$

$$> \xi \cdot \mathcal{V}_n(\pi^\dagger) - 2V_{p,n}(\pi^\dagger) \geq (\xi - 2)\mathcal{V}_n(\pi^\dagger) \geq \frac{\xi}{2}\mathcal{V}_n(\pi^\dagger) \quad \text{(since } \xi \geq 4)$$

By $\mathcal{E}_2^c$, we have $V'_{s,n}(\pi^\dagger)^2 < 4 \cdot \max\{V_{p,n}(\pi^\dagger)^2, V_{h,n}(\pi^\dagger)^2\}$. This implies $V'_{s,n}(\pi^\dagger) < 2 \cdot \max\{V_{p,n}(\pi^\dagger), V_{h,n}(\pi^\dagger)\}$. By definition, $\mathcal{V}_n(\pi^\dagger) \geq \max\{V_{s,n}(\pi^\dagger), V_{p,n}(\pi^\dagger), V_{h,n}(\pi^\dagger)\}$. Therefore, $\mathcal{V}_n(\pi^\dagger) \geq \max\{V_{s,n}(\pi^\dagger), V'_{s,n}(\pi^\dagger)/2\}$.

Chaining these inequalities and converting to $\ell_2$-norms:

$$|\hat{V}_n(\pi^\dagger) - \hat{V}'_n(\pi^\dagger)| \geq \frac{\xi}{2}\mathcal{V}_n(\pi^\dagger) \quad \text{(from (a))}$$

$$\geq \frac{\xi}{2} \cdot \max\{V_{s,n}(\pi^\dagger), V'_{s,n}(\pi^\dagger)/2\} \quad \text{(from (b))}$$

$$\geq \frac{\xi}{8}(V_{s,n}(\pi^\dagger) + V'_{s,n}(\pi^\dagger)) \quad \text{(since } \max\{a, b/2\} \geq (a+b)/4)$$

$$= \frac{\xi}{8n}(\|\mathbf{w}(\pi^\dagger)\|_2 + \|\mathbf{w}'(\pi^\dagger)\|_2)$$

$$\geq \frac{\xi}{8n}\|\mathbf{w}(\pi^\dagger) + \mathbf{w}'(\pi^\dagger)\|_2 \quad \text{(by triangle inequality, as } w_i \geq 0)$$

Multiplying by $n$, we have shown that on $\mathcal{E} \cap \mathcal{E}_1^c \cap \mathcal{E}_2^c$, the following holds:

$$\left| \sum_{i=1}^n (Y_i(\pi^\dagger) - Y'_i(\pi^\dagger)) \right| \geq \frac{\xi}{8}\|\mathbf{w}(\pi^\dagger) + \mathbf{w}'(\pi^\dagger)\|_2$$

From previous steps, and by taking the supremum inside the probability:

$$\mathbb{P}(\mathcal{E} \mid \mathcal{G}_n) \leq 2 \cdot \mathbb{P}\left( \sup_{\pi \in \Pi} \frac{|\sum_{i=1}^n (Y_i(\pi) - Y'_i(\pi))|}{\|\mathbf{w}(\pi) + \mathbf{w}'(\pi)\|_2} \geq \frac{\xi}{8} \mid \mathcal{G}_n \right)$$

Let $\mathbb{P}_{Z'}$ be the measure over the ghost sample $D'_n$ conditional on $\mathcal{G}_n$. The variables $\Delta_i(\pi) = Y_i(\pi) - Y'_i(\pi)$ are conditionally symmetric about 0. Let $\epsilon = \{\epsilon_i\}_{i=1}^n$ be independent Rademacher variables. By standard symmetrization arguments, $\mathbb{P}_{Z'}(\sup \dots)$ is equal to $\mathbb{E}_{Z'|\mathcal{G}_n}[\mathbb{E}_\epsilon[\mathbf{1}\{\sup \dots\}]]$.

Taking the expectation of $\mathbb{P}(\mathcal{E} \mid \mathcal{G}_n)$ over $\mathcal{G}_n$:

$$\mathbb{P}(\mathcal{E}) \leq 2 \cdot \mathbb{E}_{\mathcal{G}_n, Z'}\left[ \mathbb{P}_\epsilon \left( \sup_{\pi \in \Pi} \frac{|\sum \epsilon_i(Y_i(\pi) - Y'_i(\pi))|}{\|\mathbf{w}(\pi) + \mathbf{w}'(\pi)\|_2} \geq \frac{\xi}{8} \right) \right]$$

This expected probability is upper-bounded by the supremum over all data realizations $D_n, D'_n$:

$$\mathbb{P}(\mathcal{E}) \leq 2 \sup_{D_n, D'_n} \mathbb{P}_\epsilon \left( \sup_{\pi \in \Pi} \frac{|\sum_{i=1}^n \epsilon_i(Y_i(\pi) - Y'_i(\pi))|}{\|\mathbf{w}(\pi) + \mathbf{w}'(\pi)\|_2} \geq \frac{\xi}{8} \right)$$

$\square$

## A.4 Concentration of the Normalized Rademacher Process

This lemma bounds the tail probability of the self-normalized Rademacher process that emerged from the symmetrization in Lemma A.10. This process is the key to our overlap-free argument, as we show that the class of vectors being bounded is, by construction, contained within the $\ell_2$ unit ball, regardless of the magnitude of the IS weights.

**Definition A.11.** *For any given data realization $D_n = \{(X_i, A_i, R_i, \mu_i)\}_{i=1}^n$ and ghost realization $D'_n = \{(X_i, A'_i, R'_i, \mu_i)\}_{i=1}^n$, we define the class of vectors $\mathcal{F}_{D_n,D'_n} \subseteq \mathbb{R}^n$:*

$$\mathcal{F}_{D_n,D'_n} := \left\{ f \in \mathbb{R}^n \mid \exists \pi \in \Pi \ s.t. \ \|\mathbf{w}(\pi) + \mathbf{w}'(\pi)\|_2 > 0, f = \frac{\mathbf{Y}(\pi) - \mathbf{Y}'(\pi)}{\|\mathbf{w}(\pi) + \mathbf{w}'(\pi)\|_2} \right\} \cup \{\mathbf{0}\} \qquad (36)$$

*where $\mathbf{Y}(\pi), \mathbf{Y}'(\pi)$ are the vectors of weighted rewards and $\mathbf{w}(\pi), \mathbf{w}'(\pi)$ are the vectors of IS weights.*

*We also define the target process and the worst-case complexity:*

1. *$S'_n(\mathcal{F}) := \sup_{f \in \mathcal{F}} |\sum_{i=1}^n \epsilon_i f_i|$ (Unnormalized Rademacher process).*

2. *$\mathcal{I}_{sup}(\Pi, n) := \sup_{D_n, D'_n} \left\{ \int_0^2 \sqrt{\log N(\epsilon, \mathcal{F}_{D_n,D'_n}, d_{\ell_2})} d\epsilon \right\}$ (Worst-case Dudley Integral).*

We recall the auxiliary lemmas: Lemma A.4 (Dudley's Integral) and Lemma A.5 (Massart's Concentration).

## A.5 Uniform Concentration for Bounded Process

This lemma bounds "Term (ii)" of our error decomposition: the deviation of the (bounded) conditional value function $g_\pi(x) = \mathbb{E}[Y_i(\pi) \mid X_i]$ from its true expectation $V(\pi)$. This replaces the argument of [12, Lemma B.3], substituting Natarajan dimension with a data-dependent concentration bound based on empirical Rademacher complexity.

Recall the related notations and definitions:

**Definition A.12.** *Recall the dataset $D_n = \{X_i\}_{i=1}^n$. We explicitly define the following functional classes and random variables as functions of $D_n$:*

1. *Conditional Value Class: $\mathcal{F}_g := \{g_\pi : \mathcal{X} \to [0,1] \mid g_\pi(x) = \mathbb{E}[w_\pi(A|x)R \mid X = x], \pi \in \Pi\}$.*

2. *Target Empirical Process: $Z_n(\mathcal{F}_g) := \sup_{f \in \mathcal{F}_g} |(P_n - P)f|$. For concentration analysis, we denote this as $g(D_n) := Z_n(\mathcal{F}_g)$.*

3. *Empirical Rademacher Complexity: $\mathcal{R}_n(\mathcal{F}_g) := \mathbb{E}_\epsilon[\sup_{f \in \mathcal{F}_g} |\frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i)| \mid D_n]$. We denote this as $h(D_n) := \mathcal{R}_n(\mathcal{F}_g)$.*

4. *Empirical Dudley Integral: $\mathcal{I}_n(\mathcal{F}_g) := \int_0^1 \sqrt{\log N(\epsilon, \mathcal{F}_g, L_2(P_n))} d\epsilon$.*

**Lemma A.13** (Uniform Concentration for Bounded Process). *For any $\delta \in (0,1)$, the following inequality holds with probability at least $1 - \delta$ (over the draw of $D_n$):*

$$Z_n(\mathcal{F}_g) \le 2\mathcal{R}_n(\mathcal{F}_g) + 3\sqrt{\frac{\log(4/\delta)}{2n}}$$

*Furthermore, there exists a universal constant $C_2 < \infty$ such that with probability at least $1 - \delta$:*

$$Z_n(\mathcal{F}_g) \le \frac{C_2}{\sqrt{n}} \mathcal{I}_n(\mathcal{F}_g) + 3\sqrt{\frac{\log(4/\delta)}{2n}}$$

*Proof.* The proof proceeds in three steps. First, we concentrate $g(D_n)$ and $h(D_n)$ around their expectations using McDiarmid's inequality. Second, we link them via Symmetrization. Third, we bound $h(D_n)$ using Dudley's chaining argument.

We analyze the sensitivity of $g(D_n)$ and $h(D_n)$ to changing a single data point $X_j$ to $X'_j$. For $g(D_n) = Z_n(\mathcal{F}_g)$:

$$c_j = \sup_{D_n, X'_j} |g(D_n) - g(D_n \text{ with } X_j \to X'_j)| \le \sup_{f \in \mathcal{F}_g} \left| \frac{1}{n}(f(X_j) - f(X'_j)) \right| \le \frac{1}{n}$$

The last inequality holds because $\mathcal{F}_g \subseteq [0,1]$. For $h(D_n) = \mathcal{R}_n(\mathcal{F}_g)$:

$$c_j' = \sup_{D_n, X_j'} |h(D_n) - h(D_n \text{ with } X_j \to X_j')| \leq \mathbb{E}_\epsilon \left[ \sup_{f \in \mathcal{F}_g} \left| \frac{1}{n} \epsilon_j (f(X_j) - f(X_j')) \right| \right] \leq \frac{1}{n}$$

Applying McDiarmid's inequality (Lemma A.2) with $\sum c_i^2 \leq 1/n$, and setting $t_0 = \sqrt{\frac{\log(4/\delta)}{2n}}$, we have with probability at least $1 - \delta/2$:

$$Z_n \leq \mathbb{E}[Z_n] + t_0 \quad \text{and} \quad \mathbb{E}[\mathcal{R}_n] \leq \mathcal{R}_n + t_0$$

By standard symmetrization arguments, $\mathbb{E}[Z_n] \leq 2\mathbb{E}[\mathcal{R}_n]$. Combining this with the high-probability bounds from Step 1 yields the first statement:

$$Z_n \leq \mathbb{E}[Z_n] + t_0 \leq 2\mathbb{E}[\mathcal{R}_n] + t_0 \leq 2(\mathcal{R}_n + t_0) + t_0 = 2\mathcal{R}_n + 3t_0$$

We apply the empirical version of Dudley's inequality [24, Theorem 2].

$$\mathcal{R}_n(\mathcal{F}_g) \leq \inf_{\alpha > 0} \left( 4\alpha + \frac{12}{\sqrt{n}} \int_\alpha^1 \sqrt{\log N(\epsilon, \mathcal{F}_g, L_2(P_n))} d\epsilon \right)$$

Taking $\alpha \to 0$, the first term vanishes, yielding $\mathcal{R}_n(\mathcal{F}_g) \leq \frac{12}{\sqrt{n}} \mathcal{I}_n(\mathcal{F}_g)$. Substituting this into the result from Step 2 proves the second statement. $\qquad \square$

## A.6 Proof of Theorem 3.3

**Definition A.14.** *Let* $D_n = \{X_i\}_{i=1}^n$.

1. $V_n(\pi) = \mathbb{E}[\hat{V}_n(\pi) \mid X_1, \ldots, X_n]$.

2. $\mathcal{V}_n(\pi)$ *is the weight-based regularizer from Definition 3.2.*

3. $\beta_1(\Pi, n, \delta) = C_1 \left( \mathcal{I}_{sup}(\Pi, n) + \sqrt{\log(4/\delta)} \right)$ *(data-independent constant).*

4. $\beta_2(\Pi, n, \delta, D_n) = C_2 \left( \mathcal{I}_n(\mathcal{F}_g) + \sqrt{\log(4/\delta)} \right)$ *(data-dependent random variable).*

5. $\beta(D_n) = 8\beta_1 + \beta_2(D_n)$.

6. $StdErr_n(\pi) = \beta(D_n) \cdot \mathcal{V}_n(\pi)$.

We aim to prove that with probability at least $1 - \delta$, $|\hat{V}_n(\pi) - V(\pi)| \leq StdErr_n(\pi)$ holds uniformly for all $\pi \in \Pi$.

*Proof.* **Step 1: Error Decomposition.** For any $\pi \in \Pi$, we use the triangle inequality:

$$\left| \hat{V}_n(\pi) - V(\pi) \right| \leq \underbrace{\left| \hat{V}_n(\pi) - V_n(\pi) \right|}_{\text{Term (i)}} + \underbrace{|V_n(\pi) - V(\pi)|}_{\text{Term (ii)}}$$

**Step 2: Normalization.** We divide by $\mathcal{V}_n(\pi) > 0$ and take the supremum over $\pi \in \Pi$:

$$\sup_{\pi \in \Pi} \frac{|\hat{V}_n(\pi) - V(\pi)|}{\mathcal{V}_n(\pi)} \leq \sup_{\pi \in \Pi} \frac{|\text{Term (i)}|}{\mathcal{V}_n(\pi)} + \sup_{\pi \in \Pi} \frac{|\text{Term (ii)}|}{\mathcal{V}_n(\pi)}$$

We will bound the two terms on the right-hand side separately, each with a probability budget of $\delta/2$.

**Step 3: Bounding Term (ii) (Bounded Drift).** We want to bound $\sup_{\pi \in \Pi} \frac{|V_n(\pi) - V(\pi)|}{\mathcal{V}_n(\pi)}$. By the definition of $\mathcal{V}_n(\pi)$, we have $\mathcal{V}_n(\pi) \geq n^{-1/2}$ for all $\pi$. Therefore,

$$\sup_{\pi \in \Pi} \frac{|\text{Term (ii)}|}{\mathcal{V}_n(\pi)} \leq \sup_{\pi \in \Pi} \frac{|V_n(\pi) - V(\pi)|}{n^{-1/2}} = Z_n(\mathcal{F}_g) \cdot \sqrt{n}$$

We apply Lemma A.13 with a confidence level of $\delta' = \delta/2$. The lemma states that with probability at least $1 - \delta/2$,

$$Z_n(\mathcal{F}_g) \leq \frac{C_2'}{\sqrt{n}} \mathcal{I}_n(\mathcal{F}_g) + 3\sqrt{\frac{\log(4/\delta')}{2n}} = \frac{C_2'}{\sqrt{n}} \mathcal{I}_n(\mathcal{F}_g) + 3\sqrt{\frac{\log(8/\delta)}{2n}}$$

Let $\mathcal{E}_{ii}$ be this event. On $\mathcal{E}_{ii}$, we have:

$$\sup_{\pi \in \Pi} \frac{|\text{Term (ii)}|}{\mathcal{V}_n(\pi)} \leq \sqrt{n} \cdot Z_n(\mathcal{F}_g)$$

$$\leq \sqrt{n} \left( \frac{C_2'}{\sqrt{n}} \mathcal{I}_n(\mathcal{F}_g) + 3\sqrt{\frac{\log(8/\delta)}{2n}} \right)$$

$$= C_2' \mathcal{I}_n(\mathcal{F}_g) + \sqrt{\frac{9}{2} \log(8/\delta)}$$

By defining the constant $C_2$ in Theorem 3.3 as $C_2 := \max(C_2', \sqrt{\frac{9}{2}})$, this entire expression is upper-bounded by $C_2 \left( \mathcal{I}_n(\mathcal{F}_g) + \sqrt{\log(8/\delta)} \right)$, which is precisely the definition of $\beta_2(D_n)$. Thus, $\mathbb{P}(\mathcal{E}_{ii}^c) \leq \delta/2$.

**Step 4: Bounding Term (i) (Unbounded Fluctuation).** We want to bound $\sup_{\pi \in \Pi} \frac{|\hat{V}_n(\pi) - V_n(\pi)|}{\mathcal{V}_n(\pi)}$. Let $\mathcal{E}_i$ be the event:

$$\mathcal{E}_i := \left\{ \sup_{\pi \in \Pi} \frac{|\hat{V}_n(\pi) - V_n(\pi)|}{\mathcal{V}_n(\pi)} \geq 8\beta_1 \right\}$$

We apply Lemma A.10 with $\xi = 8\beta_1$. This gives:

$$\mathbb{P}(\mathcal{E}_i) \leq 2 \sup_{D_n, D_n'} \mathbb{P}_\epsilon \left( S_n'(\mathcal{F}_{D_n, D_n'}) \geq \frac{8\beta_1}{8} \right) = 2 \sup_{D_n, D_n'} \mathbb{P}_\epsilon \left( S_n'(\mathcal{F}_{D_n, D_n'}) \geq \beta_1 \right)$$

Now we apply Lemma **??** with a confidence level of $\delta_0 = \delta/4$. The lemma states:

$$\sup_{D_n, D_n'} \mathbb{P}_\epsilon \left( S_n'(\mathcal{F}_{D_n, D_n'}) \geq C_1 \left( \mathcal{I}_{\sup}(\Pi, n) + \sqrt{\log(1/\delta_0)} \right) \right) \leq \delta_0$$

By our definition of $\beta_1 = C_1 \left( \mathcal{I}_{\sup} + \sqrt{\log(4/\delta)} \right)$ and our choice of $\delta_0 = \delta/4$, the term inside the probability exactly matches $\beta_1$. Thus, $\mathbb{P}(\mathcal{E}_i) \leq 2 \cdot \delta_0 = 2 \cdot (\delta/4) = \delta/2$.

**Step 5: Union Bound.** Let $\mathcal{E} = \mathcal{E}_i^c \cap \mathcal{E}_{ii}^c$. By a union bound on the complementary events:

$$\mathbb{P}(\mathcal{E}^c) = \mathbb{P}(\mathcal{E}_i \cup \mathcal{E}_{ii}) \leq \mathbb{P}(\mathcal{E}_i) + \mathbb{P}(\mathcal{E}_{ii}) \leq \frac{\delta}{2} + \frac{\delta}{2} = \delta$$

Thus, with probability at least $1 - \delta$, the event $\mathcal{E}$ holds. On the event $\mathcal{E}$, we have for all $\pi \in \Pi$:

$$\frac{|\hat{V}_n(\pi) - V(\pi)|}{\mathcal{V}_n(\pi)} \leq \sup_\pi \frac{|\text{Term (i)}|}{\mathcal{V}_n} + \sup_\pi \frac{|\text{Term (ii)}|}{\mathcal{V}_n}$$

$$\leq (8\beta_1) + (\beta_2(D_n))$$

$$= \beta(D_n)$$

This implies $|\hat{V}_n(\pi) - V(\pi)| \leq \beta(D_n) \cdot \mathcal{V}_n(\pi) = \text{StdErr}_n(\pi)$, proving part (a).

Part (b) follows immediately from Proposition 3.1 given that event (a) holds, with the trivial upper bound of 1 coming from the bounded rewards $R \in [0, 1]$. $\qquad\square$

## A.7 Proof of Corollary 3.5

The proof aims to establish a data-independent convergence rate by showing that, under the uniform overlap condition for the optimal policy (Assumption 3.4), the random components $\beta(D_n)$ and $\mathcal{V}_n(\pi^*)$ in Theorem 3.3 concentrate around well-behaved deterministic values. We employ a union bound over three high-probability events, each allocated a failure probability budget of $\delta/3$.

*Proof.* From Theorem 3.3, we know that with probability at least $1 - \delta/3$, the following event holds:

$$\mathcal{E}_A := \{\mathcal{L}(\hat{\pi}) \leq 2 \cdot \beta(D_n) \cdot \mathcal{V}_n(\pi^*)\} \tag{37}$$

where the total complexity is $\beta(D_n) = 8\beta_1(\Pi, n, \delta/3) + \beta_2(\Pi, n, \delta/3, D_n)$.

**Step 1: Concentration of the Regularizer $\mathcal{V}_n(\pi^*)$.** We show that $\mathcal{V}_n(\pi^*)$ is of order $O_p(n^{-1/2})$. Under Assumption 3.4, the weights $w_i(\pi^*)$ are uniformly bounded almost surely by a constant $C_w < \infty$. This implies that all powers of the weights are also bounded: $w_i(\pi^*)^2 \leq C_w^2$ and $w_i(\pi^*)^4 \leq C_w^4$. We apply Bernstein's Inequality (Lemma A.1) to each empirical term with a failure probability of $\delta_B = \delta/9$:

1. **For $V_{s,n}(\pi^*)$:** Let $Z_i = w_i(\pi^*)^2$. The variables $Z_i$ are i.i.d. and bounded in $[0, C_w^2]$. Their variance is bounded by $\mathbb{E}[Z_i^2] \leq C_w^4$. By Bernstein's inequality, with probability at least $1 - \delta_B$:

$$\sum_{i=1}^{n} Z_i \leq n\mathbb{E}[Z_i] + \sqrt{2nC_w^4 \log(1/\delta_B)} + \frac{1}{3}C_w^2 \log(1/\delta_B)$$

Dividing by $n^2$ and taking the square root:

$$V_{s,n}(\pi^*) = \sqrt{\frac{1}{n^2}\sum Z_i} \leq \sqrt{\frac{C_w^2}{n} + O(n^{-3/2})} = O(n^{-1/2})$$

2. **For $V_{p,n}(\pi^*)$:** Similarly, the conditional expectations $\mathbb{E}[w_i(\pi^*)^2|X_i]$ are bounded by $C_w^2$. The same Bernstein argument yields $V_{p,n}(\pi^*) \leq O(n^{-1/2})$ with high probability.

3. **For $V_{h,n}(\pi^*)$:** The terms $\mathbb{E}[w_i(\pi^*)^4|X_i]$ are bounded by $C_w^4$. With high probability, their empirical average is $O(1)$. Thus, $V_{h,n}(\pi^*) = \frac{1}{n}(\sum \mathbb{E}[w_i^4|X_i])^{1/4} = n^{-3/4}(\frac{1}{n}\sum \mathbb{E}[w_i^4|X_i])^{1/4} \leq O(n^{-3/4})$.

Since $n^{-1/2}$ dominates $n^{-3/4}$ for large $n$, by a union bound over these three events, with probability at least $1 - \delta/3$, event $\mathcal{E}_B$ holds, where there exists a constant $C_V(C_w, \delta)$ such that:

$$\mathcal{E}_B := \left\{ \mathcal{V}_n(\pi^*) \leq \frac{C_V}{\sqrt{n}} \right\} \tag{38}$$

**Step 2: Concentration of the Complexity Term $\beta(D_n)$.** The term $\beta_1$ is structurally data-independent (it depends on the supremum over all possible datasets). The random component is $\beta_2(D_n) = C_2(\mathcal{I}_n(\mathcal{F}_g) + \sqrt{\log(24/\delta)})$, which depends on $D_n$ through the empirical Dudley integral $\mathcal{I}_n(\mathcal{F}_g)$.

We prove that $\mathcal{I}_n(\mathcal{F}_g)$ concentrates around its expectation using McDiarmid's inequality. Let $D_n$ and $D_n'$ be two datasets differing only in the $j$-th element $X_j \to X_j'$. The empirical $L_2(P_n)$ norm is $\|f\|_{L^2(P_n)} = \sqrt{\frac{1}{n}\sum_{i=1}^{n} f(X_i)^2}$. Changing one point from $X_j$ to $X_j'$ changes the squared norm by at most $\frac{1}{n}$ (since $f \in [0,1]$), thus we know $c_i = O(1/n)$. Applying McDiarmid's inequality (Lemma A.2) with failure probability $\delta/3$:

$$\mathbb{P}\left(\beta_2(D_n) > \mathbb{E}[\beta_2(D_n)] + t\right) \leq \exp\left(-\frac{2t^2}{\sum c_i^2}\right) \tag{39}$$

Setting the right-hand side to $\delta/3$ yields a deviation $t = O(\sqrt{\log(1/\delta)/n})$. Thus, with probability at least $1 - \delta/3$, the following event holds:

$$\mathcal{E}_C := \left\{ \beta_2(D_n) \leq \mathbb{E}[\beta_2(D_n)] + O\left(\sqrt{\frac{\log(1/\delta)}{n}}\right) \right\} \tag{40}$$

We define the deterministic complexity constant $\beta_C(\Pi, n, \delta) := 8\beta_1 + \mathbb{E}[\beta_2(D_n)] + O\left(\sqrt{\frac{\log(1/\delta)}{n}}\right)$ to absorb all these data-independent quantities.

By a union bound, with probability at least $1 - \delta$, events $\mathcal{E}_A, \mathcal{E}_B, \mathcal{E}_C$ hold simultaneously. On this joint event:

$$\mathcal{L}(\hat{\pi}) \leq 2 \cdot \beta(D_n) \cdot \mathcal{V}_n(\pi^*) \leq 2 \cdot \beta_C(\Pi, n, \delta) \cdot \frac{C_V}{\sqrt{n}} = O\left(\frac{\beta_C(\Pi, n, \delta)}{\sqrt{n}}\right) \tag{41}$$

$\square$

## A.8 Proof of Theorem 3.6

We first fix the notations and definitions needed in the proof:

**Definition A.15.** *We write:*

1. $d_\mu(\pi_j, \pi_k)^2 := \mathbb{E}_X\left[\int_\mathcal{A} \frac{(\pi_j(a|x) - \pi_k(a|x))^2}{\mu(a|x)} d\lambda_\mathcal{A}(a)\right]$ *(the $\chi^2$-pseudo-metric).*

2. $\langle \pi, \pi'\rangle_{\mu^{-1}} := \mathbb{E}_X\left[\int_\mathcal{A} \frac{\pi(a|x)\pi'(a|x)}{\mu(a|x)} d\lambda_\mathcal{A}(a)\right]$ *(the associated inner product).*

3. $C_w$: *A constant such that for all $\pi \in \Pi$, $\|\pi\|_{\mu^{-1}}^2 = \langle \pi, \pi\rangle_{\mu^{-1}} \leq C_w$.*

4. $M = M(\epsilon) = N_{pack}(\epsilon, \Pi, d_\mu)$. *By definition, there exists a subset* $\Pi_0 = \{\pi_1, \ldots, \pi_M\} \subseteq \Pi$ *such that* $d_\mu(\pi_j, \pi_k) \geq \epsilon$ *for all* $j \neq k$.

We construct a set of $M + 1$ "hard" problem instances $\mathcal{P}_{hard} = \{P_0, P_1, \ldots, P_M\} \subseteq \mathcal{P}(C_w, \sigma_R^2)$. We fix a behavior policy $\mu$ that satisfies the $C_w$ condition. Let $\Delta > 0$ be a perturbation magnitude to be chosen later.

1. $P_0$: $(Q_0, \mu)$, where $Q_0(x, a) = 0$ for all $(x, a)$.

2. $P_j$ (for $j = 1, .., M$): $(Q_j, \mu)$, where $Q_j(x, a) = \Delta \cdot \frac{\pi_j(a|x)}{\mu(a|x)}$.

For this construction, the data $D_n = \{(X_i, A_i, R_i)\}_{i=1}^n$ is drawn as $X_i \sim P_X$, $A_i \sim \mu(\cdot|X_i)$, and $R_i \sim \mathcal{N}(Q_j(X_i, A_i), \sigma_R^2)$.

*Proof.* We follow a standard information-theoretic argument based on Fano's inequality [12, Appendix C.6] to establish a lower bound on the minimax risk.

We bound the KL divergence between $P_j$ and $P_0$ for $j \geq 1$.

$$D_{KL}(P_j^{\otimes n} || P_0^{\otimes n}) = n \cdot D_{KL}(P_j || P_0)$$

The KL divergence between two $n$-sample distributions is $n$ times the KL divergence between the single-sample distributions.

$$D_{KL}(P_j || P_0) = \mathbb{E}_{(X,A) \sim P_0} \left[ D_{KL} \left( \mathcal{N}(Q_j(X, A), \sigma_R^2) \| \mathcal{N}(Q_0(X, A), \sigma_R^2) \right) \right]$$

Using the known formula $D_{KL}(\mathcal{N}(\mu_1, \sigma^2) \| \mathcal{N}(\mu_2, \sigma^2)) = \frac{(\mu_1 - \mu_2)^2}{2\sigma^2}$:

$$D_{KL}(P_j || P_0) = \mathbb{E}_{X \sim P_X, A \sim \mu(\cdot|X)} \left[ \frac{(Q_j(X, A) - 0)^2}{2\sigma_R^2} \right]$$

$$= \frac{\Delta^2}{2\sigma_R^2} \mathbb{E}_{X \sim P_X} \left[ \int_{\mathcal{A}} \mu(a|x) \left( \frac{\pi_j(a|x)}{\mu(a|x)} \right)^2 d\lambda_{\mathcal{A}}(a) \right]$$

$$= \frac{\Delta^2}{2\sigma_R^2} \mathbb{E}_{X \sim P_X} \left[ \int_{\mathcal{A}} \frac{\pi_j(a|x)^2}{\mu(a|x)} d\lambda_{\mathcal{A}}(a) \right] = \frac{\Delta^2}{2\sigma_R^2} \|\pi_j\|_{\mu^{-1}}^2$$

By definition of the problem class $\mathcal{P}(C_w, \sigma_R^2)$, $\|\pi_j\|_{\mu^{-1}}^2 \leq C_w$.

$$\max_{j \geq 1} D_{KL}(P_j^{\otimes n} || P_0^{\otimes n}) \leq \frac{n\Delta^2 C_w}{2\sigma_R^2}$$

The parameter of interest is the optimal policy $\pi_j^*$ for problem $P_j$. The risk is the suboptimality $\mathcal{L}_j(\hat{\pi}) = V_j(\pi_j^*) - V_j(\hat{\pi})$. Under $P_j$, the value of a policy $\pi$ is:

$$V_j(\pi) = \mathbb{E}_X \left[ \int \pi(a|x) Q_j(x, a) da \right] = \Delta \cdot \mathbb{E}_X \left[ \int \frac{\pi(a|x) \pi_j(a|x)}{\mu(a|x)} da \right] = \Delta \langle \pi, \pi_j \rangle_{\mu^{-1}}$$

The value is maximized when $\pi$ is maximally correlated with $\pi_j$ in the $\langle \cdot, \cdot \rangle_{\mu^{-1}}$ inner product. Assuming $\Pi_0 \subseteq \Pi$ and $\Pi$ is convex, the optimal policy $\pi_j^*$ is $\pi_j$ itself, as $V_j(\pi_j) = \Delta \|\pi_j\|_{\mu^{-1}}^2 \geq \Delta \langle \pi, \pi_j \rangle_{\mu^{-1}}$ by Cauchy-Schwarz. We establish the separation between $P_j$ and $P_k$ for $j, k \in \{1, \ldots, M\}, j \neq k$.

$$\mathcal{L}_j(\pi_k) + \mathcal{L}_k(\pi_j) = (V_j(\pi_j) - V_j(\pi_k)) + (V_k(\pi_k) - V_k(\pi_j)) \tag{42}$$

$$= \Delta(\langle \pi_j, \pi_j \rangle_{\mu^{-1}} - \langle \pi_k, \pi_j \rangle_{\mu^{-1}}) + \Delta(\langle \pi_k, \pi_k \rangle_{\mu^{-1}} - \langle \pi_j, \pi_k \rangle_{\mu^{-1}}) \tag{43}$$

$$= \Delta \left( \|\pi_j\|_{\mu^{-1}}^2 - 2\langle \pi_j, \pi_k \rangle_{\mu^{-1}} + \|\pi_k\|_{\mu^{-1}}^2 \right) \tag{44}$$

$$= \Delta \|\pi_j - \pi_k\|_{\mu^{-1}}^2 = \Delta d_\mu(\pi_j, \pi_k)^2 \tag{45}$$

By construction of our $\epsilon$-packing set $\Pi_0$, $d_\mu(\pi_j, \pi_k)^2 \geq \epsilon^2$.

$$\mathcal{L}_j(\pi_k) + \mathcal{L}_k(\pi_j) \geq \Delta \epsilon^2$$

This implies that $\max(\mathcal{L}_j(\pi_k), \mathcal{L}_k(\pi_j)) \geq \frac{1}{2}\Delta \epsilon^2$. For any estimator $\hat{\pi}$, if $\hat{\pi} = \pi_k$ when the true state is $j$, the risk is $\mathcal{L}_j(\pi_k)$. This forms the basis of the Fano risk bound. The minimum risk (separation) between any two hypotheses is $2\epsilon' = \frac{1}{2}\Delta \epsilon^2$.

We apply Lemma A.6 using the risk $\mathbb{E}[\mathcal{L}(\hat{\pi})]$ and the separation $\epsilon' = \frac{1}{4}\Delta\epsilon^2$.

$$\inf_{\hat{\pi}} \sup_{j \in \{0,..,M\}} \mathbb{E}_j[\mathcal{L}_j(\hat{\pi})] \geq \epsilon' \cdot \inf_{\hat{\pi}} \sup_j \mathbb{P}_j(\mathcal{L}_j(\hat{\pi}) \geq \epsilon') \tag{46}$$

$$\geq \frac{1}{4}\Delta\epsilon^2 \left( 1 - \frac{\max_{j \geq 1} D_{KL}(P_j^{\otimes n} || P_0^{\otimes n}) + \log 2}{\log M} \right) \tag{47}$$

$$\geq \frac{1}{4}\Delta\epsilon^2 \left( 1 - \frac{n\Delta^2 C_w/(2\sigma_R^2) + \log 2}{\log M} \right) \tag{48}$$

Let $\alpha = n\Delta^2 C_w/(2\sigma_R^2)$. The bound is $\frac{1}{4}\Delta\epsilon^2(1 - \frac{\alpha + \log 2}{\log M})$.

To make this bound non-vacuous, we require $\alpha < \log M - \log 2$. We select $\Delta$ to balance the terms, assuming $M \geq 4$ (so $\log M \geq 2\log 2$). Let $\alpha = \frac{1}{4}\log M$.

$$\frac{n\Delta^2 C_w}{2\sigma_R^2} = \frac{1}{4}\log M \implies \Delta^2 = \frac{2\sigma_R^2 \log M}{4nC_w} = \frac{\sigma_R^2 \log M}{2nC_w}$$

$$\Delta = \sigma_R \sqrt{\frac{\log M}{2nC_w}}$$

This choice of $\Delta$ is valid as long as $Q_j$ remains in the class (e.g., bounded). Substituting this $\Delta$ into the risk bound:

$$\inf_{\hat{\pi}} \sup_P \mathbb{E}[\mathcal{L}(\hat{\pi})] \geq \frac{1}{4}\left( \sigma_R\sqrt{\frac{\log M}{2nC_w}} \right)\epsilon^2 \left( 1 - \frac{\log M/4 + \log 2}{\log M} \right) \tag{49}$$

$$\geq \frac{1}{4}\left( \sigma_R\sqrt{\frac{\log M}{2nC_w}} \right)\epsilon^2 \left( 1 - \frac{1}{4} - \frac{1}{2} \right) \quad (\text{since } M \geq 4, \log 2 \leq \frac{1}{2}\log M) \tag{50}$$

$$= \frac{1}{16}\sigma_R\sqrt{\frac{\log M}{2nC_w}}\epsilon^2 \tag{51}$$

Let $C_3 = \frac{\sigma_R}{16\sqrt{2}}$.

$$\inf_{\hat{\pi}} \sup_P \mathbb{E}[\mathcal{L}(\hat{\pi})] \geq C_3\epsilon^2 \cdot \sqrt{\frac{\log M(\epsilon)}{n \cdot C_w}}$$

This concludes the proof. $\qquad\square$

# B    Experimental Details

This appendix provides the complete technical specifications required to reproduce the experimental results presented in Section 5. The implementation complies with the theoretical assumptions outlined in Section 3, particularly regarding reward boundedness and the structure of the policy class.

## B.1    Benchmark Specification

All benchmarks share a common foundational structure. The context space is 5-dimensional, with $X \sim U([-1,1]^5)$. The action space is $\mathcal{A} = [-1,1]$. The observed reward is $R_i = r(X_i, A_i) + \epsilon_i$, where $\epsilon_i \sim \mathcal{N}(0, \sigma_{noise}^2)$.

We define the *Truncated Normal* distribution $\mathcal{N}_\mathcal{T}(\mu, \sigma^2, [a, b])$ with probability density function:

$$\phi_\mathcal{T}(x; \mu, \sigma, a, b) = \frac{\frac{1}{\sigma}\phi(\frac{x-\mu}{\sigma})}{\Phi(\frac{b-\mu}{\sigma}) - \Phi(\frac{a-\mu}{\sigma})} \cdot \mathbb{I}(a \leq x \leq b) \tag{52}$$

where $\phi(\cdot)$ and $\Phi(\cdot)$ are the standard normal p.d.f. and c.d.f., respectively.

The specific functional forms for the three benchmarks are:

**Benchmark 1: BiasedBehaviorSharpPeak**  Designed to test learning when the optimal action lies in a low-density tail of the behavior policy.

1. **True Reward:** $r(x, a) = \exp(-50(a - x_1)^2)$.

2. **Behavior Policy:** $\mu(a|x) = \phi_{\mathcal{T}}(a; \mu_b(x), 0.35, -1, 1)$, where the mean $\mu_b(x) = -0.8x_1 - 0.15$ is systematically biased away from the optimal action $a^*(x) = x_1$.

3. **Noise:** $\sigma_{noise} = 0.05$.

**Benchmark 2: SafetyConstrainedReward**  Simulates a safety-critical application where high rewards are adjacent to catastrophic penalties.

1. **True Reward:** $r(x, a) = r_{base}(x, a) + r_{penalty}(a)$, where:
$$r_{base}(x, a) = \exp(-15(a - 0.1(x_1 + x_2))^2)$$
$$r_{penalty}(a) = -3 \cdot \mathbb{I}(a > 0.4)(a - 0.4)^2$$

2. **Behavior Policy:** A risky policy that frequently violates the safety constraint ($a > 0.4$): $\mu(a|x) = \phi_{\mathcal{T}}(a; 0.4, 0.3, -1, 1)$.

3. **Noise:** $\sigma_{noise} = 0.1$.

**Benchmark 3: SparseRewardWithNoise**  Tests the ability to recover a sparse signal amidst high aleatoric noise.

1. **True Reward:** Let $a^*_{sparse}(x) = 0.5x_1 + 0.3x_1$. The reward is non-zero only in a narrow region around $a^*_{sparse}$:
$$r(x, a) = \mathbb{I}(|a - a^*_{sparse}(x)| < 0.15) \cdot \exp(-10|a - a^*_{sparse}(x)|)$$

2. **Behavior Policy:** Uniform random policy, $\mu(a|x) = U([-1, 1])$, providing poor, unguided coverage.

3. **Noise:** $\sigma_{noise} = 0.4$ (significantly higher than other tasks).

## B.2    Data Generation and Preprocessing

For each of the $N = 15$ random seeds, we generate a unique offline dataset $D_n = \{(X_i, A_i, R_i, \mu_i)\}_{i=1}^n$ with $n = 10,000$ samples. To strictly satisfy the theoretical assumption that $R \in [0, 1]$ (crucial for the validity of the self-normalized concentration bounds), we apply min-max normalization to the training rewards:

$$\tilde{R}_i = \frac{R_i - \min_j R_j}{\max_j R_j - \min_j R_j} \tag{53}$$

The behavior density values $\mu_i = \mu(A_i|X_i)$ are recorded during data generation and provided to the learning algorithms.

## B.3    Policy Network Architecture

We parameterize the stochastic policy $\pi_\theta(a|x)$ using a Beta distribution, transformed from its standard support of $[0, 1]$ to the action space $\mathcal{A} = [-1, 1]$. The network is a Multi-Layer Perceptron (MLP) with the following structure:

1. **Input Layer:** 5 units (context dimension).

2. **Hidden Layers:** Two fully connected layers with 64 units each, using ReLU activation.

3. **Output Layer:** 2 units, corresponding to the raw parameters for the Beta distribution.

To ensure valid Beta parameters $\alpha, \beta > 1$ (enforcing a unimodal distribution conducive to optimization), we apply a Softplus activation with a bias:

$$\alpha(x) = \text{Softplus}(o_1(x)) + 1.0, \quad \beta(x) = \text{Softplus}(o_2(x)) + 1.0$$

The policy samples a raw action $a_{raw} \sim \text{Beta}(\alpha(x), \beta(x))$ and applies the affine transformation $a = 2a_{raw} - 1$ to obtain the final action $a \in [-1, 1]$. For evaluation, we use the deterministic mean of this distribution:

$$a_{det}(x) = 2\left(\frac{\alpha(x)}{\alpha(x) + \beta(x)}\right) - 1$$

## B.4 Optimization and Hyperparameters

All algorithms are implemented in PyTorch. Optimization is performed using Adam with a fixed learning rate of $1 \times 10^{-4}$.

The PPL-MM algorithm (Algorithm 1) is configured with $K = 20$ outer MM steps. Within each outer step, we perform $T_{PG} = 150$ inner Policy Gradient steps to maximize the surrogate objective. This yields a total of 3,000 gradient updates, matched by the Naive PG baseline for fair comparison.

To test robustness, we evaluate four hyperparameter variants:

Table 1: Hyperparameter Variants for Robustness Analysis

| Variant Name | IS Clip Threshold ($C_{clip}$) | Denominator Clamp ($\epsilon_\mu$) |
|---|---|---|
| Standard | 50.0 | $1 \times 10^{-6}$ |
| HighClip | 100.0 | $1 \times 10^{-6}$ |
| LowClip | 20.0 | $1 \times 10^{-6}$ |
| HighClamp | 50.0 | $1 \times 10^{-5}$ |

## B.5 Statistical Metrics Derivation and Usage

To ensure the rigorous interpretability of our empirical results, we rely on a complete suite of statistical tools designed for paired experimental designs. This subsection provides a self-contained derivation of these metrics, justifying their selection and detailing their calculation.

Let $\mathcal{D} = \{D_n^{(i)}\}_{i=1}^N$ be the set of $N = 15$ fixed offline datasets used across all experiments. For a given task and hyperparameter variant, let $V_i^{\text{PPL}}$ and $V_i^{\text{Base}}$ denote the true policy values achieved by PPL-MM and the Naive PG baseline on dataset $D_n^{(i)}$, respectively. Our primary random variable of interest is the *paired performance difference*:

$$\Delta_i = V_i^{\text{PPL}} - V_i^{\text{Base}}, \quad i = 1, \ldots, N \tag{54}$$

By design, the $\Delta_i$ are independent and identically distributed (i.i.d.) random variables with unknown true mean $\mu_\Delta$ and variance $\sigma_\Delta^2$. Our one-sided null hypothesis for superiority is $H_0 : \mu_\Delta \leq 0$, against the alternative $H_1 : \mu_\Delta > 0$.

### B.5.1 Paired t-statistic

The paired t-test [22] is the most powerful test for $\mu_\Delta$ under the assumption that the differences $\Delta_i$ are normally distributed, $\Delta_i \sim \mathcal{N}(\mu_\Delta, \sigma_\Delta^2)$. Even if this assumption is slightly violated, with $N = 15$, the Central Limit Theorem ensures that the sample mean $\bar{\Delta}$ is approximately normal.

We first compute the sample mean and sample standard deviation of the differences:

$$\bar{\Delta} = \frac{1}{N} \sum_{i=1}^N \Delta_i, \quad S_\Delta = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (\Delta_i - \bar{\Delta})^2} \tag{55}$$

The standard error of the mean difference is $SE(\bar{\Delta}) = S_\Delta / \sqrt{N}$. The t-statistic is derived as the ratio of the observed signal ($\bar{\Delta}$) to the noise ($SE(\bar{\Delta})$):

$$t = \frac{\bar{\Delta} - 0}{SE(\bar{\Delta})} = \frac{\bar{\Delta}\sqrt{N}}{S_\Delta} \tag{56}$$

Under $H_0$, this statistic follows Student's $t$-distribution with $N - 1 = 14$ degrees of freedom. We compute the one-sided $p$-value as $p = 1 - F_{t,14}(t)$, where $F_{t,14}$ is the c.d.f. of the $t_{14}$ distribution.

### B.5.2 Cohen's $d$ (Effect Size for Paired Samples)

While the $p$-value indicates statistical significance (confidence that $\mu_\Delta > 0$), it relies heavily on the sample size $N$. To quantify the *magnitude* of the improvement in a standardized, scale-free manner, we use Cohen's $d$. For paired designs, the appropriate variant is Cohen's $d$ [5], which standardizes the mean difference by

the standard deviation of the *differences* themselves (rather than the pooled standard deviation of the raw scores). This correctly accounts for the correlation between the paired runs. The estimator is given by:

$$\hat{d} = \frac{\bar{\Delta}}{S_\Delta} \tag{57}$$

Note the direct relationship $t = \hat{d}\sqrt{N}$. A value of $\hat{d} = 1.0$ indicates that the mean improvement is equal to one full standard deviation of the run-to-run variability. In our results, values of $d > 3.0$ indicate an extremely strong effect where the performance distributions of the two algorithms are almost entirely disjoint.

### B.5.3   Benjamini-Hochberg FDR Control

We perform hypothesis tests for $m = 12$ distinct conditions (3 tasks $\times$ 4 variants). Testing each at a significance level $\alpha = 0.05$ would inflate the probability of false positive findings. To address this, we control the False Discovery Rate (FDR), defined as the expected proportion of false rejections among all rejected hypotheses: $\text{FDR} = \mathbb{E}[V/R|R > 0]$, where $V$ is the number of false rejections and $R$ is the total number of rejections.

We employ the Benjamini-Hochberg (BH) procedure [10], which is more powerful than family-wise error rate methods (like Bonferroni) for exploratory analysis. Let $p_{(1)} \leq p_{(2)} \leq \cdots \leq p_{(m)}$ be the ordered $p$-values from the 12 individual t-tests, and let $H_{(1)}, \ldots, H_{(m)}$ be the corresponding null hypotheses. The BH procedure finds the largest index $k$ such that:

$$p_{(k)} \leq \frac{k}{m}\alpha \tag{58}$$

We then reject all null hypotheses $H_{(1)}, \ldots, H_{(k)}$. This guarantees that $\text{FDR} \leq \alpha$ under the assumption of independent or positively dependent test statistics.

## B.6   Detailed Statistical Results

Table 2 presents the complete, granular results of our paired statistical evaluation across all 12 experimental conditions (3 tasks $\times$ 4 variants). For each condition, we report the mean paired difference $\bar{\Delta}$, the 95% confidence interval for the mean difference, the $t$-statistic from the paired $t$-test ($df = 14$), the FDR-adjusted $p$-value, and Cohen's $d$ effect size.

Table 2: **Complete Paired Statistical Results ($N = 15$ seeds).** Mean differences ($\bar{\Delta}$) are calculated as Performance$_{\text{PPL-MM}}$ − Performance$_{\text{Naive PG}}$. Bold values indicate statistical significance after Benjamini-Hochberg FDR correction ($\alpha = 0.05$).

| Benchmark Task | Variant | Mean Diff. ($\bar{\Delta}$) | 95% CI | $t$-stat | $p_{\textbf{FDR}}$ | Cohen's $d$ |
|---|---|---|---|---|---|---|
| BiasedBehavior SharpPeak | Standard | 0.017 | $[-0.021, 0.056]$ | 0.82 | 0.427 | 0.21 |
| | HighClip | 0.047 | $[0.008, 0.087]$ | 2.26 | 0.054 | 0.58 |
| | LowClip | -0.035 | $[-0.072, 0.005]$ | -1.69 | 0.135 | -0.44 |
| | HighClamp | 0.017 | $[-0.021, 0.055]$ | 0.82 | 0.427 | 0.21 |
| SafetyConstrained Reward | Standard | **0.350** | $[0.337, 0.364]$ | 48.53 | $< 10^{-3}$ | 12.53 |
| | HighClip | **0.350** | $[0.337, 0.364]$ | 48.39 | $< 10^{-3}$ | 12.49 |
| | LowClip | **0.353** | $[0.339, 0.366]$ | 48.71 | $< 10^{-3}$ | 12.58 |
| | HighClamp | **0.350** | $[0.337, 0.364]$ | 48.53 | $< 10^{-3}$ | 12.53 |
| SparseReward WithNoise | Standard | **0.428** | $[0.399, 0.457]$ | 27.64 | $< 10^{-3}$ | 7.14 |
| | HighClip | **0.428** | $[0.399, 0.457]$ | 27.64 | $< 10^{-3}$ | 7.14 |
| | LowClip | **0.428** | $[0.398, 0.457]$ | 27.64 | $< 10^{-3}$ | 7.14 |
| | HighClamp | **0.428** | $[0.398, 0.458]$ | 27.64 | $< 10^{-3}$ | 7.14 |