
Learning Parametric Distributions from Samples and Preferences

Marc Jourdan¹ Gizem Yüce¹ Nicolas Flammarion¹

Abstract

Recent advances in language modeling have underscored the role of preference feedback in enhancing model performance. This paper investigates the conditions under which preference feedback improves parameter estimation in classes of continuous parametric distributions. In our framework, the learner observes pairs of samples from an unknown distribution along with their relative preferences depending on the same unknown parameter. We show that preference-based M-estimators achieve a better asymptotic variance than sample-only M-estimators, further improved by deterministic preferences. Leveraging the hard constraints revealed by deterministic preferences, we propose an estimator achieving an estimation error scaling of $\mathcal{O}(1/n)$ —a significant improvement over the $\Theta(1/\sqrt{n})$ rate attainable with samples alone. Next, we establish a lower bound that matches this accelerated rate; up to dimension and problem-dependent constants. While the assumptions underpinning our analysis are restrictive, they are satisfied by notable cases such as Gaussian or Laplace distributions for preferences based on the log-probability reward.

1. Introduction

Recent progress in language modeling has showcased the effectiveness of preference feedback for fine-tuning (Ziegler et al., 2019; Ouyang et al., 2022; Bai et al., 2022; Touvron et al., 2023; Dubey et al., 2024). Preference data—indicating relative quality between outcomes—consistently outperforms approaches using positive examples only like supervised fine-tuning (Iverson et al., 2024). This empirical success suggests that preference feedback introduces new, complementary information beyond the observed data. Understanding how and why preferences provide this advantage requires connecting the preference model to the

data-generating process (Ge et al., 2024).

To understand the role of preference feedback, we focus on a simpler yet illustrative problem: parameter estimation for parametric distributions and preferences. Specifically, the learner observes pairs of samples from an unknown distribution, along with preferences informed by the same parameter. For instance, preferences based on log-probabilities naturally link the preference and probability models, though other formulations are possible (Huang et al., 2024).

For continuous distributions, we uncover a significant statistical learning gap between preference-based and sample-only estimators. This paper primarily investigates this gap, taking the sample-only maximum likelihood estimator (MLE)—optimal among unbiased estimators—as a baseline. The well-established theory of M-estimators (Van der Vaart, 2000) suggests that preference-based M-estimators improve asymptotic variance under certain conditions. Yet, this improvement is modest: when samples are of similar quality, preference feedback approaches a fair coin toss, providing minimal additional information. While reducing asymptotic variance is encouraging, it does not fully explain the substantial performance gains observed empirically in large-scale language models.

For deterministic preferences, we prove a more striking result: preference-based estimators achieve a statistically significant acceleration in parameter estimation. Specifically, we show that the estimation error scales as $\mathcal{O}(1/n)$ instead of the $\mathcal{O}(1/\sqrt{n})$ rate achieved by sample-only estimators. This acceleration is supported by a matching lower bound, up to dimension and problem-dependent constants.

While this acceleration might sound surprising, the $\Theta(1/n)$ rate can already be observed in a special case of sample-only parameter estimation. For instance, consider estimating the location parameter θ of a uniform distribution on $[\theta, \theta + 1]$ based solely on samples (Wainwright, 2019). The minimax rate for estimation error is $\Theta(1/n)$. The optimal estimator achieving the accelerated rate is the minimum of uniform observations whose density is positive at θ . This improved rate arises from the accumulation of random variables having a positive density at a specific point through a minimum (or maximum) operator, in contrast to the slower aggregation inherent to averaging. Similarly, for deterministic preferences with log-likelihood rewards, we observe

¹School of Computer and Communication Sciences, EPFL, Lausanne, Vaud, Switzerland. Correspondence to: Marc Jourdan <marc.jourdan@epfl.ch>.

the true ordering between likelihoods. As it enforces hard constraints—through a minimum operator—on the admissible parameters, our preference-based estimator achieves accelerated convergence.

To illustrate this acceleration, consider the standard normal distribution with preferences based on log-probabilities. Let $n \in \mathbb{N}$ and $[n] := \{1, \dots, n\}$. For each $i \in [n]$, observe samples $(X_i, Y_i) \sim \mathcal{N}(0_2, I_2)$ along with their log-likelihood deterministic preference $Z_i := \text{sign}((Y_i - X_i)S_i)$, where $S_i := (X_i + Y_i)/2$ is their average. The triplet (X_i, Y_i, Z_i) imposes a hard constraint based on S_i on the location of candidate estimators θ that are consistent with this log-likelihood deterministic preference. Specifically, they satisfy $\theta \leq S_i$ if $S_i > 0$, and $\theta \geq S_i$ otherwise. The set of feasible parameters satisfying all constraints is thus $[\max_{i \in [n], S_i < 0} S_i, \min_{i \in [n], S_i > 0} S_i]$. Since the density of $\mathcal{N}(0, 1/2)$ is positive near zero, the length of this interval decreases as $\mathcal{O}(1/n)$ with high probability.

1.1. Contributions

For continuous parametric probability distributions, we study the statistical learning gap between preference-based estimators and sample-only estimators.

- First, we show that preference-based M-estimators achieve a better asymptotic variance than sample-only M-estimators. The variance is further improved for deterministic preference.
- Second, we introduce an estimator satisfying the constraints revealed by the deterministic preferences, and prove an accelerated estimation error rate of $\mathcal{O}(1/n)$. This constitutes a significant improvement over the $\Theta(1/\sqrt{n})$ rate achieved by M-estimators.
- Third, we provide a lower bound of $\Omega(1/n)$, matching our upper bound up to problem-specific constants.

Our results are derived under general assumptions on the distributions and the preferences. While restrictive, they are satisfied by notable cases such as Gaussian or Laplace distributions for preferences based on log-probabilities.

1.2. Related Work

Learning parametric distributions. Parametric estimation is a central approach in statistics, reducing inference about a distribution to the estimation of a finite-dimensional parameter (Lehmann & Casella, 2006; Wasserman, 2013). The maximum likelihood estimator (MLE) is the most fundamental method in this setting. Its asymptotic properties are well studied (Cramér, 1946; Ibragimov & Has’ Minskii, 2013; Van der Vaart, 2000), while non-asymptotic guarantees have been established in Birgé & Massart (1993) and Spokoiny (2012). Lower bounds in parametric estimation rely on techniques such as Le Cam’s two-point

method (LeCam, 1973), Fano’s method (Fano, 1952), and Assouad’s method (Assouad, 1983), and provide fundamental limits on estimation accuracy (Tsybakov, 2009).

Learning parametric value/preference functions. In the tabular setting, learning from pairwise comparisons aligns with the ranking problem. The performance of MLE under the Bradley-Terry model (Bradley & Terry, 1952) and its extensions has been extensively studied (Hunter, 2004; Negahban et al., 2012; Hajek et al., 2014; Rajkumar & Agarwal, 2014; Shah et al., 2016; Shah & Wainwright, 2018; Mao et al., 2018). The continuous setting, where generalization beyond observed preferences is required, has received less attention, except for linear utility functions (Zhu et al., 2023; Ge et al., 2024; Yao et al., 2025). Beyond analyzing the sample complexity of reward learning with MLE under the Bradley-Terry noise model, Zhu et al. (2023) study the performance of policies trained on the learned reward model. They show that while MLE may fail, a pessimistic variant can yield a policy with improved performance. Relaxing the noise assumption, Ge et al. (2024) show that utility parameters remain unidentifiable without strong modeling assumptions, even with noise-free query responses. However, they demonstrate that, in the active learning setting, utility can still be learned, even in the absence of noise. Their results highlight that the sampling distribution of observations must be aligned with the utility function to achieve improved sample complexity. Yao et al. (2025) leverages sparsity in the preference model and establish sharp estimation rates depending on the sparsity level. Finally, related estimation problems have also been studied in the contexts of dueling bandits and reinforcement learning (Faury et al., 2020; Saha et al., 2023).

Fine-tuning with preference data. Large language models often go through a post-training phase focusing mainly on learning from preference feedback (Lambert, 2024), to improve capabilities such as summarization, instruction following, and reasoning. The standard approach, reinforcement learning from human feedback (RLHF) (Ziegler et al., 2019), trains a reward model to align with human preferences and then optimizes the policy using reinforcement learning, typically with PPO (Schulman et al., 2017). RLHF follows three main steps: supervised fine-tuning, reward model training, and policy optimization. Another line of work has explored alternatives to PPO to simplify training. One such method, direct preference optimization (DPO), reformulates the reward function to learn a policy directly from preference data, avoiding an explicit reward model. Other preference optimization objectives have also been proposed (Meng et al., 2024). Finally, while preference data has traditionally been gathered through human annotators, the learning paradigm has recently expanded to include self-play where the model critiques its own generations (Dubey

et al., 2024; Huang et al., 2024).

2. Problem Statement

Parameter estimation. Let $\Theta \subseteq \mathbb{R}^k$ be a set of parameters for a class of continuous probability distributions \mathcal{F} over $\mathcal{X} \subseteq \mathbb{R}^d$. Let $B_\Theta := \max_{\theta \in \Theta} \|\theta\|$ be the bound on Θ for the norm $\|\cdot\|$ specific to \mathcal{F} . Let S_{k-1} be the unit sphere for this norm. Let $p_\theta^{\otimes 2}$ be the distribution of two independent observations of p_θ .

Let θ^* be an unknown parameter to estimate. Our samples are drawn from $p_{\theta^*}^{\otimes 2}$, i.e., $(X, Y) \sim p_{\theta^*}^{\otimes 2}$. We use two archetypal examples satisfying our assumptions. First, the class $\mathcal{F}_{\mathcal{N}, \Sigma}$ of multivariate Gaussian distributions with known covariance Σ , where Θ are the natural parameters with norm $\|\cdot\|_\Sigma$ where $\|x\|_\Sigma := \sqrt{x^\top \Sigma x}$. Second, the class $\mathcal{F}_{\text{Lap}, b}$ of Laplace distributions with known scale b , where Θ are the mean parameters with norm $\|\cdot\|$.

Preference feedback. Let $\ell_\theta : \mathcal{X}^2 \rightarrow \mathbb{R}$ be a parametric preference function. Given a parametric reward function r_θ , a reward-based preference function is defined as $\ell_\theta(x, y) = r_\theta(x) - r_\theta(y)$. As a concrete example for our derivations, we consider preference based on the log-probability reward $r_\theta = \log p_\theta$. Given observations (x, y) , the true preference z of x over y is governed by $\text{sign}(\ell_\theta(x, y)) \in \{\pm 1, 0\}$. In many settings, however, the observed preference Z can be stochastic due to noise or randomness in human feedback.

Conditioned on $(X, Y) \sim p_{\theta^*}^{\otimes 2}$, we denote the p.d.f. of the law of the preference Z by $h(\ell_\theta(X, Y), \cdot)$. On $\mathcal{X}^2 \times \{\pm 1, 0\}$, the p.d.f. of the law of (X, Y, Z) is denoted as $q_{\theta, h}(x, y, z) := p_{\theta^*}^{\otimes 2}(x, y)h(\ell_\theta(x, y), z)$. Under deterministic feedback, the true preferences are observed:

$$h_{\text{det}}(\cdot, z) := \mathbb{1}(z = \text{sign}(\cdot)) . \quad (1)$$

Under stochastic feedback, noisy preferences $z \in \{\pm 1\}$ are observed based on the sigmoid link:

$$h_{\text{sto}}(\cdot, z) := \sigma(z \cdot) \quad \text{with} \quad \sigma(x) := (1 + e^{-x})^{-1} . \quad (2)$$

Informative preferences. A natural question is to see when preference $Z \sim h(\ell_{\theta^*}(X, Y), \cdot)$ helps to estimate θ^* compared to using samples $(X, Y) \sim p_{\theta^*}^{\otimes 2}$ only. Intuitively, given observations with null preference gradient, parameters close to θ^* could have similar preferences. Therefore, those samples are not sufficient to discriminate between them. For that, let $\mathcal{G}_0(\theta^*) = \{(x, y) \in \mathcal{X}^2 \mid |\ell_{\theta^*}(x, y)| > 0\}$ (resp. $\mathcal{G}_1(\theta^*) = \{(x, y) \in \mathcal{G}_0(\theta^*) \mid \|\nabla_{\theta^*} \ell_{\theta^*}(x, y)\| > 0\}$) be the set of pairs with non-zero preference (resp. gradient) function. For observations in $\mathcal{G}_0(\theta^*)^c$, the preference is zero, hence uninformative. For observations in $\mathcal{G}_1(\theta^*)^c$, the preference is locally independent of the parameter. Therefore,

they do not provide gradient information to distinguish θ^* from a neighboring alternative parameter. Only the preferences of samples in $\mathcal{G}_1(\theta^*)$ can provide information on θ^* , hence preference learning is meaningful if these samples are observed, i.e., $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*)) > 0$ for all $\theta^* \in \Theta$.

Negative examples. The above condition is restrictive both on ℓ_θ and p_θ , even when considering $r_\theta = \log p_\theta$. For example, taking p_θ as the uniform distribution over $[0, \theta]$, we have $\mathbb{P}_{p_\theta^{\otimes 2}}(\mathcal{G}_1(\theta)) = 0$.

2.1. Sample-only MLE

In the absence of preference observations, a natural baseline is to estimate θ^* directly from the observations. Given $(X_i, Y_i)_{i \in [n]} \sim p_{\theta^*}^{\otimes 2n}$, the sample-only (SO) MLE is

$$\begin{aligned} \hat{\theta}_n^{\text{SO}} &\in \arg \min_{\theta} L_n^{\text{SO}}(\theta) \quad \text{with} \\ L_n^{\text{SO}}(\theta) &:= - \sum_{i \in [n]} \log p_\theta^{\otimes 2}(X_i, Y_i) . \quad (\text{SO MLE}) \end{aligned}$$

Asymptotic normality. Under enough regularity (Van der Vaart, 2000), SO MLE is asymptotically normal, i.e.,

$$\sqrt{n}(\hat{\theta}_n^{\text{SO}} - \theta^*) \rightsquigarrow_{n \rightarrow +\infty} \mathcal{N}(0_k, \mathcal{I}(p_{\theta^*}^{\otimes 2})^{-1}) ,$$

where $\mathcal{I}(p_\theta) := \mathbb{E}_{p_\theta}[-\nabla_\theta^2 \log p_\theta]$ is the Fisher information matrix of p_θ and \rightsquigarrow denote the convergence in distribution. Let \succeq denote the Loewner order on p.s.d. matrices. By the Cramér-Rao bound (Rao, 1992), SO MLE has optimal asymptotic covariance among the class of unbiased sample-only estimators, i.e., all sample-only unbiased estimator with asymptotic variance V satisfy $V \succeq \mathcal{I}(p_{\theta^*}^{\otimes 2})^{-1}$.

While asymptotic guarantees provide insight into estimator behavior as $n \rightarrow \infty$, they do not capture performance in the relevant regime of moderate sample sizes. Modern statistics gives meaningful non-asymptotic concentration results on empirical estimators, e.g., for high-dimensional statistics (Vershynin, 2018; Wainwright, 2019).

Regularity conditions. The asymptotic statistics literature has devised weak regularity conditions under which asymptotic normality holds. ‘‘Classical conditions’’ assume stronger conditions, e.g., $\theta \mapsto \log p_\theta(x)$ is three times continuously differentiable for every $x \in \mathcal{X}$ and the integral of its third derivative converges uniformly for all θ (Van der Vaart, 2000, Chapter 5.6). Those ‘‘weak’’ or ‘‘classical’’ conditions ensure that integrals and derivatives can be exchanged, and Taylor approximations around θ^* are well controlled. Throughout this paper, we use ‘‘under enough regularity’’ to refer to these regularity conditions on both p_θ and ℓ_θ .

For preferences based on the reward $r_\theta = \log p_\theta$, the regularity of p_θ implies the one of the preference ℓ_θ due to the

properties of the logarithm. Moreover, those regularity conditions are satisfied for numerous well-known distributions such as $\mathcal{F}_{\mathcal{N},\Sigma}$ and $\mathcal{F}_{\text{Lap},b}$. When studying deterministic preferences, we introduce general geometric assumptions on p_θ and ℓ_θ . Since these conditions are inherently more restrictive, our goal is not to identify the weakest possible regularity assumptions under which our derivations hold.

3. Preference-based M-estimator

In this section, we investigate when preference-based estimators can improve upon sample-only estimators. Given preference-labeled observations $\{(X_i, Y_i, Z_i)\}_{i \in [n]}$, we define the stochastic preferences MLE (**SP MLE**) as

$$\hat{\theta}_n^{\text{SP}} \in \arg \min_{\theta} L_n^{\text{SP}}(\theta) \quad \text{with} \\ L_n^{\text{SP}}(\theta) := L_n^{\text{SO}}(\theta) - \sum_{i \in [n]} \log \sigma(Z_i \ell_\theta(X_i, Y_i)). \quad (\text{SP MLE})$$

This objective extends **SO MLE** by adding a preference-based term: a binary classification loss using the logistic function $(-\log \sigma(x))$. When preferences are stochastic, this estimator corresponds to the MLE under a probabilistic preference model, justifying its name. Under sufficient regularity, M-estimators achieve asymptotic normality, so our goal is to obtain lower asymptotic covariance for **SP MLE** than for **SO MLE**. In addition, we want to show that **SP MLE** reaches a lower asymptotic covariance for deterministic preferences than for stochastic preferences.

3.1. Stochastic Preferences

Under stochastic feedback, we are given noisy preference observations $(X_i, Y_i, Z_i)_{i \in [n]} \sim q_{\theta^*, h_{\text{sto}}}^{\otimes n}$, where h_{sto} is defined in Equation (2). **SP MLE** is a specific instance of M-estimator. Under enough regularity (Van der Vaart, 2000, Chapter 5.5), **SP MLE** is asymptotically normal, i.e.,

$$\sqrt{n}(\hat{\theta}_n^{\text{SP}} - \theta^*) \rightsquigarrow_{n \rightarrow \infty} \mathcal{N}(0_k, \mathcal{I}(q_{\theta^*, h_{\text{sto}}})^{-1}),$$

where $\mathcal{I}(q_{\theta, h_{\text{sto}}}) := \mathbb{E}_{q_{\theta, h_{\text{sto}}}}[-\nabla_\theta^2 \log q_{\theta, h_{\text{sto}}}]$ denotes the Fisher information matrix of $q_{\theta, h_{\text{sto}}}$. By the Cramér-Rao bound (Rao, 1992), this variance is optimal among unbiased estimators that rely on stochastic preferences. Lemma 3.1 compares its efficiency to the sample-only MLE.

Lemma 3.1. *Let $\Delta_{\theta^*}^{\text{SP}} := \mathbb{E}_{p_{\theta^*}^{\otimes 2}}[\sigma(\ell_{\theta^*})\sigma(-\ell_{\theta^*})\nabla_{\theta^*}\ell_{\theta^*}\nabla_{\theta^*}\ell_{\theta^*}^\top]$. Then, $\mathcal{I}(q_{\theta^*, h_{\text{sto}}}) = \mathcal{I}(p_{\theta^*}^{\otimes 2}) + \Delta_{\theta^*}^{\text{SP}}$. The p.s.d. matrix $\Delta_{\theta^*}^{\text{SP}}$ is definite if $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(|\langle u, \nabla_{\theta^*}\ell_{\theta^*} \rangle| > 0) > 0$ for all $u \in \mathcal{S}_{k-1}$.*

Lemma 3.1 shows that $\mathcal{I}(q_{\theta^*, h_{\text{sto}}}) \succeq \mathcal{I}(p_{\theta^*}^{\otimes 2})$ and exhibits a condition under which $\hat{\theta}_n^{\text{SP}}$ is asymptotically better than $\hat{\theta}_n^{\text{SO}}$, meaning that incorporating preference data improves asymptotic efficiency. The condition in Lemma 3.1 ensures that $\nabla_{\theta^*}\ell_{\theta^*}$ spans all directions with some probability, making the preference-based estimator asymptotically superior to the sample-only MLE.

For preferences based on the reward $r_\theta = \log p_\theta$, this condition holds for both Laplace and Gaussian distributions: $\Delta_{\theta^*}^{\text{SP}} = \frac{4}{b^2} \Delta_{\text{Lap}(0,1)}^{\text{SP}}$ for $\mathcal{F}_{\text{Lap},b}$ (Appendix G), and $\Delta_{\theta^*}^{\text{SP}} = 2\Sigma^{1/2} \Delta_{\mathcal{N}(0_d, I_d)}^{\text{SP}} \Sigma^{1/2}$ for $\mathcal{F}_{\mathcal{N},\Sigma}$ (Appendix F).

Thus, stochastic preferences can improve parameter estimation compared to sample-only estimators. However, non-asymptotic performance can differ, and in practice, the reduction in asymptotic variance may be small, as we investigate empirically in Section 6. Next, we examine whether M-estimators based on deterministic preferences can further improve upon their stochastic counterparts.

3.2. Deterministic Preferences

We now consider the setting where true preferences are observed, meaning that the preference labels Z_i are deterministic. We observe $(X_i, Y_i, Z_i)_{i \in [n]} \sim q_{\theta^*, h_{\text{det}}}^{\otimes [n]}$, where h_{det} is defined in Equation (1). We use the same M-estimator as in the stochastic setting, $\hat{\theta}_n^{\text{SP}}(\theta) \in \arg \min_{\theta} L_n^{\text{SP}}(\theta)$, but now with deterministic preferences. To distinguish this setting, we introduce the notation SP_{det} for the preference-based estimator under deterministic feedback.

Consistency of SP_{det} . Define the population-level objective: $M(\theta) := \mathbb{E}_{p_{\theta^*}^{\otimes 2}}[\log q_{\theta, h_{\text{sto}}}(X, Y, \text{sign}(\ell_{\theta^*}(X, Y)))]$.

Under enough regularity, $\hat{\theta}_n^{\text{SP}_{\text{det}}}$ converges to a maximizer of $M(\theta)$ (Van der Vaart, 2000, Chapter 5.2). However, unlike in the stochastic setting, θ^* may not be a maximizer of M since standard regularity conditions on p_θ and ℓ_θ are insufficient. A sufficient condition for consistency is

$$\mathbb{E}_{p_{\theta^*}^{\otimes 2}}[\text{sign}(\ell_{\theta^*})\sigma(-|\ell_{\theta^*}|)\nabla_{\theta^*}\ell_{\theta^*}] = 0_k, \quad (3)$$

which holds for $\mathcal{F}_{\mathcal{N},\Sigma}$ (Appendix F) and $\mathcal{F}_{\text{Lap},b}$ (Appendix G) when using reward $r_\theta = \log p_\theta$.

Asymptotic variance of SP_{det} . If Equation (3) holds, then under additional regularity conditions (Van der Vaart, 2000, Chapter 5.3) SP_{det} is asymptotically normal with covariance $V_{\theta^*}^{\text{SP}_{\text{det}}}$ given by the following lemma.

Lemma 3.2. *Let $H_{\theta^*}^{\text{SP}_{\text{det}}} := \mathbb{E}_{p_{\theta^*}^{\otimes 2}}[u_{\theta^*}\nabla_{\theta^*}^2 \ell_{\theta^*}]$, $\Delta_{\theta^*}^{\text{SP}_{\text{det}}} := \mathbb{E}_{p_{\theta^*}^{\otimes 2}}[(2\sigma(|\ell_{\theta^*}|) - 1)\sigma(-|\ell_{\theta^*}|)\nabla_{\theta^*}\ell_{\theta^*}\nabla_{\theta^*}\ell_{\theta^*}^\top]$ and $R_{\theta^*}^{\text{SP}_{\text{det}}} := \mathbb{E}_{p_{\theta^*}^{\otimes 2}}[u_{\theta^*}(M_{\theta^*} + M_{\theta^*}^\top)]$ where $u_{\theta^*} := \text{sign}(\ell_{\theta^*})\sigma(-|\ell_{\theta^*}|)$ and $M_{\theta^*} := -\nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(\nabla_{\theta^*}\ell_{\theta^*})^\top$. Then, we have $V_{\theta^*}^{\text{SP}_{\text{det}}} := V_{1,\theta^*}^{-1}V_{2,\theta^*}V_{1,\theta^*}^{-1}$ where $V_{1,\theta^*} = \mathcal{I}(q_{\theta^*, h_{\text{sto}}}) - H_{\theta^*}^{\text{SP}_{\text{det}}}$ and $V_{2,\theta^*} = \mathcal{I}(q_{\theta^*, h_{\text{sto}}}) - \Delta_{\theta^*}^{\text{SP}_{\text{det}}} - R_{\theta^*}^{\text{SP}_{\text{det}}}$. If $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(|\ell_{\theta^*}\langle u, \nabla_{\theta^*}\ell_{\theta^*} \rangle| > 0) > 0$ for all $u \in \mathcal{S}_{k-1}$, the p.s.d. matrix $\Delta_{\theta^*}^{\text{SP}_{\text{det}}}$ is definite.*

Equation (3) and $V_{\theta^*}^{\text{SP}_{\text{det}}} \prec \mathcal{I}(q_{\theta^*, h_{\text{sto}}})^{-1}$ depend on the geometry of $p_{\theta^*}^{\otimes 2}$ and ℓ_{θ^*} . We verify that these conditions hold for $\mathcal{F}_{\mathcal{N},\Sigma}$ (Appendix F) and $\mathcal{F}_{\text{Lap},b}$ (Appendix G) when

using reward $r_\theta = \log p_\theta$. For Laplace distribution, we have $H_{\theta^*}^{\text{SP det}} = R_{\theta^*}^{\text{SP det}} = 0$ and $\Delta_{\theta^*}^{\text{SP det}} = \frac{4}{b^2} \Delta_{\text{Lap}(0,1)}^{\text{SP det}}$. For Gaussian distributions, we have $H_{\theta^*}^{\text{SP det}} = 0_{d \times d}$, $\Delta_{\theta^*}^{\text{SP det}} = 2\Sigma^{1/2} \Delta_{\mathcal{N}(0_d, I_d)}^{\text{SP det}} \Sigma^{1/2}$ and $R_{\theta^*}^{\text{SP det}} = 2\Sigma^{1/2} R_{\mathcal{N}(0_d, I_d)}^{\text{SP det}} \Sigma^{1/2}$ with $R_{\mathcal{N}(0_d, I_d)}^{\text{SP det}} \succeq 0_{d \times d}$.

Thus, deterministic preferences improve parameter estimation compared to stochastic preferences.

In conclusion, preference-based M-estimators provide asymptotic improvements in estimation efficiency. Next, we explore whether estimators beyond the M-estimation framework can achieve further gains, potentially exceeding the asymptotic normality limitations.

4. Beyond M-estimators

While computationally efficient, the SP_{det} estimator does not fully leverage the constraints imposed by deterministic preferences. Unlike in the stochastic setting, deterministic preferences provide separability: there exist parameters that classify training examples perfectly, including θ^* itself. A key limitation of SP_{det} is that, like standard logistic regression, it minimizes a convex surrogate loss (negative log-likelihood). This approach can lead to misclassification of training examples.¹ This limitation suggests an opportunity to directly minimize the 0-1 loss², potentially achieving faster rates of convergence.

0-1 loss minimization. Given $(X_i, Y_i, Z_i)_{i \in [n]} \sim q_{\theta^*, h_{\text{det}}}^{\otimes [n]}$, we consider the set \mathcal{C}_n of parameters that minimize the empirical 0 – 1 loss, i.e.,

$$\begin{aligned} \mathcal{C}_n &:= \arg \min_{\theta \in \Theta} \sum_{i \in [n]} \mathbb{1}(Z_i \ell_\theta(X_i, Y_i) < 0) \\ &= \{\theta \in \Theta \mid \forall i \in [n], Z_i \ell_\theta(X_i, Y_i) \geq 0\}, \end{aligned} \quad (4)$$

which is non-empty as $\theta^* \in \mathcal{C}_n$. Parameters $\theta \in \mathcal{C}_n$ perfectly classify all training examples. Any estimator $\hat{\theta}_n^{\text{AE}} \in \mathcal{C}_n$ is referred to as an arbitrary estimator (AE).

Alternatively, we constrain MLE to this feasible set, defining the deterministic preferences MLE (**DP MLE**), i.e.,

$$\hat{\theta}_n^{\text{DP}} \in \arg \min \{L_n^{\text{SO}}(\theta) \mid \theta \in \mathcal{C}_n\} \quad (\text{DP MLE})$$

if $\hat{\theta}_n^{\text{SO}} \notin \mathcal{C}_n$, and $\hat{\theta}_n^{\text{DP}} := \hat{\theta}_n^{\text{SO}}$ otherwise. This estimator minimizes the negative log-likelihood of the samples while ensuring perfect preference classification. For Gaussian with $r_\theta = \log p_\theta$, $\hat{\theta}_n^{\text{DP}}$ estimates θ^* better than $\hat{\theta}_n^{\text{SO}}$ for all n , i.e., **DP MLE** dominates **SO MLE** statistically.

¹For binary classification with separable data, logistic regression converges in direction toward a separating hyperplane.

²For non-separable data, the minimization of the 0-1 classification loss can be NP-hard even for the simple class of linear classifiers, e.g., [Feldman et al. \(2012\)](#).

Lemma 4.1. For all $n \in \mathbb{N}$ and almost surely, we have,

$$\text{for } \mathcal{F}_{\mathcal{N}, \Sigma}, \quad \|\hat{\theta}_n^{\text{DP}} - \theta^*\|_\Sigma \leq \|\hat{\theta}_n^{\text{SO}} - \theta^*\|_\Sigma.$$

For stochastic preferences, minimizing the 0 – 1 loss is generally NP-hard, requiring a convex surrogate like the logistic function. However, for deterministic preferences, computing \mathcal{C}_n is more tractable. If $\theta \mapsto \ell_\theta$ is affine, then \mathcal{C}_n is a convex polytope, defined by at most n half-space constraints. For Gaussian-based preferences, i.e., $r_\theta = \log p_\theta$ and $\mathcal{F}_{\mathcal{N}, \Sigma}$, we have $Z_i \ell_\theta(X_i, Y_i) \geq 0$ if and only if $Z_i \langle X_i - Y_i, \theta - \Sigma^{-1}(X_i + Y_i)/2 \rangle \geq 0$.

Consistency of 0 – 1 loss minimization. Define the disagreement probability between θ and θ^* as $m(\theta) := \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{D}(\theta^*, \theta))$ where $\mathcal{D}(\theta^*, \theta) := \{(x, y) \in \mathcal{X}^2 \mid \ell_{\theta^*}(x, y) \ell_\theta(x, y) < 0\}$ is the set of observations where θ and θ^* assign informative yet opposite preferences.

Under enough regularity ([Van der Vaart, 2000](#), Chapter 5.2), $\hat{\theta}_n^{\text{AE}}$ and $\hat{\theta}_n^{\text{DP}}$ converge in $\mathcal{C}(\theta^*) := \{\theta \in \Theta \mid m(\theta) = 0\}$, which is the non-empty set of minimizers of $m(\theta)$ as $m(\theta) \geq 0 = m(\theta^*)$. We note the set $\mathcal{C}(\theta^*)$ contains θ^* , but possibly others. To ensure consistency ($\hat{\theta}_n^{\text{AE}}, \hat{\theta}_n^{\text{DP}} \rightarrow \theta^*$), we impose the following identifiability assumption that guarantees $\mathcal{C}(\theta^*) = \{\theta^*\}$.

Assumption 4.2 (Identifiability). For all $\theta \neq \theta^*$, $m(\theta) > 0$.

When $r_\theta = \log p_\theta$, it holds for both Gaussian ($\mathcal{F}_{\mathcal{N}, \Sigma}$) (Appendix F) and Laplace ($\mathcal{F}_{\text{Lap}, b}$) (Appendix G) cases.

Fast estimation rate. Once consistency is established, the next goal is to analyze the convergence rate of the estimation errors $\|\hat{\theta}_n^{\text{AE}} - \theta^*\|$ and $\|\hat{\theta}_n^{\text{DP}} - \theta^*\|$. Since these are not M-estimators, they are not necessarily limited to the typical parametric rate $\Omega(1/\sqrt{n})$.

Theorem 4.3 states our main result for Laplace and Gaussian distributions when using log-probability rewards, i.e., a high-probability accelerated rate in $\mathcal{O}(1/n)$.

Theorem 4.3. Let $\delta \in (0, 1)$. For $\mathcal{F}_{\text{Lap}, 1}$ and $\mathcal{F}_{\mathcal{N}, 1}$, we have, for all $n \geq \mathcal{O}(\log(1/\delta))$, with probability $1 - \delta$,

$$\forall \hat{\theta}_n \in \mathcal{C}_n, \quad n|\hat{\theta}_n - \theta^*| = \mathcal{O}(\log(1/\delta)).$$

For $\mathcal{F}_{\mathcal{N}, \Sigma}$ with $d > 1$, there exists positive $A_d = d \rightarrow +\infty \mathcal{O}(\sqrt{d})$ such that, for all $n \geq \tilde{\mathcal{O}}(\log(1/\delta))$, with probability $1 - \delta$,

$$\forall \hat{\theta}_n \in \mathcal{C}_n, \quad n\|\hat{\theta}_n - \theta^*\|_\Sigma \leq \mathcal{O}(A_d \log(1/\delta) \log n).$$

Theorem 4.3 is a direct corollary of our main result, showing that $\max_{\theta \in \mathcal{C}_n} \|\theta - \theta^*\| = \mathcal{O}(1/n)$ (see Theorem 4.8 below). It directly guarantees faster convergence rates for both $\hat{\theta}_n^{\text{AE}}$ and $\hat{\theta}_n^{\text{DP}}$. Theorem 4.8 holds under general geometric conditions on p_θ and ℓ_θ that we introduce with intuitions, while sketching the proof in Section 4.1.

Negative examples. Assumption 4.2 is restrictive both on ℓ_θ and p_θ , even when considering $r_\theta = \log p_\theta$. For example, when all the distributions in \mathcal{F} agree on their preferences, $\text{sign}(\ell_\theta(x, y))$ is independent of θ . Therefore, we have $m(\theta) = 0$ for all $\theta \neq \theta^*$, since $\ell_{\theta^*}(x, y)\ell_\theta(x, y) \geq 0$. Such cases include scenarios where $p_\theta(x)$ is a monotonic function, e.g., the exponential distribution and the Pareto distribution with a known location, as well as the Laplace distribution with a known location. This motivates later assumptions on the directionality of $\nabla_{\theta^*}\ell_{\theta^*}$ for observed samples.

Link to iterative human preference alignment. Many human preference alignment methods build on the Bradley-Terry model for preference, based on rewards. Direct alignment algorithms use variants of the log-likelihood to define the implicit reward of a policy (Rafailov et al., 2023). Choosing $\ell_\theta(x, y) = \log p_\theta(x) - \log p_\theta(y)$ coincides with the optimal policy for maximum entropy RL (Swamy et al., 2025). When leveraging offline preference data, the assumption $(X, Y) \sim p_{\theta^*}^{\otimes 2}$ is unrealistic, as ℓ_{θ^*} is collected from a fixed data set of pairs of observations. However, “online” preference data has become a popular paradigm in the training of recent LLMs. Those iterative alignment procedures rely on the preference data from an earlier model (Dubey et al., 2024). At stage N , the model p_{θ_N} is trained based on the preference data for generations by the previous model, i.e., $(X, Y) \sim p_{\theta_{N-1}}^{\otimes 2}$. Under the realizability assumption and without mode collapse, this self-refinement paradigm should converge towards the true model p_{θ^*} . Our setting characterizes the limiting behavior of this iterative process, i.e., preference based on ℓ_{θ^*} for observations from p_{θ^*} . Nonetheless, we do not claim the direct applicability of DP MLE for realistic LLM training.

4.1. Upper Bound on the Estimation Error

We establish a high-probability upper bound on the estimation error $\max_{\theta \in \mathcal{C}_n} \|\theta - \theta^*\|$ in the general case. This requires grasping the geometry of \mathcal{C}_n relative to θ^* .

Linearized feasibility set. Since \mathcal{C}_n is defined by nonlinear preference constraint, analyzing its geometry is challenging, and we thus consider a linearized approximation of it. We define the linearized constraint set as

$$\tilde{\mathcal{C}}_n := \{\theta \in \Theta \mid \forall i \in [n], (X_i, Y_i) \notin \tilde{\mathcal{D}}(\theta^*, \theta)\},$$

where $\tilde{\mathcal{D}}(\theta^*, \theta) := \{(x, y) \in \mathcal{X}^2 \mid \ell_{\theta^*}(x, y)^2 + \ell_{\theta^*}(x, y)\langle \theta - \theta^*, \nabla \ell_{\theta^*}(x, y) \rangle < 0\}$. This set replaces ℓ_θ with its first-order Taylor expansion around θ^* , neglecting higher-order terms. A key assumption is that the true constraints are at least as strong as the linearized ones. This ensures $\mathcal{C}_n \subseteq \tilde{\mathcal{C}}_n$, allowing us to control \mathcal{C}_n via $\tilde{\mathcal{C}}_n$.

Assumption 4.4 (Linearization validity). For all $\theta \neq \theta^*$, $\tilde{\mathcal{D}}(\theta^*, \theta) \subseteq \mathcal{D}(\theta^*, \theta)$.

Directional analysis and informative constraints. To quantify the geometry of $\tilde{\mathcal{C}}_n$ relative to θ^* , we analyze deviations along directions $u \in \mathcal{S}_{k-1}$. Define the set of informative samples along direction u :

$$\mathcal{G}_1(\theta^*, u) := \{(x, y) \mid \ell_{\theta^*}(x, y)\langle u, \nabla_{\theta^*}\ell_{\theta^*}(x, y) \rangle < 0\}.$$

This set contains observations whose preferences give information along the direction u . Assuming that preferences are informative along all directions, we prevent degenerate cases where some directions lack preference information.

Assumption 4.5 (Informative Preferences). For all $u \in \mathcal{S}_{k-1}$, $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*, u)) > 0$.

Deviation bound via minimum informative sample. Define $R_{n,u}$ as the maximal deviation from θ^* within $\tilde{\mathcal{C}}_n$ along the direction u , i.e.,

$$R_{n,u} := \max\{\varepsilon \geq 0 \mid \theta^* + \varepsilon u \in \tilde{\mathcal{C}}_n\}.$$

We define the scaling factor

$$\forall (x, y) \in \mathcal{G}_1(\theta^*, u), V_{\theta^*,u}(x, y) := \frac{\ell_{\theta^*}(x, y)}{-\langle u, \nabla_{\theta^*}\ell_{\theta^*}(x, y) \rangle}.$$

The value $V_{\theta^*,u}(X_i, Y_i)$ quantifies the amount of information in the preference between X_i and Y_i to discriminate θ^* from other parameters on the half-line directed by u . The lower $V_{\theta^*,u}(X_i, Y_i)$ is, the more discriminative is the preference between X_i and Y_i . Since $(x, y) \in \mathcal{G}_1(\theta^*, u) \setminus \tilde{\mathcal{D}}(\theta^*, \theta^* + \varepsilon u)$ if and only if $V_{\theta^*,u}(x, y) \geq \varepsilon$, we obtain

$$R_{n,u} \leq \min_{i \in [n]} \{V_{\theta^*,u}(X_i, Y_i) \mid (X_i, Y_i) \in \mathcal{G}_1(\theta^*, u)\}.$$

Therefore, the maximal deviation $R_{n,u}$ is upper bounded by the minimum of positive random variables. It remains to upper bound the resulting value of this minimum with high probability and conclude provided some regularities hold, e.g., positive density at zero. By analyzing the distribution of $V_{\theta^*,u}$, we derive the following probabilistic bound.

Lemma 4.6. Suppose Assumption 4.5 hold. For all $u \in \mathcal{S}_{k-1}$, with probability $1 - \delta$,

$$R_{n,u} \leq F_{\theta^*,u}^{-1}(\min\{1, \log(1/\delta)/n\}),$$

with $F_{\theta^*,u}(\varepsilon) := \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(V_{\theta^*,u} \in (0, \varepsilon])$ c.d.f. of $V_{\theta^*,u}$.

Since $\max_{\theta \in \mathcal{C}_n} \|\theta - \theta^*\| \leq \max_{u \in \mathcal{S}_{k-1}} R_{n,u}$, Lemma 4.6 shows that the estimation error can be controlled by the behavior of $F_{\theta^*,u}^{-1}$ around zero, where $F_{\theta^*,u}^{-1}(0) = 0$.

Regularity assumption. To control $F_{\theta^*,u}^{-1}$ near zero, the density $F'_{\theta^*,u}$ should be positive near zero, and we assume control on $(F_{\theta^*,u}^{-1})''$.

Assumption 4.7 (Positive density at zero and regularity of inverse c.d.f.). For all $u \in \mathcal{S}_{k-1}$, $F'_{\theta^*,u}(0) \in (0, +\infty)$ and there exists $(x_{\theta^*,u}, M_{\theta^*,u}) \in (0, 1) \times \mathbb{R}_+$ such that $\sup_{x \in [0, x_{\theta^*,u}]} |(F_{\theta^*,u}^{-1})''(x)| \leq M_{\theta^*,u}$.

Using this assumption and $(F_{\theta^*,u}^{-1})'(0) = 1/F'_{\theta^*,u}(0)$, the first-order Taylor expansion with remainder yields

$$\forall x \in [0, x_{\theta^*,u}], |F_{\theta^*,u}^{-1}(x) - x/F'_{\theta^*,u}(0)| \leq M_{\theta^*,u} x^2/2.$$

This argument leads to our main theorem, directly for $k = 1$ and using a covering argument for $k > 1$.

Theorem 4.8. Suppose Assumptions 4.2, 4.4, 4.5 and 4.7 hold. Let $\delta \in (0, 1)$. Let $\gamma > 0$ and $N(\gamma)$ be the γ -covering number of Θ for the norm $\|\cdot\|$. Let $A_{\theta^*}^{-1} = \min_{u \in \mathcal{S}_{k-1}} F'_{\theta^*,u}(0)$, $B_{\theta^*}^{-1} = \min_{u \in \mathcal{S}_{k-1}} x_{\theta^*,u}$ and $C_{\theta^*} = \max_{u \in \mathcal{S}_{k-1}} M_{\theta^*,u}/2$. When $k = 1$, for all $n \geq B_{\theta^*} \log(2/\delta)$,

$$\max_{\theta \in \mathcal{C}_n} \|\theta - \theta^*\| \leq \frac{A_{\theta^*}}{n} \log(2/\delta) + \frac{C_{\theta^*}}{n^2} \log(2/\delta)^2,$$

with probability $1 - \delta$. When $k > 1$, for all $n \geq B_{\theta^*} \log(N(\gamma)/\delta)$, with probability $1 - \delta$,

$$\max_{\theta \in \mathcal{C}_n} \|\theta - \theta^*\| \leq \gamma + \frac{A_{\theta^*}}{n} \log \frac{N(\gamma)}{\delta} + \frac{C_{\theta^*}}{n^2} \log \left(\frac{N(\gamma)}{\delta} \right)^2.$$

When $k > 1$, the choice of the optimal parameter γ depends on the norm $\|\cdot\|$. Since Θ is bounded by B_{Θ} , $N(\gamma)$ is upper bounded by the covering of the ball having diameter B_{Θ} . As an example, let $N_2(\gamma)$ be the ε -covering number of the unit ball in \mathbb{R}^k for the Euclidean norm. Then, it is known that $\log N_2(\varepsilon) \approx \log(\varepsilon^2 k)/\varepsilon^2$ if $\varepsilon \gtrsim 1/\sqrt{k}$ and $\log N_2(\varepsilon) \approx k \log \frac{1}{\varepsilon^2 k}$ if $\varepsilon \lesssim 1/\sqrt{k}$.³ Therefore, optimizing over γ yields an upper bound on $\max_{\theta \in \mathcal{C}_n} \|\theta - \theta^*\|_2$ scaling as $\tilde{\mathcal{O}}(A_{\theta^*} k/n)$ when $n \gtrsim A_{\theta^*} k^{3/2}$, and $\tilde{\mathcal{O}}((A_{\theta^*}/n)^{1/3})$ otherwise, where $\tilde{\mathcal{O}}(\cdot)$ hides logarithmic terms. For large sample size compared to the dimension, i.e., $n \gtrsim A_{\theta^*} k^{3/2}$, we recover a rate of $\tilde{\mathcal{O}}(1/n)$.

In conclusion, we have derived generic assumptions under which the rate of decay of the estimation error of $\hat{\theta}_n^{\text{AE}}$ and $\hat{\theta}_n^{\text{DP}}$ is in $\tilde{\mathcal{O}}(1/n)$. This is a significant improvement compared to the asymptotic normality of the SP_{det} estimator that implies a rate of $\mathcal{O}(1/\sqrt{n})$.

Positive examples. While Assumptions 4.4, 4.5 and 4.7 are restrictive, they hold for $\mathcal{F}_{\mathcal{N},\Sigma}$ (Appendix F) and $\mathcal{F}_{\text{Lap},b}$ (Appendix G) when using reward $r_{\theta} = \log p_{\theta}$. This yields Theorem 4.3. We have $A_{\theta^*} = 2b$, $B_{\theta^*} = 8$ and $C_{\theta^*} = 16b$ for $\mathcal{F}_{\text{Lap},b}$, and $A_{\theta^*} = \frac{\pi(d-1)\Gamma(d/2)}{2\Gamma((d-1)/2)} = +\infty$ $\mathcal{O}(\sqrt{d})$ for $\mathcal{F}_{\mathcal{N},\Sigma}$, hence a rate in $\mathcal{O}(d^{3/2}/n)$ when $n \gg d^2$.

³E.g., using Gilbert-Varshamov for the lower bound (Gilbert, 1952) and Maurey's empirical method for the upper bound.

Extended discussions. In Appendix B, we discuss how to verify or weaken our assumptions (Appendix B.1), the sources of misspecification (Appendix B.2) and other reward models than log-likelihood (Appendix B.3).

5. Lower Bound for Deterministic Feedback

In this section, we show that the rate $\mathcal{O}(1/n)$ is minimax optimal (up to a logarithmic factor) by deriving a matching lower bound. The standard approach to minimax lower bounds in estimation relies on Fano-type inequalities and hypothesis testing reductions. However, due to Assumption 4.2, the Kullback-Leibler divergence and χ^2 distance between q_{θ^*} and q_{θ} are infinite for $\theta \neq \theta^*$, making these tools ineffective. Instead, we use Assouad's Lemma (Tsybakov, 2009), which provides lower bounds via the total variation distance (defined as $\text{TV}(\mathbb{P}, \mathbb{Q}) := \|\mathbb{P} - \mathbb{Q}\|_1$ for distributions \mathbb{P} and \mathbb{Q}). Since TV is not well-behaved for product distributions, we use the squared Hellinger distance, defined as $H^2(\mathbb{P}, \mathbb{Q}) := \frac{1}{2} \|\sqrt{\mathbb{P}} - \sqrt{\mathbb{Q}}\|_2^2$, which satisfies

$$\text{TV}(\mathbb{P}^{\otimes n}, \mathbb{Q}^{\otimes n}) \leq \sqrt{2H^2(\mathbb{P}^{\otimes n}, \mathbb{Q}^{\otimes n})} \leq \sqrt{2nH^2(\mathbb{P}, \mathbb{Q})}.$$

For further analytical convenience, we also employ the Bhattacharyya coefficient, $\text{BC}(\mathbb{P}, \mathbb{Q}) := \|\sqrt{\mathbb{P}\mathbb{Q}}\|_1$, which is related to the Hellinger distance by $H^2(\mathbb{P}, \mathbb{Q}) = 1 - \text{BC}(\mathbb{P}, \mathbb{Q})$. Since $q_{\theta^*} q_{\theta}$ is zero for disagreeing preferences, we define the restricted BC as

$$\widetilde{\text{BC}}(\tilde{\theta}, \theta) := \mathbb{E}_{p_{\tilde{\theta}}^{\otimes 2}} \left[\mathbb{1} \left(\mathcal{D}(\tilde{\theta}, \theta) \cup \mathcal{D}_0(\tilde{\theta}, \theta) \right) \sqrt{p_{\tilde{\theta}}^{\otimes 2}/p_{\theta}^{\otimes 2}} \right],$$

where $\mathcal{D}_0(\tilde{\theta}, \theta) := \mathcal{G}_0(\tilde{\theta})^c \triangle \mathcal{G}_0(\theta)^c$ is the set where the preferences are zero for exactly one parameter.

Lemma 5.1 decomposes the Hellinger distance between two distributions over the preference triplets into the Hellinger distance between sample-only distributions and the disagreement restricted Bhattacharyya coefficient.

Lemma 5.1. $H^2(q_{\tilde{\theta}}, q_{\theta}) = \widetilde{\text{BC}}(\tilde{\theta}, \theta) + H^2(p_{\tilde{\theta}}^{\otimes 2}, p_{\theta}^{\otimes 2})$ for all $\theta, \tilde{\theta} \in \Theta$.

As $H^2(p_{\tilde{\theta}}^{\otimes 2}, p_{\theta}^{\otimes 2}) \leq 2H^2(p_{\tilde{\theta}}, p_{\theta})$, deriving a lower bound requires controlling $H^2(p_{\tilde{\theta}}, p_{\theta})$ and $\widetilde{\text{BC}}(\tilde{\theta}, \theta)$, hence, we impose the following assumption.

Assumption 5.2. There exists positive constants c_1, c_2 independent of k and a dimension and problem-dependent scaling function $\alpha_{\mathcal{F}}(k)$ such that for all $\theta, \tilde{\theta} \in \Theta$,

$$\widetilde{\text{BC}}(\tilde{\theta}, \theta) + H^2(p_{\tilde{\theta}}^{\otimes 2}, p_{\theta}^{\otimes 2}) \leq \frac{c_1}{\alpha_{\mathcal{F}}(k)} \|\theta - \tilde{\theta}\| + c_2 \|\theta - \tilde{\theta}\|^2.$$

Theorem 5.3 bounds the minimax estimation error.

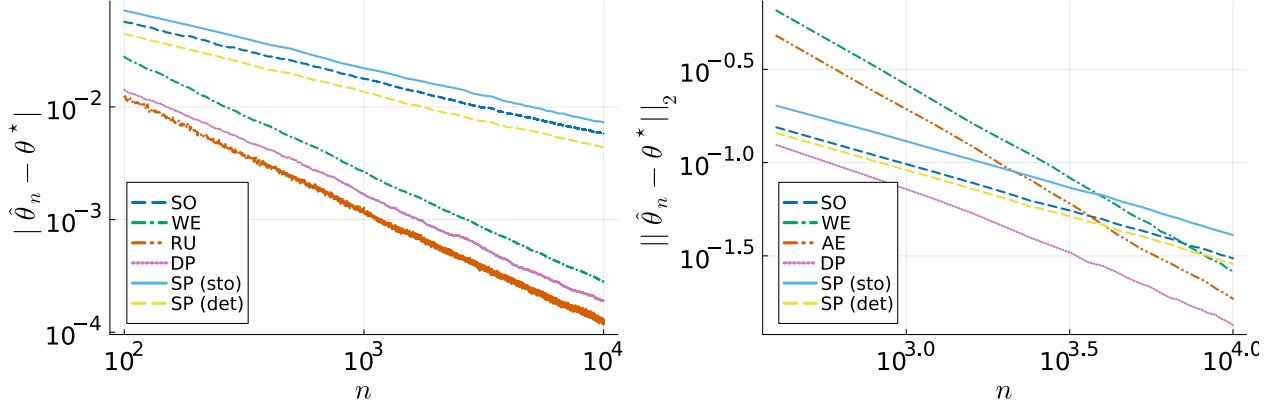


Figure 1. Estimation errors for $\mathcal{N}(\theta^*, I_d)$ where $\theta^* \sim \mathcal{U}([1, 2]^d)$ for (a) $d = 1$ with $N_{\text{runs}} = 10^3$, and (b) $d = 20$ with $N_{\text{runs}} = 10^2$.

Theorem 5.3. Let $R_{\max} := \inf_{\hat{\theta}} \sup_{\theta^* \in \Theta} \mathbb{E}_{q_{\theta^*}} [\|\hat{\theta} - \theta^*\|]$. Suppose Assumption 5.2 holds. Then,

$$R_{\max} \geq \Omega \left(\min \left\{ \frac{\alpha_{\mathcal{F}}(k)\sqrt{k}}{n}, \sqrt{\frac{k}{n}} \right\} \right).$$

This result confirms that the $O(1/n)$ rate is minimax optimal up to logarithmic factors. The scaling $\alpha_{\mathcal{F}}(k)$ comes from $\widetilde{\text{BC}}(\theta^*, \theta)$ (Assumption 5.2), yet it is challenging to link $\alpha_{\mathcal{F}}(k)$ with A_{θ^*} without further assumptions.

Positive examples. While Assumption 5.2 is restrictive, even when using rewards $r_{\theta} = \log p_{\theta}$, it holds for $\mathcal{F}_{\text{Lap}, b}$ (Appendix G), i.e.,

$$\widetilde{\text{BC}}(\theta^*, \theta) = |\theta^* - \theta|/(2b) + \mathcal{O}(|\theta^* - \theta|^2),$$

as well as for $\mathcal{F}_{\mathcal{N}, \Sigma}$ (Appendix F), i.e.,

$$\widetilde{\text{BC}}(\theta^*, \theta) = 2e^{-\|\theta^* - \theta\|_{\Sigma}^2/4} F_{\theta^*, u}(\|\theta^* - \theta\|_{\Sigma}),$$

with $F_{\theta^*, u}(\varepsilon) \leq \varepsilon/A_{\theta^*}$ and $\alpha_{\mathcal{F}}(d) = A_{\theta^*}$.

Dimensionality gap. While the lower bound in Theorem 5.3 scales as $\Omega(\alpha_{\mathcal{F}}(k)\sqrt{k}/n)$ for $n \geq \alpha_{\mathcal{F}}(k)^2$, the upper bound in Theorem 4.8 scales as $\mathcal{O}(A_{\theta^*}k/n)$ for $n \gg A_{\theta^*}k^{3/2}$. Even for the simple case of Gaussian distributions where $A_{\theta^*} = \alpha_{\mathcal{F}}(d)$, there is a dimensionality gap. Closing this gap is an important direction for future work. Improvements might come from a tighter analysis, e.g., both for the upper and lower bounds, or the derivation of better estimators based on deterministic preferences.

6. Experiments

In this section, we compare the empirical performance of the different estimators introduced in this paper. For preferences based on $r_{\theta} = \log p_{\theta}$, we conduct a set of experiments

for Gaussian distributions, and defer to Appendix H.1 for experiments on Laplace and Rayleigh distributions. In particular, we consider a uniformly drawn mean parameter $\theta^* \sim \mathcal{U}([1, 2]^d)$ and the isotropic covariance $\Sigma = I_d$. For sample size $n \in [N_{\max}]$ with $N_{\max} = 10^4$, we compute the estimation errors $\|\hat{\theta}_n - \theta^*\|_2$. We repeat this process for N_{runs} different instances and for various choices of d .

For $\mathcal{F}_{\mathcal{N}, I_d}$ (Appendix F), the M-estimators can be implemented as $\hat{\theta}_n^{\text{SO}} = \frac{1}{2n} \sum_{i \in [n]} (X_i + Y_i)$,

$$\hat{\theta}_n^{\text{SP}} = \arg \min_{\theta} \|\theta - \hat{\theta}_n^{\text{SO}}\|_2^2 - \frac{1}{n} \sum_{i \in [n]} \log \sigma(Z_i \ell_{\theta}(X_i, Y_i)),$$

where $\ell_{\theta}(X_i, Y_i) = \langle X_i - Y_i, \theta - (X_i + Y_i)/2 \rangle$. Then, the estimators based on $\mathcal{C}_n = \{\theta \mid \forall i \in [n], Z_i \ell_{\theta}(X_i, Y_i) \geq 0\}$ are $\hat{\theta}_n^{\text{DP}} = \arg \min_{\theta \in \mathcal{C}_n} \|\theta - \hat{\theta}_n^{\text{SO}}\|_2^2$ and an arbitrary estimator $\hat{\theta}_n^{\text{AE}} \in \mathcal{C}_n$. As \mathcal{C}_n is an interval for $d = 1$, we use the randomized uniform (RU) estimator, i.e., $\hat{\theta}_n^{\text{RU}} \sim \mathcal{U}(\mathcal{C}_n)$. We also consider the worst-case estimator (WE), defined as $\hat{\theta}_n^{\text{WE}} := \arg \max_{\theta \in \mathcal{C}_n} \|\theta - \theta^*\|_1$. While it is not a valid estimator due to its θ^* dependency, it serves as a proxy for the worst estimation error in \mathcal{C}_n .

Dependency on sample size. Figure 1(a) confirms empirically the difference in estimation rate between the M-estimators (SO MLE and SP MLE)—obtaining $\mathcal{O}(1/\sqrt{n})$ —and our estimators based on \mathcal{C}_n —achieving $\mathcal{O}(1/n)$.

However, Figure 1(b) also reveals that the performance of AE and WE deteriorates quickly at small sample sizes when the dimension increases. In contrast, DP MLE consistently outperforms all the other estimators, including SO MLE as theoretically shown in Lemma 4.1.

While SP_{det} outperforms SO MLE, Figure 1 also reveals that SP performs worse than SO MLE for finite sample size. Therefore, only an M-estimator based on deterministic preferences improves on sample-only M-estimators empirically.

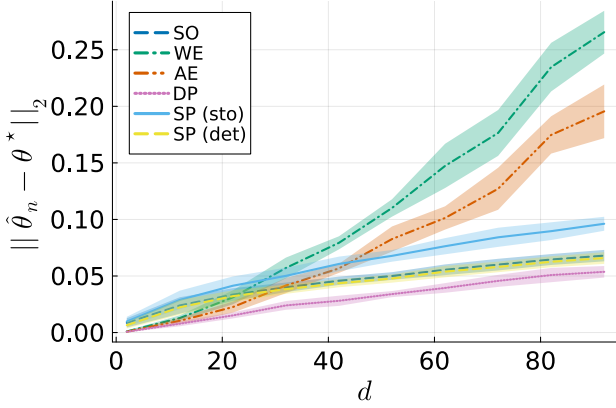


Figure 2. Estimation errors as a function of d with $\mathcal{N}(\theta^*, I_d)$ where $\theta^* \sim \mathcal{U}([1, 2]^d)$, for $n = 10^4$ and $N_{\text{runs}} = 10^3$.

This further highlights the weakness of asymptotic results compared to non-asymptotic guarantees.

While Figure 1(a) suggests that RU and WE perform on par with DP MLE, Figures 1(b) and 2 highlight that DP MLE outperforms WE and AE for larger dimensions, where the gap increases when d is nonnegligible compared to n . We conjecture that RU suffers from the same limitation as AE for larger d . As empirical evidence, we study other estimators that disentangle the effect of RU’s randomness versus its mean behavior, see Appendix H.2.

Dependency on dimension. Figure 2 strengthens the aforementioned empirical observations. For fixed sample size and increasing dimension, DP MLE is the only estimator obtaining the best-of-both world estimation error rate, i.e., $\mathcal{O}(\min\{d^{3/2}/n, \sqrt{d/n}\})$.

Covariance gap. We show that the covariance gap between SP MLE and SO MLE is relative mild: $(\Delta_{\text{Lap}(0,1)}^{\text{SP}}, \Delta_{\text{Lap}(0,1)}^{\text{SP}_{\text{det}}}) \approx (0.16, 0.08)$ and $(\Delta_{\mathcal{N}(0,1)}^{\text{SP}}, \Delta_{\mathcal{N}(0,1)}^{\text{SP}_{\text{det}}}, R_{\mathcal{N}(0,1)}^{\text{SP}_{\text{det}}}) \approx (0.17, 0.08, 0.10)$. Moreover, $\Delta_{\mathcal{N}(0_d, I_d)}^{\text{SP}}$, $\Delta_{\mathcal{N}(0_d, I_d)}^{\text{SP}_{\text{det}}}$ and $R_{\mathcal{N}(0_d, I_d)}^{\text{SP}_{\text{det}}}$ are close to $\alpha_d I_d$ where $\alpha_d > 0$ is decreasing in d (see Figure 3 in Appendix F). In addition to having a small empirical gap for a moderate value of n , the asymptotic gaps between SO MLE and SP MLE are mild.

Supplementary experiments. Following the approach of Tang et al. (2024a), we compare estimators using other convex surrogates of the 0-1 loss (Appendix H.3): they all perform similarly. For the logistic loss, we showcase the “mild” impact of normalization and regularization (Appendix H.4).

7. Perspectives

This work investigates the role of preference feedback in parameter estimation for continuous parametric distributions. We establish conditions under which preference-based estimators outperform sample-only methods. For stochastic preferences, the preference-based MLE achieves a lower asymptotic variance than its sample-only counterpart. For deterministic preferences, we demonstrate that preference-based estimators can significantly accelerate parameter estimation, achieving an improved $\mathcal{O}(1/n)$ convergence rate compared to the $\mathcal{O}(1/\sqrt{n})$ rate of M-estimators. Our lower bound analysis further confirms that this acceleration is min-max optimal up to dimension-dependent constants.

While our results provide a solid theoretical foundation, several open questions remain. A finer analysis of beyond-M-estimators and their constraint set geometry would allow to better quantify the properties of DP MLE, and provide insights for designing improved estimators that better leverage deterministic preferences. Additionally, exploring alternative preference functions beyond the log-probability gap could extend the applicability of our results.

Finally, a key challenge for future work is to quantify the benefits of preference-based estimation for discrete distributions. For distributions with small support, preference feedback may only localize the unknown parameter within a subset of the simplex, leading to diminishing information gains as the sample size increases. However, understanding how preference-based estimators perform in finite-sample settings, particularly in high-dimensional problems, remains an interesting open problem. Addressing these questions could provide further insights into the role of preferences in machine learning and statistical estimation.

Acknowledgements

We thank Jaouad Mourtada for insightful early discussions on the univariate Gaussian case. This work was supported by the Swiss National Science Foundation (grant number 212111) and by an unrestricted gift from Google.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

References

- Assouad, P. Deux remarques sur l’estimation. *C. R. Acad. Sci. Paris Sér. I Math.*, 296(23):1021–1024, 1983. ISSN 0249-6291.
- Azar, M. G., Guo, Z. D., Piot, B., Munos, R., Rowland, M.,

- Valko, M., and Calandriello, D. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pp. 4447–4455. PMLR, 2024.
- Bai, Y., Kadavath, S., Kundu, S., Asbell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Bezanson, J., Edelman, A., Karpinski, S., and Shah, V. B. Julia: A fresh approach to numerical computing. *SIAM Review*, 59(1):65–98, 2017. doi: 10.1137/141000671. URL <https://epubs.siam.org/doi/10.1137/141000671>.
- Birgé, L. and Massart, P. Rates of convergence for minimum contrast estimators. *Probability Theory and Related Fields*, 97:113–150, 1993.
- Bradley, R. A. and Terry, M. E. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- Cramér, H. *Mathematical Methods of Statistics*, volume vol. 9 of *Princeton Mathematical Series*. Princeton University Press, Princeton, NJ, 1946.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Dvoretzky, A., Kiefer, J., and Wolfowitz, J. Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator. *The Annals of Mathematical Statistics*, 27(3):642–669, 1956.
- Fano, R. Class notes for transmission of information. In *Course 6.574*. MIT Cambridge, MA, 1952.
- Faury, L., Abeille, M., Calauzènes, C., and Fercoq, O. Improved optimistic algorithms for logistic bandits. In *International Conference on Machine Learning*, pp. 3052–3060. PMLR, 2020.
- Feldman, V., Guruswami, V., Raghavendra, P., and Wu, Y. Agnostic learning of monomials by halfspaces is hard. *SIAM Journal on Computing*, 41(6):1558–1590, 2012. doi: 10.1137/120865094.
- Ge, L., Juba, B., and Vorobeychik, Y. Learning linear utility functions from pairwise comparison queries. *arXiv preprint arXiv:2405.02612*, 2024.
- Gilbert, E. N. A comparison of signalling alphabets. *The Bell System Technical Journal*, 31(3):504–522, 1952. doi: 10.1002/j.1538-7305.1952.tb01393.x.
- Gorbatovski, A., Shaposhnikov, B., Sinii, V., Malakhov, A., and Gavrilov, D. The differences between direct alignment algorithms are a blur. *arXiv preprint arXiv:2502.01237*, 2025.
- Hajek, B., Oh, S., and Xu, J. Minimax-optimal inference from partial rankings. *Advances in Neural Information Processing Systems*, 27, 2014.
- Hong, J., Lee, N., and Thorne, J. Orpo: Monolithic preference optimization without reference model. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 11170–11189, 2024.
- Huang, A., Block, A., Foster, D. J., Rohatgi, D., Zhang, C., Simchowitz, M., Ash, J. T., and Krishnamurthy, A. Self-improvement in language models: The sharpening mechanism. In *NeurIPS 2024 Workshop on Mathematics of Modern Machine Learning*, 2024.
- Huangfu, Q. and Hall, J. J. Parallelizing the dual revised simplex method. *Mathematical Programming Computation*, 10(1):119–142, 2018.
- Hunter, D. R. Mm algorithms for generalized bradley-terry models. *The annals of statistics*, 32(1):384–406, 2004.
- Ibragimov, I. A. and Has’ Minskii, R. Z. *Statistical estimation: asymptotic theory*, volume 16. Springer Science & Business Media, 2013.
- Iverson, H., Wang, Y., Liu, J., Wu, Z., Pyatkin, V., Lambert, N., Smith, N. A., Choi, Y., and Hajishirzi, H. Unpacking dpo and ppo: Disentangling best practices for learning from preference feedback. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Lambert, N. *Reinforcement Learning from Human Feedback*. Online, 2024. URL <https://rlhfbbook.com>.
- LeCam, L. Convergence of estimates under dimensionality restrictions. *The Annals of Statistics*, pp. 38–53, 1973.
- Lehmann, E. L. and Casella, G. *Theory of point estimation*. Springer Science & Business Media, 2006.
- Lubin, M., Dowson, O., Garcia, J. D., Huchette, J., Legat, B., and Vielma, J. P. Jump 1.0: Recent improvements to a modeling language for mathematical optimization. *Mathematical Programming Computation*, 15(3):581–589, 2023.
- Mao, C., Weed, J., and Rigollet, P. Minimax rates and efficient algorithms for noisy sorting. In *Algorithmic Learning Theory*, pp. 821–847. PMLR, 2018.

- mathlib Community, T. The lean mathematical library. In *Proceedings of the 9th ACM SIGPLAN International Conference on Certified Programs and Proofs*, POPL '20. ACM, January 2020. doi: 10.1145/3372885.3373824. URL <http://dx.doi.org/10.1145/3372885.3373824>.
- Meng, Y., Xia, M., and Chen, D. Simpo: Simple preference optimization with a reference-free reward. *Advances in Neural Information Processing Systems*, 37:124198–124235, 2024.
- Munos, R., Valko, M., Calandriello, D., Azar, M. G., Rowland, M., Guo, Z. D., Tang, Y., Geist, M., Mesnard, T., Fiegel, C., et al. Nash learning from human feedback. In *International Conference on Machine Learning*, pp. 36743–36768. PMLR, 2024.
- Negahban, S., Oh, S., and Shah, D. Iterative ranking from pair-wise comparisons. *Advances in neural information processing systems*, 25, 2012.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., and Finn, C. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36: 53728–53741, 2023.
- Rajkumar, A. and Agarwal, S. A statistical convergence perspective of algorithms for rank aggregation from pairwise data. In *International conference on machine learning*, pp. 118–126. PMLR, 2014.
- Rao, C. R. Information and the accuracy attainable in the estimation of statistical parameters. In *Breakthroughs in Statistics: Foundations and basic theory*, pp. 235–247. Springer, 1992.
- Saha, A., Pacchiano, A., and Lee, J. Dueling rl: Reinforcement learning with trajectory preferences. In *International Conference on Artificial Intelligence and Statistics*, pp. 6263–6289. PMLR, 2023.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Shah, N. B. and Wainwright, M. J. Simple, robust and optimal ranking from pairwise comparisons. *Journal of machine learning research*, 18(199):1–38, 2018.
- Shah, N. B., Balakrishnan, S., Bradley, J., Parekh, A., Ramch, K., Wainwright, M. J., et al. Estimation from pairwise comparisons: Sharp minimax bounds with topology dependence. *Journal of Machine Learning Research*, 17(58):1–47, 2016.
- Spokoyny, V. Parametric estimation. Finite sample theory. *The Annals of Statistics*, 40(6):2877 – 2909, 2012.
- Swamy, G., Dann, C., Kidambi, R., Wu, S., and Agarwal, A. A minimaximalist approach to reinforcement learning from human feedback. In *Forty-first International Conference on Machine Learning*, 2024.
- Swamy, G., Choudhury, S., Sun, W., Wu, Z. S., and Bagneil, J. A. All roads lead to likelihood: The value of reinforcement learning in fine-tuning. *arXiv preprint arXiv:2503.01067*, 2025.
- Tang, Y., Guo, Z. D., Zheng, Z., Calandriello, D., Munos, R., Rowland, M., Richemond, P. H., Valko, M., Avila Pires, B., and Piot, B. Generalized preference optimization: A unified approach to offline alignment. In Salakhutdinov, R., Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett, J., and Berkenkamp, F. (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 47725–47742. PMLR, 21–27 Jul 2024a.
- Tang, Y., Guo, Z. D., Zheng, Z., Calandriello, D., Munos, R., Rowland, M., Richemond, P. H., Valko, M., Pires, B. Á., and Piot, B. Generalized preference optimization: A unified approach to offline alignment. *arXiv preprint arXiv:2402.05749*, 2024b.
- The Sage Developers. *SageMath, the Sage Mathematics Software System (Version 9.7)*, 2022. <https://www.sagemath.org>.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Tsybakov, A. B. Nonparametric estimators. *Introduction to Nonparametric Estimation*, pp. 1–76, 2009.
- Van der Vaart, A. W. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- Vershynin, R. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Wächter, A. and Biegler, L. T. On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming. *Mathematical programming*, 106:25–57, 2006.

- Wainwright, M. J. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*, volume 48. Cambridge University Press, 2019.
- Wang, R., Sun, J., Hua, S., and Fang, Q. Asft: Aligned supervised fine-tuning through absolute likelihood. *arXiv preprint arXiv:2409.10571*, 2024.
- Wasserman, L. *All of statistics: a concise course in statistical inference*. Springer Science & Business Media, 2013.
- Yao, Y., He, L., and Gastpar, M. Leveraging sparsity for sample-efficient preference learning: A theoretical perspective. *arXiv preprint arXiv:2501.18282*, 2025.
- Zhu, B., Jordan, M., and Jiao, J. Principled reinforcement learning with human feedback from pairwise or k-wise comparisons. In *International Conference on Machine Learning*, pp. 43037–43067. PMLR, 2023.
- Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., Christiano, P., and Irving, G. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.

A. Outline

The appendices are organized as follows:

- In Appendix B, we provide detailed discussions on our assumptions, the sources of misspecification and other reward models.
- In Appendix C, we prove the general results presented in Section 3 such as Lemma 3.1.
- In Appendix D, we focus on Section 4 and detail the proofs of Lemma 4.6 and Theorem 4.8
- In Appendix E, we prove the results presented in Section 5.
- In Appendix F, for $\mathcal{F}_{\mathcal{N}, \Sigma}$ and preferences based on $r_\theta = \log p_\theta$, we prove all the assumptions introduced in this paper.
- In Appendix G, for $\mathcal{F}_{\text{Lap}, b}$ and preferences based on $r_\theta = \log p_\theta$, we prove all the assumptions introduced in this paper.
- In Appendix H, we provide supplementary experiments to support our theoretical findings.

B. Extended Discussions

We provide detailed discussions on how to verify or weaken our assumptions (Appendix B.1), the sources of misspecification (Appendix B.2) and other reward models than log-likelihood (Appendix B.3).

B.1. Verifying or Weakening our Assumptions

Since our assumptions are restrictive, it is natural to wonder how they can be verified or weakened.

Verifying our assumptions. Even a closed-form definition of p_θ and ℓ_θ is given, our assumptions are challenging to verify, hence we suggest using a formal verifier (e.g., Lean ([mathlib Community, 2020](#))) or software (e.g., SageMath ([The Sage Developers, 2022](#))). Empirically, they can be confirmed or rejected by sampling from $p_{\theta^*}^{\otimes 2}$. Assumption 4.4 is rejected by exhibiting $(X_i, Y_i) \in \tilde{\mathcal{D}}(\theta^*, \theta) \setminus \mathcal{D}(\theta^*, \theta)$. Assumptions 4.2 and 4.5 are confirmed by finding $(X_i, Y_i) \in \mathcal{D}(\theta^*, \theta)$ and $(X_i, Y_i) \in \mathcal{G}_1(\theta^*, u)$. Those tests’ sampling complexity scales as the inverse event’s probability. Using Dvoretzky–Kiefer–Wolfowitz inequality ([Dvoretzky et al., 1956](#)), $F_{\theta^*, u}$ can be estimated to verify that Assumption 4.7 holds.

Restrictive assumptions. When studying DP MLE only, we conjecture that the “global” Assumptions 4.2 and 4.4 can be weakened to local versions. Using time-uniform concentration results, we can build a sequence of shrinking confidence regions $(\mathcal{R}_n)_n$ around SO MLE that contains θ^* for all time n with high probability. Then, we modify DP MLE to be constrained on $\mathcal{R}_n \cap \mathcal{C}_n$. For n large enough and with high probability, $\mathcal{R}_n \cap \mathcal{C}_n$ will be included in a local neighborhood of θ^* under which the “local” Assumptions 4.2 and 4.4 are satisfied. Given that Assumption 4.4 is based on “ignoring” the reminder term in a first-order Taylor expansion, assuming a local version is a significantly weaker requirement.

B.2. Sources of Misspecification

There are several possible sources of misspecification not taken into account by our current analysis.

Preference model. The Bradley-Terry model that uses reward-based preferences has limited expressivity as it doesn’t allow for intransitive preferences. Even when individuals exhibit transitive preferences, their averaged preferences can be intransitive due to disagreements, see [Munos et al. \(2024\)](#) or [Swamy et al. \(2024\)](#).

Parameter space. When $\theta^* \notin \Theta$, the deterministic preferences might not provide separability within Θ . The definition of DP MLE should be modified to combine the cross-entropy loss and the classification 0-1 loss, i.e.,

$$\hat{\theta}_n^{\text{DP}} \in \arg \min_{\theta \in \Theta} \left\{ L_n^{\text{SO}}(\theta) + \lambda \sum_{i \in [n]} \mathbb{1}(Z_i \ell_\theta(X_i, Y_i) < 0) \right\}, \quad (5)$$

where $\lambda > 0$ is a regularization between those two losses. Equation (5) is reminiscent of single-stage alignment procedures such as ORPO ([Hong et al., 2024](#)) and ASFT ([Wang et al., 2024](#)), see, e.g., [Gorbatovski et al. \(2025\)](#). Without separability, solving Eq. (5) can be NP-hard. Under sufficient regularity, $\hat{\theta}_n^{\text{DP}}$ converges to $\theta_0 \in \arg \min_{\theta \in \Theta} \{\text{KL}(\theta^*, \theta) + \lambda m(\theta)\}$ where

$m(\theta) = \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{D}(\theta^*, \theta))$ and $\theta_0 \neq \theta^*$. As $\theta \mapsto m(\theta)$ can be non-convex, computing θ_0 might be challenging. Deriving a tractable ELBO method for this optimization is an interesting direction to obtain tractable and robust estimators. As θ_0 lies in the boundary of Θ , we should control the maximal deviation with respect to θ_0 for directions that point towards the interior of Θ to prove an accelerated rate. While parts of our analysis could be used, we believe that finer technical arguments are required.

Parametric model. The true distribution p^* of the observations might not even be a member of our class of distributions \mathcal{F} , i.e., $p^* \notin \mathcal{F}$. These situations occur when \mathcal{F} doesn't contain the true structure, e.g., other parametric or non-parametric class of distributions. Then, $L_n^{\text{SO}}(\theta)$ can be interpreted as a quasi-log-likelihood term. Let us denote by SO quasi-MLE the estimator based on **SO MLE** for this quasi-log-likelihood. Under sufficient regularity, SO quasi-MLE converges towards $\theta_0 \in \arg \min_{\theta \in \Theta} \text{KL}(p^*, p_\theta)$ where $p^* \neq p_{\theta_0} \in \mathcal{F}$. Without the separability from well-specified deterministic preference, we define DP quasi-MLE as in Eq. (5). Under sufficient regularity, DP quasi-MLE converges towards the minimizer of a similar optimization problem combining the KL term and a misspecified equivalent of $m(\theta)$.

B.3. Reward models

Except for Theorem 4.3, all the derivations in Section 4 hold for general (hence reward-based) preference models provided Assumptions 4.2, 4.4, 4.5 and 4.7 hold. Characterizing the expressivity of parametric rewards satisfying those assumptions is interesting, yet challenging. We provide two positive and one negative examples.

Positive: monotonic reward. Suppose that $\tilde{\ell}_\theta(x, y) = f(p_\theta(x)) - f(p_{\theta^*}(x))$ where f is increasing on $[0, 1]$. Since $\text{sign}(\tilde{\ell}_\theta) = \text{sign}(\ell_\theta)$, hence the parameters with zero classification loss and our estimators are the same. Therefore, our results hold for this class of rewards when our assumptions hold for the log-likelihood reward. When f is decreasing, the preferences are “reversed”, and similar arguments can be made. This example includes (1) normalization by a multiplicative constant (e.g., temperature β) and (2) the odds-ratio reward-based preference based on $f(x) = \log(x/(1-x))$ used by ORPO in Hong et al. (2024).

Positive: margin with Gaussian. Suppose that $\tilde{\ell}_\theta = \ell_\theta + c$ where c is a constant and ℓ_θ is the Gaussian log-likelihood preference. By extending our computations from Appendix F, Assumptions 4.2, 4.4, 4.5 and 4.7 hold with c -dependent positive constants. Margins are used by SimPO from Meng et al. (2024) and IPO from Azar et al. (2024).

Negative: reference model with Gaussian. Suppose that $\tilde{\ell}_\theta = \ell_\theta - \ell_{\theta_0}$ where θ_0 is known and ℓ_θ is the Gaussian log-likelihood preference. Since $\tilde{\ell}_\theta(x, y) = \langle x - y, \theta - \theta_0 \rangle$ and $\nabla_\theta \tilde{\ell}_\theta(x, y) = x - y$, Assumption 4.5 is violated for $u = \theta^* - \theta_0$, i.e., $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*, \theta^* - \theta_0)) = 0$. Not all direct alignment algorithms rely on a reference model, see, e.g., SimPO and ORPO.

C. Proofs of Section 3

C.1. Stochastic Preferences

Under enough regularity, by swapping the integration and the differentiation operators, we can show that

$$\mathbb{E}_{q_{\theta, h_{\text{sto}}}}[\nabla_\theta \log q_{\theta, h_{\text{sto}}}] = \nabla_\theta 1 = 0_k \quad \text{and} \quad \mathbb{E}_{q_{\theta, h_{\text{sto}}}}[-\nabla_\theta^2 \log q_{\theta, h_{\text{sto}}}] = \mathbb{E}_{q_{\theta, h_{\text{sto}}}}[\nabla_\theta \log q_{\theta, h_{\text{sto}}} \nabla_\theta \log q_{\theta, h_{\text{sto}}}^\top].$$

Below, we detail the proof of Lemma 3.1.

Proof. Direct computation yields that

$$\begin{aligned} \nabla_{\theta^*} \log q_{\theta^*}(x, y, z) &= \nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(x, y) + z\sigma(-z\ell_{\theta^*}(x, y))\nabla_{\theta^*}\ell_{\theta^*}(x, y), \\ \nabla_{\theta^*}^2 \log q_{\theta^*}(x, y, z) &= \nabla_{\theta^*}^2 \log p_{\theta^*}^{\otimes 2}(x, y) - \sigma(\ell_{\theta^*}(x, y))\sigma(-\ell_{\theta^*}(x, y))\nabla_{\theta^*}\ell_{\theta^*}(x, y)\nabla_{\theta^*}\ell_{\theta^*}(x, y)^\top \\ &\quad + z\sigma(-z\ell_{\theta^*}(x, y))\nabla_{\theta^*}^2\ell_{\theta^*}(x, y). \end{aligned}$$

where we used that $g'(x) = \sigma(-x)$ and $g''(x) = -\sigma'(-x) = -\sigma(x)\sigma(-x)$ with $g(x) = \log \sigma(x)$. By definition of h_{sto} ,

$$\begin{aligned} \mathbb{E}_{Z|(X, Y)}[Z\sigma(-Z\ell_{\theta^*}(X, Y))\nabla_{\theta^*}^2\ell_{\theta^*}(X, Y)] \\ = (\sigma(-\ell_{\theta^*}(X, Y))\sigma(\ell_{\theta^*}(X, Y)) - \sigma(\ell_{\theta^*}(X, Y))\sigma(-\ell_{\theta^*}(X, Y)))\nabla_{\theta^*}^2\ell_{\theta^*}(X, Y) = 0_{d \times d}. \end{aligned}$$

Therefore, we have $\mathcal{I}(q_{\theta^*, h_{\text{sto}}}) = \mathcal{I}(p_{\theta^*}^{\otimes 2}) + \Delta_{\theta^*}^{\text{SP}}$ with $\Delta_{\theta^*}^{\text{SP}} = \mathbb{E}_{p_{\theta^*}^{\otimes 2}}[\sigma(\ell_{\theta})\sigma(-\ell_{\theta})\nabla_{\theta}\ell_{\theta}(\nabla_{\theta}\ell_{\theta})^{\top}]$. For all $x \in \mathbb{R}^k$, we have

$$x^{\top}\Delta_{\theta^*}^{\text{SP}}x = \|x\|^2\mathbb{E}_{p_{\theta^*}^{\otimes 2}}[\sigma(\ell_{\theta})\sigma(-\ell_{\theta})\langle x/\|x\|, \nabla_{\theta}\ell_{\theta} \rangle^2] \geq 0.$$

It is direct to see that this inequality is strict except if $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\langle x/\|x\|, \nabla_{\theta}\ell_{\theta} \rangle^2 = 0) = 1$. Therefore, $\Delta_{\theta^*}^{\text{SP}}$ is a positive definite matrix if $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(|\langle u, \nabla_{\theta}\ell_{\theta^*} \rangle| > 0) > 0$ for all $u \in \mathcal{S}_{k-1}$. Note that this condition is implied by Assumption 4.5. \square

C.2. Deterministic Preferences

Consistency of SP. Let $M(\theta) := \mathbb{E}_{p_{\theta^*}^{\otimes 2}}[\log q_{\theta, h_{\text{sto}}}(X, Y, \text{sign}(\ell_{\theta^*}(X, Y)))]$. Under enough regularity, we obtain

$$\mathbb{E}_{p_{\theta^*}^{\otimes 2}}[\nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}] = 0_k \quad \text{and} \quad \nabla_{\theta^*} M(\theta^*) = \mathbb{E}_{p_{\theta^*}^{\otimes 2}}[\text{sign}(\ell_{\theta^*}(X, Y))\sigma(-|\ell_{\theta^*}(X, Y)|)\nabla_{\theta^*}\ell_{\theta^*}(X, Y)].$$

Therefore, θ^* is the unique maximizer of $M(\theta)$ if $\nabla_{\theta^*} M(\theta^*) = 0_k$, i.e., if (3) holds true.

Asymptotic normality of SP. Provided (3), under enough regularity, the theory of M-estimator yields that

$$\sqrt{n}(\hat{\theta}_n^{\text{SPdet}} - \theta^*) \rightsquigarrow_{n \rightarrow +\infty} \mathcal{N}(0_k, V_{1, \theta^*}^{-1} V_{2, \theta^*} V_{1, \theta^*}^{-1}),$$

where

$$\begin{aligned} V_{1, \theta^*} &= \mathbb{E}_{(X, Y) \sim p_{\theta^*}^{\otimes 2}} [-\nabla_{\theta^*}^2 \log q_{\theta^*, h_{\text{sto}}}(X, Y, \text{sign}(\ell_{\theta^*}(X, Y)))] , \\ V_{2, \theta^*} &= \mathbb{E}_{(X, Y) \sim p_{\theta^*}^{\otimes 2}} [\nabla_{\theta^*} \log q_{\theta^*, h_{\text{sto}}}(X, Y, \text{sign}(\ell_{\theta^*}(X, Y))) \nabla_{\theta^*} \log q_{\theta^*, h_{\text{sto}}}(X, Y, \text{sign}(\ell_{\theta^*}(X, Y)))^{\top}] . \end{aligned}$$

Below we detail the proof of Lemma 3.2.

Proof. Combining $z = \text{sign}(\ell_{\theta^*}(x, y))$ with the same manipulation as above yields

$$\begin{aligned} \nabla_{\theta^*} \log q_{\theta^*, h_{\text{sto}}}(x, y, z) &= \nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(x, y) + \text{sign}(\ell_{\theta^*}(x, y))\sigma(-|\ell_{\theta^*}(x, y)|)\nabla_{\theta^*}\ell_{\theta^*}(x, y) , \\ \nabla_{\theta^*} \log q_{\theta^*, h_{\text{sto}}}(x, y, z) \nabla_{\theta^*} \log q_{\theta^*, h_{\text{sto}}}(x, y, z)^{\top} &= \nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(x, y) \nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(x, y)^{\top} \\ &\quad + \sigma(-|\ell_{\theta^*}(x, y)|)^2 \nabla_{\theta^*}\ell_{\theta^*}(x, y) \nabla_{\theta^*}\ell_{\theta^*}(x, y)^{\top} \\ &\quad + \text{sign}(\ell_{\theta^*}(x, y))\sigma(-|\ell_{\theta^*}(x, y)|) (\nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(x, y) \nabla_{\theta^*}\ell_{\theta^*}(x, y)^{\top} + \nabla_{\theta^*}\ell_{\theta^*}(x, y) \nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(x, y)^{\top}) \\ \nabla_{\theta^*}^2 \log q_{\theta^*, h_{\text{sto}}}(x, y, z) &= \nabla_{\theta^*}^2 \log p_{\theta^*}^{\otimes 2}(x, y) - \sigma(\ell_{\theta^*}(x, y))\sigma(-\ell_{\theta^*}(x, y))\nabla_{\theta^*}\ell_{\theta^*}(x, y) \nabla_{\theta^*}\ell_{\theta^*}(x, y)^{\top} \\ &\quad + \text{sign}(\ell_{\theta^*}(x, y))\sigma(-|\ell_{\theta^*}(x, y)|)\nabla_{\theta^*}^2 \ell_{\theta^*}(x, y) . \end{aligned}$$

where we used that $g'(x) = \sigma(-x)$ and $g''(x) = -\sigma'(-x) = -\sigma(x)\sigma(-x)$ with $g(x) = \log \sigma(x)$. Using that $\sigma(-|\ell_{\theta^*}(x, y)|)^2 = \sigma(-|\ell_{\theta^*}(x, y)|) - \sigma(-\ell_{\theta^*}(x, y))\sigma(\ell_{\theta^*}(x, y))$, we have

$$V_{1, \theta^*} = \mathcal{I}(p_{\theta^*}^{\otimes 2}) + \Delta_{\theta^*}^{\text{SP}} - H_{\theta^*}^{\text{SPdet}} \quad \text{and} \quad V_{2, \theta^*} = \mathcal{I}(p_{\theta^*}^{\otimes 2}) + M_{2, \theta^*} - \Delta_{\theta^*}^{\text{SP}} - R_{\theta^*}^{\text{SPdet}}$$

where $\Delta_{\theta^*}^{\text{SP}} = \mathbb{E}_{p_{\theta^*}^{\otimes 2}}[\sigma(\ell_{\theta})\sigma(-\ell_{\theta})\nabla_{\theta}\ell_{\theta}(\nabla_{\theta}\ell_{\theta})^{\top}]$ as in Lemma 3.1, and we define

$$\begin{aligned} H_{\theta^*}^{\text{SPdet}} &= \mathbb{E}_{p_{\theta^*}^{\otimes 2}} [\text{sign}(\ell_{\theta^*})\sigma(-|\ell_{\theta^*}|)\nabla_{\theta^*}^2 \ell_{\theta^*}] \quad , \quad M_{2, \theta^*} = \mathbb{E}_{p_{\theta^*}^{\otimes 2}} [\sigma(-|\ell_{\theta^*}|)\nabla_{\theta^*}\ell_{\theta^*}(\nabla_{\theta^*}\ell_{\theta^*})^{\top}] \quad \text{and} \\ R_{\theta^*}^{\text{SPdet}} &= -\mathbb{E}_{p_{\theta^*}^{\otimes 2}} [\text{sign}(\ell_{\theta^*})\sigma(-|\ell_{\theta^*}|) (\nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(\nabla_{\theta^*}\ell_{\theta^*})^{\top} + \nabla_{\theta^*}\ell_{\theta^*}(\nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2})^{\top})] . \end{aligned}$$

Using that $\mathcal{I}(q_{\theta^*, h_{\text{sto}}}) = \mathcal{I}(p_{\theta^*}^{\otimes 2}) + \Delta_{\theta^*}^{\text{SP}}$ (Lemma 3.1), SP_{det} is asymptotically better than SP if and only if $V_{1, \theta^*}^{-1} V_{2, \theta^*} V_{1, \theta^*}^{-1} \prec \mathcal{I}(q_{\theta^*, h_{\text{sto}}})^{-1}$. This condition can be rewritten as

$$\mathcal{I}(p_{\theta^*}^{\otimes 2}) + M_{2, \theta^*} - \Delta_{\theta^*}^{\text{SP}} - R_{\theta^*}^{\text{SPdet}} \prec \left(\mathcal{I}(p_{\theta^*}^{\otimes 2}) + \Delta_{\theta^*}^{\text{SP}} - H_{\theta^*}^{\text{SPdet}} \right) \left(\mathcal{I}(q_{\theta^*, h_{\text{sto}}}) - H_{\theta^*}^{\text{SPdet}} \right) \mathcal{I}(q_{\theta^*, h_{\text{sto}}})^{-1} . \quad (6)$$

The condition (6) heavily depends on the geometry of $p_{\theta^*}^{\otimes 2}$ and ℓ_{θ^*} , hence it is unreasonable to assume in all generality.

In the following, we consider the special case where $H_{\theta^*}^{\text{SP det}} = 0_{d \times d}$. This occurs when $\theta \rightarrow \ell_\theta$ is linear, e.g., for $\mathcal{F}_{\mathcal{N}, \Sigma}$ and $\mathcal{F}_{\text{Lap}, b}$ and preferences based on $r_\theta = \log p_\theta$. Then, the condition (6) rewrites as

$$R_{\theta^*}^{\text{SP det}} + \Delta_{\theta^*}^{\text{SP det}} \succ 0_{d \times d} \quad \text{with} \quad \Delta_{\theta^*}^{\text{SP det}} := 2\Delta_{\theta^*}^{\text{SP}} - M_{2, \theta^*}.$$

Using that $\min_{x \in \mathbb{R}} \sigma(|x|) = 1/2$ achieved only at $x = 0$, we have directly that, for all $x \in \mathbb{R}^k$,

$$x^\top \Delta_{\theta^*}^{\text{SP det}} x = \|x\|^2 \mathbb{E}_{p_{\theta^*}^{\otimes 2}} [(2\sigma(|\ell_{\theta^*}|) - 1)\sigma(-|\ell_{\theta^*}|) \langle x/\|x\|, \nabla_{\theta^*} \ell_{\theta^*} \rangle^2] \geq 0,$$

It is direct to see that this inequality is strict except if $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\ell_{\theta^*} \langle x/\|x\|, \nabla_{\theta^*} \ell_{\theta^*} \rangle = 0) = 1$. Therefore, $\Delta_{\theta^*}^{\text{SP det}}$ is a positive definite matrix if $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(|\ell_{\theta^*} \langle u, \nabla_{\theta^*} \ell_{\theta^*} \rangle| > 0) > 0$ for all $u \in \mathcal{S}_{k-1}$. Then, a sufficient condition for the condition (6) to hold is that $R_{\theta^*}^{\text{SP det}}$ is a p.s.d. matrix, i.e., $R_{\theta^*}^{\text{SP det}} \succeq 0_{d \times d}$.

In summary, we have derived sufficient conditions for SP_{det} to be asymptotically better than SP, namely $H_{\theta^*}^{\text{SP det}} = 0_{d \times d}$, $R_{\theta^*}^{\text{SP det}} \succeq 0_{d \times d}$ and $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(|\ell_{\theta^*} \langle u, \nabla_{\theta^*} \ell_{\theta^*} \rangle| > 0) > 0$ for all $u \in \mathcal{S}_{k-1}$. Note that this last condition is implied by Assumption 4.5. \square

D. Proofs of Section 4

D.1. Proof of Lemma 4.1

Proof. For Gaussian distributions, this is a direct consequence of the following facts: $\theta^* \in \mathcal{C}_n$, $\hat{\theta}_n^{\text{DP}} \in \arg \min_{\theta \in \mathcal{C}_n} \|\theta - \hat{\theta}_n^{\text{SO}}\|_\Sigma^2$ and \mathcal{C}_n is convex. \square

D.2. Proof of Lemma 4.6

Proof. Let $u \in \mathcal{S}_{k-1}$. Let $\tilde{F}_{\theta^*, u}$ be the c.d.f. of $V_{\theta^*, u}(X, Y)$ when $(X, Y) \sim (p_{\theta^*}^{\otimes 2})|_{\mathcal{G}_1(\theta^*, u)}$, i.e., $p_{\theta^*}^{\otimes 2}$ truncated to $\mathcal{G}_1(\theta^*, u)$. Then, $\tilde{F}_{\theta^*, u}(\varepsilon) = F_{\theta^*, u}(\varepsilon)/\alpha_{\theta^*, u}$. Let $\alpha_{\theta^*, u} = \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*, u))$ and $N_{\theta^*, u} = \sum_{i \in [n]} \mathbb{1}((X_i, Y_i) \in \mathcal{G}_1(\theta^*, u)) \sim \text{Bin}(n, \alpha_{\theta^*, u})$. Let $\tilde{R}_{n, u} = \min_{i \in [n], (X_i, Y_i) \in \mathcal{G}_1(\theta^*, u)} V_{\theta^*, u}(X_i, Y_i)$. Using the derivation in Section 4.1, we have that $R_{n, u} \leq \tilde{R}_{n, u}$. Let $\varepsilon > 0$. Conditioned on $N_{\theta^*, u}$, it is direct to see that

$$\mathbb{P}(\tilde{R}_{n, u} > \varepsilon \mid N_{\theta^*, u}) = 1 - (1 - (1 - \tilde{F}_{\theta^*, u}(\varepsilon))^{N_{\theta^*, u}}) = \left(1 - \tilde{F}_{\theta^*, u}(\varepsilon)\right)^{N_{\theta^*, u}}.$$

Using that $N_{\theta^*, u} \sim \text{Bin}(n, \alpha_{\theta^*, u})$, $\mathbb{E}_{X \sim \text{Bin}(n, p)}[s^X] = (1 - p + ps)^n$ and $1 - x \leq \exp(-x)$, we obtain that

$$\mathbb{P}(R_{n, u} > \varepsilon) \leq \mathbb{P}(\tilde{R}_{n, u} > \varepsilon) \leq \left(1 - \alpha_{\theta^*, u} \tilde{F}_{\theta^*, u}(\varepsilon)\right)^n \leq \exp(-n F_{\theta^*, u}(\varepsilon)).$$

Taking $\varepsilon = F_{\theta^*, u}^{-1}(\min\{1, \log(1/\delta)/n\})$ concludes the proof. \square

D.3. Proof of Theorem 4.8

Proof. For all $u \in \mathcal{S}_{k-1}$, let $(A_{\theta^*}, B_{\theta^*}, C_{\theta^*})$ defined as in Theorem 4.8. Since $\mathcal{C}_n \subseteq \tilde{\mathcal{C}}_n$ under Assumption 4.4, we obtain that $\max_{\theta \in \mathcal{C}_n} \|\theta - \theta^*\| \leq \max_{\theta \in \tilde{\mathcal{C}}_n} \|\theta - \theta^*\|$.

Case $k = 1$. Since $|\mathcal{S}_0| = 2$, using Lemma 4.6 with a union bound yield that, with probability at least $1 - \delta$,

$$\max_{\theta \in \tilde{\mathcal{C}}_n} \|\theta - \theta^*\| \leq \max_{u \in \mathcal{S}_0} R_{n, u} \leq \max_{u \in \mathcal{S}_0} F_{\theta^*, u}^{-1}(\min\{1, \log(2/\delta)/n\}).$$

Under Assumption 4.7, for $n \geq B_{\theta^*} \log(2/\delta)$, we can conclude the proof since

$$n \max_{\theta \in \mathcal{C}_n} \|\theta - \theta^*\| \leq n \max_{\theta \in \tilde{\mathcal{C}}_n} \|\theta - \theta^*\| \leq A_{\theta^*} \log(2/\delta) + C_{\theta^*} \log(2/\delta)^2/n.$$

Case $k > 1$. Let $N(\gamma)$ be the γ -covering number of Θ for the norm $\|\cdot\|$. Let $\{\theta_j\}_{j \in [N(\gamma)]}$ be such a γ -covering. For all $j \in [N(\gamma)]$, let $\varepsilon_j = \|\theta_j - \theta^*\|$ and $u_j = (\theta_j - \theta^*)/\varepsilon_j$. Using triangular inequality, we obtain

$$\max_{\theta \in \tilde{\mathcal{C}}_n} \|\theta - \theta^*\| \leq \gamma + \max_{j \in [N(\gamma)], \theta_j \in \tilde{\mathcal{C}}_n} \|\theta_j - \theta^*\| \leq \gamma + \max_{j \in [N(\gamma)]} \mathbb{1}(\theta^* + \varepsilon_j u_j \in \tilde{\mathcal{C}}_n) \varepsilon_j \leq \gamma + \max_{j \in [N(\gamma)]} R_{n, u_j}.$$

Using Lemma 4.6 with a union bound yield that, with probability at least $1 - \delta$,

$$n \max_{\theta \in \mathcal{C}_n} \|\theta - \theta^*\| \leq n\gamma + \max_{j \in [N(\gamma)]} nF_{\theta^*, u_j}^{-1}(\log(N(\gamma)/\delta)/n) \leq n\gamma + A_{\theta^*} \log(N(\gamma)/\delta) + C_{\theta^*} \log(N(\gamma)/\delta)^2/n.$$

where the last inequality relies on Assumption 4.7 for $n \geq B_{\theta^*} \log(N(\gamma)/\delta)$. \square

E. Proofs of Section 5

E.1. Proof of Lemma 5.1

Proof. It is direct to see that

$$\sqrt{q_{\theta^*}(x, y, z)q_{\theta}(x, y, z)} = \begin{cases} 0 & \text{if } (x, y) \in \mathcal{D}(\theta^*, \theta) \cup (\mathcal{G}_0(\theta^*)^c \triangle \mathcal{G}_0(\theta)^c) \\ \sqrt{p_{\theta^*}(x)p_{\theta^*}(y)p_{\theta}(x)p_{\theta}(y)} & \text{otherwise} \end{cases}.$$

Therefore, we have

$$\text{BC}(q_{\theta^*}, q_{\theta}) = \int_{(x, y) \notin \mathcal{D}(\theta^*, \theta) \cup (\mathcal{G}_0(\theta^*)^c \triangle \mathcal{G}_0(\theta)^c)} \sqrt{p_{\theta^*}(x)p_{\theta^*}(y)p_{\theta}(x)p_{\theta}(y)} dx dy = \text{BC}(p_{\theta^*}^{\otimes 2}, p_{\theta}^{\otimes 2}) - \widetilde{\text{BC}}(\theta^*, \theta).$$

Using that $H^2(\mathbb{P}, \mathbb{Q}) = 1 - \text{BC}(\mathbb{P}, \mathbb{Q})$, we conclude the proof. \square

E.2. Proof of Theorem 5.3

Consider the hypercube $\Theta' = \{\theta_b = \delta b : b \in \{0, 1\}^d\} \subseteq \Theta$. Note that $\|\theta_b - \theta_{b'}\| \geq \frac{\delta}{\sqrt{k}} d_H(b, b')$, where $d_H(b, b')$ denotes the Hamming distance between b and b' . Then using Assouad's lemma we have

$$R_{\max} \geq \frac{\delta\sqrt{k}}{4} \left(1 - \max_{d_H(b, b')=1} \text{TV}(q_{\theta_b}^{\otimes n}, q_{\theta_{b'}}^{\otimes n}) \right).$$

Upper bounding TV with H^2 , Lemma 5.1 yields

$$\text{TV}(q_{\theta_b}^{\otimes n}, q_{\theta_{b'}}^{\otimes n}) \leq \sqrt{n(\widetilde{\text{BC}}(\theta_b, \theta_{b'}) + H^2(p_{\theta_b}^{\otimes 2}, p_{\theta_{b'}}^{\otimes 2}))}.$$

Then, Assumption 5.2 implies

$$R_{\max} \geq \frac{\delta\sqrt{k}}{4} \left(1 - \sqrt{n \left(\frac{c_1\delta}{\alpha_{\mathcal{F}}(k)} + 2c_2\delta^2 \right)} \right).$$

Picking $\delta = \frac{1}{2(c_1+2c_2)} \min\left\{\frac{\alpha_{\mathcal{F}}(k)}{n}, \frac{1}{\sqrt{n}}\right\}$ ensures that the term in parenthesis is always greater than $1/2$, hence

$$R_{\max} \geq \frac{\sqrt{k}}{8(c_1+2c_2)} \min\left\{\frac{\alpha_{\mathcal{F}}(k)}{n}, \frac{1}{\sqrt{n}}\right\}.$$

F. Multivariate Gaussian with Known Covariance

In the following, $\theta = \Sigma^{-1}\mu$ denote the natural parameter of multivariate Gaussian with known covariance matrix Σ . We have $\mathcal{X} = \mathbb{R}^d$ and $k = d$. Let $\theta \in \Theta$ and $u \in \mathcal{S}_{d-1}$ for the norm $\|\cdot\|_{\Sigma}$, i.e., $\|u\|_{\Sigma} = 1$. Let $\mathcal{S}_{2,d-1} = \{x \in \mathbb{R}^d \mid \|u\|_2 = 1\}$. It is direct to see that

$$\ell_{\theta}(x, y) = \log \frac{p_{\theta}(x)}{p_{\theta}(y)} = \langle x - y, \theta - \Sigma^{-1}(x + y)/2 \rangle \quad \text{and} \quad \nabla_{\theta^*} \ell_{\theta^*}(x, y) = x - y.$$

Therefore, we have

$$\mathcal{G}_0(\theta^*) = \{(x, y) \in (\mathbb{R}^d)^2 \mid |\langle x - y, \theta^* - (x + y)/2 \rangle| > 0\},$$

$$\mathcal{G}_1(\theta^*) = \{(x, y) \in \mathcal{G}_0(\theta^*) \mid \|x - y\| > 0\},$$

$$\mathcal{D}(\theta^*, \theta) = \{(x, y) \in (\mathbb{R}^d)^2 \mid \langle x - y, \theta^* - \Sigma^{-1}(x + y)/2 \rangle^2 + \langle x - y, \theta^* - \Sigma^{-1}(x + y)/2 \rangle \langle \theta - \theta^*, x - y \rangle < 0\},$$

$$\mathcal{G}_1(\theta^*, u) = \{(x, y) \in (\mathbb{R}^d)^2 \mid \langle x - y, \theta^* - \Sigma^{-1}(x + y)/2 \rangle \langle u, x - y \rangle < 0\},$$

$$\forall (x, y) \in \mathcal{G}_1(\theta^*, u), \quad V_{\theta^*, u}(x, y) = \frac{\langle x - y, \Sigma^{-1}(x + y)/2 - \theta^* \rangle}{\langle u, x - y \rangle}.$$

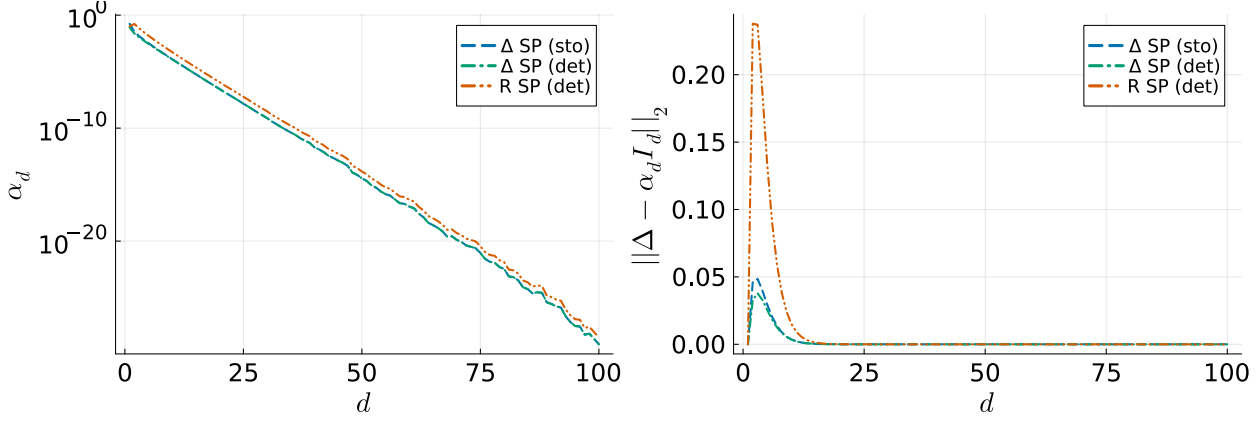


Figure 3. Approximations of $\Delta_{\mathcal{N}(0_d, I_d)}^{\text{SP}}$, $\Delta_{\mathcal{N}(0_d, I_d)}^{\text{SPdet}}$ and $R_{\mathcal{N}(0_d, I_d)}^{\text{SPdet}}$ by (a) $\alpha_d I_d$ and (b) associated error for varying d . $N_{\text{runs}} = 10^6$.

Proof that $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*)) > 0$. It is direct to see that $\dim(\mathcal{G}_0(\theta^*)^c) < 2d$ and $\dim(\mathcal{G}_0(\theta^*) \setminus \mathcal{G}_1(\theta^*)) < 2d$. Given that $p_{\theta^*}^{\otimes 2}$ is a continuous distribution on $(\mathbb{R}^d)^2$, we obtain that $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*)) = \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_0(\theta^*)) = 1$.

Condition in Lemma 3.1. The condition of Lemma 3.1 is implied by Assumption 4.5, hence we refer to the proof of this result below. Therefore, we have $\mathcal{I}(q_{\theta^*, h_{\text{sto}}}) \succ \mathcal{I}(p_{\theta^*}^{\otimes 2})$.

Consistency of SP_{det} . To study SP_{det} for $\mathcal{F}_{\mathcal{N}, \Sigma}$, we use the change of variable $D = \Sigma^{-1/2}(X - Y)/\sqrt{2}$ and $S = \sqrt{2}\Sigma^{-1/2}(\Sigma\theta^* - (X + Y)/2)$. Then, we have $(D, S) \sim \mathcal{N}(0_{2d}, I_{2d})$ and

$$\ell_{\theta^*}(X, Y) = \langle S, D \rangle, \quad \nabla_{\theta^*} \ell_{\theta^*}(X, Y) = \sqrt{2}\Sigma^{1/2}D, \quad \nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(X, Y) = 2\Sigma\theta^* - \sqrt{2}\Sigma^{1/2}S.$$

Let $M(D, S) = \text{sign}(\langle D, S \rangle)\sigma(-|\langle D, S \rangle|)D$. Then, $M(-D, -S) = -M(D, S)$ for all $(D, S) \in \mathbb{R}^{2d}$. By integration of an odd function with respect to 0_d with a symmetric distribution around 0_{2d} , we obtain $\mathbb{E}_{(D, S) \sim \mathcal{N}(0_{2d}, I_{2d})} [M(D, S)] = 0_d$. Therefore, the condition (3) is satisfied and the SP_{det} is a consistent estimator.

Asymptotic variance of SP_{det} . Let $H_{\theta^*}^{\text{SPdet}}$ and $R_{\theta^*}^{\text{SPdet}}$ defined in Lemma 3.2. Since $\ell_{\theta}(x, y) = \langle x - y, \theta - (x + y)/2 \rangle$ is linear in θ , we have $\nabla_{\theta^*}^2 \ell_{\theta^*} = 0_{d \times d}$ and $H_{\theta^*}^{\text{SPdet}} = 0_{d \times d}$. The condition $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(|\ell_{\theta^*}\langle u, \nabla_{\theta^*} \ell_{\theta^*} \rangle| > 0) > 0$ for all $u \in \mathcal{S}_{d-1}$ is implied by Assumption 4.5, hence we refer to the proof of this result below. Then, the condition $R_{\theta^*}^{\text{SPdet}} \succeq 0_{d \times d}$ is equivalent to $M_3 \succeq 0_{d \times d}$ where

$$M_3 = \mathbb{E}_{(D, S) \sim \mathcal{N}(0_{2d}, I_{2d})} [\text{sign}(\langle D, S \rangle)\sigma(-|\langle D, S \rangle|)(DS^{\top} + SD^{\top})].$$

When $d = 1$, we have $M_3 = 2\mathbb{E}_{(D, S) \sim \mathcal{N}(0_2, I_2)} [\sigma(-|D, S|)|DS|] > 0$. When $d > 1$, for all $u \in \mathcal{S}_{d-1}$, we have

$$u^{\top} M_3 u = 2\mathbb{E}_{(D, S) \sim \mathcal{N}(0_{2d}, I_{2d})} [\text{sign}(\langle D, S \rangle)\sigma(-|\langle D, S \rangle|)\langle u, D \rangle\langle u, S \rangle],$$

By rotational symmetry of $\mathcal{N}(0_{2d}, I_{2d})$ and the function to be integrated, showing that $\min_{u \in \mathcal{S}_{d-1}} u^{\top} M_3 u \geq 0$ is equivalent to showing that $e_1^{\top} M_3 e_1 \geq 0$, i.e.,

$$\mathbb{E}_{(D, S) \sim \mathcal{N}(0_{2d}, I_{2d})} [\text{sign}(\langle D, S \rangle)\sigma(-|\langle D, S \rangle|)D_1 S_1].$$

By symmetry, we conjecture that $\Delta_{\mathcal{N}(0_d, I_d)}^{\text{SP}}$, $\Delta_{\mathcal{N}(0_d, I_d)}^{\text{SPdet}}$ and $R_{\mathcal{N}(0_d, I_d)}^{\text{SPdet}}$ are of the form $\alpha_d I_d$ where $\alpha_d > 0$ is decreasing in d . Figure 3 validates this conjecture numerically.

Using the sufficient condition derived in Appendix C.2, we have shown that SP_{det} is asymptotically better than SP.

Proof of Assumption 4.4. Since $\ell_{\theta}(x, y) = \langle x - y, \theta - (x + y)/2 \rangle$ is linear in θ , we have $\mathcal{D}(\theta^*, \theta) = \tilde{\mathcal{D}}(\theta^*, \theta)$.

Proof of Assumption 4.5 For $(X, Y) \sim p_{\theta^*}^{\otimes 2}$, let $D = \Sigma^{-1/2}(X - Y)/\sqrt{2}$ and $S = \sqrt{2}\Sigma^{-1/2}((X + Y)/2 - \Sigma\theta^*)$. Then, we have $(D, S) \sim \mathcal{N}(0_{2d}, I_{2d})$. Defining $U = D/\|D\|$, we have $U \sim \mathcal{U}(\mathcal{S}_{2,d-1})$ is independent of S . Since $U = D/\|D\| \sim \mathcal{U}(\mathcal{S}_{2,d-1})$ and $\Sigma^{-1}(x + y)/2 - \theta^* = \Sigma^{-1/2}S/\sqrt{2}$, we obtain

$$\begin{aligned} \mathbb{P}_{(X,Y) \sim p_{\theta^*}^{\otimes 2}}((X, Y) \in \mathcal{G}_1(\theta^*, u)) &= \mathbb{P}_{(U,S) \sim \mathcal{U}(\mathcal{S}_{2,d-1}) \otimes \mathcal{N}(0_d, I_d)}(\langle U, S \rangle \langle u, \Sigma^{1/2}U \rangle > 0) \\ &= \mathbb{P}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})}(\langle \Sigma^{1/2}u, U \rangle > 0) / 2 + \mathbb{P}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})}(\langle \Sigma^{1/2}u, U \rangle < 0) / 2 = 1/2. \end{aligned}$$

where we used that, conditioned on U , $\langle U, S \rangle \sim \mathcal{N}(0, 1)$ and $\mathbb{P}_{X \sim \mathcal{N}(0,1)}(X < 0) = \mathbb{P}_{X \sim \mathcal{N}(0,1)}(X > 0) = 1/2$. The last equality uses that $\mathbb{P}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})}(\langle \Sigma^{1/2}u, U \rangle > 0) = \mathbb{P}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})}(\langle \Sigma^{1/2}u, U \rangle < 0) = 1/2$ by symmetry of the uniform distribution. Therefore, we have shown that $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*, u)) = 1/2$ for all $u \in \mathcal{S}_{2,d-1}$.

Proof of Assumption 4.7. Let us define $v = \Sigma^{1/2}u$, hence $v \in \mathcal{U}(\mathcal{S}_{2,d-1})$. Let Φ denote the c.d.f. of $\mathcal{N}(0, 1)$ and $\text{erf}(x) = 2\Phi(x\sqrt{2}) - 1$ be the error function. Let $\varepsilon > 0$. Similarly as above, we obtain that

$$\begin{aligned} F_{\theta^*, u}(\varepsilon) &= \mathbb{P}_{(X,Y) \sim p_{\theta^*}^{\otimes 2}}(0 < V_{\theta^*, u}(X, Y) \leq \varepsilon) \\ &= \mathbb{P}_{(U,S) \sim \mathcal{U}(\mathcal{S}_{2,d-1}) \otimes \mathcal{N}(0_d, I_d)}\left(0 < \frac{\langle U, S \rangle}{\langle v, U \rangle} \leq \sqrt{2}\varepsilon\right) \\ &= \frac{1}{2} \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} \left[2\Phi\left(\sqrt{2}\varepsilon|\langle v, U \rangle|\right) - 1 \right] = \frac{1}{2} \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [\text{erf}(\varepsilon|\langle v, U \rangle|)]. \end{aligned}$$

where we use conditioning by U as above. By change of variable, we obtain that

$$F_{\theta^*, u}(\varepsilon) = \frac{1}{\sqrt{\pi}} \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} \left[\int_0^{\varepsilon|\langle v, U \rangle|} e^{-t^2} dt \right] = \frac{\varepsilon}{\sqrt{\pi}} \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} \left[|\langle v, U \rangle| \int_0^1 e^{-x^2 \varepsilon^2 \langle v, U \rangle^2} dx \right].$$

Using that $1 - x^2 \leq e^{-x^2} \leq 1$, we obtain that

$$0 \geq \frac{\sqrt{\pi}}{\varepsilon} F_{\theta^*, u}(\varepsilon) - \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|] \geq -\frac{\varepsilon^2}{3} \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|^3].$$

Using that

$$\int_0^1 x(-2x\varepsilon^2 \langle v, U \rangle^2) e^{-x^2 \varepsilon^2 \langle v, U \rangle^2} dx = e^{-\varepsilon^2 \langle v, U \rangle^2} - \int_0^1 e^{-x^2 \varepsilon^2 \langle v, U \rangle^2} dx,$$

we obtain

$$F'_{\theta^*, u}(\varepsilon) = \frac{1}{\sqrt{\pi}} \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle| e^{-\varepsilon^2 \langle v, U \rangle^2}] \quad \text{and} \quad F''_{\theta^*, u}(\varepsilon) = \frac{2\varepsilon}{\sqrt{\pi}} \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|^3 e^{-\varepsilon^2 \langle v, U \rangle^2}].$$

Therefore, using Lemma F.1, we have

$$F'_{\theta^*, u}(0) = \frac{1}{\sqrt{\pi}} \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|] = \frac{2}{d-1} \frac{\Gamma(d/2)}{\pi \Gamma((d-1)/2)} =_{d \rightarrow +\infty} \mathcal{O}(1/\sqrt{d})$$

Let us define

$$\varepsilon_{\theta^*, u} = \sqrt{\frac{\mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|]}{2 \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|^3]}} \quad \text{and} \quad M_{\theta^*, u} = 4\pi \sqrt{\frac{\mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|^3]}{\mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|]^5}}.$$

Then, for all $\varepsilon \in (0, \varepsilon_{\theta^*, u}]$, we obtain that

$$\begin{aligned} \frac{F''_{\theta^*, u}(\varepsilon)}{F'_{\theta^*, u}(\varepsilon)^3} &= \frac{\pi \varepsilon \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|^3 e^{-\varepsilon^2 \langle v, U \rangle^2}]}{2 \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle| e^{-\varepsilon^2 \langle v, U \rangle^2}]^3} \\ &\leq \frac{\pi \varepsilon}{2} \frac{\mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|^3]}{(\mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|] - \varepsilon^2 \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle v, U \rangle|^3])^3} \leq M_{\theta^*, u}. \end{aligned}$$

Since we have $(F_{\theta^*,u}^{-1})''(x) = -\frac{F_{\theta^*,u}''(F_{\theta^*,u}^{-1}(x))}{F_{\theta^*,u}'(F_{\theta^*,u}^{-1}(x))^3}$, we obtain

$$\sup_{x \in (0, x_{\theta^*,u}]} |(F_{\theta^*,u}^{-1})''(x)| \leq M_{\theta^*,u} \quad \text{where} \quad x_{\theta^*,u} = F_{\theta^*,u}(\varepsilon_{\theta^*,u}) \leq \sqrt{\frac{\mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle \Sigma^{1/2}u, U \rangle|^3]}{2\pi \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle \Sigma^{1/2}u, U \rangle|^3]}}.$$

Proof of Assumption 4.2. Let $\varepsilon = \|\theta^* - \theta\|$ and $u = (\theta^* - \theta)/\varepsilon$. Then, we have

$$\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{D}(\theta^*, \theta)) = \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{D}(\theta^*, \theta^* + \varepsilon u)) \geq \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{D}(\theta^*, \theta^* + \varepsilon u) \cap \mathcal{G}_1(\theta^*, u)) = \mathbb{P}_{(X,Y) \sim p_{\theta^*}^{\otimes 2}}(0 < V_{\theta^*,u}(X, Y) < \varepsilon)$$

Using the above computation, we obtain that $\mathbb{P}_{(X,Y) \sim p_{\theta^*}^{\otimes 2}}(0 < V_{\theta^*,u}(X, Y) < \varepsilon) > 0$, hence $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{D}(\theta^*, \theta)) > 0$.

Proof of Assumption 5.2. Using that $1 - e^{-x} \leq x$, we obtain

$$H^2(p_{\theta^*}, p_{\theta}) = 1 - \exp\left(-\frac{1}{8}\|\theta^* - \theta\|_{\Sigma}^2\right) \leq \frac{1}{8}\|\theta^* - \theta\|_{\Sigma}^2.$$

First, we notice that $\dim(\mathcal{G}_0(\theta^*)^c \triangle \mathcal{G}_0(\theta)^c) < 2d$, hence we can show that

$$\int_{(x,y) \in \mathcal{G}_0(\theta^*)^c \triangle \mathcal{G}_0(\theta)^c} \sqrt{p_{\theta^*}(x)p_{\theta^*}(y)p_{\theta}(x)p_{\theta}(y)} dx dy = 0.$$

Second, we see that

$$\begin{aligned} & \|x - \Sigma\theta^*\|_{\Sigma^{-1}}^2 + \|y - \Sigma\theta^*\|_{\Sigma^{-1}}^2 + \|x - \Sigma\theta\|_{\Sigma^{-1}}^2 + \|y - \Sigma\theta\|_{\Sigma^{-1}}^2 \\ &= \|x - y\|_{\Sigma^{-1}}^2 + \|\theta - \theta^*\|_{\Sigma}^2 + \|x + y - \Sigma(\theta^* + \theta)\|_{\Sigma^{-1}}^2, \\ \mathcal{D}(\theta^*, \theta) &= \{(x, y) \in (\mathbb{R}^d)^2 \mid \langle x - y, \theta^* - \Sigma^{-1}(x + y)/2 \rangle^2 + \langle x - y, \theta^* - \Sigma^{-1}(x + y)/2 \rangle \langle \theta - \theta^*, x - y \rangle < 0\}, \end{aligned}$$

Then, we consider the change of variable $u = \Sigma^{-1/2}(x - y)$ and $v = \Sigma^{-1/2}(x + y)$, whose Jacobian has $\det(\Sigma)2^{-d}$ as absolute value of its determinant. Therefore, we obtain

$$\begin{aligned} e^{\frac{1}{4}\|\theta - \theta^*\|_{\Sigma}^2} \widetilde{\text{BC}}(\theta^*, \theta) &= \frac{1}{(4\pi)^d} \int_{(u,v)} \mathbf{1}\left(0 < -\frac{\langle u, \Sigma^{1/2}\theta^* - v/2 \rangle}{\langle \Sigma^{1/2}(\theta - \theta^*), u \rangle} < 1\right) e^{-\frac{1}{4}\|u\|^2 - \frac{1}{4}\|v - \Sigma^{1/2}(\theta + \theta^*)\|^2} du dv \\ &= \frac{1}{(2\pi)^d} \int_{(\tilde{u}, \tilde{v})} \mathbf{1}\left(\left|\frac{\langle \tilde{u}, \tilde{v} \rangle}{\langle \Sigma^{1/2}(\theta - \theta^*), \tilde{u} \rangle}\right| < \sqrt{2}\right) e^{-\frac{1}{2}\|\tilde{u}\|^2 - \frac{1}{2}\|\tilde{v}\|^2} d\tilde{u} d\tilde{v} \\ &= \mathbb{P}_{(X,Y) \sim \mathcal{N}(0_d, I_d)^{\otimes 2}}\left(\left|\frac{\langle X, Y \rangle}{\langle \Sigma^{1/2}(\theta - \theta^*), X \rangle}\right| < \sqrt{2}\right) \\ &= \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})}\left[\text{erf}\left(|\langle \Sigma^{1/2}(\theta - \theta^*), U \rangle|\right)\right] = 2F_{\theta^*,u}(\varepsilon) \end{aligned}$$

where the second equality uses the change of variable $\tilde{u} = u/\sqrt{2}$ and $\tilde{v} = (v - \Sigma^{1/2}(\theta + \theta^*))/\sqrt{2}$, whose Jacobian has determinant 2^d . The third and the fourth re-uses computation done previously with $\varepsilon = \|\theta - \theta^*\|_{\Sigma}$ and $u = (\theta - \theta^*)/\varepsilon$. Using Lemma F.1 and the above upper bound on $F_{\theta^*,u}(\varepsilon)$, we obtain

$$\widetilde{\text{BC}}(\theta^*, \theta) = 2e^{-\varepsilon^2/4} F_{\theta^*,u}(\varepsilon) \leq \frac{4}{d-1} \frac{\Gamma(d/2)}{\pi \Gamma((d-1)/2)} \|\theta^* - \theta\|_{\Sigma}.$$

Lemma F.1. Let Γ be the Γ function. Then,

$$\forall u \in \mathcal{S}_{2,d-1}, \quad \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle u, U \rangle|] = \frac{2}{d-1} \frac{\Gamma(d/2)}{\sqrt{\pi} \Gamma((d-1)/2)} =_{d \rightarrow +\infty} \mathcal{O}(1/\sqrt{d}).$$

Proof. Due to rotational symmetry of the distribution, for any unit vector u ,

$$\mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle u, U \rangle|] = \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|\langle e_1, U \rangle|] = \mathbb{E}_{U \sim \mathcal{U}(\mathcal{S}_{2,d-1})} [|U_1|].$$

The density of U_1 is given by

$$f_{U_1}(x) = \frac{\Gamma(\frac{d}{2})}{\sqrt{\pi} \Gamma(\frac{d-1}{2})} (1-x^2)^{\frac{d-3}{2}}, \quad x \in [-1, 1],$$

and the expectation can be computed as

$$\begin{aligned} \mathbb{E}[|U_1|] &= \int_{-1}^1 |x| f_{U_1}(x) dx = 2 \int_0^1 x \frac{\Gamma(\frac{d}{2})}{\sqrt{\pi} \Gamma(\frac{d-1}{2})} (1-x^2)^{\frac{d-3}{2}} dx = \frac{2\Gamma(\frac{d}{2})}{\sqrt{\pi} \Gamma(\frac{d-1}{2})} \int_0^1 x (1-x^2)^{\frac{d-3}{2}} dx \\ &= \frac{\Gamma(\frac{d}{2})}{\sqrt{\pi} \Gamma(\frac{d-1}{2})} \int_0^1 (1-u)^{\frac{d-3}{2}} du = \frac{\Gamma(\frac{d}{2})}{\sqrt{\pi} \Gamma(\frac{d-1}{2})} \cdot \frac{1}{\frac{d-1}{2}} = \frac{2}{d-1} \frac{\Gamma(\frac{d}{2})}{\sqrt{\pi} \Gamma(\frac{d-1}{2})}. \end{aligned}$$

Therefore, for large d , $\mathbb{E}[|U_1|] =_{d \rightarrow +\infty} \mathcal{O}(1/\sqrt{d})$. \square

G. Laplace with Known Scale

In the following, θ denote the mean parameter of Laplace distribution with known scale b . We have $\mathcal{X} = \mathbb{R}$ and $k = d = 1$. Let $\theta \in \Theta$ and $u \in \{\pm 1\}$. It is direct to see that

$$\ell_\theta(x, y) = \log \frac{p_\theta(x)}{p_\theta(y)} = |y - \theta|/b - |x - \theta|/b = \frac{1}{b} \begin{cases} y - x & \text{if } \theta < \min\{x, y\} \\ x - y & \text{if } \theta > \max\{x, y\} \\ (2\theta - (x + y))\text{sign}(x - y) & \text{if } \theta \in [\min\{x, y\}, \max\{x, y\}] \end{cases},$$

$$\text{and } \nabla_{\theta^*} \ell_{\theta^*}(x, y) = \begin{cases} 0 & \text{if } \theta^* < \min\{x, y\} \text{ or } \theta^* > \max\{x, y\} \\ \frac{2}{b} \text{sign}(x - y) & \text{if } \theta^* \in [\min\{x, y\}, \max\{x, y\}] \end{cases}.$$

Therefore, we have

$$\begin{aligned} \mathcal{G}_0(\theta^*) &= \{(x, y) \in \mathbb{R}^2 \mid ||y - \theta^*| - |x - \theta^*|| > 0\}, \\ \mathcal{G}_1(\theta^*) &= \{(x, y) \in \mathbb{R}^2 \mid \theta^* \in [\min\{x, y\}, \max\{x, y\}]\}, \\ \mathcal{G}_1(\theta^*, u) &= \{(x, y) \in \mathbb{R}^2 \mid \theta^* \in [\min\{x, y\}, \max\{x, y\}] \wedge u((x + y)/2 - \theta^*) > 0\}, \\ \tilde{\mathcal{D}}(\theta^*, \theta) &= \{(x, y) \in \mathbb{R}^2 \mid \theta^* \in [\min\{x, y\}, \max\{x, y\}] \wedge 0 < \text{sign}(\theta - \theta^*)((x + y)/2 - \theta^*) < |\theta - \theta^*|\}, \\ \forall (x, y) \in \mathcal{G}_1(\theta^*, u), \quad V_{\theta^*, u}(x, y) &= u((x + y)/2 - \theta^*). \end{aligned}$$

When $\theta^* > \theta$, we have

$$\begin{aligned} \mathcal{D}(\theta^*, \theta) &= \{(x, y) \mid \{\theta^*, \theta\} \subset [\min\{x, y\}, \max\{x, y\}] \wedge \theta < (x + y)/2 < \theta^*\} \\ &\cup \{(x, y) \mid \theta < \min\{x, y\} \wedge \theta^* \in ((x + y)/2, \max\{x, y\}]\} \\ &\cup \{(x, y) \mid \theta < \min\{x, y\} \wedge \theta^* > \max\{x, y\}\} \\ &\cup \{(x, y) \mid \theta^* > \max\{x, y\} \wedge \theta \in [\min\{x, y\}, (x + y)/2)\}. \end{aligned}$$

When $\theta^* < \theta$, we have

$$\begin{aligned} \mathcal{D}(\theta^*, \theta) &= \{(x, y) \mid \{\theta^*, \theta\} \subset [\min\{x, y\}, \max\{x, y\}] \wedge \theta^* < (x + y)/2 < \theta\} \\ &\cup \{(x, y) \mid \theta > \max\{x, y\} \wedge \theta^* \in [\min\{x, y\}, (x + y)/2)\} \\ &\cup \{(x, y) \mid \theta^* < \min\{x, y\} \wedge \theta > \max\{x, y\}\} \\ &\cup \{(x, y) \mid \theta^* < \min\{x, y\} \wedge \theta \in ((x + y)/2, \max\{x, y\}]\}. \end{aligned}$$

Proof that $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*)) > 0$. It is direct to see that $\dim(\mathcal{G}_0(\theta^*)^\complement) < 2$. Given that $p_{\theta^*}^{\otimes 2}$ is a continuous distribution on $(\mathbb{R})^2$, we obtain that $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_0(\theta^*)) = 1$. Using the symmetry of the Laplace distribution around its mean, we have that

$$\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*)) = \mathbb{P}_{p_{\theta^*}^{\otimes 2}}((-\infty, \theta^*) \times (\theta^*, +\infty)) + \mathbb{P}_{p_{\theta^*}^{\otimes 2}}((\theta^*, +\infty) \times (-\infty, \theta^*)) = 1/2.$$

Condition in Lemma 3.1. The condition of Lemma 3.1 is implied by Assumption 4.5, hence we refer to the proof of this result below. Therefore, we have $\mathcal{I}(q_{\theta^*, h_{\text{sto}}}) \succ \mathcal{I}(p_{\theta^*}^{\otimes 2})$.

Consistency of SP_{det} . To study SP_{det} for $\mathcal{F}_{\text{Lap}, b}$, we use the change of variable $D = \theta^* - X$ and $S = \theta^* - Y$. For all $(D, S) \in \mathcal{G}_1(0)$, we have

$$\ell_{\theta^*}(X, Y) = \frac{1}{b}(D + S)\text{sign}(S - D), \quad \nabla_{\theta^*} \ell_{\theta^*}(X, Y) = \frac{2}{b}\text{sign}(S - D), \quad \nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(X, Y) = 0.$$

For all $(D, S) \notin \mathcal{G}_1(0)$, we have $\nabla_{\theta^*} \ell_{\theta^*}(X, Y) = 0$ and $\nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(X, Y) \neq 0$. Let $M(D, S) = \mathbb{1}((D, S) \in \mathcal{G}_1(0)) \sigma(-|D + S|/b) \text{sign}(D + S)$. Then, $M(-D, -S) = -M(D, S)$ for all $(D, S) \in \mathbb{R}^2$. By integration of an odd function with respect to 0 with a symmetric distribution around 0, we obtain $\mathbb{E}_{(D, S) \sim \mathcal{N}(0_{2d}, I_{2d})} [M(D, S)] = 0$. Therefore, the condition (3) is satisfied and SP_{det} is a consistent estimator.

Asymptotic variance of SP_{det} . Let $H_{\theta^*}^{\text{SP}_{\text{det}}}$ and $R_{\theta^*}^{\text{SP}_{\text{det}}}$ defined in Lemma 3.2. By definition of ℓ_{θ^*} , we obtain $\nabla_{\theta^*}^2 \ell_{\theta^*} = 0$ and $H_{\theta^*}^{\text{SP}_{\text{det}}} = 0$. Moreover, using the above formula, we have $\nabla_{\theta^*} \ell_{\theta^*}(X, Y) \nabla_{\theta^*} \log p_{\theta^*}^{\otimes 2}(X, Y) = 0$ for all $(D, S) \in \mathcal{G}_1(0)$, hence we obtain $R_{\theta^*}^{\text{SP}_{\text{det}}} = 0$. The condition $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(|\ell_{\theta^*} \langle u, \nabla_{\theta^*} \ell_{\theta^*} \rangle| > 0) > 0$ for all $u \in \mathcal{S}_{d-1}$ is implied by Assumption 4.5, hence we refer to the proof of this result below. Using the sufficient condition derived in Appendix C.2, we have shown that SP_{det} is asymptotically better than SP.

Proof of Assumption 4.4. Using that $\tilde{\mathcal{D}}(\theta^*, \theta) \subseteq \mathcal{G}_1(\theta^*)$, we simply need to show that $\tilde{\mathcal{D}}(\theta^*, \theta) \subseteq \mathcal{G}_1(\theta^*) \cap \mathcal{D}(\theta^*, \theta)$. Let us consider the case $\theta^* > \theta$. Then, we have

$$\begin{aligned} \tilde{\mathcal{D}}(\theta^*, \theta) &= \{(x, y) \in \mathbb{R}^2 \mid \theta^* \in [\min\{x, y\}, \max\{x, y\}] \wedge \theta < (x + y)/2 < \theta^*\} \\ &= \{(x, y) \mid \{\theta^*, \theta\} \subset [\min\{x, y\}, \max\{x, y\}] \wedge \theta < (x + y)/2 < \theta^*\} \\ &\cup \{(x, y) \mid \theta < \min\{x, y\} \wedge \theta^* \in ((x + y)/2, \max\{x, y\}]\} = \mathcal{G}_1(\theta^*) \cap \mathcal{D}(\theta^*, \theta). \end{aligned}$$

The same result follows when $\theta^* < \theta$ by using the same argument. In summary, we have shown that $\tilde{\mathcal{D}}(\theta^*, \theta) = \mathcal{G}_1(\theta^*) \cap \mathcal{D}(\theta^*, \theta) \subseteq \mathcal{D}(\theta^*, \theta)$.

Proof of Assumption 4.5. Using the symmetry of the Laplace distribution around its mean, we have $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*, u)) = \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*, 1))$ for all $u \in \{\pm 1\}$. Then, by integrating for $x < y$, we obtain

$$\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*, 1)) = \frac{1}{2b^2} \int_{x \in (-\infty, \theta^*)} e^{x/b} \left(\int_{y \in (2\theta^* - x, +\infty)} e^{-y/b} dy \right) dx = \frac{1}{2b} \int_{x \in (-\infty, \theta^*)} e^{2x - 2\theta^*/b} dx = \frac{1}{4}.$$

Proof of Assumption 4.7. Let $\varepsilon > 0$. Using the symmetry of the Laplace distribution around its mean, we have $F_{\theta^*, u}(\varepsilon) = F_{\theta^*, 1}(\varepsilon)$ for all $u \in \{\pm 1\}$. Similarly as above, by integrating for $x < y$, we obtain that

$$\begin{aligned} F_{\theta^*, 1}(\varepsilon) &= \mathbb{P}_{(X, Y) \sim p_{\theta^*}^{\otimes 2}}(0 < V_{\theta^*, 1}(X, Y) \leq \varepsilon) \\ &= \frac{1}{2b^2} \int_{x \in (-\infty, \theta^*)} e^{x/b} \left(\int_{y \in (2\theta^* - x, 2\varepsilon + 2\theta^* - x)} e^{-y/b} dy \right) dx \\ &= \frac{1}{2b} \left(\int_{x \in (-\infty, \theta^*)} e^{(2x - 2\theta^*)/b} dx - \int_{x \in (-\infty, \theta^*)} e^{(2x - 2\theta^* - 2\varepsilon)/b} dx \right) = \frac{1}{4} (1 - e^{-2\varepsilon/b}). \end{aligned}$$

Therefore, we have

$$F'_{\theta^*,u}(x) = \frac{1}{2b}e^{-2\varepsilon/b} \quad , \quad F_{\theta^*,u}^{-1}(x) = -\frac{b}{2}\log(1-4x) \quad \text{and} \quad (F_{\theta^*,u}^{-1})''(x) = \frac{8b}{(1-4x)^2}.$$

Then, we obtain $F'_{\theta^*,u}(0) = \frac{1}{2b}$ and we can take $x_{\theta^*,u} = 1/8$ and $M_{\theta^*,u} = 32b$.

Proof of Assumption 4.2. Let $\varepsilon = |\theta^* - \theta|$ and $u = \text{sign}(\theta^* - \theta)$. Using the above computation, we have

$$\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{D}(\theta^*, \theta)) \geq \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\tilde{\mathcal{D}}(\theta^*, \theta^* + \varepsilon u)) \geq \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\tilde{\mathcal{D}}(\theta^*, \theta^* + \varepsilon u) \cap \mathcal{G}_1(\theta^*, u)) = \mathbb{P}_{(X,Y) \sim p_{\theta^*}^{\otimes 2}}(0 < V_{\theta^*,u}(X, Y) < \varepsilon)$$

Using the above computation, we obtain that $\mathbb{P}_{(X,Y) \sim p_{\theta^*}^{\otimes 2}}(0 < V_{\theta^*,u}(X, Y) < \varepsilon) > 0$, hence $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{D}(\theta^*, \theta)) > 0$.

Proof of Assumption 5.2. Using that $f(x) = x^2 - 1 + (1+x)e^{-x}$ is positive on \mathbb{R}_+ , we obtain

$$H^2(p_{\theta^*}, p_{\theta}) = 1 - \left(1 + \frac{|\theta^* - \theta|}{2b}\right) \exp\left(-\frac{|\theta^* - \theta|}{2b}\right) \leq \frac{(\theta^* - \theta)^2}{4b^2}.$$

First, we notice that $\dim(\mathcal{G}_0(\theta^*)^{\mathbb{G}} \triangle \mathcal{G}_0(\theta)^{\mathbb{G}}) < 2$, hence we can show that

$$\int_{(x,y) \in \mathcal{G}_0(\theta^*)^{\mathbb{G}} \triangle \mathcal{G}_0(\theta)^{\mathbb{G}}} \sqrt{p_{\theta^*}(x)p_{\theta^*}(y)p_{\theta}(x)p_{\theta}(y)} dx dy = 0.$$

We consider the case $\theta^* < \theta$ since $\theta^* > \theta$ is done similarly as $\widetilde{\text{BC}}(\theta^*, \theta) = \widetilde{\text{BC}}(\theta, \theta^*)$. Let $\varepsilon = \theta - \theta^*$. By integrating for $x < y$, we have

$$\begin{aligned} \widetilde{\text{BC}}(\theta^*, \theta) &= \frac{1}{2b^2} \int_x e^{x/b} \left(\int_y \mathbb{1}(x \leq \theta^* < (x+y)/2 < \theta^* + \varepsilon \leq y) e^{-y/b} dy \right) dx \\ &+ \frac{e^{-(\varepsilon+\theta^*)/b}}{2b^2} \int_x e^{x/b} \left(\int_y \mathbb{1}(y < \theta^* + \varepsilon \wedge x \leq \theta^* < (x+y)/2) dy \right) dx \\ &+ \frac{e^{-\varepsilon/b}}{2b^2} \int_x \left(\int_y \mathbb{1}(\theta^* < x < y < \theta^* + \varepsilon) dy \right) dx \\ &+ \frac{e^{-\theta^*/b}}{2b^2} \int_y e^{-y/b} \left(\int_x \mathbb{1}(\theta^* < x \wedge (x+y)/2 < \theta^* + \varepsilon \leq y) dx \right) dy \end{aligned}$$

Direct computation yields

$$\begin{aligned} \int_{x \in (\theta^* - \varepsilon, \theta^*)} e^{x/b} \left(\int_{y \in (2\theta^* - x, \theta^* + \varepsilon)} 1 dy \right) dx &= \int_{x \in (\theta^* - \varepsilon, \theta^*)} e^{x/b} (x + \varepsilon - \theta^*) dx = e^{(\theta^* - \varepsilon)/b} \int_{u \in (0, \varepsilon)} u e^{u/b} du, \\ \int_{u \in (0, \varepsilon)} u e^{u/b} du &= b \left(e^{\varepsilon/b} (\varepsilon - b) + b \right), \\ \int_x \left(\int_y \mathbb{1}(\theta^* < x < y < \theta^* + \varepsilon) dy \right) dx &= \int_{x \in (\theta^*, \theta^* + \varepsilon)} (\theta^* + \varepsilon - x) dx = \frac{\varepsilon^2}{2}, \\ \int_{y \in (\theta^* + \varepsilon, \theta^* + 2\varepsilon)} e^{-y/b} \left(\int_{x \in (\theta^*, 2\theta^* + 2\varepsilon - y)} 1 dx \right) dy &= \int_{y \in (\theta^* + \varepsilon, \theta^* + 2\varepsilon)} e^{-y/b} (\theta^* + 2\varepsilon - y) dy \\ &= e^{-(\theta^* + 2\varepsilon)/b} \int_{u \in (0, \varepsilon)} u e^{u/b} du. \end{aligned}$$

Moreover, we have

$$\begin{aligned}
 & \int_x e^{x/b} \left(\int_y \mathbb{1}(x \leq \theta^* < (x+y)/2 < \theta^* + \varepsilon \leq y) e^{-y/b} dy \right) dx \\
 &= \int_{x \in (-\infty, \theta^* - \varepsilon)} e^{x/b} \left(\int_{y \in (2\theta^* - x, 2\theta^* + 2\varepsilon - x)} e^{-y/b} dy \right) dx + \int_{x \in (\theta^* - \varepsilon, \theta^*)} e^{x/b} \left(\int_{y \in (\theta^* + \varepsilon, 2\theta^* + 2\varepsilon - x)} e^{-y/b} dy \right) dx \\
 &= b \int_{x \in (-\infty, \theta^* - \varepsilon)} \left(e^{-(2\theta^* - 2x)/b} - e^{-(2\theta^* + 2\varepsilon - 2x)/b} \right) dx + b \int_{x \in (\theta^* - \varepsilon, \theta^*)} \left(e^{-(\theta^* + \varepsilon - x)/b} - e^{-(2\theta^* + 2\varepsilon - 2x)/b} \right) dx \\
 &= b \left(\frac{b}{2} e^{-2\varepsilon/b} - \frac{b}{2} e^{-4\varepsilon/b} + b \left(e^{-\varepsilon/b} - e^{-2\varepsilon/b} \right) + \frac{b}{2} \left(e^{-4\varepsilon/b} - e^{-2\varepsilon/b} \right) \right) = b^2 \left(e^{-\varepsilon/b} - e^{-2\varepsilon/b} \right)
 \end{aligned}$$

Therefore, we have

$$\begin{aligned}
 \widetilde{\text{BC}}(\theta^*, \theta^* + \varepsilon) &= \frac{1}{2} \left(e^{-\varepsilon/b} - e^{-2\varepsilon/b} \right) + \frac{1}{2b} \left(e^{-2\varepsilon/b} + e^{-2(\theta^* + \varepsilon)/b} \right) \left(e^{\varepsilon/b}(\varepsilon - b) + b \right) + e^{-\varepsilon/b} \frac{\varepsilon^2}{4b^2} \\
 &= \frac{1}{2} \left(e^{-\varepsilon/b} - e^{-2\varepsilon/b} \right) + \frac{1}{2} \left(e^{-2\varepsilon/b} + e^{-2(\theta^* + \varepsilon)/b} \right) \left(\frac{\varepsilon}{b} e^{\varepsilon/b} - e^{\varepsilon/b} + 1 \right) + e^{-\varepsilon/b} \frac{\varepsilon^2}{4b^2} \\
 &= \frac{1}{2} \left(e^{-2(\theta^* + \varepsilon)/b} - e^{-(2\theta^* + \varepsilon)/b} \right) + \frac{1}{2} \left(e^{-\varepsilon/b} + e^{-(2\theta^* + \varepsilon)/b} \right) \frac{\varepsilon}{b} + e^{-\varepsilon/b} \frac{\varepsilon^2}{4b^2} \\
 &= \frac{1}{2} e^{-\varepsilon/b} \left(e^{-2\theta^*/b} (e^{-\varepsilon/b} - 1 + \varepsilon/b) + \frac{\varepsilon}{b} + \frac{\varepsilon^2}{2b^2} \right).
 \end{aligned}$$

Then, we can conclude that

$$\widetilde{\text{BC}}(\theta^*, \theta^* - \varepsilon) = \widetilde{\text{BC}}(\theta^* - \varepsilon, \theta^*) = \frac{1}{2} e^{-\varepsilon/b} \left(e^{-2(\theta^* - \varepsilon)/b} (e^{-\varepsilon/b} - 1 + \varepsilon/b) + \frac{\varepsilon}{b} + \frac{\varepsilon^2}{2b^2} \right).$$

Using that $f(x) = 1 - x + x^2/2 - e^{-x}$ is positive on \mathbb{R}_+ , we obtain

$$\widetilde{\text{BC}}(\theta^*, \theta) \leq \frac{|\theta^* - \theta|}{2b} \left(1 + \frac{|\theta^* - \theta|}{2b} \left(1 + e^{-2 \min\{\theta^*, \theta\}/b} \right) \right).$$

H. Supplementary Experiments

Using the same empirical setup as in Section 6, we conduct additional experiments to support our theoretical claims for other distributions (Appendix H.1), other estimators for Gaussian distributions based on \mathcal{C}_n (Appendix H.2), other convex surrogates of the 0-1 loss (Appendix H.3) or normalized/regularized versions of the logistic loss (Appendix H.4).

Reproducibility. Code for reproducing our empirical results is available at <https://github.com/tml-epfl/learning-parametric-distributions-from-samples-and-preferences>. Our code is implemented in Julia (Bezanson et al., 2017), version 1.11.5. The plots are generated with StatsPlots. The optimization problems defining some of our estimators are solved numerically with JuMP (Lubin et al., 2023), by using the Ipopt (Wächter & Biegler, 2006) and HiGHS (Huangfu & Hall, 2018) solvers. Other dependencies are listed in the Readme.md that provides detailed julia instructions to reproduce our experiments, as well as a script.sh to run them all at once. Our experiments are conducted on 12 Intel(R) Core(TM) Ultra 7 165U 4.9GHz CPU.

Gaussian distribution with known variance. For $\mathcal{F}_{\mathcal{N},1}$, the SP_{det} and SP estimators are computed with the Ipopt solver. For $\mathcal{F}_{\mathcal{N},I_d}$, the SP_{det} , SP , DP and WE estimators are computed with the Ipopt solver, and the AE estimator uses the HiGHS solver.

H.1. Accelerated Rates for Other Distributions

H.1.1. LAPLACE DISTRIBUTION WITH KNOWN SCALE

Estimators. For $\mathcal{F}_{\text{Lap},1}$ (Appendix G), we have

$$\hat{\theta}_n^{\text{SO}} = \text{median}(\{X_i\}_{i \in [n]} \cup \{Y_i\}_{i \in [n]}) \quad \text{and} \quad \mathcal{C}_n = \{\theta \mid \forall i \in [n], Z_i(|Y_i - \theta| - |X_i - \theta|) \geq 0\}.$$

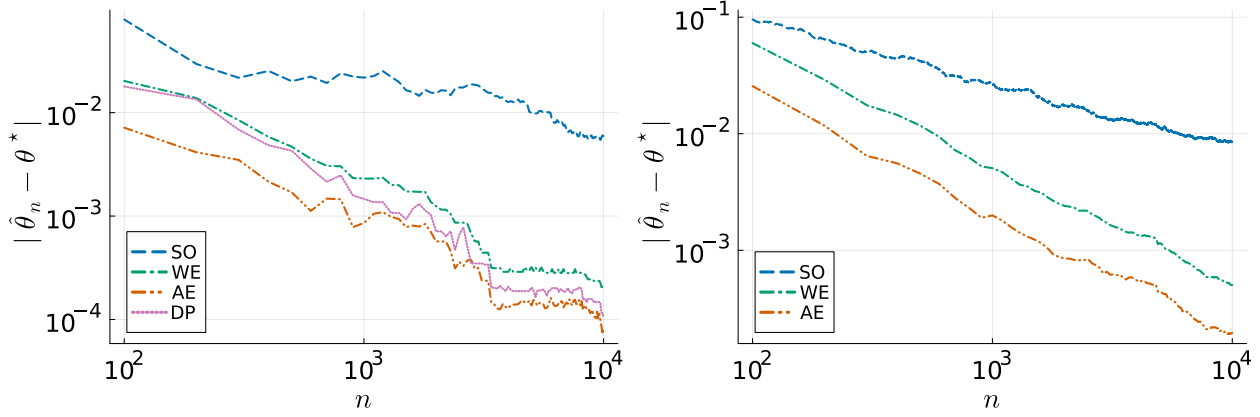


Figure 4. Estimation errors for (a) $\text{Lap}(\theta^*, 1)$ where $\theta^* \sim \mathcal{U}([1, 2])$ with $N_{\text{runs}} = 10$ and (b) $\text{Rayleigh}(\sqrt{\theta^*})$ where $\theta^* \sim \mathcal{U}([1, 2])$ with $N_{\text{runs}} = 10^2$.

The estimators based on \mathcal{C}_n are $\hat{\theta}_n^{\text{AE}} \in \mathcal{C}_n$, $\hat{\theta}_n^{\text{WE}} := \arg \max_{\theta \in \mathcal{C}_n} |\theta - \theta^*|$ and

$$\hat{\theta}_n^{\text{DP}} = \arg \max_{\theta \in \mathcal{C}_n} \sum_{i \in [n]} (|Y_i - \theta| + |X_i - \theta|) .$$

Those three estimators are computed with the Ipopt solver.

Experiments. Figure 4(a) confirms empirically the difference in estimation rate between the M-estimators (**SO MLE**)—obtaining $\mathcal{O}(1/\sqrt{n})$ —and our estimators based on \mathcal{C}_n —achieving $\mathcal{O}(1/n)$. Moreover, AE and WE perform on par with **DP MLE**.

H.1.2. RAYLEIGH DISTRIBUTION

Let $\sigma > 0$ be the scale parameter characterizing a Rayleigh distribution. In the following, let $\theta = -\frac{1}{2\sigma^2} < 0$ denote the natural parameter of a Rayleigh distribution. We have $\Theta \subseteq \mathbb{R}_+^*$, $\mathcal{X} = \mathbb{R}_+$ and $k = d = 1$. The probability density function is defined as

$$\forall x \in \mathbb{R}_+, \quad p_\theta(x) = \exp(x^2\theta + \log(x) + \log(2\theta)) .$$

Let $\theta \in \Theta$ and $u \in \{\pm 1\}$. It is direct to see that, for all $(x, y) \in \mathbb{R}_+^2$,

$$\ell_\theta(x, y) = \log \frac{p_\theta(x)}{p_\theta(y)} = (x^2 - y^2)\theta + \log(x/y) \quad \text{and} \quad \frac{d\ell_{\theta^*}}{d\theta^*}(x, y) = x^2 - y^2 = (x - y)(x + y) .$$

Therefore, we have

$$\begin{aligned} \mathcal{G}_0(\theta^*) &= \{(x, y) \in \mathbb{R}_+^2 \mid |(x^2 - y^2)\theta^* + \log(x/y)| > 0\} , \\ \mathcal{G}_1(\theta^*) &= \{(x, y) \in \mathbb{R}_+^2 \mid |x - y| > 0\} , \\ \mathcal{G}_1(\theta^*, u) &= \{(x, y) \in \mathbb{R}_+^2 \mid u((x^2 - y^2)^2\theta^* + (x^2 - y^2)\log(x/y)) < 0\} , \\ \mathcal{D}(\theta^*, \theta) &= \{(x, y) \in \mathbb{R}_+^2 \mid ((x^2 - y^2)\theta^* + \log(x/y))^2 + (x^2 - y^2)(\theta - \theta^*)((x^2 - y^2)\theta^* + \log(x/y))\} , \\ \forall (x, y) \in \mathcal{G}_1(\theta^*, u), \quad V_{\theta^*, u}(x, y) &= -u \left(\theta^* + \frac{1}{x + y} \frac{\log(x) - \log(y)}{x - y} \right) . \end{aligned}$$

Proof that $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*)) > 0$. It is direct to see that $\dim(\mathcal{G}_0(\theta^*)^c) < 2$ and $\dim(\mathcal{G}_0(\theta^*) \setminus \mathcal{G}_1(\theta^*)) < 2$. Given that $p_{\theta^*}^{\otimes 2}$ is a continuous distribution on $(\mathbb{R}_+)^2$, we obtain that $\mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_0(\theta^*)) = \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*)) = 1$.

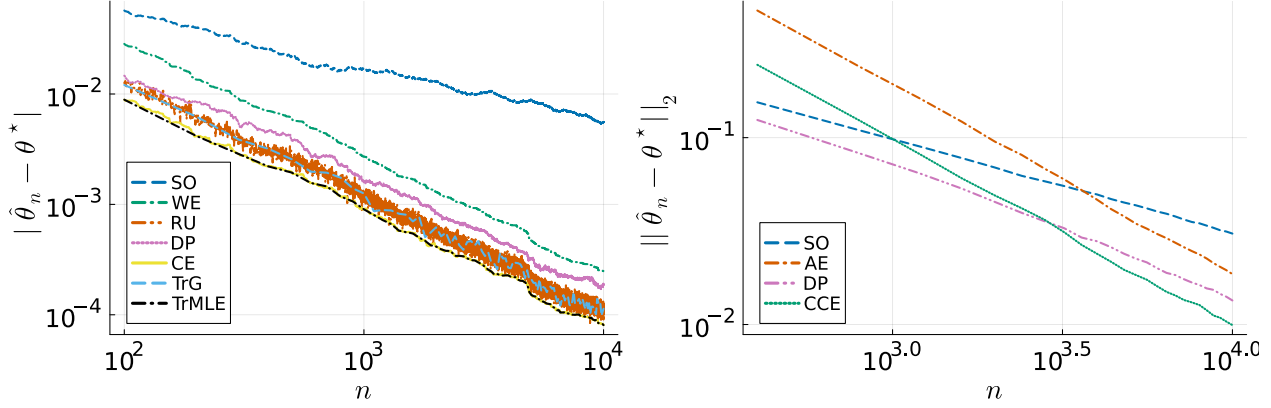


Figure 5. Estimation errors for $\mathcal{N}(\theta^*, I_d)$ where $\theta^* \sim \mathcal{U}([1, 2]^d)$ with $N_{\text{runs}} = 10^2$ for (a) $d = 1$ and (b) $d = 20$.

Proof of Assumption 4.4 and 4.5. Since $\ell_\theta(x, y) = (x^2 - y^2)\theta + \log(x/y)$ is linear in θ , we have $\mathcal{D}(\theta^*, \theta) = \tilde{\mathcal{D}}(\theta^*, \theta)$. Let $(X, Y) \sim p_{\theta^*}^{\otimes 2}$. Then, we have

$$\begin{aligned} \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_1(\theta^*, 1)) &= \mathbb{P}_{(X, Y) \sim p_{\theta^*}^{\otimes 2}}\left(\theta^* < \frac{1}{X^2} \frac{\log(Y/X)}{1 - (Y/X)^2}\right) > 0, \\ \mathbb{P}_{p_{\theta^*}^{\otimes 2}}(\mathcal{G}_0(\theta^*, -1)) &= \mathbb{P}_{(X, Y) \sim p_{\theta^*}^{\otimes 2}}\left(\theta^* > \frac{1}{X^2} \frac{\log(Y/X)}{1 - (Y/X)^2}\right) > 0. \end{aligned}$$

Estimators. We have

$$\hat{\theta}_n^{\text{SO}} = \frac{1}{4n} \sum_{i \in [n]} (X_i^2 + Y_i^2) \quad \text{and} \quad \mathcal{C}_n = \{\theta \mid \forall i \in [n], Z_i((X_i^2 - Y_i^2)\theta + \log(X_i/Y_i)) \geq 0\}.$$

The estimators based on \mathcal{C}_n are $\hat{\theta}_n^{\text{AE}} \in \mathcal{C}_n$, $\hat{\theta}_n^{\text{WE}} := \arg \max_{\theta \in \mathcal{C}_n} |\theta - \theta^*|$. Those two estimators are computed with the Ipopt solver.

Experiments. Figure 4(b) confirms empirically the difference in estimation rate between the M-estimators (SO MLE)—obtaining $\mathcal{O}(1/\sqrt{n})$ —and our estimators based on \mathcal{C}_n —achieving $\mathcal{O}(1/n)$. Moreover, AE and WE perform similarly.

H.2. Other Estimators for Gaussian Distributions

To better understand the surprising performance of the RU estimator, we consider other estimators that disentangle the effect of RU’s randomness versus its mean behavior.

Univariate Gaussian. The center estimator (CE) returns the center of the interval \mathcal{C}_n . The truncated Gaussian estimator (TrG) returns a realization from a Gaussian distribution with mean CE and variance $4/n$, which is truncated to \mathcal{C}_n . The truncated MLE (TrMLE) returns the average of the observations $(\{X_i\}_{i \in [n]} \cup \{Y_i\}_{i \in [n]}) \cap \mathcal{C}_n$.

Figure 5(a) reveals that TrG performs on par with RU, yet CE and TrMLE outperform both TrG and RU. This suggests that being far away from the boundary of \mathcal{C}_n improves performance compared to DP that lies on the boundary of \mathcal{C}_n (as observed empirically). Moreover, randomization on \mathcal{C}_n worsens performance compared to CE.

Using the derivation in the introduction on univariate Gaussian, it is coherent that CE improves on DP by a multiplicative constant: the average of those two (non-independent) random variables decreases faster. Formally, this could be proven by refining the proof of Lemma 4.6 to account for the property that $n = N_{\theta^*, -1} + N_{\theta^*, 1}$.

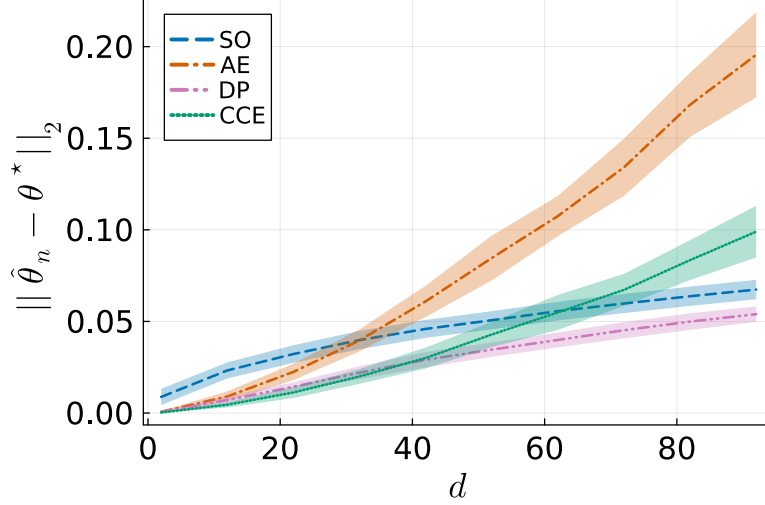


Figure 6. Estimation errors as a function of d with $\mathcal{N}(\theta^*, I_d)$ where $\theta^* \sim \mathcal{U}([1, 2]^d)$, for $n = 10^4$ and $N_{\text{runs}} = 10^2$.

Multivariate Gaussian. For $d > 1$, multiple centers exist. We use the Chebyshev center estimator (CCE) of \mathcal{C}_n .

Figures 5(b) and 6 shows that CCE outperforms AE by a constant margin. It only outperforms DP in the regime of large n compared to d and performs worse than **SO MLE** for small n . Geometrically, for small n and large d , we conjecture that the random polytope \mathcal{C}_n is more likely to be “spiky” along some directions. Due to those distant vertices, the center would become a worse estimator than DP, since the “average” is intuitively less robust to outliers. In contrast, **DP MLE** dominates **SO MLE** statistically (Lemma 4.1), hence it achieves rate $O(\sqrt{d/n})$ when n is small compared to d .

H.3. Estimators Based on Convex Surrogate of the 0-1 Loss

While **DP MLE** minimizes an objective that minimizes the 0-1 loss, **SP MLE** minimizes an objective involving the logistic loss $f_{\text{Log}}(x) = \log(1 + \exp(-x))$. As in Tang et al. (2024b), we can generalize this approach to f any convex surrogate of the 0-1 loss, see Figure 7(a). For example, we consider the Hinge loss (Hin), i.e., $f_{\text{Hin}}(x) := \max\{0, 1 - x\}$, the square loss (Squ), i.e., $f_{\text{Squ}}(x) := (1 - x)^2$, the truncated square loss (TrS), i.e., $f_{\text{TrS}}(x) := \max\{0, 1 - x\}^2$, the Savage loss (Sav), i.e., $f_{\text{Sav}}(x) := (1 + \exp(x))^{-2}$, and the exponential loss (Exp), i.e., $f_{\text{Exp}}(x) := \exp(-x)$.

Given $(X_i, Y_i, Z_i)_{i \in [n]} \sim q_{\theta^*, h_{\text{det}}}^{\otimes [n]}$ and a loss f , we consider the estimator

$$\hat{\theta}_n^f \in \arg \min_{\theta \in \Theta} \left\{ L_n^{\text{SO}}(\theta) + \sum_{i \in [n]} f(Z_i \ell_{\theta}(X_i, Y_i)) \right\}.$$

All those estimators are computed with the `Ipopt` solver.

Figure 7(b) shows that all estimators perform on par with **SP MLE**, i.e., the one based on the logistic loss.

H.4. Impact of Normalization and Regularization

The estimator defined in Appendix H.3 can be further generalized by introducing a regularization parameter $\lambda \geq 0$ and a normalization parameter $\beta > 0$, see, e.g., Gorbatovski et al. (2025). Given $(X_i, Y_i, Z_i)_{i \in [n]} \sim q_{\theta^*, h_{\text{det}}}^{\otimes [n]}$, a loss f and regularization/normalization (λ, β) , we consider the estimator

$$\hat{\theta}_n^{f, \lambda, \beta} \in \arg \min_{\theta \in \Theta} \left\{ L_n^{\text{SO}}(\theta) + \lambda \sum_{i \in [n]} f(\beta Z_i \ell_{\theta}(X_i, Y_i)) \right\}.$$

While similar modifications could be made for other losses, we focus on the logistic loss $f_{\text{Log}}(x) = \log(1 + \exp(-x))$. In particular, we recover **SP MLE** by taking $\lambda = \beta = 1$.

Figures 8(a) and (b) showcase the “mild” impact of normalization and regularization.

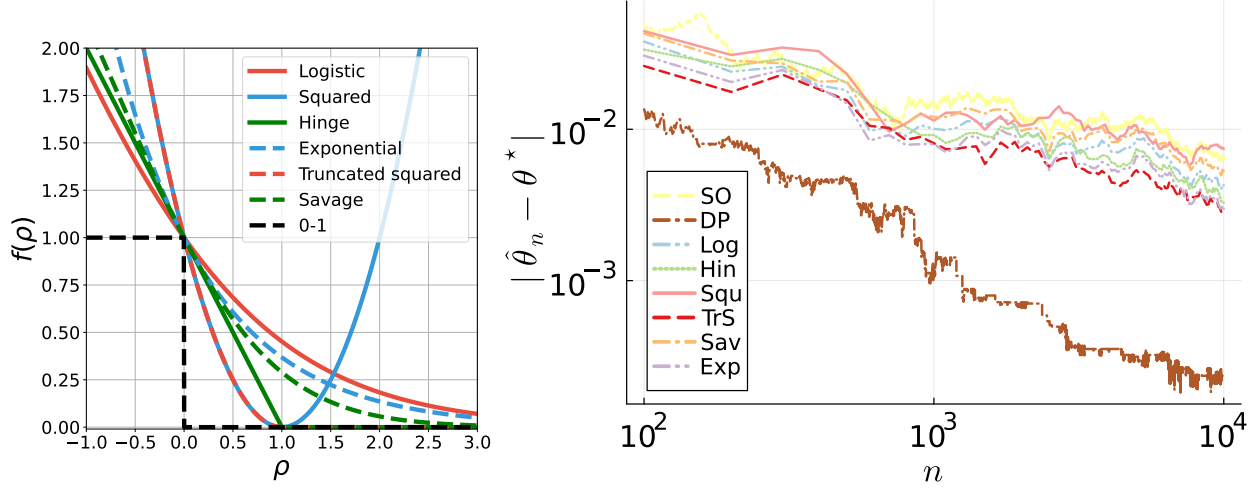


Figure 7. (a) Figure 2 in Tang et al. (2024b): notable examples of binary classification loss functions. (b) Estimation errors when minimizing the empirical losses for $\mathcal{N}(\theta^*, 1)$ where $\theta^* \sim \mathcal{U}([1, 2])$ with $N_{\text{runs}} = 10$.

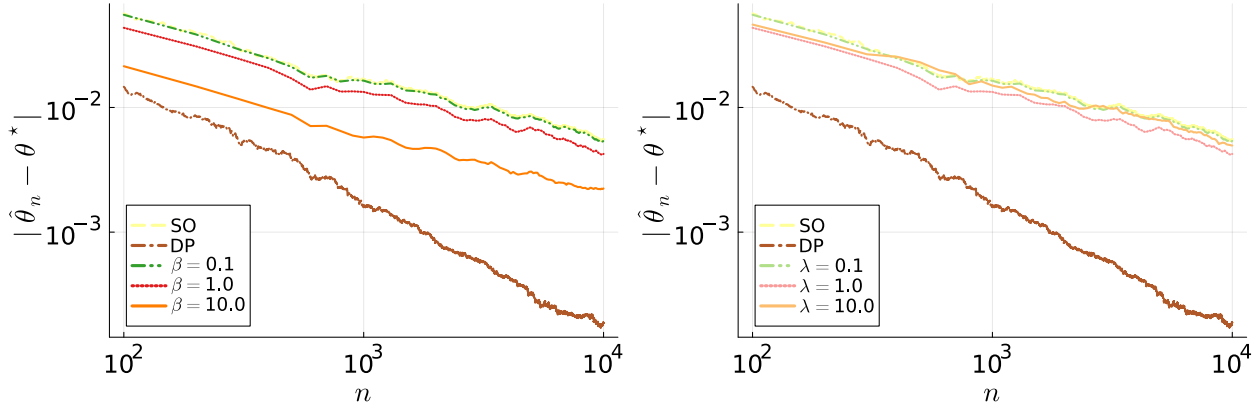


Figure 8. Estimation errors when minimizing the empirical losses for $\mathcal{N}(\theta^*, 1)$ where $\theta^* \sim \mathcal{U}([1, 2])$ with $N_{\text{runs}} = 10^2$ when (a) normalizing by β with regularization $\lambda = 1$ and (b) regularizing by λ with normalization $\beta = 1$.