# IMPROVED STOCHASTIC OPTIMIZATION OF LOGSUMEXP

**Anonymous authors** 

000

001

002003004

010 011

012

013

014

015

016

018

019

021

022

024

025

026

027 028 029

031

032

033

034

037 038

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

## **ABSTRACT**

The LogSumExp function, also known as the free energy, plays a central role in many important optimization problems, including entropy-regularized optimal transport and distributionally robust optimization (DRO). It is also the dual to the Kullback-Leibler (KL) divergence, which is widely used in machine learning. In practice, when the number of exponential terms inside the logarithm is large or infinite, optimization becomes challenging since computing the gradient requires differentiating every term. Previous approaches that replace the full sum with a small batch introduce significant bias. We propose a novel approximation to Log-SumExp that can be efficiently optimized using stochastic gradient methods. This approximation is rooted in a sound modification of the KL divergence in the dual, resulting in a new f-divergence called the safe KL divergence. The accuracy of the approximation is controlled by a tunable parameter and can be made arbitrarily small. Like the LogSumExp, our approximation preserves convexity. Moreover, when applied to an L-smooth function bounded from below, the smoothness constant of the resulting objective scales linearly with L. Experiments in DRO and continuous optimal transport demonstrate the advantages of our approach over state-of-the-art baselines and the effective treatment of numerical issues associated with the standard LogSumExp and KL.

#### 1 Introduction

Optimization problems arising in various fields involve the LogSumExp function, or, more generally, the log-partition functional

$$F(\varphi; \mu) := \ln \int e^{\varphi(x)} \, \mathrm{d}\mu(x) \in (-\infty, \infty]$$
 (1)

mapping a measurable function  $\varphi$  to  $(-\infty, \infty]$  based on a probability measure  $\mu$ . The goal in such optimization problems is to minimize an objective involving F w.r.t.  $\varphi$  over some class.

LogSumExp function appears commonly in optimization objectives, e.g., multiclass classification with softmax probabilities (Bishop & Nasrabadi, 2006), semi-dual formulation of entropy-regularized optimal transport (OT) (Peyré & Cuturi, 2019; Genevay et al., 2016), minimax problems (Pee & Royset, 2011), distributionally robust optimization (DRO) (Hu & Hong; Ben-Tal et al., 2013; Kuhn et al., 2024), maximum likelihood estimation (MLE) for exponential families and graphical models (Wainwright et al., 2008), variational Bayesian methods (Khan & Nielsen, 2018; Khan & Rue, 2023), information geometry (Amari & Nagaoka, 2000), KL-regularized Markov decision processes (Tiapkin et al., 2024). Such optimization problems are characterized by two challenges. First, the decision variable  $\varphi$  (or a parameter  $\theta$  defining  $\varphi$ ) often has large or infinite dimension. Second, the support of the measure  $\mu$  can be large or even infinite. The first challenge is usually addressed by the use of first-order methods, especially stochastic gradient descent (SGD), with cheap iterations. Unfortunately, there is no universal way to tackle the second challenge. If the number of exponential terms under the logarithm is large or infinite, i.e., the Monte Carlo estimate  $\log \sum_{i=1}^N e^{\varphi(x)}$  for a large N, the gradient computation requires differentiating each term, and, to the best of our knowledge, no cheap unbiased stochastic gradient have been proposed.

In the current work, we propose a general-purpose approach to tackle both mentioned challenges. To that end, we use a SoftPlus approximation of  $F(\varphi; \mu)$  that allows using stochastic gradient methods

while remaining close to the original objective. We start with a variational formulation analogous to the one in the Gibbs principle, but with the KL-divergence replaced with another f-divergence – the safe KL divergence. The corresponding variational problem can be of interest itself, as it possesses some properties which can be beneficial compared to the KL penalty — e.g., uniform density bound. Moreover, it can be also viewed as an approximation of a conditional value at risk functional (CVaR). In fact, the same functional (with different parameters) appeared in Soma & Yoshida (2020) in the context of smooth CVaR approximation. Thus, we demonstrate that it generates a family of problems including CVaR and LogSumExp minimization as limit cases.

**Related works.** We are not aware of any universal way to efficiently treat optimization objectives involving the LogSumExp functional (1), especially in infinite-dimensional settings. This functional appeared in different applications and was treated on a case-by-case basis. Bouchard (2007) studies three upper bounds on LogSumExp for approximate Bayesian inference. One of them is a particular case of the approximation proposed in the present work. Titsias (2016) constructs a bound on softmax probabilities and shows that it leads to a bound on LogSumExp in the context of multiclass classification. Nielsen & Sun (2016) approximate LogSumExp in the context of estimating divergences between mixture models. Their approximation combines LogSumExp bounds based on min and max. Hu & Hong and, susequently, Levy et al. (2020) study DRO problems with f-divergences. They propose a batch-based approximation. When the ambiguity set is the unit simplex, and KL divergence penalty is used, the original objective is the LogSumExp of the losses over the entire dataset, while the approximation replaces it with the average of LogSumExp terms computed on individual batches. This approximation introduces bias, which can only be reduced by using large batch sizes.

## Contributions. Our main contributions are as follows:

- 1. We introduce a general-purpose and computationally efficient approach for handling the Log-SumExp function in large-scale optimization problems by proposing a novel relaxation of this function. The proposed relaxation preserves key properties of the original LogSum-Exp function, such as convexity and smoothness, and enables the use of stochastic gradient methods for machine learning tasks. Furthermore, our method only requires a simple and tunable scalar parameter, allowing the relaxation to be made arbitrarily close to the original LogSumExp objective as desired.
- 2. We provide the theoretical backbone of this approximation, demonstrating that it is due to a modified version of the KL-divergence in the dual formulation. We termed the resulting *f*-divergence the (*Overflow-*)*Safe KL* divergence. It can be applied to various applications where KL-divergence is used.
- 3. We empirically demonstrate the effectiveness of our approach on tasks, including computing continuous entropy-regularized OT and various DRO formulations. Our method outperforms existing state-of-the-art baselines in these applications and circumvents the overflow issue (Remark 3.1). It can also be combined with existing techniques. Therefore, it serves as a versatile tool for solving large-scale optimization problems.
- 4. Additionally, we provide insights into a few remarkable connections between the proposed approximation and existing notions such as the conditional value-at-risk.

**Notation.** Given  $a, a_1, \ldots, a_n \in \mathbb{R}$ , we define  $\operatorname{LogSumExp}(a_1, \ldots, a_n) \coloneqq \operatorname{log}(\sum_{i=1}^n e^{a_i})$  and  $\operatorname{SoftPlus}(a) \coloneqq \operatorname{log}(1 + e^a)$ . Given a measurable space  $\mathcal{X}$ , by  $\mathcal{P}(\mathcal{X})$  we denote the space of probability measures on  $\mathcal{X}$ , and by  $\mathcal{C}(\mathcal{X})$  the space of continuous functions on  $\mathcal{X}$ . Let  $\mu, \nu \in \mathcal{P}(\mathcal{X})$ . Define Kullback–Leibler (KL) divergence as

$$D_{KL}(\mu,\nu) \coloneqq \begin{cases} \int_{\mathcal{X}} \log \frac{\mathrm{d}\mu}{\mathrm{d}\nu}(x) \, \mathrm{d}\mu(x) & \mu \ll \nu, \\ +\infty & \text{otherwise}, \end{cases}$$

where  $\log$  is the natural logarithm and  $\mu \ll \nu$  denotes that  $\mu$  is absolutely continuous w.r.t.  $\nu$ .

## 2 SOFTPLUS APPROXIMATION OF LOG-PARTITION FUNCTION

In this section, we present our approximation to the log-partition function (1) and describe its theoretical properties. Recall that by the Gibbs variational principle

$$F(\varphi; \mu) = \sup \left\{ \int \varphi(x) \, d\nu(x) - D_{KL}(\nu, \mu) : \nu \in \mathcal{P}(\mathcal{X}), \int |\varphi(x)| \, d\nu(x) < \infty \right\}$$
 (2)

with the minimum attained at the Gibbs measure  $d\nu^*(x) = e^{\varphi(x) - F(\varphi;\mu)} d\mu(x)$ , once  $F(\varphi;\mu) < \infty$  (Polyanskiy & Wu, 2025, Proposition 4.7).

We are going to construct an approximation of F with better regularity properties by changing  $D_{KL}$  to another f-divergence. Specifically, for any  $0 < \rho < 1$ , let us define the following.

**Definition 2.1** (Safe KL entropy). We define the safe KL entropy generator  $f_{\rho} \colon [0, \infty) \to \mathbb{R}$  by

$$f_{\rho}(t) := \begin{cases} t \log t + 1 + \frac{1 - \rho t}{\rho} \log(1 - \rho t), & 0 \le t \le \frac{1}{\rho}, \\ +\infty, & otherwise. \end{cases}$$
(3)

The resulting  $f_{\rho}$ -divergence, which we refer to as the safe KL divergence, is given by

$$D_{\rho}(\nu,\mu) := \begin{cases} \int f_{\rho} \left( \frac{d\nu}{d\mu}(x) \right) d\mu(x), & \nu \ll \mu, \\ +\infty, & \text{otherwise.} \end{cases}$$
 (4)

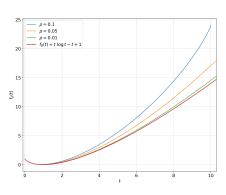


Figure 1:  $f_{\rho}(t)$  for different values of  $\rho$ .

It is easy to see that  $f_{\rho}(t) \to f_0(t) \coloneqq t \log t + 1 - t$  as  $\rho \to 0$ . Since  $f_0$  induces the standard KL-divergence,  $D_{\rho}$  is its approximation with accuracy regulated by the parameter  $\rho$ .

Using the variational representation, we define

$$F_{\rho}(\varphi; \mu) := \sup \left\{ \int \varphi(x) \, \mathrm{d}\nu(x) - D_{\rho}(\nu, \mu) : \right.$$

$$\nu \in \mathcal{P}(\mathcal{X}), \int |\varphi(x)| \, \mathrm{d}\nu(x) < \infty \right\}. \tag{5}$$

(i.e.,  $F_{\rho}(\cdot;\mu)$  is the convex conjugate of  $D_{\rho}(\cdot,\mu)$ ). Note that the last term in  $f_{\rho}$  prevents the density  $\frac{d\nu}{d\mu}$  from being too large. In particular, it can not be greater than  $\frac{1}{\rho}$ . This can make the safe KL divergence a reasonable choice for unbalanced OT or DRO, as it imposes a hard constraint

on the reweighting unlike the standard  $D_{KL}$ . Moreover, it can also be used instead of the entropy penalization in regularized OT (cf. capacity constrained transport in (Benamou et al., 2015, section 5.2)).

Again, by the convex duality and the variational principle (see Birrell et al., 2022, Theorem 6), we state the following properties.

**Lemma 2.2.** The functional  $F_o$  defined by (5) has an equivalent variational representation

$$F_{\rho}(\varphi; \mu) = \inf_{\alpha \in \mathbb{R}} \alpha + \int f_{\rho}^{*}(\varphi(x) - \alpha) d\mu(x).$$

It is straightforward to check the following.

**Lemma 2.3.** The conjugate function to  $f_{\rho}$  is a rescaled SoftPlus, specifically,

$$f_{\rho}^{*}(s) := \sup_{t \in \mathbb{R}_{+}} st - f_{\rho}(t) = \frac{1}{\rho} \log (1 + \rho e^{s}) - 1.$$

Therefore, we obtain

$$F_{\rho}(\varphi; \mu) = \inf_{\alpha \in \mathbb{R}} \alpha - 1 + \frac{1}{\rho} \int \log\left(1 + \rho e^{\varphi(x) - \alpha}\right) d\mu(x). \tag{6}$$

In essence, we have replaced the exponential function with a rescaled SoftPlus. Furthermore, it is easy to see that the optimal  $\alpha^*$  satisfies

$$\int \frac{e^{\varphi(x)-\alpha^*}}{1+\rho e^{\varphi(x)-\alpha^*}} \,\mathrm{d}\mu(x) = 1,\tag{7}$$

in particular,  $\alpha^* < F(\varphi; \mu)$ . Moreover, the maximum in (5) is attained at  $d\nu_\rho^*(s) = \frac{e^{\varphi(x) - \alpha^*}}{1 + \rho e^{\varphi(x) - \alpha^*}} \, \mathrm{d}\mu(x)$ . Note that  $0 < \frac{d\nu_\rho^*(x)}{d\mu(x)} < \frac{1}{\rho}$ , which is due to the fact that the derivative of  $t \log t$  explodes at 0, preventing reaching the constraint.

The next proposition (proved in Appendix A) ensures that  $F_{\rho}$  is a valid approximation of F.

**Proposition 2.4.** Let  $\mu \in \mathcal{P}(\mathcal{X})$  and  $\varphi$  be a measurable function on  $\mathcal{X}$ .

- (i) For all  $0 < \rho \le \rho' < 1$ , it holds  $F_{\rho'}(\varphi; \mu) \le F_{\rho}(\varphi; \mu)$ .
- (ii) As  $\rho \to 0+$ ,  $F_{\rho}(\varphi; \mu) \to F_0(\varphi; \mu) := F(\varphi; \mu)$ .
- (iii) If  $F(2\varphi;\mu) < \infty$ , then for all  $0 < \rho \le \frac{1}{4}e^{2F(\varphi;\mu) F(2\varphi;\mu)}$

$$F_{\rho}(\varphi;\mu) \ge F(\varphi;\mu) + \frac{\rho}{2} - 4\rho e^{F(2\varphi;\mu) - 2F(\varphi;\mu)}.$$
 (8)

(iv) If 
$$\varphi(x) \leq M$$
 for all  $x \in \mathcal{X}$ , then  $F_{\rho}(\varphi; \mu) \geq F(\varphi; \mu) - \rho e^{M - F(\varphi; \mu)}$  for  $\rho \in (0, e^{F(\varphi; \mu) - M})$ .

In particular, (i) and (iii) show that  $F_{\rho} - O(\rho) \leq F \leq F_{\rho}$ , and thus the parameter  $\rho$  allows one to control the approximation accuracy. In the case of LogSumExp, the above proposition yields the following simple bounds.

**Corollary 2.5.** Let  $a_1, \ldots, a_n \in \mathbb{R}$ . Then for any  $0 < \rho < 1$ 

$$\operatorname{LogSumExp}(a_1, \dots, a_n) - \rho \leq \inf_{\alpha \in \mathbb{R}} \alpha - 1 + \frac{1}{\rho} \sum_{i=1}^n \log(1 + \rho e^{a_i - \alpha}) \leq \operatorname{LogSumExp}(a_1, \dots, a_n).$$

For  $\rho = 1$  our approximation coincides with Bouchard's bound for LogSumExp (Bouchard, 2007).

#### 2.1 LINKS TO CVAR

Recall that the conditional value at risk (CVaR) w.r.t. a probability measure  $\mu \in \mathcal{P}(\mathcal{X})$  at level  $\rho \in (0, 1)$ , associated with a function  $\varphi$ , can be defined (in the case of continuous distribution) as

$$\mathrm{CVaR}_{\rho}(\varphi;\mu) \coloneqq \mathbb{E}_{X \sim \mu} \left[ \varphi(X) | \varphi(X) \ge Q_{1-\rho} \right] = \frac{1}{\rho} \int_{\varphi(x) \ge Q_{1-\rho}} \varphi(x) \, \mathrm{d}\mu(x),$$

where  $Q_{1-\rho}$  is the  $(1-\rho)$ -quantile of  $\varphi(X)$ ,  $X \sim \mu$  (Rockafellar et al., 2000). Moreover, by Theorem 1 in Rockafellar et al. (2000) CVaR also has the following variational formulation:

$$CVaR_{\rho}(\varphi; \mu) = \inf_{\alpha \in \mathbb{R}} \alpha + \frac{1}{\rho} \int (\varphi(x) - \alpha)_{+} d\mu(x).$$
 (9)

Remarkably, in Soma & Yoshida (2020) the authors obtained a smooth approximation to CVaR which, up to an additive constant, has the same form as  $F_{\rho}$ . However, they considered the approximation w.r.t. a different parameter—a "temperature" inside SoftPlus. Finally, Levy et al. (2020) proposed another similar smoothed version of CVaR (KL-regularized CVaR) in the context of DRO. For our approximation, we obtain the following bounds.

**Proposition 2.6.** For all  $0 < \rho < 1$  and  $\lambda > 0$ 

$$CVaR_{\rho}(\varphi;\mu) + \lambda(\log \rho - 1) \le \lambda F_{\rho}(\varphi/\lambda;\mu) \le CVaR_{\rho}(\varphi;\mu) + \lambda\left(\log \rho - 1 + \frac{1}{\rho}\right).$$
 (10)

#### 2.2 The case of parametric models

In some applications, the function  $\varphi$  is defined as the parametric loss function  $L(x,\theta)$  and the goal is to minimize objective involving (1) w.r.t. parameter  $\theta$  to find the best model from the parametric family. In this section, we study our approximation to (1) in this parametric setting. Fix some closed parameter set  $\Theta \subset \mathbb{R}^d$  and a loss function  $L \colon \mathcal{X} \times \Theta \to \mathbb{R}$ . Combining our approximation (6) and the minimization w.r.t. parameter  $\theta$ , we obtain the following minimization problem (note that we shifted  $\alpha$  by  $\log \rho$  compared to (6))

$$\min_{\theta \in \Theta, \alpha \in \mathbb{R}} G_{\rho}(\theta, \alpha) := \alpha + \log \rho - 1 + \frac{1}{\rho} \int \log \left( 1 + e^{L(x, \theta) - \alpha} \right) d\mu(x).$$

Clearly,  $G_{\rho}$  is convex in  $\alpha$ . Moreover, if L is convex in  $\theta$  for  $\mu$ -a.e. x, then  $G_{\rho}$  is jointly convex, meaning that our approximation preserves convexity.

Note that  $f_{\rho}(t) = \frac{1}{\rho} \left( (\rho t) \log(\rho t) + (1 - \rho t) \log(1 - \rho t) \right) + 1 - t \log \rho$ . Thus, unlike the KL entropy function  $t \log t + 1 - t$ ,  $f_{\rho}$  possesses the following favorable properties:

**Lemma 2.7.** The entropy function  $f_{\rho}$  is  $\rho$ -strongly convex. Its conjugate function  $f_{\rho}^*$  is  $\frac{1}{\rho}$ -smooth.

The above properties are important from the computational optimization point of view. Recall that  $\frac{d}{dt}\log(1+e^t)=\frac{e^t}{1+e^t}=:\sigma(t)$ . Thus, we immediately obtain the following formulas for the gradient:

$$\nabla_{\theta} G_{\rho}(\theta, \alpha) = \frac{1}{\rho} \int \sigma \left( L(x, \theta) - \alpha \right) \nabla_{\theta} L(x, \theta) \, \mathrm{d}\mu(x),$$
$$\partial_{\alpha} G_{\rho}(\theta, \alpha) = 1 - \frac{1}{\rho} \int \sigma \left( L(x, \theta) - \alpha \right) \, \mathrm{d}\mu(x).$$

This yields, in particular, that the variance of the (naïve) stochastic gradient is bounded by  $\frac{1}{\rho}$  and the second moment of  $\nabla_{\theta}L(X,\theta)$ ,  $X \sim \mu$ . In the same way one can calculate the Hessian of  $G_{\rho}$ , see Appendix C. By Proposition C.2, if  $L(x,\theta)$  is bounded from below, then  $G_{\rho}$  is smooth on  $\Theta \times (-\infty, a]$  for any  $a \in \mathbb{R}$  meaning that our approximation preserves smoothness of the loss L.

## 3 APPLICATIONS

In this section we consider several particular applications involving the objective (1) and show numerically, that our general-purpose approach based on approximation (6) leads to better performance of SGD-type algorithms than the baseline algorithms designed specifically for these applications. Source code for all experiments can be found in supplementary material.

## 3.1 CONTINUOUS ENTROPY-REGULARIZED OT

The classical optimal transport (Monge–Kantorovich) problem consists in finding a coupling of two probability measures  $\mu, \nu \in \mathcal{P}(\mathcal{X})$  which minimizes the integral of a given measurable cost function  $c \colon \mathcal{X} \times \mathcal{X} \to \mathbb{R}_+$  (e.g., a distance), i.e.,  $W(\mu, \nu) \coloneqq \inf_{\pi \in \Pi(\nu, \mu)} \int c(x, z) \, \mathrm{d}\pi(x, z)$ , where  $\Pi(\mu, \nu) \subset \mathcal{P}(\mathcal{X} \times \mathcal{X})$  is the set of couplings (transport plans) of  $\mu$  and  $\nu$  (see Kantorovich, 1942; Villani, 2008; Santambrogio, 2015). For simplicity of demonstration, we assume that the measures are defined on the same space  $\mathcal{X}$ , but the results extend trivially to the case of two different spaces. Following Cuturi (2013), we consider entropy-regularized optimal transport (eOT) problem:

$$\min_{\pi \in \Pi(\nu,\mu)} \int c(x,z) \, \mathrm{d}\pi(x,z) + \varepsilon D_{KL}(\pi,\nu \otimes \mu)$$
 (11)

where  $\nu \otimes \mu$  is the product measure. It is known that eOT admits the following dual and semi-dual formulations (see, e.g., Genevay et al. (2016)):

$$W_{\varepsilon}(\mu,\nu) = \underbrace{\max_{u,v \in \mathcal{C}(\mathcal{X})} \iint f_{\varepsilon}(x,y,u,v) \, \mathrm{d}\mu(x) \, \mathrm{d}\nu(y)}_{\text{dual}} = \underbrace{\max_{v \in \mathcal{C}(\mathcal{X})} \int h_{\varepsilon}(x,v) \, \mathrm{d}\mu(x)}_{\text{semi-dual}},$$

where

$$f_{\varepsilon}(x, y, u, v) := u(x) + v(y) - \varepsilon \exp\left(\frac{u(x) + v(y) - c(x, y)}{\varepsilon}\right),$$
 (12)

$$h_{\varepsilon}(x,v) := \int v(y) d\nu(y) - \varepsilon \log \left( \int \exp \left( \frac{v(y) - c(x,y)}{\varepsilon} \right) d\nu(y) \right) - \varepsilon, \tag{13}$$

and  $\varepsilon > 0$  is the regularization coefficient. Genevay et al. (2016) consider a reproducing kernel Hilbert space (RKHS)  $\mathcal{H}$  defined on  $\mathcal{X}$ , with a kernel  $\kappa$ , and apply SGD to solve the dual problem.

**Remark 3.1** (The overflow issue). The main drawback of this approach is the presence of the exponent in the dual objective (and consequently in the SGD updates, see Appendix B.1). Specifically, exponents are prone to floating-point exceptions (Goldberg, 1991), especially if the regularization parameter  $\varepsilon$  is relatively small, which is often the case. For example, if  $\varepsilon = 0.01$  and  $z \ge 7.1$ , then  $e^{z/\varepsilon}$  exceeds the representable range of a double-precision (float64) floating-point number — an overflow happens. When single precision (float32) is used, an overflow happens even for  $z \ge 0.89$ .

272

273

274 275

281

282

283

284 285

287

288

289

290 291 292

293

295

296

297

298

299 300 301

302

303

304 305 306

307

308 309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

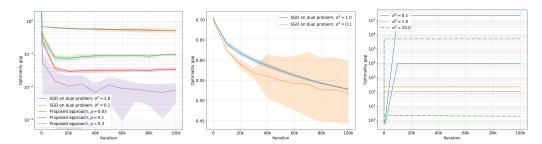


Figure 2: Left: convergence of kernel SGD applied to the dual objective (12) (blue and orange) and approximate semi-dual problem (14) (green, red and purple). Solid lines show average optimality gap across 20 runs, shaded regions indicate  $\pm$  one standard deviation. Y-axis uses logarithmic scale. Middle: a zoomed-in view of blue and orange curves from the plot on the left. Right: examples of divergent optimality gap curves obtained by running the baseline approach with the stepsize parameter  $C=10^{-2}$ .

Our approach. If we consider instead the semi-dual formulation and use the approximation (6), we arrive at the problem

$$\max_{\substack{v \in \mathcal{C}(\mathcal{X})\\ \alpha \in \mathbb{R}}} \iint \tilde{h}_{\varepsilon}(x, y, v, \alpha) \, d\mu(x) \, d\nu(y) \tag{14}$$
with  $\tilde{h}_{\varepsilon}(x, y, v, \alpha) \coloneqq v(y) - \alpha - \frac{\varepsilon}{\rho} \log \left(1 + \rho e^{(v(y) - c(x, y) - \alpha)/\varepsilon}\right) - \varepsilon,$ 

with 
$$\tilde{h}_{\varepsilon}(x, y, v, \alpha) := v(y) - \alpha - \frac{\varepsilon}{\rho} \log(1 + \rho e^{(v(y) - c(x, y) - \alpha)/\varepsilon}) - \varepsilon,$$
 (15)

which can also be solved by SGD. Analytic form of SGD iterates for objectives (12) and (15) can be found in Appendix B.1. One can show, in the same way as in Genevay et al. (2016), that this corresponds to the regularized OT problem (11) with Safe KL divergence  $D_{\rho}$  rather than the usual KL, i.e.

$$\min_{\pi \in \Pi(\nu,\mu)} \int c(x,z) \, \mathrm{d}\pi(x,z) + \varepsilon D_{\rho}(\pi,\nu \otimes \mu).$$

Note that this problem, in turn, can be viewed as a combination of the entropy-regularized and the capacity-constrained optimal transport. For  $\rho > 0$ , this approach is much more stable than the previous one when used in SGD. We illustrate this in the following experiments.

**Experiments.** Consider a setup analogous to the one described in Section 5 of Genevay et al. (2016). Specifically,  $\mu$  is a 1D Gaussian, and  $\nu$  is a mixture of two Gaussians (see Appendix B for a plot of densities). Gaussian kernel  $\kappa(x,x') = \exp\left(-\frac{\|x-x'\|^2}{\sigma^2}\right)$  with a bandwidth hyperparameter  $\sigma^2 > 0$  is used. The regularization coefficient is set to  $\varepsilon = 0.01$ . We consider kernel SGD applied to the dual objective (12) as a baseline approach (Genevay et al., 2016). We compare it to the proposed approach, namely, kernel SGD applied to the approximate semi-dual problem (14). For details on how the optimality gap is estimated, see Appendix B.

When applying kernel SGD to the dual and approximate semi-dual formulations, we consider hyperparameters  $\sigma^2 \in \{0.1, 1, 10\}$  (kernel bandwidth),  $C \in \{10^{-4}, 10^{-3}, \dots, 10\}$  (stepsize parameter), and  $\rho \in \{0.03, 0.1, 0.3\}$  (approximation accuracy). Double floating-point precision is used. In the experiment, the proposed approach works best with  $\sigma^2 = 10$ , and C = 1 for  $\rho \in \{0.03, 0.1\}$ , C=10 for  $\rho=0.3$ . Baseline works best with  $\sigma^2\in\{0.1,1\}$  and  $C=10^{-3}$ . Figure 2 (left) shows performance of the two approaches. For clarity, we provide a zoomed-in view of the curves generated by the baseline in the middle. As seen from the figures, the baseline is extremely slow, which happens due to the small stepsize. Larger values of C lead to numerical instabilities as illustrated by the plot on the right. Apparently, exponent causes a large magnitude of the gradient at a certain step, which brings an iterate to a region where it stagnates. On the contrary, our approximate semidual formulation permits larger stepsizes, which results in faster convergence. Indeed, the method usually achieves a relatively low optimality gap in about  $2 \cdot 10^4$  iterations, and plateaus after that.

#### 3.2 DISTRIBUTIONALLY ROBUST OPTIMIZATION WITH KL DIVERGENCE

One of the approaches to training a model that is robust to data distribution shifts and noisy observations is called Distributionally Robust Optimization (DRO) (Kuhn et al., 2024). In contrast to the standard Empirical Risk Minimization (ERM) approach, which minimizes the average loss on the training sample, DRO minimizes the risk for the worst-case distribution among those close to a reference measure (e.g., empirical distribution). A prominent example is KL divergence DRO (Hu & Hong), which is formulated as the saddle-point problem

$$\min_{\theta \in \Theta} \max_{p \in \Delta^n} \sum_{i=1}^n p_i \ell_i(\theta) - \lambda D_{KL}(p, \hat{p}), \tag{16}$$

where  $\theta \in \Theta$  is the model parameters,  $\ell_i(\theta)$  is the respective loss on the i-th training example,  $\Delta^n$  is the unit simplex in  $\mathbb{R}^n$ ,  $\hat{p} \in \Delta^n$  is the weight vector defining the empirical distribution (typically  $\hat{p} = \frac{1}{n}\mathbf{1}$ ), and  $D_{KL}$  is the Kullback–Leibler divergence which discourages distributions that are too far from the empirical one,  $\lambda>0$  is the penalty coefficient. For fixed  $\theta$ , the solution of the maximization problem is given by  $p_i^*(\theta)\coloneqq \frac{e^{\ell_i(\theta)/\lambda}}{\sum_j e^{\ell_j(\theta)/\lambda}}$ , which reduces the problem to

$$\min_{\theta \in \Theta} \mathcal{L}(\theta) := \lambda \log \left( \frac{1}{n} \sum_{i=1}^{n} e^{\ell_i(\theta)/\lambda} \right). \tag{17}$$

However, when n is large, computing the full gradient  $\nabla \mathcal{L}(\theta) = \sum_{i=1}^n p_i^*(\theta) \nabla \ell_i(\theta)$  becomes costly. A straightforward approach (Levy et al., 2020) is to sample a batch D, compute the respective softmax weights  $p_i^D(\theta) \coloneqq \frac{e^{\ell_i(\theta)/\lambda}}{\sum_{j\in D} e^{\ell_j(\theta)/\lambda}}$ , and define a gradient estimator by

$$\tilde{\nabla}_D \mathcal{L}(\theta) = \sum_{i \in D} p_i^D(\theta) \nabla \ell_i(\theta). \tag{18}$$

However, this introduces a bias and requires using large batch sizes to keep it sufficiently small.

**Our approach.** Instead, we propose to use the approximation (6), which results in the problem

$$\min_{\begin{subarray}{c} \theta \in \Theta \\ \alpha \in \mathbb{R} \end{subarray}} G(\theta, \alpha) \coloneqq \frac{1}{n} \sum_{i=1}^n \Bigl\{ \alpha + \frac{\lambda}{\rho} \log \bigl( 1 + \rho e^{(\ell_i(\theta) - \alpha)/\lambda} \bigr) \Bigr\}.$$

Like in the previous subsection, this can be interpreted as switching from  $D_{KL}$  penalty in (16) to Safe KL  $D_{\rho}$ . The respective gradient estimators are

$$\tilde{\nabla}_{\theta}^{D}G(\theta,\alpha) := \frac{1}{|D|} \sum_{i \in D} \sigma_{\rho} \left( \frac{\ell_{i}(\theta) - \alpha}{\lambda} \right) \nabla \ell_{i}(\theta), \quad \tilde{\nabla}_{\alpha}^{D}G(\theta,\alpha) := 1 - \frac{1}{|D|} \sum_{i \in D} \sigma_{\rho} \left( \frac{\ell_{i}(\theta) - \alpha}{\lambda} \right). \tag{19}$$

**Experiments.** Consider California housing dataset (Pace & Barry, 1997) consisting of 20640 objects represented by 8 features. Let  $\ell_i$  be the squared error of a linear model,  $\ell_i(\theta) = (y_i - \theta^\top x_i)^2$ . We use SGD with Nesterov acceleration and the gradient estimator (18) (Levy et al., 2020) as the baseline approach for solving (17). As an alternative, accelerated SGD is used with our proposed gradient estimator (19). We considered values of the stepsize  $\eta \in \{10^{-8}, 10^{-7}, \dots, 10^{-4}\}$  and approximation accuracy parameter  $\rho \in \{10^{-4}, 10^{-3}, 10^{-2}\}$ . Momentum was set to 0.9 (without tuning) and entropy penalty  $\lambda$  to 1. Least squares solution was used as a starting point for optimization

Convergence curves for both approaches corresponding to the best hyperparameters (i.e., resulting in smallest objective values) are displayed in Figure 3. Y-axis depicts the value of  $\mathcal{L}(\theta)$ , X-axis displays the number of epochs (i.e., full passes over the dataset), error bars denote  $\pm$  one standard deviation across 20 runs. We omit first few epochs to zoom in on smaller objective values. As seen from the figure, convergence of the proposed approach is consistently good across all batch sizes. As for the baseline, |D|=10 requires small stepsize  $\eta=10^{-6}$  to avoid divergence, and thus results in slow convergence. Larger batches enable the increase of  $\eta$  up to  $10^{-5}$ . The baseline still converges slower than the proposed approach with |D|=100, but almost matches its performance for |D|=1000.

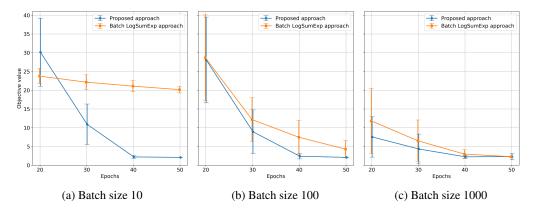


Figure 3: Convergence of SGD with Nesterov acceleration on the California housing dataset for the baseline gradient estimator (18) and the proposed estimator (19). The plots show  $\mathcal{L}(\theta)$  versus the number of epochs for the best hyperparameters for different batch sizes. Error bars indicate  $\pm$  one standard deviation over 20 runs.

#### 3.3 DISTRIBUTIONALLY ROBUST OPTIMIZATION WITH UNBALANCED OT

In the KL divergence DRO described in the previous subsection, uncertainty set is limited to distributions with the same support as the empirical measure  $\mu = \frac{1}{n} \sum_i \delta_{x_i}$ . Another popular approach, Wasserstein DRO (WDRO) (Mohajerin Esfahani & Kuhn, 2018; Sinha et al., 2020), considers the worst-case risk over shifts within a Wasserstein (OT) ball around a reference measure  $\mu$  instead of the KL-ball in (16), thus including continuous probability measures. Unfortunately, this approach is not resilient to outliers that are geometrically far from the clean distribution since OT metric is sensitive to them (Nietert et al., 2023). A natural generalization is to switch to semi-balanced OT (Liero et al., 2018; Chizat et al., 2019; Kondratyev et al., 2016), which replaces a hard constraint on one of the marginals with a mismatch penalty function, e.g.,

$$W_{\beta}(\nu,\mu) = \inf_{\substack{\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) \\ \pi_1 = \nu}} \int c(x,z) \, \mathrm{d}\pi(x,z) + \beta \, D_{KL}(\pi_2,\mu),$$

where  $\pi_1$  and  $\pi_2$  are first and second marginals of  $\pi$ , respectively,  $\beta > 0$  is the marginal penalty parameter. Intuitively, this discrepancy measure allows to ignore some points (e.g., outliers) by paying a small price for mismatch in marginals. The (penalty-form) DRO problem can be written as

$$\min_{\theta \in \Theta} \max_{\nu \in \mathcal{P}(\mathcal{X})} \int \ell(\theta, x) \, d\nu(x) - \lambda W_{\beta}(\nu, \mu),$$

where  $\lambda > 0$  is the Lagrangian penalty parameter. Using standard duality, Wang et al. (2024) showed that when  $\mu = \frac{1}{n} \sum_{i} \delta_{x_i}$  is the empirical distribution, this is equivalent to

$$\min_{\theta \in \Theta} F(\theta) := \lambda \beta \log \left( \frac{1}{n} \sum_{i=1}^{n} e^{\hat{\ell}_i(\theta)/(\lambda \beta)} \right) \quad \text{with} \quad \hat{\ell}_i(\theta) := \sup_{z \in \mathcal{X}} \{ \ell(\theta; z) - \lambda c(z, x_i) \}, \tag{20}$$

To avoid the costly gradient computation of LogSumExp, the authors drop the logarithm and use SGD to optimize the sum of exponents,

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^{n} e^{\hat{\ell}_i(\theta)/(\lambda \beta)}.$$
 (21)

The major downside of this approach is that the exponent terms have a large variance, and SGD is prone to floating-point exceptions (overflow) unless a very small stepsize is tuned, which slows down the convergence and can be time-consuming and unstable in practice.

**Our approach.** To overcome this issue, we propose leveraging the approximation (6), which leads to the problem

$$\min_{\substack{\theta \in \Theta \\ \alpha \in \mathbb{R}}} \frac{1}{n} \sum_{i=1}^{n} \left\{ \alpha + \frac{\lambda \beta}{\rho} \log \left( 1 + \rho e^{(\hat{\ell}_i(\theta) - \alpha)/(\lambda \beta)} \right) \right\},\tag{22}$$

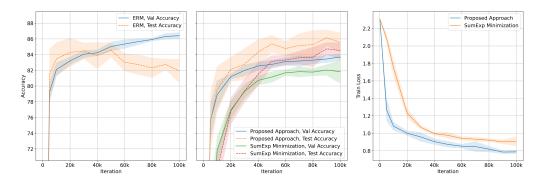


Figure 4: Performance of ERM and two DRO approaches on MNIST with noisy labels. Left: ERM accuracy on the noisy validation set vs. clean test set. Middle: validation vs. test accuracy for DRO approaches. Right: training loss  $F(\theta)$  from (20).

where  $\rho > 0$  is a parameter controlling the accuracy of the approximation. This approximation can be efficiently optimized with SGD. Note that our method can also be applied to other DRO algorithms such as Sinkhorn DRO (Wang et al., 2021), which we omit to avoid redundancy.

**Experiments.** We consider MNIST dataset (Deng, 2012) with train and validation labels corrupted by feature-dependent noise (see Algan & Ulusoy, 2020) (noise ratio 25%), and original (clean) test labels. Let  $\theta$  denote weights of a CNN with two convolutional layers (32 and 64 channels, kernel size 3, ReLU activations, and  $2\times 2$  max pooling), followed by a fully connected classifier with one hidden layer of 128 units, and let  $\ell(\theta;z)$  be its cross entropy loss on object z. In the experiment, SGD (with batch size 1) is applied to problems (21) (baseline) and (22) (proposed approach). We consider values of the stepsize  $\eta \in \{10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}\}$ . Parameters  $\lambda$  and  $\beta$  were set to 1 since smaller values required a smaller stepsize and resulted in a slow convergence, while larger values pushed the model towards fitting the noisy distribution instead of the true one. Approximation accuracy parameter  $\rho$  in (22) was set to 0.1. For the inner maximization problem in (20), just 5 iterations of Nesterov's accelerated gradient method were sufficient to reach plateau in terms of the objective value. Additionally, we used SGD for the usual empirical risk minimization (ERM) to observe the effects of conventional (non-robust) training on noisy data.

Figure 4 demonstrates the performance of different approaches with best hyperparameter  $\eta$  ( $10^{-3}$  for ERM,  $10^{-4}$  for the proposed approach, and  $10^{-5}$  for the baseline). Shaded regions indicate  $\pm$  one standard deviation across 10 runs, except that for the baseline approach (21) we excluded a single run that caused a floating-point exception. The plot on the left illustrates that ERM fits to the corrupted data well (accuracy on the noisy validation set is increasing) which results in decreasing accuracy on the clean test set. In contrast, the plot in the middle shows that an increase in validation accuracy results in the increase in test accuracy for both DRO approaches, which indicates that they are more capable at learning the underlying clean distribution. The plot on the right shows the train loss  $F(\theta)$  from (20). As seen from the figure, the proposed approach converges faster than the baseline. This is caused by the fact that the baseline requires a small stepsize to avoid overflows.

### 4 Conclusion

We introduce a novel approximation to the log partition function (and in particular, to LogSum-Exp), which arises in numerous applications across machine learning and optimization. In the dual formulation, it corresponds to the safe KL divergence. Our LogSumExp approximation preserves convexity and smoothness, and can be efficiently minimized using stochastic gradient methods. Importantly, the respective gradient estimator has controllable bias independent of batch size, in contrast to prior approaches. Our empirical results highlight the practical advantages of the proposed approximation across tasks in continuous entropy-regularized OT and DRO. An important direction for future work is to leverage the approximation for other applications, where the LogSumExp function and duality of the KL divergence play a role.

## REFERENCES

- Görkem Algan and Ilkay Ulusoy. Label noise types and their effects on deep learning. *arXiv preprint arXiv:2003.10471*, 2020.
- Shun-ichi Amari and Hiroshi Nagaoka. *Methods of Information Geometry*. American Mathematical Soc., 2000. ISBN 978-0-8218-4302-4.
  - Aharon Ben-Tal, Dick den Hertog, Anja De Waegenaere, Bertrand Melenberg, and Gijs Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59(2):341–357, February 2013. ISSN 0025-1909, 1526-5501. doi: 10.1287/mnsc.1120.1641. URL http://pubsonline.informs.org/doi/abs/10.1287/mnsc.1120.1641.
  - Jean-David Benamou, Guillaume Carlier, Marco Cuturi, Luca Nenna, and Gabriel Peyré. Iterative bregman projections for regularized transportation problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–A1138, 2015.
  - Jeremiah Birrell, Paul Dupuis, Markos A. Katsoulakis, Yannis Pantazis, and Luc Rey-Bellet. (f,gamma)-divergences: Interpolating between f-divergences and integral probability metrics. *Journal of Machine Learning Research*, 23(39):1–70, 2022. URL http://jmlr.org/papers/v23/21-0100.html.
  - Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
  - Guillaume Bouchard. Efficient bounds for the softmax function and applications to approximate inference in hybrid models. In NIPS 2007 workshop for approximate Bayesian inference in continuous/hybrid systems, volume 6, 2007.
  - Lenaic Chizat, Gabriel Peyré, Bernhard Schmitzer, and François-Xavier Vialard. Unbalanced optimal transport: Dynamic and Kantorovich formulation. *arXiv:1508.05216 [math]*, February 2019. URL http://arxiv.org/abs/1508.05216. arXiv: 1508.05216.
  - Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26, 2013.
  - Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
  - Aude Genevay, Marco Cuturi, Gabriel Peyré, and Francis Bach. Stochastic optimization for large-scale optimal transport. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems* 29, pp. 3440–3448. Curran Associates, Inc., 2016.
  - David Goldberg. What every computer scientist should know about floating-point arithmetic. *ACM computing surveys (CSUR)*, 23(1):5–48, 1991.
  - Zhaolin Hu and L Jeff Hong. Kullback-Leibler Divergence Constrained Distributionally Robust Optimization.
  - Leonid Kantorovich. On the translocation of masses. (Doklady) Acad. Sci. URSS (N.S.), 37:199–201, 1942.
  - Mohammad Emtiyaz Khan and Didrik Nielsen. Fast yet simple natural-gradient descent for variational inference in complex models. In 2018 International Symposium on Information Theory and Its Applications (ISITA), pp. 31–35. IEEE, 2018.
  - Mohammad Emtiyaz Khan and Håvard Rue. The Bayesian learning rule. (arXiv:2107.04562), June 2023.
- Stanislav Kondratyev, Léonard Monsaingeon, and Dmitry Vorotnikov. A new optimal transport distance on the space of finite Radon measures. *Advances in Differential Equations*, 21(11/12): 1117 1164, 2016. doi: 10.57262/ade/1476369298. URL https://doi.org/10.57262/ade/1476369298.

- Daniel Kuhn, Soroosh Shafiee, and Wolfram Wiesemann. Distributionally Robust Optimization. (arXiv:2411.02549), November 2024. doi: 10.48550/arXiv.2411.02549.
- Daniel Levy, Yair Carmon, John C Duchi, and Aaron Sidford. Large-scale methods for distributionally robust optimization. *Advances in neural information processing systems*, 33:8847–8860, 2020.
  - Matthias Liero, Alexander Mielke, and Giuseppe Savaré. Optimal entropy-transport problems and a new Hellinger–Kantorovich distance between positive measures. *Inventiones mathematicae*, 211(3):969–1117, March 2018. ISSN 0020-9910, 1432-1297. doi: 10.1007/s00222-017-0759-8. URL http://link.springer.com/10.1007/s00222-017-0759-8.
  - Peyman Mohajerin Esfahani and Daniel Kuhn. Data-driven distributionally robust optimization using the Wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1-2):115–166, September 2018. ISSN 0025-5610, 1436-4646. doi: 10.1007/s10107-017-1172-1.
  - Frank Nielsen and Ke Sun. Guaranteed bounds on information-theoretic measures of univariate mixtures using piecewise log-sum-exp inequalities. *Entropy*, 18(12):442, 2016.
  - Sloan Nietert, Ziv Goldfeld, and Soroosh Shafiee. Outlier-robust wasserstein dro. *Advances in Neural Information Processing Systems*, 36:62792–62820, 2023.
  - R Kelley Pace and Ronald Barry. Sparse spatial autoregressions. *Statistics & Probability Letters*, 33 (3):291–297, 1997.
  - EY Pee and Johannes O Royset. On solving large-scale finite minimax problems using exponential smoothing. *Journal of optimization theory and applications*, 148(2):390–421, 2011.
  - Gabriel Peyré and Marco Cuturi. Computational optimal transport: With applications to data science. Foundations and Trends® in Machine Learning, 11(5-6):355–607, 2019. ISSN 1935-8237. doi: 10.1561/2200000073. URL http://dx.doi.org/10.1561/2200000073. arXiv:1803.00567.
  - Yury Polyanskiy and Yihong Wu. *Information theory: From coding to learning*. Cambridge university press, 2025.
  - R Tyrrell Rockafellar, Stanislav Uryasev, et al. Optimization of conditional value-at-risk. *Journal of risk*, 2:21–42, 2000.
  - Filippo Santambrogio. *Optimal Transport for Applied Mathematicians: Calculus of Variations*, *PDEs*, *and Modeling*. Springer International Publishing, 2015. ISBN 9783319208282. doi: 10.1007/978-3-319-20828-2. URL http://dx.doi.org/10.1007/978-3-319-20828-2.
  - Aman Sinha, Hongseok Namkoong, Riccardo Volpi, and John Duchi. Certifying some distributional robustness with principled adversarial training. *arXiv:1710.10571 [cs, stat]*, May 2020.
  - Tasuku Soma and Yuichi Yoshida. Statistical learning with conditional value at risk. *arXiv preprint arXiv:2002.05826*, 2020.
  - Daniil Tiapkin, Denis Belomestny, Daniele Calandriello, Eric Moulines, Alexey Naumov, Pierre Perrault, Michal Valko, and Pierre Menard. Demonstration-regularized RL. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=lF2aip4Scn.
- Michalis K. Titsias. One-vs-each approximation to softmax for scalable estimation of probabilities. Advances in Neural Information Processing Systems, 29, 2016.
  - Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
  - Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.

Jie Wang, Rui Gao, and Yao Xie. Sinkhorn distributionally robust optimization. *arXiv* preprint arXiv:2109.11926, 2021.

Zifan Wang, Yi Shen, Michael Zavlanos, and Karl H Johansson. Outlier-robust distributionally robust optimization via unbalanced optimal transport. *Advances in Neural Information Processing Systems*, 37:52189–52214, 2024.

## A Proofs for Section 2

*Proof of Proposition 2.4.* (i,ii) Consider the function  $g(t) := \frac{\ln(1+t)}{t}$ . It is decreasing and convex on  $(0,\infty)$ ,  $g(t) \to 1$  and  $g'(t) \to -\frac{1}{2}$  as  $t \to 0+$ . Note that

$$F_{\rho}(\varphi;\mu) = \inf_{\alpha \in \mathbb{R}} \alpha - 1 + \int e^{\varphi(x) - \alpha} g\left(\rho e^{\varphi(x) - \alpha}\right) d\mu(x).$$

Then (i) follows immediately from (6) and the monotonicity of g. The monotone convergence theorem yields (ii) since

$$F(\varphi; \mu) = \inf_{\alpha \in \mathbb{R}} \alpha - 1 + \int e^{\varphi(x) - \alpha} d\mu(x).$$

Now, let us prove (iii). Consider the optimal  $\alpha_{\rho}$ , satisfying (7). By Jensen's inequality

$$\int \ln\left(1 + \rho e^{\varphi(x) - \alpha_{\rho}}\right) d\mu(x) = -\int \ln\left(1 - \frac{\rho e^{\varphi(x) - \alpha_{\rho}}}{1 + \rho e^{\varphi(x) - \alpha_{\rho}}}\right) d\mu(x)$$

$$\geq -\ln\left(1 - \int \frac{\rho e^{\varphi(x) - \alpha_{\rho}}}{1 + \rho e^{\varphi(x) - \alpha_{\rho}}} d\mu(x)\right) = -\ln(1 - \rho),$$

thus

$$F_{\rho}(\varphi;\mu) = \alpha_{\rho} - 1 + \frac{1}{\rho} \int \ln\left(1 + \rho e^{\varphi(x) - \alpha_{\rho}}\right) d\mu(x) \ge \alpha_{\rho} - 1 - \frac{\ln(1 - \rho)}{\rho} \ge \alpha_{\rho} + \frac{\rho}{2}. \quad (23)$$

It remains to get a lower bound on  $\alpha_{\rho}$ . By the monotonicity of  $\frac{t}{1+t}$  we deduce that  $\alpha_{\rho} \geq \alpha$  for any  $\alpha$  such that

$$\int \frac{\rho e^{\varphi(x) - \alpha}}{1 + \rho e^{\varphi(x) - \alpha}} \, \mathrm{d}\mu(x) \ge \rho.$$

Since  $\frac{t}{1+t} \ge t - t^2$ ,

$$\int \frac{\rho e^{\varphi(x) - \alpha}}{1 + \rho e^{\varphi(x) - \alpha}} \, \mathrm{d}\mu(x) \ge \int \left( \rho e^{\varphi(x) - \alpha} - \rho^2 e^{2\varphi(x) - 2\alpha} \right) \, \mathrm{d}\mu(x) = \rho e^{F(\varphi; \mu) - \alpha} - \rho^2 e^{F(2\varphi; \mu) - 2\alpha}.$$

Denoting  $u := e^{F(\varphi;\mu)-\alpha}$ , it is enough to find u such that

$$u - au^2 \ge 1$$
, where  $a := \rho e^{F(2\varphi;\mu) - 2F(\varphi;\mu)} \le \frac{1}{4}$ .

Thus, taking

$$u := \frac{1}{2a} \left( 1 - \sqrt{1 - 4a} \right) \le 1 + 4a,$$

we obtain

$$\alpha_{\rho} \ge F(\varphi; \mu) - \ln u \ge F(\varphi; \mu) - \ln (1 + 4a) \ge F(\varphi; \mu) - 4a.$$

Combining this with (23), we get (8).

(iv) Finally, let  $\varphi(x) \leq M$  for all  $x \in \mathcal{X}$ . Then by concavity

$$\int \ln\left(1 + \rho e^{\varphi(x) - \alpha}\right) d\mu(x) \ge \int e^{\varphi(x) - M} \ln\left(1 + \rho e^{M - \alpha}\right) d\mu(x) = e^{F(\varphi; \mu) - M} \ln\left(1 + \rho e^{M - \alpha}\right)$$

for all  $\alpha \in \mathbb{R}$ . Therefore,

$$F_{\rho}(\varphi; \mu) \ge \min_{\alpha} \alpha - 1 + \frac{e^{F(\varphi; \mu) - M}}{\rho} \ln \left( 1 + \rho e^{M - \alpha} \right)$$

$$= F(\varphi; \mu) - 1 - \frac{1 - \rho e^{M - F(\varphi; \mu)}}{\rho e^{M - F(\varphi; \mu)}} \ln \left( 1 - \rho e^{M - F(\varphi; \mu)} \right)$$

$$\ge F(\varphi; \mu) - \rho e^{M - F(\varphi; \mu)}.$$

Here we used the inequality

$$\frac{1-t}{t}\ln(1-t) \le t - 1, \quad 0 < t < 1.$$

*Proof of Corollary 2.5.* Set  $\mu_n := \frac{1}{n} \sum_{i=1}^n \delta_{a_i} \in \mathcal{P}(\mathbb{R})$ . Then

$$LogSumExp(a_1, ..., a_n) = \ln n + \ln \left( \int x d\mu_n(x) \right) = \ln n + F(id; \mu_n)$$

and

$$\inf_{\alpha \in \mathbb{R}} \alpha - 1 + \frac{1}{\rho} \sum_{i=1}^{n} \ln(1 + \rho e^{a_i - \alpha}) = \inf_{\alpha \in \mathbb{R}} \alpha - 1 + \frac{n}{\rho} \int \ln(1 + \rho e^{x - \alpha}) \, \mathrm{d}\mu_n(x)$$

$$= \inf_{\alpha \in \mathbb{R}} \alpha - 1 + \frac{n}{\rho} \int \ln\left(1 + \frac{\rho}{n} e^{x - \alpha + \ln n}\right) \, \mathrm{d}\mu_n(x)$$

$$= \ln n + F_{\rho/n}(id; \mu_n).$$

Since

$$e^{F(id;\mu_n) - \max_i a_i} = \frac{\sum_{i=1}^n e^{a_i}}{n \max_i e^{a_i}} \ge \frac{1}{n} > \frac{\rho}{n},$$

Proposition 2.4(i,iv) yields

$$F(id; \mu_n) - \rho \le F_{\rho/n}(id; \mu_n) \le F(id; \mu_n).$$

The claim follows.

*Proof of Proposition 2.6.* As  $\lambda \ln(1 + e^{t/\lambda}) > t_+ := \max\{0, t\}$ , we get

$$\lambda F_{\rho}(\varphi/\lambda; \mu) \ge \lambda \inf_{\alpha \in \mathbb{R}} \alpha - 1 + \frac{1}{\rho} \int \left( \ln \rho + \frac{\varphi(x)}{\lambda} - \alpha \right)_{+} d\mu(x)$$
$$= \lambda (\ln \rho - 1) + \inf_{\alpha \in \mathbb{R}} \alpha + \frac{1}{\rho} \int (\varphi(x) - \alpha)_{+} d\mu(x).$$

The infimum in the r.h.s. is the variational formula for CVaR (9), thus we get the first inequality in (10). The second inequality can be obtained in a similar way using that  $\lambda \ln(1 + e^{t/\lambda}) < t_+ + \lambda$ .

## B ADDITIONAL DETAILS ON THE EOT EXPERIMENT

## B.1 SGD ITERATES FOR THE DUAL AND APPROXIMATE SEMI-DUAL FORMULATIONS

By the property of RKHS, if  $u \in \mathcal{H}$ , then  $u(x) = \langle u, \kappa(\cdot, x) \rangle_{\mathcal{H}}$ , and the derivatives of  $f_{\varepsilon}$  are as follows:

$$\nabla_u f_{\varepsilon}(x, y, u, v) = \kappa(\cdot, x) - \exp\left(\frac{u(x) + v(y) - c(x, y)}{\varepsilon}\right) \kappa(\cdot, x),$$
$$\nabla_v f_{\varepsilon}(x, y, u, v) = \kappa(\cdot, y) - \exp\left(\frac{u(x) + v(y) - c(x, y)}{\varepsilon}\right) \kappa(\cdot, y).$$



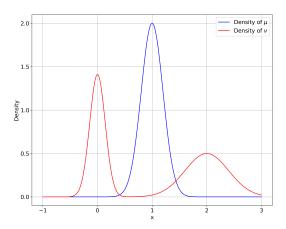


Figure 5: Densities of source and target distributions in the eOT experiment.

Thus, the SGD iterates for the dual objective (12) take the simple form

$$(u_k, v_k) = (u_0, v_0) + \sum_{i=1}^k \beta_i \left( \kappa \left( \cdot, x_i \right), \kappa \left( \cdot, y_i \right) \right)$$
(24)

with 
$$\beta_i \coloneqq \frac{C}{\sqrt{i}} \left( 1 - e^{\frac{u_{i-1}(x_i) + v_{i-1}(y_i) - c(x_i, y_i)}{\varepsilon}} \right),$$
 (25)

where  $(x_i, y_i)$  are i.i.d. samples from  $\mu \otimes \nu$ , and C > 0 is the initial stepsize. Similarly, SGD iterates for (15) are computed as follows:

$$\begin{split} v_k &= v_0 + \sum_{i=1}^k \tilde{\beta}_i \kappa\left(\cdot, y_i\right), \\ \alpha_k &= \alpha_0 - \sum_{i=1}^k \tilde{\beta}_i \quad \text{with } \beta_i \coloneqq \frac{C}{\sqrt{i}} \Big(1 - \sigma_\rho \left(\frac{u_{i-1}(x_i) + v_{i-1}(y_i) - c(x_i, y_i)}{\varepsilon}\right)\Big), \end{split}$$

where  $\sigma_{\rho}(t) \coloneqq \frac{e^t}{1+\rho e^t}$ .

#### B.2 Densities of source and target distributions

Figure 5 shows the densities of source and target distributions ( $\mu$  and  $\nu$ , respectively) used in the experiment in Section 3.1. Specifically,  $\mu$  is a 1D Gaussian, and  $\nu$  is a mixture of two Gaussians.

### B.3 COMPUTING A PROXY FOR OPTIMALITY GAP

Optimality gap in the experiment is estimated as follows:

- 1. Test sets  $\{x_i\}_{i=1}^N$  and  $\{y_i\}_{i=1}^N$  of size  $N=10^4$  are sampled from  $\mu$  and  $\nu$ . The corresponding empirical distributions are denoted  $\hat{\mu}$  and  $\hat{\nu}$ , respectively.
- 2. Similarly to Genevay et al. (2016), we obtain a proxy  $\hat{W}$  for  $W(\mu, \nu)$  by solving the semi-discrete eOT problem

$$\label{eq:local_energy} \begin{split} \max_{\mathbf{v} \in \mathbb{R}^N} \; \mathbb{E}_{X \sim \mu} \, \hat{h}_{\varepsilon}(X, \mathbf{v}) \\ \text{with } \; \hat{h}_{\varepsilon}(x, \mathbf{v}) \coloneqq \frac{1}{N} \sum_{i=1}^N v_i - \varepsilon \log \Bigl( \frac{1}{N} \sum_{i=1}^N e^{\frac{v_i - c(x, y_i)}{\varepsilon}} \Bigr) - \varepsilon, \end{split}$$

which corresponds to replacing the expectation  $\mathbb{E}_{Y \sim \nu}$  in (13) with the average over the test set  $\mathbb{E}_{Y \sim \hat{\nu}}$ . We preform 10 runs of SGD, each consisting of  $2 \cdot 10^5$  iterations, and define  $\hat{W}$  as largest achieved value on the test set, i.e., the largest  $\mathbb{E}_{X \sim \hat{\mu}} \hat{h}_{\varepsilon}(X, \mathbf{v})$ .

3. Finally, given a potential  $v \in \mathcal{C}(\mathcal{X})$ , we estimate the optimality gap as  $\hat{W} - \mathbb{E}_{X \sim \hat{\mu}} \hat{h}_{\varepsilon}(X, \mathbf{v})$ , where  $\mathbf{v} = (v(y_1), \dots, v(y_N))^{\top}$  is the evaluation of v on the test set.

## C PROPERTIES OF SOFTPLUS

Let 
$$F(x) = \log(1 + e^{f(x)})$$
, then

$$\nabla F(x) = \sigma(f(x))\nabla f(x),\tag{26}$$

$$\nabla^2 F(x) = \sigma(f(x)) \nabla^2 f(x) + \sigma(f(x)) (1 - \sigma(f(x))) \nabla f(x) \nabla f(x)^{\top}. \tag{27}$$

Suppose f(x) is L-smooth (possibly non-convex). Let us derive smoothness constant of F. We will use the following

**Lemma C.1.** Consider function  $f_a(x) = \sigma(x) + 2\sigma'(x)(x-a)$ ,  $x \ge a$  with parameter  $a \le 0$ . It holds  $f_a(x) \le 2 - \frac{a}{2}$ .

*Proof.* By the properties of the sigmoid function  $\sigma(x)$ ,  $\sigma'(x) \leq \frac{1}{4}$  and  $\sigma(x) \leq 1$ . Therefore,  $f_a(x) \leq 1 + \frac{x-a}{2}$ . If  $x \leq 2$ , the result follows. Let us now show that the derivative

$$\frac{d}{dx}f_a(x) = \sigma'(x)[3 + 2(1 - 2\sigma(x))(x - a)]$$

is negative if x > 2. Indeed, due to monotonicity of the sigmoid function  $\sigma(x)$ ,

$$\sigma(x) > \sigma(2) > 0.88 \Rightarrow 2(1 - 2\sigma(x)) < -\frac{3}{2}$$

Moreover, x-a>2, so  $3+2(1-2\sigma(x))(x-a)<0$  and  $\frac{d}{dx}f_a(x)<0$ . Therefore, if x>2, then  $f_a(x)< f_a(2)\leq 2-\frac{a}{2}$ .

**Proposition C.2.** Let  $f \in C^1(\mathbb{R}^d)$  be L-smooth and bounded from below by  $f_* \in \mathbb{R}$ , then  $F(x) = \log(1 + e^{f(x)})$  is smooth with parameter

$$\begin{cases}
\frac{4}{3}L & \text{if } f_* \ge 0, \\
\left(\frac{4}{3} - \frac{f_*}{2}\right)L & \text{if } f_* < 0.
\end{cases}$$
(28)

*Proof.* W.l.o.g., we can assume that  $f \in C^2$ . From (27) and Lemma C.1 we get

$$\|\nabla^{2} F(x)\| \leq \sigma(f(x)) \|\nabla^{2} f(x)\| + \sigma'(f(x)) \|\nabla f(x)\|^{2}$$

$$\leq L\sigma(f(x)) + 2L\sigma'(f(x))(f(x) - f_{*})$$

$$= L(\sigma(f(x)) + 2\sigma'(f(x))f(x)) - 2L\sigma'(f(x))f_{*}.$$

Analyzing the function  $h(t) := (\sigma(t) + 2t\sigma'(t))$ , one can show that  $\max_t h(t) < \frac{4}{3}$ . Thus, in the case  $f_* \ge 0$ , using the fact that  $\sigma'(t) > 0$  we obtain

$$\|\nabla^2 F(x)\| \le Lh(f(x)) \le \frac{4}{3}L.$$

Now, consider the case  $f_* < 0$ . Since  $\sigma'(t) = \sigma(t)(1 - \sigma(t)) \le \frac{1}{4}$ ,

$$\|\nabla^2 F(x)\| \le Lh(f(x)) - 2L\sigma'(f(x))f_* \le \frac{4}{3}L - \frac{L}{2}f_*.$$

The claim follows.

**Remark C.3.** The factor  $\frac{1}{2}$  in front of  $-f_*$  in (28) can't be improved. Indeed, consider  $f(x) = \frac{1}{2}(x-a)^2 - \frac{1}{2}a^2$  with  $f_* = -\frac{1}{2}a^2$ . The second derivative of  $F(x) = \log(1+e^{f(x)})$  is

$$F''(x) = \sigma(f(x)) + \sigma(f(x)) \left(1 - \sigma(f(x))\right) (x - a)^{2},$$

$$F''(0) = \sigma(0) + \sigma(0) (1 - \sigma(0)) a^2 = \frac{1}{2} + \frac{a^2}{4} = \frac{1}{2} - \frac{f_*}{2}.$$

**Proposition C.4.** If f is convex, then  $F(x) = \log(1 + e^{f(x)})$  is also convex.

*Proof.* Trivially follows from (27).

# D LLM USAGE DISCLOSURE

In the preparation of this manuscript, large language models (LLMs) were used to improve the readability. All substantive contributions are solely by the authors.