
Optimal algorithms for group distributionally robust optimization and beyond

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Distributionally robust optimization (DRO) can improve the robustness and fairness
2 of learning methods. In this paper, we devise stochastic algorithms for a class
3 of DRO problems including group DRO, subpopulation fairness, and empirical
4 conditional value at risk (CVaR) optimization. Our new algorithms achieve faster
5 convergence rates than existing algorithms for multiple DRO settings. We also
6 provide a new information-theoretic lower bound that implies our bounds are tight
7 for group DRO. Empirically, too, our algorithms outperform known methods.

8 1 Introduction

9 Commonly, machine learning models are trained to optimize the average performance. However,
10 such models may not perform equally well among all demographic subgroups due to a hidden bias in
11 the training set or distribution shift in training and test phases [Hovy and Søgaard, 2015; Hashimoto
12 *et al.*, 2018; Martinez *et al.*, 2021; Duchi and Namkoong, 2021]. Biases in datasets are also directly
13 related to fairness concerns in machine learning [Buolamwini and Gebru, 2018; Jurgens *et al.*, 2017].

14 Recently, various algorithms based on distributionally robust optimization (DRO) have been proposed
15 to address these problems [Hovy and Søgaard, 2015; Hashimoto *et al.*, 2018; Hu *et al.*, 2018; Oren *et al.*,
16 2019; Williamson and Menon, 2019; Sagawa *et al.*, 2020; Curi *et al.*, 2020; Zhang *et al.*, 2021;
17 Martinez *et al.*, 2021; Duchi and Namkoong, 2021]. However, these algorithms are often highly
18 tailored to each specific DRO formulation. Furthermore, it is often unclear whether these proposed
19 algorithms are optimal in terms of the convergence rate. Are there a unified algorithmic methodology
20 and a lower bound for these problems?

21 **Contributions.** In this paper, we study a general class of DRO problems, which includes group
22 DRO [Hu *et al.*, 2018; Oren *et al.*, 2019; Sagawa *et al.*, 2020], subpopulation fairness [Martinez *et al.*,
23 2021], conditional value at risk (CVaR) optimization [Curi *et al.*, 2020], and many others. Let
24 $\Theta \subseteq \mathbb{R}^n$ be a convex set of model parameters and $\ell(\theta; z) : \Theta \rightarrow \mathbb{R}_+$ be a convex loss of the model
25 with parameter θ with respect to data point z . The data point z may be drawn from one out of m
26 distributions P_1, \dots, P_m which are accessible via a stochastic oracle that returns an i.i.d. sample
27 $z \sim P_i$. Let Q be a convex subset of the probability simplex in \mathbb{R}^m that contains the uniform vector,
28 i.e., $(1/m, \dots, 1/m) \in Q$. Our DRO formulation is as follows:

$$\min_{\theta \in \Theta} \max_{q \in Q} \sum_{i=1}^m q_i \mathbf{E}_{z \sim P_i} [\ell(\theta; z)]. \quad (1)$$

29 If Q are the probability simplex and scaled k -set polytope, we can recover group DRO [Sagawa *et al.*,
30 2020] and subpopulation fairness [Martinez *et al.*, 2021], respectively. Moreover, we formulate
31 a new, more general fairness concept based on weighted rankings with Q being a permutahedron,
32 which includes these special cases; see Section 2 for details.

Table 1: Summary of convergence results for group DRO. Here, m denotes the number of groups, n the dimension of θ , G the Lipschitz constant of loss function ℓ , D the diameter of feasible set Θ , M the range of loss function ℓ , and T the number of calls to stochastic oracle.

reference	convergence rate $\mathbf{E}[\varepsilon_T]$	iteration complexity	lower bound
[Sagawa <i>et al.</i> , 2020]	$O\left(m\sqrt{\frac{G^2D^2+M^2\log m}{T}}\right)$	$O(m+n) + \text{proj. onto } \Theta$	$\Omega\left(\sqrt{\frac{G^2D^2+M^2m}{T}}\right)$ (Theorem 5)
Ours (Theorem 2)	$O\left(\sqrt{\frac{G^2D^2+M^2m\log m}{T}}\right)$	$O(m+n) + \text{proj. onto } \Theta$	
Ours (Theorem 3)	$O\left(\sqrt{\frac{G^2D^2+M^2m}{T}}\right)$	$O(m+n) + \text{proj. onto } \Theta$ + solving scalar equation	

33 For our general DRO, we devise an efficient stochastic gradient algorithm. Furthermore, we show that
 34 it achieves the information-theoretic optimal convergence rate for group DRO. Our main technical
 35 contributions are as follows;

- 36 • We provide a generic stochastic gradient algorithm for our general DRO. By specializing it
 37 in the group DRO setting, we provide two algorithms (GDRO-EXP3 and GDRO-TINF)
 38 that improve the rate of Sagawa *et al.* [2020] by a factor of $\Omega(\sqrt{m})$ with the almost same
 39 complexity per iteration; see Table 1. Furthermore, our generic algorithm can be specialized
 40 to improve the convergence rate of Curi *et al.* [2020] for subpopulation fairness (a.k.a.
 41 empirical CVaR optimization). Finally, we show that our algorithm runs efficiently if Q is a
 42 permutahedron, which includes all aforementioned subclasses.
- 43 • We prove a matching information-theoretic lower bound for the convergence rate of group
 44 DRO. This implies that no algorithm can improve the convergence rate of GDRO-TINF
 45 (up to a constant factor). To the best of our knowledge, this is the first information-theoretic
 46 lower bound for group DRO.
- 47 • Our experiments on real-world and synthetic datasets show that our algorithms also empiri-
 48 cally outperform the known algorithm, supporting our theoretical analysis.

49 1.1 Our techniques

50 **Algorithms.** The core idea of our algorithms is *stochastic no-regret dynamics* [Hazan, 2016]. We
 51 regard DRO (1) as a two-player zero-sum game between a player who picks $\theta \in \Theta$ and another player
 52 who picks $q \in Q$. The two players iteratively update their solution using online learning algorithms;
 53 in particular, we will use online gradient descent (OGD) [Zinkevich, 2003] and online mirror descent
 54 (OMD) [Cesa-Bianchi and Lugosi, 2006] for the θ -player and q -player, respectively. In addition, we
 55 need to estimate gradients for both players, since the objective function of our DRO is stochastic and
 56 we cannot obtain exact gradients.

57 The convergence rate of stochastic no-regret dynamics depends on the expected regret of OGD and
 58 OMD. To obtain the optimal convergence rate, we must carefully choose the regularizer in OMD as
 59 well as gradient estimators, exploiting the structure of our DRO. In particular, we need to balance
 60 the variance of gradient estimators and the diameter terms in *both* OGD and OMD. This is the most
 61 challenging part of the algorithm design. Inspired by adversarial multi-armed bandit algorithms,
 62 we design gradient estimators for no-regret dynamics of OGD and OMD in our DRO. Indeed, our
 63 algorithms for group DRO (GDRO-EXP3 and GDRO-TINF) are based on adversarial multi-armed
 64 bandit algorithms, EXP3 [Auer *et al.*, 2003] and Tsallis-INF [Zimmert and Seldin, 2021], respectively,
 65 hence the name. Although each building block (OGD, OMD, and gradient estimators) is fairly known
 66 in the literature, we need to put them together in the right combination to obtain the optimal rate.

67 **Lower bound.** For the lower bound, we carefully design a family of group DRO instances for which
 68 any algorithm requires a certain number of queries to achieve a good objective value. To bound
 69 the number of queries, we use information-theoretic tools such as Le Cam’s lemma and bound the
 70 Kullback-Leibler divergence between Bernoulli distributions. Such tools are also used at the heart of
 71 lower bounds for stochastic convex optimization [Agarwal *et al.*, 2012] and adversarial multi-armed
 72 bandits [Auer *et al.*, 2003], but the connection to those settings is much more subtle here, and our
 73 construction is specifically designed for group DRO-type problems.

74 **1.2 Related work**

75 DRO is a wide field ranging from robust optimization to machine learning and statistics [Goh and Sim,
76 2010; Bertsimas *et al.*, 2018], whose original idea dates back to Scarf [1958]. Popular choices of the
77 uncertainty set in DRO include balls around an empirical distribution in Wasserstein distance [Esfahani
78 and Kuhn, 2018; Blanchet *et al.*, 2019], f -divergence [Namkoong and Duchi, 2016; Duchi and
79 Namkoong, 2021], χ^2 -divergence [Staib *et al.*, 2019], and maximum mean discrepancy [Staib and
80 Jegelka, 2019; Kirschner *et al.*, 2020].

81 DRO algorithms have been mainly studied for the offline setting, i.e., algorithms can access all data
82 points of the empirical distribution. Note that our DRO is not offline because the group distributions
83 are given by the stochastic oracles. Namkoong and Duchi [2016] proposed stochastic gradient
84 algorithms for offline DRO with f -divergence uncertainty sets. Curi *et al.* [2020] used no-regret
85 dynamics for empirical CVaR minimization. Their algorithm invokes sampling from k -DPP in each
86 iteration, which is more computationally demanding than our algorithm. Furthermore, our algorithm
87 gets rid of an $O(\log m)$ factor in the convergence rate using the Tsallis entropy regularizer; see
88 Theorem 4. Qi *et al.* [2021]; Jin *et al.* [2021] devised stochastic gradient algorithms for several DRO
89 with non-convex losses.

90 Agarwal *et al.* [2012] gave a lower bound for stochastic convex optimization, which is a special case
91 of our DRO with only one distribution. Recently, Carmon *et al.* [2021] showed a lower bound for
92 minimax problem $\min_x \max_{i=1}^m f_i(x)$ for non-stochastic Lipschitz convex f_i . Our lower bound deals
93 with the stochastic functions, so this result does not apply.

94 In this paper, we assume that the group information is given in advance. However, the group
95 information might not be easy to define in practice. Bao *et al.* [2021] propose a simple method to
96 define groups for classification problems based on mistakes of models in the training phase. Their
97 method often generates group DRO instances with large m . Our algorithms are more efficient for
98 such group DRO thanks to the better dependence on m in the convergence rate.

99 No-regret dynamics is a well-studied method for solving two-player zero-sum games [Cesa-Bianchi
100 and Lugosi, 2006]. For non-stochastic convex-concave games, one can achieve $O(1/T)$ convergence
101 via predictable sequences [Rakhlin and Sridharan, 2013]. This result does not apply to our setting
102 because our DRO is a stochastic game.

103 **Notations.** Throughout the paper, m denotes the number of distributions (groups) and n denotes
104 the dimension of a variable θ . For a positive integer m , we write $[m] := \{1, \dots, m\}$. The orthogonal
105 projection onto set Θ is denoted by proj_Θ . The i th standard unit vector is denoted by e_i and the
106 all-one vector is denoted by $\mathbf{1}$. The probability simplex in \mathbb{R}^m is denoted by Δ_m .

107 **2 Examples contained in our general DRO**

108 In this section, we show how several DRO formulations in the literature can be phrased in our general
109 DRO formulation (1). In addition, we propose a novel fairness constraint based on weighted rankings
110 using our general DRO.

111 **Group DRO.** When Q equals the probability simplex, we obtain group DRO [Hu *et al.*, 2018;
112 Oren *et al.*, 2019; Sagawa *et al.*, 2020]:

$$\min_{\theta \in \Theta} \max_{i=1}^m \mathbf{E}_{z \sim P_i} [\ell(\theta; z)]. \quad (2)$$

113 That is, group DRO aims to minimize the expected loss in the worst group, thereby ensuring better
114 performance across all groups.

115 **Empirical CVaR, Subpopulation fairness, Average top- k worst group loss.** Group DRO may
116 yield overly pessimistic solutions. For instance, the groups might be automatically generated by other
117 algorithms (such as one in Bao *et al.* [2021]) and there might exist a few “outlier” groups that make
118 the group DRO objective trivial.

119 For such a case, we can restrict Q to a small subset of the probability simplex so that the solution
120 cannot put large weights on a few outlier groups. Especially, let $Q = \left\{ q \in \Delta_m : 0 \leq q_i \leq \frac{1}{pm} \right\}$ for

121 some parameter $p \in (0, 1)$, i.e., Q is a scaled k -set polytope. The intuition behind the choice of Q
 122 is that, by limiting the largest entry of q to $1/pm$, DRO would optimize the expected loss over the
 123 worst p -fraction subgroups of m groups. Therefore, if the fraction of outlier groups is sufficiently
 124 small compared to p , then p -fraction subgroups must contain “inlier” groups as well. Therefore, it is
 125 likely that DRO with Q finds solutions z more robust than group DRO.

126 When P_i is the Dirac measure of data z_i , then the resulting DRO is empirical CVaR optimization [Curi
 127 *et al.*, 2020]. In the fairness context, the same problem is called subpopulation fairness [Williamson
 128 and Menon, 2019; Martinez *et al.*, 2021; Duchi and Namkoong, 2021].

129 If $p = m/k$ for some positive integer k , the resulting DRO is the average top- k worst group
 130 loss [Zhang *et al.*, 2021]:

$$\min_{\theta \in \Theta} \frac{1}{k} \sum_{i=1}^k L_i^\downarrow(\theta),$$

131 where $L_i^\downarrow(\theta)$ denotes the the i th largest population group loss of θ . More precisely, let $L_i(\theta) =$
 132 $\mathbf{E}_{z \sim P_i}[\ell(\theta; z)]$ for $i \in [m]$ and sort them in the non-increasing order: $L_1^\downarrow(\theta) \geq \dots \geq L_m^\downarrow(\theta)$.

133 **Weighted ranking of group losses.** The aforementioned DRO formulations are special cases of the
 134 following DRO, which we call the *weighted ranking of group losses*. Let $\alpha \in \Delta^m$ be a fixed vector
 135 with non-increasing entries. Let Q be the permutahedron of α , the convex hull of $(\alpha_{\sigma(1)}, \dots, \alpha_{\sigma(m)})$
 136 for all permutations σ of $[m]$. Then, the resulting DRO is

$$\min_{\theta \in \Theta} \sum_{i=1}^m \alpha_i L_i^\downarrow(\theta).$$

137 Group DRO corresponds to $\alpha = (1, 0, \dots, 0)$ and the average top- k worst group losses corresponds
 138 to $\alpha = (\underbrace{1/k, \dots, 1/k}_{k \text{ times}}, 0, \dots, 0)$. Another example that is contained in none of the above examples is

139 *lexicographic minimax fairness* [Diana *et al.*, 2021]. The goal of lexicographical minimax fairness
 140 is to find $\theta \in \Theta$ such that the sequence $(L_1^\downarrow(\theta), \dots, L_m^\downarrow(\theta))$ is lexicographically minimum. This
 141 corresponds to α with sufficiently varied entries, i.e., $\alpha_1 \gg \alpha_2 \gg \dots \gg \alpha_m$.

142 3 Algorithms

143 In this section, we describe our algorithms. First, we present a generic algorithm for our general
 144 DRO (1) and provide a unified convergence analysis in Section 3.1. Then, we specialize it into two
 145 concrete algorithms for group DRO (2) in Section 3.2. We sketch algorithms for the average of top- k
 146 group losses and weighted ranking of group loss in Section 3.3.

147 3.1 Algorithm for the general case

148 We present our algorithm for general DRO (1). At a high level, our algorithm can be regarded as
 149 stochastic no-regret dynamics. Let us denote $L(\theta, q) := \sum_{i=1}^m q_i \mathbf{E}_{z \sim P_i}[\ell(\theta; z)]$. Imagine that the
 150 θ -player and q -player run online algorithms \mathcal{A}_θ and \mathcal{A}_q , respectively, to solve the minimax problem
 151 $\min_{\theta \in \Theta} \max_{q \in Q} L(\theta, q)$. That is, for $t = 1, \dots, T$,

- 152 • $\theta_t \in \Theta$ and $q_t \in Q$ are determined by \mathcal{A}_θ and \mathcal{A}_q , respectively.
- 153 • Both players feed gradient estimators $\hat{\nabla}_{\theta,t}$ and $\hat{\nabla}_{q,t}$ to \mathcal{A}_θ and \mathcal{A}_q , respectively. Here,
 154 $\mathbf{E}[\hat{\nabla}_{\theta,t}] = \nabla_\theta L(\theta_t, q_t)$ and $\mathbf{E}[\hat{\nabla}_{q,t}] = \nabla_q L(\theta_t, q_t)$.

155 Let

$$\varepsilon_T := \max_{q \in Q} L(\bar{\theta}_{1:T}, q) - \min_{\theta \in \Theta} \max_{q \in Q} L(\theta, q)$$

156 be the optimality gap of the averaged iterate $\bar{\theta}_{1:T} = \frac{1}{T} \sum_{t=1}^T \theta_t$. We can bound the expected
 157 convergence rate $\mathbf{E}[\varepsilon_T]$ via regrets R_θ and R_q of these online algorithms (see Appendix A for a
 158 formal definition), i.e.,

$$\mathbf{E}[\varepsilon_T] \leq \frac{\mathbf{E}[R_\theta(T)] + \mathbf{E}[R_q(T)]}{T}. \quad (3)$$

159 We can obtain hence the convergence rate of the above algorithms by investigating the expected regret
160 bounds of these online algorithms.

161 To get a concrete algorithm, we must specify the online algorithms $\mathcal{A}_\theta, \mathcal{A}_q$ as well as the gradient
162 estimators $\hat{\nabla}_{\theta,t}, \hat{\nabla}_{q,t}$. We use OGD and OMD as \mathcal{A}_θ and \mathcal{A}_q , respectively. We construct the gradient
163 estimators by sampling $i_t \sim q_t$ and $z \sim P_{i_t}$ and setting $\hat{\nabla}_{\theta,t} = \nabla_\theta \ell(\theta_t; z)$ and $\hat{\nabla}_{q,t} = \frac{\ell(\theta_t; z)}{q_{t,i_t}} \mathbf{e}_{i_t}$.
164 This leads to Algorithm 1. There, $\Psi : Q \rightarrow \mathbb{R}$ denotes the regularizer of OMD and $\eta_{\theta,t}$ and η_q
165 denote the step sizes of OGD and OMD, respectively.¹ It turns out that this combination of online
166 algorithms and gradient estimators yields the best convergence rate (for group DRO) because the
167 expected regrets of both players are optimal.

Algorithm 1 Algorithm for general DRO (1)

Require: initial solution $\theta_1 \in \Theta$, number of iterations T , step sizes $\eta_{\theta,t} > 0$ ($t \in [T]$), $\eta_q > 0$, and
a strictly convex function $\Psi : Q \rightarrow \mathbb{R}$.
1: Let $q_1 = (1/m, \dots, 1/m)$.
2: **for** $t = 1, \dots, T$ **do**
3: Sample $i_t \sim q_t$.
4: Call the stochastic oracle to obtain $z \sim P_{i_t}$.
5: $\theta_{t+1} \leftarrow \text{proj}_\Theta(\theta_t - \eta_{\theta,t} \nabla_\theta \ell(\theta_t; z))$
6: $\nabla \Psi(\tilde{q}_{t+1}) \leftarrow \nabla \Psi(q_t) - \frac{\eta_q}{q_{t,i_t}} \ell(\theta_t; z) \mathbf{e}_{i_t}$; $q_{t+1} \leftarrow \text{argmin}_{q \in Q} D_\Psi(q, \tilde{q}_{t+1})$, where
 $D_\Psi(x, y) = \Psi(x) - \Psi(y) - \nabla \Psi(x)^\top (y - x)$ is the Bregman divergence with respect to Ψ .
7: **return** $\frac{1}{T} \sum_{t=1}^T \theta_t$.

168 We now analyze the convergence rate of Algorithm 1. We make the following standard assumptions.

169 **Assumption 1.** *The loss function $\ell(\theta; z)$ is continuously differentiable and G -Lipchitz in θ , and has
170 range $[0, M]$ for all z . The Euclidean diameter of the feasible region Θ is at most D .*

171 The following theorem follows from plugging regret bounds of OGD and OGD, and the construction
172 of the gradient estimators into (3).

173 **Theorem 1.** *If $\eta_{\theta,t}$ is nonincreasing, Algorithm 1 achieves the expected convergence rate*

$$\mathbf{E}[\varepsilon_T] \leq \frac{1}{T} \left(\frac{G^2}{2} \sum_{t=1}^T \eta_{\theta,t} + \frac{D^2}{2\eta_{\theta,T}} + \frac{M^2}{2} \eta_q \sum_{t=1}^T \mathbf{E}_{i_t} \left[\frac{(\nabla^2 \Psi(q_t))_{i_t, i_t}^{-1}}{q_{t, i_t}^2} \right] + \frac{\max_{q^* \in Q} D_\Psi(q^*, \mathbf{1}/m)}{\eta_q} \right).$$

174 A formal proof can be found in Appendix B. We will see how specific choices of the regularizer Ψ
175 yield various algorithms and convergence rates for group DRO and others in the next subsections. A
176 few remarks on the regularizers, step sizes, and projection step are in order.

177 **Regularizer.** Although Algorithm 1 works with general Ψ , we can choose a specific regularizer for
178 Q appearing in applications, e.g, the probability simplex, scaled k -set polytope, or a permutahedron.
179 In the next subsections, we show that the entropy regularizer $\Psi(x) = \sum_i (x_i \log x_i - x_i)$ and Tsallis
180 entropy regularizer $\Psi(x) = 2(1 - \sum_i \sqrt{x_i})$ yield efficient algorithms with improved convergence
181 rates for these cases.

182 **Step sizes.** The theorem includes decreasing step sizes such as $\eta_{\theta,t} = \frac{D}{mG\sqrt{t}}$ in addition to fixed
183 step sizes. Decreasing step sizes have the advantage that we do not require the knowledge of T
184 at the beginning of the algorithm but come at the cost of an extra constant factor in the expected
185 convergence rate. Since both step size policies give the asymptotically same convergence rate, we
186 describe only fixed step sizes in the theorems in the next subsections. In practice, decreasing step
187 sizes stabilize the algorithm and often outperform fixed step sizes.

188 **Projection step.** In general, the Bregman projection $\text{argmin}_{q \in Q} D_\Psi(q, \tilde{q}_{t+1})$ is convex, but may
189 be costly to compute. For the applications described in Section 2, Q is a permutahedron. In this case,

¹We make a standard assumption that the regularizer Ψ is differentiable and strictly convex, and satisfies $\|\nabla \Psi(x)\| \rightarrow +\infty$ as x tends to the boundary of Q .

Algorithm 2 GDRO-EXP3

Require: initial solution $\theta_1 \in \Theta$, number of iterations T , and step sizes $\eta_{\theta,t} > 0$ ($t \in [T]$), $\eta_q > 0$.

- 1: Let $q_t = (1/m, \dots, 1/m)$.
- 2: **for** $t = 1, \dots, T$ **do**
- 3: Sample $i_t \sim q_t$.
- 4: Call the stochastic oracle to obtain $z \sim P_{i_t}$.

$$5: \quad \theta_{t+1} \leftarrow \text{proj}_{\Theta}(\theta_t - \eta_{\theta,t} \nabla_{\theta} \ell(\theta_t; z))$$

$$6: \quad \tilde{q}_{t+1} \leftarrow q_t \exp\left(\frac{\eta_q \ell(\theta_t; z) \mathbf{e}_{i_t}}{q_{t,i_t}}\right)$$

$$7: \quad q_{t+1} \leftarrow \frac{\tilde{q}_{t+1}}{\sum_i \tilde{q}_{t+1,i}}$$

$$8: \quad \textbf{return} \quad \frac{1}{T} \sum_{t=1}^T \theta_t.$$

Algorithm 3 GDRO-TINF

Require: initial solution $\theta_1 \in \Theta$, number of iterations T , and step sizes $\eta_{\theta,t} > 0$ ($t \in [T]$), $\eta_q > 0$.

- 1: Let $q_t = (1/m, \dots, 1/m)$.

- 2: **for** $t = 1, \dots, T$ **do**

- 3: Sample $i_t \sim q_t$.

- 4: Call the stochastic oracle to obtain $z \sim P_{i_t}$.

$$5: \quad \theta_{t+1} \leftarrow \text{proj}_{\Theta}(\theta_t - \eta_{\theta,t} \nabla_{\theta} \ell(\theta_t; z))$$

$$6: \quad \tilde{q}_{t+1} \leftarrow q_t \left(\mathbf{1} - \frac{\eta_q \sqrt{q_t}}{q_{t,i_t}} \ell(\theta_t; z) \mathbf{e}_{i_t} \right)^{-2}$$

- 7: Compute $\alpha \in \mathbb{R}$ such that $\sum_{i=1}^m (\sqrt{\tilde{q}_{t+1,i}} - \alpha)^{-2} = 1$.

$$8: \quad q_{t+1} \leftarrow (\sqrt{\tilde{q}_{t+1}} - \alpha \mathbf{1})^{-2}$$

$$9: \quad \textbf{return} \quad \frac{1}{T} \sum_{t=1}^T \theta_t.$$

190 it is known that the Bregman projection with respect to the entropy and Tsallis entropy regularizers
 191 can be done in $O(m \log m)$ time [Lim and Wright, 2016]. If Q is the probability simplex, we even
 192 have a closed form for the Bregman projection.

193 3.2 Algorithms for Group DRO

194 As applications of our generic algorithm, we now describe two concrete algorithms for group DRO (2)
 195 and their convergence rates.

196 **GDRO-EXP3.** Let Ψ be the entropy regularizer, which corresponds to the EXP3 algorithm for
 197 the q -player. The resulting algorithm, GDRO-EXP3, is shown in Algorithm 2. The update is in a
 198 closed formula and its complexity is $O(m + n)$ time. The convergence rate follows from Theorem 1.

199 **Theorem 2.** *If $\eta_{\theta,t}$ is nonincreasing, GDRO-EXP3 (Algorithm 2) achieves the expected convergence*
 200 *rate*

$$\mathbf{E}[\varepsilon_T] \leq \frac{1}{T} \left(\frac{G^2}{2} \sum_{t=1}^T \eta_{\theta,t} + \frac{D^2}{2\eta_{\theta,T}} + \frac{mM^2}{2} \eta_q T + \frac{\log m}{\eta_q} \right). \quad (4)$$

201 For $\eta_{\theta,t} = \frac{D}{G\sqrt{T}}$ and $\eta_q = \sqrt{\frac{2\log m}{mM^2T}}$, we obtain

$$\mathbf{E}[\varepsilon_T] \leq \sqrt{2} \frac{\sqrt{G^2 D^2 + 2M^2 m \log m}}{\sqrt{T}}.$$

202 **GRDO-TINF.** We can further improve the convergence rate using the Tsallis entropy regularizer
 203 at the cost of a slightly higher iteration complexity. The update of q_t is then

$$\tilde{q}_{t+1} = q_t \left(\mathbf{1} - \frac{\eta_q \sqrt{q_t}}{q_{t,i_t}} \ell(\theta_t; z) \mathbf{e}_{i_t} \right)^{-2}, \quad q_{t+1} := \left(\sqrt{\tilde{q}_{t+1}} - \alpha \mathbf{1} \right)^{-2},$$

204 where the multiplication, square-root, and power operations are entry-wise and $\alpha \in \mathbb{R}$ is the unique
 205 solution of equation $\sum_{i=1}^m (\sqrt{\tilde{q}_{t+1,i}} - \alpha)^{-2} = 1$. The solution α can be computed via the Newton
 206 method. Practically, one can use α in the previous iteration to warm start the Newton method. In
 207 each iteration, the algorithm performs a single orthogonal projection onto Θ , the Newton method for
 208 finding α , and $O(m + n)$ operations to update θ_t, q_t . The pseudocode is given in Algorithm 3. From
 209 Theorem 1, we obtain the following convergence rate.

210 **Theorem 3.** *If $\eta_{\theta,t}$ is nonincreasing, GDRO-TINF (Algorithm 3) achieves the expected convergence*
 211 *rate*

$$\mathbf{E}[\varepsilon_T] \leq \frac{1}{T} \left(\frac{G^2}{2} \sum_{t=1}^T \eta_{\theta,t} + \frac{D^2}{2\eta_{\theta,T}} + \sqrt{m} M^2 \eta_q T + \frac{\sqrt{m}}{\eta_q} \right). \quad (5)$$

212 For $\eta_{\theta,t} = \frac{D}{G\sqrt{T}}$ and $\eta_q = \frac{1}{M\sqrt{T}}$, we obtain

$$\mathbf{E}[\varepsilon_T] \leq \sqrt{2} \frac{\sqrt{G^2 D^2 + 4M^2 m}}{\sqrt{T}}.$$

213 **Comparison to Sagawa *et al.* [2020].** Our algorithms improve the convergence rate of Sagawa *et al.* [2020] by a factor of $O(\sqrt{m})$; see Table 1. The reason lies in the choice of gradient estimator. All
 214 algorithms are stochastic no-regret dynamics. As outlined above, their convergence hence can be
 215 bounded by the regrets of the players, which depend on the variance of the local norm of the gradient
 216 estimators. Their strategy is based on uniform sampling that yields a variance of $O(m)$ for both
 217 players, whereas our bound is $O(\sqrt{m})$ thanks to the gradient estimators tailored to the regularizer of
 218 OMD. More details may be found in Appendix D.

220 3.3 Algorithm for weighted ranking of group losses

221 We now consider a more general case that Q is a permutahedron. Applying Algorithm 1 with the
 222 Tsallis entropy regularizer, we obtain the following result.

223 **Theorem 4.** *If $\eta_{\theta,t}$ is nonincreasing and Q is a permutahedron, Algorithm 1 with the Tsallis entropy
 224 regularizer achieves the same expected convergence rate as Theorem 3. Furthermore, the iteration
 225 complexity is $O(m \log m + n)$.*

226 This implies a convergence rate of $O(\sqrt{\frac{G^2 D^2 + M^2 m}{T}})$ for empirical CVaR optimization, which
 227 improves $O(\sqrt{\frac{G^2 D^2 + M^2 m \log m}{T}})$ convergence by Curi *et al.* [2020]. Furthermore, their iteration
 228 complexity is $O(m^3)$ due to the k -DPP sampling step, so our algorithm is even faster in terms of
 229 iteration complexity.

230 4 Lower bound

231 Theorem 3 states that we can find an ε -optimal solution for group DRO in $O(\frac{G^2 D^2 + M^2 m}{\varepsilon^2})$ calls to
 232 stochastic oracles. Next, we show that this query complexity is information-theoretically optimal.

233 Let \mathcal{L} be a class of convex G -Lipschitz loss functions $\ell : \Theta \rightarrow [0, M]$. Given a loss function $\ell \in \mathcal{L}$,
 234 and an m -set $\mathcal{P} = \{P_1, \dots, P_m\}$ of distributions, denote the optimality gap of $\theta \in \Theta$ by

$$R(\theta, \ell, \mathcal{P}) = \max_{P \in \mathcal{P}} \mathbf{E}_{z \sim P} [\ell(\theta; z)] - \min_{\theta^* \in \Theta} \max_{P \in \mathcal{P}} \mathbf{E}_{z \sim P} [\ell(\theta^*; z)].$$

235 Let \mathcal{A}_T be the set of algorithms that outputs $\hat{\theta} \in \Theta$ making T queries to the stochastic oracle.

Theorem 5 (Lower Bound).

$$\inf_{\hat{\theta} \in \mathcal{A}_T} \sup_{\ell \in \mathcal{L}, \Theta, \mathcal{P}} \mathbf{E}_{\mathcal{P}} [R(\hat{\theta}, \ell, \mathcal{P})] \geq \Omega \left(\max \left\{ \frac{GD}{\sqrt{T}}, M \sqrt{\frac{m}{T}} \right\} \right),$$

236 where Θ runs over convex sets with diameter D and \mathcal{P} over m -sets of distributions, and $\mathbf{E}_{\mathcal{P}}$ denotes
 237 the expectation over outcomes of the stochastic oracle in \mathcal{P} .

238 As $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y} \leq \sqrt{2(x+y)}$ for $x, y \geq 0$, this theorem immediately implies that the
 239 minimax convergence rate is $\Omega \left(\sqrt{\frac{G^2 D^2 + M^2 m}{T}} \right)$, which equals the convergence rate achieved by
 240 Algorithm 3 up to a constant factor.

241 **Proof Sketch.** It suffices to show two lower bounds $\frac{GD}{\sqrt{T}}$ and $M \sqrt{\frac{m}{T}}$ independently. The former is
 242 a well-known lower bound for stochastic convex optimization [Agarwal *et al.*, 2012]. To illustrate the
 243 latter, we take an algorithmic dependent point of view via the Le cam’s method. For any algorithm
 244 in \mathcal{A}_T , we need to construct instances $\mathcal{P}_0, \mathcal{P}_1$ such that the total variation distance between the
 245 distributions over the query outcomes (they depend on both the behavior of the algorithm and the
 246 instance) with respect to \mathcal{P}_0 and \mathcal{P}_1 is small. On the other hand, the objective function of the two

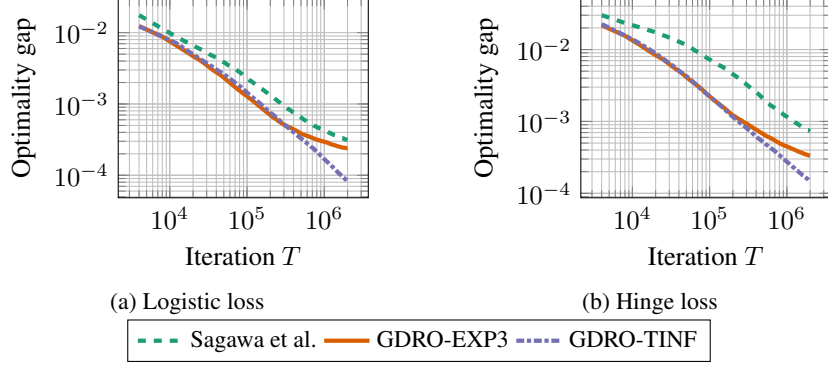


Figure 1: Results on Adult dataset

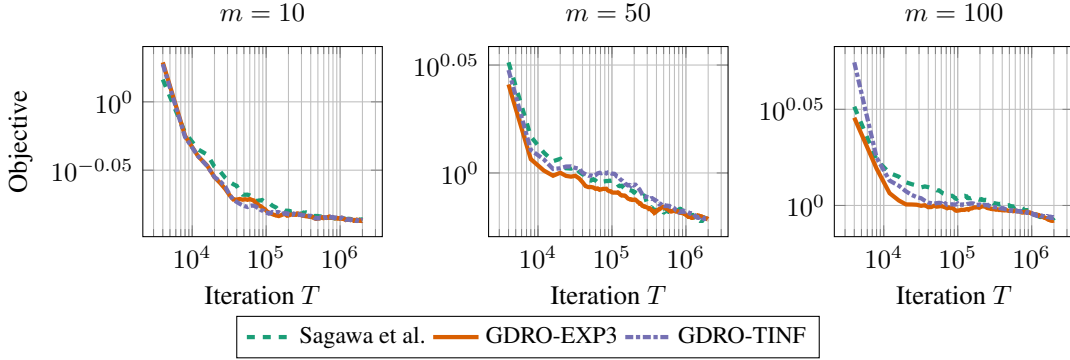


Figure 2: Results on synthetic dataset

247 instances must be well-separated, i.e., any fixed θ is δ sub-optimal for either \mathcal{P}_0 or \mathcal{P}_1 . So, any
 248 algorithm that solves group DRO up to error δ needs to distinguish two instances \mathcal{P}_0 and \mathcal{P}_1 . This
 249 implies a query lower bound because the total variation distance of the outcome distributions of
 250 these instances is small. The challenge is how to construct such instances for the regime of small
 251 dimensions of θ , e.g, $n = 1$. To this end, we carefully construct linear functions for m groups using
 252 opposite slopes. Then, based on the behavior of the algorithm, we tweak the noise bias in one of the
 253 groups with a positive slope, in a way that any fixed θ is $\Theta(\delta)$ sub-optimal for one of these instances.
 254 For the detailed proof, see Appendix C.

255 5 Experiments

256 In this section, we compare our algorithms with the known algorithm using real-world and synthetic
 257 datasets. We follow the setup in [Namkoong and Duchi, 2016].

258 **Adult dataset.** For the real-world dataset, we use Adult dataset [Dua and Graff, 2017]. The dataset
 259 consists of age, gender, race, educational background, and many other attributes of 48,842 individuals
 260 from the US census. The task is to predict whether the person’s income is greater than 50,000 USD
 261 or not. We set up 6 groups based on the race and gender attributes: each group corresponds to
 262 a combination of {black, white, others} \times {female, male}. Converting the categorical features to
 263 dummy variables, we obtain a 101-dimensional feature vector $a \in \mathbb{R}^n$ ($n = 101$) for each individual.
 264 We train the linear model with the logistic loss and hinge loss functions. The group-DRO objective is
 265 the worst empirical loss over the 6 groups:

$$\max_{i=1}^6 \frac{1}{|I_i|} \sum_{(a,b) \in I_i} \ell(\theta; a, b),$$

266 where I_i is the set of data points in the i th group. The feasible region is set to the Euclidean ball of
 267 radius $D = 10$.

268 **Synthetic dataset.** To observe the performance of the algorithms over the regime of high-dimension
 269 model parameters and the larger number of groups, we also conducted experiments using the following
 270 synthetic instances. First, we set $n = 500$ and varied $m \in \{10, 50, 100\}$. For each group $i \in [m]$, we
 271 generated the true classifier $\theta_i^* \in \mathbb{R}^n$ from the uniform distribution over the unit sphere in \mathbb{R}^n . The i th
 272 group distribution P_i was the empirical distribution of 1,000 data points, where each data point (a, b)
 273 was drawn as $a \sim N(0, I_n)$ and $b = \text{sign}(a^\top \theta_i^*)$ with probability 0.9 and $b = -\text{sign}(a^\top \theta_i^*)$ with
 274 probability 0.1. We trained the linear model with the hinge loss function. Finally, the group-DRO
 275 objective is

$$\max_{i=1}^m \mathbf{E}_{(a,b) \sim P_i} [\ell(\theta; a, b)].$$

276 The feasible region is set to the Euclidean ball of radius $D = 10$.

277 5.1 Algorithms

278 We implemented GDRO-EXP3, GDRO-TINF, and the algorithm in [Sagawa *et al.*, 2020] in Python.
 279 We ran our algorithms for $T = 2,000,000$ iterations.

280 **Inner online algorithms.** It is known that EXP3 has a variance as large as $O(T^2)$ [Lattimore and
 281 Szepesvári, 2020]. Therefore, vanilla EXP3 often fails to achieve a sublinear regret even though it
 282 achieves $O(\sqrt{T})$ regret *in expectation*. This large variance makes it difficult to reliably evaluate the
 283 performance of the algorithms. To stabilize the algorithms, we replaced EXP3 with its variation,
 284 EXP3P [Auer *et al.*, 2003], which achieves $O(\sqrt{T})$ regret *with high probability*. Note that this change
 285 does not harm our expected convergence bounds.

286 **Step sizes.** The choice of step sizes is crucial to the practical performance of first-order methods.
 287 We found that the decreasing step size $\eta_{\theta,t} \sim 1/\sqrt{t}$ for θ_t and the fixed step size $\eta_q \sim 1/\sqrt{T}$ for
 288 q_t gave the best results. More precisely, we set $\eta_{\theta,t} = \frac{C_\theta D}{\sqrt{t}}$ ($t \in [T]$) and $\eta_q = C_q \sqrt{\frac{\log m}{mT}}$, where
 289 $C_\theta \in [0.1, 5.0]$ and $C_q \in [0.1, 3.0]$ are hyper-parameters tuned for each algorithm. We used the best
 290 hyper-parameter found by Optuna [Akiba *et al.*, 2019] for the shown results.

291 **Mini-batch.** The use of mini-batch often improves the stability of stochastic gradient algorithms.
 292 In our experiments, we used mini-batches of size 10 to evaluate stochastic gradients. Neither the
 293 objective values of outputs nor the stability was improved with larger mini-batch sizes. The group
 294 DRO objective is evaluated using the entire dataset.

295 **Initialization.** For both datasets, we initialized the algorithms with $\theta_1 = \mathbf{0}$.

296 5.2 Results

297 We show the results of our experiment in Figures 1 and 2.

298 **Adult dataset.** In Figure 1, we plot the optimality gap of the averaged iterate $\frac{1}{T} \sum_{t=1}^T \theta_t$ against the
 299 number of iteration T . We observe that all the algorithms converge with a rate roughly $T^{-0.5}$ for both
 300 loss functions, consistent with our convergence bound. Furthermore, our algorithms (GDRO-EXP3
 301 and GDRO-TINF) achieve faster convergence compared to the algorithm by Sagawa *et al.* [2020].
 302 Interestingly, GDRO-TINF achieves a 10^{-4} optimality gap in $T = 10^6$ iterations, which is faster
 303 than the theoretical $T^{-0.5}$ rate in Theorem 3.

304 **Synthetic dataset.** In Figure 2, we plot the objective values of the averaged iterate against the
 305 number of iterations. For all the values of m , our algorithms (especially GDRO-EXP3) consistently
 306 achieve smaller loss values faster than the known algorithm. The performance gap between our
 307 algorithms and the known algorithm increased as m grows, which verifies that our algorithms have
 308 better dependence on m in the convergence rate.

309 **References**

- 310 Alekh Agarwal, Peter L. Bartlett, Pradeep Ravikumar, and Martin J. Wainwright. Information-
311 theoretic lower bounds on the oracle complexity of stochastic convex optimization. *IEEE Transac-*
312 *tions on Information Theory*, pages 3235–3249, 2012.
- 313 Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna:
314 A next-generation hyperparameter optimization framework. In *Proceedings of the 25rd ACM*
315 *SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019.
- 316 Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed
317 bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2003.
- 318 Yujia Bao, Shiyu Chang, and Regina Barzilay. Predict then interpolate: A simple algorithm to learn
319 stable classifiers. In *Proceedings of the 38th International Conference on Machine Learning*,
320 volume 139, pages 640–650, 2021.
- 321 Dimitris Bertsimas, Vishal Gupta, and Nathan Kallus. Data-driven robust optimization. *Mathematical*
322 *Programming*, 167(2):235–292, 2018.
- 323 Jose Blanchet, Yang Kang, and Karthyek Murthy. Robust wasserstein profile inference and applica-
324 tions to machine learning. *Journal of Applied Probability*, 56(3):830–857, 2019.
- 325 Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial
326 gender classification. In *Proceedings of the 1st Conference on Fairness, Accountability and*
327 *Transparency*, pages 77–91, 2018.
- 328 Yair Carmon, Arun Jambulapati, Yujia Jin, and Aaron Sidford. Thinking inside the ball: Near-optimal
329 minimization of the maximal loss. In *Proceedings of 34th Conference on Learning Theory*, volume
330 134 of *Proceedings of Machine Learning Research*, pages 866–882, 2021.
- 331 Nicolo Cesa-Bianchi and Gabor Lugosi. *Prediction, Learning, and Games*. Cambridge University
332 Press, 2006.
- 333 Sebastian Curi, Kfir Y. Levy, Stefanie Jegelka, and Andreas Krause. Adaptive sampling for stochastic
334 risk-averse learning. In *Advances in Neural Information Processing Systems*, pages 1036–1047,
335 2020.
- 336 Emily Diana, Wesley Gill, Ira Globus-Harris, Michael Kearns, Aaron Roth, and Saeed Sharifi-
337 Malvajerdi. Lexicographically fair learning: Algorithms and generalization. In *Proceedings of the*
338 *2nd Symposium on Foundations of Responsible Computing*, pages 6:1–6:23, 2021.
- 339 Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- 340 John C. Duchi and Hongseok Namkoong. Learning models with uniform performance via distribu-
341 tionally robust optimization. *The Annals of Statistics*, 49(3):1378 – 1406, 2021.
- 342 Peyman Mohajerin Esfahani and Daniel Kuhn. Data-driven distributionally robust optimization
343 using the wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical*
344 *Programming*, 171(1):115–166, 2018.
- 345 Joel Goh and Melvyn Sim. Distributionally robust optimization and its tractable approximations.
346 *Operations Research*, 58(4-part-1):902–917, 2010.
- 347 Tatsunori Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. Fairness without
348 demographics in repeated loss minimization. In *Proceedings of the 35th International Conference*
349 *on Machine Learning*, pages 1929–1938, 2018.
- 350 Elad Hazan. *Introduction to Online Convex Optimization*. 2016.
- 351 Dirk Hovy and Anders Søgaard. Tagging performance correlates with author age. In *Proceedings of*
352 *the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International*
353 *Joint Conference on Natural Language Processing*, pages 483–488, 2015.

- 354 Weihua Hu, Gang Niu, Issei Sato, and Masashi Sugiyama. Does distributionally robust supervised
355 learning give robust classifiers? In *Proceedings of the 35th International Conference on Machine*
356 *Learning*, pages 2029–2037, 2018.
- 357 Jikai Jin, Bohang Zhang, Haiyang Wang, and Liwei Wang. Non-convex distributionally robust
358 optimization: Non-asymptotic analysis. In *Advances in Neural Information Processing Systems*,
359 volume 34, pages 2771–2782, 2021.
- 360 David Jurgens, Yulia Tsvetkov, and Dan Jurafsky. Incorporating dialectal variability for socially
361 equitable language identification. In *Proceedings of the 55th Annual Meeting of the Association*
362 *for Computational Linguistics*, pages 51–57, 2017.
- 363 Johannes Kirschner, Ilija Bogunovic, Stefanie Jegelka, and Andreas Krause. Distributionally ro-
364 bust bayesian optimization. In *Proceedings of the 33rd International Conference on Artificial*
365 *Intelligence and Statistics*, pages 2174–2184, 2020.
- 366 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 367 Cong Han Lim and Stephen J. Wright. Efficient bregman projections onto the permutahedron and
368 related polytopes. In *Proceedings of the 19th International Conference on Artificial Intelligence*
369 *and Statistics*, pages 1205–1213, 2016.
- 370 Natalia L Martinez, Martin A Bertran, Afroditi Papadaki, Miguel Rodrigues, and Guillermo Sapiro.
371 Blind pareto fairness and subgroup robustness. In *Proceedings of the 38th International Conference*
372 *on Machine Learning*, pages 7492–7501, 2021.
- 373 Hongseok Namkoong and John C Duchi. Stochastic gradient methods for distributionally robust
374 optimization with f -divergences. In *Advances in Neural Information Processing Systems*, 2016.
- 375 Yonatan Oren, Shiori Sagawa, Tatsunori Hashimoto, and Percy Liang. Distributionally robust
376 language modeling. In *Proceedings of the Conference on Empirical Methods in Natural Language*
377 *Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-*
378 *IJCNLP)*, pages 4227–4237, 2019.
- 379 Qi Qi, Zhishuai Guo, Yi Xu, Rong Jin, and Tianbao Yang. An online method for a class of
380 distributionally robust optimization with non-convex objectives. In *Advances in Neural Information*
381 *Processing Systems*, volume 34, pages 10067–10080, 2021.
- 382 Alexander Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable
383 sequences. In *Advances in Neural Information Processing Systems*, 2013.
- 384 Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust
385 neural networks for group shifts: On the importance of regularization for worst-case generalization.
386 In *The 8th International Conference on Learning Representations*, 2020.
- 387 Herbert Scarf. A min-max solution of an inventory problem. *Studies in the mathematical theory of*
388 *inventory and production*, 1958.
- 389 Matthew Staib and Stefanie Jegelka. Distributionally robust optimization and generalization in kernel
390 methods. In *Advances in Neural Information Processing Systems*, 2019.
- 391 Matthew Staib, Bryan Wilder, and Stefanie Jegelka. Distributionally robust submodular maximization.
392 In *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, pages
393 506–516, 2019.
- 394 Robert Williamson and Aditya Menon. Fairness risk measures. In *Proceedings of the 36th Interna-*
395 *tional Conference on Machine Learning*, pages 6786–6797, 2019.
- 396 Jingzhao Zhang, Aditya Krishna Menon, Andreas Veit, Srinadh Bhojanapalli, Sanjiv Kumar, and
397 Suvrit Sra. Coping with label shift via distributionally robust optimisation. In *The 9th International*
398 *Conference on Learning Representations*, 2021.
- 399 Julian Zimmert and Yevgeny Seldin. Tsallis-inf: An optimal algorithm for stochastic and adversarial
400 bandits. *Journal of Machine Learning Research*, 22(28):1–49, 2021.

401 Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Pro-*
402 *ceedings of the 20th International Conference on International Conference on Machine Learning*,
403 pages 928–935, 2003.

404 Checklist

- 405 1. For all authors...
 - 406 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
407 contributions and scope? [Yes]
 - 408 (b) Did you describe the limitations of your work? [No]
 - 409 (c) Did you discuss any potential negative societal impacts of your work? [No] This paper
410 is a theoretical paper.
 - 411 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
412 them? [Yes]
- 413 2. If you are including theoretical results...
 - 414 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Assump-
415 tion 1.
 - 416 (b) Did you include complete proofs of all theoretical results? [Yes] Omitted Proof can
417 be found in the supplemental material.
- 418 3. If you ran experiments...
 - 419 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
420 mental results (either in the supplemental material or as a URL)? [Yes] The experiment
421 code can be found in the supplemental material.
 - 422 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
423 were chosen)? [Yes] See Section 5.
 - 424 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
425 ments multiple times)? [No]
 - 426 (d) Did you include the total amount of compute and the type of resources used (e.g., type
427 of GPUs, internal cluster, or cloud provider)? [Yes]
- 428 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 429 (a) If your work uses existing assets, did you cite the creators? [Yes]
 - 430 (b) Did you mention the license of the assets? [Yes]
 - 431 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
 - 432 (d) Did you discuss whether and how consent was obtained from people whose data you’re
433 using/curating? [No] We used a public dataset.
 - 434 (e) Did you discuss whether the data you are using/curating contains personally identifiable
435 information or offensive content? [No]
- 436 5. If you used crowdsourcing or conducted research with human subjects...
 - 437 (a) Did you include the full text of instructions given to participants and screenshots, if
438 applicable? [No] Not applicable
 - 439 (b) DidNo describe any potential participant risks, with links to Institutional Review Board
440 (IRB) approvals, if applicable? [No]
 - 441 (c) DidNo include the estimated hourly wage paid to participants and the total amount
442 spent on participant compensation? [No]