

机器学习理论研究导引

大作业

庄镇华 502022370071 人工智能学院

作业提交注意事项

- (1) 本次作业提交截止时间为 **2023/07/02 23:59:59**, 截止时间后不再接收作业, 本次作业记零分;
- (2) 作业提交方式: 使用此 LaTeX 模板书写解答, 只需提交编译生成的 pdf 文件, 将 pdf 文件提交至南大网盘:
<https://box.nju.edu.cn/u/d/c1944351a2de4e93ada5/>
- (3) pdf 文件命名方式: 学号-姓名-作业号-v 版本号, 例 MG1900000-张三-5-v1; 如果需要更改已提交的解答, 请在截止时间之前提交新版本的解答, 并将版本号加一;
- (4) 未按照要求提交作业, 或 pdf 命名方式不正确, 将会被扣除部分作业分数.

摘要

本文研究了组分布鲁棒优化 (GDRO), 其目的是学习一个在 m 个不同分布表现良好的模型。首先, 我们将 GDRO 形式化为随机凸凹鞍点问题, 并提出了两种随机镜像下降 (SMD) 算法, 第一个算法在每次迭代中使用 1 个样本, 第二个算法每次使用 m 个样本并实现了用 $O(m(\log m)/\epsilon^2)$ 的样本复杂度来找到最优解, 这几乎与该问题的下界 $\Omega(m/\epsilon^2)$ 相匹配。然后, 我们考虑一个更实际的问题, 即从每个分布中采样的数量是不同的, 并提出了一种新的加权组分布鲁棒优化 (weighted DRO) 情景, 这使我们能够获得分布相关的收敛速度。用 n_i 表示第 i 个分布的样本预算, 并假设 $n_1 \geq n_2 \geq \dots \geq n_m$ 。为了解决这个问题, 我们在 SMD 中引入非均匀抽样, 使得样本预算满足期望, 并证明了第 i 个分布的超额风险以 $O(\sqrt{n_i \log m}/n_i)$ 的速率收敛。最后, 我们给出了未来展望, 即将 SA 算法扩展到具有类似形式的问题, 例如联邦学习和极小极大遗憾优化等问题。

1 引言

在经典统计机器学习中，我们的目标是 minimize 固定分布 \mathcal{P}_0 的风险，即

$$\min_{\mathbf{w} \in \mathcal{W}} \{R_0(\mathbf{w}) = \mathbb{E}_{\mathbf{z} \sim \mathcal{P}_0} [\ell(\mathbf{w}; \mathbf{z})]\}, \quad (1.1)$$

其中 $\mathbf{z} \in \mathcal{Z}$ 是从 \mathcal{P}_0 中抽取的样本， \mathcal{W} 代表假设类， $\ell(\mathbf{w}; \mathbf{z})$ 是度量模型 \mathbf{w} 在 \mathbf{z} 上预测误差的损失函数。近年来，用于优化 (1.1) 的各种算法相继被提出，并且可以分为两类：样本平均近似 (SAA) 和随机近似 (SA) 方法。在 SAA 方法中，我们最小化经验风险，即从 \mathcal{P}_0 中采样的一组样本的平均风险，而在 SA 方法中，我们利用目标函数 $R_0(\cdot)$ 的随机观测值直接解决原始问题。

然而，在单一分布上训练的模型可能缺乏鲁棒性，这是因为 (i) 尽管平均损失很低，但其仍可能在少数类中表现糟糕；(ii) 当在另一个分布上测试时，其性能可能会急剧下降。分布式鲁棒优化 (DRO) 提供了一种原理性的方法通过最小化分布 \mathcal{P}_0 邻域最坏情况下的损失来解决这些问题 [Ben-Tal et al., 2013]。最近，其已经被广泛应用于优化、统计学、运筹学和机器学习等领域。在本文中，我们考虑最新出现的问题，即组分布鲁棒优化 (GDRO) 问题，其优化有限分布上的最大风险 [Sagawa et al., 2020]。GDRO 可以表述为一个随机极小极大问题：

$$\min_{\mathbf{w} \in \mathcal{W}} \max_{i \in [m]} \{R_i(\mathbf{w}) = \mathbb{E}_{\mathbf{z} \sim \mathcal{P}_i} [\ell(\mathbf{w}; \mathbf{z})]\} \quad (1.2)$$

这里 $\mathcal{P}_1, \dots, \mathcal{P}_m$ 表示 m 个分布。一个启发性的例子是联邦学习，即在多个客户端部署一个集中式模型，并且每个客户端都面临着不同的数据分布。

由 [Nemirovski et al., 2009, §3.2] 所启发，我们可以将 (1.2) 转化为随机凸凹鞍点问题：

$$\min_{\mathbf{w} \in \mathcal{W}} \max_{\mathbf{q} \in \Delta_m} \left\{ \phi(\mathbf{w}, \mathbf{q}) = \sum_{i=1}^m q_i R_i(\mathbf{w}) \right\} \quad (1.3)$$

其中 $\Delta_m = \{\mathbf{q} \in \mathbb{R}^m : \mathbf{q} \geq 0, \sum_{i=1}^m q_i = 1\}$ 是 $(m-1)$ 维单纯形，然后可以利用随机镜像下降方法 (SMD) 来求解 (1.3)。

2 相关工作

分布鲁棒优化 (DRO) 源于 [Scarf, 1958] 的开创性工作，并且随着鲁棒优化的发展而得到了广泛的关注。它已经广泛地应用于各种机器学习任务，包括对抗训练 [Sinha et al., 2018]、算法公平性 [Hashimoto et al., 2018]、类别不均衡 [Xu et al., 2020]、长尾分布 [Samuel and Chechik, 2021]、标签移位 [Zhang et al., 2021] 等问题。

一般来说，DRO 反应了我们对目标分布的不确定性。为确保在分布扰动下的良好性能，它最小化不确定集中最坏分布的风险，即

$$\min_{\mathbf{w} \in \mathcal{W}} \sup_{\mathcal{P} \in \mathcal{S}(\mathcal{P}_0)} \{\mathbb{E}_{\mathbf{z} \sim \mathcal{P}} [\ell(\mathbf{w}; \mathbf{z})]\} \quad (2.1)$$

其中 $\mathcal{S}(\mathcal{P}_0)$ 表示 \mathcal{P}_0 周围概率分布的集合。

本文的研究重点是 (1.2)/(1.3) 中的 GDRO 问题，而不是 (2.1) 中的传统 DRO 问题。[Sagawa et al., 2020] 已经应用了 SMD [Nemirovski et al., 2009] 到 (1.3)，但由于梯度的方差上界较大，仅获得次优的样本复杂度。在后续工作中，[Haghtalab et al., 2022] 和 [Soma et al., 2022] 已经

通过重用样本和应用来自 MAB 的技术来降低样本复杂度，但他们的分析存在依赖性问题。为了处理不同分布的非均匀噪声，[Agarwal and Zhang, 2022] 提出了 GDRO 算法的变种——即极小极大遗憾优化算法 (MRO)，它将风险 $R_i(\mathbf{w})$ 替换为“超额风险” $R_i(\mathbf{w}) - \min_{\mathbf{w} \in \mathcal{W}} R_i(\mathbf{w})$ 。更一般地，可以在 DRO 中引入校准项，以防止任何单一分布支配最大值 [Słowiak and Bottou, 2022]。

3 组分布鲁棒优化 (GDRO) 的随机近似 (SA) 方法

在本节中，我们提供了针对 GDRO 问题的两种有效 SA 方法，它们的区别在于每轮采样的数目不同，即随机梯度的构建方式不同，进而导致不同的样本复杂度。

3.1 预备知识

首先，我们重温随机镜像下降 (SMD) 方法的设定 [Nemirovski et al., 2009]。我们为定义域 \mathcal{W} 配置距离生成函数 $\nu_w(\cdot)$ ，其关于范数 $\|\cdot\|_w$ 是 1-强凸的。我们定义与 $\nu_w(\cdot)$ 相关的 Bregman 距离为

$$B_w(\mathbf{u}, \mathbf{v}) = \nu_w(\mathbf{u}) - [\nu_w(\mathbf{v}) + \langle \nabla \nu_w(\mathbf{v}), \mathbf{u} - \mathbf{v} \rangle].$$

对于单纯形 Δ_m ，我们选择熵函数 $\nu_q(\mathbf{q}) = \sum_{i=1}^m q_i \ln q_i$ 作为距离生成函数，其关于向量 ℓ_1 -norm $\|\cdot\|_1$ 是 1-强凸的。类似的， $B_q(\cdot, \cdot)$ 是与 $\nu_q(\cdot)$ 关联的 Bregman 距离。

然后，我们介绍关于定义域和损失函数的一般假设。

假设 1. 定义域 \mathcal{W} 是凸的并且根据 $\nu_w(\cdot)$ 计算的直径上界为 D ，即

$$\max_{\mathbf{w} \in \mathcal{W}} \nu_w(\mathbf{w}) - \min_{\mathbf{w} \in \mathcal{W}} \nu_w(\mathbf{w}) \leq D^2. \quad (3.1)$$

对于 Δ_m ，容易验证其由熵函数计算的直径以 $\sqrt{\ln m}$ 为界。为了简化表示，我们假设损失属于 $[0, 1]$ ，并且其梯度也是有界的。

假设 2. 对所有 $i \in [m]$ ，我们有

$$0 \leq \ell(\mathbf{w}; \mathbf{z}) \leq 1, \quad \forall \mathbf{w} \in \mathcal{W}, \mathbf{z} \sim \mathcal{P}_i. \quad (3.2)$$

假设 3. 对所有 $i \in [m]$ ，我们有

$$\|\nabla_{\mathbf{w}} \ell(\mathbf{w}; \mathbf{z})\|_{w,*} \leq G, \quad \forall \mathbf{w} \in \mathcal{W}, \mathbf{z} \sim \mathcal{P}_i \quad (3.3)$$

其中 $\|\cdot\|_{w,*}$ 是 $\|\cdot\|_w$ 的对偶范数。

最后，我们讨论性能度量标准。为了分析收敛性，我们通过如下误差来衡量 (1.3) 的近似解 $(\bar{\mathbf{w}}, \bar{\mathbf{q}})$ 的质量

$$\epsilon_\phi(\bar{\mathbf{w}}, \bar{\mathbf{q}}) = \max_{\mathbf{q} \in \Delta_m} \phi(\bar{\mathbf{w}}, \mathbf{q}) - \min_{\mathbf{w} \in \mathcal{W}} \phi(\mathbf{w}, \bar{\mathbf{q}}) \quad (3.4)$$

它可以支配 $\bar{\mathbf{w}}$ 对于原始问题 (1.2) 的最优性，因为

$$\begin{aligned} \max_{i \in [m]} R_i(\bar{\mathbf{w}}) - \min_{\mathbf{w} \in \mathcal{W}} \max_{i \in [m]} R_i(\mathbf{w}) &= \max_{\mathbf{q} \in \Delta_m} \sum_{i=1}^m q_i R_i(\bar{\mathbf{w}}) - \min_{\mathbf{w} \in \mathcal{W}} \max_{\mathbf{q} \in \Delta_m} \sum_{i=1}^m q_i R_i(\mathbf{w}) \\ &\leq \max_{\mathbf{q} \in \Delta_m} \sum_{i=1}^m q_i R_i(\bar{\mathbf{w}}) - \min_{\mathbf{w} \in \mathcal{W}} \sum_{i=1}^m \bar{q}_i R_i(\mathbf{w}) = \epsilon_\phi(\bar{\mathbf{w}}, \bar{\mathbf{q}}). \end{aligned} \quad (3.5)$$

3.2 GDRO 的在线优化方法

为了应用 SMD，关键是构造 (1.3) 中函数 $\phi(\mathbf{w}, \mathbf{q})$ 的随机梯度。我们首先给出它关于 \mathbf{w} 和 \mathbf{q} 的真实梯度：

$$\nabla_{\mathbf{w}}\phi(\mathbf{w}, \mathbf{q}) = \sum_{i=1}^m q_i \nabla R_i(\mathbf{w}), \text{ and } \nabla_{\mathbf{q}}\phi(\mathbf{w}, \mathbf{q}) = [R_1(\mathbf{w}), \dots, R_m(\mathbf{w})]^\top. \quad (3.6)$$

在第 t 轮中，用 \mathbf{w}_t 和 \mathbf{q}_t 表示当前解。我们首先从 $[m]$ 中均匀选择一个索引 i_t ，然后从分布 \mathcal{P}_{i_t} 中抽取一个样本 $\mathbf{z}_t^{(i_t)}$ ，并将随机梯度定义为

$$\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t) = q_{t,i_t} m \nabla \ell(\mathbf{w}_t; \mathbf{z}_t^{(i_t)}), \text{ and } \mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t) = [0, \dots, m \ell(\mathbf{w}_t; \mathbf{z}_t^{(i_t)}), \dots, 0]^\top \quad (3.7)$$

显然，它们是真实梯度的无偏估计：

$$\mathbb{E}_{t-1}[\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t)] = \nabla_{\mathbf{w}}\phi(\mathbf{w}_t, \mathbf{q}_t), \text{ and } \mathbb{E}_{t-1}[\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)] = \nabla_{\mathbf{q}}\phi(\mathbf{w}_t, \mathbf{q}_t)$$

其中 $\mathbb{E}_{t-1}[\cdot]$ 表示直到第 $t-1$ 轮以随机性为条件的期望值。

然后，我们应用 SMD 来更新 \mathbf{w}_t 和 \mathbf{q}_t ：

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \{ \eta_w \langle \mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t), \mathbf{w} - \mathbf{w}_t \rangle + B_w(\mathbf{w}, \mathbf{w}_t) \}, \quad (3.8)$$

$$\mathbf{q}_{t+1} = \operatorname{argmin}_{\mathbf{q} \in \Delta_m} \{ \eta_q \langle -\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t), \mathbf{q} - \mathbf{q}_t \rangle + B_q(\mathbf{q}, \mathbf{q}_t) \} \quad (3.9)$$

其中 $\eta_w > 0$ 和 $\eta_q > 0$ 是两个步长。 \mathbf{w}_t 的更新规则取决于距离生成函数 $\nu_w(\cdot)$ 的选择。例如，如果 $\nu_w(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2$ ，则 (3.8) 变为随机梯度下降 (SGD)，即

$$\mathbf{w}_{t+1} = \Pi_{\mathcal{W}}[\mathbf{w}_t - \eta_w \mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t)]$$

其中 $\Pi_{\mathcal{W}}[\cdot]$ 表示到 \mathcal{W} 中最近点的欧式投影。由于 $B_q(\mathbf{q}, \mathbf{q}_t)$ 是根据负熵定义的，因此 (3.9) 等价于

$$q_{t+1,i} = \frac{q_{t,i} \exp(\eta_q \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}))}{\sum_{j=1}^m q_{t,j} \exp(\eta_q \ell(\mathbf{w}_t; \mathbf{z}_t^{(j)}))}, \quad \forall i \in [m] \quad (3.10)$$

即为应用于最大化问题的 Hedge 算法 [Freund and Schapire, 1997]。一开始，我们设定 $\mathbf{w}_1 = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \nu_w(\mathbf{w})$ ， $\mathbf{q}_1 = \frac{1}{m} \mathbf{1}_m$ ，其中 $\mathbf{1}_m$ 是由 1 组成的 m 维向量。在最后一步，我们返回平均迭代 $\bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$ 和 $\bar{\mathbf{q}} = \frac{1}{T} \sum_{t=1}^T \mathbf{q}_t$ 作为最终解。整体的程序流程在算法 1 中。

算法 1 GDRO 的在线优化算法

输入： 两个步长： η_w 和 η_q

- 1: 初始化 $\mathbf{w}_1 = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \nu_w(\mathbf{w})$ ， $\mathbf{q}_1 = [1/m, \dots, 1/m]^\top \in \mathbb{R}^m$
 - 2: **for** $t = 1$ to T **do**
 - 3: $i_t \sim \text{Uniform}(1, \dots, m)$
 - 4: 从分布 \mathcal{P}_{i_t} 中抽取一个样本 $\mathbf{z}_t^{(i_t)}$
 - 5: 构造公式 (3.7) 中定义的随机梯度
 - 6: 分别根据 (3.8) 和 (3.9) 更新 \mathbf{w}_t 和 \mathbf{q}_t
 - 7: **end for**
 - 8: **return** $\bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$ 和 $\bar{\mathbf{q}} = \frac{1}{T} \sum_{t=1}^T \mathbf{q}_t$
-

基于随机凸凹优化的 SMD 理论保证 [Nemirovski et al., 2009, §3.1], 对于算法 1, 我们有如下定理:

定理 1. 基于假设 1, 2 和 3, 并且设定算法 1 中的 $\eta_w = \frac{D^2}{m} \sqrt{\frac{8}{5T(D^2G^2 + \ln m)}}$, $\eta_q = \frac{\ln m}{m} \sqrt{\frac{8}{5T(D^2G^2 + \ln m)}}$, 我们有

$$\mathbb{E}[\epsilon_\phi(\bar{\mathbf{w}}, \bar{\mathbf{q}})] \leq 2m \sqrt{\frac{10(D^2G^2 + \ln m)}{T}}$$

并且以至少 $1 - \delta$ 的概率有,

$$\epsilon_\phi(\bar{\mathbf{w}}, \bar{\mathbf{q}}) \leq 2m \left(\sqrt{10} + 2\sqrt{2 \ln 1/\delta} \right) \sqrt{\frac{D^2G^2 + \ln m}{T}}.$$

备注 1. 定理 1 表明算法 1 达到 $O(m\sqrt{(\log m)/T})$ 的收敛速度。由于每次迭代消耗 1 个样本, 因此样本复杂度为 $O(m^2(\log m)/\epsilon^2)$ 。

3.3 GDRO 的随机镜像下降方法

接下来, 我们采用另一种方式构造 (1.3) 中函数 $\phi(\mathbf{w}, \mathbf{q})$ 的随机梯度。在第 t 轮中, 用 \mathbf{w}_t 和 \mathbf{q}_t 表示当前解。我们从每个分布 \mathcal{P}_i 中抽取一个样本 $\mathbf{z}_t^{(i)}$, 并将随机梯度定义为

$$\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t) = \sum_{i=1}^m q_{t,i} \nabla \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}), \text{ and } \mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t) = [\ell(\mathbf{w}_t; \mathbf{z}_t^{(1)}), \dots, \ell(\mathbf{w}_t; \mathbf{z}_t^{(m)})]^\top. \quad (3.11)$$

显然, 它们也是真实梯度的无偏估计:

$$\mathbb{E}_{t-1}[\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t)] = \nabla_{\mathbf{w}} \phi(\mathbf{w}_t, \mathbf{q}_t), \text{ and } \mathbb{E}_{t-1}[\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)] = \nabla_{\mathbf{q}} \phi(\mathbf{w}_t, \mathbf{q}_t)$$

其中 $\mathbb{E}_{t-1}[\cdot]$ 表示直到第 $t-1$ 轮以随机性为条件的期望值。

然后, 我们根据 (3.8) 和 (3.9) 应用 SMD 来更新 \mathbf{w}_t 和 \mathbf{q}_t 。初始时, 我们仍设定 $\mathbf{w}_1 = \arg\min_{\mathbf{w} \in \mathcal{W}} \nu_w(\mathbf{w})$, $\mathbf{q}_1 = \frac{1}{m} \mathbf{1}_m$ 。在最后一步, 我们返回平均迭代 $\bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$ 和 $\bar{\mathbf{q}} = \frac{1}{T} \sum_{t=1}^T \mathbf{q}_t$ 作为最终解。整体的程序流程在算法 2 中。

算法 2 GDRO 的随机镜像下降算法

输入: 两个步长: η_w 和 η_q

- 1: 初始化 $\mathbf{w}_1 = \arg\min_{\mathbf{w} \in \mathcal{W}} \nu_w(\mathbf{w})$, $\mathbf{q}_1 = [1/m, \dots, 1/m]^\top \in \mathbb{R}^m$
 - 2: **for** $t = 1$ to T **do**
 - 3: 对于每个 $i \in [m]$, 从分布 \mathcal{P}_i 中抽取一个样本 $\mathbf{z}_t^{(i)}$
 - 4: 构造公式 (3.11) 中定义的随机梯度
 - 5: 分别根据 (3.8) 和 (3.9) 更新 \mathbf{w}_t 和 \mathbf{q}_t
 - 6: **end for**
 - 7: **return** $\bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$ 和 $\bar{\mathbf{q}} = \frac{1}{T} \sum_{t=1}^T \mathbf{q}_t$
-

类似的, 对于算法 2, 我们有如下定理:

定理 2. 基于假设 1, 2 和 3, 并且设定算法 2 中的 $\eta_w = D^2 \sqrt{\frac{8}{5T(D^2G^2 + \ln m)}}$, $\eta_q = (\ln m) \sqrt{\frac{8}{5T(D^2G^2 + \ln m)}}$, 我们有

$$\mathbb{E}[\epsilon_\phi(\bar{\mathbf{w}}, \bar{\mathbf{q}})] \leq 2 \sqrt{\frac{10(D^2G^2 + \ln m)}{5T}}$$

并且以至少 $1 - \delta$ 的概率有,

$$\epsilon_\phi(\bar{\mathbf{w}}, \bar{\mathbf{q}}) \leq 2 \left(\sqrt{10} + 2\sqrt{2 \ln 1/\delta} \right) \sqrt{\frac{D^2 G^2 + \ln m}{T}}.$$

备注 2. 定理 2 表明算法 2 达到 $O(\sqrt{(\log m)/T})$ 的收敛速度。由于每次迭代消耗 m 个样本, 因此样本复杂度为 $O(m(\log m)/\epsilon^2)$, 几乎与下界 $\Omega(m/\epsilon^2)$ 匹配 [Soma et al., 2022, Theorem 5]。

备注 3. 算法 1 的收敛速度比算法 2 慢, 并且其随机梯度 (3.7) 的对偶范数是算法 2 的 m 倍, 导致其样本复杂度是算法 2 的 m 倍, 为 $O(m^2(\log m)/\epsilon^2)$ 。

4 加权组分布鲁棒优化 (Weighted GDRO) 和随机近似 (SA) 方法

正如我们在第 3 节中所做的那样, 在为 GDRO 问题设计 SA 方法时, 通常假设算法可以自由地从每个分布中抽取样本。然而, 由于从不同分布收集数据的成本不同, 这种假设在实践中可能不成立 [Radivojac et al., 2004]。在本节中, 我们将研究从每个分布中采样的样本数量可能不同的情况, 用 n_i 表示可以从分布 \mathcal{P}_i 中采样的数量, 不失一般性, 我们假设 $n_1 \geq n_2 \geq \dots \geq n_m$ 。注意, 我们有一个简单的**基线**, 它只运行 n_m 次算法 2, 并且优化误差为 $\epsilon_\phi(\bar{\mathbf{w}}, \bar{\mathbf{q}}) = O(\sqrt{(\log m)/n_m})$ 。

4.1 非均匀采样的随机镜像下降算法

为了满足预算, 我们建议将非均匀采样纳入 SMD 算法。在进入技术细节之前, 我们首先解释使用非均匀采样的主要思想。一种方法是在每次迭代中以 $p_i = n_i/n_1$ 的概率从每个分布 \mathcal{P}_i 中抽取 1 个样本。那么, 在 n_1 次迭代之后, 从 \mathcal{P}_i 中采样的期望数目将是 $n_1 p_i = n_i$, 因此预算在期望意义上满足。

具体地, 在第 t 轮中, 我们首先生成一组伯努利随机变量 $\{b_t^{(1)}, \dots, b_t^{(m)}\}$, 其中 $\Pr[b_t^{(i)} = 1] = p_i$, 用于确定是否从第 i 分布采样。如果 $b_t^{(i)} = 1$, 我们从 \mathcal{P}_i 中提取样本 $\mathbf{z}_t^{(i)}$ 。现在, 问题变成如何利用这些样本构建随机梯度, 令 $\mathcal{C}_t = \{i | b_t^{(i)} = 1\}$ 表示所选分布的索引集合。如果我们求解 (1.3) 中的原始问题, 那么随机梯度应按以下方式构造以确保无偏性。

$$\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t) = \sum_{i \in \mathcal{C}_t} \frac{q_{t,i}}{p_i} \nabla \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}), \text{ and } [\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)]_i = \begin{cases} \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)})/p_i, & i \in \mathcal{C}_t \\ 0, & \text{otherwise} \end{cases} \quad (4.1)$$

然后, 可以利用 SMD 来更新 \mathbf{w}_t 和 \mathbf{q}_t 。为了分析优化误差, 我们需要限制 (4.1) 中随机梯度的范数。为此, 我们有 $\|\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t)\|_{w,*} \leq G n_1/n_m$ 和 $\|\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)\|_\infty \leq n_1/n_m$ 。根据定理 2 的论证, 我们可以证明误差为 $\epsilon_\phi(\bar{\mathbf{w}}, \bar{\mathbf{q}}) = O(\sqrt{(\log m)/n_1} \cdot n_1/n_m) = O(\sqrt{n_1 \log m}/n_m)$, 甚至大于基线的 $O(\sqrt{(\log m)/n_m})$ 误差。

在下文中, 我们证明了上述过程的简单转化仍可以产生与基线互补的有意义的结果。我们观察到 (4.1) 中随机梯度的范数较大是由概率的倒数 $1/p_i$ 引起的。一个自然的想法是忽略 $1/p_i$, 并定义以下随机梯度:

$$\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t) = \sum_{i \in \mathcal{C}_t} q_{t,i} \nabla \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}), \text{ and } [\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)]_i = \begin{cases} \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}), & i \in \mathcal{C}_t \\ 0, & \text{otherwise.} \end{cases} \quad (4.2)$$

通过这种方式，它们不再是 (1.3) 的随机梯度，而是可以被视为加权 GDRO 问题的随机梯度：

$$\min_{\mathbf{w} \in \mathcal{W}} \max_{\mathbf{q} \in \Delta_m} \left\{ \varphi(\mathbf{w}, \mathbf{q}) = \sum_{i=1}^m q_i p_i R_i(\mathbf{w}) \right\} \quad (4.3)$$

其中每个风险 $R_i(\cdot)$ 通过因子 p_i 缩放。基于 (4.2) 中的梯度，我们仍然使用 (3.8) 和 (3.9) 来更新 \mathbf{w}_t 和 \mathbf{q}_t ，并在算法 3 中总结了完整的过程。

算法 3 加权 GDRO 的随机镜像下降算法

输入： 两个步长: η_w 和 η_q

- 1: 初始化 $\mathbf{w}_1 = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \nu_w(\mathbf{w})$, $\mathbf{q}_1 = [1/m, \dots, 1/m]^\top \in \mathbb{R}^m$
 - 2: **for** $t = 1$ to n_1 **do**
 - 3: 对于每个 $i \in [m]$, 生成一个伯努利随机变量 $b_t^{(i)}$, 其中 $\Pr[b_t^{(i)} = 1] = p_i$, 如果 $b_t^{(i)} = 1$, 则从分布 \mathcal{P}_i 中抽取一个样本 $\mathbf{z}_t^{(i)}$
 - 4: 构造 (4.2) 中定义的随机梯度
 - 5: 分别根据 (3.8) 和 (3.9) 更新 \mathbf{w}_t 和 \mathbf{q}_t
 - 6: **end for**
 - 7: **return** $\bar{\mathbf{w}} = \frac{1}{n_1} \sum_{t=1}^{n_1} \mathbf{w}_t$ 和 $\bar{\mathbf{q}} = \frac{1}{n_1} \sum_{t=1}^{n_1} \mathbf{q}_t$
-

我们省略算法 3 对于 (4.3) 的优化误差，因为它具有与定理 2 完全相同的形式。我们真正感兴趣的是它在多重分布上解的理论保证。为此，我们有以下定理：

定理 3. 基于假设 1, 2 和 3, 并且设定算法 3 中的 $\eta_w = D^2 \sqrt{\frac{8}{5n_1(D^2G^2 + \ln m)}}$, $\eta_q = (\ln m) \sqrt{\frac{8}{5n_1(D^2G^2 + \ln m)}}$, 则以至少 $1 - \delta$ 的概率, 有

$$R_i(\bar{\mathbf{w}}) - \frac{n_1}{n_i} p_\varphi^* \leq \frac{1}{p_i} \mu(\delta) \sqrt{\frac{D^2 G^2 + \ln m}{n_1}} = \mu(\delta) \frac{\sqrt{(D^2 G^2 + \ln m) n_1}}{n_i}, \quad \forall i \in [m]$$

其中 p_φ^* 是式 (4.3) 的最优值, 并且 $\mu(\delta) = 2 \left(\sqrt{10} + 2\sqrt{2 \ln 1/\delta} \right)$.

备注 4. 我们观察到算法 3 表现出与分布相关的收敛性: 样本数 n_i 越大, 目标风险 $n_1 p_\varphi^*/n_i$ 越小, 收敛速度 $O(\sqrt{n_1 \log m}/n_i)$ 越快。注意, 它的收敛率总是优于使用 (4.1) 作为梯度的 SMD 方法 $O(\sqrt{n_1 \log m}/n_m)$ 的收敛率。此外, 当 $n_i \geq \sqrt{n_1 n_m}$ 时, 它比基线收敛得更快。特别的, 对于分布 \mathcal{P}_1 , 算法 3 获得了 $O(\sqrt{(\log m)/n_1})$ 的收敛速率, 几乎和从单个分布学习的最优 $O(\sqrt{1/n_1})$ 收敛速率相匹配。我们注意到, 很难描述 p_φ^* 的阶, 以前的研究通常将这些量视为极小值甚至为 0 [Agarwal and Zhang, 2022, Corollary 9]。

5 定理证明

在这一节中, 我们给出主要定理的证明。

5.1 随机镜像下降算法定理的证明

定理 1 和 2 的证明基于 [Nemirovski et al., 2009] 的 §2.3 节——随机镜像下降算法, 因此为了完整性, 我们首先证明 §2.3 节中的主要结果。

算法 4 随机镜像下降算法**输入:** 步长: η

- 1: 初始化 $\mathbf{w}_1 = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \nu_w(\mathbf{w})$
- 2: **for** $t = 1$ to T **do**
- 3: 从分布 \mathcal{P} 中抽取一个样本 \mathbf{z}_t
- 4: 构造随机梯度 $\nabla \ell(\mathbf{w}_t; \mathbf{z}_t)$
- 5: 根据下式更新 \mathbf{w}_t

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \{ \eta \langle \nabla \ell(\mathbf{w}_t; \mathbf{z}_t), \mathbf{w} - \mathbf{w}_t \rangle + B_w(\mathbf{w}, \mathbf{w}_t) \}$$

6: **end for**7: **return** $\bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

定理 4. 基于假设 1, 2 和 3, 在算法 4 中设定 $\eta = \frac{\sqrt{2D}}{G\sqrt{T}}$, 我们有

$$\mathbb{E}[R(\bar{\mathbf{w}}) - R^*] \leq DG\sqrt{\frac{2}{T}}. \quad (5.1)$$

此外, 以至少 $1 - \delta$ 的概率有,

$$R(\bar{\mathbf{w}}) - R^* \leq DG\sqrt{\frac{2}{T}} \left(1 + 4\sqrt{\ln \frac{1}{\delta}} \right). \quad (5.2)$$

证明: 设 $\mathbf{w}_* \in \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} R(\mathbf{w})$ 是最小化 $R(\cdot)$ 的最优解. 根据镜像下降的性质, 即 Nemirovski et al. [2009] 的引理 2.1, 我们有

$$\begin{aligned} & \eta \langle \nabla \ell(\mathbf{w}_t; \mathbf{z}_t), \mathbf{w}_t - \mathbf{w}_* \rangle \\ & \leq B_w(\mathbf{w}_*, \mathbf{w}_t) - B_w(\mathbf{w}_*, \mathbf{w}_{j+1}) + \frac{\eta^2}{2} \|\nabla \ell(\mathbf{w}_t; \mathbf{z}_t)\|_{w,*}^2 \\ & \stackrel{(3.3)}{\leq} B_w(\mathbf{w}_*, \mathbf{w}_t) - B_w(\mathbf{w}_*, \mathbf{w}_{j+1}) + \frac{\eta^2 G^2}{2}. \end{aligned} \quad (5.3)$$

因此, 我们有

$$\begin{aligned} & \eta \langle \nabla R(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}_* \rangle \\ & = \eta \langle \nabla \ell(\mathbf{w}_t; \mathbf{z}_t), \mathbf{w}_t - \mathbf{w}_* \rangle + \eta \langle \nabla R(\mathbf{w}_t) - \nabla \ell(\mathbf{w}_t; \mathbf{z}_t), \mathbf{w}_t - \mathbf{w}_* \rangle \\ & \stackrel{(5.3)}{\leq} B_w(\mathbf{w}_*, \mathbf{w}_t) - B_w(\mathbf{w}_*, \mathbf{w}_{j+1}) + \frac{\eta^2 G^2}{2} + \eta \langle \nabla R(\mathbf{w}_t) - \nabla \ell(\mathbf{w}_t; \mathbf{z}_t), \mathbf{w}_t - \mathbf{w}_* \rangle. \end{aligned}$$

对于 $t = 1, \dots, T$, 将上述不等式求和, 我们有

$$\begin{aligned} & \sum_{t=1}^T \eta \langle \nabla R(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}_* \rangle \\ & \leq B_w(\mathbf{w}_*, \mathbf{w}_1) + \frac{G^2}{2} \sum_{t=1}^T \eta^2 + \sum_{t=1}^T \eta \langle \nabla R(\mathbf{w}_t) - \nabla \ell(\mathbf{w}_t; \mathbf{z}_t), \mathbf{w}_t - \mathbf{w}_* \rangle \\ & \stackrel{(3.1)}{\leq} D^2 + \frac{G^2}{2} \sum_{t=1}^T \eta^2 + \sum_{t=1}^T \eta \langle \nabla R(\mathbf{w}_t) - \nabla \ell(\mathbf{w}_t; \mathbf{z}_t), \mathbf{w}_t - \mathbf{w}_* \rangle. \end{aligned} \quad (5.4)$$

由于风险函数的凸性，我们有

$$\begin{aligned}
R(\bar{\mathbf{w}}) - R(\mathbf{w}_*) &= R\left(\sum_{t=1}^T \frac{\mathbf{w}_t}{T}\right) - R(\mathbf{w}_*) \\
&\leq \left(\sum_{t=1}^T \frac{1}{T} R(\mathbf{w}_t)\right) - R(\mathbf{w}_*) = \sum_{t=1}^T \frac{1}{T} (R(\mathbf{w}_t) - R(\mathbf{w}_*)) \\
&\leq \sum_{t=1}^T \frac{1}{T} \langle \nabla R(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w}_* \rangle \\
&\stackrel{(5.4)}{\leq} \frac{D^2 + \frac{G^2}{2} \sum_{t=1}^T \eta^2 + \sum_{t=1}^T \eta \langle \nabla R(\mathbf{w}_t) - \nabla \ell(\mathbf{w}_t; \mathbf{z}_t), \mathbf{w}_t - \mathbf{w}_* \rangle}{\eta T}.
\end{aligned} \tag{5.5}$$

定义

$$\delta_t = \eta \langle \nabla R(\mathbf{w}_t) - \nabla \ell(\mathbf{w}_t; \mathbf{z}_t), \mathbf{w}_t - \mathbf{w}_* \rangle. \tag{5.6}$$

由于 \mathbf{w}_t 和 \mathbf{w}_* 不依赖于 \mathbf{z}_t ，因此

$$\mathbb{E}_{t-1}[\delta_t] = \eta \langle \mathbb{E}_{t-1}[\nabla R(\mathbf{w}_t) - \nabla \ell(\mathbf{w}_t; \mathbf{z}_t)], \mathbf{w}_t - \mathbf{w}_* \rangle = 0.$$

因此， $\delta_1, \dots, \delta_t$ 是鞅差序列。

5.1.1 期望上界

对 (5.5) 式左右两边取期望，我们有

$$\mathbb{E}[R(\bar{\mathbf{w}}) - R(\mathbf{w}_*)] \leq \frac{2D^2 + TG^2\eta^2}{2\eta T}. \tag{5.7}$$

设 $\eta = \frac{\sqrt{2}D}{G\sqrt{T}}$ ，根据基本不等式，我们有

$$\mathbb{E}[R(\bar{\mathbf{w}}) - R(\mathbf{w}_*)] \leq DG\sqrt{\frac{2}{T}}. \tag{5.8}$$

因而 (5.1) 得证。

5.1.2 高概率上界

我们不加证明的给出如下引理，即鞅的 Hoeffding-Azuma 不等式。

引理 1. 设 V_1, V_2, \dots 是关于某个序列 X_1, X_2, \dots 的鞅差序列，其中 $V_i \in [A_i, A_i + c_i]$ ， A_i 是随机变量， c_i 是常数。如果 $S_n = \sum_{i=1}^n V_i$ ，那么对于任意 $t > 0$ ，

$$\Pr[S_n > t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right).$$

为了应用上面的引理，我们需要证明 $|\delta_t|$ 是有界的：

$$\begin{aligned}
\|\nabla R(\mathbf{w}_t) - \nabla \ell(\mathbf{w}_t; \mathbf{z}_t)\|_{w,*} &\leq \|\nabla R(\mathbf{w}_t)\|_{w,*} + \|\nabla \ell(\mathbf{w}_t; \mathbf{z}_t)\|_{w,*} \\
&\leq \mathbb{E}_{t-1} \left[\|\nabla \ell(\mathbf{w}_t; \mathbf{z}_t)\|_{w,*} \right] + \|\nabla \ell(\mathbf{w}_t; \mathbf{z}_t)\|_{w,*} \stackrel{(3.3)}{\leq} 2G,
\end{aligned} \tag{5.9}$$

$$\begin{aligned}
\|\mathbf{w}_t - \mathbf{w}_*\|_w &\leq \|\mathbf{w}_t - \mathbf{w}_1\|_w + \|\mathbf{w}_1 - \mathbf{w}_*\|_w \\
&\leq \left(\sqrt{B_w(\mathbf{w}_t, \mathbf{w}_1)} + \sqrt{B_w(\mathbf{w}_*, \mathbf{w}_1)} \right) \stackrel{(3.1)}{\leq} 2D.
\end{aligned} \tag{5.10}$$

因此,

$$\begin{aligned}
|\delta_t| &= \eta |\langle \nabla R(\mathbf{w}_t) - \nabla \ell(\mathbf{w}_t; \mathbf{z}_t), \mathbf{w}_t - \mathbf{w}_* \rangle| \\
&\leq \eta \|\nabla R(\mathbf{w}_t) - \nabla \ell(\mathbf{w}_t; \mathbf{z}_t)\|_{w,*} \|\mathbf{w}_t - \mathbf{w}_*\|_w \stackrel{(5.9), (5.10)}{\leq} 4\eta DG.
\end{aligned} \tag{5.11}$$

根据引理 1, 以至少 $1 - \delta$ 的概率, 我们有

$$\sum_{t=1}^T \delta_t \leq 4DG \sqrt{2 \sum_{t=1}^T \eta^2 \ln \frac{1}{\delta}}. \tag{5.12}$$

将 (5.12) 代入 (5.5), 以至少 $1 - \delta$ 的概率, 我们有

$$\begin{aligned}
R(\bar{\mathbf{w}}) - R(\mathbf{w}_*) &\leq \frac{D^2 + \frac{G^2}{2} \sum_{t=1}^T \eta^2 + 4DG \sqrt{2 \sum_{t=1}^T \eta^2 \ln \frac{1}{\delta}}}{\sum_{t=1}^T \eta} \\
&\leq DG \sqrt{\frac{2}{T}} + 4DG \sqrt{\frac{2 \ln \frac{1}{\delta}}{T}} \\
&= DG \sqrt{\frac{2}{T}} \left(1 + 4 \sqrt{\ln \frac{1}{\delta}} \right)
\end{aligned}$$

5.2 定理 1 和 2 的证明

首先证明定理 2。证明基于 [Nemirovski et al., 2009] 的引理 3.1 和命题 3.2, 即 §3.1 节——求解随机凸凹鞍点问题的镜像下降算法。

5.2.1 合并 (3.8) 和 (3.9) 中的更新规则

虽然 \mathbf{w} 和 \mathbf{q} 是通过 SMD 的两个实例来更新的, 但可以通过将 \mathbf{w} 和 \mathbf{q} 连接为单个变量 $[\mathbf{w}; \mathbf{q}] \in \mathcal{W} \times \Delta_m$, 并重新定义范数和距离生成函数, 以达到只使用 1 个 SMD 实例的目的 [Nemirovski et al., 2009, §3.1]。设 \mathcal{E} 是 \mathcal{W} 所在的空间。我们给笛卡尔积 $\mathcal{E} \times \mathbb{R}^m$ 配备以下范数和对偶范数:

$$\|[\mathbf{w}; \mathbf{q}]\| = \sqrt{\frac{1}{2D^2} \|\mathbf{w}\|_w^2 + \frac{1}{2 \ln m} \|\mathbf{q}\|_1^2}, \text{ and } \|[\mathbf{u}, \mathbf{v}]\|_* = \sqrt{2D^2 \|\mathbf{u}\|_{w,*}^2 + 2 \|\mathbf{v}\|_\infty^2 \ln m}. \tag{5.13}$$

我们使用符号 $\mathbf{x} = [\mathbf{w}; \mathbf{q}]$, 并为集合 $\mathcal{W} \times \Delta_m$ 选择距离生成函数

$$\nu(\mathbf{x}) = \nu([\mathbf{w}; \mathbf{q}]) = \frac{1}{2D^2} \nu_w(\mathbf{w}) + \frac{1}{2 \ln m} \nu_q(\mathbf{q}). \tag{5.14}$$

易知 $\nu(\mathbf{x})$ 关于范数 $\|\cdot\|$ 是 1-强凸的。设 $B(\cdot, \cdot)$ 是与 $\nu(\cdot)$ 相关联的 Bregman 距离:

$$\begin{aligned}
B(\mathbf{x}_1, \mathbf{x}_2) &= \nu(\mathbf{x}_1) - [\nu(\mathbf{x}_2) + \langle \nabla \nu(\mathbf{x}_2), \mathbf{x}_1 - \mathbf{x}_2 \rangle] \\
&= \frac{1}{2D^2} (\nu_w(\mathbf{w}_1) - [\nu_w(\mathbf{w}_2) + \langle \nabla \nu_w(\mathbf{w}_2), \mathbf{w}_1 - \mathbf{w}_2 \rangle]) \\
&\quad + \frac{1}{2 \ln m} (\nu_q(\mathbf{q}_1) - [\nu_q(\mathbf{q}_2) + \langle \nabla \nu_q(\mathbf{q}_2), \mathbf{q}_1 - \mathbf{q}_2 \rangle]) \\
&= \frac{1}{2D^2} B_w(\mathbf{w}_1, \mathbf{w}_2) + \frac{1}{2 \ln m} B_q(\mathbf{q}_1, \mathbf{q}_2)
\end{aligned} \tag{5.15}$$

其中 $\mathbf{x}_1 = [\mathbf{w}_1, \mathbf{q}_1]$ 并且 $\mathbf{x}_2 = [\mathbf{w}_2, \mathbf{q}_2]$.

然后, 我们可以证明域 $\mathcal{W} \times \Delta_m$ 是有界的, 因为

$$\max_{(\mathbf{w}, \mathbf{q}) \in \mathcal{W} \times \Delta_m} B([\mathbf{w}, \mathbf{q}], [\mathbf{w}_1, \mathbf{q}_1]) = \frac{1}{2D^2} \max_{\mathbf{w} \in \mathcal{W}} B_w(\mathbf{w}, \mathbf{w}_1) + \frac{1}{2 \ln m} \max_{\mathbf{q} \in \Delta_m} B_q(\mathbf{q}, \mathbf{q}_1) \stackrel{(3.1)}{\leq} 1. \quad (5.16)$$

接着, 我们考虑用以下 SMD 来更新 \mathbf{x}_t :

$$\mathbf{x}_{t+1} = \underset{\mathbf{x} \in \mathcal{W} \times \Delta_m}{\operatorname{argmin}} \left\{ \eta \langle [\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t); -\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)], \mathbf{x} - \mathbf{x}_t \rangle + B(\mathbf{x}, \mathbf{x}_t) \right\} \quad (5.17)$$

这里 $\eta > 0$ 是步长。初始时, 我们设置 $\mathbf{x}_1 = \underset{\mathbf{x} \in \mathcal{W} \times \Delta_m}{\operatorname{argmin}} \nu(\mathbf{x}) = [\mathbf{w}_1, \mathbf{q}_1]$ 。从 (5.15) 中的 Bregman 距离的分解过程, 我们得出通过设置如下步长, (5.17) 等价于 (3.8) 和 (3.9)。

$$\eta_w = 2\eta D^2, \text{ and } \eta_q = 2\eta \ln m. \quad (5.18)$$

5.2.2 无偏随机梯度的 SMD 算法收敛性分析

为了简化符号, 我们定义

$$\begin{aligned} F(\mathbf{w}_t, \mathbf{q}_t) &= [\nabla_{\mathbf{w}} \phi(\mathbf{w}_t, \mathbf{q}_t), -\nabla_{\mathbf{q}} \phi(\mathbf{w}_t, \mathbf{q}_t)] \\ &= \left[\sum_{i=1}^m q_{t,i} \nabla R_i(\mathbf{w}_t), -[R_1(\mathbf{w}_t), \dots, R_m(\mathbf{w}_t)]^\top \right] \end{aligned} \quad (5.19)$$

其包含在 $(\mathbf{w}_t, \mathbf{q}_t)$ 处的 $\phi(\cdot, \cdot)$ 的真实梯度, 并且

$$\begin{aligned} \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t) &= [\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t); -\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)] \\ &\stackrel{(3.11)}{=} \left[\sum_{i=1}^m q_{t,i} \nabla \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}), -[\ell(\mathbf{w}_t; \mathbf{z}_t^{(1)}), \dots, \ell(\mathbf{w}_t; \mathbf{z}_t^{(m)})]^\top \right] \end{aligned} \quad (5.20)$$

即式 (5.17) 中的随机梯度, 且其范数的上界为:

$$\begin{aligned} \|\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t)\|_{w,*} &= \left\| \sum_{i=1}^m q_{t,i} \nabla \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}) \right\|_{w,*} \leq \sum_{i=1}^m q_{t,i} \left\| \nabla \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}) \right\|_{w,*} \stackrel{(3.3)}{\leq} \sum_{i=1}^m q_{t,i} G = G, \\ \|\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)\|_\infty &= \left\| [\ell(\mathbf{w}_t; \mathbf{z}_t^{(1)}), \dots, \ell(\mathbf{w}_t; \mathbf{z}_t^{(m)})]^\top \right\|_\infty \stackrel{(3.2)}{\leq} 1 \end{aligned}$$

因而

$$\|\mathbf{g}(\mathbf{w}_t, \mathbf{q}_t)\|_* = \sqrt{2D^2 \|\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t)\|_{w,*}^2 + 2 \|\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)\|_\infty^2 \ln m} \leq \underbrace{\sqrt{2D^2 G^2 + 2 \ln m}}_{:=M}. \quad (5.21)$$

由 $\phi(\cdot, \cdot)$ 的凸凹性, 我们有 [Nemirovski et al., 2009, (3.9)]

$$\begin{aligned}
& \max_{\mathbf{q} \in \Delta_m} \phi(\bar{\mathbf{w}}, \mathbf{q}) - \min_{\mathbf{w} \in \mathcal{W}} \phi(\mathbf{w}, \bar{\mathbf{q}}) \\
&= \max_{\mathbf{q} \in \Delta_m} \phi\left(\sum_{t=1}^T \frac{\mathbf{w}_t}{T}, \mathbf{q}\right) - \min_{\mathbf{w} \in \mathcal{W}} \phi\left(\mathbf{w}, \sum_{t=1}^T \frac{\mathbf{q}_t}{T}\right) \\
&\leq \frac{1}{T} \left[\max_{\mathbf{q} \in \Delta_m} \sum_{t=1}^T \phi(\mathbf{w}_t, \mathbf{q}) - \min_{\mathbf{w} \in \mathcal{W}} \sum_{t=1}^T \phi(\mathbf{w}, \mathbf{q}_t) \right] \\
&\leq \frac{1}{T} \max_{(\mathbf{w}, \mathbf{q}) \in \mathcal{W} \times \Delta_m} \sum_{t=1}^T \left[\langle \nabla_{\mathbf{w}} \phi(\mathbf{w}_t, \mathbf{q}_t), \mathbf{w}_t - \mathbf{w} \rangle - \langle \nabla_{\mathbf{q}} \phi(\mathbf{w}_t, \mathbf{q}_t), \mathbf{q}_t - \mathbf{q} \rangle \right] \\
&\stackrel{(5.19)}{=} \frac{1}{T} \max_{\mathbf{x} \in \mathcal{W} \times \Delta_m} \sum_{t=1}^T \langle F(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{x} \rangle.
\end{aligned} \tag{5.22}$$

因此，我们可以将优化误差分解如下：

$$\begin{aligned}
& \max_{\mathbf{q} \in \Delta_m} \phi(\bar{\mathbf{w}}, \mathbf{q}) - \min_{\mathbf{w} \in \mathcal{W}} \phi(\mathbf{w}, \bar{\mathbf{q}}) \\
&\leq \underbrace{\frac{1}{T} \max_{\mathbf{x} \in \mathcal{W} \times \Delta_m} \sum_{t=1}^T \langle \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{x} \rangle}_{:= E_1} \\
&\quad + \underbrace{\frac{1}{T} \max_{\mathbf{x} \in \mathcal{W} \times \Delta_m} \sum_{t=1}^T \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{x} \rangle}_{:= E_2}.
\end{aligned} \tag{5.23}$$

我们继续约束 (5.23) 中的两项。为了约束第一项 E_1 ，我们仿照 (5.3) 的推导，有

$$\begin{aligned}
\eta \langle \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{x} \rangle &\leq B(\mathbf{x}, \mathbf{x}_t) - B(\mathbf{x}, \mathbf{x}_{t+1}) + \frac{\eta^2}{2} \|\mathbf{g}(\mathbf{w}_t, \mathbf{q}_t)\|_*^2 \\
&\stackrel{(5.21)}{\leq} B(\mathbf{x}, \mathbf{x}_t) - B(\mathbf{x}, \mathbf{x}_{t+1}) + \frac{M^2}{2} \eta^2.
\end{aligned} \tag{5.24}$$

对于 $t = 1, \dots, T$ ，对上述不等式求和，我们有

$$\sum_{t=1}^T \eta \langle \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{x} \rangle \leq B(\mathbf{x}, \mathbf{x}_1) + \frac{M^2}{2} \sum_{t=1}^T \eta^2 \stackrel{(5.16)}{\leq} 1 + \frac{M^2}{2} \sum_{t=1}^T \eta^2. \tag{5.25}$$

接下来，我们考虑第二项 E_2 。由于最大化运算，变量 \mathbf{x} 受到算法随机性的影响，因此我们不能将 $\eta \langle E_{t-1}[\mathbf{g}(\mathbf{w}_t, \mathbf{q}_t)] - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{x} \rangle$, $t = 1, \dots, T$ 视为鞅差序列。为了解决这一挑战，我们使用 Nemirovski et al. [2009] 中的“幽灵迭代”技术，得出以下引理：

引理 2. 在定理 2 的条件下，我们有

$$\mathbb{E} \left[\max_{\mathbf{x} \in \mathcal{W} \times \Delta_m} \sum_{t=1}^T \eta \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{x} \rangle \right] \leq 1 + 2M^2 \sum_{t=1}^T \eta^2. \tag{5.26}$$

此外，以至少 $1 - \delta$ 的概率，我们有

$$\begin{aligned}
& \max_{\mathbf{x} \in \mathcal{W} \times \Delta_m} \sum_{t=1}^T \eta \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{x} \rangle \\
&\leq 1 + 2M^2 \sum_{t=1}^T \eta^2 + 4M \sqrt{2 \sum_{t=1}^T \eta^2 \left(\ln \frac{1}{\delta} \right)}.
\end{aligned} \tag{5.27}$$

证明：我们通过以 $F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t)$ 作为梯度执行 SMD 来创建虚拟序列：

$$\mathbf{y}_{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{W} \times \Delta_m} \left\{ \eta \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x} - \mathbf{y}_t \rangle + B(\mathbf{x}, \mathbf{y}_t) \right\} \quad (5.28)$$

其中 $\mathbf{y}_1 = \mathbf{x}_1$ 。然后，我们进一步将误差项分解为

$$\begin{aligned} & \max_{\mathbf{x} \in \mathcal{W} \times \Delta_m} \sum_{t=1}^T \eta \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{x} \rangle \\ & \leq \underbrace{\max_{\mathbf{x} \in \mathcal{W} \times \Delta_m} \sum_{t=1}^T \eta \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{y}_t - \mathbf{x} \rangle}_{:=A} \\ & \quad + \underbrace{\sum_{t=1}^T \eta \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{y}_t \rangle}_{:=B}. \end{aligned} \quad (5.29)$$

为了约束项 A ，我们重复 (5.24) 的分析，有

$$\begin{aligned} & \eta \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{y}_t - \mathbf{x} \rangle \\ & \leq B(\mathbf{x}, \mathbf{y}_t) - B(\mathbf{x}, \mathbf{y}_{t+1}) + \frac{\eta^2}{2} \|F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t)\|_*^2 \\ & \leq B(\mathbf{x}, \mathbf{y}_t) - B(\mathbf{x}, \mathbf{y}_{t+1}) + 2M^2\eta^2 \end{aligned} \quad (5.30)$$

其中最后一步利用以下不等式

$$\|F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t)\|_* \leq \|F(\mathbf{w}_t, \mathbf{q}_t)\|_* + \|\mathbf{g}(\mathbf{w}_t, \mathbf{q}_t)\|_* \stackrel{(5.21)}{\leq} 2M. \quad (5.31)$$

对 $t = 1, \dots, T$ 累加 (5.30)，我们有

$$\begin{aligned} A &= \sum_{t=1}^T \eta \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{y}_t - \mathbf{x} \rangle \\ &\leq B(\mathbf{x}, \mathbf{y}_1) + 2M^2 \sum_{t=1}^T \eta^2 \stackrel{(5.16)}{\leq} 1 + 2M^2 \sum_{t=1}^T \eta^2. \end{aligned} \quad (5.32)$$

为了约束 (5.29) 中的项 B ，我们定义

$$\delta_t = \eta \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{y}_t \rangle.$$

由于 \mathbf{x}_t 和 \mathbf{y}_t 独立于用于在 (5.20) 中构造梯度 $\mathbf{g}(\mathbf{w}_t, \mathbf{q}_t)$ 的随机样本 $\mathbf{z}_t^{(1)}, \dots, \mathbf{z}_t^{(m)}$ ， $\delta_1, \dots, \delta_T$ 构成鞅差序列。因此，我们有

$$\mathbb{E}[B] = \mathbb{E} \left[\sum_{t=1}^T \delta_t \right] = 0. \quad (5.33)$$

对 (5.29) 两边求期望，我们有

$$\begin{aligned} & \mathbb{E} \left[\max_{\mathbf{x} \in \mathcal{W} \times \Delta_m} \sum_{t=1}^T \eta \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{x} \rangle \right] \\ & \leq \mathbb{E} \left[\max_{\mathbf{x} \in \mathcal{W} \times \Delta_m} A \right] + \mathbb{E}[B] \stackrel{(5.32), (5.33)}{\leq} 1 + 2M^2 \sum_{t=1}^T \eta^2 \end{aligned}$$

为了建立高概率上界, 我们遵循 5.1.2 节中的分析, 并利用引理 1 来确定 B 的上界。为此, 我们首先证明 $|\delta_t|$ 有界:

$$\begin{aligned}
 |\delta_t| &= \left| \eta \langle F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t), \mathbf{x}_t - \mathbf{y}_t \rangle \right| \\
 &\leq \eta \|F(\mathbf{w}_t, \mathbf{q}_t) - \mathbf{g}(\mathbf{w}_t, \mathbf{q}_t)\|_* \|\mathbf{x}_t - \mathbf{y}_t\| \\
 &\stackrel{(5.31)}{\leq} 2M\eta (\|\mathbf{x}_t - \mathbf{x}_1\| + \|\mathbf{x}_1 - \mathbf{y}_t\|) \\
 &\leq 2M\eta \left(\sqrt{B(\mathbf{x}_t, \mathbf{x}_1)} + \sqrt{B(\mathbf{y}_t, \mathbf{x}_1)} \right) \stackrel{(5.16)}{\leq} 4M\eta.
 \end{aligned}$$

根据引理 1 和联合界不等式, 以至少 $1 - \delta$ 的概率有

$$B = \sum_{t=1}^T \delta_t \leq 4M \sqrt{2 \sum_{t=1}^T \eta^2 \left(\ln \frac{1}{\delta} \right)}. \quad (5.34)$$

我们把 (5.32) 和 (5.34) 代入 (5.29) 即可得到 (5.27)。

至此, 引理2证明完毕。

结合 (5.23), (5.25) 和 (5.26), 我们有

$$\mathbb{E} \left[\max_{\mathbf{q} \in \Delta_m} \phi(\bar{\mathbf{w}}, \mathbf{q}) - \min_{\mathbf{w} \in \mathcal{W}} \phi(\mathbf{w}, \bar{\mathbf{q}}_t) \right] \leq \frac{2}{\eta T} + \frac{5M^2}{2} \eta.$$

通过设定

$$\eta = \frac{2}{M\sqrt{5T}} = \sqrt{\frac{2}{5T(D^2G^2 + \ln m)}}, \quad (5.35)$$

我们有

$$\mathbb{E} \left[\max_{\mathbf{q} \in \Delta_m} \phi(\bar{\mathbf{w}}, \mathbf{q}) - \min_{\mathbf{w} \in \mathcal{W}} \phi(\mathbf{w}, \bar{\mathbf{q}}_t) \right] \leq 2M\sqrt{\frac{5}{T}} = 2\sqrt{\frac{10(D^2G^2 + \ln m)}{T}}.$$

结合 (5.23), (5.25) 和 (5.27), 以至少 $1 - \delta$ 的概率, 我们有

$$\max_{\mathbf{q} \in \Delta_m} \phi(\bar{\mathbf{w}}, \mathbf{q}) - \min_{\mathbf{w} \in \mathcal{W}} \phi(\mathbf{w}, \bar{\mathbf{q}}_t) \leq 2M\sqrt{\frac{5}{T}} + 4M\sqrt{\frac{\ln 1/\delta}{T}} = 2 \left(\sqrt{10} + 2\sqrt{2 \ln 1/\delta} \right) \sqrt{\frac{D^2G^2 + \ln m}{T}}.$$

至此, 定理2证明完毕。

接下来证明定理 1, 易知算法 1 和算法 2 的不同之处仅在于随机梯度的构建方式, 因此根据算法 1 的假设, 我们有

$$\|\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t)\|_{w,*} = \left\| q_{t,i_t} m \nabla \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}) \right\|_{w,*} \leq q_{t,i_t} m \left\| \nabla \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}) \right\|_{w,*} \stackrel{(3.3)}{\leq} q_{t,i_t} m G \leq mG,$$

$$\|\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)\|_\infty = \left\| [0, \dots, m\ell(\mathbf{w}_t; \mathbf{z}_t^{(i_t)}), \dots, 0]^\top \right\|_\infty \stackrel{(3.2)}{\leq} m.$$

$$\begin{aligned}
 \left\| [\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t); -\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)] \right\|_* &= \sqrt{2D^2 \|\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t)\|_{w,*}^2 + 2\|\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)\|_\infty^2 \ln m} \\
 &\leq \underbrace{\sqrt{2m^2 D^2 G^2 + 2m^2 \ln m}}_{:=M}.
 \end{aligned}$$

通过设定

$$\eta = \frac{2}{M\sqrt{5T}} = \sqrt{\frac{2}{5Tm^2(D^2G^2 + \ln m)}},$$

并仿照定理 2 的证明过程即可完成证明。

至此, 定理1证明完毕。

5.3 定理 3 的证明

对于 (4.2) 中的随机梯度，其范数上界可以通过和 (3.11) 类似方式得出，即

$$\begin{aligned}\|\mathbf{g}_w(\mathbf{w}_t, \mathbf{q}_t)\|_{w,*} &= \left\| \sum_{i \in C_t} q_{t,i} \nabla \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}) \right\|_{w,*} \leq \sum_{i \in C_t} q_{t,i} \left\| \nabla \ell(\mathbf{w}_t; \mathbf{z}_t^{(i)}) \right\|_{w,*} \stackrel{(3.3)}{\leq} \sum_{i \in C_t} q_{t,i} G = G, \\ \|\mathbf{g}_q(\mathbf{w}_t, \mathbf{q}_t)\|_\infty &= \max_{i \in C_t} |\ell(\mathbf{w}_t; \mathbf{z}_t^{(i)})| \stackrel{(3.2)}{\leq} 1.\end{aligned}$$

因此，通过与定理 2 完全相同的分析，我们有

$$\mathbb{E}[\epsilon_\varphi(\bar{\mathbf{w}}, \bar{\mathbf{q}})] \leq 2 \sqrt{\frac{10(D^2 G^2 + \ln m)}{n_1}}$$

以至少为 $1 - \delta$ 的概率，有

$$\epsilon_\varphi(\bar{\mathbf{w}}, \bar{\mathbf{q}}) \leq 2 \left(\sqrt{10} + 2\sqrt{2 \ln 1/\delta} \right) \sqrt{\frac{D^2 G^2 + \ln m}{n_1}}. \quad (5.36)$$

接下来，我们将讨论如何在每个分布 \mathcal{P}_i 上限定 $\bar{\mathbf{w}}$ 的风险，即 $R_i(\bar{\mathbf{w}})$ 。在 (3.5) 中的推导后，我们知道

$$\max_{i \in [m]} p_i R_i(\bar{\mathbf{w}}) - \min_{\mathbf{w} \in \mathcal{W}} \max_{i \in [m]} p_i R_i(\mathbf{w}) \leq \epsilon_\varphi(\bar{\mathbf{w}}, \bar{\mathbf{q}}).$$

因此，对于所有分布 \mathcal{P}_i ， $R_i(\bar{\mathbf{w}})$ 的上界为：

$$R_i(\bar{\mathbf{w}}) \leq \frac{1}{p_i} \min_{\mathbf{w} \in \mathcal{W}} \max_{i \in [m]} p_i R_i(\mathbf{w}) + \frac{1}{p_i} \epsilon_\varphi(\bar{\mathbf{w}}, \bar{\mathbf{q}}).$$

以 (5.36) 中的高概率界为例，我们以至少 $1 - \delta$ 的概率有

$$\begin{aligned}R_i(\bar{\mathbf{w}}) &\leq \frac{1}{p_i} \min_{\mathbf{w} \in \mathcal{W}} \max_{i \in [m]} p_i R_i(\mathbf{w}) + \frac{1}{p_i} 2 \left(\sqrt{10} + 2\sqrt{2 \ln 1/\delta} \right) \sqrt{\frac{D^2 G^2 + \ln m}{n_1}} \\ &= \frac{n_1}{n_i} \min_{\mathbf{w} \in \mathcal{W}} \max_{i \in [m]} p_i R_i(\mathbf{w}) + 2 \left(\sqrt{10} + 2\sqrt{2 \ln 1/\delta} \right) \frac{\sqrt{(D^2 G^2 + \ln m) n_1}}{n_i}.\end{aligned} \quad (5.37)$$

至此，定理3证明完毕。

6 总结与展望

对于 GDRO 问题，本文提出了两种基于 SMD 的 SA 算法。这两种算法每轮使用的样本数不同（分别为 1 和 m ），随机梯度的构建方式不同，其中第二个算法取得了接近最优的 $O(m(\log m)/\epsilon^2)$ 样本复杂度。然后，本文形式化了加权 GDRO 问题来处理各分布样本预算不均衡的情况，将非均匀采样纳入到 SMD 方法，以满足样本预算期望，并提供了分布相关的收敛速度。

关于未来展望，虽然本文主要关注 GDRO，但可能的未来工作是将 SA 算法扩展到具有类似形式的问题，例如 PAC 协作学习 [Blum et al., 2017]、联邦学习 [Mohri et al., 2019] 和极小极大遗憾优化 [Agarwal and Zhang, 2022] 问题。

参考文献

- Aharon Ben-Tal, Dick den Hertog, Anja De Waegenare, Bertrand Melenberg, and Gijb Renen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59(2):341–357, 2013.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *International Conference on Learning Representations*, 2020.
- A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- Herbert Scarf. A min-max solution of an inventory problem. *Studies in the Mathematical Theory of Inventory and Production*, pages 201–209, 1958.
- Aman Sinha, Hongseok Namkoong, and John Duchi. Certifying some distributional robustness with principled adversarial training. In *International Conference on Learning Representations*, 2018.
- Tatsunori Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. Fairness without demographics in repeated loss minimization. In *Proceedings of the 35th International Conference on Machine Learning*, pages 1929–1938, 2018.
- Ziyu Xu, Chen Dan, Justin Khim, and Pradeep Ravikumar. Class-weighted classification: Trade-offs and robust approaches. In *Proceedings of the 37th International Conference on Machine Learning*, pages 10544–10554, 2020.
- Dvir Samuel and Gal Chechik. Distributional robustness loss for long-tail learning. In *Proceedings of the 2021 IEEE/CVF International Conference on Computer Vision*, pages 9475–9484, 2021.
- Jingzhao Zhang, Aditya Krishna Menon, Andreas Veit, Srinadh Bhojanapalli, Sanjiv Kumar, and Suvrit Sra. Coping with label shift via distributionally robust optimisation. In *International Conference on Learning Representations*, 2021.
- Nika Haghtalab, Michael I. Jordan, and Eric Zhao. On-demand sampling: Learning optimally from multiple distributions. In *Advances in Neural Information Processing Systems 35*, pages 406–419, 2022.
- Tasuku Soma, Khashayar Gatmiry, and Stefanie Jegelka. Optimal algorithms for group distributionally robust optimization and beyond. *ArXiv e-prints*, arXiv:2212.13669, 2022.
- Alekh Agarwal and Tong Zhang. Minimax regret optimization for robust machine learning under distribution shift. In *Proceedings of 35th Conference on Learning Theory*, pages 2704–2729, 2022.

- Agnieszka Słowik and Léon Bottou. On distributionally robust optimization and data rebalancing. In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, pages 1283–1297, 2022.
- Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.
- Predrag Radivojac, Nitesh V. Chawla, A. Keith Dunker, and Zoran Obradovic. Classification and knowledge discovery in protein databases. *Journal of Biomedical Informatics*, 37(4):224–239, 2004.
- Alekh Agarwal and Tong Zhang. Minimax regret optimization for robust machine learning under distribution shift. In *Proceedings of 35th Conference on Learning Theory*, pages 2704–2729, 2022.
- L. Zhang, P. Zhao, T. Yang, and Z.-H. Zhou. Stochastic approximation approaches to group distributionally robust optimization, 2023.
- Avrim Blum, Nika Haghtalab, Ariel D. Procaccia, and Mingda Qiao. Collaborative PAC learning. In *Advances in Neural Information Processing Systems 30*, pages 2389–2398, 2017.
- Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In *Proceedings of the 36th International Conference on Machine Learning*, pages 4615–4625, 2019.