



B站 丸一口



FedAWARE^[1,2]

Tackling Hybrid Heterogeneity on Federated
Optimization via Gradient Diversity Maximization

作者: Dun Zeng 曾趯 电子科技大学UESTC



[1]<https://arxiv.org/abs/2310.02702>



[2]<https://github.com/dunzeng/FedAWARE>



01/17

CONTENTS



Part 01

背景介绍

Introduction of Background

Part 02

算法介绍

Introduction of Algorithm

Part 03

写作赏析

Wrighting

Part 04

交流评价

Communications





PART 01

背景介绍

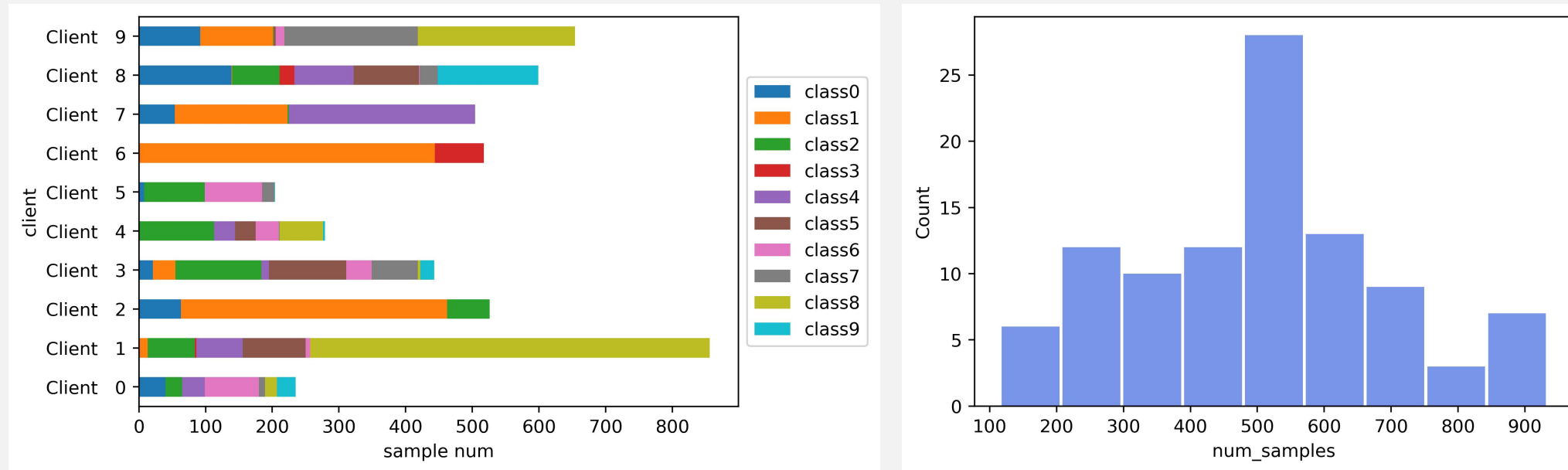
Introduction of Background



两种异构

Heterogeneities

Statistical Heterogeneity: Data Non-IID(数据非独立同分布)



[3]

System Heterogeneity:

- hardware specifications
- operating systems
- software configurations

This can result in communication overheads, computational disparities, and bias

Hybrid Heterogeneity

使用梯度差异性来衡量统计异质性状态的程度

Tackling Hybrid Heterogeneity on Federated Optimization via **Gradient Diversity** Maximization

Definition 3.1 (Gradient Diversity). We define the gradient diversity as the following ratio:

$$D(\mathbf{x}) := \sqrt{\frac{\sum_{i=1}^N \lambda_i \|\nabla f_i(\mathbf{x})\|^2}{\|\nabla f(\mathbf{x})\|^2}} \geq 1. \quad (3)$$

Assumption 3.1 Bounded global variance
Assumption 3.2 Bounded gradient dissimilarity

Corollary 3.1 (Bounded gradient diversity). Let Assumption 3.1 hold, it induces that

$$D(\mathbf{x}) \leq \sqrt{1 + \frac{\sigma_g^2}{\|\nabla f(\mathbf{x})\|^2}},$$

which is also connected to Assumption 3.2 with $G = 0$. In this case, the B is the upper bound of gradient diversity.

Bigger Gradient Diversity \leftrightarrow Smaller Norm of Global Gradient

Assumption 3.1 (Bounded global variance). We assume the averaged global variance is bounded, i.e.,
 $\sum_{i=1}^N \lambda_i \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2 \leq \sigma_g^2$ for all $\mathbf{x} \in \mathcal{X}$.

Assumption 3.2 (Bounded gradient dissimilarity). There exist constants $B \geq 1, G \geq 0$ such that
 $\sum_{i=1}^N \lambda_i \|\nabla f_i(\mathbf{x})\|^2 \leq B^2 \|\nabla f(\mathbf{x})\|^2 + G^2$ for all $\mathbf{x} \in \mathcal{X}$.



混合异质性效应的上限

Tackling Hybrid Heterogeneity on Federated Optimization via Gradient Diversity Maximization

Lemma 3.2 (Upper bound of balanced local drift). *Let Assumptions 2.1 2.2 3.1 hold. For any client $i \in [N]$ with balanced local iteration steps $k \in [K]$ with local learning rate $\eta_l \leq \frac{1}{K}$, the average of local drift can be bounded by:*

$$\sum_{i=1}^N \lambda_i \mathbb{E} \left\| \mathbf{x}_i^{t,k} - \mathbf{x}^t \right\|^2 \leq 5\eta_l(\sigma_l^2 + 6K\sigma_g^2 + 6K\mathbb{E} \left\| \nabla f(\mathbf{x}^t) \right\|^2). \quad (4)$$

Corollary 3.2 (Loose upper bound of unbalanced local drift). *We denote the local drift of client i from the global gradient as $\zeta_i(\mathbf{x}) = \left\| \nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x}) \right\|^2$. Let Assumption 2.1 2.2 3.1 hold. For all client $i \in [N]$ with arbitrary local iteration steps $k \in [K_i]$ with local learning rate $\eta_l \leq \frac{1}{K_i}$, the average of local drift can be bounded as follows:*

$$\sum_{i=1}^N \lambda_i \mathbb{E} \left\| \mathbf{x}_i^{t,k} - \mathbf{x}^t \right\|^2 \leq \Phi_{Hetero} + 5\eta_l(\sigma_l^2 + 6K_{\min}^2\sigma_g^2 + 6\tilde{K}\mathbb{E} \left\| \nabla f(\mathbf{x}^t) \right\|^2), \quad (5)$$

Smaller Norm of Global Gradient \leftrightarrow Lower drift

Assumption 2.1 (Smoothness). *Each objective $f_i(\mathbf{x})$ for all $i \in [N]$ is L -smooth, inducing that for all $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, it holds $\left\| \nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y}) \right\| \leq L\left\| \mathbf{x} - \mathbf{y} \right\|$.*

Assumption 2.2 (Unbiasedness and Bounded Local Variance). *For each $i \in [N]$ and $\mathbf{x} \in \mathbb{R}^d$, we assume the access to an unbiased stochastic gradient $\nabla F_i(\mathbf{x}, \xi_i)$ of client's true gradient $\nabla f_i(\mathbf{x})$, i.e., $\mathbb{E}_{\xi_i \sim \mathcal{D}_i} [\nabla F_i(\mathbf{x}, \xi_i)] = \nabla f_i(\mathbf{x})$. The function f_i have σ_l -bounded (local) variance i.e., $\mathbb{E}_{\xi_i \sim \mathcal{D}_i} \left[\left\| \nabla F_i(\mathbf{x}, \xi_i) - \nabla f_i(\mathbf{x}) \right\|^2 \right] \leq \sigma_l^2$.*



Gradient Diversity

Tackling Hybrid Heterogeneity on Federated Optimization via **Gradient Diversity** Maximization

Pre-Experiments

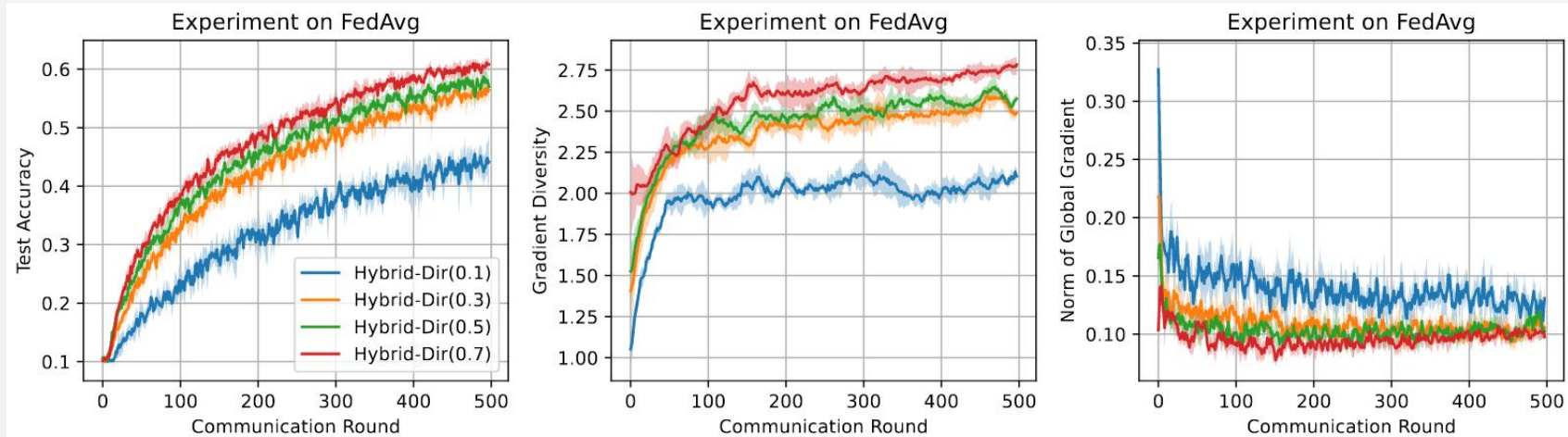


Figure 2. Observation: gradient diversity is related to convergence quality.

Bigger Acc \leftrightarrow Bigger Gradient Diversity \leftrightarrow Smaller Norm of Global Gradient



PART 02

算 法 介 绍

Introduction of Algorithm



都是梯度惹的祸

All heterogeneities can be reduce by decreasing Norm of Global Gradient

重构一组权重，通过自适应权重来减小 Norm of Global Gradient 以获取更好的ACC

formance. To this end, we use an adjustable $\tilde{\lambda}$ to study a *surrogate global objective*:

$$\tilde{f}(\mathbf{x}) = \sum_{i=1}^N \tilde{\lambda}_i f_i(\mathbf{x}).$$

$$\tilde{\lambda} = \min_{\lambda} \left\| \sum_{i=1}^N \lambda_i \nabla f_i(\mathbf{x}) \right\|^2, \quad (8)$$

As the dimension of gradients can be millions, we use the **Frank-Wolfe algorithm** (Jaggi, 2013) to solve it.

Solving Equation 8 typically **requires access to all local first-order gradients from clients**, which is often infeasible in FL systems. To overcome this limitation, FEDAWARE approximates utilizes the local updates **using the history momentum** of clients

Methodology

FedAWARE

Algorithm 1 FEDAWARE

Require: x^0, m^0, α

- 1: **for** round $t \in [T]$ **do**
- 2: Server sample clients S^t and broadcast model x^t
- 3: **for** client $i \in S^t$ in parallel **do**
- 4: $x_i^{t,0} = x^t$
- 5: **for** local update step $k \in [K_i]$ **do**
- 6: $x_i^{t,k} = x_i^{t,k-1} - \eta_l \nabla F_i(x_i^{t,k-1})$
- 7: **end for**
- 8: Client uploads local updates $g_i^t = x_i^{t,0} - x_i^{t,K_i}$
- 9: **end for**
- 10: Server updates local momentum
$$m_i^t = \begin{cases} \alpha m_i^{t-1} + (1 - \alpha) g_i^t, & \text{if } i \in S^t \\ m_i^{t-1}, & \text{if } i \notin S^t \end{cases}$$
- 11: Server computes $\tilde{\lambda}^t$ by (8) with m^t
- 12: Server computes global estimates $d^t = \sum_{i \in S^t} \tilde{\lambda}_i^t m_i^t$ and updates $x^{t+1} = x^t - \eta d^t$
- 13: **end for**

As the dimension of gradients can be millions, we use the **Frank-Wolfe algorithm** (Jaggi, 2013) to solve it.

Solving Equation 8 typically **requires access to all local first-order gradients from clients**, which is often infeasible in FL systems. To overcome this limitation, FEDAWARE approximates utilizes the local updates **using the history momentum** of clients

$$\tilde{\lambda} = \min_{\lambda} \left\| \sum_{i=1}^N \lambda_i \nabla f_i(x) \right\|^2, \quad (8)$$





PART 03

写作赏析

Awesome Wrihting



Wrighting

Awesome

3. Analyses on Hybrid Heterogeneity

3.1. 测量统计异质性影响

3.2. 混合异质性效应的上限

3.3. 减轻混合异质性影响

Corollary 3.1 (Bounded gradient diversity). *Let Assumption 3.1 hold, it induces that*

$$D(\mathbf{x}) \leq \sqrt{1 + \frac{\sigma_g^2}{\|\nabla f(\mathbf{x})\|^2}},$$

which is also connected to Assumption 3.2 with $G = 0$. In this case, the B is the upper bound of gradient diversity.

Corollary 3.2 (Loose upper bound of unbalanced local drift). *We denote the local drift of client i from the global gradient as $\zeta_i(\mathbf{x}) = \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2$. Let Assumption 2.1 2.2 3.1 hold. For all client $i \in [N]$ with arbitrary local iteration steps $k \in [K_i]$ with local learning rate $\eta_l \leq \frac{1}{K_i}$, the average of local drift can be bounded as follows:*

$$\sum_{i=1}^N \lambda_i \mathbb{E} \left\| \mathbf{x}_i^{t,k} - \mathbf{x}^t \right\|^2 \leq \Phi_{Hetero} + 5\eta_l(\sigma_l^2 + 6K_{min}^2\sigma_g^2 + 6\tilde{K} \mathbb{E} \|\nabla f(\mathbf{x}^t)\|^2), \quad (5)$$

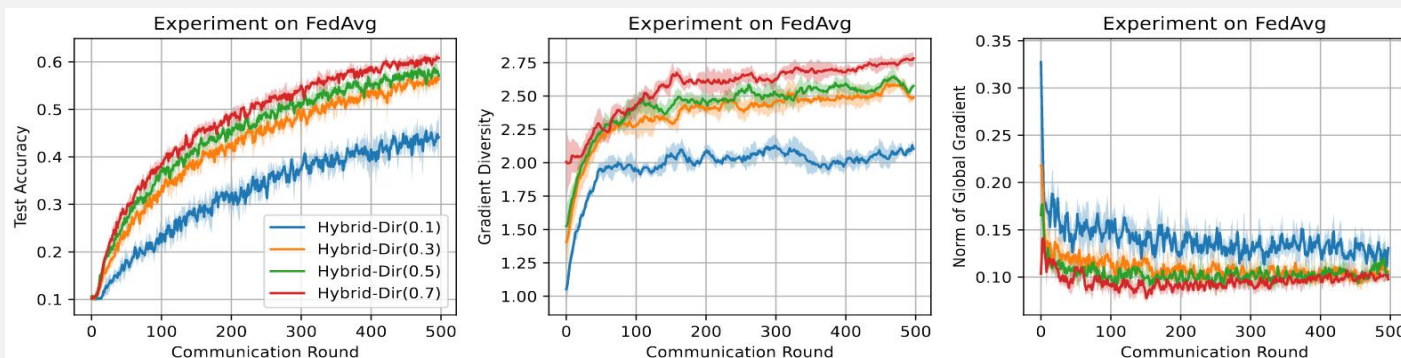


Figure 2. Observation: gradient diversity is related to convergence quality.

Bigger Acc \leftrightarrow Bigger Gradient Diversity \leftrightarrow Smaller Norm of Global Gradient

两节理论 + 一节实验得到优化对象

Wrighting

Awesome



推论 3.2 (非平衡本地漂移的松散上界)。我们标记客户端 i 与全局梯度的本地漂移为 $\zeta_i(\mathbf{x}) = \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2$ 。假设 2.1、2.2、3.1 成立。对于所有客户端 $i \in [N]$ ，其任意本地迭代步数 $k \in [K_i]$ ，且本地学习率 $\eta_l \leq \frac{1}{K_i}$ ，本地漂移的平均值可以被限制如下：

$$\sum_{i=1}^N \lambda_i \mathbb{E} \|\mathbf{x}_i^{t,k} - \mathbf{x}^t\|^2 \leq \Phi_{\text{Hetero}} + 5\eta_l \left(\sigma_l^2 + 6K_{\min}^2 \sigma_g^2 + 6\tilde{K} \mathbb{E} \|\nabla f(\mathbf{x}^t)\|^2 \right), \quad (5)$$

其中 $\tilde{K} = \sum_{i=1}^N \lambda_i K_i$, $K_{\min} = \min(K_1, \dots, K_N)$ ，以及 $\Phi_{\text{Hetero}} = \sum_{i=1}^N 30\eta_l (K_i - K_{\min}) \zeta_i(\mathbf{x}^t)$ 。

评注 3.1 (推论 3.2 的解释)。本文没有假设 Φ_{Hetero} 中本地差异性 $\zeta_i(\mathbf{x}^t)$ 的界限。因此，当本地更新步数变得不平衡时，上界 (4) 被 (5) 替代，并引入了混合异质性项 Φ_{Hetero} 。这是因为至少会有一个客户端 i 使得 $K_i - K_{\min} = 0$ ，使得假设 3.1 不适用。此外，根据系统异质性（非平衡本地步骤），由于额外的本地步骤， Φ_{Hetero} 项被放大，使得 (5) 成为一个非常松散的界限。因此，这可能会对联邦优化的性能产生负面影响。

对混合异质性的洞察。混合异质性影响是由局部差异和不平衡局部步骤协同引起的。联邦优化的先前工作隐含地最小化上限 (5) 以提高优化性能。例如，`fedprox` (li et al., 2020c) 利用惩罚项来减少推论 3.2 中的局部漂移 $\mathbb{E} \|\mathbf{x}_i^{t,K_i} - \mathbf{x}^t\|^2$ 。类似地，`SCAFFLOD` (karimireddy et al., 2020)、`fedavgm` (hsu et al., 2019) 和 `feddyn` (acar et al., 2020) 使用方差正则化项校正本地更新以缩小方差 σ_l 和 σ_g 。`fednova` (wang et al., 2020) 基于局部步骤对本地更新进行剪辑，以减少 $K_i, \forall i$ 的尺度效应。总之，以前的工作主要是通过操纵本地更新来减轻异质性。

这段写得是真好啊，把其他人的工作全部分析清楚了

理论分析出前人工作的优化方式



Wrighting

Awesome

重构一组权重，通过自适应权重来减小 Norm of Global Gradient 以获取更好的ACC

formance. To this end, we use an adjustable $\tilde{\lambda}$ to study a *surrogate global objective*:

$$\tilde{f}(\mathbf{x}) = \sum_{i=1}^N \tilde{\lambda}_i f_i(\mathbf{x}).$$

提出异于前人工作的优化方式：自适应权重

$$\tilde{\lambda} = \min_{\lambda} \left\| \sum_{i=1}^N \lambda_i \nabla f_i(\mathbf{x}) \right\|^2, \quad (8)$$

As the dimension of gradients can be millions, we use the **Frank-Wolfe algorithm** (Jaggi, 2013) to solve it.

Solving Equation 8 typically **requires access to all local first-order gradients from clients**, which is often infeasible in FL systems. To overcome this limitation, FEDAWARE approximates utilizes the local updates **using the history momentum** of clients





PART 04

交 流 评 价

Outline, Contribution, Argue Comments



交流评价

Communications

写作是真牛逼👍 阅读体验很丝滑



学长，这两天在做PPT准备汇报你的FedAWARE，发现我之前的“用实验（Figure 2）去说明 Gradient Diversity大的Acc好，这边感觉逼格就没前面高了”这个其实不对，你这边的写法是没问题的。3.3节用实验去和3.1+3.2照应了，是先理论、再实验验证的路子，这没问题的。我之前没看清，是我的问题。后面Frank-Wolfe algorithm和momentum的引入，还是觉得有点唐突哈，这两点可能还是需要修一下。



🤔 你说的问题确实，我自己写的时候也不太确定该怎么说这个衔接



这个文章在nips和icml都拿了个borderline 看缘分吧 看来路还长



嗯嗯👍 多谢帮我宣传 如果有啥反馈也可以和我说



- ① The connection between theory and methodology needs to be improved.
- ② 写得很好，中稿只是时间问题。



B站 丸一口

感谢您的三连指导

2024/04/22

