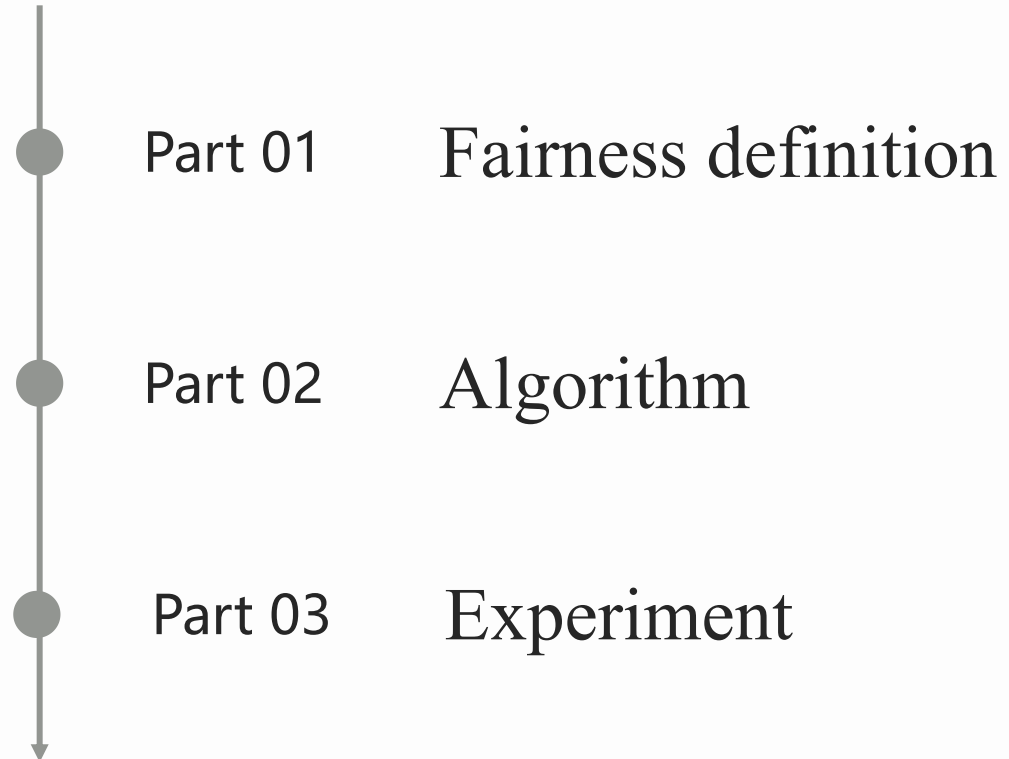# Fair Resource Allocation in Federated Learning

Tian Li et.     CMU     ICLR2020

Power by 丸一口

CONTENT

# Part 01

# Fairness definition
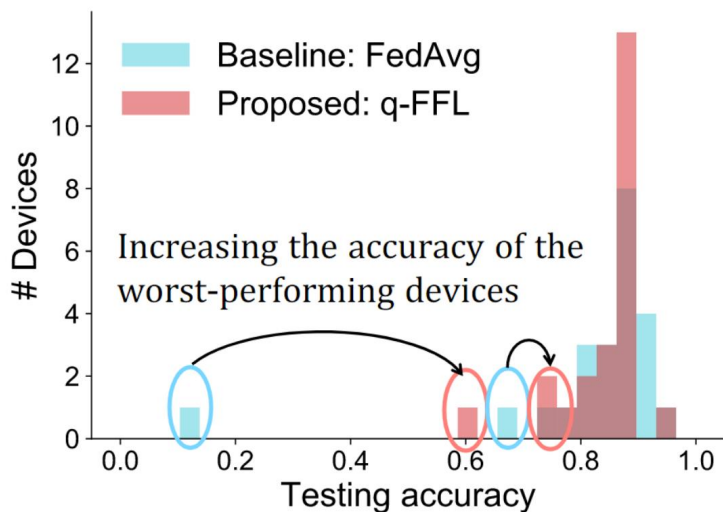
Two reasonable definition of fairness.

# Fairness

More Average: Emphasize that everyone **has equal opportunities** and care about **poor** clients;

均衡公平性[1]：强调 **人人平等有机会** ，关心 **表现差** 的客户；

Contribution: emphasizing distribution according to work, where those who **work more receive more.**

贡献公平性[1]：强调 **按劳分配，多劳多得，优胜劣汰** 。



**Definition 1** (*Fairness of performance distribution*). For trained models $w$ and $\tilde{w}$, we informally say that model $w$ provides a more *fair* solution to the federated learning objective (1) than model $\tilde{w}$ if the performance of model $w$ on the $m$ devices, $\{a_1, \ldots a_m\}$, is more *uniform* than the performance of model $\tilde{w}$ on the $m$ devices.

Uniformity：In this work, we mainly use the *variance* of the *performance* distribution as a measure of uniformity

Performance: In this work, we take 'performance', $a_k$, to be the *testing accuracy* of applying the trained model $w$ on the test data for device k.

[1] https://zhuanlan.zhihu.com/p/601227861

# Part 02

---

## Algorithm

How to improve the fairness of Federated Learning?

# Objection function

FedAvg

$$\min_{w} f(w) = \sum_{k=1}^{m} p_k F_k(w)$$

$m$ is the total number of devices, $p_k \geq 0$, and $\sum_k p_k = 1$

q-Fair Federated Learning (q-FFL)

$$\min_{w} f_q(w) = \sum_{k=1}^{m} \frac{p_k}{q+1} F_k^{q+1}(w)$$

| | | |
|---|---|---|
| q = 0 | → | FedAvg |
| q increase... | → | Focus on the big one of $F_k$ (performace bad) |
| q = ∞ | → | Agnostic Federated Learning (AFL)[2] |

[2] Mohri, Mehryar, Gary Sivek, and Ananda Theertha Suresh. "Agnostic federated learning." International Conference on Machine Learning. PMLR, 2019. 不可知论的联邦学习

# Extra parameter: q

$$\min_{w} \; f_q(w) = \sum_{k=1}^{m} \frac{p_k}{q+1} F_k^{q+1}(w)$$

Effectiveness of q: larger q inducing more fairness

How to get a proper q: grid search

Reduce the parameter search: auto choice learning rate

# Auto choice learning rate

One concern with solving such a family of objectives is that it requires step-size tuning for every value of $q$. In particular, in gradient-based methods, the step-size inversely depends on the Lipschitz constant of the function's gradient, which will change as we change $q$. This can quickly cause the

$$q \rightarrow L(q) \rightarrow \eta \propto 1/L(q)$$

实际上是二阶导的上界作为L

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla F(\mathbf{w}_t)$$

$$\rightarrow \eta \nabla F(\mathbf{w}_t) = \mathbf{w}_t - \mathbf{w}_{t+1}$$

$$\rightarrow \frac{1}{\eta} = \frac{\nabla F(\mathbf{w}_t)}{\mathbf{w}_t - \mathbf{w}_{t+1}}$$

$$\rightarrow \frac{\nabla F(\mathbf{w}_{t+1}) - \nabla F(\mathbf{w}_t)}{\mathbf{w}_{t+1} - \mathbf{w}_t} = \frac{\nabla F(\mathbf{w}_t) - \nabla F(\mathbf{w}_{t+1})}{\mathbf{w}_t - \mathbf{w}_{t+1}} \leq \frac{\nabla F(\mathbf{w}_t)}{\mathbf{w}_t - \mathbf{w}_{t+1}} = \frac{1}{\eta} = L$$

LV9 管理员 丸一口-M-DPFL

佬们，TianLi这篇ICLR2020《FAIR RESOURCE ALLOCATION IN FEDERATEDLEARNING》里说，"在基于梯度的方法中，步长成反比取决于函数梯度的Lipschitz常数"，这有啥依据吗，没看到她引文献呢

LV5 Zed-D-联邦优化-异构计算 (Author of FedLab)

分析里有学习率的地方一般会乘上 L 这时候取反比可以缓解L的影响

这个操作在优化分析里是基操 都是这么做的

# Compute L

$q \rightarrow L(q) \rightarrow \eta \propto 1/L(q)$

$$\min_{w} \ f_q(w) = \sum_{k=1}^{m} \frac{p_k}{q+1} F_k^{q+1}(w)$$

**Lemma 3.** *If the non-negative function $f(\cdot)$ has a Lipschitz gradient with constant L, then for any $q \geq 0$ and at any point $w$,*

$$L_q(w) = Lf(w)^q + qf(w)^{q-1}\|\nabla f(w)\|^2 \qquad (3)$$

*is an upper-bound for the local Lipschitz constant of the gradient of $\frac{1}{q+1}f^{q+1}(\cdot)$ at point $w$.*

*Proof.* At any point $w$, we can compute the Hessian $\nabla^2 \left(\frac{1}{q+1}f^{q+1}(w)\right)$ as: $\boxed{\nabla \left(\frac{1}{q+1}f^{q+1}(w)\right) = f^q(w) \cdot \nabla f(w)}$

$$\nabla^2 \left(\frac{1}{q+1}f^{q+1}(w)\right) = qf^{q-1}(w)\underbrace{\nabla f(w)\nabla^T f(w)}_{\preceq\|\nabla f(w)\|^2 \times I} + f^q(w)\underbrace{\nabla^2 f(w)}_{\preceq L \times I}. \qquad (4)$$

As a result, $\|\nabla^2 \frac{1}{q+1}f^{q+1}(w)\|_2 \leq L_q(w) = Lf(w)^q + qf(w)^{q-1}\|\nabla f(w)\|^2$. $\qquad \square$

[3] 非凸优化基石：Lipschitz Condition - Zeap的文章 - 知乎 https://zhuanlan.zhihu.com/p/27554191

# q-FedSGD

$$\min_{w} \ f_q(w) = \sum_{k=1}^{m} \frac{p_k}{q+1} F_k^{q+1}(w), \qquad (2)$$

$$L_q(w) = Lf(w)^q + qf(w)^{q-1}\|\nabla f(w)\|^2 \qquad (3)$$

---

**Algorithm 1** $q$-FedSGD

---

1: **Input:** $K, T, q, 1/L, w^0, p_k, k = 1, \cdots, m$
2: **for** $t = 0, \cdots, T-1$ **do**
3:      Server selects a subset $S_t$ of $K$ devices at random (each device $k$ is chosen with prob. $p_k$)
4:      Server sends $w^t$ to all selected devices
5:      Each selected device $k$ computes:
$$\Delta_k^t = F_k^q(w^t)\nabla F_k(w^t) \quad \boxed{\text{grad of (2)}}$$
$$h_k^t = qF_k^{q-1}(w^t)\|\nabla F_k(w^t)\|^2 + LF_k^q(w^t) \quad \boxed{\text{compute (3)}}$$
6:      Each selected device $k$ sends $\Delta_k^t$ and $h_k^t$ back to the server
7:      Server updates $w^{t+1}$ as:
$$w^{t+1} = w^t - \frac{\sum_{k \in S_t} \Delta_k^t}{\sum_{k \in S_t} h_k^t} \quad \boxed{\text{use 1/L as } \eta}$$
8: **end for**

---

# q-FedAvg

$$\min_w \ f_q(w) = \sum_{k=1}^m \frac{p_k}{q+1} F_k^{q+1}(w), \tag{2}$$

$$L_q(w) = Lf(w)^q + qf(w)^{q-1}\|\nabla f(w)\|^2 \tag{3}$$

E次梯度变化值，L是当做1/η在用

**Algorithm 2** $q$-FedAvg

1: **Input:** $K, E, T, q, 1/L, \eta, w^0, p_k, k = 1, \cdots, m$
2: **for** $t = 0, \cdots, T-1$ **do**
3:    Server selects a subset $S_t$ of $K$ devices at random (each device $k$ is chosen with prob. $p_k$)
4:    Server sends $w^t$ to all selected devices
5:    Each selected device $k$ updates $w^t$ for $E$ epochs of SGD on $F_k$ with step-size $\eta$ to obtain $\bar{w}_k^{t+1}$
6:    Each selected device $k$ computes:

$$\Delta w_k^t = L(w^t - \bar{w}_k^{t+1}) \quad \text{?}$$

$$\Delta_k^t = F_k^q(w^t)\Delta w_k^t \quad \boxed{\text{grad of (2)}}$$

$$h_k^t = qF_k^{q-1}(w^t)\|\Delta w_k^t\|^2 + LF_k^q(w^t) \quad \boxed{\text{compute (3)}}$$

$$\Delta_k^t = F_k^q(w^t)\nabla F_k(w^t) \quad \text{FedSGD}$$
$$h_k^t = qF_k^{q-1}(w^t)\|\nabla F_k(w^t)\|^2 + LF_k^q(w^t)$$

7:    Each selected device $k$ sends $\Delta_k^t$ and $h_k^t$ back to the server
8:    Server updates $w^{t+1}$ as:

$$w^{t+1} = w^t - \frac{\sum_{k\in S_t}\Delta_k^t}{\sum_{k\in S_t}h_k^t} \quad \boxed{\text{use 1/L as }\eta}$$
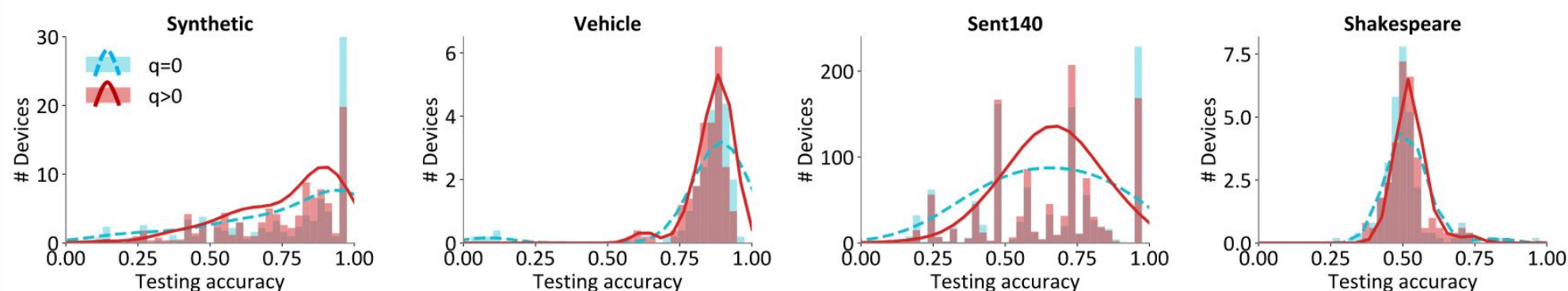
9: **end for**

[3] 非凸优化基石：Lipschitz Condition - Zeap的文章 - 知乎 https://zhuanlan.zhihu.com/p/27554191

# Part 03

# Experiment

How to construct a wonderful paper

# Extra parameter: q



| Dataset | Objective | Average (%) | Worst 10% (%) | Best 10% (%) | Variance |
|---|---|---|---|---|---|
| Synthetic | $q = 0$ | $80.8 \pm .9$ | $18.8 \pm 5.0$ | $100.0 \pm 0.0$ | $724 \pm 72$ |
|  | $q = 1$ | $79.0 \pm 1.2$ | $\mathbf{31.1} \pm 1.8$ | $100.0 \pm 0.0$ | $\mathbf{472} \pm 14$ |
| Vehicle | $q = 0$ | $87.3 \pm .5$ | $43.0 \pm 1.0$ | $\mathbf{95.7} \pm 1.0$ | $291 \pm 18$ |
|  | $q = 5$ | $87.7 \pm .7$ | $\mathbf{69.9} \pm .6$ | $94.0 \pm .9$ | $\mathbf{48} \pm 5$ |
| Sent140 | $q = 0$ | $65.1 \pm 4.8$ | $15.9 \pm 4.9$ | $100.0 \pm 0.0$ | $697 \pm 132$ |
|  | $q = 1$ | $66.5 \pm .2$ | $\mathbf{23.0} \pm 1.4$ | $100.0 \pm 0.0$ | $\mathbf{509} \pm 30$ |
| Shakespeare | $q = 0$ | $51.1 \pm .3$ | $39.7 \pm 2.8$ | $\mathbf{72.9} \pm 6.7$ | $82 \pm 41$ |
|  | $q = .001$ | $52.1 \pm .3$ | $\mathbf{42.1} \pm 2.1$ | $69.0 \pm 4.4$ | $\mathbf{54} \pm 27$ |

# Thanks~

Power by 丸一口