# Google

# Federated Composite Optimization

*(arXiv: 2011.08474)*

**Honglin Yuan**, Manzil Zaheer, Sashank Reddi

# Federated Learning

[Konečný et al., '15]

Notation: $w_{r,k}^m$
→ $m^{th}$ client
→ $r^{th}$ round, $k^{th}$ local iteration

# FEDAVG: the *de facto* standard of FL

[Konečný et al., '15]

Notation: $w_{r,k}^m$
- $m^{\text{th}}$ client
- $r^{\text{th}}$ round, $k^{\text{th}}$ local iteration



$w_{0,0}^1$ — Local data, $K$ steps of SGD → $w_{0,K}^1$

$w_{0,0}^2$ → $w_{0,K}^2$

$w_{0,0}^3$ → $w_{0,K}^3$

$w_{0,0}^M$ → $w_{0,K}^M$

Average & Broadcast

$$w_{1,0}^m \equiv \frac{1}{M} \sum_{m=1}^{M} w_{0,K}^m$$

$w_{1,0}^1$ → $w_{1,K}^1$ ...

$w_{1,0}^2$ → $w_{1,K}^2$ ...

$w_{1,0}^3$ → $w_{1,K}^3$ ...

$w_{1,0}^M$ → $w_{1,K}^M$ ...

Google

# FEDAVG: Generalized Formulation

[Karimireddy et al., ICML'20, Reddi et al., '20, etc]

Notation: $w_{r,k}^m$ — $m^{\text{th}}$ client

$r^{\text{th}}$ round, $k^{\text{th}}$ local iteration

**Algorithm 1** Federated Averaging (FEDAVG)

1: **procedure** FEDAVG$(w_0, \eta_c, \eta_s)$
2:     **for** $r = 0, \ldots, R-1$ **do**
3:        sample a subset of clients $\mathcal{S}_r \subseteq [M]$    **Client sampling**
4:        **on client** $m \in \mathcal{S}_r$ **in parallel do**
5:           client initialization $w_{r,0}^m \leftarrow w_r$    **Client update**
6:           **for** $k = 0, \ldots, K-1$ **do**
7:              $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$
8:              $w_{r,k+1}^m \leftarrow w_{r,k}^m - \eta_c \cdot g_{r,k}^m$
9:        $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r} (w_{r,K}^m - w_{r,0}^m)$    **Average client deltas (as pseudo anti-gradient)**
10:      $w_{r+1} \leftarrow w_r + \eta_s \cdot \Delta_r$    **Server update with server learning rate $\eta_s$**

Google

# Introducing Federated Composite Optimization (FCO)

- FedAvg (and other existing FL algorithms) solves **unconstrained** (smooth) problem only

  - $\min_{w \in \mathbb{R}^d} \frac{1}{M} \sum_{m=1}^{M} F_m(w)$ , where $F_m(w) := \mathbb{E}_{\xi \sim \mathcal{D}_m}[f(w; \xi)]$     *[e.g., Woodworth et al., NeurIPS'20]*

    Data distribution of the $m^{\text{th}}$ client

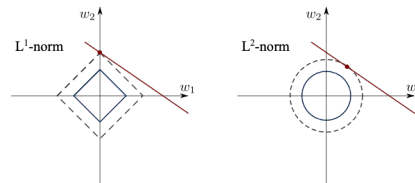- We propose Federated **composite** optimization (FCO)

  - $\min_{w \in \mathbb{R}^d} \Phi(w) := \frac{1}{M} \sum_{m=1}^{M} [F_m(w) + \psi_m(w)]$ , where $\psi_m$ is convex composite functions

Google

# Example of $\psi_m$: FL with Regularization

$$\min_{w \in \mathbb{R}^d} \frac{1}{M} \sum_{m=1}^{M} [F_m(w) + \psi_m(w)]$$

- Let $\psi_m(w)$ be regularizers

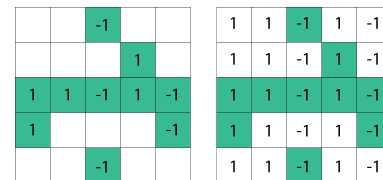- Federated Lasso for sparsity representations

$$\min_w \frac{1}{M} \sum_{m=1}^{M} \mathbb{E}_{(x,y) \sim \mathcal{D}_m} \|x^T w - y\|_2^2 + \lambda \|w\|_1$$

  Potential application: cross-silo distributed biomedical data

- Federated matrix completion for recommendation system

$$\min_W \frac{1}{M} \sum_{m=1}^{M} F_m(W) + \lambda \|W\|_*$$ → Matrix nuclear norm promotes low-rank

Google

# Example of $\psi_m$ : FL with (Personalized) Constraints

- Let $\psi_m(w)$ be convex indicator $\begin{cases} 0 & \text{if } w \in C_m \\ +\infty & \text{if } w \notin C_m \end{cases}$

$$\min_{w \in \mathbb{R}^d} \frac{1}{M} \sum_{m=1}^{M} [F_m(w) + \psi_m(w)]$$

- Problem becomes

$$\min_w \quad \frac{1}{M} \sum_{m=1}^{M} F_m(w)$$
$$\text{s.t.} \quad w \in \bigcap_{m=1}^{M} C_m \quad \longrightarrow \quad \text{Fulfill all constraints}$$

- **Budgeting**, each customer has a budget constraint

- FL with monotonic constraints $\longrightarrow$ Improve interpretability

- Inputs welcome!

Google

# Mix & Match of Setups

$$\min_{w\in\mathbb{R}^d} \frac{1}{M}\sum_{m=1}^{M}[F_m(w)+\psi_m(w)]$$

- **Homogeneous** vs **heterogeneous objective** $F_m$: standard "heterogeneity" in FL

  *[e.g., Li et al., MLSys'20, Karimireddy et al., ICML'20, Woodworth et al., NeurIPS'20]*

- **Homogeneous** vs **heterogeneous composite** $\psi_m$

- **Client** and/or **server** access to composite oracle $\psi_m$
  - Client-side oracle: better convergence? Privacy for personalized constraints?
  - Server-side oracle: computationally light

- **In this work, we focus on homogeneous $\psi_m \equiv \psi$ but allowing for heterogeneous $F_m$**

$$\min_{w\in\mathbb{R}^d} \Phi(w) := \frac{1}{M}\sum_{m=1}^{M} F_m(w) + \psi(w)$$

Google

# Composite 101: Proximal Gradient Descent

- Consider sequential $\min\ F(w) + \psi(w)$, where $F$ smooth, $\psi$ "simple" and convex

- Proximal Gradient Descent (PGD)

$$w_{t+1} \leftarrow \mathbf{prox}_{\eta\psi}\left(w_t - \eta\nabla F(w_t)\right)$$

*Proximal additive*

$$:= \underset{w}{\operatorname{argmin}}\left\{ \boxed{F\left(w_t\right) + \langle\nabla F\left(w_t\right), w - w_t\rangle} + \boxed{\frac{1}{2\eta}\|w - w_t\|_2^2} + \boxed{\psi(w)} \right\}$$

*First-order Taylor expansion of F*　　*Smoothness estimation*

- **prox** operator can often be computed analytically

$$\psi(w) = \chi_{\mathcal{C}}(w) := \begin{cases} 0 & \text{if } w \in \mathcal{C} \\ +\infty & \text{if } w \notin \mathcal{C} \end{cases}$$ → *Projected GD*

$$\psi(w) = \frac{1}{2}\lambda\|w\|_2^2$$ → *Weight decay (variant)*
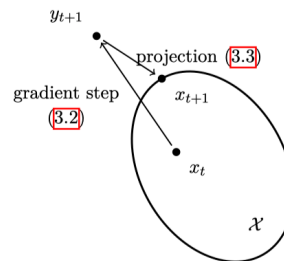
$$\psi(w) = \lambda\|w\|_1$$ → *Soft-thresholding*



*Image source: [Bubeck, 2015]*

# First Attempt: FEDAVG + Proximal Gradient Descent

**Algorithm 1** Federated Averaging (FEDAVG)

1: **procedure** FEDAVG($w_0, \eta_c, \eta_s$)
2:     **for** $r = 0, \ldots, R-1$ **do**
3:         sample a subset of clients $\mathcal{S}_r \subseteq [M]$
4:         **on client** $m \in \mathcal{S}_r$ **in parallel do**
5:             client initialization $w_{r,0}^m \leftarrow w_r$
6:             **for** $k = 0, \ldots, K-1$ **do**
7:                 $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$
8:                 $w_{r,k+1}^m \leftarrow w_{r,k}^m - \eta_c \cdot g_{r,k}^m$
9:         $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r} (w_{r,K}^m - w_{r,0}^m)$
10:        $w_{r+1} \leftarrow w_r + \eta_s \cdot \Delta_r$
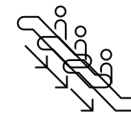
**Algorithm 2** Federated PGD

1: **procedure** FEDPGD($w_0, \eta_c, \eta_s$)
2:     **for** $r = 0, \ldots, R-1$ **do**
3:         sample a subset of clients $\mathcal{S}_r \subseteq [M]$
4:         **on client** $m \in \mathcal{S}_r$ **in parallel do**
5:             client initialization $w_{r,0}^m \leftarrow w_r$
6:             **for** $k = 0, \ldots, K-1$ **do**
7:                 $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$
8:                 $w_{r,k+1}^m \leftarrow \mathbf{prox}_{\eta_c \psi}(w_{r,k}^m - \eta_c g_{r,k}^m)$
9:         $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r} (w_{r,K}^m - w_{r,0}^m)$
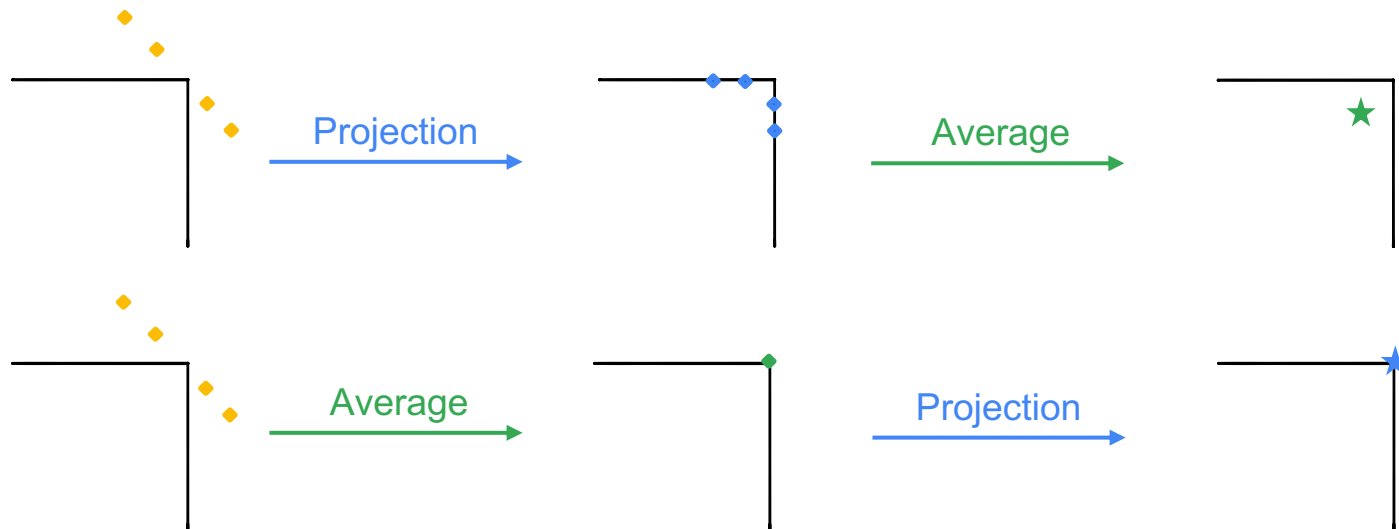10:        $w_{r+1} \leftarrow \mathbf{prox}_{\eta_s \eta_c K \psi}(w_r + \eta_s \Delta_r)$

$$\min_{w \in \mathbb{R}^d} \Phi(w) := \frac{1}{M} \sum_{m=1}^{M} F_m(w) + \psi(w)$$

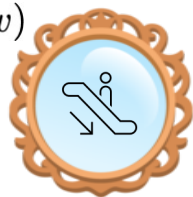# First Attempt: FEDAVG + Proximal Gradient Descent

- Challenge: Averaging and proximal operations discord

  - Averaging and (nonlinear) proximal operators **do not commute**

  - Intuition: Averaging on post-projected points "blunt" the sharpness of projection

$$\min \; F(w) + \psi(w)$$

# Composite 201: (composite) Mirror Descent

*[Nemirovski et al., '83, Duchi et al., COLT'10]*

$$\text{PGD}: \quad w_{t+1} = \underset{w}{\arg\min}\left\{ F(w_t) + \langle \nabla F(w_t), w - w_t \rangle + \psi(w) + \frac{1}{2\eta}\|w - w_t\|_2^2 \right\}$$

Arbitrary distance-generating $h$

$$\text{MD}: \quad w_{t+1} = \underset{w}{\arg\min}\left\{ F(w_t) + \langle \nabla F(w_t), w - w_t \rangle + \psi(w) + \frac{1}{\eta}D_h(w, w_t) \right\}$$
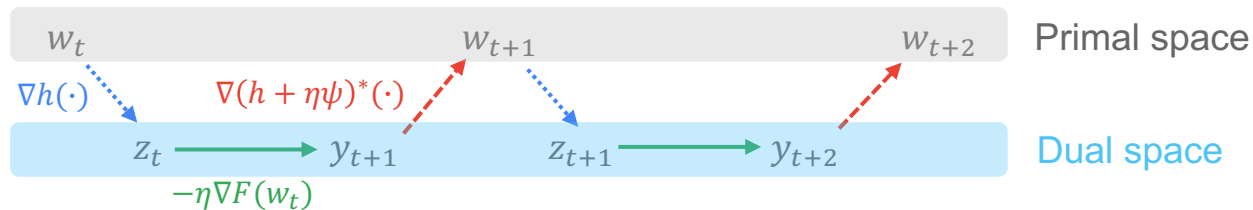
$$D_h(w, w_t) := h(w) - h(w_t) - \langle \nabla h(w_t), w - w_t \rangle$$

*(reduces to PGD if $h(w) = \frac{1}{2}\|w\|^2$ )*

Primal-dual interpretation of MD

- $z_t = \nabla h(w_t)$    Forward mirror (Primal -> Dual)
- $y_{t+1} = z_t - \eta \cdot \nabla F(w_t)$    Gradient step (in dual space)
- $w_{t+1} = \nabla(h + \eta\psi)^*(y_{t+1})$    Backward mirror (Dual -> Primal)

*\* indicates convex conjugate*



$w_t$     $w_{t+1}$     $w_{t+2}$    Primal space

$\nabla h(\cdot)$    $\nabla(h + \eta\psi)^*(\cdot)$

$z_t$    $y_{t+1}$    $z_{t+1}$    $y_{t+2}$    Dual space

$-\eta\nabla F(w_t)$

# Federated Mirror Descent (FEDMID)

- Federated Mirror Descent (FEDMID) generalizes Federated PGD

**Algorithm 2** Federated PGD

1: **procedure** FEDPGD($w_0, \eta_c, \eta_s$)
2:     **for** $r = 0, \dots, R-1$ **do**
3:         sample a subset of clients $\mathcal{S}_r \subseteq [M]$
4:         **on client** $m \in \mathcal{S}_r$ **in parallel do**
5:           client initialization $w_{r,0}^m \leftarrow w_r$
6:           **for** $k = 0, \dots, K-1$ **do**
7:             $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$
8:             $w_{r,k+1}^m \leftarrow \mathbf{prox}_{\eta_c \psi}(\nabla h(w_{r,k}^m) - \eta_c g_{r,k}^m)$
9:         $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r}(w_{r,K}^m - w_{r,0}^m)$
10:       $w_{r+1} \leftarrow \mathbf{prox}_{\eta_s \eta_c K \psi}(\nabla h(w_r) + \eta_s \Delta_r)$

**Algorithm 2** Federated Mirror Descent (FEDMID)

1: **procedure** FEDMID($w_0, \eta_c, \eta_s$)
2:     **for** $r = 0, \dots, R-1$ **do**
3:         sample a subset of clients $\mathcal{S}_r \subseteq [M]$
4:         **on client** $m \in \mathcal{S}_r$ **in parallel do**
5:           client initialization $w_{r,0}^m \leftarrow w_r$
6:           **for** $k = 0, \dots, K-1$ **do**
7:             $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$
8:             $w_{r,k+1}^m \leftarrow \nabla(h + \eta_c \psi)^*(\nabla h(w_{r,k}^m) - \eta_c g_{r,k}^m)$
9:         $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r}(w_{r,K}^m - w_{r,0}^m)$
10:       $w_{r+1} \leftarrow \nabla(h + \eta_s \eta_c K \psi)^*(\nabla h(w_r) + \eta_s \Delta_r)$

# Composite 202: Dual Averaging

[Nesterov et al., '09, Xiao et al., '10, Flammarion et al., COLT'17]

$$\min \ F(w) + \psi(w)$$

Dual Averaging (a.k.a. Lazy Mirror Descent)

$$w_t = \nabla(h + \eta t \psi)^*(z_t)$$

Backward mirror (Dual -> Primal) – retrieve primal

$$= \arg \min_{w} \{\langle -z_t, w \rangle + \eta t \psi(w) + h(w)\}$$

$$z_{t+1} = z_t - \eta \cdot \nabla F(w_t)$$

Gradient step (in dual space)



Recall MD:



Google

# Mirror Descent vs Dual Averaging

### Mirror Descent

$$w_t \qquad\qquad w_{t+1}$$

$\nabla h(\cdot)$

$\nabla (h + \eta\psi)^*(\cdot)$

$$z_t \longrightarrow y_{t+1}$$

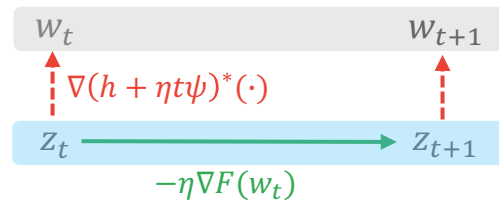$-\eta \nabla F(w_t)$

- Forward **and** backward mirror
- Persistent **primal** states

### Dual Averaging

$$w_t \qquad\qquad w_{t+1}$$

$\nabla (h + \eta t\psi)^*(\cdot)$

$$z_t \longrightarrow z_{t+1}$$

$-\eta \nabla F(w_t)$

- Backward mirror **only**
- Persistent **dual** states

# Federated Dual Averaging (FEDDUALAVG)

Notation: $w_{r,k}^m$ → $m^{th}$ client
→ $r^{th}$ round, $k^{th}$ local iteration

# Federated Dual Averaging (FEDDUALAVG)

Notation: $w_{r,k}^m$ → $m^{th}$ client

→ $r^{th}$ round, $k^{th}$ local iteration
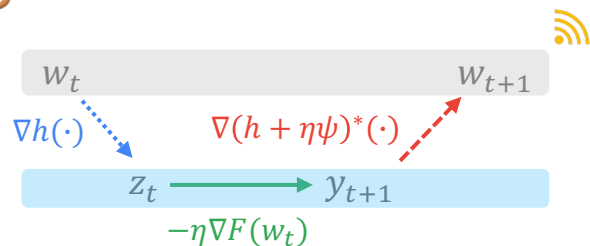
---

**Algorithm 3** Federated Dual Averaging

1: **procedure** FEDDUALAVG($w_0, \eta_c, \eta_s$)
2:   server initialization $z_0 \leftarrow \nabla h(w_0)$
3:   **for** $r = 0, \dots, R-1$ **do**
4:     sample a subset of clients $\mathcal{S}_r \subseteq [M]$
5:     **on client $m \in \mathcal{S}_r$ in parallel do**
6:       client initialization $z_{r,0}^m \leftarrow z_r$
7:       **for** $k = 0, \dots, K-1$ **do**
8:         $\tilde{\eta}_{r,k} \leftarrow \eta_s \eta_c r K + \eta_c k$
9:         $w_{r,k}^m \leftarrow \nabla(h + \tilde{\eta}_{r,k}\psi)^*(z_{r,k}^m)$
10:        $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$
11:        $z_{r,k+1}^m \leftarrow z_{r,k}^m - \eta_c g_{r,k}^m$
12:    $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r} (z_{r,K}^m - z_{r,0}^m)$
13:    $z_{r+1} \leftarrow z_r + \eta_s \Delta_r$
14:    $w_{r+1} \leftarrow \nabla(h + \eta_s \eta_c (r+1) K \psi)^*(z_{r+1})$

→ Compute primal point

→ Client **dual** update

→ Average client **dual** deltas

→ Server **dual** update

→ (Optional) primal output

Google

# FEDMID (a.k.a. FEDPGD) vs FEDDUALAVG

## Mirror Descent

$$w_t \qquad\qquad\qquad w_{t+1}$$

$$\nabla h(\cdot) \qquad \nabla(h + \eta\psi)^*(\cdot)$$

$$z_t \longrightarrow y_{t+1}$$

$$-\eta\nabla F(w_t)$$
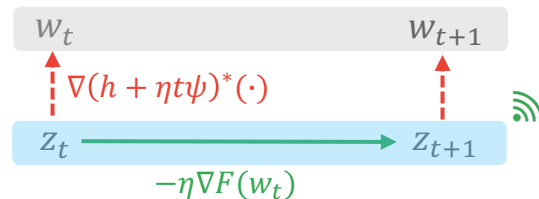
- Forward **and** backward mirror

- Persistent **primal** updates

- FEDMID: average the **primal**

- theoretically challenging due to the **nonlinearity of mirror map**.

## Dual Averaging

$$w_t \qquad\qquad\qquad w_{t+1}$$

$$\nabla(h + \eta t\psi)^*(\cdot)$$

$$z_t \longrightarrow z_{t+1}$$

$$-\eta\nabla F(w_t)$$

- Backward mirror **only**

- Persistent **dual** updates

- FEDDUALAVG: average the **dual**

- Enjoys nice theoretical interpretation via **dual shadow sequence**.

- outperforms FEDMID empirically.

Google

# Theory: Blanket Assumptions

$$\min_{w \in \mathbb{R}^d} \quad \Phi(w) := \frac{1}{M} \sum_{m=1}^{M} F_m(w) + \psi(w)$$

$$\text{where} \quad F_m(w) := \mathbb{E}_{\xi \sim \mathcal{D}_m}[f(w; \xi)]$$

**Assumption 1.** *Let* $\|\cdot\|$ *be an arbitrary norm and* $\|\cdot\|_*$ *be its dual norm.*

(a) $\psi : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ *is a closed convex function with closed* $\mathbf{dom}\,\psi$. *Assume* $\Phi(w) = F(w) + \psi(w)$ *attains a finite optimum at* $\theta^\star \in \mathbf{dom}\,\psi$.

(b) $h : \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$ *is a Legendre function that is 1-strongly convex with respect to* $\|\cdot\|$. *Assume* $\mathbf{dom}\,h \supset \mathbf{dom}\,\psi$.

*(a) & (b): standard regularity assumptions for composite setup*

(c) $f(\cdot, \xi) : \mathbb{R}^d \to \mathbb{R}$ *is a closed convex function that is differentiable on* $\mathbf{dom}\,\psi$ *for any fixed* $\xi$. *In addition,* $f(\cdot, \xi)$ *is L-smooth on* $\mathbf{dom}\,\psi$, *namely for any* $u, w \in \mathbf{dom}\,\psi$,

$$f(u; \xi) \leq f(w; \xi) + \langle \nabla f(w; \xi), u - w \rangle + \frac{1}{2} L \|u - w\|^2.$$

*(c): smoothness of f*

(d) $\nabla f$ *has* $\sigma^2$-*bounded variance under* $\|\cdot\|_*$ *norm within* $\mathbf{dom}\,\psi$, *namely for any* $w \in \mathbf{dom}\,\psi$,

$$\mathbb{E}_{\xi \sim \mathcal{D}_m} \|\nabla f(w, \xi) - \nabla F_m(w)\|_*^2 \leq \sigma^2.$$

*(d): additive bounded variance*

(e) *Assume all the M clients participate in client updates for every round, namely* $\mathcal{S}_r = [M]$.

*(e): full participation (for simplicity of exposition)*

Google

# Theorem 1: Small Client Learning Rate $\eta_c$ Regime

In small $\eta_c$ regime, both FEDMID and FEDDUALAVG can match minibatch rate

**Theorem 1.** Assuming A1, for **sufficiently small** $\eta_c$, and appropriate $\eta_s$, both FEDMID and FEDDUALAVG can output $\hat{w}$ such that

$$\mathbb{E}\left[\Phi(\hat{w})\right] - \Phi(w^\star) \lesssim \frac{LB}{R} + \frac{\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}}$$

where $B := D_h(w^\star, w_0)$ is the Bregman divergence distance between optimum $w^\star$ and initial $w_0$

*L: smoothness*
*σ: variance bound*
*M: # of clients*
*K: # of local steps*
*R: # of rounds*

# Stronger Guarantee for FEDDUALAVG (bounded gradient)

We establish (possibly) stronger guarantee for **FEDDUALAVG** with larger $\eta_c$ and unit $\eta_s = 1$

**Theorem 2.** Assuming A1, and in addition assume $\displaystyle\sup_{w \in \mathbf{dom}\psi} \|\nabla f(w, \xi)\|_* \leq G$, then for $\eta_s = 1$ and $\eta_c \leq \frac{1}{4L}$, FEDDUALAVG can output $\hat{w}$ such that

$$\mathbb{E}\left[\Phi\left(\hat{w}\right)\right] - \Phi(w^\star) \lesssim \frac{B}{\eta_c KR} + \frac{\eta_c \sigma^2}{M} + \eta_c^2 LK^2G^2$$

Moreover for appropriate $\eta_c$

*faster convergence (usefulness of client step)*

*Overhead for infrequent communication*

$$\mathbb{E}\left[\Phi\left(\hat{w}\right)\right] - \Phi(w^\star) \lesssim \frac{LB}{KR} + \frac{\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}} + \frac{L^{\frac{1}{3}} B^{\frac{2}{3}} G^{\frac{2}{3}}}{R^{\frac{2}{3}}}.$$

*matches [Stich ICLR'19] bound on smooth unconstrained FEDAVG*

$B := D_h(w^\star, w_0)$
$L$: smoothness
$\sigma$: variance bound
$M$: # of clients
$K$: # of local steps
$R$: # of rounds

Google

# Stronger Guarantee for FEDDUALAVG (quadratic $F$)

We can relax the bounded gradient assumption if $F$ is quadratic, and heterogeneity is bounded.

**Theorem 3.** Assuming A1, and in addition assume $\displaystyle\sup_{w \in \mathbf{dom}\,\psi} \|\nabla F_m(w) - \nabla F(w)\|_* \leq \zeta^2$

and $F$ is quadratic, then FEDDUALAVG can output $\hat{w}$ such that

$$\mathbb{E}\left[\Phi\left(\hat{w}\right)\right] - \Phi(w^\star) \lesssim \frac{B}{\eta_{\mathrm{c}}KR} + \frac{\eta_{\mathrm{c}}\sigma^2}{M} + \eta_{\mathrm{c}}^2 LK\sigma^2 + \eta_{\mathrm{c}}^2 LK^2\zeta^2,$$

moreover for appropriate $\eta_c$

*faster convergence
(usefulness of client step)*

*Overhead for infrequent communication*

$$\mathbb{E}\left[\Phi\left(\hat{w}\right)\right] - \Phi(w^\star) \lesssim \frac{LB}{KR} + \frac{\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}}K^{\frac{1}{2}}R^{\frac{1}{2}}} + \frac{L^{\frac{1}{3}}B^{\frac{2}{3}}\sigma^{\frac{2}{3}}}{K^{\frac{1}{3}}R^{\frac{2}{3}}} + \frac{L^{\frac{1}{3}}B^{\frac{2}{3}}\zeta^{\frac{2}{3}}}{R^{\frac{2}{3}}}.$$

*matches best known bound on
<u>smooth unconstrained</u> FEDAVG
[Khaled AISTATS'20,
Woodworth NeurIPS'20 etc]*

$B := D_h(w^\star, w_0)$
*L: smoothness*
*σ: variance bound*
*M: # of clients*
*K: # of local steps*
*R: # of rounds*

Google

# Summary of Theoretical Results

- FEDMID & FEDDUALAVG, small $\eta_c$:

$$\frac{LB}{R} + \frac{\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}}$$

- FEDDUALAVG, larger $\eta_c$:

$$\frac{LB}{KR} + \frac{\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}} + \frac{L^{\frac{1}{3}} B^{\frac{2}{3}} G^{\frac{2}{3}}}{R^{\frac{2}{3}}}$$

$$\frac{LB}{KR} + \frac{\sigma B^{\frac{1}{2}}}{M^{\frac{1}{2}} K^{\frac{1}{2}} R^{\frac{1}{2}}} + \frac{L^{\frac{1}{3}} B^{\frac{2}{3}} \sigma^{\frac{2}{3}}}{K^{\frac{1}{3}} R^{\frac{2}{3}}} + \frac{L^{\frac{1}{3}} B^{\frac{2}{3}} \zeta^{\frac{2}{3}}}{R^{\frac{2}{3}}}$$

$B := D_h(w^\star, w_0)$    K: # of local steps
L: smoothness    R: # of rounds
σ: variance bound    G: gradient bound
M: # of clients    ζ: heterogeneity bound

# Proof Sketch -- FedDualAvg

**Main observation**: the averaged dual $\overline{z_{r,k}} := \frac{1}{M} \sum_{m=1}^{M} z_{r,k}^m$ "almost" does **centralized dual averaging**

$$\overline{z_{r,k+1}} = \overline{z_{r,k}} - \eta_c \cdot \frac{1}{M} \sum_{m=1}^{M} \nabla f(w_{r,k}^m; \xi_{r,k}^m)$$

*Variance-reduced* but **biased** *stochastic gradient oracle*

**Step 1:** convergence of the averaged dual (a.k.a. perturbed iterate analysis)

$$\mathbb{E}\left[\Phi\left(\frac{1}{KR}\sum_{r=0}^{R-1}\sum_{k=1}^{K}\nabla\left(h+\tilde{\eta}_{r,k}\psi\right)^*\left(\overline{z_{r,k}}\right)\right)\right] - \Phi(w^\star) \leq \underbrace{\frac{B}{\eta_c KR} + \frac{\eta_c \sigma^2}{M}}_{\substack{Rate\ if\ synchronize \\ every\ iterations}} + \frac{L}{MKR}\underbrace{\left[\sum_{r=0}^{R-1}\sum_{k=0}^{K-1}\sum_{m=1}^{M}\mathbb{E}\left\|\overline{z_{r,k}} - z_{r,k}^m\right\|_*^2\right]}_{Discrepancy\ overhead}$$

**Step 2:** bound $\mathbb{E}\|\overline{z_{r,k}} - z_{r,k}^m\|_*^2$ by stability analysis

Google

# Experiments

- Platform setup: *TensorFlow/Federated* & *google-research/federated*

- We evaluate the following 4 algorithms:

  1. Federated Dual Averaging (FEDDUALAVG)

  2. Federated Mirror Descent (FEDMID)

  3. FEDDUALAVG-OSP (only-server-proximal)

  4. FEDMID-OSP (only-server-proximal)

*potential light computation but less principled - for ablation study purpose*

# FEDMID VS FEDMID-OSP

**Algorithm 2** Federated Mirror Descent (FEDMID)

1: **procedure** FEDMID($w_0, \eta_c, \eta_s$)
2:    **for** $r = 0, \ldots, R-1$ **do**
3:       sample a subset of clients $\mathcal{S}_r \subseteq [M]$
4:       **on client** $m \in \mathcal{S}_r$ **in parallel do**
5:         client initialization $w_{r,0}^m \leftarrow w_r$
6:         **for** $k = 0, \ldots, K-1$ **do**
7:           $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$
8:           $w_{r,k+1}^m \leftarrow \nabla(h + \eta_c \psi)^*(\nabla h(w_{r,k}^m) - \eta_c g_{r,k}^m)$
9:       $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r}(w_{r,K}^m - w_{r,0}^m)$
10:     $w_{r+1} \leftarrow \nabla(h + \eta_s \eta_c K \psi)^*(\nabla h(w_r) + \eta_s \Delta_r)$

**Algorithm 4** Federated Mirror Descent Only Server

1: **procedure** FEDMID-OSP($w_0, \eta_c, \eta_s$)
2:    **for** $r = 0, \ldots, R-1$ **do**
3:       sample a subset of clients $\mathcal{S}_r \subseteq [M]$
4:       **on client** $m \in \mathcal{S}_r$ **in parallel do**
5:         client initialization $w_{r,0}^m \leftarrow w_r$
6:         **for** $k = 0, \ldots, K-1$ **do**
7:           $g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$
8:           $w_{r,k+1}^m \leftarrow \nabla h^*(\nabla h(w_{r,k}^m) - \eta_c g_{r,k}^m)$    ▷
9:       $\Delta_r = \frac{1}{|\mathcal{S}_r|} \sum_{m \in \mathcal{S}_r}(w_{r,K}^m - w_{r,0}^m)$
10:     $w_{r+1} \leftarrow \nabla(h + \eta_s \eta_c K \psi)^*(\nabla h(w_r) + \eta_s \Delta_r)$

*Proximal $\psi$ skipped*

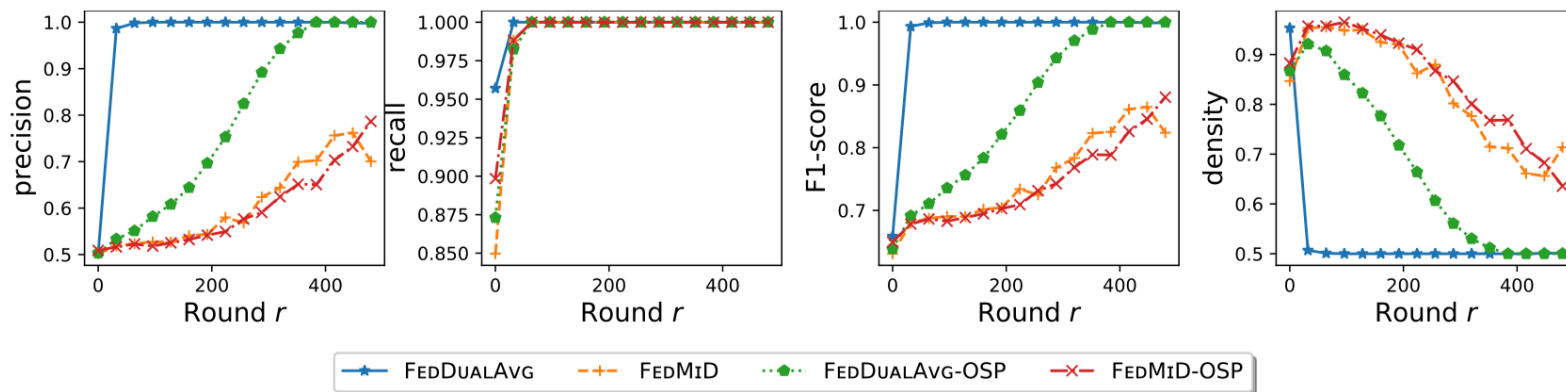*Reduces to $w_{r,k+1}^m \leftarrow w_{r,k}^m - \eta_c g_{r,k}^m$ if $h = \frac{1}{2}||\cdot||^2$*

Google

# Experiment 1: Federated Lasso on Synthetic Dataset

- **Synthetic dataset**: $y = x^T w^\star + b^\star + \varepsilon$; known sparse ground truth $w^*$

    *(64 clients, 128 samples per client, ground truth density 512/1024)*

- **Problem**: $\min\limits_{w \in \mathbb{R}^d, b \in \mathbb{R}} \dfrac{1}{M} \sum\limits_{m=1}^{M} \mathbb{E}_{(x,y) \sim \mathcal{D}_m} (x^\top w + b - y)_2^2 + \lambda \|w\|_1$
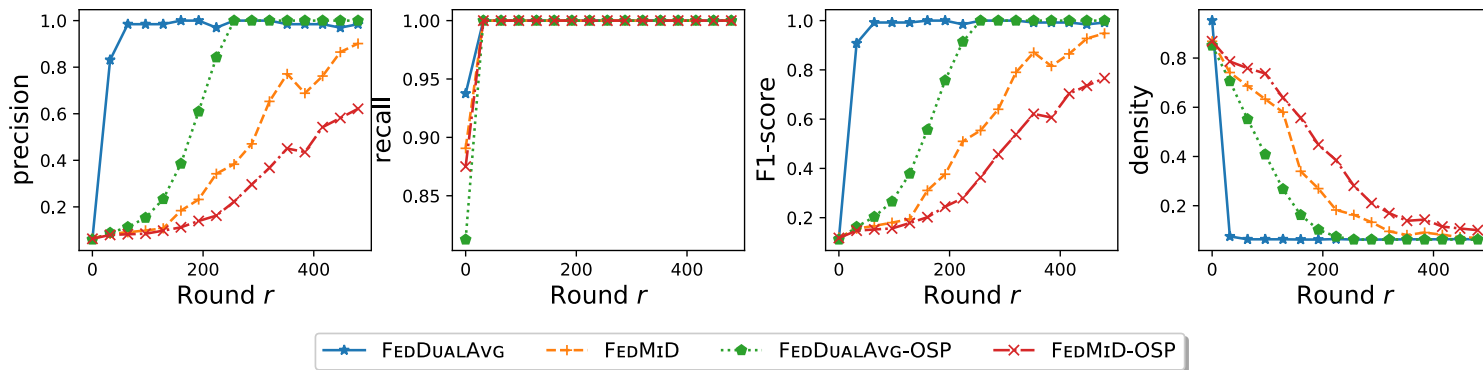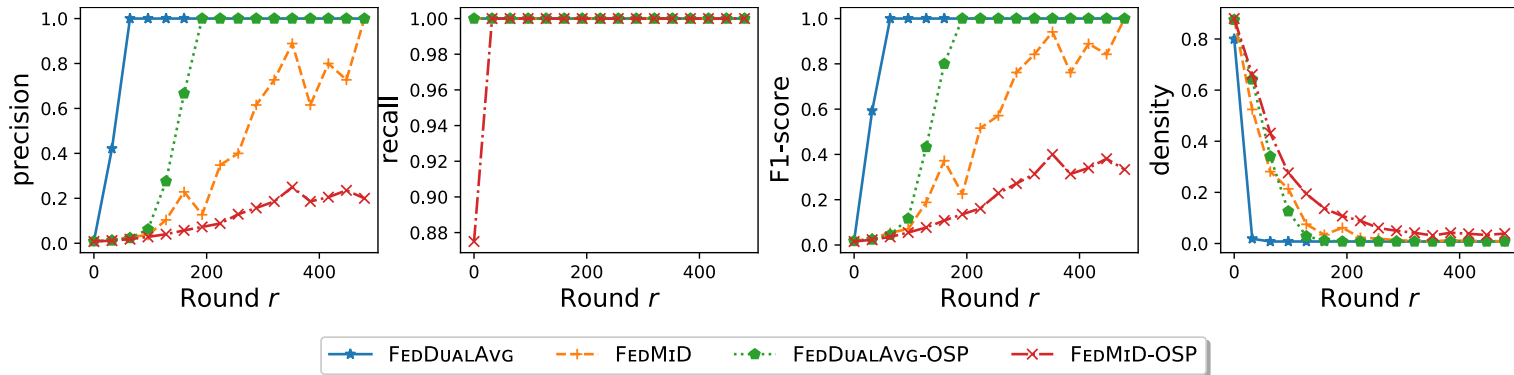
- **Metric**: F1-score of the estimated sparsity, precision, recall, density



*For all algorithms, we tune only $\eta_s$*
*and $\eta_c$ to attain the best F1-score*

Google

# Experiment 1: Sparser Ground Truth

- **Sparser dataset:** *(64 clients, 128 samples per client, ground truth density **64**/1024)*



- **Even sparser dataset:** *(64 clients, 128 samples per client, ground truth density **8**/1024)*

# Experiment 1: More Distributed Data

- **Even more distributed:** *(256 clients, 32 samples per client, ground truth density 512/1024)*
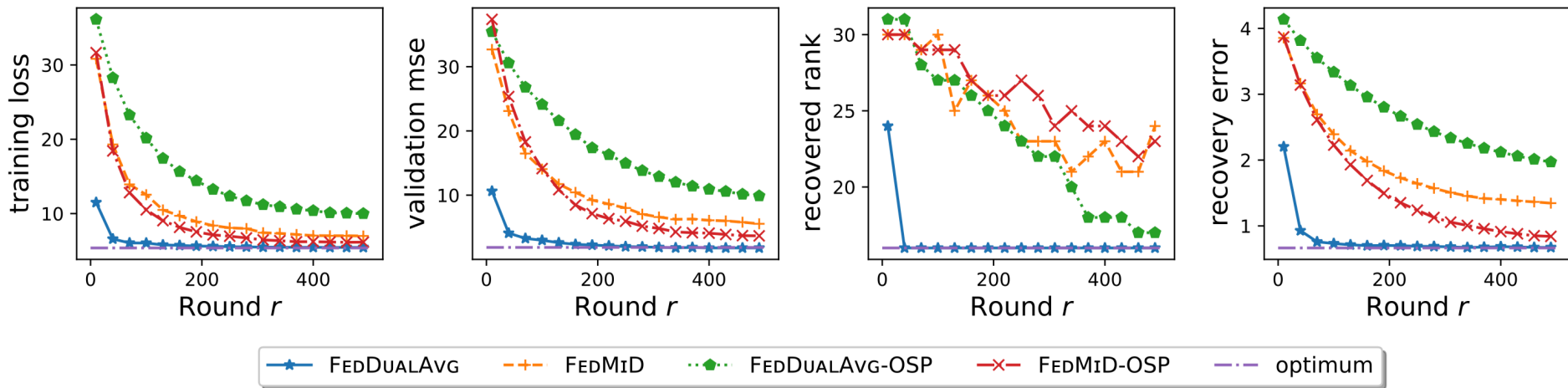
# Experiment 2: Low-Rank Matrix Estimation

- **Synthetic dataset**: $y = \langle X, W^\star \rangle + b^\star + \varepsilon$; known **low-rank** ground truth $W^*$

    *(64 clients, 128 samples per client, ground truth rank **16**/32)*

- **Problem**: $\displaystyle \min_{W \in \mathbb{R}^{d_1 \times d_2}, b \in \mathbb{R}} \frac{1}{M} \sum_{m=1}^{M} \mathbb{E}_{(X,y) \sim \mathcal{D}_m} \left( \langle X, W \rangle + b - y \right)^2 + \lambda \|W\|_{\text{nuc}}$

- **Metric**: training loss, validation mse, recovered rank, recovered error (in Frobenius norm)

# Experiment 2: Sparser Ground Truth

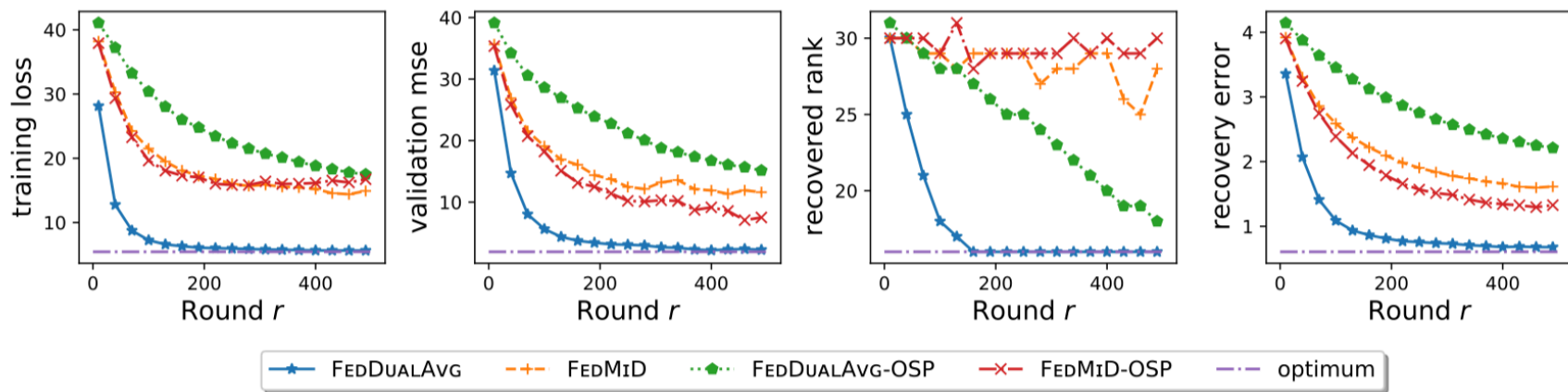- **Lower rank dataset:** *(64 clients, 128 samples per client, ground truth rank **4**/32)*



- **Even lower rank dataset:** *(64 clients, 128 samples per client, ground truth rank **1/32**)*

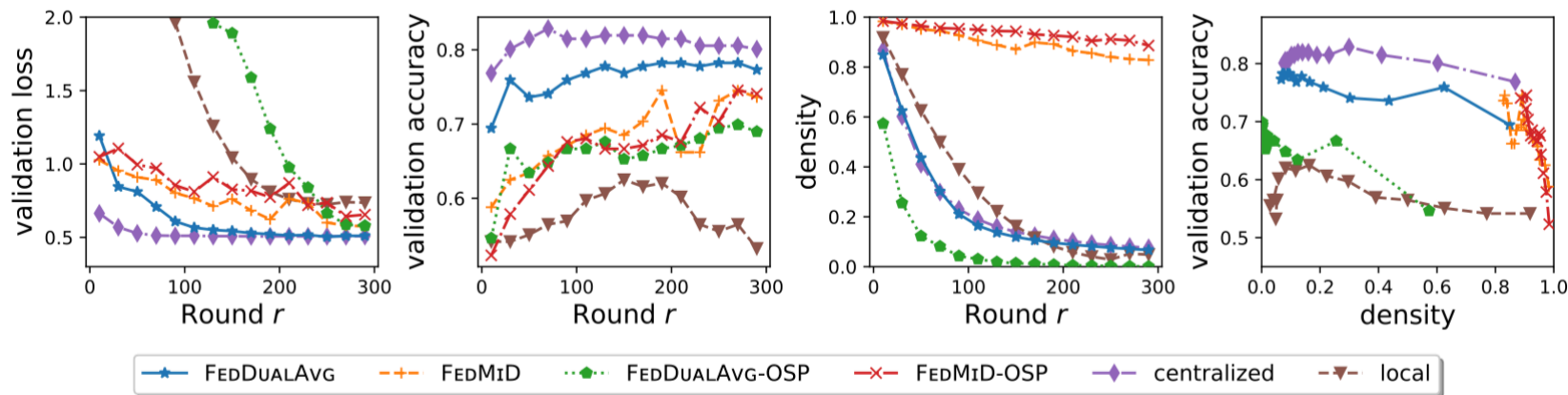

Google

# Experiment 2: More Distributed Data

● **More distributed:** *(256 clients, 32 samples per client, ground truth density 512/1024)*
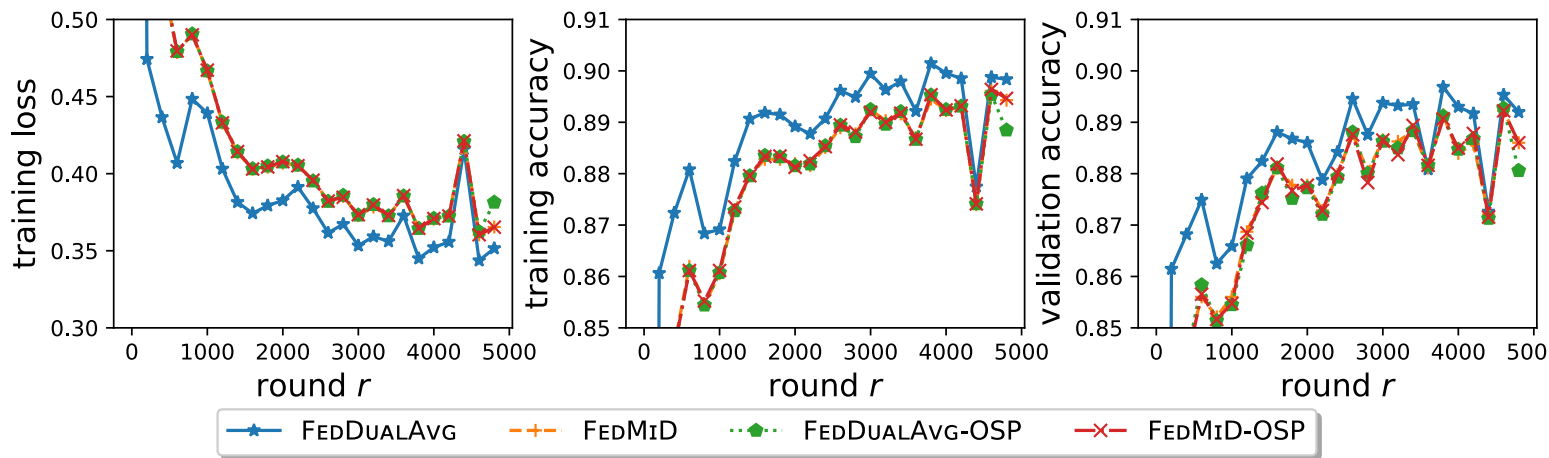
# Experiment 3: Sparse Logistic Regression for fMRI

- **Dataset**: fMRI scans on response to binary image recognition

  *(6 subjects, 11-12 sessions per subject, 18 scans per session, 39,912 voxels)*

- **Federated Setup:** Each client possesses the data of a **session**. (59 training clients in total)

- **Problem**: l1-regularized logistic regression
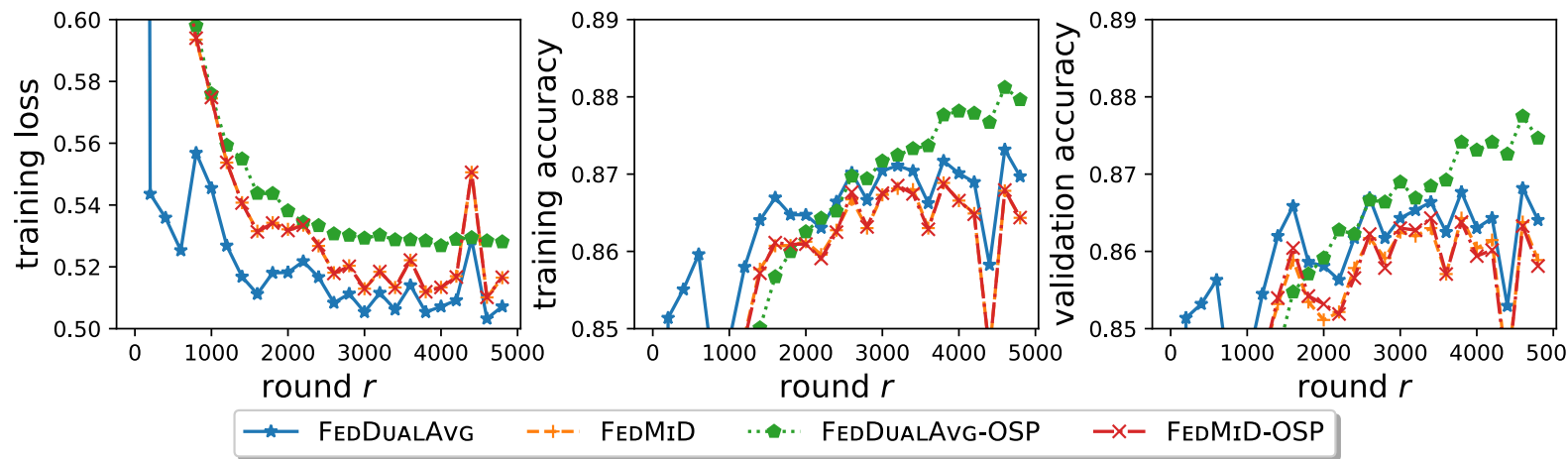
- **Metric**: density, validation accuracy

# Experiment 4: norm-ball constrained FL

- **Dataset**: Federated EMNIST (10 classes or 62 classes)

- **Metric**: Training loss, training accuracy, test accuracy

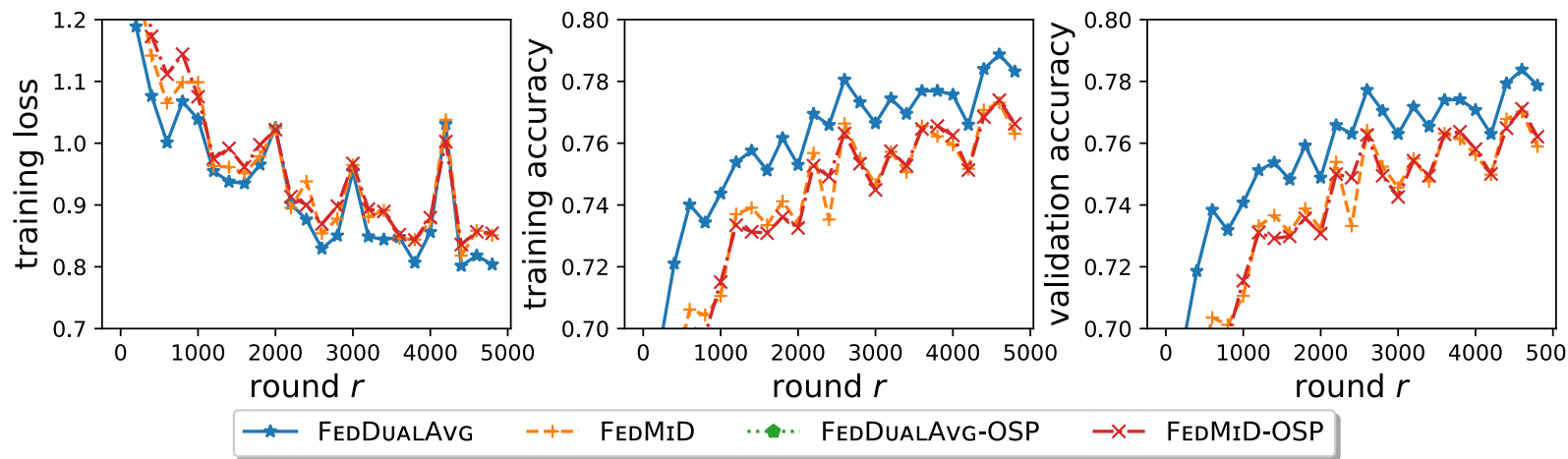- *L1-constrained logistic regression for EMNIST-10*

# Experiment 4: norm-ball constrained FL

- *L2-constrained logistic regression for EMNIST-10*

# Experiment 4: norm-ball constrained FL

- *L1-constrained 2-hidden-layer NN on EMNIST-62*

# Thank you!

Paper: https://arxiv.org/abs/2011.08474

Email: yuanhl@cs.stanford.edu