

FedNest: Federated Bilevel, Minimax, and Compositional Optimization

Davoud Ataee Tarzanagh¹ Mingchen Li² Christos Thrampoulidis³
Samet Oymak²

¹ University of Michigan

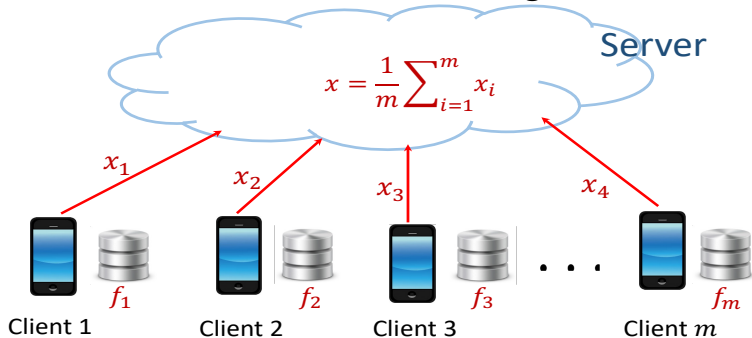
² University of California, Riverside

³ University of British Columbia

ICML, July 17-23 2022

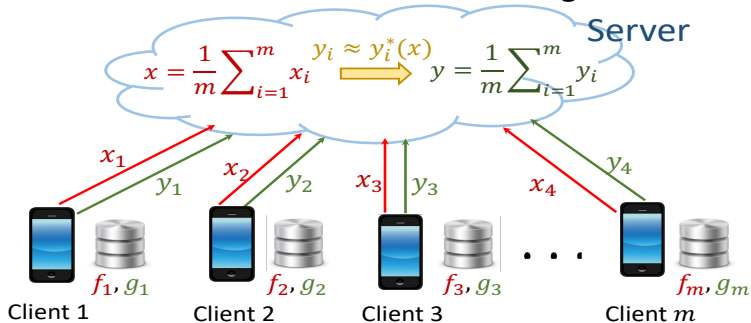


Federated Learning



Goal:
$$\min_{x \in \mathbb{R}^{d_1}} \frac{1}{m} \sum_{i=1}^m f_i(x)$$

Federated Bilevel Learning



Goal:

$$\begin{aligned}
 \min_{\mathbf{x} \in \mathbb{R}^{d_1}} \quad & f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) && \text{(outer)} \\
 \text{subj. to} \quad & \mathbf{y}^*(\mathbf{x}) \in \underset{\mathbf{y} \in \mathbb{R}^{d_2}}{\text{argmin}} \quad \frac{1}{m} \sum_{i=1}^m g_i(\mathbf{x}, \mathbf{y}) && \text{(inner)}
 \end{aligned}$$

Federated Bilevel Optimization (FBO)



Our Setting:

- **Stochastic:** Access to (f_i, g_i) is via stochastic sampling:

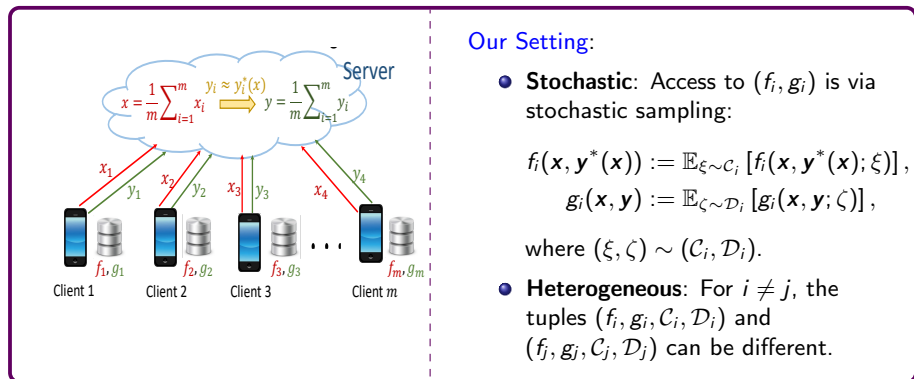
$$f_i(x, y^*(x)) := \mathbb{E}_{\xi \sim \mathcal{C}_i} [f_i(x, y^*(x); \xi)],$$

$$g_i(x, y) := \mathbb{E}_{\zeta \sim \mathcal{D}_i} [g_i(x, y; \zeta)],$$

where $(\xi, \zeta) \sim (\mathcal{C}_i, \mathcal{D}_i)$.

- **Heterogeneous:** For $i \neq j$, the tuples $(f_i, g_i, \mathcal{C}_i, \mathcal{D}_i)$ and $(f_j, g_j, \mathcal{C}_j, \mathcal{D}_j)$ can be different.

Federated Bilevel Optimization (FBO)



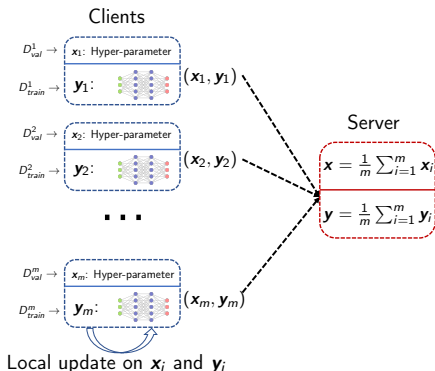
Applications of FBO: meta-learning, hyperparameter optimization, neural network architecture search, actor-critic reinforcement learning, GANs, . . .

Motivating Example

- **Federated Hyper-Parameter Optimization:** Collaboratively find **machine learning (ML) model** and **the hyper-parameters** while keeping the data decentralized

Motivating Example

- **Federated Hyper-Parameter Optimization:** Collaboratively find **machine learning (ML) model** and **the hyper-parameters** while keeping the data decentralized



- **Inner objective:**
$$\frac{1}{m} \sum_{i=1}^m g(x, y; D_{train}^i)$$
- **Outer objective:**
$$\frac{1}{m} \sum_{i=1}^m f(y^*(x); D_{val}^i)$$
- y is an **ML model** such as neural network.
- x contains **hyper-parameters** such as
 - regularization parameters,
 - learning rates, and
 - batch size.

Two Special Cases

- FBO subsumes two popular problem classes with the nested structure.

Federated Minimax Optimization

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^{d_1}} \quad & f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x}); \xi) \\ \text{subj. to} \quad & \mathbf{y}^*(\mathbf{x}) = \underset{\mathbf{y} \in \mathbb{R}^{d_2}}{\operatorname{argmin}} - \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}, \mathbf{y}; \xi) \end{aligned}$$

FBO with

$$g_i(\mathbf{x}, \mathbf{y}; \zeta) = -f_i(\mathbf{x}, \mathbf{y}; \xi).$$

Application:

- Training Generative-adversarial Networks (GANs)

Two Special Cases

- FBO subsumes two popular problem classes with the nested structure.

Federated Minimax Optimization

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^{d_1}} \quad & f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x}); \xi) \\ \text{subj. to } \quad & \mathbf{y}^*(\mathbf{x}) = \underset{\mathbf{y} \in \mathbb{R}^{d_2}}{\operatorname{argmin}} - \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}, \mathbf{y}; \xi) \end{aligned}$$

FBO with

$$g_i(\mathbf{x}, \mathbf{y}; \zeta) = -f_i(\mathbf{x}, \mathbf{y}; \xi).$$

Application:

- Training Generative-adversarial Networks (GANs)

Federated Compositional Optimization

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^{d_1}} \quad & f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{y}; \xi) \\ \text{subj. to } \quad & \mathbf{y}^*(\mathbf{x}) = \underset{\mathbf{y} \in \mathbb{R}^{d_2}}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m \|\mathbf{y} - \mathbf{r}_i(\mathbf{x}; \zeta)\|^2 \end{aligned}$$

FBO with

$$\begin{aligned} f_i(\mathbf{x}, \mathbf{y}; \xi) &= f_i(\mathbf{y}; \xi) \text{ and} \\ g_i(\mathbf{x}, \mathbf{y}; \zeta) &= \|\mathbf{y} - \mathbf{r}_i(\mathbf{x}; \zeta)\|^2. \end{aligned}$$

Application:

- Model Agnostic Meta-Learning (MAML)

Federated (Single-Level) Optimization

$$\min_{\mathbf{x} \in \mathbb{R}^{d_1}} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\xi \sim \mathcal{C}_i} [f_i(\mathbf{x}; \xi)]$$

Federated (Single-Level) Optimization

$$\min_{\mathbf{x} \in \mathbb{R}^{d_1}} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\xi \sim \mathcal{C}_i} [f_i(\mathbf{x}; \xi)]$$

Gradient-Type Federated Optimization

For $k = 0, \dots, K - 1$:

① i -th client:

- For $\nu = 0, \dots, \tau_i - 1$:

$$\mathbf{x}_{i,\nu+1}^k = \mathbf{x}_{i,\nu}^k + \alpha_i^k \mathbf{h}_{i,\nu}^k$$

② Server:


$$\mathbf{x}^{k+1} = 1/m \sum_{i=1}^m \mathbf{x}_{i,\tau_i}^k$$

- τ_i is number of local iterations
- α_i^k is the stepsize
- **FedAvg** (McMahan et al., 2017):
 $\mathbf{h}_{i,\nu}^k = -\nabla f_i(\mathbf{x}_{i,\nu}^k; \xi_{i,\nu}^k)$
- **FedSVRG** (Konečný et al., 2018):
 $\mathbf{h}_{i,\nu}^k = -\nabla f_i(\mathbf{x}_{i,\nu}^k; \xi_{i,\nu}^k) +$
 $\nabla f_i(\mathbf{x}^k; \xi_{i,\nu}^k) - \frac{1}{m} \sum_{i=1}^m \nabla f_i(\mathbf{x}^k; \xi_i^k)$

Challenges in FBO

- **FedAvg** can lead to convergence to a point different from $\mathbf{y}^*(\mathbf{x})$.
- Each client i requires access to the **global Hessian inverse**:

$$\nabla f_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) = \nabla_{\mathbf{x}} f_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) - \nabla_{\mathbf{x}\mathbf{y}}^2 g(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) \cdot \underbrace{\left[\sum_{i=1}^m \nabla_{\mathbf{y}}^2 g_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) \right]^{-1} \nabla_{\mathbf{y}} f_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x}))}_{\mathbf{p}_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x}))}$$

- 
- **FedAvg** can lead to convergence to a point different from $\mathbf{y}^*(\mathbf{x})$.
 - Each client i requires access to the **global Hessian inverse**:

$$\nabla f_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) = \nabla_{\mathbf{x}} f_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) - \nabla_{\mathbf{x}\mathbf{y}}^2 g(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) \cdot \underbrace{\left[\sum_{i=1}^m \nabla_{\mathbf{y}}^2 g_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) \right]^{-1} \nabla_{\mathbf{y}} f_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x}))}_{\mathbf{p}_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x}))}$$

Our approaches:

- Use **FedSVRG** to solve the inner problem.
- Estimate $\mathbf{p}(\mathbf{x}, \mathbf{y}^*(\mathbf{x}))$ via a Federated Inverse Hessian-Gradient-Product (**FedIHGP**).

- N -Neumann series approximation (Ghadimi & Wang, 2018):

$$\begin{aligned} \mathbf{p}(\mathbf{x}, \mathbf{y}) &:= \sum_{i=1}^m \left[\sum_{i=1}^m \nabla_{\mathbf{y}}^2 g_i(\mathbf{x}, \mathbf{y}) \right]^{-1} \nabla_{\mathbf{y}} f_i(\mathbf{x}, \mathbf{y}) \\ &\approx \sum_{i=1}^m \left[\frac{N}{\ell_{g,1}} \prod_{n=1}^{N'} \sum_{i=1}^m \left(I - \frac{1}{\ell_{g,1}} \nabla_{\mathbf{y}}^2 g_i(\mathbf{x}, \mathbf{y}; \zeta_n) \right) \right] \nabla_{\mathbf{y}} f_i(\mathbf{x}, \mathbf{y}; \xi) \end{aligned}$$

- FedIHGP** provides a **federated recursive strategy** to estimate \mathbf{p} .

$\mathbf{p}_{N'}$ = FedIHGP ($\mathbf{x}, \mathbf{y}, N$)

Select $N' \in \{0, \dots, N-1\}$, $S_0, \dots, S_{N'} \in \mathcal{S}$ UAR. Set

- i -th client: $\mathbf{p}_{i,0} = \nabla_{\mathbf{y}} f_i(\mathbf{x}, \mathbf{y}; \xi_{i,0})$
- server: $\mathbf{p}_0 = \frac{N}{\ell_{g,1}} |S_0|^{-1} \sum_{i \in S_0} \mathbf{p}_{i,0}$

If $N' = 0$ Return $\mathbf{p}_{N'}$.

For $n = 1, \dots, N'$:

- i -th client: $\mathbf{p}_{i,n} = \left(I - \frac{1}{\ell_{g,1}} \nabla_{\mathbf{y}}^2 g_i(\mathbf{x}, \mathbf{y}; \zeta_{i,n}) \right) \mathbf{p}_{i,n-1}$
- server: $\mathbf{p}_n = |S_n|^{-1} \sum_{i \in S_n} \mathbf{p}_{i,n}$

- FedIHGP** avoids explicit Hessian:

- matrix-vector products
- vector communications

- $\| \mathbf{p}(\mathbf{x}, \mathbf{y}) - \mathbb{E}[\mathbf{p}_{N'}] \| \leq \mathcal{O} \left(\underbrace{\left(\frac{\kappa_g - 1}{\kappa_g} \right)^N}_{< 1} \right)$.

- $\kappa_g := \frac{\ell_{g,1}}{\mu_g}$ (condition number).

Proposed Algorithm: FedNest

For $k = 0, \dots, K - 1$

- ① $\mathbf{y}^{k+1} = \mathbf{FedInn}(\mathbf{x}^k, \mathbf{y}^k, \beta^k)$ // one or multiple FedSVRGs on \mathbf{y}
- ② $\mathbf{x}^{k+1} = \mathbf{FedOut}(\mathbf{x}^k, \mathbf{y}^{k+1}, \alpha^k)$ // FedSVRG + FedIHGP on \mathbf{x}

Proposed Algorithm: FedNest

For $k = 0, \dots, K - 1$

- ① $\mathbf{y}^{k+1} = \mathbf{FedInn}(\mathbf{x}^k, \mathbf{y}^k, \beta^k)$ // one or multiple FedSVRGs on \mathbf{y}
- ② $\mathbf{x}^{k+1} = \mathbf{FedOut}(\mathbf{x}^k, \mathbf{y}^{k+1}, \alpha^k)$ // FedSVRG + FedIHGP on \mathbf{x}

Inner Optimizer: FedInn

- $\mathbf{y}_k \approx \mathbf{y}^*(\mathbf{x}_k)$
- It avoids inner client drift:
 $\|\mathbf{y}_{i,\nu}^k - \mathbf{y}^k\|^2 \leq O(\tau_i(\beta_i^k)^2)$
- The global convergence of **FedInn** ensures **accurate hypergradient computation**.

Proposed Algorithm: FedNest

For $k = 0, \dots, K - 1$

- ① $\mathbf{y}^{k+1} = \mathbf{FedInn}(\mathbf{x}^k, \mathbf{y}^k, \beta^k)$ // one or multiple FedSVRGs on \mathbf{y}
- ② $\mathbf{x}^{k+1} = \mathbf{FedOut}(\mathbf{x}^k, \mathbf{y}^{k+1}, \alpha^k)$ // FedSVRG + FedIHGP on \mathbf{x}

Inner Optimizer: FedInn

- $\mathbf{y}_k \approx \mathbf{y}^*(\mathbf{x}_k)$
- It avoids inner client drift:
 $\|\mathbf{y}_{i,\nu}^k - \mathbf{y}^k\|^2 \leq O(\tau_i(\beta_i^k)^2)$
- The global convergence of **FedInn** ensures **accurate hypergradient computation**.

Outer Optimizer: FedOut

- It avoids outer client drift:

$$\|\mathbf{x}_{i,\nu}^k - \mathbf{x}^k\|^2 \leq O(\tau_i(\alpha_i^k)^2) + \|\mathbf{y}^{k+1} - \mathbf{y}^*(\mathbf{x}^k)\|^2$$

- It gives new convergence guarantees for federated bilevel, minimax, compositional, and single-level optimization.

Assumptions

- $f_i(\mathbf{z}), \nabla f_i(\mathbf{z}), \nabla g_i(\mathbf{z}), \nabla^2 g_i(\mathbf{z})$ are $\ell_{f,0}, \ell_{f,1}, \ell_{g,1}, \ell_{g,2}$ -Lipschitz continuous, respectively.
 - $g_i(\mathbf{x}, \mathbf{y})$ is μ_g -strongly convex in \mathbf{y} for all $\mathbf{x} \in \mathbb{R}^{d_1}$.
 - $\nabla f_i(\mathbf{z}; \xi), \nabla g_i(\mathbf{z}; \zeta), \nabla^2 g_i(\mathbf{z}; \zeta)$ are unbiased estimators of $\nabla f_i(\mathbf{z}), \nabla g_i(\mathbf{z}), \nabla^2 g_i(\mathbf{z})$; and their variances are bounded.
- These assumptions are common in the (non-federated) BO literature.

Assumptions

- $f_i(\mathbf{z}), \nabla f_i(\mathbf{z}), \nabla g_i(\mathbf{z}), \nabla^2 g_i(\mathbf{z})$ are $\ell_{f,0}, \ell_{f,1}, \ell_{g,1}, \ell_{g,2}$ -Lipschitz continuous, respectively.
 - $g_i(\mathbf{x}, \mathbf{y})$ is μ_g -strongly convex in \mathbf{y} for all $\mathbf{x} \in \mathbb{R}^{d_1}$.
 - $\nabla f_i(\mathbf{z}; \xi), \nabla g_i(\mathbf{z}; \zeta), \nabla^2 g_i(\mathbf{z}; \zeta)$ are unbiased estimators of $\nabla f_i(\mathbf{z}), \nabla g_i(\mathbf{z}), \nabla^2 g_i(\mathbf{z})$; and their variances are bounded.
- These assumptions are common in the (non-federated) BO literature.

Theorem (Informal)

Under the above assumptions, if we choose the stepsize properly, then the iterates of FedNest satisfy

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} \left[\left\| \nabla f(\mathbf{x}^k) \right\|^2 \right] = \mathcal{O}\left(\frac{1}{\sqrt{K}}\right) \quad \text{and} \quad \mathbb{E} \left[\left\| \mathbf{y}^k - \mathbf{y}^*(\mathbf{x}^k) \right\|^2 \right] = \mathcal{O}\left(\frac{1}{\sqrt{K}}\right).$$

Theory: Comparison with Previous Results

- Sample complexity of FedNest and comparable non-FL methods to find an ϵ -stationary point of f , i.e., $1/K \sum_{k=1}^K \mathbb{E}[\|\nabla f(\mathbf{x}^k)\|^2] \leq \epsilon$:
- $\kappa_g = \ell_{g,1}/\mu_g$ (condition number).

		Non-Federated		
	FedNest	ALSET	BSA	TTSA
batch size	$\mathcal{O}(1)$			
samples in ξ	$\mathcal{O}(\kappa_g^5 \epsilon^{-2})$	$\mathcal{O}(\kappa_g^5 \epsilon^{-2})$	$\mathcal{O}(\kappa_g^6 \epsilon^{-2})$	$\mathcal{O}(\kappa_g^p \epsilon^{-2.5})$
samples in ζ	$\mathcal{O}(\kappa_g^9 \epsilon^{-2})$	$\mathcal{O}(\kappa_g^9 \epsilon^{-2})$	$\mathcal{O}(\kappa_g^9 \epsilon^{-3})$	$\mathcal{O}(\kappa_g^p \epsilon^{-2.5})$

ALSET(Chen et al., 2021), **BSA**(Ghadimi & Wang, 2018), **TTSA**(Hong et al., 2020).

- Main takeaways:
 - **FedNest** enjoys the same convergence as non-federated alternating SGD (**ALSET**), despite objective heterogeneity.

LFedNest: Communication Efficiency via Local Hypergradient

Light-FedNest (LFedNest):

- computes hypergradients locally
- only needs a single communication round for the outer update

	definition		properties			
	outer optimizer	inner optimizer	global outer gradient	global IHGP	global inner gradient	# communication rounds
FedNest	SVRG on x	SVRG on y	yes	yes	yes	$2T + N + 3$
LFedNest	SGD on x	SGD on y	no	no	no	$T + 1$
FedNest_{SGD}	SVRG on x	SGD on y	yes	yes	no	$T + N + 3$

- T : # inner iterations (y update)
- N : # terms of Neumann series

Minimax Experiment

- Minimax saddle point problem (on non-i.i.d. synthetic dataset):

$$\min_{\mathbf{x} \in \mathbb{R}^{d_1}} f(\mathbf{x}) := \frac{1}{m} \max_{\mathbf{y} \in \mathbb{R}^{d_2}} \sum_{i=1}^m f_i(\mathbf{x}, \mathbf{y}),$$

where

$$f_i(\mathbf{x}, \mathbf{y}) := - \left[\frac{1}{2} \|\mathbf{y}\|^2 - \mathbf{b}_i^\top \mathbf{y} + \mathbf{y}^\top \mathbf{A}_i \mathbf{x} \right] + \frac{\lambda}{2} \|\mathbf{x}\|^2,$$

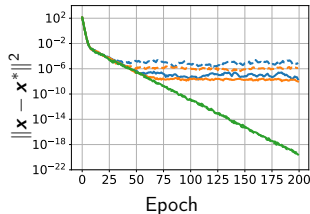
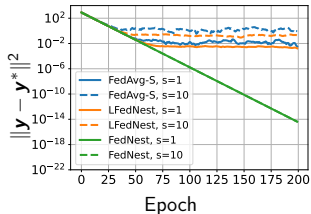
Minimax Experiment

- Minimax saddle point problem (on non-i.i.d. synthetic dataset):

$$\min_{\mathbf{x} \in \mathbb{R}^{d_1}} f(\mathbf{x}) := \frac{1}{m} \max_{\mathbf{y} \in \mathbb{R}^{d_2}} \sum_{i=1}^m f_i(\mathbf{x}, \mathbf{y}),$$

where

$$f_i(\mathbf{x}, \mathbf{y}) := - \left[\frac{1}{2} \|\mathbf{y}\|^2 - \mathbf{b}_i^\top \mathbf{y} + \mathbf{y}^\top \mathbf{A}_i \mathbf{x} \right] + \frac{\lambda}{2} \|\mathbf{x}\|^2,$$



- LFedNest performs slightly better than FedAvg-S (Hou et al., 2021).
- FedNest converges linearly despite heterogeneity.

Hyperparameter Tuning for Label Imbalance

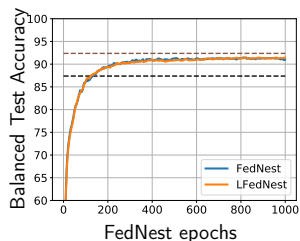
- Imbalanced classification is the problem of classification when there is an unequal distribution of classes in the training dataset.
- **Goal:** Design fairness-seeking loss functions via bilevel optimization

$$\begin{aligned} & \underset{\mathbf{x}=(\Delta, l)}{\text{minimize}} \quad \sum_{i=1}^m \overbrace{\sum_{(\mathbf{u}, t) \in \mathcal{D}_{val}^i} \frac{1}{p_t} \log(1 + \sum_{s \neq t} e^{y_s(\mathbf{u}) - y_t(\mathbf{u})})}^{=: f_i(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) \leftarrow \text{Balanced Val Loss}} \\ \text{s.t.} \quad & \mathbf{y}^*(\mathbf{x}) = \underset{\mathbf{y}}{\text{arg min}} \underbrace{\sum_{i=1}^m \sum_{(\mathbf{u}, t) \in \mathcal{D}_{train}^i} \log(1 + \sum_{s \neq t} e^{l_s - l_t} \cdot e^{\Delta_s y_s(\mathbf{u}) - \Delta_t y_t(\mathbf{u})})}_{=: g_i(\mathbf{x}, \mathbf{y}) \leftarrow \text{Parametric Train Loss}} \end{aligned}$$

- **Outer optimization** aims to maximize the class-balanced validation accuracy.
- **Inner optimization** trains model parameter \mathbf{y} to minimize $g = \sum_{i=1}^m g_i$.

Hyperparameter Tuning for Label Imbalance

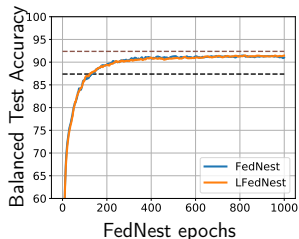
- **Brown dashed line:** Non-Federated bilevel training
- **Black dashed line:** Non-Federated accuracy without bilevel tuning



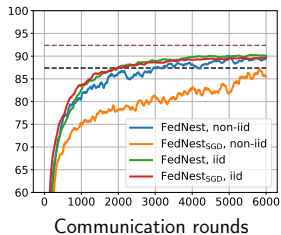
- **I.I.D. setup:** FedNest behaves similarly to LFedNest.

Hyperparameter Tuning for Label Imbalance

- **Brown dashed line:** Non-Federated bilevel training
- **Black dashed line:** Non-Federated accuracy without bilevel tuning



- **I.I.D. setup:** FedNest behaves similarly to LFedNest.



- FedNest (with SVRG) is significantly better than FedNest_{SGD} (with SGD in FedInn).

Federated (**Single-Level**) Optimization

FedAvg (McMahan et al., 2017)

FedSVRG (Konečný et al., 2018)

FedProx (Li et al., 2020)

SCAFFOLD (Karimireddy et al., 2020)

FedNova (Wang et al., 2020)

FedLin (Mitra et al., 2021)

Stochastic **Bilevel** Optimization

BSA (Ghadimi & Wang, 2018)

TTSA (Hong et al., 2020)

ALSET (Chen et al., 2021)

stocBiO (Ji et al., 2020)

Stochastic **MiniMax** Optimization

SGDA (Lin et al., 2020)

SMD (Rafique et al., 2021)

Stochastic **Compositional** Optimization

SCGD (Wang et al., 2017)

NASA (Ghadimi et al., 2020)

Conclusion and Future Work

- **Conclusion:**

- **FedNest** gives a new framework for federated bilevel, minimax, and compositional optimization.
- **FedNest** matches the sample complexity of the alternating SGD.

	FedNest			
	Bilevel	Minimax	Compositional	Single-Level
batch size	$\mathcal{O}(1)$			
samples complexity	$\mathcal{O}(\epsilon^{-2})$			

- **Future Work:**

- Other applications and properties of FedNest such as federated actor-critic reinforcement learning.
- Sparsification or quantization for communication-efficient FBO

Paper: [ICML 2022](#), [arXiv](#)

Code: github.com/ucr-optml/FedNest

References

- Chen, T., Sun, Y., and Yin, W. Closing the gap: Tighter analysis of alternating stochastic gradient methods for bilevel problems. *Advances in Neural Information Processing Systems*, 34, 2021.
- Ghadimi, S. and Wang, M. Approximation methods for bilevel programming. *arXiv preprint arXiv:1802.02246*, 2018.
- Ghadimi, S., Ruszczyński, A., and Wang, M. A single timescale stochastic approximation method for nested stochastic optimization. *SIAM Journal on Optimization*, 30(1):960–979, 2020.
- Hong, M., Wai, H.-T., Wang, Z., and Yang, Z. A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic. *arXiv preprint arXiv:2007.05170*, 2020.
- Hou, C., Thekumparampil, K. K., Fantì, G., and Oh, S. Efficient algorithms for federated saddle point optimization. *arXiv preprint arXiv:2102.06333*, 2021.
- Ji, K., Yang, J., and Liang, Y. Provably faster algorithms for bilevel optimization and applications to meta-learning. *ArXiv*, abs/2010.07962, 2020.
- Karimireddy, S. P., Kale, S., Mohri, M., Reddi, S., Stich, S., and Suresh, A. T. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pp. 5132–5143. PMLR, 2020.
- Konečný, J., McMahan, H. B., Ramage, D., and Richtárik, P. Federated optimization: Distributed machine learning for on-device intelligence. *International Conference on Learning Representations*, 2018.
- Li, T., Sahu, A. K., Zaheer, M., Sanjabi, M., Talwalkar, A., and Smith, V. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems*, 2:429–450, 2020.
- Lin, T., Jin, C., and Jordan, M. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pp. 6083–6093. PMLR, 2020.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, AISTATS 2017, 20-22 April 2017, Fort Lauderdale, FL, USA*, pp. 1273–1282, 2017. URL <http://proceedings.mlr.press/v54/mcmahan17a.html>.
- Mitra, A., Jaafar, R., Pappas, G., and Hassani, H. Linear convergence in federated learning: Tackling client heterogeneity and sparse gradients. *Advances in Neural Information Processing Systems*, 34, 2021.
- Rafique, H., Liu, M., Lin, Q., and Yang, T. Weakly-convex-concave min-max optimization: provable algorithms and applications in machine learning. *Optimization Methods and Software*, pp. 1–35, 2021.
- Wang, J., Liu, Q., Liang, H., Joshi, G., and Poor, H. V. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in neural information processing systems*, 2020.
- Wang, M., Fang, E. X., and Liu, H. Stochastic compositional gradient descent: algorithms for minimizing compositions of expected-value functions. *Mathematical Programming*, 161(1-2):419–449, 2017. 