

BACKGROUND

Conditional stochastic optimization has applications in a wide range of machine learning tasks, such as invariant learning, AUPRC maximization, and meta-learning.

- Model-Agnostic Meta-Learning (MAML)

In meta-learning, we attempt to train models that can efficiently adapt to unseen tasks via learning with meta data from similar tasks. When the tasks are distributed at different clients, a federated version of MAML would be beneficial to leverage information from all workers to improve the model performance on the downstream tasks.

$$\min_x \mathbb{E}_{i \sim \mathcal{P}_{\text{task}}, a \sim D_{\text{query}}^i} \mathcal{L}_i \left(\mathbb{E}_{b \sim D_{\text{support}}^i} (x - \lambda \nabla \mathcal{L}_i(x, b)), a \right)$$

- Online AUPRC Maximization In AUPRC maximization, AP loss is used as the loss function, instead of cross-entropy loss since AP is the surrogate function of AUPRC and cross-entropy is corresponding to accuracy. By directly optimizing AP loss, model performance is improved with the metric of AUPRC.

$$\hat{\text{AP}} = \mathbb{E}_{\xi \sim \mathcal{D}^+} \frac{\mathbb{E}_{\xi \sim \mathcal{D}} \mathbf{I}(y=1) \ell(x; z^+, z)}{\mathbb{E}_{\xi \sim \mathcal{D}} \ell(x; z^+, z)}$$

CHALLENGES

To address the large-scale distributed data challenges across multiple clients with communication-efficient distributed training, federated learning (FL) is gaining popularity. Many optimization algorithms for CSO problems have been developed in the centralized setting. Nonetheless, the algorithm for CSO problems under FL is still underexplored.

$$\min_{x \in \mathcal{X}} F(x) := \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\xi^n} f_{\xi^n}^n(\mathbb{E}_{\eta^n | \xi^n} g_{\eta^n}^n(x, \xi^n)), \quad (1)$$

Federated conditional stochastic optimization subsumes the standard federated learning optimization as a special case when the inner-layer function $g_{\eta^n}^n(x, \xi^n) = x$. In addition, federated stochastic compositional optimization is similar to federated conditional stochastic optimization given that both problems contain two-layer nested expectations. However, they are fundamentally different. In federated stochastic compositional optimization, the inner randomness η and the outer randomness ξ are independent and data samples of the inner layer are available directly from η (instead of a conditional distribution as in Problem (1)). Therefore, when randomnesses η and ξ are independent and $g_{\eta^n}^n(x, \cdot) = g_{\eta^n}^n(x)$, (1) is converted into federated stochastic compositional optimization.

EXISTING METHODS

Summary of complexity results of proposed federated conditional stochastic optimization algorithms to reach an ϵ -stationary point. Sample complexity is defined as the number of calls to the First-order Oracle (FO) by clients to reach an ϵ -stationary point. Communication complexity denotes the total number of back-and-forth communication rounds between each client and the central server required to reach an ϵ -stationary point.

Algorithm	Sample	Communication
FCSG	$O(\epsilon^{-6})$	$O(\epsilon^{-3})$
FCSG-M	$O(\epsilon^{-6})$	$O(\epsilon^{-3})$
Theoretical Lower Bound	$O(\epsilon^{-5})$	-
Acc-FCSG-M	$O(\epsilon^{-5})$	$O(\epsilon^{-2})$

FEDSGDA+ ALGORITHM

Algorithm 1 FCSG and FCSG-M Algorithm

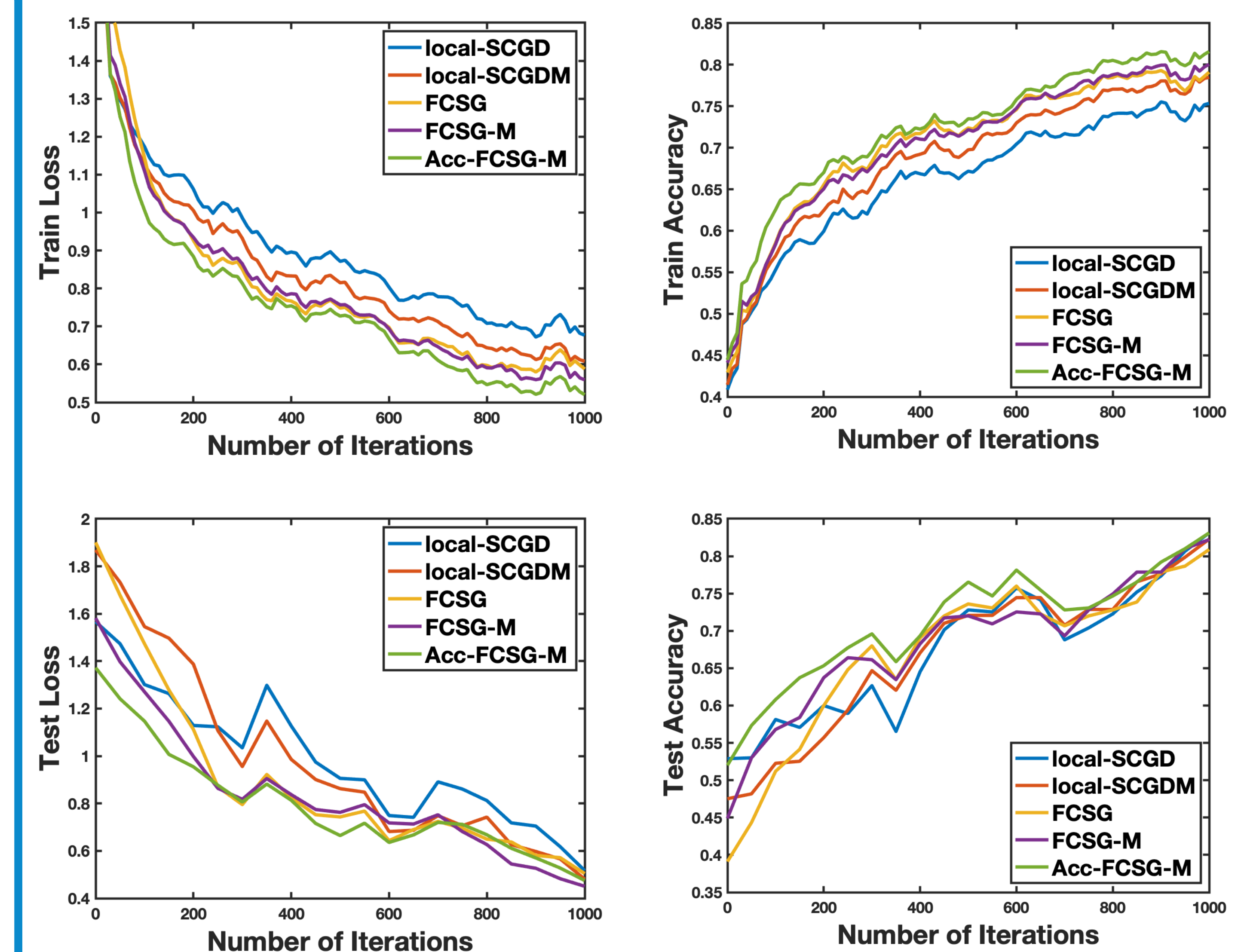
- 1: **Input:** Parameters: T , momentum weight β , learning rate α , the number of local updates q , inner batch size m and outer batch size b , as well as the initial outer batch size B ;
- 2: **Initialize:** $x_0^n = \bar{x}_0 = \frac{1}{N} \sum_{k=1}^N x_0^k$. Draw B samples of $\{\xi_{t,1}^n, \dots, \xi_{t,B}^n\}$ and draw m samples $\mathcal{B}_{0,i}^n = \{\eta_{ij}^n\}_{j=1}^m$ from $P(\eta^n | \xi_{0,i}^n)$ for each $\xi_{0,i}^n \in \{\xi_{t,1}^n, \dots, \xi_{t,B}^n\}$; $u_1^n = \frac{1}{B} \sum_{i=1}^B \nabla \hat{F}^n(x_0^n; \xi_{0,i}^n, \mathcal{B}_{0,i}^n)$.
- 3: **for** $t = 1, 2, \dots, T$ **do**
- 4: **for** $n = 1, 2, \dots, N$ **do**
- 5: **if** $\text{mod}(t, q) = 0$ **then**
- 6: **Server Update:**
- 7: $u_t^n = \bar{u}_t = \frac{1}{N} \sum_{i=1}^N u_t^i$
- 8: $x_t^n = \bar{x}_t = \frac{1}{N} \sum_{n=1}^N (x_{t-1}^n - \alpha u_t^n)$
- 9: **else**
- 10: $x_t^n = x_{t-1}^n - \alpha u_t^n$
- 11: **end if**
- 12: Draw b samples of $\{\xi_{t,1}^n, \dots, \xi_{t,b}^n\}$
- 13: Draw m samples $\mathcal{B}_{t,n}^n = \{\eta_{ij}^n\}_{j=1}^m$ from $P(\eta^n | \xi_{t,i}^n)$ for each $\xi_{t,i}^n \in \{\xi_{t,1}^n, \dots, \xi_{t,b}^n\}$,
- 14: $u_{t+1}^n = \frac{1}{b} \sum_{i=1}^b \nabla \hat{F}^n(x_t^n; \xi_{t,i}^n, \mathcal{B}_{t,i}^n)$
- 15: $u_{t+1}^n = (1 - \beta)u_t^n + \frac{\beta}{b} \sum_{i=1}^b \nabla \hat{F}^n(x_t^n; \xi_{t,i}^n, \mathcal{B}_{t,i}^n)$
- 16: **end for**
- 17: **end for**
- 18: **Output:** x chosen uniformly random from $\{\bar{x}_t\}_{t=1}^T$.

FEDSGDA-M ALGORITHM

Algorithm 2 Acc-FCSG-M Algorithm

- 1: **Initialize:** Draw B samples of $\{\xi_1^n, \dots, \xi_B^n\}$ and draw m samples $\mathcal{B}_{0,i}^n = \{\eta_{ij}^n\}_{j=1}^m$ from $P(\eta^n | \xi_i^n)$ for each $\xi_i^n \in \{\xi_1^n, \dots, \xi_B^n\}$, then $u_1^n = \frac{1}{B} \sum_{i=1}^B \nabla \hat{F}^n(x_0^n; \xi_{0,i}^n, \mathcal{B}_{0,i}^n)$ for $n \in [N]$.
- 2: **for** $t = 1, 2, \dots, T$ **do**
- 3: **for** $n = 1, 2, \dots, N$ **do**
- 4: **if** $\text{mod}(t, q) = 0$ **then**
- 5: $u_t^n = \bar{u}_t = \frac{1}{N} \sum_{i=1}^N u_t^i$, $x_t^n = \bar{x}_t = \frac{1}{N} \sum_{n=1}^N (x_{t-1}^n - \alpha u_t^n)$
- 6: **else**
- 7: $x_{t,n} = x_{t-1}^n - \alpha u_t^n$
- 8: **end if**
- 9: Draw b samples of $\{\xi_{t,1}^n, \dots, \xi_{t,b}^n\}$, draw m samples $\mathcal{B}_{t,n}^n = \{\eta_{ij}^n\}_{j=1}^m$ from $P(\eta^n | \xi_{t,i}^n)$ for each $\xi_{t,i}^n \in \{\xi_{t,1}^n, \dots, \xi_{t,b}^n\}$,
- 10: $u_{t+1}^n = \frac{1}{b} \sum_{i=1}^b \nabla \hat{F}^n(x_t^n; \xi_{t,i}^n, \mathcal{B}_{t,i}^n) + (1 - \beta)(u_t^n - \frac{1}{b} \sum_{i=1}^b \nabla \hat{F}^n(x_{t-1}^n; \xi_{t,i}^n, \mathcal{B}_{t,i}^n))$
- 11: **end for**
- 12: **end for**

EXPERIMENTS



The 5-way-1-shot case over Omniglot Dataset in MAML (Train loss; Train accuracy; Test loss; Test accuracy)