# Federated Multi-Objective Learning

Haibo Yang[1], Zhuqing Liu[2], Jia Liu[2], Chaosheng Dong[3], Michinari Momma[3]

[1]Rochester Institute of Technology, [2]The Ohio State University, [3]Amazon

## Overview

In this work, we propose a new **federated multi-objective learning (FMOL)** framework with multiple clients distributively and collaboratively solving an MOO problem while keeping their training data private.

❖ Our FMOL framework allows a different set of objective functions across different clients to support a wide range of applications, which advances and generalizes the MOO formulation to the federated learning paradigm.

❖ For this FMOL framework, we propose two new federated multi-objective optimization (FMOO) algorithms called federated multi-gradient descent averaging (FMGDA) and federated stochastic multi-gradient descent averaging (FSMGDA).

❖ Both algorithms allow local updates to significantly reduce communication costs, while achieving the same convergence rates as those of their algorithmic counterparts in centralized learning.

## Problem Formulation

For a system with M clients and S tasks (objectives) in total, our FMOL framework can be written as follows:

$$\min_{\mathbf{x}} \quad \mathrm{Diag}(\mathbf{FA}^\top),$$

$$\mathbf{F} \triangleq \begin{bmatrix} f_{1,1} & \cdots & f_{1,M} \\ \vdots & \ddots & \vdots \\ f_{S,1} & \cdots & f_{S,M} \end{bmatrix}_{S \times M}, \mathbf{A} \triangleq \begin{bmatrix} a_{1,1} & \cdots & a_{1,M} \\ \vdots & \ddots & \vdots \\ a_{S,1} & \cdots & a_{S,M} \end{bmatrix}_{S \times M},$$

where $\mathbf{F}$ groups all potential objectives $f_{s,i}(x)$ for each task s at each client i, and $A \in \{0,1\}^{S \times M}$ is a binary objective indicator matrix, with each element $a_{s,i} = 1$ if task s is of client i's interest and $a_{s,i} = 0$ otherwise.

1. Each client has only one distinct objective: $A = I_M, S = M, Diag(FA^T) = [f_1(x), f_2(x), \ldots, f_S(x)]$. E.g., multi-task learning, classic federated learning as MOO problem.
2. All clients share the same S objectives: A is an all-one matrix. E.g., distributed MOO with decentralized data.
3. Each client has a different subset of objectives.

## Algorithm



- Local update: communication efficient
- Two-sided learning rates: local learning rate controls the derivations and noise, global learning rate manage the learning progress
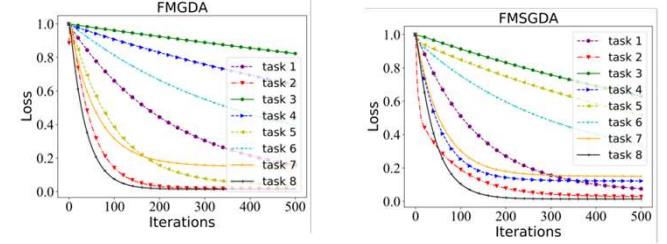- Convex quadratic optimization for common direction

## Convergence Rates

Table 1: Convergence rate results (shaded parts are our results) comparisons.

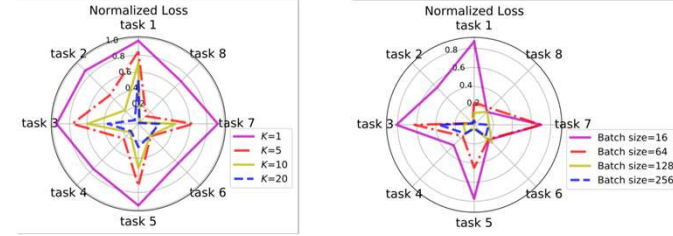| Methods | Strongly Convex | | Non-convex | |
|---|---|---|---|---|
| | Rate | Assumption* | Rate | Assumption* |
| MGD [7] | $\mathcal{O}(r^T)$ # | Linear search & sequence convergence | $\mathcal{O}(1/T)$ | Linear search & sequence convergence |
| SMGD [8] | $\mathcal{O}(1/T)$ | First moment bound & Lipschitz continuity of $\lambda$ | Not provided | Not provided |
| FMGDA | $\mathcal{O}(\exp(-\mu T))$ # | Not needed | $\mathcal{O}(1/T)$ | Not needed |
| FSMGDA | $\tilde{\mathcal{O}}(1/T)$ | $(\alpha,\beta)$-Lipschitz continuous stochastic gradient | $\mathcal{O}(1/\sqrt{T})$ | $(\alpha,\beta)$-Lipschitz continuous stochastic gradient |

# Notes on constants: $\mu$ is the strong convexity modulus; r is a constant depends on $\mu$, s.t., $r \in (0,1)$.
* Assumption short-hands: "Linear search": learning rate linear search [7]; "Sequence convergence": $\{x_t\}$ converges to $x^*$ [7]; "First moment bound" (Asm. 5.2(b) [8]): $\mathbb{E}[\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{x})\|] \le \eta(a + b\|\nabla f(\mathbf{x})\|)$;"Lipschitz continuity of $\lambda$" (Asm. 5.4 [8]): $\|\boldsymbol{\lambda}_k - \boldsymbol{\lambda}_t\| \le \beta \|[(\nabla f_1(\mathbf{x}_k) - \nabla f_1(\mathbf{x}_t))^T, \ldots, (\nabla f_S(\mathbf{x}_k) - \nabla f_S(\mathbf{x}_t))^T]\|$; "$(\alpha,\beta)$-Lipschitz continuous stochastic gradient": see Asm. 4.

**Assumption 4** $((\alpha,\beta)$-Lipschitz Continuous Stochastic Gradient). A function f has $(\alpha,\beta)$-Lipschitz continuous stochastic gradients if there exist two constants $\alpha, \beta > 0$ such that, for any two independent training samples $\xi$ and $\xi'$, $\mathbb{E}[\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{y}, \xi')\|^2] \le \alpha\|\mathbf{x} - \mathbf{y}\|^2 + \beta\sigma^2$.

## Numerical Results



Effectiveness of FMOL algorithms: River Flow dataset with 8 tasks.



Ablation study: local steps and batch size

## Conclusion

- Proposed a general Federated Multi-Objective Learning (FMOL) framework..
- Proposed two federated (stochastic) multi-gradient descent averaging algorithms with theoretical guarantees.

## Acknowledgments