Federated Composite Optimization

Overview

New Setting: Federated Composite Optimization (FCO)

- FL with (possibly non-smooth) regularizers or constraints
- Arises naturally in FL applications that involve sparsity, low-rank, or constraints.
- Standard FL algorithms (e.g., FedAvg) are for smooth unconstrained settings.

Straightforward Extension of FedAvg suffers from "curse of primal averaging"

Our proposal: Federated Dual Averaging

- Novel server dual averaging procedure
- Theoretical and empirical advantages

Problem definition and Examples

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

- $\circ F_m(w) := \mathbb{E}_{\xi \sim \mathcal{D}_m}[f(w;\xi)]$ is the loss function of the m-th client
- ψ is a (possibly non-smooth, non-finite) convex regularizer Ο
- **Federated Lasso** for sparsity representations

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

• Federated matrix completion via nuclear norm

$$\min_{W} \frac{1}{M} \sum_{m=1}^{M} F_m(W) + \lambda \|W\|_*$$



 $+\infty$ if $w \notin \mathcal{C}$.

Background of (non-federated) composite optimization

Composite 101: **ProxGD** is the standard algorithm for solving non-federated CO:

$$w_{t+1} \leftarrow \operatorname{prox}_{\eta\psi} (w_t - \eta \nabla F(w_t)) \qquad \qquad \text{Proximal additive}$$

$$:= \underset{w}{\operatorname{argmin}} \left\{ F(w_t) + \langle \nabla F(w_t), w - w_t \rangle + \frac{1}{2\eta} \|w - w_t\|_2^2 + \frac{\psi(w)}{2} \right\}$$
First-order Taylor expansion of F Smoothness estimation

Composite 201: Mirror descent generalizes ProxGD to general Bregman divergence

$$w_{t+1} = \underset{w}{\operatorname{argmin}} \left\{ F\left(w_{t}\right) + \left\langle \nabla F\left(w_{t}\right), w - w_{t} \right\rangle + \psi(w) + \frac{1}{\eta} D_{h}(w, w_{t}) \right\} \underset{if \ h(w) = \frac{1}{2} ||w||^{2}}{\operatorname{reduces to PGD}}$$

Primal-Dual interpretation of Mirror Descent [Nemirovski and Yudin, 1983, Flammarion and Bach, 2017]

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





FedMiD: a straightforward extension

Algorithm 1 Federated Averaging (FEDAVG)			Algorithm 2 Federated Mirror Descent (FEDMID)	
1: procedure FEDAVG (w_0, η_c, η_s)		_	1: procedure FEDMID (w_0, η_c, η_s)	
2: for $r = 0,, R - 1$ do			2: for $r = 0,, R - 1$ do	
3:	sample a subset of clients $\mathcal{S}_r \subseteq [M]$		3:	sample a subset of clients $\mathcal{S}_r \subseteq [M]$
4:	on client $m \in S_r$ in parallel do		4:	on client $m \in \mathcal{S}_r$ in parallel do
5:	client initialization $w_{r,0}^m \leftarrow w_r$		5:	client initialization $w_{r,0}^m \leftarrow w_r$
6:	for $k = 0, \ldots, K - 1$ do		6:	for $k = 0, \ldots, K - 1$ do
7:	$g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$		7:	$g_{r,k}^m \leftarrow \nabla f(w_{r,k}^m; \xi_{r,k}^m)$
8:	$w_{r,k+1}^m \leftarrow w_{r,k}^m - \eta_{\mathrm{c}} \cdot g_{r,k}^m$	\rightarrow	8:	$w_{r,k+1}^m \leftarrow \nabla (h + \eta_{\mathrm{c}} \psi)^* (\nabla h(w_{r,k}^m) - \eta_{\mathrm{c}} g_{r,k}^m)$
9:	$\Delta_r = rac{1}{ \mathcal{S}_r } \sum_{m \in \mathcal{S}_r} (w^m_{r,K} - w^m_{r,0})$		9:	$\Delta_r = rac{1}{ \mathcal{S}_r } \sum_{m \in \mathcal{S}_r} (w^m_{r,K} - w^m_{r,0})$
10:	$w_{r+1} \leftarrow w_r + \eta_{\rm s} \cdot \Delta_r$	\rightarrow	10:	$w_{r+1} \leftarrow \nabla (h + \eta_{\mathrm{s}} \eta_{\mathrm{c}} K \psi)^* (\nabla h(w_r) + \eta_{\mathrm{s}} \Delta_r)$

Issue of FedMiD: curse of primal averaging



- Persistent **primal** states
- Persistent dual states

Our main proposal: FedDualAvg

Algorithm 3 Federated Dual Averaging						
1: procedure FEDDUALAVG (w_0, η_c, η_s)						
2: S	2: server initialization $z_0 \leftarrow \nabla h(w_0)$					
3: for $r = 0,, R - 1$ 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 $k = 0, \ldots, K - 1$ do					
8:	$ ilde{\eta}_{r,k} \leftarrow \eta_{ m s} \eta_{ m c} r K + \eta_{ m 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:	$z_{r,k+1}^{m} \leftarrow z_{r,k}^{m} - \eta_{ ext{c}} g_{r,k}^{m}$					
12:	$\Delta_r = rac{1}{ \mathcal{S}_r } \sum_{m \in \mathcal{S}_r} (z^m_{r,K} - z^m_{r,0})$					
13:	$z_{r+1} \leftarrow z_r + \eta_{\mathrm{s}} \Delta_r$					
14:	$w_{r+1} \leftarrow \nabla (h + \eta_{\mathrm{s}} \eta_{\mathrm{c}} (r+1) K \psi)^* (z_{r+1})$					

Locally: each client runs dual averaging, tracking a pair of primal and dual states.

Communication: dual states are aggregated across clients.

- --- Compute primal point
- → Average client **dual** deltas
- → Server **dual** update
- → (Optional) primal output

Stanford Google Research



Main theoretical results



(e) Assume all the M clients participate in client updates for every round, namely $S_r = [M]$.

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

Theorem 4.3. Assuming A1, and in addition assume $\sup_{w \in \operatorname{dom}\psi} \|\nabla F_m(w) - \nabla F(w)\|_* \leq \zeta^2$ and *F* is quadratic, then FedDualAvg can output \hat{w} such that

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

aster convergence moreover for appropriate η_{i} (usefulness of client step) Overhead for infrequent communication

 $B \coloneqq D_h(w^*, w_0)$ L: smoothness matches best known bound on σ : variance bound smooth unconstrained FedAvg *M: # of clients* K: # of local steps c.f. [Woodworth et al., 2020] *R: # of rounds*

Experiments



Figure 1: Results on sparse (ℓ_1 -regularized) logistic regression for a federated fMRI dataset based on (Haxby, 2001). centralized corresponds to training on the centralized dataset gathered from all the training clients. local corresponds to training on the local data from only one training client without communication. FEDAVG (∂) corresponds to running FEDAVG algorithms with subgradient in lieu of SGD to handle the non-smooth ℓ_1 -regularizer. FEDMID is another straightforward extension of FEDAVG running local proximal gradient method (see Section 3.1 for details). We show that using our proposed algorithm FEDDUALAVG, one can 1) achieve performance comparable to the centralized baseline without the need to gather client data, and 2) significantly outperforms the local baseline on the isolated data and the FEDAVG baseline. See Section 5.3 for details.

Main References

- A.S. Nemirovski and D.B. Yudin "Problem complexity and method efficiency in optimization." 1983
- Yurii Nesterov "Primal-dual subgradient methods for convex problems" In: Mathematical Programming, 120(1), 2009
- Nicolas Flammarion and Francis Bach "Stochastic composite least-squares regression with convergence rate O(1/n)" In: COLT 2017.
- Brendan McMahan et al. "Communicationefficient learning of deep networks from decentralized data" In: AISTATS 2017
- Blake Woodworth et al "Minibatch vs Local SGD for Heterogeneous Distributed Learning" In: NeurIPS 2020.



[[]McMahan et al., 2017]