



THE UNIVERSITY OF
SYDNEY



A-FedPD: Aligning Dual-Drift is All Federated Primal-Dual Learning Needs

Yan Sun¹, Li Shen², Dacheng Tao³

¹The University of Sydney ²Sun Yat-sen University ³Nanyang Technological University

Introduction

As a popular paradigm for juggling data privacy and collaborative training, federated learning (FL) is flourishing to distributively process the large scale of heterogeneous datasets on edged clients. From an optimization perspective, two mainstream frameworks have gained widespread applications.

	Federated Primal Methods	Federated Primal-dual Methods
Basis	mandatory aggregation	constraints optimization
Sub-problems splitting	natural operator splitting	Lagrangian operator splitting
Advantages	1. easy to implement	1. support longer local training 2. fewer communications

Federated ADMM (Primal-dual) Family

The consensus finite-sum minimization:

$$\min_{\theta, \theta_i} \frac{1}{C} \sum_{i \in [C]} f_i(\theta_i), \quad s.t. \theta_i = \theta \text{ for } \forall i.$$

Augmented Lagrangian function:

$$L = \frac{1}{C} \sum_{i \in [C]} f_i(\theta_i) + \langle \lambda_i, \theta_i - \theta \rangle + \frac{\rho}{2} \|\theta_i - \theta\|^2.$$

Alternating multiplier optimization: $\{\theta_i \rightarrow \theta \rightarrow \lambda_i \rightarrow \dots \dots\}$

When Primal Dual meets Partial Participation

Partial participation has led to the incompleteness of local problem solutions, and we have summarized the following three frameworks.

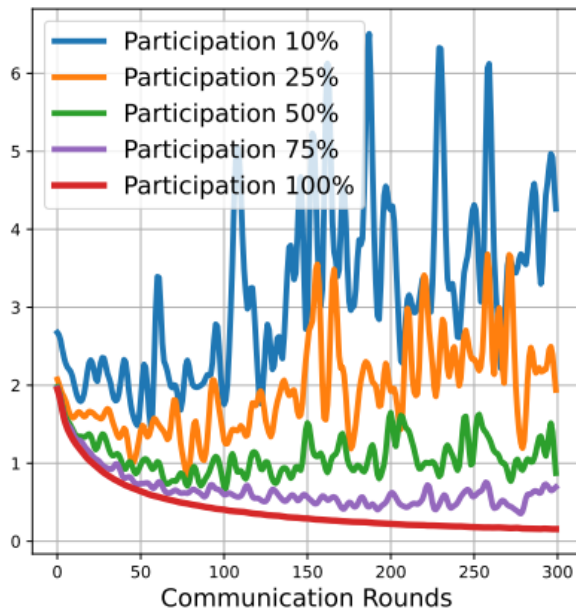
	FedPD	FedADMM	FedDyn
Global consensus	Average in [C]	Average in [P]	Average in [P]
Global dual	Average in [C]	Average in [P]	Average in [C]
Local primal	[C]	[P]	[P]

Average in [C]: average the parameters of all clients

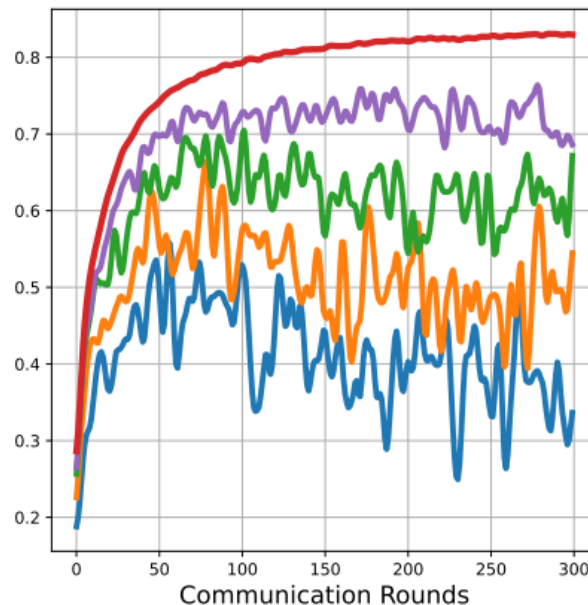
Average in [P]: average the parameters of participated clients

Instability in Federated Primal Dual Family

- significant biases caused by partial participation



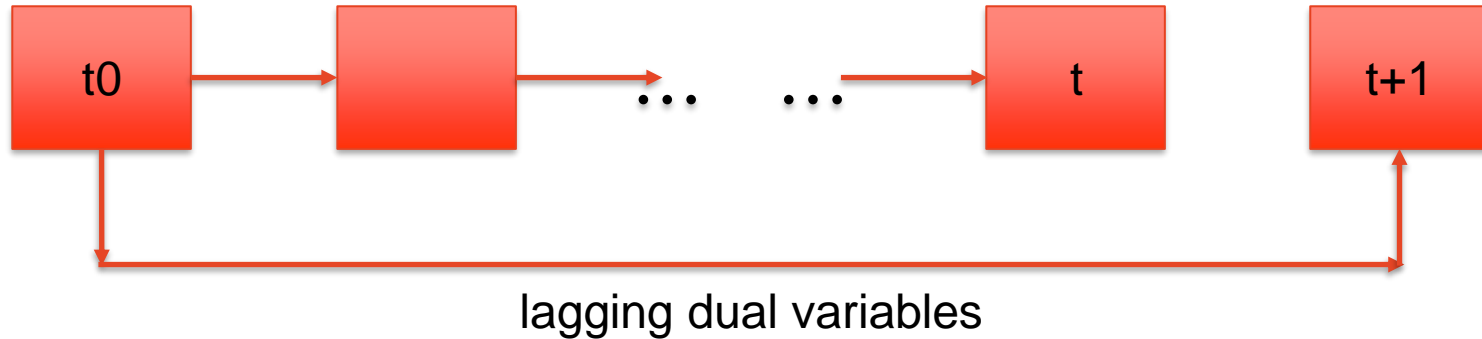
(a) Train loss.



(b) Test accuracy.

"Dual Drift"

- At round t , $\theta^{t+1} = \operatorname{argmin}_{\theta} L(\theta^t, \lambda_i^t)$ for $i \in [P^t]$. At round $t + 1$, when a client $i \notin [P^{\tau}]_{\tau=t_0+1}^t$ ($t_0 \ll t$) is activated, the local sub-problem $L(\theta^t, \lambda_i^{t_0})$ will fall into a very unstable state due to the mismatch between the primal and dual variables.



Our Method: virtual dual update

Algorithm 1 A-FedPD Algorithm

Input: $\theta^0, \theta_i^0, T, K, \lambda_i^0, \rho$

Output: global average model

```
1: Initialization :  $\theta_i^0 = \theta^0, \lambda_i^0 = 0$ .
2: for  $t = 0, 1, 2, \dots, T - 1$  do
3:   randomly select active clients set  $\mathcal{P}^t$  from  $\mathcal{C}$ 
4:   for client  $i \in \mathcal{P}^t$  in parallel do
5:     receive  $\lambda_i^t, \theta^t$  from the global server
6:      $\theta_{i,k}^{t+1} = \text{LocalTrain}(\lambda_i^t, \theta^t, \eta^t, K)$ 
7:     send  $\theta_{i,k}^{t+1}$  to the global server
8:   end for
9:    $\bar{\theta}^{t+1} = \frac{1}{P} \sum_{i \in \mathcal{P}^t} \theta_i^{t+1}$ 
10:   $\lambda_i^{t+1} = D\text{-Update}(\lambda_i^t, \theta^t, \theta_i^{t+1}, \bar{\theta}^{t+1}, \mathcal{P}^t)$ 
11:   $\bar{\lambda}^{t+1} = \frac{1}{C} \sum_{i \in \mathcal{C}} \lambda_i^{t+1}$ 
12:   $\theta^{t+1} = \bar{\theta}^{t+1} + \frac{1}{\rho} \bar{\lambda}^{t+1}$ 
13: end for
14: return global average model
```

\diamond *LocalTrain*: (Optimize Eq.(4))

Input: $\lambda_i^t, \theta^t, \eta^t, K$

Output: $\theta_{i,K}^t$

```
1: for  $k = 0, 1, 2, \dots, K - 1$  do
2:   calculate the stochastic gradient  $g_{i,k}^t$ 
3:    $\theta_{i,k+1}^t = \theta_{i,k}^t - \eta^t (g_{i,k}^t + \lambda_i^t + \rho(\theta_{i,k}^t - \theta^t))$ 
4: end for
```

\diamond *D-Update*: (update dual variables)

Input: $\lambda_i^t, \theta^t, \theta_i^{t+1}, \bar{\theta}^{t+1}, \mathcal{P}^t$

Output: λ_i^{t+1}

```
1: if  $i \in \mathcal{P}^t$  then
2:    $\lambda_i^{t+1} = \lambda_i^t + \rho_t (\theta_i^{t+1} - \theta^t)$ 
3: else
4:    $\lambda_i^{t+1} = \lambda_i^t + \rho_t (\bar{\theta}^{t+1} - \theta^t)$ 
5: end if
```

General updates via
the true model θ_i

Virtual updates via
the constructed θ

Optimization and Generalization

Theorem 1 Let non-convex objective f satisfies Assumption 1, let ρ be selected as a non-zero positive constant, $\{\bar{\theta}^t\}_{t=0}^T$ sequence generated by algorithm 1 satisfies:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f(\bar{\theta}^t)\|^2 \leq \frac{\rho [f(\bar{\theta}^1) - f^*] + R_0}{T} + \mathcal{O}(\epsilon),$$

where f^* is the optimum and $R_0 = \frac{1}{C} \sum_{i \in \mathcal{C}} \mathbb{E}_t \|\theta_i^1 - \theta^0\|^2$ is the first local training volumes.

Theorem 2 Let non-convex objective f satisfies Assumption 1 and 2 and $H = \sup_{\theta, \xi} f(\theta, \xi)$, after T communication rounds training with Algorithm 1, the generalization error bound achieves:

$$\mathbb{E} [F(\theta^T) - f(\theta^T)] \leq \frac{\kappa_c}{CS} (HPT)^{\frac{\mu L}{1+\mu L}},$$

where μ is a constant related to the learning rate and $\kappa_c = 4 (G^2/L)^{\frac{1}{1+\mu L}}$ is a constant.

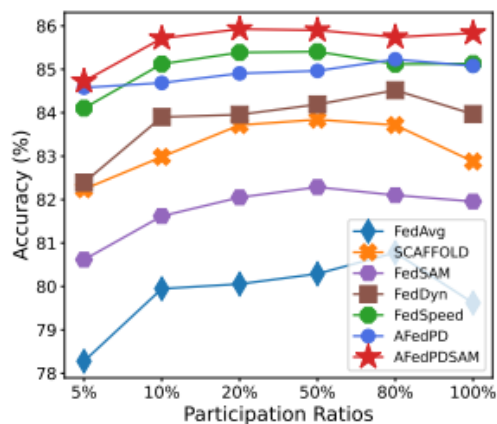
Experimental Results

Table 2: Test accuracy on the CIFAR-10 / 100 dataset. We fix the total client $C = 100$ and $P = 10$ under training local 50 iterations. We test 3 setups of IID, Dir-1.0, and Dir-0.1 on each dataset. Each group is tested on LeNet (upper portion) and ResNet-18 (lower portion) models. Each results are tested with 4 different random seeds. “-” means can not stably converge. “Family” distinguishes whether the algorithm is a primal method (P) or a primal dual method (PD) and “Local Opt” distinguishes whether the algorithm adopts SGD-based or SAM-based local optimizer.

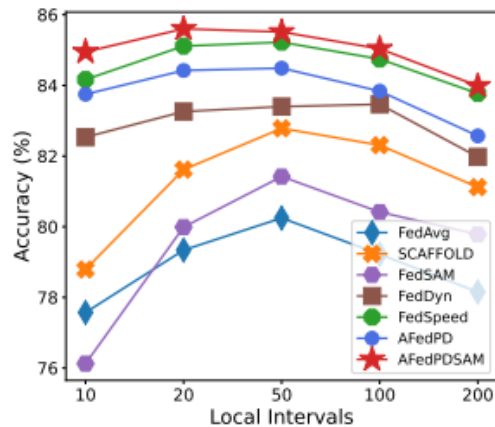
	FAMILY	LOCAL OPT	CIFAR-10			CIFAR-100		
			IID	DIR-1.0	DIR-0.1	IID	DIR-1.0	DIR-0.1
FEDAVG	P	SGD	81.87 \pm .12	80.58 \pm .15	75.57 \pm .27	40.11 \pm .17	39.65 \pm .07	38.37 \pm .14
FEDCM	P	SGD	80.34 \pm .14	79.31 \pm .33	72.89 \pm .37	43.33 \pm .13	42.35 \pm .25	37.11 \pm .51
SCAFFOLD	P	SGD	84.25 \pm .16	83.61 \pm .14	78.66 \pm .29	49.65 \pm .06	49.11 \pm .14	46.36 \pm .30
FEDSAM	P	SAM	83.22 \pm .09	81.94 \pm .13	77.41 \pm .36	43.02 \pm .09	42.83 \pm .29	42.29 \pm .23
FEDDYN	PD	SGD	84.49 \pm .22	84.20 \pm .14	79.51 \pm .13	50.27 \pm .11	49.64 \pm .21	46.30 \pm .26
FEDSPEED	PD	SAM	86.01 \pm .18	85.11 \pm .21	80.86 \pm .18	54.01 \pm .15	53.45 \pm .23	51.28 \pm .18
A-FEDPD	PD	SGD	85.31 \pm .14	84.94 \pm .13	80.28 \pm .20	51.41 \pm .15	51.17 \pm .17	48.15 \pm .28
A-FEDPDSAM	PD	SAM	86.47 \pm .18	85.90 \pm .29	81.96 \pm .19	55.56 \pm .27	54.62 \pm .16	53.15 \pm .19
FEDAVG	P	SGD	81.67 \pm .21	80.94 \pm .17	76.24 \pm .35	44.68 \pm .21	44.27 \pm .25	41.64 \pm .27
FEDCM	P	SGD	84.22 \pm .17	82.85 \pm .21	76.93 \pm .32	50.04 \pm .16	48.66 \pm .28	44.07 \pm .30
SCAFFOLD	P	SGD	84.31 \pm .14	83.70 \pm .11	78.70 \pm .21	50.69 \pm .21	50.28 \pm .21	47.12 \pm .34
FEDSAM	P	SAM	83.79 \pm .28	82.58 \pm .19	77.83 \pm .27	48.66 \pm .29	48.42 \pm .19	45.03 \pm .22
FEDDYN	PD	SGD	83.71 \pm .26	82.66 \pm .15	79.44 \pm .25	-	-	-
FEDSPEED	PD	SAM	86.90 \pm .18	85.92 \pm .24	81.47 \pm .19	53.22 \pm .28	52.75 \pm .16	49.66 \pm .13
A-FEDPD	PD	SGD	85.11 \pm .12	84.33 \pm .16	81.05 \pm .28	48.15 \pm .22	48.02 \pm .29	46.24 \pm .26
A-FEDPDSAM	PD	SAM	87.44 \pm .13	86.46 \pm .25	82.48 \pm .21	55.30 \pm .23	53.49 \pm .17	50.31 \pm .23



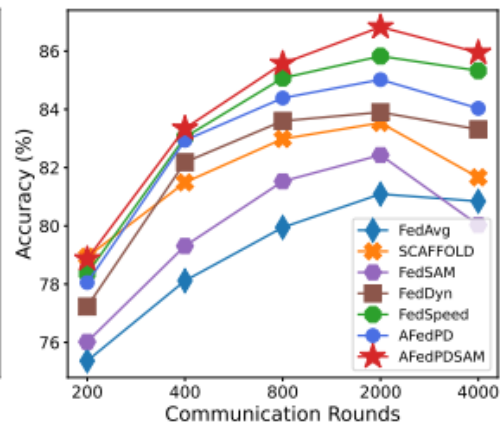
Ablation Study



(a) Different Participation Ratios.



(b) Different Local Intervals.



(c) Different Rounds.

Figure 2: Test of the proposed *A-FedPD* method on setups of different participation ratios, different local intervals, and different rounds. In these experiments, we fix the total training data samples and total training iterations and then learn their variation trends.



THE UNIVERSITY OF
SYDNEY

Thank you!