

A New Formulation of Federated Learning

• Cross-silo Federated Learning (FL) is typically formulated as the optimization problem

$$\min_{x \in \mathbb{R}^d} \left[f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right], \qquad (\text{ERM})$$

where i = 1, 2, ..., n are the clients/silos, and f_i is the loss defined by the data owned by client i.

2 We propose a **new formulation** of FL, which we call **FedMix**:

$$x^* = \arg\min_{x \in \mathbb{R}^d} \left[\tilde{f}(x) = \frac{1}{n} \sum_{i=1}^n f_i(\alpha_i x + (1 - \alpha_i) x_i) \right]. \quad (\mathsf{FedMix})$$

- $x_i = \operatorname{argmin} f_i(x)$ is the pure model trained entirely on data owned by device i
- $\alpha_1, \ldots, \alpha_n \in [0, 1]$ are explicit personalization parameters.
- At node i we deploy the personalized model

$$T_i(x) = \alpha_i x^* + (1 - \alpha_i) x_i.$$

3 Key properties of **FedMix**:

- Efficiently solvable as a finite-sum problem
- Adaptive to communication constraints
- Adaptive to personalization

Motivation 1: From Local GD to FedMix

1 Local GD (LGD) is a simple form of Federated Averaging for solving (ERM), used in the hope that its local steps address the communication bottleneck.

• Local GD:

$$\dot{x}_{t+1}^{i} = \begin{cases} x_t^{i} - \gamma \nabla f_i(x_t^{i}) & \text{if } t \mod H \neq 0\\ \frac{1}{n} \sum_{i=1}^{n} \left[x_t^{i} - \gamma \nabla f_i(x_t^{i}) \right] & \text{if } t \mod H = 0 \end{cases},$$

where $H \geq 1$ is the number of local steps.

• However, LGD does not improve on GD in terms of communication complexity when solving ERM [2, 1] with heterogeneous data!

2 In light of this, Hanzely and Richtárik [2] recently proposed the implicitly personalized FL formulation

$$\min_{x \in \mathbb{R}^d} \left[f_{\lambda}(y_1, y_2, \dots, y_n) = \frac{1}{n} \sum_{i=1}^n f_i(y_i) + \frac{\lambda}{2n} \|y_i - \overline{y}\|^2 \right], \quad (\text{IPFL})$$

where $\overline{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$ and $\lambda > 0$ is a regularization parameter.

- As the λ varies between 0 and ∞ , the solutions of (IPFL) interpolate between the pure local optimal models (i.e., $x_i = \operatorname{argmin} f_i(x)$) and the minimizer of (ERM). The solutions y_1, \ldots, y_n found by Local GD are an *implicit mixture* of the local minimizers x_1, \ldots, x_n and the solution of (ERM) $\operatorname{argmin} f(x)$.
- Key observation: When applying SGD to (IPFL), seen as a 2-sum problem, one recovers a (variant) of LGD.
- Key result of [2]: However, when LGD is seen as a method for solving (IPFL) as opposed to (ERM), its communication complexity improves even in the heterogeneous case, and diminishes to 0 as $\lambda \to 0$. So, with increased personalization (= smaller λ), they get better communication complexity, resolving an important issue in FL.
- ³ Motivation 1: Can we design a different formulation of FL, one in which we would have the same benefits, but where local steps would not be needed? **FedMix** is the answer!

FedMix: A Simple and Communication-Efficient Alternative to Local Methods in Federated Learning

Motivation 2: from MAML to FedMix • Model Agnostic Meta Learning (MAML) objective is $\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(x - \gamma \nabla f_i(x)).$ (1)**2** Key observation: Suppose that we run GD for $H \ge 1$ steps on the quadratic objective $f_i(x) = \frac{1}{2}x^T A_i x - b_i^T x + c$, starting from some $x^0 \in \mathbb{R}^d$. If the stepsize satisfies $\gamma \leq \frac{1}{L_c}$, where $L_i = \lambda_{\max}(\mathbf{A}_i)$, then the final iterate x_i^H can be written as $x_i^H = \left(\boldsymbol{I} - \boldsymbol{J}_i^H \right) x_i + \boldsymbol{J}_i^H x^0,$ where $x_i = \arg \min f_i$ and $J_i \in \mathbb{R}^{d \times d}$ is a matrix with maximum eigenvalue smaller than 1, i.e. $\lambda_{\max}(\mathbf{J}) < 1$. • Plugging this result into Equation (1), observe that in MAML we find the initial model x^0 by solving the problem $\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i((\boldsymbol{I} - \boldsymbol{J}_i)x_i + \boldsymbol{J}_i x).$ • Hence, MAML is optimizing or a specific weighted average of the initial model x^0 and the local solutions x_1, x_2, \ldots, x_n . **3** Motivation 2: FedMix can be can be seen as dispensing with the specific matrix J_i , and instead optimizing the average weighted with an arbitrary constant $\alpha_i \in [0, 1]$. Theoretical Properties of FedMix • FedMix preserves smoothness and convexity. In particular, suppose that each objective f_i is L_i -smooth, i.e., $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L_i \|x - y\|$ for all $x, y \in \mathbb{R}^d$. Then for the **FedMix** objective f, we have (i) \tilde{f} is L_{α} -smooth with $L_{\alpha} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^{n} \alpha_i^2 L_i$ (ii) If each f_i is convex, then \tilde{f} is also convex (iii) If each f_i is μ_i -strongly convex, then \tilde{f} is μ_{α} -strongly convex with $\mu_{lpha} \stackrel{\mathrm{def}}{=} rac{1}{n} \sum_{i=1}^{n} lpha_{i}^{2} \mu_{i}$ • FedMix regularizes model variance. Let $T_1(x), T_2(x), \ldots, T_n(x)$ be the deployed personalized models. For $y_1, \ldots, y_n \in \mathbb{R}^d$ define $V(y_1,\ldots,y_n) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \|y_i - \bar{y}\|^2,$ where $\bar{y} = \frac{1}{n} \sum_{i} y_{i}$. Then, $V(T_1(x), T_2(x), \dots, T_n(x)) = (1 - \beta)^2 V(x_1, x_2, \dots, x_n),$

if we assume that $\beta := \alpha_1 = \alpha_2 = \ldots = \alpha_n$. So, the variance of the deployed models is smaller if β is smaller.

• FedMix enjoys a one-shot learning property. Suppose that each objective f_i is L_i -smooth, and let $\hat{L} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n L_i$. Given the pure local models x_1, x_2, \ldots, x_n , define the weighted average

$$x^{\text{avg}} \stackrel{\text{def}}{=} \sum_{i=1}^{n} w_i x_i, \qquad w_i \stackrel{\text{def}}{=} \frac{\alpha_i^2 L_i}{n L_\alpha}, \qquad L_\alpha \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^{n} \alpha_i^2 L_i.$$
(2)

We further define the constants

$$D \stackrel{\text{def}}{=} \max_{i,j=1,\dots,n, i \neq j} \|x_i - x_j\|, \text{ and } V \stackrel{\text{def}}{=} \sum_{i=1}^n w_i \|x_i - x^{\text{avg}}\|^2.$$
(3)

Fix any $\epsilon > 0$. Assume that either $\max_{i=1,\dots,n} \alpha_i \leq \sqrt{2\epsilon}/\sqrt{\hat{L}D}$, or $\alpha_i = \beta$ for all *i* and $\beta \leq \sqrt{2\epsilon}/\sqrt{\hat{L}D}$. Then x^{avg} is an ϵ -approximate minimizer of (FedMix).

Elnur Gasanov¹ Ahmed Khaled² Samuel Horváth¹ Peter Richtárik¹

¹King Abdullah University of Science and Technology ² Cairo University

Distributed GD Applied to FedMix

Distributed GD applied to FedMix:

 r^{k+1}

$$= x^{k} - \frac{\gamma}{n} \sum_{i=1}^{n} \alpha_{i} \nabla f_{i} \left(\alpha_{i} x^{k} + (1 - \alpha_{i}) x_{i} \right).$$
 (DGD)

Suppose that each f_i in (FedMix) is L_i -smooth and μ_i -strongly convex. Define x^{avg} , L_{α} , and \hat{L} by (2) and V, D by (3). Suppose that we run DGD for K iterations starting from $x^0 = x^{\text{avg}}$. Then the following hold:

i) If the α_i are allowed to be arbitrary, then for $\alpha_{\max} \stackrel{\text{def}}{=} \max_{i=1,\dots,n} \alpha_i$ we have

$$\tilde{f}(x^k) - \min_{x \in \mathbb{R}^d} \tilde{f}(x) \le \left(1 - \frac{\mu_\alpha}{L_\alpha}\right)^K \frac{\alpha_{\max}^2 \hat{L}D}{2}.$$

ii) If $\alpha_i = \beta$ for all *i*, then

$$\tilde{f}(x^k) - \min_{x \in \mathbb{R}^d} \tilde{f}(x) \le \left(1 - \frac{\hat{\mu}}{\hat{L}}\right)^K \frac{\beta^2 \hat{L} V}{2},\tag{4}$$

where $\hat{\mu} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^{n} \mu_i$.

Performance of Compressed Gradient Methods on FedMix

Theoretical performance of Distributed Compressed Gradient Descent (DCGD) and the DIANA method of Mishchenko et al [4] in different settings on FedMix. Note that in all cases, complexity improves as α decreases.

Algorithm	Assumption	Convergence guarantee
DCGD	smoothness scvx	$\mathbb{E}_{\tilde{x}} \ x^k - x^*\ ^2 \le (1 - \gamma_0 \overline{\mu})^k C_1 + C_2$
DOGD		$\mathbb{E}[f(x^k) - f^*] \le ((1 - \gamma_0 \overline{\mu})^k C_1 + C_2) \alpha^2$
DIANA	smoothness, scvx	$\mathbb{E} \ x^k - x^* \ ^2 \le (1 - \rho)^k C$
		$\mathbb{E}[\tilde{f}(x^k) - \tilde{f}^*] \le (1 - \rho)^k C \alpha^2$
DCGD	smoothness, cvx	$\mathbb{E}[\tilde{f}(\overline{x}^k) - \tilde{f}(x^*)] \le \frac{1}{k}C_1\alpha^2 + C_2\alpha$
DIANA	smoothness, cvx	$\mathbb{E}[\tilde{f}(\overline{x}^k) - \tilde{f}(x^*)] \le \frac{1}{k}(C_1\alpha^2 + C_2\alpha)$
DCGD	smoothness	$\min_{0 \le t \le k-1} \mathbb{E} \ \nabla \tilde{f}(x^t)\ ^2 \le \frac{\left(1 + C_1 \gamma_0^2\right)^k C_2}{\gamma_0 k} \alpha^2$
DIANA	smoothness	$\mathbb{E} \ \nabla \tilde{f}(\hat{x})\ ^2 \leq \frac{C}{k} \alpha^2$

Table: Abbreviations: cvx = convex, scvx = strongly convex. All the constants are independent of α .

Generalization of Sine Functions



Figure: Average MSE over clients as a function of the personalization parameter α . In all cases, $0 < \alpha < 1$ is optimal; that is, it is optimal to use a mixture, and not to rely on either the maximally personalized models $x_i = \arg \min f_i$ (this corresponds to $\alpha = 0$), or on the single global model $x^* = \arg \min f$ (this corresponds to $\alpha = 1$).

To test the generalization of FedMix on real data we adopt Stackoverflow dataset [3], and compare its performance with two baselines, FOMAML and Reptile.





Figure: Test accuracy of our FedMix model for different personalization parameter values, versus FOMAML and Reptile. We choose $\alpha_i = \alpha$ for all *i*, (see horizontal axis). FOMAML and Reptile are independent from the personalization parameter α . Plots correspond to different data splittings.





Generalization on Real Clients' Data



(a) 100 workers created out of two (b) 50 workers with distinct data distributions

Optimization Experiments

References

- **1** Reese Pathak and Martin J. Wainwright. FedSplit: An algorithmic framework for fast federated optimization. arXiv:2005.05238, 2020.
- 2 Filip Hanzely and Peter Richtárik. Federated learning of a mixture of global and local models. arXiv:2002.05516, 2020.
- **3** TensorFlow Developers. TensorFlow Federated Stack Overflow dataset, 2021.
- 4 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. ICML 2017.
- **6** Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč and Peter Richtárik. Distributed learning with compressed gradient differences. arXiv:1901.09269, 2019.