

IBM Research ETHzürich

Improving Vertical Federated Learning by Efficient Communication with ADMM

Chulin Xie¹ Pin-Yu Chen² Ce Zhang³ Bo Li¹

¹UIUC, ²IBM Research, ³ETH Zurich



Introduction

Vertical Federated Learning (VFL):

• Features of the samples are partitioned across clients and the labels are owned by server



VIMADMM Workflow

• Key steps of VIMADMM at each round *t*:

(1) Communication from client to server: a batch of local embeddings $\{h_j^{k(t)}\}_{j \in B(t)}$

(2) Sever updates auxiliary variables:

$$z_{j}^{(t)} = \underset{z_{j}}{\operatorname{argmin}} \quad \ell(z_{j}, y_{j}) - \lambda_{j}^{(t-1)^{\top}} z_{j} + \frac{\rho}{2} \left\| \sum_{k=1}^{M} h_{j}^{k(t)} W_{k}^{(t)} - z_{j} \right\|_{F}^{2}, \forall j \in B(t)$$
(3) Sever updates dual variables:

$$\lambda_{j}^{(t)} = \lambda_{j}^{(t-1)} + \rho\left(\sum_{k=1}^{M} h_{j}^{k(t)} W_{k}^{(t)} - z_{j}^{(t)}\right), \forall j \in B(t)$$

(4) Sever updates linear heads: $W_{k}^{(t+1)} = \underset{W_{i}}{\operatorname{argmin}} \quad \beta \mathcal{R}(W_{k}) + \sum_{j \in B(t)} \lambda_{j}^{(t)^{\top}} h_{j}^{k(t)} W_{k} + \sum_{j \in B(t)} \frac{\rho}{2} \left\| \sum_{i \in [M], i \neq k} h_{j}^{i(t)} W_{i}^{(t)} + h_{j}^{k(t)} W_{k} - z_{j}^{(t)} \right\|_{F}^{2}, \forall k \in [M]$

(5) Communication from server to each client: residual variable, dual variables, and one *corresponding* linear head

 ${s_{j}^{k}}^{(t+1)} \triangleq {z_{j}}^{(t)} - \sum_{j=1}^{k} {h_{j}^{i}}^{(t)} W_{i}^{(t+1)}, \forall j \in B(t), \forall k \in [M]$

(6) Client updates local model parameters: $\theta_{k}^{(t+1)} = \underset{\rho}{\operatorname{argmin}} \quad \beta \mathcal{R}(\theta_{k}) + \sum_{j \in B(t)} \lambda_{j}^{(t+1)^{\top}} f(x_{j}^{k};\theta_{k}) W_{k}^{(t+1)} + \frac{\rho}{2} \sum_{j \in B(t)} \left\| s_{j}^{k(t+1)} - f(x_{j}^{k};\theta_{k}) W_{k}^{(t+1)} \right\|_{F}^{2}$ Solved by **multiple** local SGD steps. Reduce **frequency** by

- Aggregation: averaging local embeddings fails to capture the unique properties of each client.
- **Communication**: communicating gradients for *each* training step incurs high costs.

Our contributions:

- ✓ **Framework**: an effective framework with multiple heads (VIM).
- ✓ Algorithm: an ADMM-based method (VIMADMM) reducing communication costs by allowing multiple local updates at each step. ✓ Empirical study: 1) VIMADMM converges faster and achieves higher
- accuracy; 2) client-level explanation under VIM based on the linear heads.

VFL with Multiple Heads (VIM)

- VFL setting (with model splitting setting)
- The features of one sample $\{x_j^1, x_j^2, \dots, x_j^M\}$ is distributed to M clients
- Each client has a local feature set $X_k = \{x_j^k\}_{j=1}^N$ and uploads local embeddings $f(x_j^k; \theta_k)$
- Server aggregates local embeddings and computes gradient with labels $\{y_j\}_{j=1}^N$
- Challenges: using averaged local embeddings as server model input [1] might lose the unique aspects of each local feature set.
- **Our idea**: server learns a model with multiple linear heads W_1, W_2, \ldots, W_M corresponding to local clients, taking their separate contribution into account $\min_{\{W_k\}_{k=1}^M, \{\theta_k\}_{k=1}^M} \sum_{j=1}^N \ell(\sum_{k=1}^M f(x_j^k; \theta_k) W_k, y_j) + \sum_{k=1}^M \beta_k \mathcal{R}_k(\theta_k) + \sum_{k=1}^M \beta_k \mathcal{R}_k(W_k)$

Local embedding Server heads Regularization

- Comparing to SGD-based VFL methods [1,2] where server sends gradients to clients at every training step of the local models, VIMADMM has lower communication costs.
- allowing multiple local updating steps at each round;
- Reduce dimensionality of information by **exchanging ADMM-related variables.**

Experiments

VIMADMM vs SOTA [1,2]

• VFL classification on four datasets: MNIST, CIFAR, NUS-WIDE and



* Faster convergence; improved communication efficiency; higher accuracy.

Client-level Explainability of VIM

• The weights norm of linear heads under clean setting (row 2), under one noisy client (row 4), and test accuracy when each client's test input features are perturbed (row 3).









Solving VIM with ADMM-based method

- Key idea: Multiple heads in VIM enable the alternating direction method of multipliers (ADMM) via a special decomposition into simpler sub-problems that can be solved in a distributed manner.
- **Formulation:** an equivalent **constrained** optimization problem ▲ auxiliary variables

$$\sum_{j=1}^{N} \ell(z_j, y_j) + \sum_{k=1}^{M} \beta_k \mathcal{R}_k(\theta_k) + \sum_{k=1}^{M} \beta_k \mathcal{R}_k(W_k)$$

s.t.
$$\sum_{k=1}^{M} f(x_j^k; \theta_k) W_k - z_j = 0, \forall j \in [N]$$

a **consensus** between the server's output and auxiliary variable • Augmented Lagrangian:

 $z_{j}^{(t+1)} = \operatorname{argmin}_{\mathcal{L}}(\{ heta_{k}^{(t+1)}\}, \{W_{k}^{(t+1)}\}, z_{j}, \{\lambda_{j'}^{(t)}\}), orall j \in [N],$

 $\lambda_{j}^{(t+1)} = \operatorname*{argmin}_{\lambda_{j}} \mathcal{L}(\{ heta_{k}^{(t+1)}\}, \{W_{k}^{(t+1)}\}, \{z_{j'}^{(t+1)}\}, \lambda_{j}), \forall j \in [N],$

 $\theta_{k}^{(t+1)} = \underset{k}{\operatorname{argmin}} \mathcal{L}(\theta_{k}, \{W_{k'}^{(t+1)}\}, \{z_{j}^{(t)}\}, \{\lambda_{j}^{(t)}\}), \forall k \in [M], \mathsf{lent}$

 $W_k^{(t+1)} = \operatorname{argmin} \mathcal{L}(\{ heta_{k'}^{(t)}\}, W_k, \{z_j^{(t)}\}, \{\lambda_j^{(t)}\}), orall k \in [M],$

 $+\sum_{j=1}^N\ell(z_j,y_j) + \sum_{k=1}^Meta_k\left(\mathcal{R}_k(heta_k) + \mathcal{R}_k(W_k)
ight) + \sum_{j=1}^N\lambda_j^ op\left(\sum_{k=1}^Mf(x_j^k; heta_k)W_k - z_j
ight)$ • Solution: dual variables constant **penalty** factor

- Decomposes the problem into four sets of sub-problems over $\{W_k\}, \{\theta_k\}, \{z_j\}, \{\lambda_j\}$
- Alternatively updating in server and clients
- Each sub-problem set can be solved in parallel across M clients or N samples.

The weights of linear heads reflect the importance of local clients; Perturbing the client with high weights has higher impact on test accuracy; VIM enables client denoising by lowering their weights.

More details and results are in our paper:

- Details for algorithm VIMADMM-J for VFL without model splitting setting.
- Results on communication costs comparison, effect of penalty factor and local steps, client summarization and the visualization of the local embedding.

References:

server

[1] Tianyi Chen, Xiao Jin, Yuejiao Sun, and Wotao Yin. Vafl: a method of vertical asynchronous federated learning. arXiv, 2020.

[2] Yaochen Hu, Di Niu, Jianming Yang, and Shengping Zhou. Fdml: A collaborative machine learning framework for distributed features. KDD, 2019.