

DEMYSTIFYING MODEL AVERAGING FOR COMMUNICATION-EFFICIENT FEDERATED MATRIX FACTORIZATION



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



Shuai Wang^{1,2}, Richard Cornelius Suwandi², Tsung-Hui Chang^{1,2}
¹ Shenzhen Research Institute of Big Data, ² The Chinese Univ. of Hong Kong, Shenzhen

Summary

The technique of model averaging (MA) has not been considered for the important matrix factorization (MF) model under the scenario of federated learning (FL).

- Propose a new MA based algorithm, named FedMAvg, by judiciously combining the alternating minimization technique and MA.
- Local GD with diminishing steps and partial client communication can greatly reduce the communication cost, even under non-i.i.d. data.

Federated Matrix Factorization Model

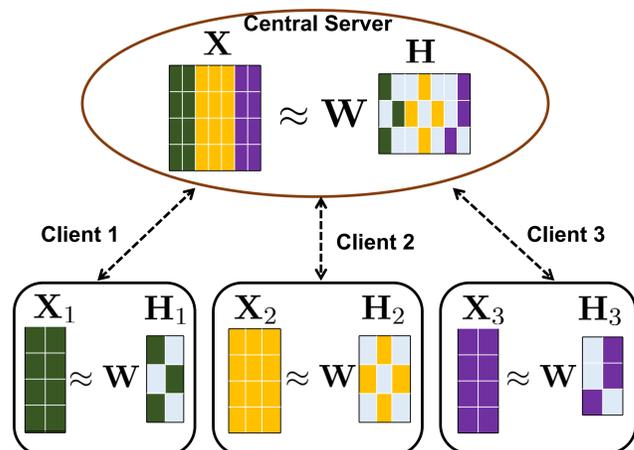
The data samples are partitioned as $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_P]$ and respectively owned by P distributed clients. Each client p owns non-overlapping data $\mathbf{X}_p \in \mathbb{R}^{M \times N_p}$, where N_p is the number of samples of client p and $\sum_{p=1}^P N_p = N$.

$$\min_{\mathbf{W}, \mathbf{H}_p, p=1, \dots, P} F(\mathbf{W}, \mathbf{H}) \triangleq \sum_{p=1}^P \omega_p F_p(\mathbf{W}, \mathbf{H}_p) \quad (1a)$$

$$\text{s.t. } \mathbf{W} \in \mathcal{W}, \mathbf{H}_p \in \mathcal{H}_p, \forall p \in \mathcal{P}, \quad (1b)$$

where $F_p(\mathbf{W}, \mathbf{H}_p) = \frac{1}{N_p} \Phi_p(\mathbf{X}_p, \mathbf{W}\mathbf{H}_p)$, $p \in \mathcal{P}$.

- $\Phi_p(\mathbf{X}_p, \mathbf{W}\mathbf{H}_p)$ measures the quality of the approximation $\mathbf{X}_p \approx \mathbf{W}\mathbf{H}_p$, e.g. $\frac{1}{N_p} \|\mathbf{X}_p - \mathbf{W}\mathbf{H}_p\|_F^2$.
- P could be large, N_p , $p = 1, \dots, P$, could be unbalanced, and \mathbf{X}_p , $p \in \mathcal{P}$ could be non-i.i.d.
- Problem (1) is challenging to solve since it is non-convex and non-smooth, and involves two blocks of variables \mathbf{W} and \mathbf{H} .



Algorithm Development

Alternating Minimization:

Given \mathbf{W}^{s-1} , each client p performs

$$\mathbf{H}_p^s = \arg \min_{\mathbf{H}_p \in \mathcal{H}_p} F_p(\mathbf{W}^{s-1}, \mathbf{H}_p), \quad (2a)$$

$$\mathbf{W}_p^s = \arg \min_{\mathbf{W}} F_p(\mathbf{W}, \mathbf{H}_p^s). \quad (2b)$$

The server does $\mathbf{W}^s = \mathcal{P}_{\mathcal{W}}(\sum_{p=1}^P \omega_p \mathbf{W}_p^s)$.

Local GD with Diminishing Q_2 :

- (2a) via $Q_1 \geq 1$ consecutive steps of PGD with respect to \mathbf{H}_p .
- (2b) via $Q_2^s \geq 1$ ($Q_2^s = \lfloor \frac{Q_2}{s} \rfloor + 1$) consecutive steps of GD with respect to \mathbf{W}_p .

Partial Client Communication (PCC):

- For each round, m clients in \mathcal{A}^s are selected by the server.
- All clients perform updating but **only the clients in \mathcal{A}^s upload their models** to the server for averaging.

Proposed FedMAvg Method

Algorithm 1 Proposed FedMAvg algorithm

Input: initial values of $\mathbf{W}_1^0 = \dots = \mathbf{W}_P^0$ at the server side, initial values of $\{\mathbf{H}_p^0\}_{p=1}^P$ at the clients, $\mathcal{A}^0 = \{1, \dots, P\}$ and \hat{Q} .
for round $s = 1$ **to** S **do**

Server side: Compute

$$\mathbf{W}^s = \mathcal{P}_{\mathcal{W}}\left(\frac{1}{m} \sum_{p \in \mathcal{A}^{s-1}} \mathbf{W}_p^{s-1}\right),$$

and select a set of clients \mathcal{A}^s (with size $|\mathcal{A}^s| = m$) by sampling with replacement according to probabilities $\{\omega_1, \dots, \omega_P\}$, and broadcast \mathbf{W}^s to all clients.

Client side:

for client $p = 1$ **to** P **in parallel do**

Set $\mathbf{H}_p^{s,0} = \mathbf{H}_p^{s-1}$ and $\mathbf{W}_p^{s,0} = \mathbf{W}^s$.

for epoch $t = 1$ **to** Q_1 **do**

$$\mathbf{H}_p^{s,t} = \mathcal{P}_{\mathcal{H}_p}\left(\mathbf{H}_p^{s,t-1} - \frac{\nabla_{\mathbf{H}_p} F_p(\mathbf{W}_p^{s,t-1}, \mathbf{H}_p^{s,t-1})}{c_p^s}\right)$$

$$\mathbf{W}_p^{s,t} = \mathbf{W}_p^{s,t-1}.$$

end for

for epoch $t = Q_1 + 1$ **to** $Q_2^s = Q_1 + Q_2^s$ **do**

$$\mathbf{W}_p^{s,t} = \mathbf{W}_p^{s,t-1} - \frac{\nabla_{\mathbf{W}} F_p(\mathbf{W}_p^{s,t-1}, \mathbf{H}_p^{s,t-1})}{d^s},$$

$$\mathbf{H}_p^{s,t} = \mathbf{H}_p^{s,t-1}.$$

end for

Denote $\mathbf{W}_p^s = \mathbf{W}_p^{s,Q^s}$ and $\mathbf{H}_p^s = \mathbf{H}_p^{s,Q^s}$.

if client $p \in \mathcal{A}^s$ **then**

Upload \mathbf{W}_p^s to the server.

end if

end for

Convergence Analysis

Bounds:

$$\|\nabla_{\mathbf{W}} F_p(\mathbf{W}, \mathbf{H}_p) - \nabla_{\mathbf{W}} F(\mathbf{W}, \mathbf{H})\|_F^2 \leq \zeta^2, \quad (3)$$

$$\|\nabla_{\mathbf{W}} F(\mathbf{W}, \mathbf{H})\|_F^2 \leq \phi^2, \quad (4)$$

Virtual Sequences: $\forall t = 1, \dots, Q$,

$$\tilde{\mathbf{W}}^{s,t} = \mathcal{P}_{\mathcal{W}}\left(\frac{1}{m} \sum_{p \in \mathcal{A}^s} \mathbf{W}_p^{s,t}\right), \tilde{\mathbf{W}}^{s,0} = \mathbf{W}^s, \quad (5)$$

Proximal Gradient:

$$G_H^{s,t} \triangleq \sum_{p=1}^P \omega_p (c_p^s)^2 \|\mathbf{H}_p^{s,t} - \mathcal{P}_{\mathcal{H}_p}(\mathbf{H}_p^{s,t})\|_F^2, \quad (6)$$

$$G_W^{s,t} \triangleq (d^s)^2 \|\tilde{\mathbf{W}}^{s,t} - \mathcal{P}_{\mathcal{W}}(\tilde{\mathbf{W}}^{s,t})\|_F^2, \quad (7)$$

Theorem 1 Let $Q_2^s = \lfloor \frac{Q_2}{s} \rfloor + 1$, and let T be the total number of iterations. Moreover, let $c_p^s = \frac{\gamma_1}{2} L_{H_p}$, $d^s = \gamma_2 L_W$, where $\gamma_1 > 1$ and $\gamma_2 \geq Q_2^s \sqrt{2(7 + 4L_W^2/L_H^2)}$. Then, under Assumptions, the sequence $\{(\tilde{\mathbf{W}}^{s,t}, \mathbf{H}^{s,t})\}$ satisfies

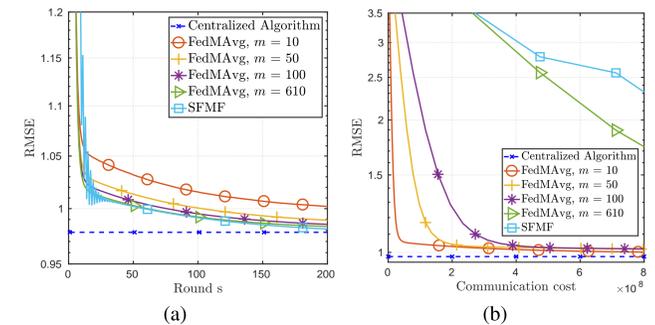
$$\begin{aligned} & \frac{1}{T} \left[\sum_{s=1}^S \sum_{t=1}^{Q_1} \mathbb{E}[G_H^{s,t-1}] + \sum_{s=1}^S \sum_{t=Q_1+1}^{Q_2^s} \mathbb{E}[G_W^{s,t-1}] \right] \\ & \leq \frac{D}{T} \left(F(\tilde{\mathbf{W}}^{1,0}, \mathbf{H}^{1,0}) - F \right) + \left(\frac{8D\zeta^2}{m\gamma_2 L_W} + \frac{96\zeta^2}{m} \right) \\ & + \frac{2D(1+8/m)(\frac{11}{3}\zeta^2 + \phi^2) \sum_{s=1}^S C_1^s}{T\gamma_2^3 L_W} \\ & + \frac{(\frac{11}{3}\zeta^2 + \phi^2) \sum_{s=1}^S C_2^s}{T\gamma_2^2} + \frac{3(\zeta^2 + \phi^2) \sum_{s=1}^S C_1^s}{2T}, \quad (8) \end{aligned}$$

where $D \triangleq \frac{\gamma_1^2 L_H}{2(\gamma_1 - 1)} + \frac{6(\gamma_2^2 + 1)L_W^2}{(\gamma_2 - 1)L_W}$, $C_1^s \triangleq Q_2^s(Q_2^s - 1)(2Q_2^s - 1)$, and $C_2^s \triangleq 6(3Q_2^s(Q_2^s - 1)/2 + 4 + 32/m)C_1^s$.

Numerical Results II

Application to Item Recommendation:

Recommendation performance:



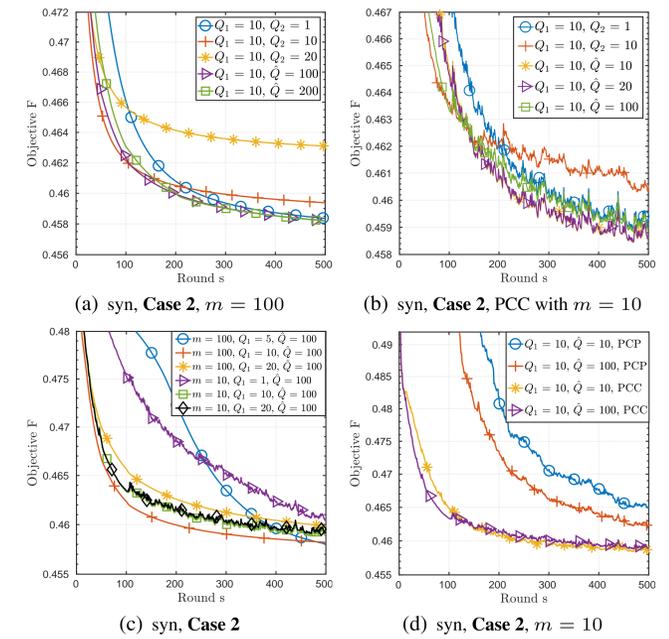
References

- B. Bahmani, B. Moseley, A. Vattani, R. Kumar, and S. Vassilvitskii, "Scalable k-means++," in Proc. VLDB 2012, Istanbul, Turkey, Aug. 27-31 2012, pp. 622-633.
- M.-F. Balcan, S. Ehrlich, and Y. Liang, "Distributed k-means and kmedian clustering on general topologies," in Proc. NeuIPS 2013, Lake Tahoe, USA, Dec. 5-10 2013, pp. 1995-2003.
- J. Chen, E. S. Azer, and Q. Zhang, "A practical algorithm for distributed clustering and outlier detection," in Proc. NeuIPS 2018, Montreal, Quebec, Canada, Dec. 2-8 2018, pp. 2248-2256.
- D. Chai, L. Wang, K. Chen, and Q. Yang, "Secure federated matrix factorization," IEEE Intelligent Systems, vol. 1, No. 1, pp. 1-8, Aug. 2020.

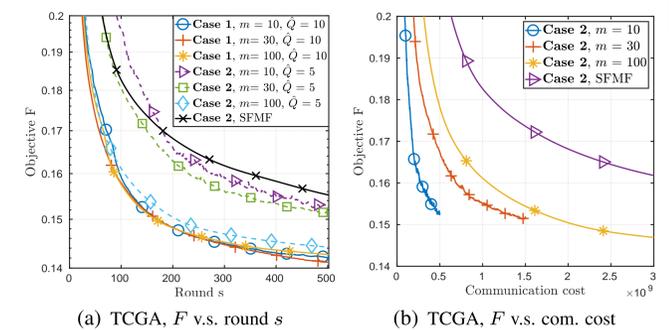
Numerical Results I

Application to Data Clustering:

Convergence behavior:



Effect of non-i.i.d and PCC (Case 2: non-i.i.d):



Clustering performance versus distributed clustering (KM|[1], BEL[2], CAL[3]) and MF (SFMF[4]):

