

Collaborative Learning by Detecting Collaboration Partners

Shu Ding, Wei Wang

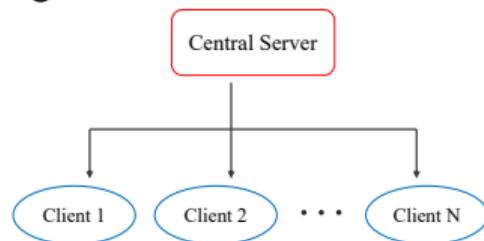
{dings, wangw}@lamda.nju.edu.cn



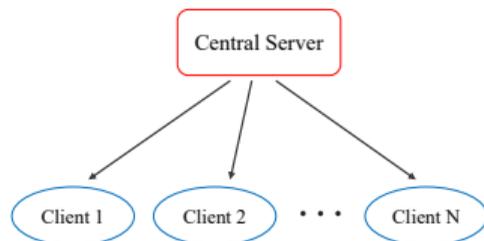
NeurIPS 2022

- ▶ Massive amounts of data are naturally dispersed over numerous clients. Each client only has limited data.
- ▶ **Collaborative learning** is a promising paradigm that enables the clients to learn models through collaboration.

- ▶ Two settings



- *Centralized model*: return one single model for all clients



- *Personalized model*: return different models for different clients

- ▶ Centralized model
One single model may perform badly on clients whose distributions are different from the average distribution.
- ▶ Personalized model
Learning personalized models is impractical when the number of clients N is very large since this costs unaffordable computational resources.
- ▶ Can we return K ($K \ll N$) appropriate models for N heterogeneous clients and expect that the returned models have comparable performance to personalized models?

► Preliminaries

- Clients $\{C_1, \dots, C_N\}$ with distributions $\{\mathcal{D}_1, \dots, \mathcal{D}_N\}$
- Each client C_i has access to m_i examples $S_i = \{(\mathbf{x}_1^i, y_1^i), \dots, (\mathbf{x}_{m_i}^i, y_{m_i}^i)\}$ drawn from \mathcal{D}_i
- Total number of examples $M = \sum_{i=1}^N m_i$

► Collaborative learning scenario

- Train the model over the weighted union of all samples $S_\alpha = \sum_{j=1}^N \alpha_j S_j$
- The model for C_i can be learned by minimizing $\hat{\mathcal{L}}_{\alpha_i}(h) = \sum_{j=1}^N \alpha_{ij} \hat{\mathcal{L}}_{S_j}(h)$ with collaboration vector $\alpha_i = (\alpha_{i1}, \dots, \alpha_{iN}) \in \Delta^N$

Theorem (Generalization Bound)

Let \mathcal{H} be the hypothesis space with VC-dimension d . Denote $h_i^* = \arg \min_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}_i}(h)$ and $\hat{h}_{\alpha_i} = \arg \min_{h \in \mathcal{H}} \hat{\mathcal{L}}_{\alpha_i}(h)$. For any given $\delta \in (0, 1)$ and $\forall i \in \{1, \dots, N\}$, with probability at least $1 - \delta$:

$$\mathcal{L}_{\mathcal{D}_i}(\hat{h}_{\alpha_i}) - \mathcal{L}_{\mathcal{D}_i}(h_i^*) \leq 2 \sum_{j=1}^N \alpha_{ij} d_{\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_j) + 2\mu \sqrt{\sum_{j=1}^N \frac{\alpha_{ij}^2}{m_j}} \sqrt{8(d \log(2M) + \log \frac{8}{\delta})}.$$

Here $d_{\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_j) = \sup_{h \in \mathcal{H}} |\mathcal{L}_{\mathcal{D}_i}(h) - \mathcal{L}_{\mathcal{D}_j}(h)|$ is the Integral Probability Metrics (IPM).

Theorem (Optimal Collaboration Vector)

Let $\Xi_i^j = d_{\mathcal{H}}(\mathcal{D}_i, \mathcal{D}_j)$ and $\lambda = \mu \sqrt{8(d \log(2M) + \log \frac{8}{\delta})}$. For client C_i , sort $\{\Xi_i^1, \dots, \Xi_i^N\}$ in ascending order to get $\{\Xi_i^{\sigma(1)}, \dots, \Xi_i^{\sigma(N)}\}$. The optimal α_i^* for client C_i is given by

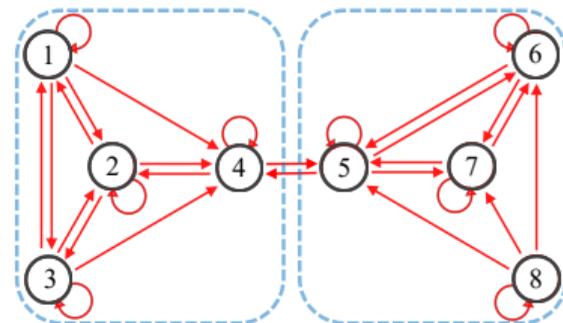
$$\alpha_{ij}^* = \left[\frac{m_j(\zeta - \Xi_i^j)}{\sum_{q \leq q_i} m_{\sigma(q)}(\zeta - \Xi_i^{\sigma(q)})} \right]_+.$$

Here $[\cdot]_+ = \max(\cdot, 0)$, ζ is the larger root of equation $\sum_{q \leq q_i} m_{\sigma(q)} (\zeta - \Xi_i^{\sigma(q)})^2 = \lambda^2$, and $q_i = \arg \max_t \left\{ t \mid \zeta \geq \Xi_i^{\sigma(t)} \wedge \left(\sum_{q \leq t} m_{\sigma(q)} \Xi_i^{\sigma(q)} \right)^2 \geq \left(\sum_{q \leq t} m_{\sigma(q)} \right) \left(\sum_{q \leq t} m_{\sigma(q)} (\Xi_i^{\sigma(q)})^2 - \lambda^2 \right) \right\}$.

- ▶ $\hat{h}_{\alpha_i^*}$ with respect to the optimal α_i^* is referred as the *personalized* model for client C_i

Collaboration Partners

- ▶ In the directed graph A , $\alpha_{ij}^* > 0$ means C_j is beneficial to C_i . Clients with similar incoming edges are called **collaboration partners** since they need similar contribution from other clients.
- ▶ Intuitively, collaboration partners should be in the same group. We could probably return the same model for C_1, C_2, C_3, C_4 while it is inappropriate to return the same model for C_4, C_5 .



In graph $A = (V, E)$, $|V| = N$, node i denotes C_i and the weight of edge from j to i is α_{ij}^* .

► Collaboration with Modularity Maximization

- Construct matrix \mathbf{U} to evaluate the **incoming-edge similarity** among clients

$$\mathbf{U} = \mathbf{D}_{in}^{-\beta} \mathbf{A} \mathbf{A}^T \mathbf{D}_{in}^{-\beta}$$

- Use **Modularity** as the objective function to evaluate the quality of group partitions

$$Q(\mathcal{G}) = \frac{1}{2W} \sum_{i,j} \left[w_{ij} - \frac{d_i d_j}{2W} \right] \delta(g_i, g_j)$$

- Relax the modularity maximization problem as a **SemiDefinite Programming**

$$\begin{aligned} \max \quad & \sum_{\mathcal{M}^+} \mathcal{M}_{ij} \boldsymbol{\nu}_i \cdot \boldsymbol{\nu}_j + \sum_{\mathcal{M}^-} -\mathcal{M}_{ij} (1 - \boldsymbol{\nu}_i \cdot \boldsymbol{\nu}_j) \\ \text{s.t.} \quad & \boldsymbol{\nu}_i \cdot \boldsymbol{\nu}_i = 1, \forall i \in \{1, \dots, N\}; \quad \boldsymbol{\nu}_i \cdot \boldsymbol{\nu}_j \geq 0, \quad \forall i \neq j, \\ & \boldsymbol{\nu}_i \in \mathbb{R}^K, \forall i \in \{1, \dots, N\}. \end{aligned}$$

▶ Collaboration with Modularity Maximization

- **Find reasonable group partitions** by solving the SDP

Given matrix \mathbf{U} , let $Q(\mathcal{G})$ be the modularity value of the group partition \mathcal{G} obtained by solving the SDP using rounding techniques. Then $Q(\mathcal{G}) > \kappa \text{OPT}_{Q(\mathcal{G})} - (1 - \kappa)$ where $\kappa = 0.766$ is the approximation factor.

- **Detect bad clients**

Edge $e_{ij} \in \mathbf{U}$ is a *weak edge* if its weight $w_{ij} < \frac{1}{N}$. A group is divided into several disjoint parts after removing all weak edges within the group. Clients do not belong to the largest part are *bad clients*. Bad clients cannot be provided with good performance guarantee.

► Collaboration with Modularity Maximization

• Number of bad clients

Given the group partition $\mathcal{G} = \{G_1, \dots, G_K\}$ returned by Algorithm ACLMM, assume

$N_k \geq 2\sqrt{Z_{in}}, \forall k \in \{1, \dots, K\}$. Let $N_{min} = \min_k N_k$, then $|\mathcal{B}| \leq \frac{N_{min} - \sqrt{N_{min}^2 - 4Z_{in}}}{2}$, where $Z_{in} \leq \frac{N}{2(N-1)} \left[\frac{N^2 - KN}{K} - 2W \left((\kappa + 1) \text{OPT}_{Q(\mathcal{G})} - \frac{K-1}{K} \right) \right]$.

• Theoretical guarantee

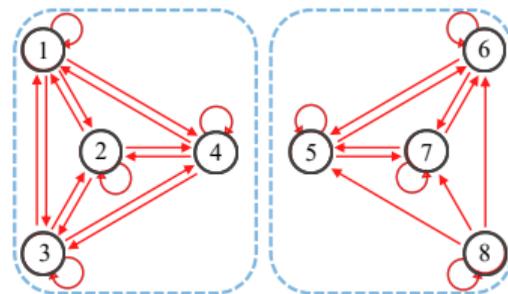
Let $\mathcal{G} = \{G_1, \dots, G_K\}$ be the group partition returned by solving the SDP. $\hat{h}_{\alpha_{G_k}}$ is the model returned by Algorithm ACLMM for client C_i in group G_k . $\text{upp}(\hat{h}_{\alpha_{G_k}})$ is the upper bound of the expected risk of $\hat{h}_{\alpha_{G_k}}$ and $\text{upp}(\hat{h}_{\alpha_i^*})$ is the upper bound of the expected risk of the personalized model $\hat{h}_{\alpha_i^*}$. The following result holds except for the bad clients in \mathcal{B} :

$$\text{upp}(\hat{h}_{\alpha_{G_k}}) - \text{upp}(\hat{h}_{\alpha_i^*}) \leq O \left(\eta(1 - \tau) \sqrt{\frac{N}{N-1}} \right).$$

- ▶ Collaboration with Clustering
- **Potential structures**

There exists a potential partition $\mathcal{P}^* = \{P_1^*, \dots, P_K^*\}$ s.t. $\Phi(\mathcal{P}) = \sum_{k=1}^K \sum_{C_i \in P_k} d(\alpha_i^*, \bar{\alpha}_k)$ is small. Assume that $\{\alpha_1^*, \dots, \alpha_N^*\}$ satisfy $(1 + \gamma, \epsilon)$ -approximation-stability property.
- **Detect bad clients**

$\bar{d} = \frac{1}{N} \text{OPT}_{\Phi(\mathcal{P})}$ is the average distance. $d^* = \frac{\gamma \bar{d}}{\epsilon t}$ is the critical distance. C_i is the *bad client* if $d_1(\alpha_i^*) \geq d^*$ or $d_2(\alpha_i^*) - d_1(\alpha_i^*) \leq \frac{t}{2} d^*$.



The example here has better structures than the aforementioned example.

- ▶ Collaboration with Clustering

- **Number of bad clients**

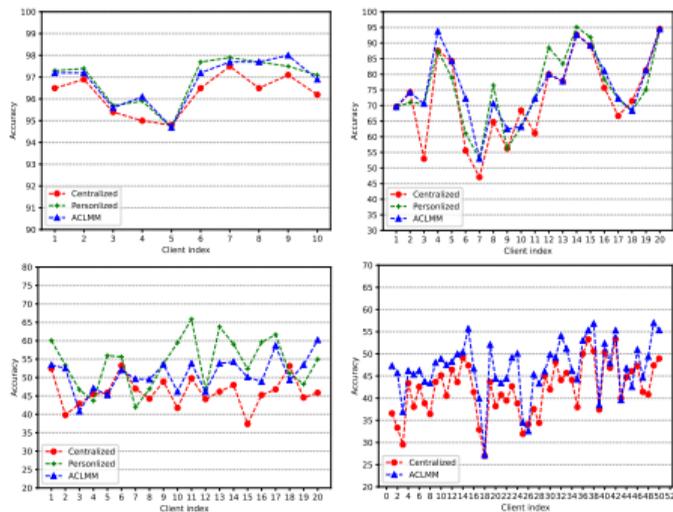
Let $\mathcal{P} = \{P_1, \dots, P_K\}$ be the group partition produced by Algorithm ACLC. Then $|\mathcal{B}| < (6 + \frac{t}{\gamma})\beta\epsilon N$ where $t > 2$ and $\beta > 1$ are given constants.

- **Theoretical guarantee**

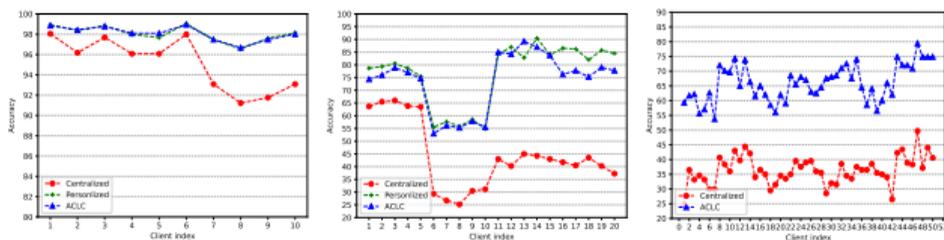
Let $\mathcal{P} = \{P_1, \dots, P_K\}$ be the group partition produced by Algorithm ACLC. $\hat{h}_{\alpha_{P_k}}$ is the model returned by Algorithm ACLC for client C_i in group P_k . $\text{upp}(\hat{h}_{\alpha_{P_k}})$ is the upper bound of the expected risk of $\hat{h}_{\alpha_{P_k}}$ and $\text{upp}(\hat{h}_{\alpha_i^*})$ is the upper bound of the expected risk of the personalized model $\hat{h}_{\alpha_i^*}$. The following result holds except for the bad clients in \mathcal{B} :

$$\text{upp}(\hat{h}_{\alpha_{P_k}}) - \text{upp}(\hat{h}_{\alpha_i^*}) \leq O\left(\frac{\gamma \text{OPT}_{\Phi(\mathcal{P})}}{\epsilon t N}\right).$$

Experimental Results



- ▶ The model learned with ACLMM performs better than the centralized model and is comparable to the personalized model.



- ▶ The model learned with ACLC performs much better than the centralized model and is comparable to the personalized model.
- ▶ The gap between the model return by ACLC and the personalized model is small.

Thank you!