

Federated Spectral Clustering via Secure Similarity Reconstruction

Dong Qiao^{1,2}, Chris Ding¹, Jicong Fan^{1,2*}

¹The Chinese University of Hong Kong, Shenzhen

²Shenzhen Research Institute of Bigdata, Shenzhen

December 9, 2023

Motivation

- Spectral Clustering (SC)
 - Step 1: Construct a similarity matrix;
 - Step 2: Perform normalized cut and Kmeans [Shi and Malik, 2000]
- Limitations of SC
 - SC heavily relies on the construction of similarity matrix;
 - SC cannot be directly performed on distributed dataset;
 - Vanilla SC cannot preserve the user's privacy.

$$K_{xx} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & k(\mathbf{x}_n, \mathbf{x}_2) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

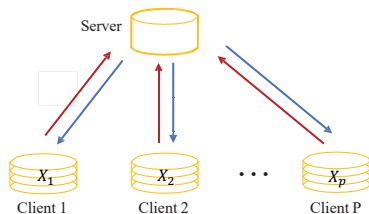


Figure: Distributed data

Similarity Reconstruction via Feature Space Factorization

- Gaussian kernel: $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = \exp(-\frac{1}{r^2} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2)$
- Nonlinear approximation: $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) \simeq \widehat{\phi}(\mathbf{x}_i)^T \widehat{\phi}(\mathbf{x}_j)$ for $\widehat{\phi}(\mathbf{x}_i) = \phi(\mathbf{Z})\mathbf{c}_i$
- For a client p with $\mathbf{X}_p \in \mathbb{R}^{m \times N_p}$, the nonlinear approximation has a matrix form of

$$\widehat{\phi}(\mathbf{X}_p) \simeq \phi(\mathbf{Z})\mathbf{C}_p$$

- To do the nonlinear approximation for P clients, we consider the distributed optimization problem:

$$\min_{\mathbf{Z}, \mathbf{C}} \triangleq \sum_{p=1}^P \omega_p f_p(\mathbf{Z}, \mathbf{C}_p)$$

where $f_p(\mathbf{Z}, \mathbf{C}_p) = \frac{1}{2} \|\phi(\mathbf{X}_p) - \phi(\mathbf{Z})\mathbf{C}_p\|_F^2 + \frac{\lambda}{2} \|\mathbf{C}_p\|_F^2$

Our design: FedSC

- **Stage I: Federated similarity reconstruction** (See left part in Fig 1)

- (1) Local update scheme:

$$\mathbf{C}_p^S = \arg \min_{\mathbf{C}_p} f_p(\mathbf{Z}^{S-1}, \mathbf{C}_p)$$

$$\mathbf{Z}_p^S = \arg \min_{\mathbf{Z}_p} f_p(\mathbf{Z}_p, \mathbf{C}_p^S)$$

- (2) Aggregation scheme:

$$\mathbf{Z}^S = \frac{1}{|\mathcal{A}^{S-1}|} \sum_{p \in \mathcal{A}^{S-1}} \mathbf{Z}_p^S$$

- **Stage II: Spectral Clustering** (See right part in Fig 1)

$$y = \text{SpectralClustering}(\mathbf{C}^T \mathcal{K}(\mathbf{Z}, \mathbf{Z}) \mathbf{C}, \mathbf{K})$$

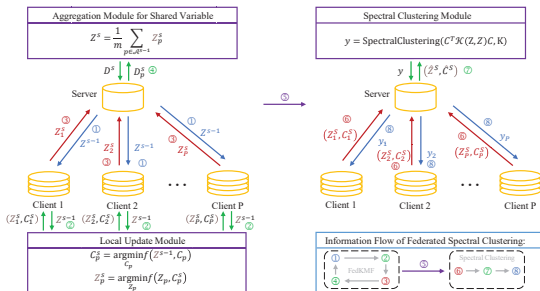


Figure: FedSC

Our design: Secure-Enhanced FedSC

- **Alternative 1:** FedSC with perturbed data

Raw data $\mathbf{X} \in \mathbb{R}^{m \times n}$ is perturbed by a noise matrix $\mathbf{E} \in \mathbb{R}^{m \times n}$ to form the noisy data matrix:

$$\tilde{\mathbf{X}} = \mathbf{X} + \mathbf{E}$$

- **Alternative 2:** FedSC with perturbed factors

We added Gaussian noise to \mathbf{Z} in every round but added to Gaussian noise to \mathbf{C} in the last round. That is,

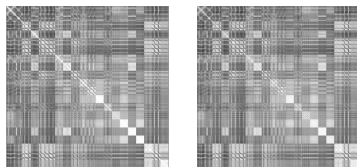
$$\begin{cases} \tilde{\mathbf{Z}}_p^s = \mathbf{Z}_p^s + \mathbf{E}_{Z_p}^s & s = 1, \dots, S; p = 1, \dots, P \\ \tilde{\mathbf{C}}_p^S = \mathbf{C}_p^S + \mathbf{E}_{C_p} & p = 1, \dots, P \end{cases}$$

where $\mathbf{E}_{Z_p}^s$ and \mathbf{E}_{C_p} are drawn from $\mathcal{N}(0, \sigma_Z^2)$ and $\mathcal{N}(0, \sigma_C^2)$. Then, we perform SC on

$$\tilde{K}_{xx} = \tilde{\mathbf{C}}^T \mathcal{K}(\tilde{\mathbf{Z}}, \tilde{\mathbf{Z}}) \tilde{\mathbf{C}}$$

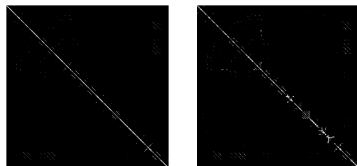
Numerical results

Part I: Similarity Reconstruction



(a) Vanilla SC

(b) FedSC

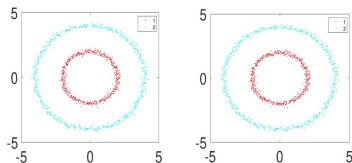


(c) Vanilla SC

(d) FedSC

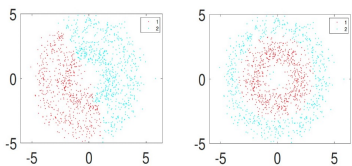
Figure: COIL20

Part II: Clustering performance of FedSC



(a) Vanilla SC

(b) FedSC



(c) Vanilla SC

(d) FedSC

Figure: Concentric circles

Numerical results

- Part III: Comparison with baselines

Table: Comparison of clustering accuracy

		Kmeans	SC	DSC	FedSC
X	Iris	0.8933 \pm 0.0000	0.9000 \pm 0.0000	0.5480 \pm 0.0679	0.9000 \pm 0.0031
	COIL20	0.6113 \pm 0.0534	0.8025 \pm 0.0009	0.1009 \pm 0.0100	0.7828 \pm 0.0231
	MNIST	0.5110 \pm 0.0424	0.5872 \pm 0.0468	0.1340 \pm 0.0113	0.6002 \pm 0.0410
	CIFAR10	0.2171 \pm 0.0132	0.2182 \pm 0.0133	0.1235 \pm 0.0062	0.2134 \pm 0.0131
X with 0.3σ	Iris	0.8420 \pm 0.0274	0.8327 \pm 0.0267	0.4533 \pm 0.0674	0.8427 \pm 0.0404
	COIL20	0.6422 \pm 0.0366	0.7997 \pm 0.0029	0.0981 \pm 0.0084	0.7793 \pm 0.0240
	MNIST	0.5530 \pm 0.0364	0.6246 \pm 0.0492	0.1355 \pm 0.0114	0.5908 \pm 0.0510
	CIFAR10	0.2209 \pm 0.0154	0.2198 \pm 0.0172	0.1194 \pm 0.0051	0.2187 \pm 0.0125
X with 0.7σ	Iris	0.6500 \pm 0.0420	0.6087 \pm 0.0468	0.3927 \pm 0.0349	0.6120 \pm 0.0455
	COIL20	0.6220 \pm 0.0627	0.7662 \pm 0.0172	0.0893 \pm 0.0055	0.6803 \pm 0.0198
	MNIST	0.5317 \pm 0.0506	0.5763 \pm 0.0385	0.1377 \pm 0.0148	0.5181 \pm 0.0462
	CIFAR10	0.2202 \pm 0.0147	0.2205 \pm 0.0084	0.1209 \pm 0.0045	0.2134 \pm 0.0152

Table: Clustering accuracy (average over 10 trials) of FedSC with perturbed factors on COIL20.

		α_c			
		0	0.05	0.1	0.15
		Mean \pm Std	Mean \pm Std	Mean \pm Std	Mean \pm Std
α_z	0	0.7831 \pm 0.0268	0.6573 \pm 0.0291	0.6012 \pm 0.0391	0.4835 \pm 0.0526
	0.05	0.7824 \pm 0.0126	0.6573 \pm 0.0243	0.6165 \pm 0.0286	0.4998 \pm 0.0567
	0.1	0.7817 \pm 0.0299	0.6625 \pm 0.0258	0.6453 \pm 0.0186	0.5508 \pm 0.0415
	0.15	0.7881 \pm 0.0223	0.6526 \pm 0.0258	0.6219 \pm 0.0467	0.5901 \pm 0.0326

Conclusion

- Secure kernelized factorization method for federated spectral clustering on distributed data
- Theoretical guarantees for optimization convergence, correct clustering, and differential privacy
- Numerical experiments on synthetic and real image datasets

THANK YOU