

# Heterogeneous Federated Learning on a Graph

Huiyuan Wang, Xuyang Zhao and Wei Lin\*

School of Mathematical Sciences, Peking University

huiyuan.wang@pku.edu.cn, xuyangzhao@pku.edu.cn and weilin@math.pku.edu.cn

## Abstract

Federated learning, where algorithms are trained across multiple decentralized devices without sharing local data, is increasingly popular in distributed machine learning practice. Typically, a graph structure  $G$  exists behind local devices for communication. In this work, we consider parameter estimation in federated learning with *heterogeneity* in *communication* and *data distribution*, coupled with a *limited computational capacity* of local devices. We encode the distribution heterogeneity by parametrizing distributions on local devices with a set of distinct  $p$ -dimensional vectors. We then propose to jointly estimate parameters of all devices under the  *$M$ -estimation framework with the fused Lasso regularization*, encouraging an equal estimate of parameters on connected devices in  $G$ . We provide a general statistical guarantee for our estimator, which can be further calibrated to obtain convergence rates for various specific problem setups. Surprisingly, our estimator attains *the optimal rate under certain graph fidelity condition on  $G$ , as if we could aggregate all samples sharing the same distribution*. If the graph fidelity condition is not met, we propose an edge selection procedure via multiple testing to ensure the optimality. To ease the burden of local computation, a *decentralized stochastic version of ADMM is provided, with convergence rate  $O(T^{-1} \log T)$*  where  $T$  denotes the number of iterations. We highlight that our algorithm transmits only parameters along edges of  $G$  at each iteration, without requiring a central machine, which *preserves privacy*. To address communication heterogeneity, we further extend it to the case where devices are *randomly inaccessible* during the training process, with a similar algorithmic convergence guarantee. The computational and statistical efficiency of our method is evidenced by simulation experiments and the 2020 US presidential election data set.

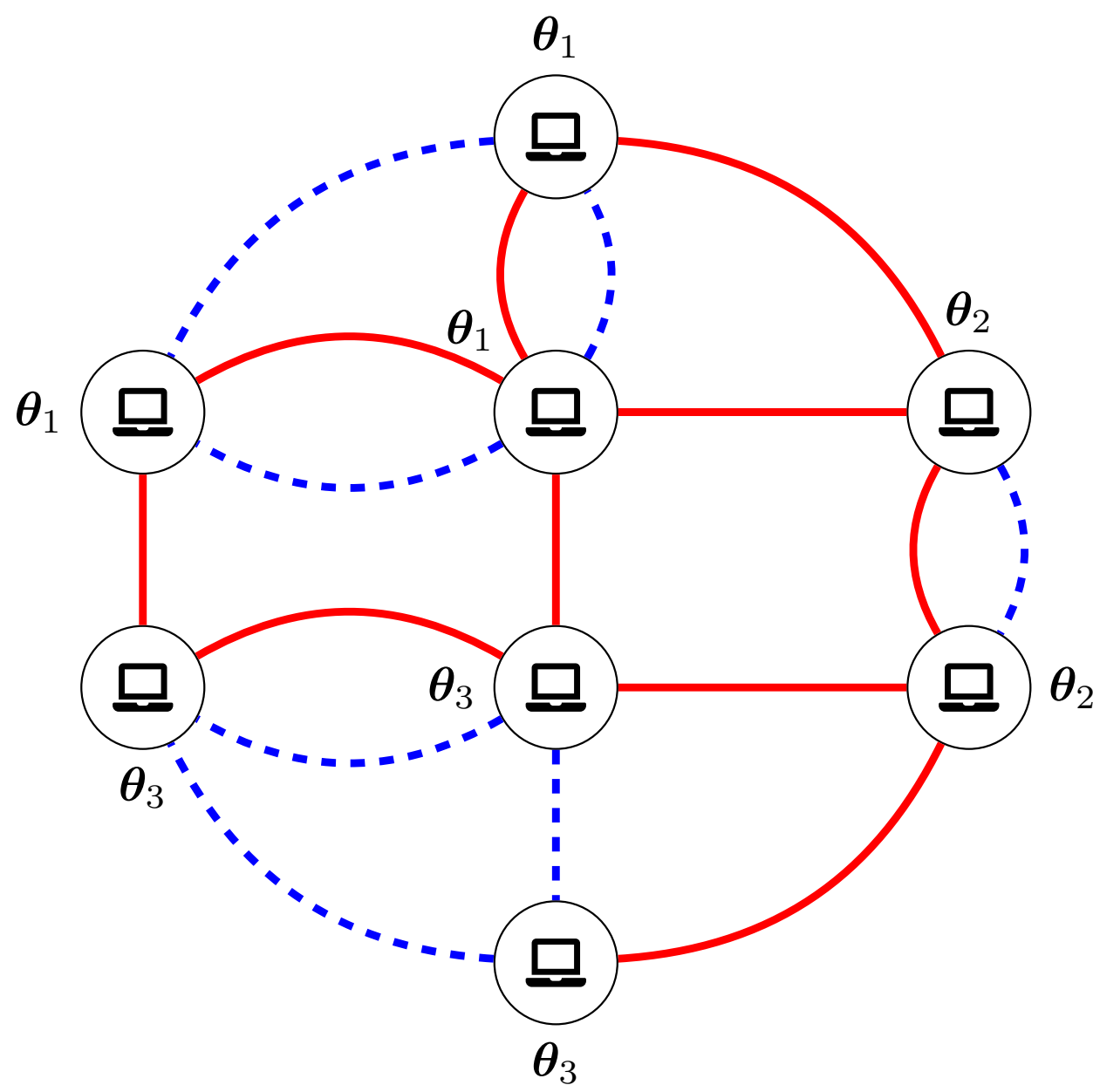
## 1. Problem setup

- (i) **Heterogeneous data distribution.** The target parameter  $\theta_u^*$  of the local device  $u$  is defined as the unique minimizer

$$\theta_u^* = \arg \min_{\theta \in \Xi} M_u(\theta) \equiv \mathbb{E}\{m_u(\mathbf{z}; \theta)\}, \quad \forall u \in V,$$

where the expectation is taken with respect to the distribution of local samples. To model distribution heterogeneity, we introduce the *characteristic graph*  $G_0$  such that device  $i$  and device  $j$  are connected in  $G_0$  *if and only if*  $\theta_i^* = \theta_j^*$ .

- (ii) **Heterogeneous communication among local devices.** To avoid the expensive cost of communicating between local devices and the central server, we consider *decentralized federated learning*; specifically, we assume that a *communication graph*  $G$  is given a priori, along whose edges local devices can transmit information with negligible costs. Here,  $G$  is required to be similar but not necessarily equal to  $G_0$ ; see Figure 1 for an illustration. Moreover, we allow *local devices to be off-line* randomly during the communication process.



**Figure 1:** Decentralized Federated Learning network illustrating communication (red solid lines) and distribution heterogeneity (blue dashed lines) among local devices with varying target parameters ( $\theta_1$ – $\theta_3$ ).

- (iii) **Limited computing power of local devices.** We consider the case where local devices have such a weak computing power that it will cost prohibitive time to derive local estimators. As such, our goal is to develop a *real-time* algorithm so that local devices are only required to perform simple iterations.

## 2. Methodology

To exploit the similarity between  $G$  and  $G_0$ , we propose the network fusion penalized estimator

$$\hat{\Theta} = (\hat{\theta}_u : u \in V) = \arg \min_{\Theta} F(\Theta) = \frac{1}{|V|} \sum_{u \in V} \widehat{M}_u(\theta_u) + \lambda R(\mathbf{D}\Theta), \quad (1)$$

where  $\Theta = (\theta_u : u \in V) \in \mathbb{R}^{|V| \times p}$ ,  $\mathbf{D}\Theta = (\theta_i - \theta_j : (i, j) \in E) \in \mathbb{R}^{|E| \times p}$ ,  $\lambda$  is a tuning parameter, and  $R(\mathbf{D}\Theta) = \sum_{(i, j) \in E} \phi(\theta_i - \theta_j)$  with  $\phi(\cdot)$  being a norm defined on  $\mathbb{R}^p$ . The fused lasso regularization  $R(\mathbf{D}\Theta)$  encourages an *equal* estimate for devices connected in  $G$ . When  $G$  deviates significantly from  $G_0$ , such an estimator may be erroneous. We then propose an edge selection procedure.

- (i) **Edge selection procedure.** Since  $(i, j) \in E_0 = E(G_0)$  is equivalent to  $\theta_i^* = \theta_j^*$ , we consider the simultaneous testing of the following null hypotheses  $H_{0,(i,j)} : \theta_i^* = \theta_j^*$  versus  $H_{1,(i,j)} : \theta_i^* \neq \theta_j^*$ ,  $(i, j) \in E$ . We construct a test statistic via  $\widehat{W}_{(i,j)} = \left\{ (\widehat{\theta}_i^{\text{loc}} - \widehat{\theta}_j^{\text{loc}})^T (\widehat{\Omega}_i + \widehat{\Omega}_j)^{-1} (\widehat{\theta}_i^{\text{loc}} - \widehat{\theta}_j^{\text{loc}}) \right\}^{1/2}$ , where  $\widehat{\Omega}_i = \{n_i \widehat{\mathbf{H}}_i(\widehat{\theta}_i^{\text{loc}})\}^{-1}$  denotes the asymptotic variance of  $\widehat{\theta}_i^{\text{loc}}$ . Adopting Bonferroni correction, we select  $E \cap E_0$  by

$$\widehat{E} = \{(i, j) \in E : |\widehat{W}_{(i,j)}|^2 \leq \chi_p^2(\alpha/|E|)\}, \quad (2)$$

where  $\chi_p^2(\alpha)$  is the upper  $\alpha$ -quantile of the  $\chi_p^2$  distribution. The edge selection procedure can eliminate wrong edges in  $E \setminus E_0$ , and thus ensure the global estimator in (1) achieve the *optimal* performance.

- (ii) **Decentralized real-time optimization procedure.** To solve the optimization problem (1), we consider the augmented Lagrangian defined as

$$\begin{aligned} L(\Theta, \mathbf{B}, \mathbf{A}) = & \frac{1}{|V|} \sum_{i \in V} \widehat{M}_i(\theta_i) + \lambda \sum_{(i,j) \in E} \phi(\beta_{ij} - \beta_{ji}) \\ & - \sum_{(i,j) \in E} \left\{ \alpha_{ij}^T (\theta_i - \beta_{ij}) + \alpha_{ji}^T (\theta_j - \beta_{ji}) \right\} + \frac{\rho}{2} \sum_{(i,j) \in E} \left\{ \|\theta_i - \beta_{ij}\|_2^2 + \|\theta_j - \beta_{ji}\|_2^2 \right\}, \end{aligned} \quad (3)$$

where  $\mathbf{B} = (\beta_{ij}, \beta_{ji} : (i, j) \in E)$  and  $\mathbf{A} = (\alpha_{ij}, \alpha_{ji} : (i, j) \in E)$ . To accommodate the *weak computing power of local devices*, we adopt one-step stochastic gradient update in the  $t$ -th iteration:

$$\theta_i(t+1) = \theta_i(t) - \eta(t) \left\{ \widetilde{\mathbf{g}}_i(t) + \rho \sum_{j \in N_i} (\theta_i(t) - \beta_{ij}(t) - \rho^{-1} \alpha_{ij}(t)) \right\}, \quad (4)$$

where  $\widetilde{\mathbf{g}}_i(t) = |\mathcal{B}_i(t)|^{-1} \sum_{b \in \mathcal{B}_i(t)} \psi_i(\mathbf{z}_b^{(i)}; \theta_i(t))$ ,  $\eta(t)$  denotes the learning rate,  $\mathcal{B}_i(t)$  denotes the mini-batch randomly sampled from  $\{\mathbf{z}_k^{(i)}\}_{k=1}^{n_i}$  on device  $i$  in the  $t$ -th iteration, and  $\widetilde{\mathbf{g}}_i(t)$  is an unbiased estimator of  $\nabla_{\theta} M_i(\theta)$  evaluated at  $\theta_i(t)$ . Except for local samples, the update equation (4) only requires  $\beta_{ij}(t)$  and  $\alpha_{ij}(t)$  which can be transmitted from device  $j$ . Thus, (4) *can be executed in parallel for all devices*. With  $\theta_i(t+1)$ ,  $i \in V$  at hand, we then update  $\beta_{ij}(t)$  and  $\beta_{ji}(t)$  by

$$\begin{aligned} \begin{pmatrix} \beta_{ij}(t+1) \\ \beta_{ji}(t+1) \end{pmatrix} = & \arg \min_{\beta_{ij}, \beta_{ji}} \left\{ \lambda \phi(\beta_{ij} - \beta_{ji}) \right. \\ & \left. + \frac{\rho}{2} \left( \|\theta_i(t+1) - \beta_{ij} - \rho^{-1} \alpha_{ij}(t)\|_2^2 + \|\theta_j(t+1) - \beta_{ji} - \rho^{-1} \alpha_{ji}(t)\|_2^2 \right) \right\}. \end{aligned} \quad (5)$$

The update equation (5) can be implemented on either device  $i$  or device  $j$ , as long as  $(\theta_j(t+1), \beta_{ji}(t), \alpha_{ji}(t))$  or  $(\theta_i(t+1), \beta_{ij}(t), \alpha_{ij}(t))$  is transmitted to the corresponding device. *For specific choice of  $\phi(\cdot)$ , for example,  $\phi(\cdot) = \|\cdot\|_1$  and  $\phi(\cdot) = \|\cdot\|_2$ , we can obtain an explicit update equation from (5).* Finally, we update  $\alpha_{ij}(t)$  and  $\alpha_{ji}(t)$  by

$$\begin{pmatrix} \alpha_{ij}(t+1) \\ \alpha_{ji}(t+1) \end{pmatrix} = \begin{pmatrix} \alpha_{ij}(t) \\ \alpha_{ji}(t) \end{pmatrix} - \rho \begin{pmatrix} \theta_j(t+1) - \beta_{ji}(t+1) \\ \theta_i(t+1) - \beta_{ij}(t+1) \end{pmatrix}. \quad (6)$$

Notice that update equation (6) also only requires parameter communication among connected devices. Both (5) and (6) can be performed in parallel across edges.

*We refer to (4) as node optimization step, and (5) and (6) as edge communication step.*

## 3. Assumptions and theoretic guarantees

We maintain the following regularity assumptions.

- (A1) **Identifiability.**  $\mathbb{E}\{m_u(\mathbf{z}; \theta)\}$  is convex and twice differentiable with respect to  $\theta$ , the Hessian matrix  $\mathbf{H}_u(\theta) = \nabla_{\theta}^2 \mathbb{E}\{m_u(\mathbf{z}; \theta)\}$  is Lipschitz continuous at  $\theta_u^*$ .
- (A2) **Bounded conditional number.** The conditional number of  $\widehat{\mathbf{H}}_u(\theta)$  is bounded by  $\kappa$ .
- (A3) **Compatibility factor.** For  $S = E \setminus E_0 \neq \emptyset$ ,  $\kappa_S(\mathbf{D}) \equiv \inf_{\Theta} \frac{\sqrt{|S|} \|\Theta\|_F}{R(\mathbf{D}\Theta)_S} \geq \kappa_0 > 0$ , where  $R\{(\mathbf{D}\Theta)_S\} = \sum_{(u,v) \in S} \phi(\theta_u - \theta_v)$ .

Based on Assumptions (A1)–(A3), we obtain the following *non-asymptotic bounds* regarding both statistical and algorithmic errors.

- (i) **Deterministic estimation bound.** For noises from *any* distribution, the penalized  $M$ -estimator  $\widehat{\Theta}$  satisfies

$$\frac{1}{|V|} \|\widehat{\Theta} - \Theta^*\|_F^2 \leq 2\kappa^2 \left( \rho^2 + \frac{4|S|}{\kappa_0} \lambda^2 \right),$$

where  $\rho = \frac{1}{\sqrt{|V|}} \|\Pi_{\text{Ker}(\mathbf{D})} \widehat{\Psi}(\Theta^*)\|_F$ , and  $\lambda = \frac{1}{\sqrt{|V|}} R^*\{(\mathbf{D}^+)^T \widehat{\Psi}(\Theta^*)\}$ ,  $\widehat{\Psi}$  is the gradient of the empirical risk function and  $R^*(\cdot)$  denotes the Fréchet dual of  $R(\cdot)$ .

- (ii) **Probabilistic estimation bound.** If the score function  $\psi_u(\mathbf{z}_i^{(u)}; \theta_u^*) = \nabla_{\theta} m_u(\mathbf{z}_i^{(u)}; \theta_u^*)$  is sub-Gaussian with parameter  $\sigma^2$ , then we obtain

$$\frac{1}{|V|} \|\widehat{\Theta} - \Theta^*\|_F^2 = O_p \left\{ \frac{\sigma^2}{\kappa_0} \left( \frac{pK(G)}{n|V|} + \frac{p|E \setminus E_0|}{n|V|} \right) \right\},$$

where  $K(G)$  is the number of connected components of  $G$ . We remark that the optimal estimation error scales as  $O_p \left( \frac{pK(E_0)}{n|V|} \right)$ .

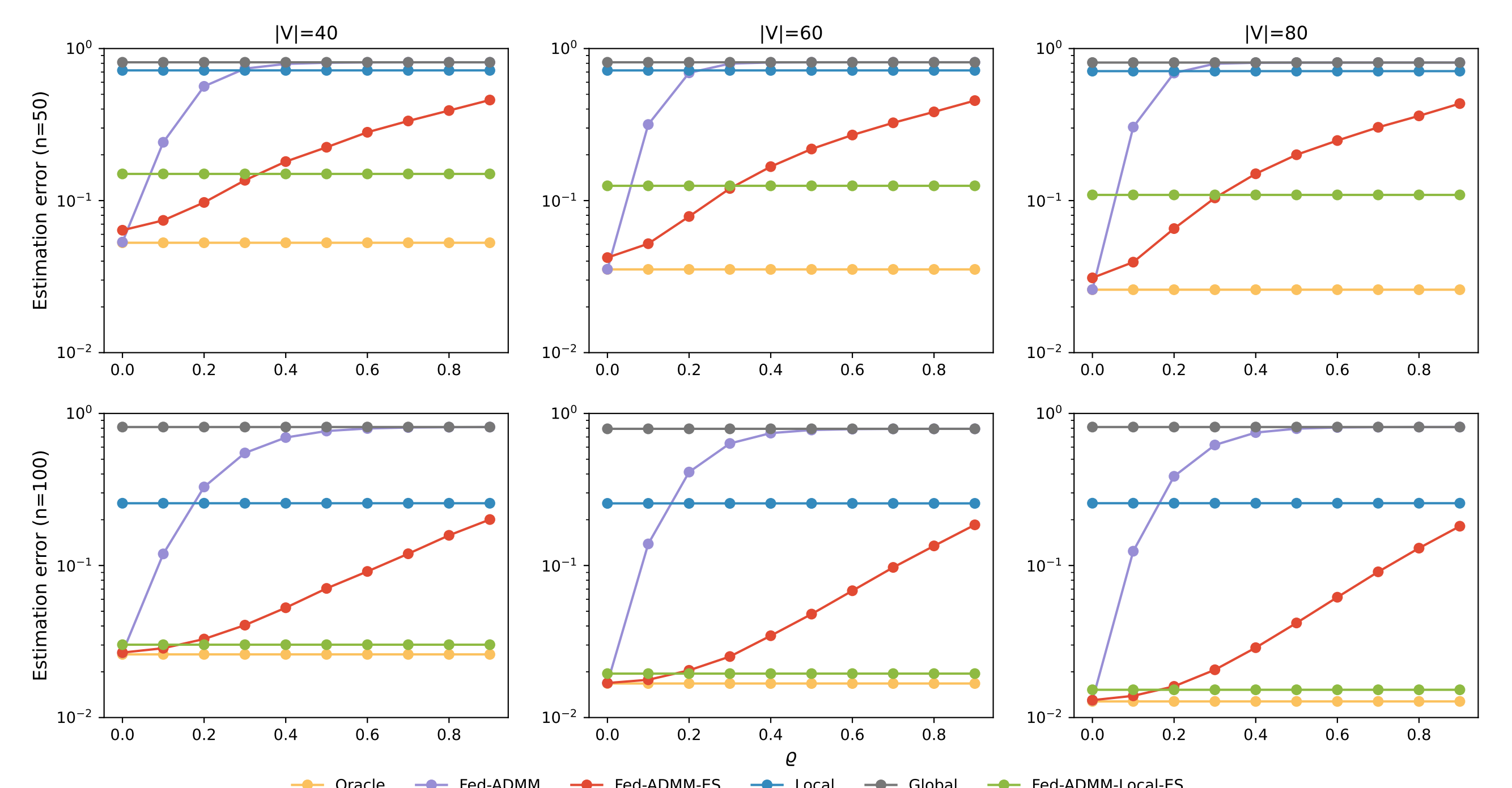
- (iii) **Edge selection ensures the optimality of the network fusion estimator.** Under mild conditions, the output of the edge selection procedure maximizes the *graph fidelity*  $\text{GF}_{G_0}(G) \equiv \frac{K(E_0)}{K(E) + |E \setminus E_0|}$ .

- (iv) **Algorithmic convergence.** Under mild conditions, by choosing  $\eta(t) = \kappa/t$ , we have

$$\frac{1}{|V|} \mathbb{E} \left( \|\widehat{\Theta} - \Theta^*\|_F^2 \mid \mathbf{z}_k^{(u)}, 1 \leq k \leq n_u, u \in V \right) \leq \frac{2\kappa^2 C_{\psi} \log T}{T},$$

for sufficiently large  $T$  such that  $\kappa C_{\psi} |V| \log T \geq C|E|$  for some constant  $C$ , where the expectation is taken with respect to the choice of mini-batches  $\{\mathcal{B}_u(t) : u \in V, 1 \leq t \leq T\}$ .

## 4. Selected simulation



**Figure 2:** Sensitivity to graph corruption.  $G$  is obtained by randomly flipping the status of edge in  $G_0$  with probability  $q$ . Fed-ADMM-ES is obtained by performing the edge selection procedure on  $G$  and Fed-ADMM-Local-ES is obtained by performing the edge selection procedure on fully connected graph.