

Byzantine-Resilient Federated Alternating Gradient Descent and Minimization for Partly-Decoupled Low Rank Matrix Learning

Problem Setting

Learn a low rank $r \ll n, q$ matrix $\Theta^* \in \mathbb{R}^{n \times q}$ from measurements of the form

$$\mathbf{y}_k := \mathbf{X}_k \boldsymbol{\theta}_k^*, k \in [q].$$

In a federated setting, while being resilient to **Byzantine Attacks**.

Factor $\Theta = \mathbf{U}\mathbf{B}$. Solving this problem requires solving (AltGDmin [Nayer and Vaswani 2022] and FedRep [Collins et al. 2021]),

$$\min_{\substack{\tilde{\mathbf{U}} \in \mathbb{R}^{n \times r} \\ \tilde{\mathbf{B}} \in \mathbb{R}^{r \times q}}} f(\tilde{\mathbf{U}}, \tilde{\mathbf{B}}) = \min_{\substack{\tilde{\mathbf{U}} \in \mathbb{R}^{n \times r} \\ \tilde{\mathbf{B}} \in \mathbb{R}^{r \times q}}} \sum_{k=1}^q \|\mathbf{y}_k - \mathbf{X}_k \tilde{\mathbf{U}} \tilde{\mathbf{b}}_k\|^2$$

Theorem 1: Byz-Fed-AltGDmin-Learn

Bounded heterogeneity:

$$\max_{\ell, \ell' \in [L]} \|\mathbf{B}_\ell^* - \mathbf{B}_{\ell'}^*\|_F^2 \leq G_B^2 \sigma_{\max}^{*2}$$

Assume RSV incoherence, Bounded heterogeneity Assumption holds, and $\frac{L_{\text{byz}}}{L} < 0.4$. If

$$n\tilde{q}p \geq C\tilde{\kappa}^{10} \mu^2 \tilde{q} r^2 \log \tilde{q} \log\left(\frac{1}{\epsilon}\right)$$

then, w.h.p. after $T = C\tilde{\kappa}^2 \log\left(\frac{1}{\epsilon}\right)$ iterations,

$$SD_F(\mathbf{U}^*, \mathbf{U}_T) \leq \max(\epsilon, 14C\tilde{\kappa}^2 G_B)$$

Experiments

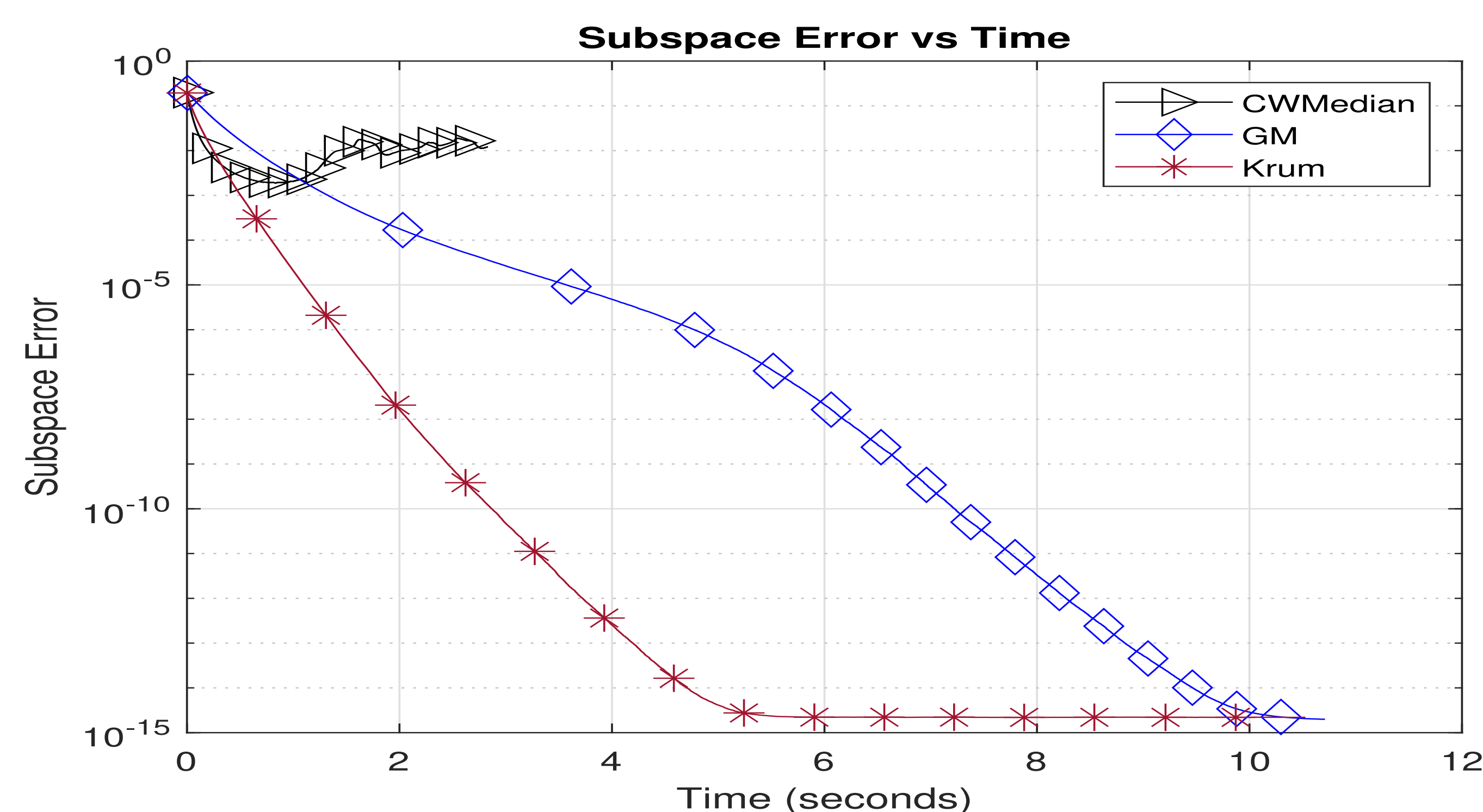


Figure 1. LRMC experiment.

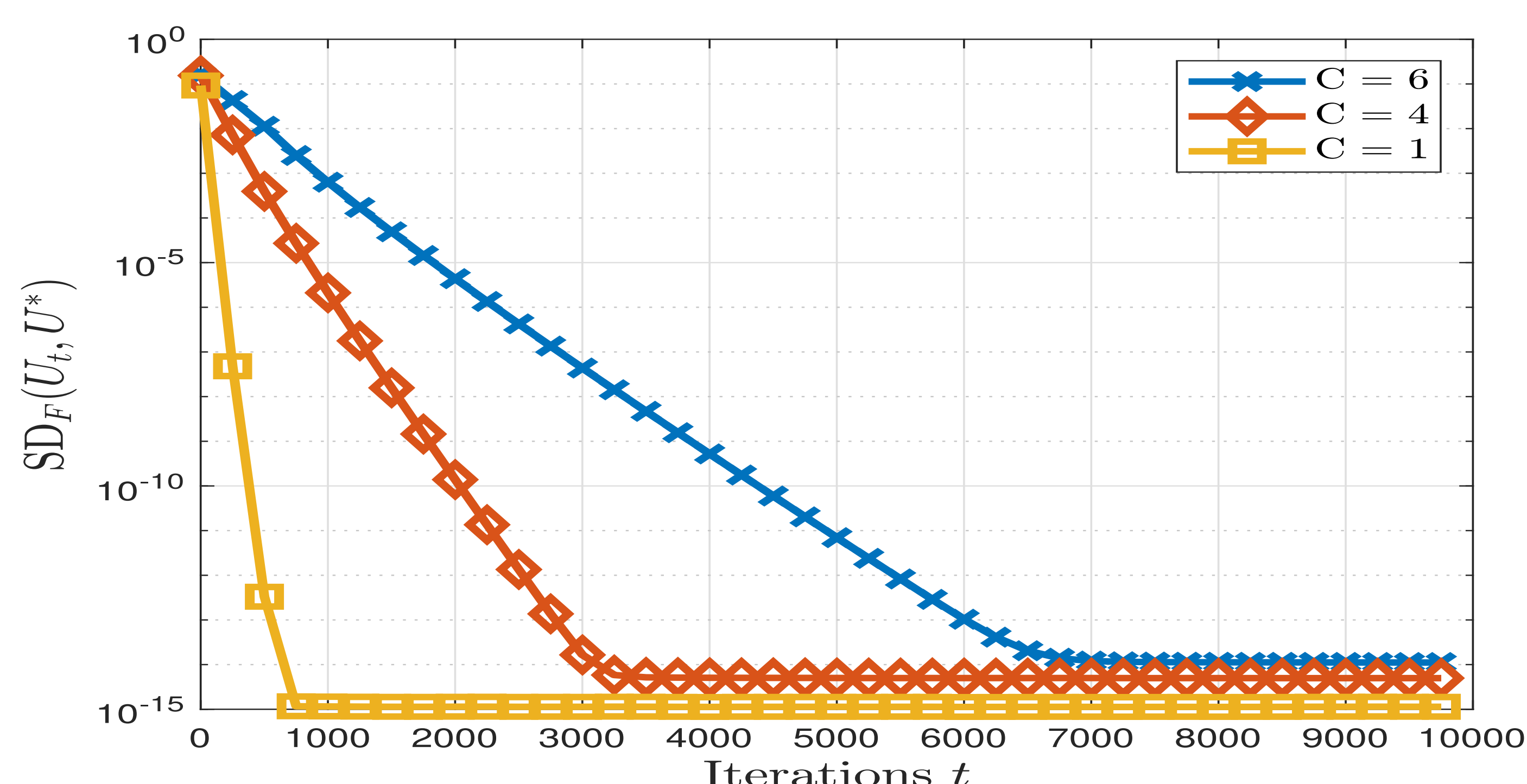


Figure 2. Heterogeneity Effect.

Applications

The way matrix \mathbf{X}_k is defined it gives variety of differently problems:

- **LoRA** technique for **large language models** PEFT, \mathbf{X}_k being identity matrices.
- For **collaborative filtering in recommendation systems** a.k.a (LRMC) problem it is a diagonal 1-0 matrix.
- For **linear multi task learning** a.k.s LRCS problem it is a random Gaussian matrix.
- For **nonlinear multi task representation** a.k.s LRPR problem it is a random Gaussian matrix, but only the magnitudes of the measurements are observed i.e., $\mathbf{z}_k = |\mathbf{y}_k|$.
- **Initialization:** All these problems require initialization, which reduces to a subspace estimation or PCA problem.

Byz-AltGDmin-LRMC algorithm

AltGDmin Initialization:

Nodes $\ell = 1, \dots, L$

Calculate and Push $\mathbf{U}_{0\ell}$ to center

Central Server

Define set $\mathcal{I}_0 = \{\}$

Key Idea 1: Subspace Median with filtering

for $\ell = 1$ to L **do**

if $\|\mathbf{u}_{0\ell}^j\| \leq 1.5\mu\sqrt{\frac{r}{n}}$ for all $j \in [n]$ **then**

Add ℓ to set \mathcal{I}_0

end for

$\mathbf{U}_0 \leftarrow \text{Byz-SubspaceEstimation}\{\mathbf{U}_{0\ell}\}_{\ell \in \mathcal{I}_0}$

Push \mathbf{U}_0 to nodes.

AltGDmin Iterations:

for $t = 1$ to T **do**

Nodes $\ell = 1, \dots, L$

Calculate and Push ∇_ℓ to center

Central Server

Define set $\mathcal{I}_t = \{\}$

Key Idea 2: Filtering with robust aggregation

for $\ell = 1$ to L **do**

Compute $\mathbf{U}_{\text{temp}} \leftarrow \mathbf{U}_{t-1} - \eta \nabla_\ell$

if $\|\mathbf{u}_{\text{temp}}^j\| \leq (1 - \frac{0.4}{\tilde{\kappa}^2})\|\mathbf{u}_{t-1}^j\| + 1.4\mu\sqrt{\frac{r}{n}}$ for all $j \in [n]$ **then**

Add ℓ to set \mathcal{I}_t

end for

$\nabla_{Kr/GM} = \text{Krum/GM}\{\nabla_\ell\}_{\ell \in \mathcal{I}_t}$

Compute $\mathbf{U}_t \leftarrow QR(\mathbf{U}_{t-1} - \eta \nabla_{Kr/GM})$

Push \mathbf{U}_t to nodes.

end for

Output \mathbf{U}_T .