Recursive Importance Sketching for Rank Constrained Least Squares: Algorithms and High-order Convergence

#### Wen Huang

Joint meeting of Wuhan University and Xiamen University

1



Yuetian Luo UW-Madison



Xudong Li Fudan University



Anru Zhang UW-Madison

Thank Yuetian Luo for creating these slides.

## **Problem of Interest**

$$\begin{split} \min_{\mathbf{X} \in \mathbb{R}^{p_1 \times p_2}} f(\mathbf{X}) &:= \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2, \quad \text{subject to} \quad \operatorname{rank}(\mathbf{X}) = r, \\ \text{where } \mathbf{y} \in \mathbb{R}^n, \mathcal{A}(\mathbf{X}) = [\langle \mathbf{A}_1, \mathbf{X} \rangle, \dots, \langle \mathbf{A}_n, \mathbf{X} \rangle]^\top. \end{split}$$

#### Problem of Interest

$$\begin{split} \min_{\mathbf{X} \in \mathbb{R}^{p_1 \times p_2}} f(\mathbf{X}) &:= \frac{1}{2} \| \mathbf{y} - \mathcal{A}(\mathbf{X}) \|_2^2, \quad \text{subject to} \quad \operatorname{rank}(\mathbf{X}) = r, \\ \text{where } \mathbf{y} \in \mathbb{R}^n, \mathcal{A}(\mathbf{X}) = [\langle \mathbf{A}_1, \mathbf{X} \rangle, \dots, \langle \mathbf{A}_n, \mathbf{X} \rangle]^\top. \end{split}$$

#### Motivation: Low rank matrix recovery

• Observe  $\mathbf{y}, \mathcal{A}$  from  $\mathbf{y} = \mathcal{A}(\mathbf{X}^*) + \epsilon$ . Goal: recover  $\mathbf{X}^*$  from  $\mathbf{y}, \mathcal{A}$ 

#### Specific problems:

- Matrix regression: A<sub>i</sub> <sup>i.i.d.</sup> ∼ N(0, 1) [Candès and Plan, 2011, Recht et al., 2010]
- Matrix Completion: A<sub>i</sub> has one entry to be 1, others are 0 [Candès and Tao, 2010]
- Phase retrieval:  $\mathbf{A}_i = \mathbf{a}_i \mathbf{a}_i^{\top}$  [Shechtman et al., 2015]
- Rank-one sensing: A<sub>i</sub> = a<sub>i</sub>b<sub>i</sub><sup>⊤</sup>
   [Cai and Zhang, 2015, Chen et al., 2015]

# **Prior Work**

• Convex relaxation:  $\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_{2}^{2} + \lambda \|\mathbf{X}\|_{*}$ 

[Recht et al., 2010, Candès and Plan, 2011]

#### Theoretical properties 🖌 computation can be intensive

- Non-convex methods: enforce rank r constraint
  - Factorize X = RL<sup>⊤</sup> + Gradient descent or Alternating Minimization on R ∈ ℝ<sup>p1×r</sup>, L ∈ ℝ<sup>p2×r</sup>

[Ma et al., 2019, Park et al., 2018, Sun and Luo, 2015, Tu et al., 2016, Wang et al., 2017,

Zhao et al., 2015, Zheng and Lafferty, 2015, Jain et al., 2013, Hardt, 2014]...

 Projected gradient descent (Singular value projection (SVP), Iterative Hard Thresholding (IHT))

[Goldfarb and Ma, 2011, Jain et al., 2010, Tanner and Wei, 2013]...

Manifold optimization [Boumal and Absil, 2011, Keshavan et al., 2009,

Vandereycken, 2013, Wei et al., 2016, Huang and Hand, 2018]

• ...

# **Prior Work**

• Convex relaxation:  $\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_{2}^{2} + \lambda \|\mathbf{X}\|_{*}$ 

[Recht et al., 2010, Candès and Plan, 2011]

#### Theoretical properties 🖌 computation can be intensive

- Non-convex methods: enforce rank r constraint
  - Factorize  $\mathbf{X} = \mathbf{R}\mathbf{L}^{\top}$  + Gradient descent or Alternating Minimization on  $\mathbf{R} \in \mathbb{R}^{p_1 \times r}$ ,  $\mathbf{L} \in \mathbb{R}^{p_2 \times r}$

[Ma et al., 2019, Park et al., 2018, Sun and Luo, 2015, Tu et al., 2016, Wang et al., 2017,

Zhao et al., 2015, Zheng and Lafferty, 2015, Jain et al., 2013, Hardt, 2014]...

 Projected gradient descent (Singular value projection (SVP), Iterative Hard Thresholding (IHT))

[Goldfarb and Ma, 2011, Jain et al., 2010, Tanner and Wei, 2013]...

Manifold optimization [Boumal and Absil, 2011, Keshavan et al., 2009,

Vandereycken, 2013, Wei et al., 2016, Huang and Hand, 2018]

• ...

#### Most of existing algorithms

- require careful tuning or
- have a convergence rate no faster than linear.

 $\implies$  Can we do better?

Meeting of WHU-XMU

<u>Recursive</u> <u>Importance</u> <u>Sketching</u> algorithm for <u>Rank</u> constrained least squares <u>Optimization</u> (RISRO).

#### Advantages

- Tuning free
- High-order convergence guarantees under proper assumptions

- Input **y**, *A*, and initialization **X**<sup>0</sup> with (economic) SVD **U**<sup>0</sup>Σ<sup>0</sup>V<sup>0</sup><sup>T</sup>
   For *t* = 0, 1, ...
  - Perform importance sketching on A.



Update sketching matrices.

- Input **y**, A, and initialization **X**<sup>0</sup> with (economic) SVD **U**<sup>0</sup>Σ<sup>0</sup>V<sup>0</sup><sup>T</sup>
   For t = 0, 1, ...
  - Perform importance sketching on  $\mathcal{A}$ . Construct importance covariates  $\mathbf{A}_i^B := \mathbf{U}^{t\top} \mathbf{A}_i \mathbf{V}^t, \mathbf{A}_i^{D_1} := \mathbf{U}_{\perp}^{t\top} \mathbf{A}_i \mathbf{V}^t, \mathbf{A}_i^{D_2} := \mathbf{U}^{t\top} \mathbf{A}_i \mathbf{V}_{\perp}^t$



Solve a dimension reduced least squares.

Update sketching matrices.

- Input **y**, *A*, and initialization **X**<sup>0</sup> with (economic) SVD **U**<sup>0</sup>Σ<sup>0</sup>V<sup>0</sup><sup>T</sup>
   For *t* = 0, 1, ...
  - Perform importance sketching on  $\mathcal{A}$ . Construct importance covariates  $\mathbf{A}_i^B := \mathbf{U}^{t\top} \mathbf{A}_i \mathbf{V}^t, \mathbf{A}_i^{D_1} := \mathbf{U}_{\perp}^{t\top} \mathbf{A}_i \mathbf{V}^t, \mathbf{A}_i^{D_2} := \mathbf{U}^{t\top} \mathbf{A}_i \mathbf{V}_{\perp}^t$



Solve a dimension reduced least squares.

$$(\mathbf{B}^{t+1}, \mathbf{D}_1^{t+1}, \mathbf{D}_2^{t+1}) = \operatorname*{arg\,min}_{\mathbf{B}, \mathbf{D}_1, \mathbf{D}_2} \sum_{i=1}^n \left( \mathbf{y}_i - \langle \mathbf{A}_i^{\mathcal{B}}, \mathbf{B} \rangle - \langle \mathbf{A}_i^{\mathcal{D}_1}, \mathbf{D}_1 \rangle - \langle \mathbf{A}_i^{\mathcal{D}_2}, \mathbf{D}_2^{\mathsf{T}} \rangle \right)^2$$

Update sketching matrices.

- Input **y**, A, and initialization **X**<sup>0</sup> with (economic) SVD **U**<sup>0</sup>**\Sigma**<sup>0</sup>**V**<sup>0</sup><sup> $\top$ </sup> For t = 0, 1, ...
  - Perform importance sketching on A. Construct importance covariates  $\mathbf{A}_{i}^{B} := \mathbf{U}^{t\top}\mathbf{A}_{i}\mathbf{V}^{t}, \mathbf{A}_{i}^{D_{1}} := \mathbf{U}_{\perp}^{t\top}\mathbf{A}_{i}\mathbf{V}^{t}, \mathbf{A}_{i}^{D_{2}} := \mathbf{U}^{t\top}\mathbf{A}_{i}\mathbf{V}_{\perp}^{t}$



Solve a dimension reduced least squares.

$$(\mathbf{B}^{t+1}, \mathbf{D}_1^{t+1}, \mathbf{D}_2^{t+1}) = \operatorname*{arg\,min}_{\mathbf{B}, \mathbf{D}_1, \mathbf{D}_2} \sum_{i=1}^{n} \left( \mathbf{y}_i - \langle \mathbf{A}_i^B, \mathbf{B} \rangle - \langle \mathbf{A}_i^{D_1}, \mathbf{D}_1 \rangle - \langle \mathbf{A}_i^{D_2}, \mathbf{D}_2^\top \rangle \right)^2$$

Update sketching matrices. Let  $\mathbf{X}_{tt}^{t+1} = (\mathbf{U}^t \mathbf{B}^{t+1} + \mathbf{U}_1^t \mathbf{D}_1^{t+1}),$  $\mathbf{X}_{V}^{t+1} = (\mathbf{V}^{t}\mathbf{B}^{t+1\top} + \mathbf{V}_{\perp}^{t}\mathbf{D}_{2}^{t+1}). \text{ Update} \\ \mathbf{U}^{t+1} = QR(\mathbf{X}_{U}^{t+1}), \mathbf{V}^{t+1} = QR(\mathbf{X}_{V}^{t+1}).$ 

$$\blacksquare \text{ (Optional) } \mathbf{X}^{t+1} = \mathbf{X}_U^{t+1} \left( \mathbf{B}^{t+1} \right)^{\dagger} \mathbf{X}_V^{t+1}$$

 $QR(\cdot)$  is the Q part in QR decomposition and  $(\cdot)^{\dagger}$  is the Moore-Penrose inverse Meeting of WHU-XMU

## **RISRO-Intuition**

Suppose  $\mathbf{y}_i = \langle \mathbf{A}_i, \overline{\mathbf{X}} \rangle + \overline{\epsilon}_i$  where  $\overline{\mathbf{X}}$  is a rank *r* target matrix. Rewritten

 $\mathbf{y}_i = \langle \mathbf{A}_i^B, \mathbf{U}^{t\top} \bar{\mathbf{X}} \mathbf{V}^t \rangle + \langle \mathbf{A}_i^{D_1}, \mathbf{U}_{\perp}^{t\top} \bar{\mathbf{X}} \mathbf{V}^t \rangle + \langle \mathbf{A}_i^{D_2}, \mathbf{U}^{t\top} \bar{\mathbf{X}} \mathbf{V}_{\perp}^t \rangle + \boldsymbol{\epsilon}_i^t,$ 

where  $\boldsymbol{\epsilon}_{\boldsymbol{i}}^{t} = \langle \boldsymbol{\mathsf{U}}_{\perp}^{t\top} \boldsymbol{\mathsf{A}}_{\boldsymbol{i}} \boldsymbol{\mathsf{V}}_{\perp}^{t}, \boldsymbol{\mathsf{U}}_{\perp}^{t\top} \bar{\boldsymbol{\mathsf{X}}} \boldsymbol{\mathsf{V}}_{\perp}^{t} \rangle + \bar{\boldsymbol{\epsilon}}_{\boldsymbol{i}}.$ 

## **RISRO-Intuition**

Suppose  $\mathbf{y}_i = \langle \mathbf{A}_i, \overline{\mathbf{X}} \rangle + \overline{\epsilon}_i$  where  $\overline{\mathbf{X}}$  is a rank *r* target matrix. Rewritten

 $\mathbf{y}_{i} = \langle \mathbf{A}_{i}^{B}, \mathbf{U}^{t\top} \bar{\mathbf{X}} \mathbf{V}^{t} \rangle + \langle \mathbf{A}_{i}^{D_{1}}, \mathbf{U}_{\perp}^{t\top} \bar{\mathbf{X}} \mathbf{V}^{t} \rangle + \langle \mathbf{A}_{i}^{D_{2}}, \mathbf{U}^{t\top} \bar{\mathbf{X}} \mathbf{V}_{\perp}^{t} \rangle + \boldsymbol{\epsilon}_{i}^{t},$ 

where  $\boldsymbol{\epsilon}_{i}^{t} = \langle \boldsymbol{\mathsf{U}}_{\perp}^{t\top} \boldsymbol{\mathsf{A}}_{i} \boldsymbol{\mathsf{V}}_{\perp}^{t}, \boldsymbol{\mathsf{U}}_{\perp}^{t\top} \bar{\boldsymbol{\mathsf{X}}} \boldsymbol{\mathsf{V}}_{\perp}^{t} \rangle + \bar{\boldsymbol{\epsilon}}_{i}.$ 

If  $\epsilon^t = 0$ . Then

 $\mathbf{B}^{t+1} = \mathbf{U}^{t\top} \bar{\mathbf{X}} \mathbf{V}^t, \quad \mathbf{D}_1^{t+1} = \mathbf{U}_{\perp}^{t\top} \bar{\mathbf{X}} \mathbf{V}^t, \quad \mathbf{D}_2^{t+1} = \mathbf{U}^{t\top} \bar{\mathbf{X}} \mathbf{V}_{\perp}^t$ 

is a solution of the least squares. Moreover if  $\mathbf{B}^{t+1}$  is invertible

$$\mathbf{X}^{t+1} = \mathbf{X}_U^{t+1} \left( \mathbf{B}^{t+1} 
ight)^{-1} \mathbf{X}_V^{t+1 op} = \mathbf{\overline{X}}$$

In general  $\epsilon^t \neq 0$ , but we hope  $\mathbf{X}^t \rightarrow \overline{\mathbf{X}}$ .

# Importance Sketching in RISRO

Sketching: do dimension reduction to speed up the computation



 Comparison of Importance Sketching and Randomized Sketching

	Importance Sketching	Randomized Sketching	
	Importance Sketching	[Mahoney, 2011, Woodruff, 2014]	
Sketching Matrix	Deterministic, $\mathbf{U}^t$ , $\mathbf{V}^t$ (with supervision)	Random	
Dimension reduction	Reduce <i>p</i> , hold <i>n</i>	Reduce <i>n</i> , hold <i>p</i>	
Statistical efficiency	High	Low	

 Alternating Minimization (Alter Mini) [Jain et al., 2013, Zhao et al., 2015]

$$\widehat{\mathbf{V}}^{t+1} = \underset{\mathbf{V}\in\mathbb{R}^{p_{2}\times r}}{\arg\min} \sum_{i=1}^{n} \left( \mathbf{y}_{i} - \langle \mathbf{A}_{i}, \mathbf{U}^{t}\mathbf{V}^{\top} \rangle \right)^{2} = \underset{\mathbf{V}\in\mathbb{R}^{p_{2}\times r}}{\arg\min} \sum_{i=1}^{n} \left( \mathbf{y}_{i} - \langle \mathbf{U}^{t\top}\mathbf{A}_{i}, \mathbf{V}^{\top} \rangle \right)^{2},$$
$$\mathbf{V}^{t+1} = \operatorname{QR}(\widehat{\mathbf{V}}^{t+1})$$



• Alternating Minimization (Alter Mini) [Jain et al., 2013, Zhao et al., 2015]

 Rank 2r iterative least squares (R2RILS) for matrix completion [Bauch and Nadler, 2020]

$$\min_{\mathbf{M} \in \mathbb{R}^{p_1 \times r}, \mathbf{N} \in \mathbb{R}^{p_2 \times r}} \sum_{(i,j) \in \Omega} \left\{ \left( \mathbf{U}^t \mathbf{N}^\top + \mathbf{M} \mathbf{V}^{t\top} - \mathbf{Y} \right)_{[i,j]} \right\}^2,$$

 $\boldsymbol{\Omega}$  is the observed entry indices.

• Alternating Minimization (Alter Mini) [Jain et al., 2013, Zhao et al., 2015]

 Rank 2r iterative least squares (R2RILS) for matrix completion [Bauch and Nadler, 2020]

$$\sum_{(i,j)\in\Omega} \left\{ \left( \mathbf{U}^{t}\mathbf{N}^{\top} + \mathbf{M}\mathbf{V}^{t\top} - \mathbf{Y} \right)_{[i,j]} \right\}^{2} \iff \sum_{(i,j)\in\Omega} \left( \mathbf{Y}_{[i,j]} - \langle \mathbf{U}^{t\top}\mathbf{A}^{ij}, \mathbf{N}^{\top} \rangle - \langle \mathbf{M}, \mathbf{A}^{ij}\mathbf{V}^{t} \rangle \right)$$

$$A_{i}$$

$$A_{i$$



- Alter Mini: Miss one set of covariates ⇒ large iteration error
- R2RILS: Double core sketch  $\Rightarrow \begin{cases} Rank \text{ deficiency in the least squares} \\ Hard in theory and implementation} \end{cases}$
- ★ RISRO: resolve both issues ⇒ High-order convergence!

#### **Convergence Analysis**



Meeting of WHU-XMU

## **Convergence Analysis: Assumption**

 Restricted isometry property (RIP) [Candès, 2008, Recht et al., 2010] :

 $\mathcal{A}: \mathbb{R}^{p_1 \times p_2} \to \mathbb{R}^n$  satisfies the *r*-RIP with *r*-restricted isometry constant  $\delta \in [0, 1)$  if

 $(1-\delta) \|\mathbf{Z}\|_{\mathbf{F}}^2 \le \|\mathcal{A}(\mathbf{Z})\|_2^2 \le (1+\delta) \|\mathbf{Z}\|_{\mathbf{F}}^2$ 

holds for all **Z** of rank at most *r*.

- Widely used [Cai and Zhang, 2013, Candès and Plan, 2011]...
- Easy to satisfy in random subgaussian design [Candès and Plan, 2011]
- Key quantity in landscape analysis [Bhojanapalli et al., 2016, Ge et al., 2017, Uschmajew and Vandereycken, 2018]

## **Convergence Analysis: Assumption**

• Restricted isometry property (RIP)

[Candès, 2008, Recht et al., 2010] :

 $\mathcal{A}: \mathbb{R}^{p_1 \times p_2} \to \mathbb{R}^n \text{ satisfies the } r-\text{RIP with } r\text{-restricted isometry constant } \delta \in [0, 1) \text{ if }$ 

 $(1-\delta) \|\mathbf{Z}\|_{\mathbf{F}}^2 \le \|\mathcal{A}(\mathbf{Z})\|_2^2 \le (1+\delta) \|\mathbf{Z}\|_{\mathbf{F}}^2$ 

holds for all **Z** of rank at most *r*.

- Widely used [Cai and Zhang, 2013, Candès and Plan, 2011]...
- Easy to satisfy in random subgaussian design [Candès and Plan, 2011]
- Key quantity in landscape analysis [Bhojanapalli et al., 2016, Ge et al., 2017, Uschmajew and Vandereycken, 2018]
   Caveat: RIP may not hold in some scenarios such as matrix completion.

Let  $\overline{\mathbf{X}}$  be a rank *r* stationary point and  $\overline{\epsilon} := \mathbf{y} - \mathcal{A}(\overline{\mathbf{X}})$ . Assume

- $\mathcal{A}$  satisfies 3*r*-restricted isometry property (RIP) with RIP constant  $\delta$
- Initialization condition:  $\|\mathbf{X}^0 \bar{\mathbf{X}}\|_{\mathbf{F}} \leq C(\delta)\sigma_r(\bar{\mathbf{X}})$
- Small residual (gradient) condition:  $\|\mathcal{A}^*(\bar{\epsilon})\|_{\mathbf{F}} \leq C'(\delta)\sigma_r(\bar{\mathbf{X}})$ .

 $\sigma_r(\bar{\mathbf{X}})$  is the *r*-th largest singular value of  $\bar{\mathbf{X}}$ .  $\mathcal{A}^*(\mathbf{b}) := \sum_{i=1}^n \mathbf{b}_i \mathbf{A}_i$  is the adjoint operator of  $\mathcal{A}$ .

Let  $\bar{\mathbf{X}}$  be a rank *r* stationary point and  $\bar{\boldsymbol{\epsilon}} := \mathbf{y} - \mathcal{A}(\bar{\mathbf{X}})$ .

<u>Theorem 1</u>: Under the assumptions above,  $\mathbf{X}^t$  generated by RISRO converges Q-linearly to  $\mathbf{\bar{X}}$ :

$$\|\mathbf{X}^{t+1} - \bar{\mathbf{X}}\|_{\mathbf{F}} \leq \frac{3}{4}\|\mathbf{X}^t - \bar{\mathbf{X}}\|_{\mathbf{F}}, \quad \forall t \geq 0.$$

Let  $\bar{\mathbf{X}}$  be a rank *r* stationary point and  $\bar{\epsilon} := \mathbf{y} - \mathcal{A}(\bar{\mathbf{X}})$ .

<u>Theorem 1</u>: Under the assumptions above,  $\mathbf{X}^t$  generated by RISRO converges Q-linearly to  $\mathbf{\bar{X}}$ :

$$\|\mathbf{X}^{t+1} - \bar{\mathbf{X}}\|_{\mathbf{F}} \leq \frac{3}{4}\|\mathbf{X}^t - \bar{\mathbf{X}}\|_{\mathbf{F}}, \quad \forall t \geq 0.$$

$$\begin{aligned} \|\mathbf{X}^{t+1} - \bar{\mathbf{X}}\|_{\mathbf{F}}^2 &\leq \frac{c_1(\delta) \|\mathbf{X}^t - \bar{\mathbf{X}}\|^2}{\sigma_t^2(\bar{\mathbf{X}})} \left( \|\mathbf{X}^t - \bar{\mathbf{X}}\|_{\mathbf{F}}^2 + \|\mathcal{A}^*(\bar{\boldsymbol{\epsilon}})\|_{\mathbf{F}} \|\mathbf{X}^t - \bar{\mathbf{X}}\|_{\mathbf{F}} + \|\mathcal{A}^*(\bar{\boldsymbol{\epsilon}})\|_{\mathbf{F}}^2 \right), \\ &\forall t \geq 0. \end{aligned}$$

Let  $\bar{\mathbf{X}}$  be a rank *r* stationary point and  $\bar{\epsilon} := \mathbf{y} - \mathcal{A}(\bar{\mathbf{X}})$ .

<u>Theorem 1</u>: Under the assumptions above,  $\mathbf{X}^t$  generated by RISRO converges Q-linearly to  $\bar{\mathbf{X}}$ :

$$\|\mathbf{X}^{t+1} - \bar{\mathbf{X}}\|_{\mathbf{F}} \leq \frac{3}{4}\|\mathbf{X}^t - \bar{\mathbf{X}}\|_{\mathbf{F}}, \quad \forall t \geq 0.$$

$$\|\mathbf{X}^{t+1} - \bar{\mathbf{X}}\|_{\mathsf{F}}^{2} \leq \frac{c_{1}(\delta) \|\mathbf{X}^{t} - \bar{\mathbf{X}}\|^{2}}{\sigma_{r}^{2}(\bar{\mathbf{X}})} \left( \|\mathbf{X}^{t} - \bar{\mathbf{X}}\|_{\mathsf{F}}^{2} + \|\mathcal{A}^{*}(\bar{\boldsymbol{\epsilon}})\|_{\mathsf{F}} \|\mathbf{X}^{t} - \bar{\mathbf{X}}\|_{\mathsf{F}} + \|\mathcal{A}^{*}(\bar{\boldsymbol{\epsilon}})\|_{\mathsf{F}}^{2} \right),$$
  
$$\forall t \geq 0.$$

If  $\bar{\boldsymbol{\epsilon}} = \boldsymbol{0}$ , then  $\{\boldsymbol{X}^t\}$  converges quadratically to  $\bar{\boldsymbol{X}}$  as

$$\|\mathbf{X}^{t+1} - \bar{\mathbf{X}}\|_{\mathbf{F}} \leq \frac{\sqrt{c_1(\delta)} \|\mathbf{X}^t - \bar{\mathbf{X}}\|_{\mathbf{F}}^2}{\sigma_r(\bar{\mathbf{X}})}, \quad \forall \, t \geq 0.$$

★ Quadratic-linear convergence

$$\|\mathbf{X}^{t+1} - \bar{\mathbf{X}}\|_{\mathbf{F}}^2 \leq \frac{c_1(\delta) \|\mathbf{X}^t - \bar{\mathbf{X}}\|^2}{\sigma_r^2(\bar{\mathbf{X}})} \left( \|\mathbf{X}^t - \bar{\mathbf{X}}\|_{\mathbf{F}}^2 + \|\mathcal{A}^*(\bar{\boldsymbol{\epsilon}})\|_{\mathbf{F}} \|\mathbf{X}^t - \bar{\mathbf{X}}\|_{\mathbf{F}} + \|\mathcal{A}^*(\bar{\boldsymbol{\epsilon}})\|_{\mathbf{F}}^2 \right)$$

- when ||X<sup>t</sup> X
   ||<sub>F</sub> ≫ ||A<sup>\*</sup>(ē)||<sub>F</sub> ⇒ quadratic convergence
   when ||X<sup>t</sup> X
   ||<sub>F</sub> ≤ c||A<sup>\*</sup>(ē)||<sub>F</sub> ⇒ reduce to linear convergence
- $ar{\epsilon}\downarrow\Longrightarrow$  Longer period of quadratic convergence.

#### ★ Quadratic-linear convergence

$$\|\mathbf{X}^{t+1} - \bar{\mathbf{X}}\|_{\mathbf{F}}^2 \leq \frac{c_1(\delta) \|\mathbf{X}^t - \bar{\mathbf{X}}\|^2}{\sigma_r^2(\bar{\mathbf{X}})} \left( \|\mathbf{X}^t - \bar{\mathbf{X}}\|_{\mathbf{F}}^2 + \|\mathcal{A}^*(\bar{\boldsymbol{\epsilon}})\|_{\mathbf{F}} \|\mathbf{X}^t - \bar{\mathbf{X}}\|_{\mathbf{F}} + \|\mathcal{A}^*(\bar{\boldsymbol{\epsilon}})\|_{\mathbf{F}}^2 \right)$$

 $ar{\epsilon}\downarrow\Longrightarrow$  Longer period of quadratic convergence.

★ 
$$\bar{\epsilon} = 0 \implies \mathbf{y} = \mathcal{A}(\bar{\mathbf{X}}) \implies$$
 matrix sensing  
[Recht et al., 2010]  
RISRO achieves quadratic convergence

#### Simulation

$$\begin{split} \mathbf{y}_i &= \langle \mathbf{A}_i, \mathbf{X}^* \rangle + \epsilon_i \text{ for } 1 \leq i \leq n, \ \mathbf{A}_i \overset{i.i.d.}{\sim} \mathcal{N}(0,1) \text{ and } \epsilon_i \overset{i.i.d.}{\sim} \mathcal{N}(0,\sigma^2). \\ \mathbf{X}^* &\in \mathbb{R}^{p \times p} \text{ with } p = 100, r = 3, \kappa(\mathbf{X}^*) = 1 \text{ and } \mathbf{X}^0 = \mathrm{SVD}_r(\mathcal{A}^*(\mathbf{y})). \end{split}$$

• (Quadratic-linear) n = 5pr,  $\sigma = 10^{\alpha}$  for  $\alpha \in \{0, -1, -2, -14\}$ 



Quadratic) n/(pr) ∈ {4,5,6,7,8}, σ = 0



Suppose  $p_1 = p_2 = p$  and  $n \ge pr$ . Under similar assumptions as in Theorem 1:

	GD	PGD (SVP / IHT)	Alter Mini	RISRO (this work)
Iteration Complexity	O(np²r)	O(np <sup>2</sup> )	$O(np^2r^2)$	$O(np^2r^2)$
Tuning	Yes	Yes	No	No
Convergence	Linear	Linear	Linear	Quadratic-(linear)

★ Improve upon Alter Mini for free

#### Comparison Simulation $\sigma = 0$



 $\kappa = 500$ 



Meeting of WHU-XMU

# Any connection of RISRO to existing optimization algorithms?



# Connection to Riemannian Manifold Optimization

#### Iteration t of RISRO:

- Perform importance sketching.
- Perform a dimension reduced least squares.

**③** Update sketching matrices and  $\mathbf{X}^{t+1}$ .

# Connection to Riemannian Manifold Optimization

Iteration t of RISRO:

- Perform importance sketching.
- Perform a dimension reduced least squares.

→ Implicitly solves "Fisher Scoring" or "Riemannian Gauss-Newton" equation in Riemannian optimization on fixed rank matrices.

Opdate sketching matrices and X<sup>t+1</sup>.

 $\implies$  Perform a type of retraction in Riemannian optimization literature



# **Riemannian Manifold Optimization**

- Target: optimize a function *f* defined on a Riemannian manifold *M*. [Absil et al., 2009]
- Common Riemannian manifolds, embedded submanifold: a smooth subset of R<sup>n</sup> + a Riemannian metric.

- Target: optimize a function *f* defined on a Riemannian manifold *M*. [Absil et al., 2009]
- Common Riemannian manifolds, embedded submanifold: a smooth subset of R<sup>n</sup> + a Riemannian metric.

• 
$$\mathcal{M}_r = \{\mathbf{X} \in \mathbb{R}^{p_1 \times p_2} : \operatorname{rank}(\mathbf{X}) = r\}$$

Riemannian metric: Euclidean inner product,  $\langle \boldsymbol{U}, \boldsymbol{V} \rangle = \operatorname{trace}(\boldsymbol{U}^\top \boldsymbol{V})$ 

## Retraction

 Iterative algorithm: x<sup>t+1</sup> = x<sup>t</sup> + ξ.
 Manifold optimization: x<sup>t+1</sup> may not lie in the manifold Solution: retraction!

### Retraction

- Iterative algorithm: x<sup>t+1</sup> = x<sup>t</sup> + ξ.
   Manifold optimization: x<sup>t+1</sup> may not lie in the manifold Solution: retraction!
- Retraction: a smooth map that brings the vector in the tangent space back to the manifold. Denote T<sub>x</sub>M as the tangent space at x



 $R: T\mathcal{M} \to \mathcal{M}, x \times \xi \to R_x(\xi) \in \mathcal{M}.$ 

## Retraction

- Iterative algorithm: x<sup>t+1</sup> = x<sup>t</sup> + ξ.
   Manifold optimization: x<sup>t+1</sup> may not lie in the manifold Solution: retraction!
- Retraction: a smooth map that brings the vector in the tangent space back to the manifold. Denote T<sub>x</sub>M as the tangent space at x
- ★ Let  $\eta^t$  be the update direction such that  $\mathbf{X}^t + \eta^t$  has the following representation,

$$\mathbf{X}^t + \boldsymbol{\eta}^t = \begin{bmatrix} \mathbf{U}^t & \mathbf{U}_{\perp}^t \end{bmatrix} \begin{bmatrix} \mathbf{B}^{t+1} & \mathbf{D}_2^{t+1\top} \\ \mathbf{D}_1^{t+1} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}^t & \mathbf{V}_{\perp}^t \end{bmatrix}^{\top}.$$

★  $X^t + \eta^t \implies X^{t+1}$ . Retraction is:

$$\mathbf{X}^{t+1} = R_{\mathbf{X}^{t}}(\boldsymbol{\eta}^{t}) = \begin{bmatrix} \mathbf{U}^{t} & \mathbf{U}_{\perp}^{t} \end{bmatrix} \begin{bmatrix} \mathbf{B}^{t+1} & \mathbf{D}_{2}^{t+1\top} \\ \mathbf{D}_{1}^{t+1} & \mathbf{D}_{1}^{t+1} (\mathbf{B}^{t+1})^{-1} \mathbf{D}_{2}^{t+1\top} \end{bmatrix} \begin{bmatrix} \mathbf{V}^{t} & \mathbf{V}_{\perp}^{t} \end{bmatrix}^{\top}$$

#### $\star$ $\eta^t$ solves the Fisher Scoring or Riemannian Gauss-Newton

Meeting of WHU-XMU

Recall  $f(\mathbf{X}) := \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2$ .

• Riemannian Gradient: grad f(X)

• Riemannian Hessian: Hessf(X)

Riemannian Newton direction η<sub>Newton</sub>

 $-\operatorname{grad} f(\mathbf{X}) = \operatorname{Hess} f(\mathbf{X})[\eta_{\operatorname{Newton}}]$ 

Recall  $f(\mathbf{X}) := \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2$ .

- Riemannian Gradient: grad  $f(\mathbf{X}) = P_{T_{\mathbf{X}}}(\mathcal{A}^*(\mathcal{A}(\mathbf{X}) \mathbf{y}))$ .  $P_{T_{\mathbf{X}}}(\cdot)$  is the orthogonal projector onto the tangent space at  $\mathbf{X}$ .
- Riemannian Hessian:  $\operatorname{Hess} f(\mathbf{X})$   $[\eta] = P_{T_{\mathbf{X}}} (\mathcal{A}^*(\mathcal{A}(\eta))) + h(\mathbf{y} - \mathcal{A}(\mathbf{X})).$  $h(\cdot)$  depends on  $\mathbf{X}, \eta$ .
- Riemannian Newton direction η<sub>Newton</sub>

 $-\operatorname{grad} f(\mathbf{X}) = \operatorname{Hess} f(\mathbf{X})[\eta_{\operatorname{Newton}}]$ 

 $\iff -\operatorname{grad} f(\mathbf{X}) = P_{\mathcal{T}_{\mathbf{X}}} \left( \mathcal{A}^*(\mathcal{A}(\eta_{\operatorname{Newton}})) \right) + h(\mathbf{y} - \mathcal{A}(\mathbf{X}))$ 

Recall  $f(\mathbf{X}) := \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2$ .

- Riemannian Gradient: grad  $f(\mathbf{X}) = P_{T_{\mathbf{X}}}(\mathcal{A}^*(\mathcal{A}(\mathbf{X}) \mathbf{y}))$ .  $P_{T_{\mathbf{X}}}(\cdot)$  is the orthogonal projector onto the tangent space at  $\mathbf{X}$ .
- Riemannian Hessian:  $\operatorname{Hess} f(\mathbf{X})$   $[\eta] = P_{T_{\mathbf{X}}} (\mathcal{A}^*(\mathcal{A}(\eta))) + h(\mathbf{y} - \mathcal{A}(\mathbf{X})).$  $h(\cdot)$  depends on  $\mathbf{X}, \eta$ .
- Riemannian Newton direction η<sub>Newton</sub>

 $-\operatorname{grad} f(\mathbf{X}) = \operatorname{Hess} f(\mathbf{X})[\eta_{\operatorname{Newton}}]$ 

 $\iff -\operatorname{grad} f(\mathbf{X}) = \boldsymbol{P}_{T_{\mathbf{X}}} \left( \mathcal{A}^*(\mathcal{A}(\eta_{\operatorname{Newton}})) \right) + \boldsymbol{h}(\mathbf{y} - \mathcal{A}(\mathbf{X}))$ 

Update in RISRO:

$$\mathbf{X}^{t} + \boldsymbol{\eta}^{t} = \begin{bmatrix} \mathbf{U}^{t} & \mathbf{U}_{\perp}^{t} \end{bmatrix} \begin{bmatrix} \mathbf{B}^{t+1} & \mathbf{D}_{2}^{t+1} \\ \mathbf{D}_{1}^{t+1} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}^{t} & \mathbf{V}_{\perp}^{t} \end{bmatrix}^{\top}.$$

Recall  $f(\mathbf{X}) := \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2$ .

- Riemannian Gradient: grad  $f(\mathbf{X}) = P_{T_{\mathbf{X}}}(\mathcal{A}^*(\mathcal{A}(\mathbf{X}) \mathbf{y}))$ .  $P_{T_{\mathbf{X}}}(\cdot)$  is the orthogonal projector onto the tangent space at  $\mathbf{X}$ .
- Riemannian Hessian:  $\operatorname{Hess} f(\mathbf{X})$   $[\eta] = P_{T_{\mathbf{X}}} (\mathcal{A}^*(\mathcal{A}(\eta))) + h(\mathbf{y} - \mathcal{A}(\mathbf{X})).$  $h(\cdot)$  depends on  $\mathbf{X}, \eta$ .
- Riemannian Newton direction η<sub>Newton</sub>

 $-\operatorname{grad} f(\mathbf{X}) = \operatorname{Hess} f(\mathbf{X})[\eta_{\operatorname{Newton}}]$ 

 $\iff -\operatorname{grad} f(\mathbf{X}) = \boldsymbol{P}_{T_{\mathbf{X}}} \left( \mathcal{A}^*(\mathcal{A}(\eta_{\operatorname{Newton}})) \right) + \boldsymbol{h}(\mathbf{y} - \mathcal{A}(\mathbf{X}))$ 

Update in RISRO:

$$\mathbf{X}^t + \boldsymbol{\eta}^t = \begin{bmatrix} \mathbf{U}^t & \mathbf{U}_{\perp}^t \end{bmatrix} \begin{bmatrix} \mathbf{B}^{t+1} & \mathbf{D}_2^{t+1\top} \\ \mathbf{D}_1^{t+1} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}^t & \mathbf{V}_{\perp}^t \end{bmatrix}^{\top}.$$

<u>Theorem 2</u>:  $\eta^t$  solves

$$-\operatorname{grad} f(\mathbf{X}^{t}) = \boldsymbol{P}_{\mathcal{T}_{\mathbf{X}^{t}}}\left(\mathcal{A}^{*}(\mathcal{A}(\eta))\right)$$

Meeting of WHU-XMU

# Connection of RISRO and Riemannian optimization

Suppose  $\mathbf{y} = \mathcal{A}(\mathbf{X}) + \epsilon$ , where **X** is a fixed matrix and  $\epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$ . Then for any  $\eta$ ,

$$\{\mathbb{E}(\operatorname{Hess} f(\mathbf{X})[\boldsymbol{\eta}])\}|_{\mathbf{X}=\mathbf{X}^{t}} = \boldsymbol{P}_{T_{\mathbf{X}^{t}}}\left(\mathcal{A}^{*}(\mathcal{A}(\boldsymbol{\eta}))\right).$$

# Connection of RISRO and Riemannian optimization

Suppose  $\mathbf{y} = \mathcal{A}(\mathbf{X}) + \epsilon$ , where **X** is a fixed matrix and  $\epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$ . Then for any  $\eta$ ,

$$\{\mathbb{E}(\text{Hess} f(\mathbf{X})[\boldsymbol{\eta}])\}|_{\mathbf{X}=\mathbf{X}^{t}} = \boldsymbol{P}_{\mathcal{T}_{\mathbf{X}^{t}}}\left(\mathcal{A}^{*}(\mathcal{A}(\boldsymbol{\eta}))\right).$$

By Theorem 2,  $\eta^t$  solves

 $-\operatorname{grad} f(\mathbf{X}^{t}) = \{\mathbb{E}(\operatorname{Hess} f(\mathbf{X})[\boldsymbol{\eta}])\} |_{\mathbf{X} = \mathbf{X}^{t}}.$ 

# Connection of RISRO and Riemannian optimization

Suppose  $\mathbf{y} = \mathcal{A}(\mathbf{X}) + \epsilon$ , where **X** is a fixed matrix and  $\epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$ . Then for any  $\eta$ ,

$$\{\mathbb{E}(\text{Hess} f(\mathbf{X})[\boldsymbol{\eta}])\}|_{\mathbf{X}=\mathbf{X}^{t}} = \boldsymbol{P}_{\mathcal{T}_{\mathbf{X}^{t}}}\left(\mathcal{A}^{*}(\mathcal{A}(\boldsymbol{\eta}))\right).$$

By Theorem 2,  $\eta^t$  solves

```
-\operatorname{grad} f(\mathbf{X}^{t}) = \{\mathbb{E}(\operatorname{Hess} f(\mathbf{X})[\boldsymbol{\eta}])\}|_{\mathbf{X}=\mathbf{X}^{t}}.
```

This algorithm is called Fisher Scoring in literature [Lange, 2010].



#### **Applications to Statistics and Machine Learning**



## Applications to Statistics and Machine Learning

• Low-rank matrix trace regression model:

$$\mathbf{y}_i = \langle \mathbf{A}_i, \mathbf{X}^* \rangle + \epsilon_i, \quad \text{ for } 1 \leq i \leq n,$$

 $\mathbf{X}^* \in \mathbb{R}^{p_1 \times p_2}$  is the true model parameter and  $\operatorname{rank}(\mathbf{X}^*) = r$ . • Phase retrieval

$$\mathbf{y}_i = |\langle \mathbf{a}_i, \mathbf{x}^* \rangle|^2$$
 for  $1 \le i \le n$ ,

 $\mathbf{X}^* \in \mathbb{R}^p$ .

Goal: estimate or recovery **X**<sup>\*</sup> (or **x**<sup>\*</sup>).

<u>Theorem</u>: Suppose A satisfies the 3r-RIP with RIP constant  $\delta$  and

• 
$$\|\mathbf{X}^0 - \mathbf{X}^*\|_{\mathbf{F}} \leq C(\delta) \cdot \sigma_r(\mathbf{X}^*)$$

• 
$$\sigma_r(\mathbf{X}^*) \geq C'(\delta) \cdot \sqrt{r} \| \mathcal{A}^*(\boldsymbol{\epsilon}) \|.$$

Then iterations generated by RISRO satisfy

$$\begin{aligned} \|\mathbf{X}^{t+1} - \mathbf{X}^*\|_{\mathbf{F}}^2 \leq c_1(\delta) \|\mathbf{X}^t - \mathbf{X}^*\|^2 \left(\frac{\|\mathbf{X}^t - \mathbf{X}^*\|_{\mathbf{F}}^2}{\sigma_r^2(\mathbf{X}^*)} + \frac{\sqrt{r}\|\mathcal{A}^*(\epsilon)\|}{\sigma_r(\mathbf{X}^*)}\right) \\ + c_2(\delta)r\|\mathcal{A}^*(\epsilon)\|^2, \end{aligned}$$

for all  $t \ge 0$ .

★ First term: Decreases quadratic-linearly.

**\star** Second term: Statistical error independent of *t*.

- Introduce a new algorithm, RISRO, for rank constrained least squares.
  - $\implies$  Tuning free, fast and has high-order convergence
- Introduce the recursive importance sketching framework
   Provide a platform to compare different algorithms from a sketching perspective
- Connect RISRO with Riemannian optimization

# Thank you! Questions?

Luo, Y., Huang, W., Li, X., & Zhang, A. R. (2020). Recursive Importance Sketching for Rank Constrained Least Squares: Algorithms and High-order Convergence. arXiv preprint arXiv:2011.08360.

# Bibliography

- Absil, P.-A., Mahony, R., and Sepulchre, R. (2009). Optimization algorithms on matrix manifolds. Princeton University Press.
- Bauch, J. and Nadler, B. (2020). Rank 2r iterative least squares: efficient recovery of ill-conditioned low rank matrices from few entries. arXiv preprint arXiv:2002.01849.
- Bhojanapalli, S., Neyshabur, B., and Srebro, N. (2016). Global optimality of local search for low rank matrix recovery.

In Advances in Neural Information Processing Systems, pages 3873–3881.

Boumal, N. and Absil, P.-a. (2011). Rtrmc: A Riemannian trust-region method for low-rank matrix completion.

In Advances in neural information processing systems. Meeting of WHU-XMU December, :