

How to improve first-order optimization methods

Geometry, block structure, and randomization

Hanbaek Lyu

Department of Mathematics
University of Wisconsin – Madison

Partially supported by NSF DMS #2206296 and #2010035

Deep Learning seminar @ KRAFTON AI
May. 29, 2024

Students and Collaborators



Yuchen Li
(4th year Ph.D.)

Block Optimization
Riemannian Optimization



William Powell
(3rd year Ph.D.)

Stochastic Optimization
Reinforcement Learning



Danny Duan
(3rd year Ph.D.)

Quasi-Newton Methods



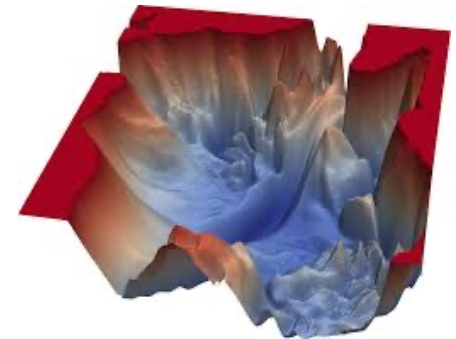
Joowon Lee
(UW Stats)

Causal Inference
Matrix Factorization

Introduction

$$\theta^* \in \arg \min_{\theta \in \Theta \subseteq \mathbb{R}^p} f(\theta)$$

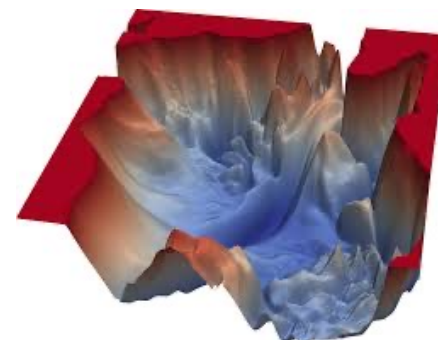
- f : Objective function
 - Nonconvex, smooth (for this talk)
 - Often represents model fitness to training data
 - e.g., DNN, LLM




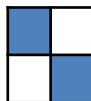
$$\theta^* \in \arg \min_{\theta \in \Theta \subseteq \mathbb{R}^p} f(\theta)$$

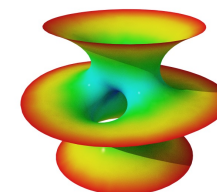
- f : Objective function

- Nonconvex, smooth (for this talk)
- Often represents model fitness to training data
 - e.g., DNN, LLM



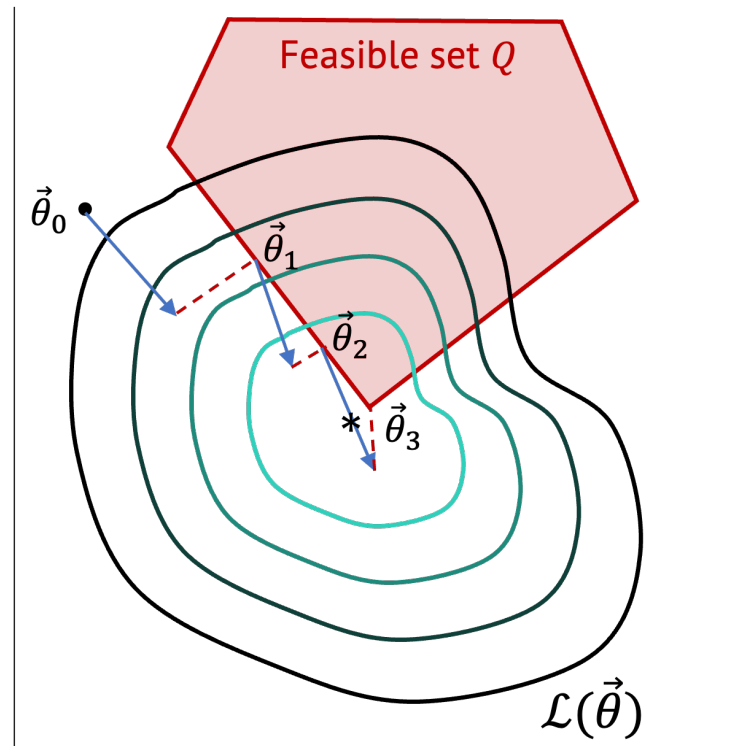
- Θ : Parameter space with various structures:

- Convex 
- block structure: $\theta = (\theta_1, \dots, \theta_m)$ 
- Riemannian manifold (e.g., Θ =set of orthogonal matrices)



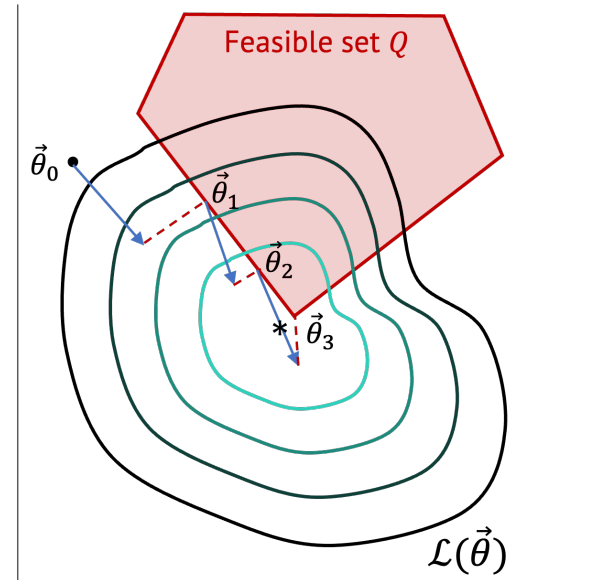
The simplest yet commonly used method

(PGD) $\theta_{n+1} \leftarrow \Pi_{\Theta} (\theta_n - \alpha_n \nabla f(\theta_n))$



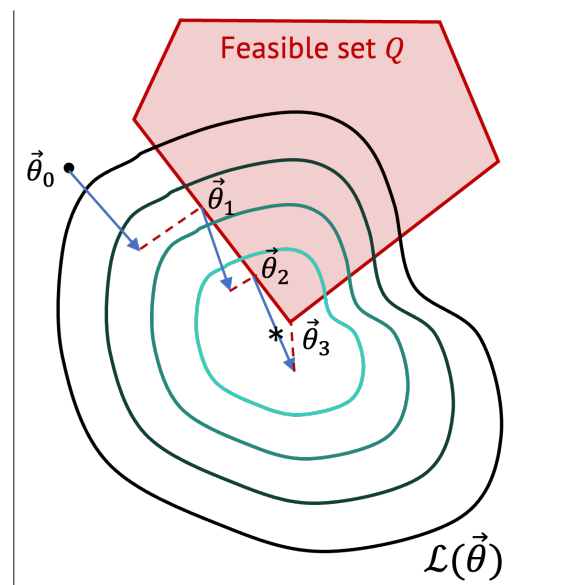
$$\text{(PGD)} \quad \theta_{n+1} \leftarrow \Pi_{\Theta} (\theta_n - \alpha_n \nabla f(\theta_n))$$

- Projected Gradient Descent (PGD)
 - Widely used in ML, Statistics, Scientific Computing
 - Reduces to Gradient Descent for unconstrained problems
 - Easy to implement, cheap to run
 - Several of modifications (e.g., momentum, adaptive stepsizes)



$$\text{(PGD)} \quad \theta_{n+1} \leftarrow \Pi_{\Theta} (\theta_n - \alpha_n \nabla f(\theta_n))$$

- Projected Gradient Descent (PGD)
 - Widely used in ML, Statistics, Scientific Computing
 - Reduces to Gradient Descent for unconstrained problems
 - Easy to implement, cheap to run
 - Several of modifications (e.g., momentum, adaptive stepsizes)



$$\text{(PSGD)} \quad \theta_{n+1} \leftarrow \Pi_{\Theta} (\theta_n - \alpha_n \widehat{\nabla} f(\theta_n))$$

Stochastic gradient $\approx \nabla$

$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} (\boldsymbol{\theta}_n - \mathbf{H}_n \nabla f(\boldsymbol{\theta}_n))$$

- Π_{Θ} : Projection onto Θ
 - Simple for nonnegativity or norm constraints
- $\nabla f(\boldsymbol{\theta}_n)$: Gradient of the objective at current parameter $\boldsymbol{\theta}_n$
 - Assume it is cheap to compute; Requires $O(p)$ storage
- \mathbf{H}_n : $p \times p$ Preconditioning matrix
 - $\mathbf{H}_n = \alpha_n \mathbf{I}_p \rightarrow$ Gradient Descent with stepsize α_n
 - $\mathbf{H}_n = \nabla^2 f(\boldsymbol{\theta}_n)^{-1} \rightarrow$ Newton's Method
 - $\mathbf{H}_n \approx \nabla^2 f(\boldsymbol{\theta}_n)^{-1} \rightarrow$ Quasi-Newton Method

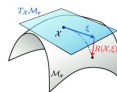
- Classical 1st-order methods like P(S)GD has some limitations
 - Difficult/expensive to project onto the constrained set
 - Still expensive with high dimensionality ($p \gg 1$)
 - e.g., overparametrization
 - Slow convergence for ill-conditioned problems
 - Mathematical guarantees are too pessimistic vs. practice
 - e.g., Usually assumes the worst-case, requires very small stepsizes, i.i.d. data sampling,

- Classical 1st-order methods like P(S)GD has some limitations
 - Difficult/expensive to project onto the constrained set
 - Still expensive with high dimensionality ($p \gg 1$)
 - e.g., overparametrization
 - Slow convergence for ill-conditioned problems
 - Mathematical guarantees are too pessimistic vs. practice
 - e.g., Usually assumes the worst-case, requires very small stepsizes, i.i.d. data sampling,
- We will propose several techniques to improve them, each based on some mathematical principles

Plan of the talk

- Classical 1st-order methods like P(S)GD has some limitations
 - Difficult/expensive to project onto the constrained set
 - Still expensive with high dimensionality ($p \gg 1$)
 - e.g., overparametrization
 - Slow convergence for ill-conditioned problems
 - Mathematical guarantees are too pessimistic vs. practice
 - e.g., Usually assumes the worst-case, requires very small stepsizes, i.i.d. data sampling,
- We will propose several techniques to improve them, each based on some mathematical principles

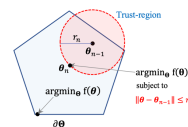
- Riemannian OPT



- Adaptive BCD



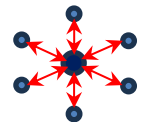
- Non-adaptive TR



- Rand. Hessian Approx.



- Gradient Consensus



A Math Warm-up:

*Basic Ideas behind nonconvex
optimization algorithms*

Taylor expansion!

Frank-Wolfe (1956)

- Taylor expand the objective f near θ_n :

$$f(\theta) \approx f(\theta_n) + \langle \nabla f(\theta_n), \theta - \theta_n \rangle$$

- Minimize the 1st-order Taylor approximation

$$\theta'_n = \arg \min_{\theta \in \Theta} \langle \nabla f(\theta_n), \theta - \theta_n \rangle$$

$$\theta_{n+1} = \text{Convex combination of } \theta_n \text{ and } \theta'_n$$

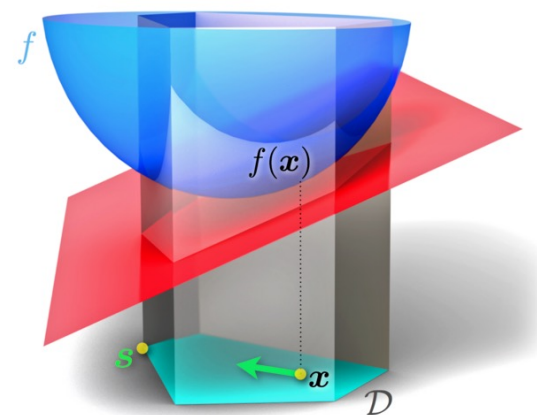


Image credit: Wikipedia

Newton's Method (1690) (Assume $\Theta = \mathbb{R}^p$)

- Taylor expand the objective f near θ_n :

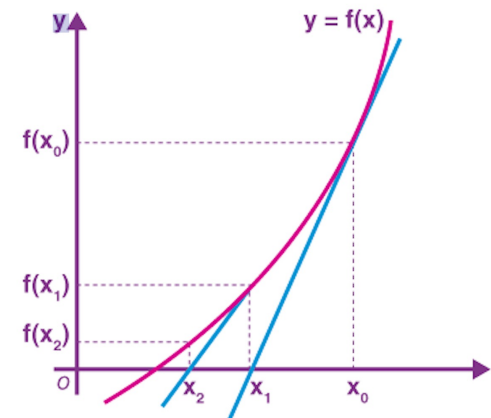
$$f(\theta) \approx f(\theta_n) + \langle \nabla f(\theta_n), \theta - \theta_n \rangle + \frac{1}{2} \langle \theta - \theta_n, \nabla^2 f(\theta_n)(\theta - \theta_n) \rangle$$

- Minimize the 2nd-order Taylor expansion

$$\begin{aligned} \theta_{n+1} &= \arg \min_{\theta \in \mathbb{R}^p} f(\theta_n) + \langle \nabla f(\theta_n), \theta - \theta_n \rangle + \frac{1}{2} \langle \theta - \theta_n, \nabla^2 f(\theta_n)(\theta - \theta_n) \rangle \\ &= \theta_n - (\nabla^2 f(\theta_n))^{-1} \nabla f(\theta_n) \end{aligned}$$

- Undefined if the Hessian is not PD
- Inverting the Hessian is expensive

Super-fast, quadratic convergence if works ☺



A Quasi-Newton Method (Assume $\Theta = \mathbb{R}^p$)

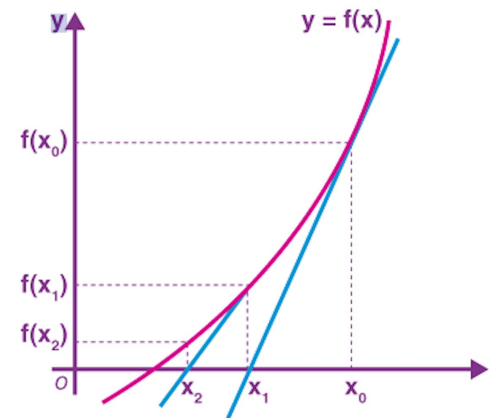
- Taylor expand the objective f near θ_n :

$$f(\theta) \approx f(\theta_n) + \langle \nabla f(\theta_n), \theta - \theta_n \rangle + \frac{1}{2} \langle \theta - \theta_n, \nabla^2 f(\theta_n)(\theta - \theta_n) \rangle$$

- Minimize the 2nd-order Taylor expansion with regularization

$$\begin{aligned} \theta_{n+1} &= \arg \min_{\theta \in \mathbb{R}^p} \langle \nabla f(\theta_n), \theta - \theta_n \rangle + \frac{1}{2} \langle \theta - \theta_n, \nabla^2 f(\theta_n)(\theta - \theta_n) \rangle + \frac{\gamma_n}{2} \|\theta - \theta_n\|^2 \\ &= \theta_n - (\nabla^2 f(\theta_n) + \gamma_n \mathbf{I}_p)^{-1} \nabla f(\theta_n) \end{aligned}$$

- Levenberg-Marquardt regularization ('44 '63)
- Can make the regularized Hessian PD
- Matrix inversion still expensive



Trust-Region (1970's)

- Taylor expand the objective f near θ_n :

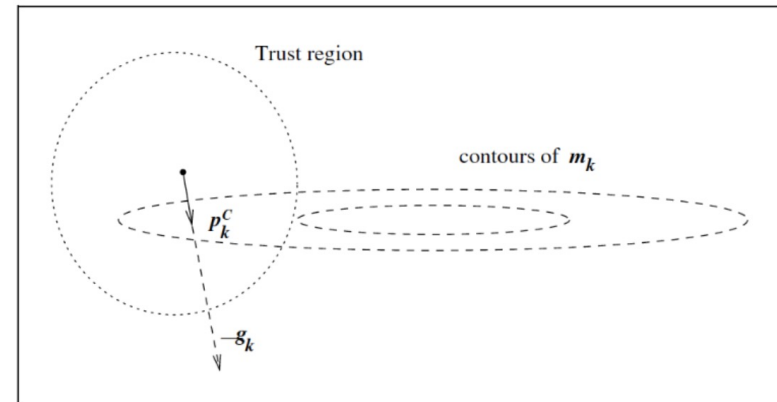
$$f(\theta) \approx f(\theta_n) + \langle \nabla f(\theta_n), \theta - \theta_n \rangle + \frac{1}{2} \langle \theta - \theta_n, \nabla^2 f(\theta_n)(\theta - \theta_n) \rangle$$

- Minimize a “quadratic model” within a trust-region

$$\theta_{n+1} = \underset{\theta \in \Theta, \|\theta - \theta_n\| \leq r_n}{\operatorname{arg\,min}} \langle \nabla f(\theta_n), \theta - \theta_n \rangle + \frac{1}{2} \langle \theta - \theta_n, \mathbf{B}_n(\theta - \theta_n) \rangle$$

- Trust-region with radius r_n
- Next radius r_{n+1} is computed adaptively based on the performance of the previous update

- PD, and no need to be close to the Hessian



PGD (Late 1950s)

- Taylor expand the objective f near θ_n :

$$f(\theta) \approx f(\theta_n) + \langle \nabla f(\theta_n), \theta - \theta_n \rangle + \frac{1}{2} \langle \theta - \theta_n, \nabla^2 f(\theta_n)(\theta - \theta_n) \rangle$$

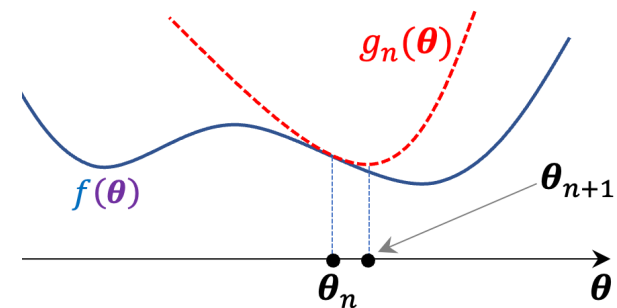
- Replace the Hessian with L times the identity

$$f(\theta) \leq f(\theta_n) + \langle \nabla f(\theta_n), \theta - \theta_n \rangle + \frac{L}{2} \|\theta - \theta_n\|^2$$

L = largest absolute eigenvalue of $\nabla^2 f(\theta)$ over all θ

- Minimize the quadratic surrogate

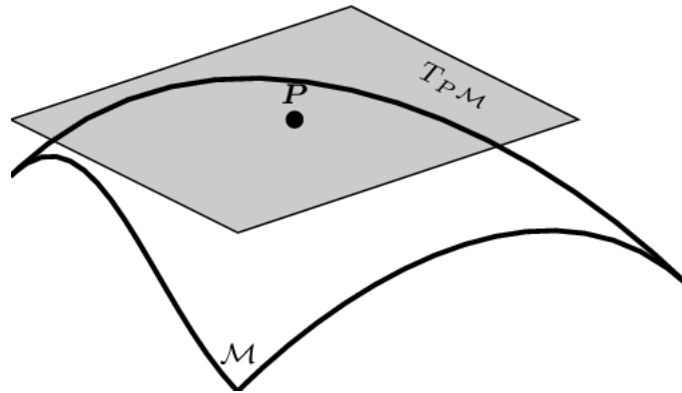
$$\begin{aligned} \theta_{n+1} &= \arg \min_{\theta \in \Theta} \langle \nabla f(\theta_n), \theta - \theta_n \rangle + \frac{L}{2} \|\theta - \theta_n\|^2 \\ &= \arg \min_{\theta \in \Theta} \left\| \theta - \left(\theta_n - \frac{1}{L} \nabla f(\theta_n) \right) \right\|^2 \\ &= \Pi_{\Theta} \left(\theta_n - \frac{1}{L} \nabla f(\theta_n) \right) \end{aligned}$$



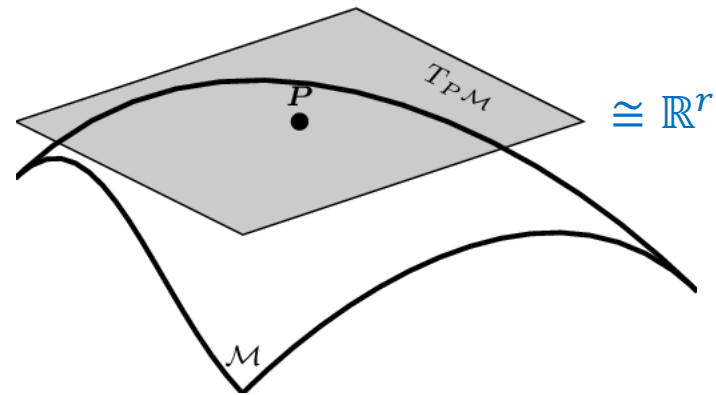
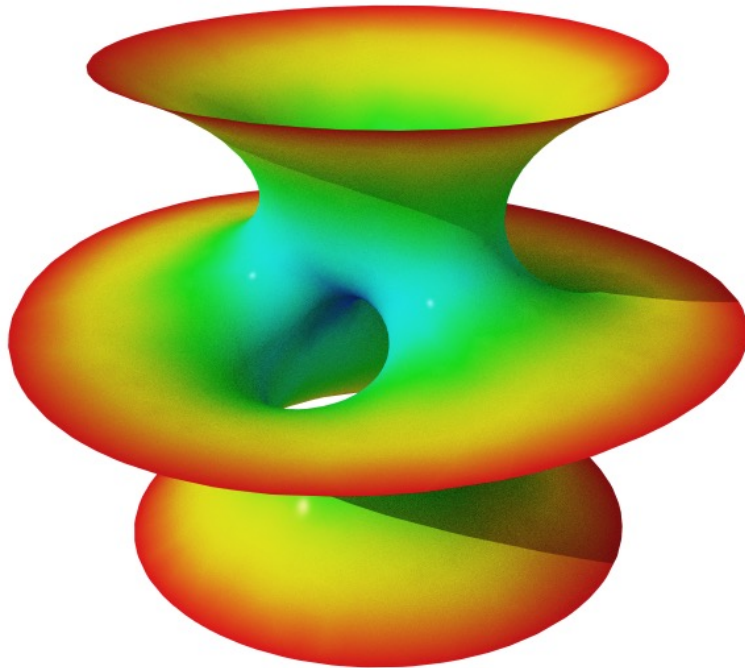
A geometric approach:

Riemannian Optimization

Tangent planes!

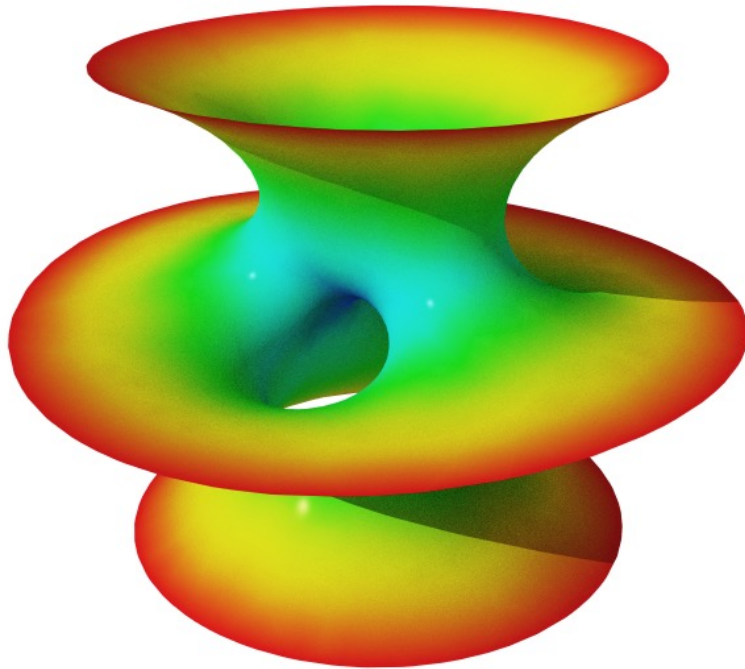


- Parameter spaces often has some intrinsic, ***low-dimensional geometric structures***
- **(r -dim) Riemannian manifold**
 - = Smooth surface in \mathbb{R}^p where every point has r -dim tangent spaces



- e.g.,
- Low-rank matrix manifolds
 - Orthogonal frames
 - PSD matrices
 - DNNs with spectrum-constrained matrix weights

- Parameter spaces often has some intrinsic, ***low-dimensional geometric structures***
- **(r -dim) Riemannian manifold**
 - = Smooth surface in \mathbb{R}^p where every point has r -dim tangent spaces

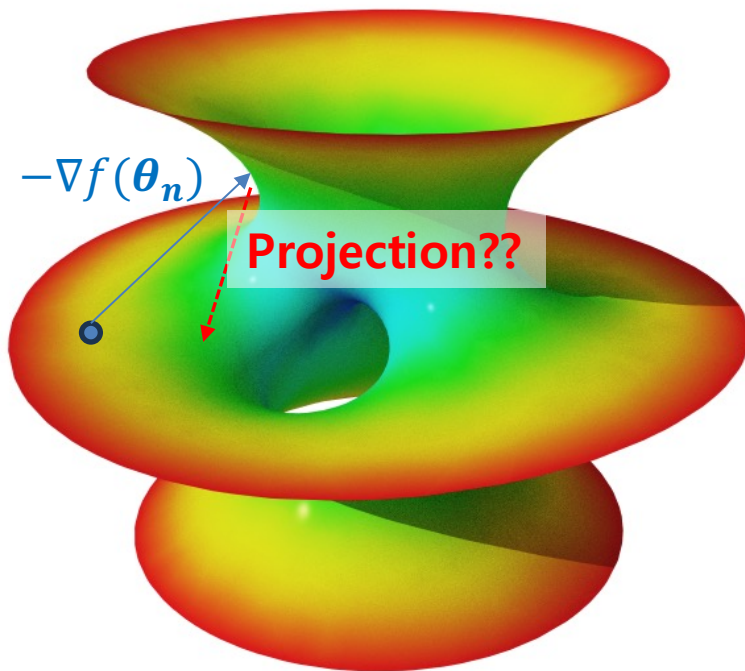


Riemann
(1826-1866)



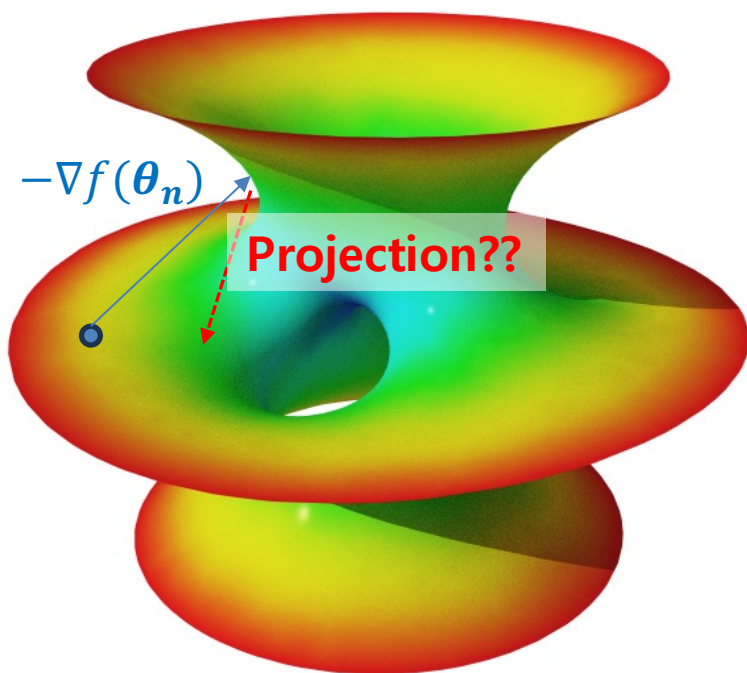
Gauss
(1777-1855)

$$\text{(PGD)} \quad \boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} (\boldsymbol{\theta}_n - \alpha_n \nabla f(\boldsymbol{\theta}_n))$$



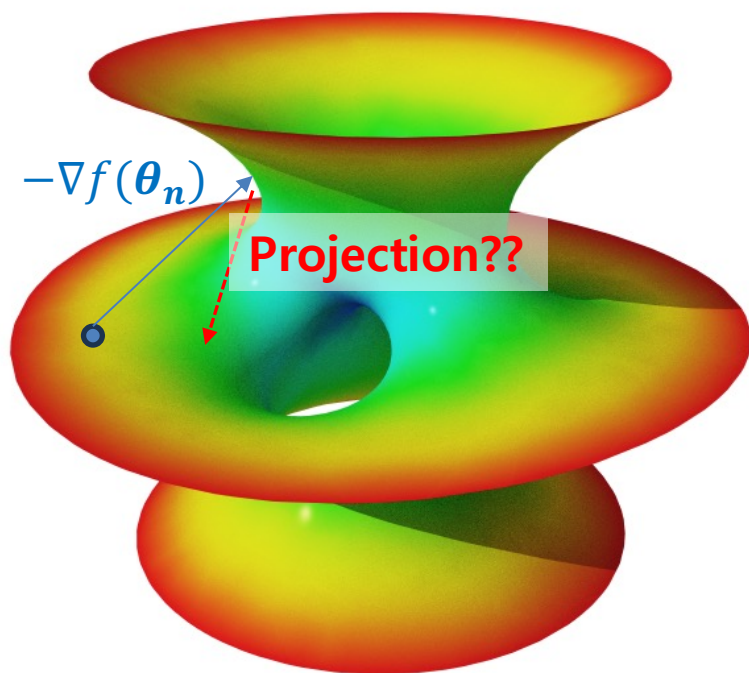
- Π_{Θ} : Projection onto Θ :
 - Not always easy!

$$\text{(PGD)} \quad \boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} (\boldsymbol{\theta}_n - \alpha_n \nabla f(\boldsymbol{\theta}_n))$$



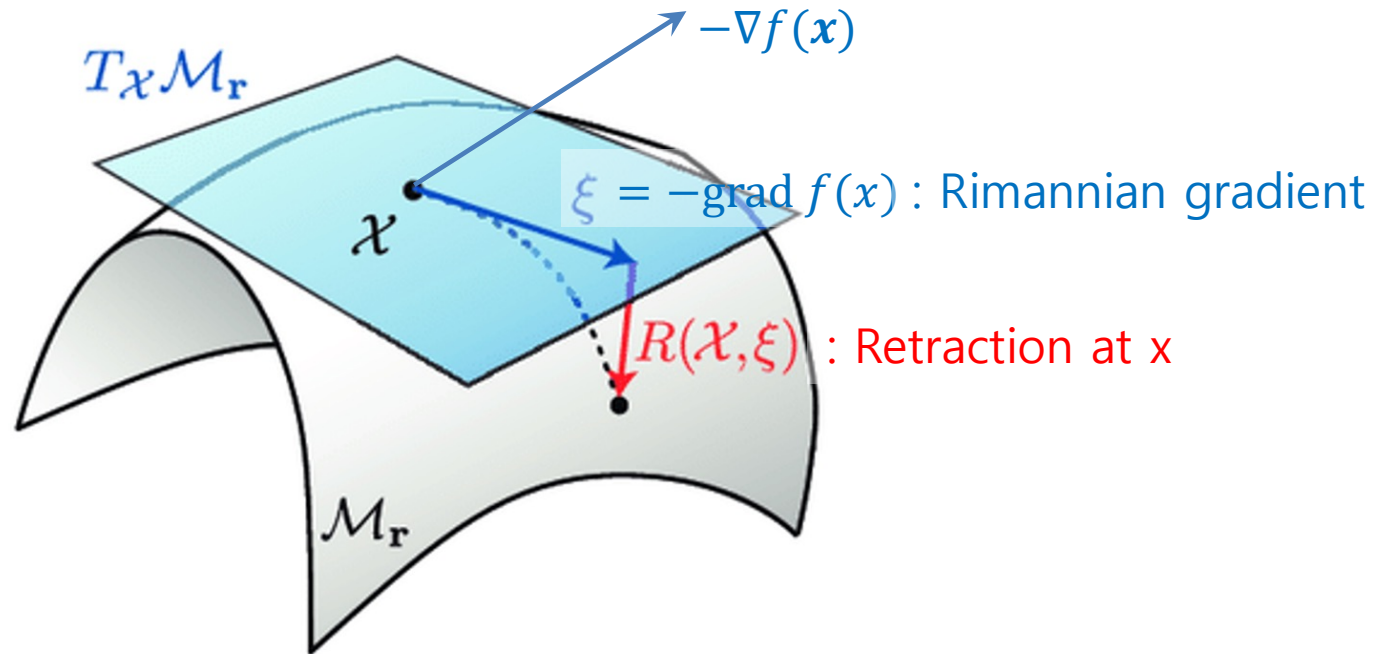
- Π_{Θ} : Projection onto Θ :
 - Not always easy!
- If Θ is low-dim (say $r \ll p$)
 - Work only with r-dim stuffs?

$$\text{(PGD)} \quad \boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} (\boldsymbol{\theta}_n - \alpha_n \nabla f(\boldsymbol{\theta}_n))$$



- Π_{Θ} : Projection onto Θ :
 - Not always easy!
- If Θ is low-dim (say $r \ll p$)
 - Work only with r-dim stuffs?
- *For certain low-rank matrix factorization problems, PGD with rSVD as projection can be useful*
[Lee, Lyu, Yao NeurIPS '23]

$$\text{(Riemannian GD)} \quad \theta_{n+1} \leftarrow \text{Rtr}_{\theta_n} (-\alpha_n \text{grad } f(\theta_n))$$



(Proj from $T_x \rightarrow \mathcal{M}$)

- RGD = gradient descent within the tangent space + Retraction

(RGD)

$$\boldsymbol{\theta}_{n+1} \leftarrow \text{Rtr}_{\boldsymbol{\theta}_n} (-\alpha_n \text{grad } f(\boldsymbol{\theta}_n))$$

- Limitation: **grad** and **Retraction** can be expensive to compute

$$\text{(RGD)} \quad \boldsymbol{\theta}_{n+1} \leftarrow \text{Rtr}_{\boldsymbol{\theta}_n} \left(-\alpha_n \text{grad} f(\boldsymbol{\theta}_n) \right)$$

- Limitation: **grad** and **Retraction** can be expensive to compute

[Li, Lyu, Balzano, Needell ICML '24]

$$\text{(Inexact RGD)} \quad \boldsymbol{\theta}_{n+1} \leftarrow \widehat{\text{Rtr}}_{\boldsymbol{\theta}_n} \left(-\alpha_n \widehat{\text{grad}} f(\boldsymbol{\theta}_n) \right)$$

- **It's OK to use approximate grad and Retraction**
(Both mathematically and practically)

(RGD)

$$\theta_{n+1} \leftarrow \text{Rtr}_{\theta_n} (-\alpha_n \text{grad } f(\theta_n))$$

- Limitation: **grad** and **Retraction** can be expensive to compute

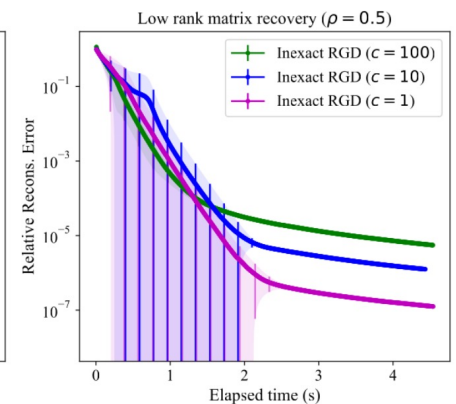
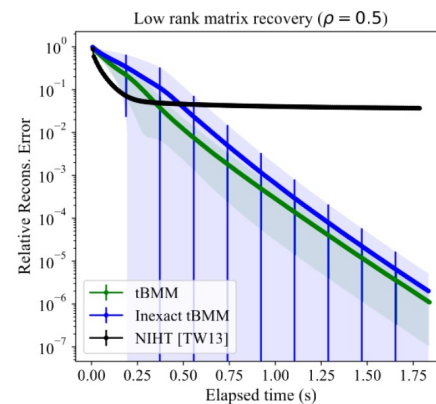
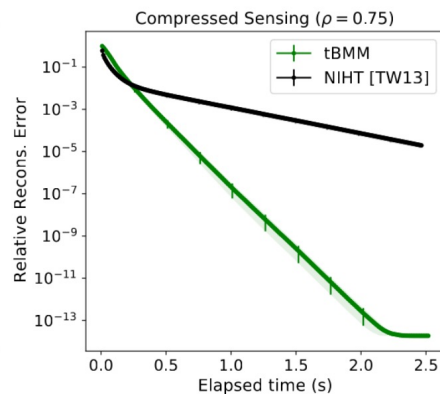
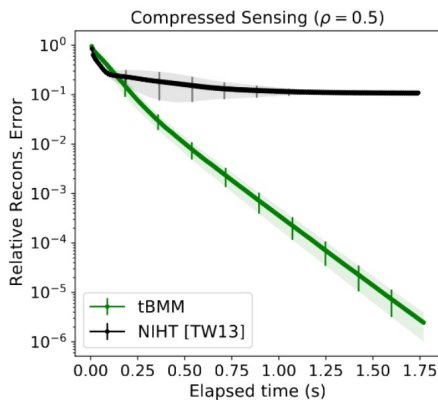
[Li, Lyu, Balzano, Needell ICML '24]

(Inexact RGD)

$$\theta_{n+1} \leftarrow \widehat{\text{Rtr}}_{\theta_n} (-\alpha_n \widehat{\text{grad}} f(\theta_n))$$

- **It's OK to use approximate grad and Retraction**

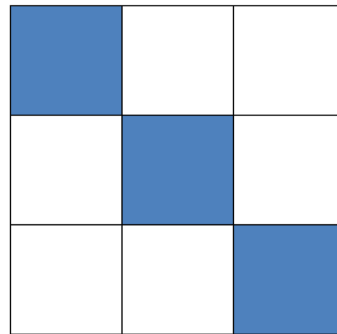
(Both mathematically and practically)



A coordinate-wise approach:

Adaptive Block Coordinate Descent

Optimize one block at a time!



$$\text{(PGD)} \quad \boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} (\boldsymbol{\theta}_n - \alpha_n \nabla f(\boldsymbol{\theta}_n))$$

- α_n : Stepsizes. How to choose them?

$$\text{(PGD)} \quad \boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} (\boldsymbol{\theta}_n - \alpha_n \nabla f(\boldsymbol{\theta}_n))$$

- α_n : Stepsizes. How to choose them?
 - "Small enough stepsize": $\alpha_n \leq 1/L$, where
 - $L =$ **Lipschitz constant for ∇f over Θ**
 \approx **Largest absolute eigenvalue of $\nabla^2 f$ over Θ**

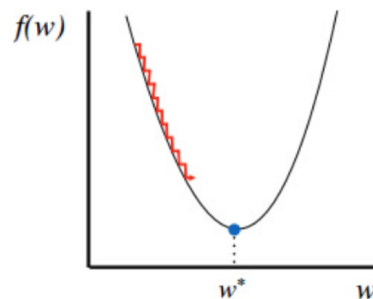
$$\text{(PGD)} \quad \boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} (\boldsymbol{\theta}_n - \alpha_n \nabla f(\boldsymbol{\theta}_n))$$

- α_n : Stepsizes. How to choose them?
 - "Small enough stepsize": $\alpha_n \leq 1/L$, where
 - $L =$ **Lipschitz constant for ∇f over Θ**
 \approx **Largest absolute eigenvalue of $\nabla^2 f$ over Θ**
 - But $\alpha_n = 1/L$ is TOO SMALL!
 - In practice use "line search" to find larger α_n that works
 - L might be unknown and hard to estimate

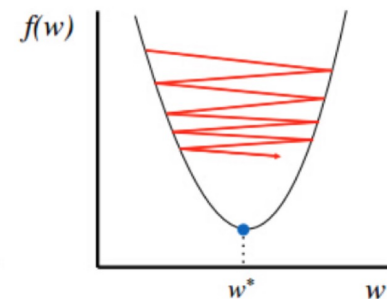
$$\text{(PGD)} \quad \theta_{n+1} \leftarrow \Pi_{\Theta} (\theta_n - \alpha_n \nabla f(\theta_n))$$

- α_n : Stepsizes. How to choose them?
 - "Small enough stepsize": $\alpha_n \leq 1/L$, where
 - $L =$ Lipschitz constant for ∇f over Θ
 \approx Largest absolute eigenvalue of $\nabla^2 f$ over Θ
 - But $\alpha_n = 1/L$ is TOO SMALL!
 - In practice use "line search" to find larger α_n that works
 - L might be unknown and hard to estimate

In practice, performance of P(S)GD depends **very sensitively** on α_n s



Too small: converge very slowly



Too big: overshoot and even diverge

Two-Block structure : $\theta = (A, B)$, $\Theta = \Theta_A \times \Theta_B$

$$\begin{array}{ll} \text{(Block PGD)} & A_{n+1} \leftarrow \Pi_{\Theta_A} (A_n - \alpha_n \nabla_A f(A_n, B_n)) \\ \text{(or BCD)} & B_{n+1} \leftarrow \Pi_{\Theta_B} (B_n - \beta_n \nabla_B f(A_{n+1}, B_n)) \end{array}$$

Two-Block structure : $\theta = (A, B)$, $\Theta = \Theta_A \times \Theta_B$

$$\begin{array}{ll} \text{(Block PGD)} & A_{n+1} \leftarrow \Pi_{\Theta_A} (A_n - \alpha_n \nabla_A f(A_n, B_n)) \\ \text{(or BCD)} & B_{n+1} \leftarrow \Pi_{\Theta_B} (B_n - \beta_n \nabla_B f(A_{n+1}, B_n)) \end{array}$$

- Known to be much more robust against stepsize choices than PGD (**Why??**)
 - e.g., low-rank matrix factorization, dictionary learning, tensor factorization, kernel learning, linear tranciever design

Two-Block structure : $\theta = (A, B)$, $\Theta = \Theta_A \times \Theta_B$

$$\begin{array}{ll} \text{(Block PGD)} & A_{n+1} \leftarrow \Pi_{\Theta_A} (A_n - \alpha_n \nabla_A f(A_n, B_n)) \\ \text{(or BCD)} & B_{n+1} \leftarrow \Pi_{\Theta_B} (B_n - \beta_n \nabla_B f(A_{n+1}, B_n)) \end{array}$$

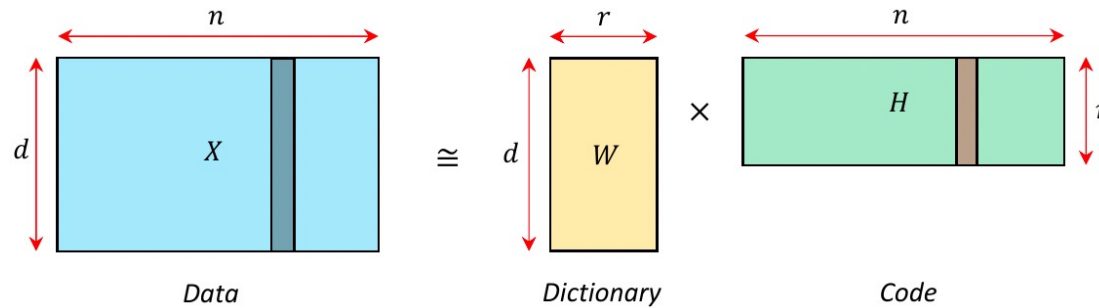
- Known to be much more robust against stepsize choices than PGD (**Why??**)
 - e.g., low-rank matrix factorization, dictionary learning, tensor factorization, kernel learning, linear tranciever design
- [Large $L \Leftrightarrow \nabla f$ changes wildly]
 - **Sensitive dependence on stepsize**

Two-Block structure : $\theta = (A, B)$, $\Theta = \Theta_A \times \Theta_B$

$$\begin{array}{ll} \text{(Block PGD)} & A_{n+1} \leftarrow \Pi_{\Theta_A} (A_n - \alpha_n \nabla_A f(A_n, B_n)) \\ \text{(or BCD)} & B_{n+1} \leftarrow \Pi_{\Theta_B} (B_n - \beta_n \nabla_B f(A_{n+1}, B_n)) \end{array}$$

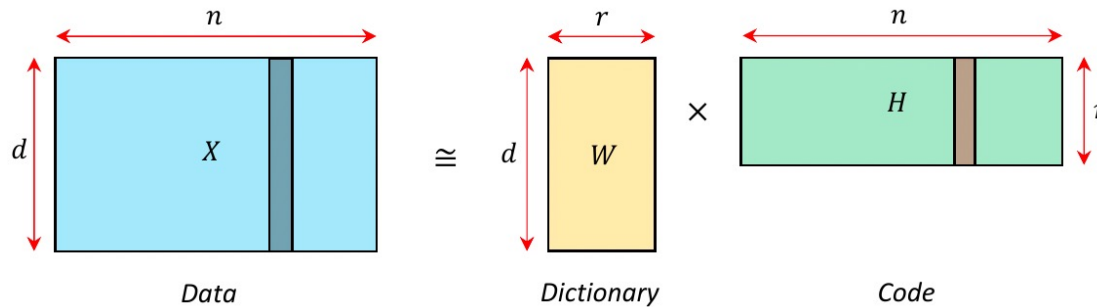
- Known to be much more robust against stepsize choices than PGD (**Why??**)
 - e.g., low-rank matrix factorization, dictionary learning, tensor factorization, kernel learning, linear tranciever design
- [Large $L \Leftrightarrow \nabla f$ changes wildly]
 - **Sensitive dependence on stepsize**
- Exploiting **block structure** → **The “effective L” is reduced**

Motivating example : **Nonnegative Matrix Factorization (Lee & Seung, Nature '99)**



$$\left\{ \begin{array}{l} \text{minimize} \quad f(W, H) = \|X - WH\|_F^2 \\ \text{subject to} \quad W \in \mathbb{R}_{\geq 0}^{d \times r}, H \in \mathbb{R}_{\geq 0}^{r \times n} \end{array} \right. \quad \begin{array}{l} \text{(Reconstruction error)} \\ \text{(Constraints)} \end{array}$$

Motivating example : Nonnegative Matrix Factorization (Lee & Seung, Nature '99)



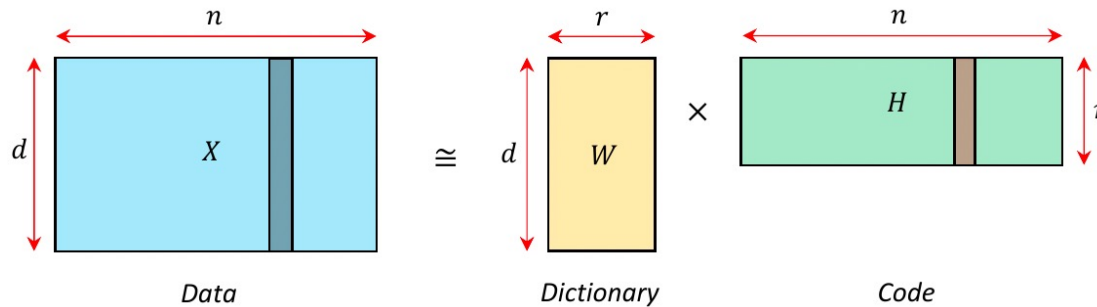
$$\begin{cases} \text{minimize} & f(W, H) = \|X - WH\|_F^2 & \text{(Reconstruction error)} \\ \text{subject to} & W \in \mathbb{R}_{\geq 0}^{d \times r}, H \in \mathbb{R}_{\geq 0}^{r \times n} & \text{(Constraints)} \end{cases}$$

PGD $\begin{cases} (W_{n+1}, H_{n+1}) \leftarrow (W_n, H_n) - \alpha_n (\nabla_W f(W_n, H_n), \nabla_H f(W_n, H_n)) \\ (W_{n+1}, H_{n+1}) \leftarrow \max(0, (W_{n+1}, H_{n+1})) \end{cases}$

- $f(W, H) = \text{Bi-convex}$
- $\nabla_W f(W, H) = (WH - X)H^T$
- $\nabla_H f(W, H) = W^T(WH - X)$

(Almost no one uses this 😊)

Motivating example : Nonnegative Matrix Factorization (Lee & Seung, Nature '99)



$$\begin{cases} \text{minimize} & f(W, H) = \|X - WH\|_F^2 & \text{(Reconstruction error)} \\ \text{subject to} & W \in \mathbb{R}_{\geq 0}^{d \times r}, H \in \mathbb{R}_{\geq 0}^{r \times n} & \text{(Constraints)} \end{cases}$$

PGD
$$\begin{cases} (W_{n+1}, H_{n+1}) \leftarrow (W_n, H_n) - \alpha_n (\nabla_W f(W_n, H_n), \nabla_H f(W_n, H_n)) \\ (W_{n+1}, H_{n+1}) \leftarrow \max(0, (W_{n+1}, H_{n+1})) \end{cases}$$

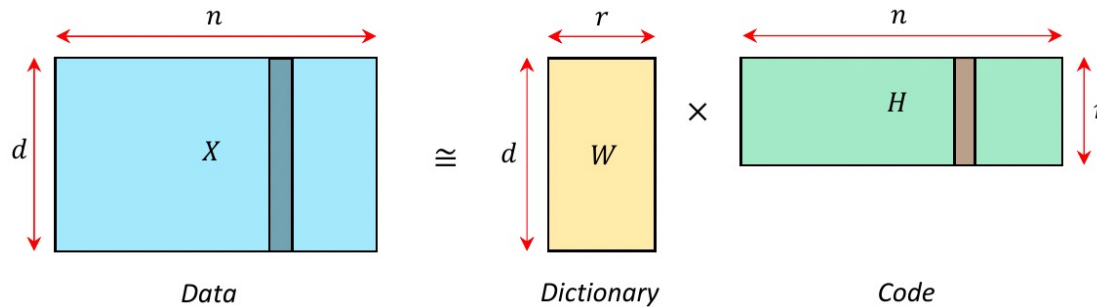
- $f(W, H) = \text{Bi-convex}$
- $\nabla_W f(W, H) = (WH - X)H^T$
- $\nabla_H f(W, H) = W^T(WH - X)$

(Almost no one uses this 😊)

$$\bullet \quad \nabla^2 f = \begin{matrix} \text{vec}(\mathbf{W}) \\ \text{vec}(\mathbf{H}) \end{matrix} \begin{bmatrix} \text{vec}(\mathbf{W})^T & \text{vec}(\mathbf{H})^T \\ \mathbf{H}\mathbf{H}^T \otimes \mathbf{I}_p & A_{12} \\ A_{12}^T & \mathbf{I}_n \otimes \mathbf{W}^T \mathbf{W} \end{bmatrix}$$

$$A_{12} = [(\mathbf{H} \otimes \mathbf{W}) + \mathbf{I}_r \otimes (\mathbf{W}\mathbf{H} - \mathbf{X})] \mathbf{C}^{(r,n)}$$

Motivating example : Nonnegative Matrix Factorization (Lee & Seung, Nature '99)



$$\begin{cases} \text{minimize} & f(W, H) = \|X - WH\|_F^2 & \text{(Reconstruction error)} \\ \text{subject to} & W \in \mathbb{R}_{\geq 0}^{d \times r}, H \in \mathbb{R}_{\geq 0}^{r \times n} & \text{(Constraints)} \end{cases}$$

PGD $\begin{cases} (W_{n+1}, H_{n+1}) \leftarrow (W_n, H_n) - \alpha_n (\nabla_W f(W_n, H_n), \nabla_H f(W_n, H_n)) \\ (W_{n+1}, H_{n+1}) \leftarrow \max(0, (W_{n+1}, H_{n+1})) \end{cases}$

- $f(W, H) = \text{Bi-convex}$
- $\nabla_W f(W, H) = (WH - X)H^T$
- $\nabla_H f(W, H) = W^T(WH - X)$

(Almost no one uses this 😊)

$$\bullet \nabla^2 f = \begin{matrix} \text{vec}(\mathbf{W}) \\ \text{vec}(\mathbf{H}) \end{matrix} \begin{bmatrix} \text{vec}(\mathbf{W})^T & \text{vec}(\mathbf{H})^T \\ \mathbf{H}\mathbf{H}^T \otimes \mathbf{I}_p & A_{12} \\ A_{12}^T & \mathbf{I}_n \otimes \mathbf{W}^T \mathbf{W} \end{bmatrix}$$

$$A_{12} = [(\mathbf{H} \otimes \mathbf{W}) + \mathbf{I}_r \otimes (\mathbf{W}\mathbf{H} - \mathbf{X})] \mathbf{C}^{(r,n)}$$

$L = \text{Max eval over all } (W, H)$

$= \text{Unbounded!}$

Adaptive BCD:

$$\begin{aligned}\mathbf{W} &\leftarrow \Pi \left(\mathbf{W} - \frac{1}{\lambda_{\max}(\mathbf{H}\mathbf{H}^T) + \varepsilon} (\mathbf{W}\mathbf{H} - \mathbf{X})\mathbf{H}^T \right), \\ \mathbf{H} &\leftarrow \Pi' \left(\mathbf{H} - \frac{1}{\lambda_{\max}(\mathbf{W}\mathbf{W}^T) + \varepsilon} \mathbf{W}^T(\mathbf{W}\mathbf{H} - \mathbf{X}) \right)\end{aligned}$$

Adaptive BCD:

$$\begin{aligned} \mathbf{W} &\leftarrow \Pi \left(\mathbf{W} - \frac{1}{\lambda_{\max}(\mathbf{H}\mathbf{H}^T) + \varepsilon} (\mathbf{W}\mathbf{H} - \mathbf{X})\mathbf{H}^T \right), \\ \mathbf{H} &\leftarrow \Pi' \left(\mathbf{H} - \frac{1}{\lambda_{\max}(\mathbf{W}\mathbf{W}^T) + \varepsilon} \mathbf{W}^T(\mathbf{W}\mathbf{H} - \mathbf{X}) \right) \end{aligned}$$

$$\bullet \quad \nabla^2 f = \begin{matrix} \text{vec}(\mathbf{W}) \\ \text{vec}(\mathbf{H}) \end{matrix} \begin{bmatrix} \text{vec}(\mathbf{W})^T & \text{vec}(\mathbf{H})^T \\ \mathbf{H}\mathbf{H}^T \otimes \mathbf{I}_p & A_{12} \\ A_{12}^T & \mathbf{I}_n \otimes \mathbf{W}^T\mathbf{W} \end{bmatrix}$$

$\nabla_{\mathbf{W}}^2 f(\cdot, \mathbf{H})$ $\nabla_{\mathbf{H}}^2 f(\mathbf{W}, \cdot)$

Adaptive BCD:

$$\begin{aligned} \mathbf{W} &\leftarrow \Pi \left(\mathbf{W} - \frac{1}{\lambda_{\max}(\mathbf{H}\mathbf{H}^T) + \varepsilon} (\mathbf{W}\mathbf{H} - \mathbf{X})\mathbf{H}^T \right), \\ \mathbf{H} &\leftarrow \Pi' \left(\mathbf{H} - \frac{1}{\lambda_{\max}(\mathbf{W}\mathbf{W}^T) + \varepsilon} \mathbf{W}^T(\mathbf{W}\mathbf{H} - \mathbf{X}) \right) \end{aligned}$$

$$\bullet \quad \nabla^2 f = \begin{matrix} \text{vec}(\mathbf{W}) \\ \text{vec}(\mathbf{H}) \end{matrix} \begin{bmatrix} \text{vec}(\mathbf{W})^T & \text{vec}(\mathbf{H})^T \\ \mathbf{H}\mathbf{H}^T \otimes \mathbf{I}_p & A_{12} \\ A_{12}^T & \mathbf{I}_n \otimes \mathbf{W}^T\mathbf{W} \end{bmatrix}$$

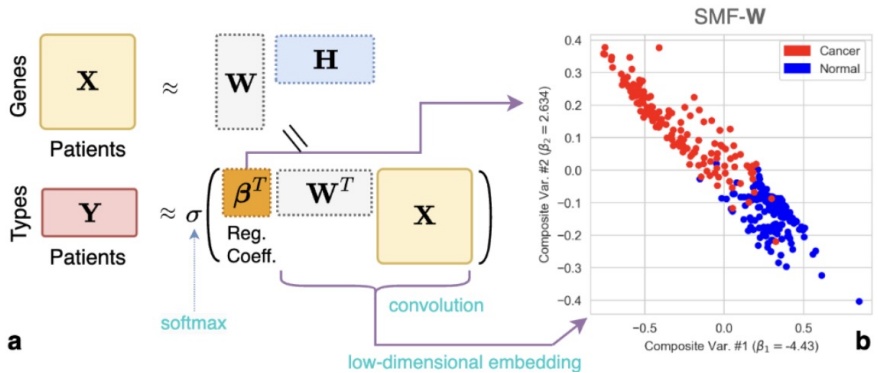
$\nabla_{\mathbf{W}}^2 f(\cdot, \mathbf{H})$ $\nabla_{\mathbf{H}}^2 f(\mathbf{W}, \cdot)$

Largest Eval of *diagonal blocks* of the Hessian

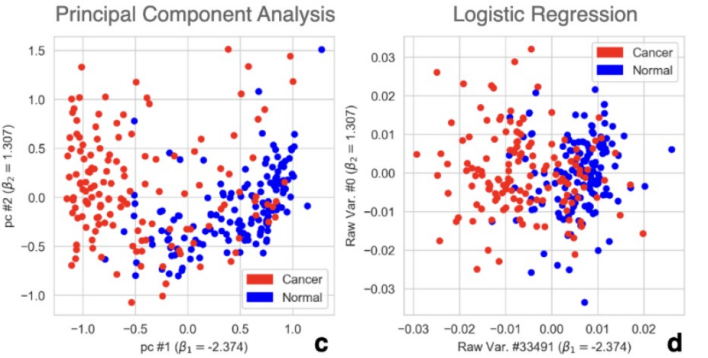
<< Largest Eval of the entire Hessian

Case study: Supervised Matrix Factorization (Lee, Lyu, Yao ICML '24)

Supervised Matrix Factorization (SMF-W)

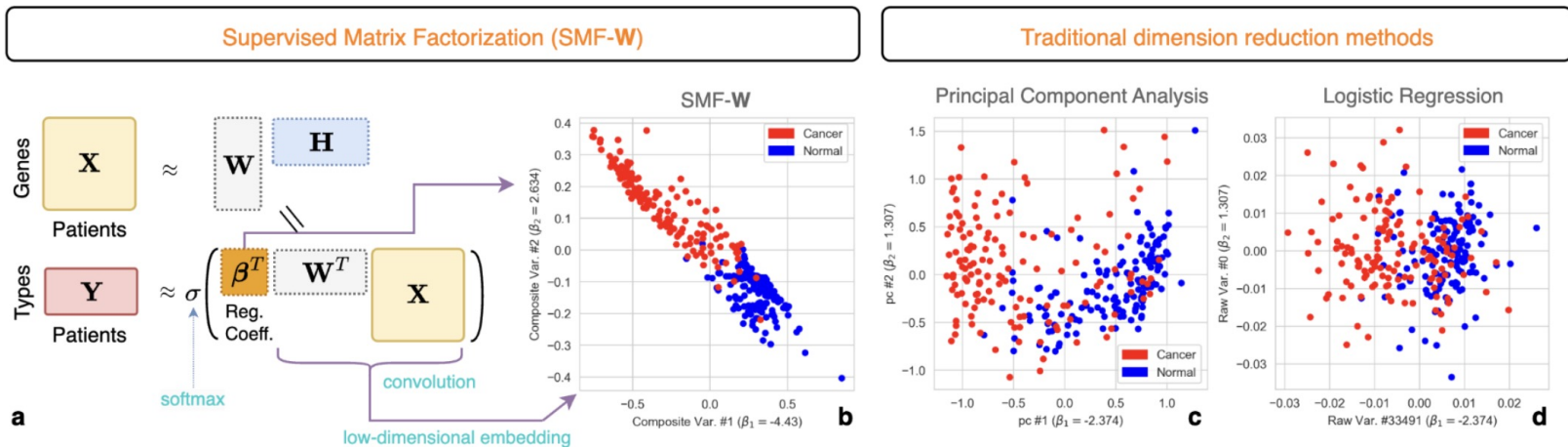


Traditional dimension reduction methods



Simultaneous dimension reduction and classification

Case study: Supervised Matrix Factorization (Lee, Lyu, Yao ICML '24)



Simultaneous dimension reduction and classification

$$\min_{\mathbf{W}, \mathbf{H}, \beta} f(\mathbf{W}, \mathbf{H}, \beta) := \sum_{i=1}^n \underbrace{\ell(y_i, \beta^T \mathbf{W}^T \mathbf{x}_i)}_{\text{Classification loss}} + \xi \underbrace{\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2}_{\text{Dimension reduction loss}}$$

Negative log-likelihood under logistic model

$$\ell(y, a) = \log(1 + \exp(a)) - \mathbf{1}_{\{y=1\}} a.$$

Parameters = Three factor matrices $\mathbf{W}, \mathbf{H}, \beta$

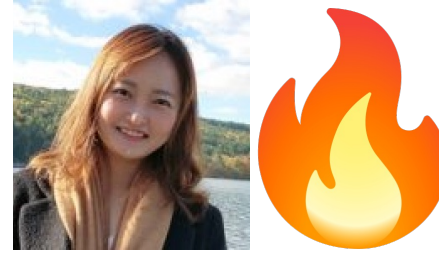
Block Coordinate Descent: Optimize one by one, not all at once

Case study: Supervised Matrix Factorization (Lee, Lyu, Yao ICML '24)

Algorithm 1 BCD for SMF-W

- 1: **Input:** $\mathbf{X} \in \mathbb{R}^{p \times n}$ (Data); $\mathbf{X}_{\text{aux}} \in \mathbb{R}^{q \times n}$ (Auxiliary covariate); $\mathbf{Y}_{\text{label}} \in \{0, \dots, \kappa\}^{1 \times n}$ (Label);
- 2: **Constraints:** Convex subsets $\mathcal{C}_1 \subseteq \mathbb{R}^{p \times r}$, $\mathcal{C}_2 \subseteq \mathbb{R}^{r \times n}$, $\mathcal{C}_3 \subseteq \mathbb{R}^{r \times \kappa}$, $\mathcal{C}_4 \subseteq \mathbb{R}^{q \times \kappa}$
- 3: **Parameters:** $\xi \geq 0$ (Tuning parameter); $T \in \mathbb{N}$ (number of iterations); $(\eta_{k;i})_{k \geq 1, 1 \leq i \leq 4}$ (step-sizes)
- 4: Initialize $\mathbf{W} \in \mathcal{C}_1$, $\mathbf{H} \in \mathcal{C}_2$, $\beta \in \mathcal{C}_3$, $\Gamma \in \mathcal{C}_4$
- 5: **For** $k = 1, 2, \dots, T$ **do:** (\triangleright For α^+ see 4.3.)
- 6: **(Update W)**
- 7: Update activation $\mathbf{a}_1, \dots, \mathbf{a}_n$ and \mathbf{K}
- 8: $\nabla_{\mathbf{W}} f(\mathbf{Z}) \leftarrow \mathbf{X}\mathbf{K}^T \beta^T + 2\xi(\mathbf{W}\mathbf{H} - \mathbf{X})\mathbf{H}^T$
- 9: Choose $\eta_{k,1}^{-1} > L_1 := \alpha^+ \|\beta\|_2^2 \cdot \|\mathbf{X}\|_2^2 + 2\xi \|\mathbf{H}\|_2^2$
- 10: $\mathbf{W} \leftarrow \Pi_{\mathcal{C}_1}(\mathbf{W} - \eta_{k,1} \nabla_{\mathbf{W}} f(\mathbf{Z}))$
- 11: **(Update H)**
- 12: $\nabla_{\mathbf{H}} f(\mathbf{Z}) \leftarrow 2\xi \mathbf{W}^T (\mathbf{W}\mathbf{H} - \mathbf{X})$
- 13: Choose $\eta_{k,2}^{-1} > L_2 := 2\xi \|\mathbf{W}\|_2^2$
- 14: $\mathbf{H} \leftarrow \Pi_{\mathcal{C}_2}(\mathbf{H} - \eta_{k,2} \nabla_{\mathbf{H}} f(\mathbf{Z}))$
- 15: **(Update β)**
- 16: Update activation $\mathbf{a}_1, \dots, \mathbf{a}_n$ and \mathbf{K}
- 17: $\nabla_{\beta} f(\mathbf{Z}) \leftarrow \mathbf{W}^T \mathbf{X}\mathbf{K}^T$
- 18: Choose $\eta_{k,3}^{-1} > L_3 := \alpha^+ \|\mathbf{W}\|_2^2 \cdot \|\mathbf{X}\|_2^2$
- 19: $\beta \leftarrow \Pi_{\mathcal{C}_3}(\beta - \eta_{k,3} \nabla_{\beta} f(\mathbf{Z}))$
- 20: **(Update Γ)**
- 21: Update activation $\mathbf{a}_1, \dots, \mathbf{a}_n$ and \mathbf{K}
- 22: $\nabla_{\Gamma} f(\mathbf{Z}) \leftarrow \mathbf{X}_{\text{aux}} \mathbf{K}^T$
- 23: Choose $\eta_{k,4}^{-1} > L_4 := \alpha^+ \|\mathbf{X}_{\text{aux}}\|_2^2$
- 24: $\Gamma \leftarrow \Pi_{\mathcal{C}_4}(\Gamma - \eta_{k,4} \nabla_{\Gamma} f(\mathbf{Z}))$
- 25: **End for**
- 26: **Output:** $\mathbf{Z} = (\mathbf{W}, \mathbf{H}, \beta, \Gamma)$

$$\nabla^2 f = \begin{matrix} \text{vec}(\mathbf{W}) \\ \text{vec}(\mathbf{H}) \\ \text{vec}(\beta) \\ \text{vec}(\Gamma) \end{matrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} & O \\ A_{21} & A_{22} & O & O \\ A_{31} & O & A_{33} & A_{34} \\ O & O & A_{43} & A_{44} \end{bmatrix}$$



Lemma D.2 (Derivatives of the SMF-W objective). Let $f(\mathbf{Z})$ denote the objective of SMF-W in (40). Suppose Assumption 4.3 holds. Recall \mathbf{h} and $\tilde{\mathbf{H}}$ defined in (44). Then the gradients of $f(\mathbf{Z})$ are given by

$$\nabla_{\mathbf{W}} f(\mathbf{Z}) = \mathbf{X}\mathbf{K}^T \beta^T + 2\xi(\mathbf{W}\mathbf{H} - \mathbf{X})\mathbf{H}^T, \quad (68)$$

$$\nabla_{\mathbf{H}} f(\mathbf{Z}) = 2\xi \mathbf{W}^T (\mathbf{W}\mathbf{H} - \mathbf{X}), \quad (69)$$

$$\nabla_{\beta} f(\mathbf{Z}) = \mathbf{W}^T \mathbf{X}\mathbf{K}^T, \quad (70)$$

$$\nabla_{\Gamma} f(\mathbf{Z}) = \mathbf{X}_{\text{aux}} \mathbf{K}^T. \quad (71)$$

The block-diagonal terms in the Hessian are given by

$$\nabla_{\text{vec}(\mathbf{W})} \nabla_{\text{vec}(\mathbf{W})}^T f(\mathbf{Z}) = (\beta \otimes \mathbf{X}) \mathbf{C}^{(\kappa, n)} \mathbf{M}\mathbf{C}^{(n, \kappa)} (\beta \otimes \mathbf{X})^T + 2\xi(\mathbf{H}\mathbf{H}^T \otimes \mathbf{I}_p), \quad (72)$$

$$\nabla_{\text{vec}(\mathbf{H})} \nabla_{\text{vec}(\mathbf{H})}^T f(\mathbf{Z}) = 2\xi(\mathbf{I}_r \otimes \mathbf{W}^T \mathbf{W}), \quad (73)$$

$$\nabla_{\text{vec}(\beta)} \nabla_{\text{vec}(\beta)}^T f(\mathbf{Z}) = (\mathbf{I}_{\kappa} \otimes \mathbf{W}^T \mathbf{X}) \mathbf{C}^{(\kappa, n)} \mathbf{M}\mathbf{C}^{(n, \kappa)} (\mathbf{I}_{\kappa} \otimes \mathbf{W}^T \mathbf{X})^T, \quad (74)$$

$$\nabla_{\text{vec}(\Gamma)} \nabla_{\text{vec}(\Gamma)}^T f(\mathbf{Z}) = (\mathbf{I}_{\kappa} \otimes \mathbf{X}_{\text{aux}}) \mathbf{C}^{(\kappa, n)} \mathbf{M}\mathbf{C}^{(n, \kappa)} (\mathbf{I}_{\kappa} \otimes \mathbf{X}_{\text{aux}})^T. \quad (75)$$

The block-off-diagonal terms in the Hessian are given by

$$\nabla_{\text{vec}(\mathbf{H})} \nabla_{\text{vec}(\mathbf{W})}^T f(\mathbf{Z}) = 2\xi \mathbf{C}^{(n, r)} [(\mathbf{H}^T \otimes \mathbf{W}^T) + \mathbf{I}_r \otimes (\mathbf{W}\mathbf{H} - \mathbf{X})^T], \quad (76)$$

$$\nabla_{\text{vec}(\beta)} \nabla_{\text{vec}(\mathbf{W})}^T f(\mathbf{Z}) = \mathbf{C}^{(\kappa, r)} (\mathbf{I}_r \otimes \mathbf{X}\mathbf{K}^T)^T + (\mathbf{I}_{\kappa} \otimes \mathbf{W}^T \mathbf{X}) \mathbf{C}^{(\kappa, n)} \mathbf{M}\mathbf{C}^{(n, \kappa)} (\beta \otimes \mathbf{X})^T, \quad (77)$$

$$\nabla_{\text{vec}(\Gamma)} \nabla_{\text{vec}(\mathbf{W})}^T f(\mathbf{Z}) = \nabla_{\text{vec}(\beta)} \nabla_{\text{vec}(\mathbf{H})}^T f(\mathbf{Z}) = \nabla_{\text{vec}(\Gamma)} \nabla_{\text{vec}(\mathbf{H})}^T f(\mathbf{Z}) = \mathbf{O}, \quad (78)$$

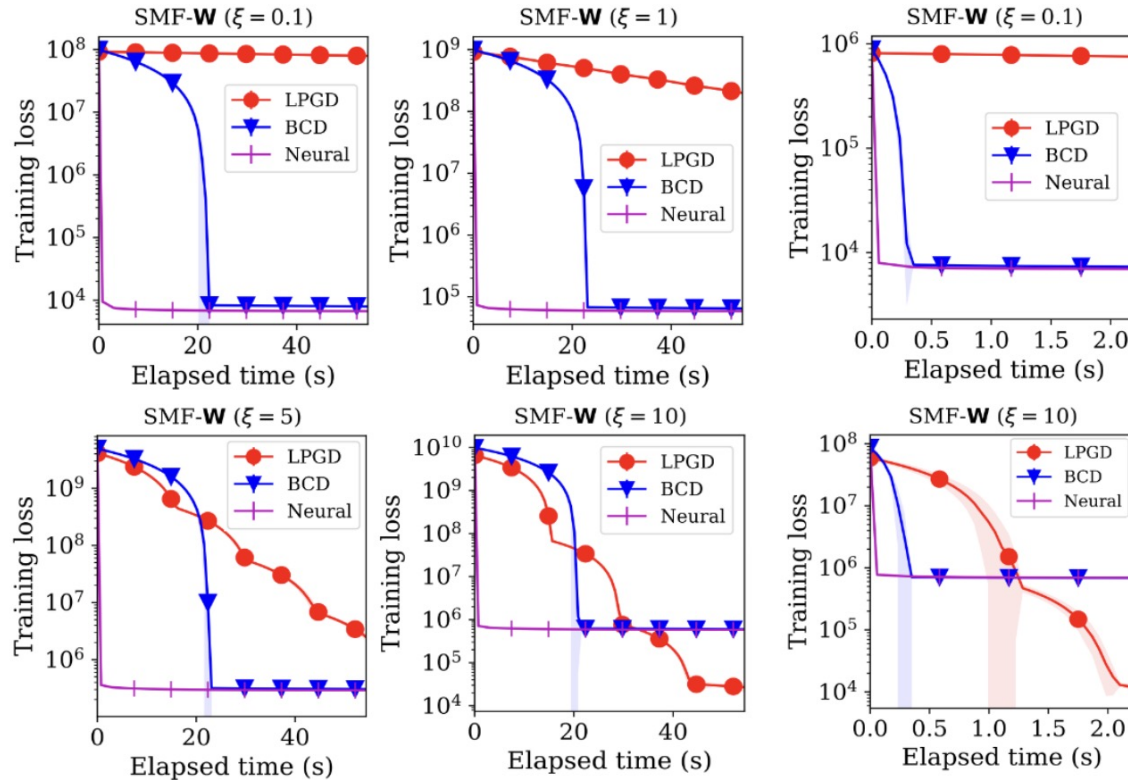
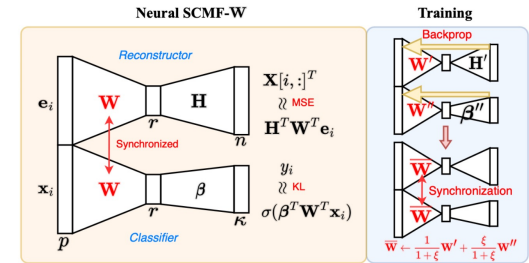
$$\nabla_{\text{vec}(\Gamma)} \nabla_{\text{vec}(\beta)}^T f(\mathbf{Z}) = (\mathbf{I}_r \otimes \mathbf{X}_{\text{aux}}) \mathbf{C}^{(\kappa, n)} \mathbf{M}\mathbf{C}^{(n, \kappa)} (\mathbf{I}_{\kappa} \otimes \mathbf{W}^T \mathbf{X})^T. \quad (79)$$

Case study: Supervised Matrix Factorization (Lee, Lyu, Yao ICML '24)

LPGD = Low-rank PGD [Lee, Lyu, Yao NeurIPS '23]

BCD = Adaptive BCD [Lee, Lyu, Yao ICML '24]

Neural = 2-year coupled NN implementation of BCD [Lee, Lyu, Yao ICML '24]

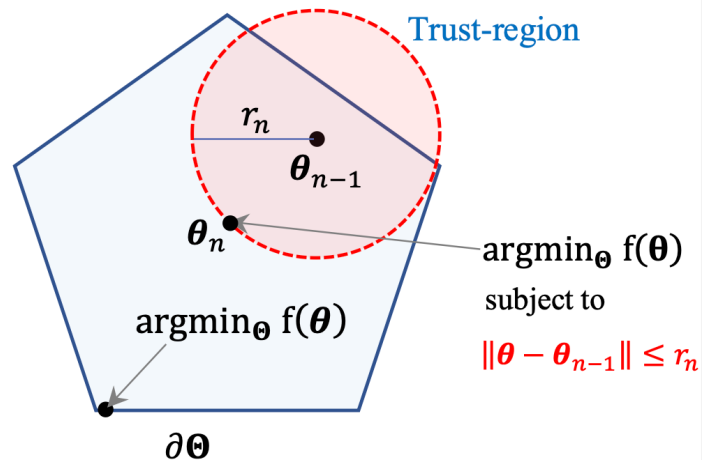


a Job postings

b Semi-synthetic MNIST data

A safeguarding approach:
Non-adaptive Trust-Region

Disciplined Agressiveness



How fast is BCD?

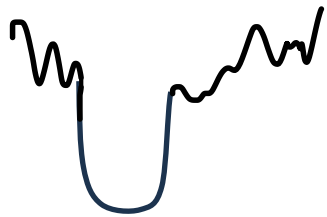
$$\begin{aligned} & A_{n+1} \leftarrow \Pi_{\Theta_A} (A_n - \alpha_n \nabla_A f(A_n, B_n, C_n)) \\ \text{(BCD)} \quad & B_{n+1} \leftarrow \Pi_{\Theta_B} (B_n - \beta_n \nabla_B f(A_{n+1}, B_n, C_n)) \\ & C_{n+1} \leftarrow \Pi_{\Theta_C} (C_n - \gamma_n \nabla_C f(A_{n+1}, B_{n+1}, C_n)) \end{aligned}$$

- Pros: **Robust against stepsize choices**

How fast is BCD?

$$\begin{aligned} \text{(BCD)} \quad & A_{n+1} \leftarrow \Pi_{\Theta_A} (A_n - \alpha_n \nabla_A f(A_n, B_n, C_n)) \\ & B_{n+1} \leftarrow \Pi_{\Theta_B} (B_n - \beta_n \nabla_B f(A_{n+1}, B_n, C_n)) \\ & C_{n+1} \leftarrow \Pi_{\Theta_C} (C_n - \gamma_n \nabla_C f(A_{n+1}, B_{n+1}, C_n)) \end{aligned}$$

- Pros: **Robust against stepsize choices**
- Cons: **Convergence rate could be slow for FLAT landscape**
 - (# of iterations to reach an ϵ -stationary points for PGD) = $O(L\epsilon^{-2})$



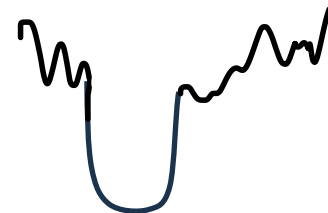
|

$$" \|\nabla f(\theta_n)\| \leq \epsilon "$$

How fast is BCD?

$$\begin{aligned} \text{(BCD)} \quad & A_{n+1} \leftarrow \Pi_{\Theta_A} (A_n - \alpha_n \nabla_A f(A_n, B_n, C_n)) \\ & B_{n+1} \leftarrow \Pi_{\Theta_B} (B_n - \beta_n \nabla_B f(A_{n+1}, B_n, C_n)) \\ & C_{n+1} \leftarrow \Pi_{\Theta_C} (C_n - \gamma_n \nabla_C f(A_{n+1}, B_{n+1}, C_n)) \end{aligned}$$

- Pros: **Robust against stepsize choices**
- Cons: **Convergence rate could be slow for FLAT landscape**



- (# of iterations to reach an ϵ -stationary points for PGD) = $O(L\epsilon^{-2})$

$$\text{" } \|\nabla f(\theta_n)\| \leq \epsilon \text{"}$$

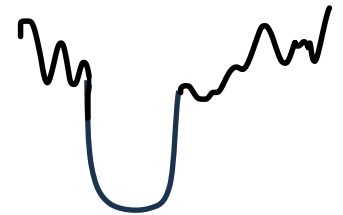
[Razaviyayn et al. '14, Lyu, Li '24+]

- (# of iterations to reach an ϵ -stationary points for **BCD**) = $O(L(1 + \rho^{-1})\epsilon^{-2})$

How fast is BCD?

$$\begin{aligned} \text{(BCD)} \quad & A_{n+1} \leftarrow \Pi_{\Theta_A} (A_n - \alpha_n \nabla_A f(A_n, B_n, C_n)) \\ & B_{n+1} \leftarrow \Pi_{\Theta_B} (B_n - \beta_n \nabla_B f(A_{n+1}, B_n, C_n)) \\ & C_{n+1} \leftarrow \Pi_{\Theta_C} (C_n - \gamma_n \nabla_C f(A_{n+1}, B_{n+1}, C_n)) \end{aligned}$$

- Pros: **Robust against stepsize choices**
- Cons: **Convergence rate could be slow for FLAT landscape**



- (# of iterations to reach an ϵ -stationary points for PGD) = $O(L\epsilon^{-2})$

$$\text{" } \|\nabla f(\theta_n)\| \leq \epsilon \text{"}$$

[Razaviyayn et al. '14, Lyu, Li '24+]

- (# of iterations to reach an ϵ -stationary points for **BCD**) = $O(L(1 + \rho^{-1})\epsilon^{-2})$

Largest allowed stepsize

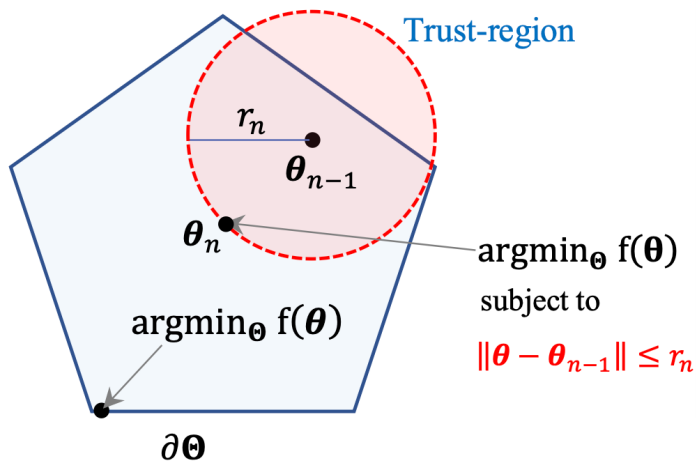
(ρ =Local smoothness parameter)

(Small $\rho \Leftrightarrow$ Flat local landscape \Leftrightarrow Large allowed stepsize)

$$\begin{aligned} & A_{n+1} \leftarrow \Pi_{\Theta_A \cap \{A: \|A - A_n\| \leq r_n\}} (A_n - \alpha_n \nabla_A f(A_n, B_n, C_n)) \\ \text{(BCD-DR)} \quad & B_{n+1} \leftarrow \Pi_{\Theta_B \cap \{B: \|B - B_n\| \leq r_n\}} (B_n - \beta_n \nabla_B f(A_{n+1}, B_n, C_n)) \\ & C_{n+1} \leftarrow \Pi_{\Theta_C \cap \{C: \|C - C_n\| \leq r_n\}} (C_n - \gamma_n \nabla_C f(A_{n+1}, B_{n+1}, C_n)) \end{aligned}$$

BCD with Diminishing Radius

$$\begin{aligned} A_{n+1} &\leftarrow \Pi_{\Theta_A \cap \{A: \|A - A_n\| \leq r_n\}} (A_n - \alpha_n \nabla_A f(A_n, B_n, C_n)) \\ \text{(BCD-DR)} \quad B_{n+1} &\leftarrow \Pi_{\Theta_B \cap \{B: \|B - B_n\| \leq r_n\}} (B_n - \beta_n \nabla_B f(A_{n+1}, B_n, C_n)) \\ C_{n+1} &\leftarrow \Pi_{\Theta_C \cap \{C: \|C - C_n\| \leq r_n\}} (C_n - \gamma_n \nabla_C f(A_{n+1}, B_{n+1}, C_n)) \end{aligned}$$



DR = Diminishing Radius: $r_n \sim c/\sqrt{n}$

(In classic Trust-Region literature, r_n is chosen adaptively)

[Lyu, Strohmeier, Needell JMLR '22]

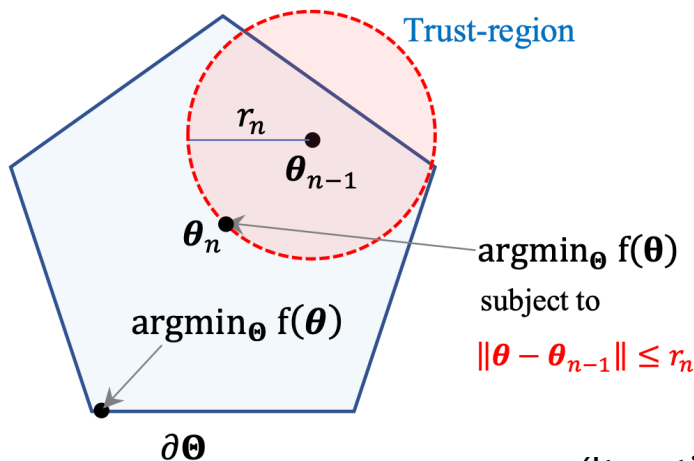
[Lyu, JMLR '24] [Powell, Lyu ICML'24]

(BCD-DR)

$$A_{n+1} \leftarrow \Pi_{\Theta_A \cap \{A: \|A - A_n\| \leq r_n\}} (A_n - \alpha_n \nabla_A f(A_n, B_n, C_n))$$

$$B_{n+1} \leftarrow \Pi_{\Theta_B \cap \{B: \|B - B_n\| \leq r_n\}} (B_n - \beta_n \nabla_B f(A_{n+1}, B_n, C_n))$$

$$C_{n+1} \leftarrow \Pi_{\Theta_C \cap \{C: \|C - C_n\| \leq r_n\}} (C_n - \gamma_n \nabla_C f(A_{n+1}, B_{n+1}, C_n))$$



DR = Diminishing Radius: $r_n \sim c/\sqrt{n}$

(In classic Trust-Region literature, r_n is chosen adaptively)

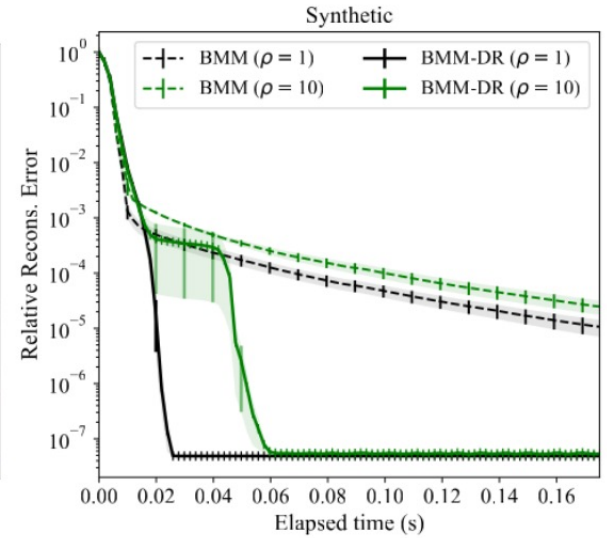
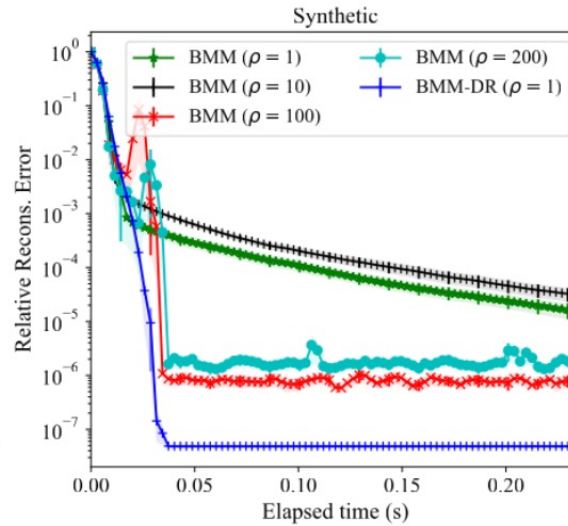
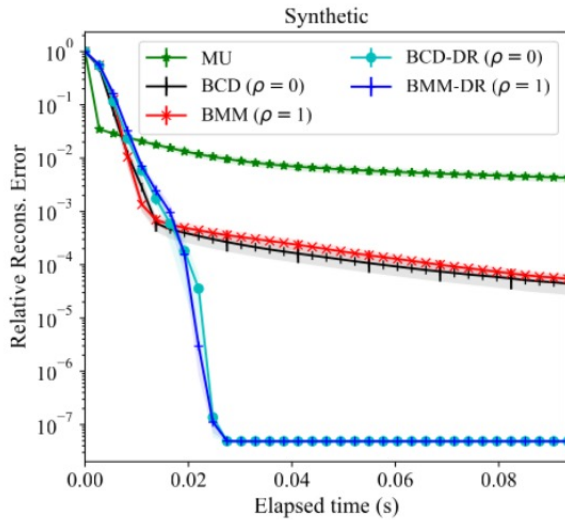
[Lyu, Strohmeier, Needell JMLR '22]

[Lyu, JMLR '24] [Powell, Lyu ICML'24]

- (Iteration complexity for PGD) = $O(L\epsilon^{-2})$
- (iteration complexity for BCD) = $O(L(1 + \rho^{-1})\epsilon^{-2})$
- (Iteration complexity for **BCD-DR**) = $O(L\epsilon^{-2})$ [Li, Lyu '24+]

No dependence on ρ

A very flat nonconvex problem



MU=Multiplicative Update

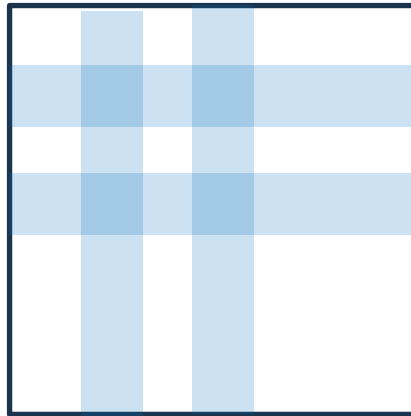
BCD = Block Coordinate Descent

BMM= Block Majorization-Minimization (in this case BCD+prox. Reg.)

A randomization approach:

Randomized Hessian Approximation

Flip the right coins!



What's known for the simplest algorithm?

$$\text{(QN)} \quad \boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}_n - \overbrace{(\nabla^2 f(\boldsymbol{\theta}_n) + \gamma_n \mathbf{I}_p)^{-1}}^{\mathbf{H}_n} \nabla f(\boldsymbol{\theta}_n)$$

What's known for the simplest algorithm?

$$\text{(QN)} \quad \boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}_n - \overbrace{(\nabla^2 f(\boldsymbol{\theta}_n) + \gamma_n \mathbf{I}_p)^{-1}}^{\mathbf{H}_n} \nabla f(\boldsymbol{\theta}_n)$$

- $\boldsymbol{\theta}_{n+1}$ can be found by solving the following linear eq.

$$(\nabla^2 f(\boldsymbol{\theta}_n) + \gamma_n \mathbf{I}_p) (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n) = \nabla f(\boldsymbol{\theta}_n)$$

What's known for the simplest algorithm?

$$\text{(QN)} \quad \theta_{n+1} \leftarrow \theta_n - \overbrace{(\nabla^2 f(\theta_n) + \gamma_n \mathbf{I}_p)^{-1}}^{\mathbf{H}_n} \nabla f(\theta_n)$$

- θ_{n+1} can be found by solving the following linear eq.

$$(\nabla^2 f(\theta_n) + \gamma_n \mathbf{I}_p) (\theta_{n+1} - \theta_n) = \nabla f(\theta_n)$$

- If we had a **Cholesky Factorization** $\mathbf{F}_n \mathbf{F}_n^T = \nabla^2 f(\theta_n) + \gamma_n \mathbf{I}_p$, then backward-forward substitution solves the above

What's known for the simplest algorithm?

$$\text{(QN)} \quad \theta_{n+1} \leftarrow \theta_n - \overbrace{(\nabla^2 f(\theta_n) + \gamma_n \mathbf{I}_p)^{-1}}^{\mathbf{H}_n} \nabla f(\theta_n)$$

- θ_{n+1} can be found by solving the following linear eq.

$$(\nabla^2 f(\theta_n) + \gamma_n \mathbf{I}_p) (\theta_{n+1} - \theta_n) = \nabla f(\theta_n)$$

- If we had a **Cholesky Factorization** $\mathbf{F}_n \mathbf{F}_n^T = \nabla^2 f(\theta_n) + \gamma_n \mathbf{I}_p$, then backward-forward substitution solves the above
- Cholesky factorization is super expensive $O(p^3)$ (same as inversion)

What's known for the simplest algorithm?

$$\text{(QN)} \quad \theta_{n+1} \leftarrow \theta_n - \overbrace{(\nabla^2 f(\theta_n) + \gamma_n \mathbf{I}_p)^{-1}}^{\mathbf{H}_n} \nabla f(\theta_n)$$

- θ_{n+1} can be found by solving the following linear eq.

$$(\nabla^2 f(\theta_n) + \gamma_n \mathbf{I}_p) (\theta_{n+1} - \theta_n) = \nabla f(\theta_n)$$

- If we had a **Cholesky Factorization** $\mathbf{F}_n \mathbf{F}_n^T = \nabla^2 f(\theta_n) + \gamma_n \mathbf{I}_p$, then backward-forward substitution solves the above
- Cholesky factorization is super expensive $O(p^3)$ (same as inversion)
- A recent, **randomized low-rank (k) Cholesky** is cheap $O(k^2 p)$

[Chen, Epperly, Tropp, and Webber '22+]

What's known for the simplest algorithm?

$$\text{(QN)} \quad \boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}_n - \overbrace{(\nabla^2 f(\boldsymbol{\theta}_n) + \gamma_n \mathbf{I}_p)^{-1}}^{\mathbf{H}_n} \nabla f(\boldsymbol{\theta}_n)$$

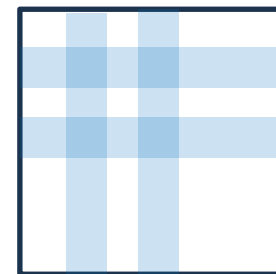
- $\boldsymbol{\theta}_{n+1}$ can be found by solving the following linear eq.

$$(\nabla^2 f(\boldsymbol{\theta}_n) + \gamma_n \mathbf{I}_p) (\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n) = \nabla f(\boldsymbol{\theta}_n)$$

- If we had a **Cholesky Factorization** $\mathbf{F}_n \mathbf{F}_n^T = \nabla^2 f(\boldsymbol{\theta}_n) + \gamma_n \mathbf{I}_p$, then backward-forward substitution solves the above
- Cholesky factorization is super expensive $O(p^3)$ (same as inversion)
- A recent, **randomized low-rank (k) Cholesky** is cheap $O(k^2 p)$

[Chen, Epperly, Tropp, and Webber '22+]

Idea: Sample rows and columns of the Hessian with prob. \propto diag(Hessian)



$$\text{(RLQN)} \quad \begin{cases} \mathbf{F}_n \mathbf{F}_n^T \approx \nabla^2 f(\boldsymbol{\theta}_n) + \tau_n \mathbf{I}_N & (\triangleright \text{Randomized low-rank Hessian approx.}) \\ \mathbf{H}_n \leftarrow \alpha_n (\mathbf{F}_n \mathbf{F}_n^T + \delta_n \mathbf{I}_N)^{-1} & (\triangleright \text{Preconditioning matrix}) \\ \boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}_n - \mathbf{H}_n \nabla f(\boldsymbol{\theta}_n) & (\triangleright \text{Parameter update}). \end{cases} \quad [\text{Duan, Lyu '24+}]$$

$$\text{(RLQN)} \quad \begin{cases} \mathbf{F}_n \mathbf{F}_n^T \approx \nabla^2 f(\boldsymbol{\theta}_n) + \tau_n \mathbf{I}_N & (\triangleright \text{Randomized low-rank Hessian approx.}) \\ \mathbf{H}_n \leftarrow \alpha_n (\mathbf{F}_n \mathbf{F}_n^T + \delta_n \mathbf{I}_N)^{-1} & (\triangleright \text{Preconditioning matrix}) \\ \boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}_n - \mathbf{H}_n \nabla f(\boldsymbol{\theta}_n) & (\triangleright \text{Parameter update}). \end{cases} \quad [\text{Duan, Lyu '24+}]$$

- Per-iteration computational cost = $O(k^2 p)$ --- Nearly 1st-order for $k \ll p$

$$\text{(RLQN)} \quad \begin{cases} \mathbf{F}_n \mathbf{F}_n^T \approx \nabla^2 f(\boldsymbol{\theta}_n) + \tau_n \mathbf{I}_N & (\triangleright \text{Randomized low-rank Hessian approx.}) \\ \mathbf{H}_n \leftarrow \alpha_n (\mathbf{F}_n \mathbf{F}_n^T + \delta_n \mathbf{I}_N)^{-1} & (\triangleright \text{Preconditioning matrix}) \\ \boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}_n - \mathbf{H}_n \nabla f(\boldsymbol{\theta}_n) & (\triangleright \text{Parameter update}). \end{cases} \quad [\text{Duan, Lyu '24+}]$$

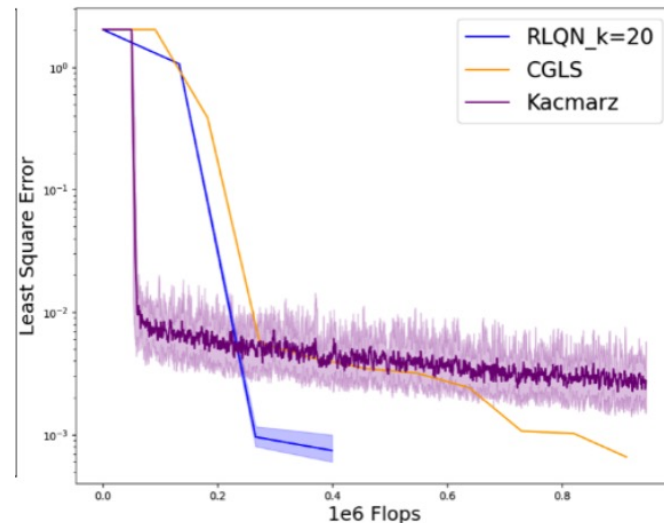
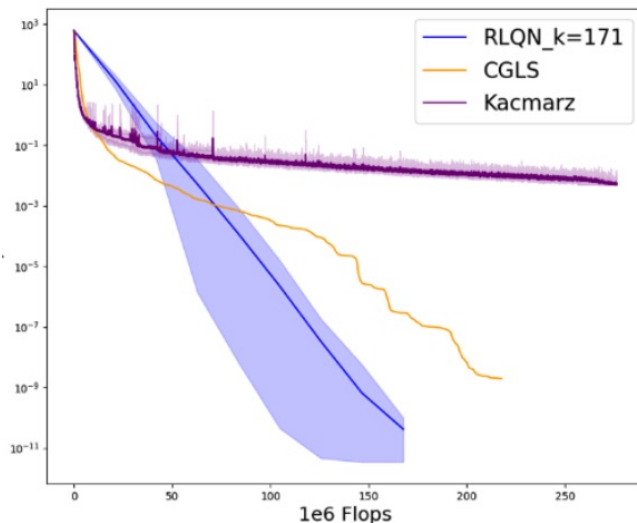
- Per-iteration computational cost = $O(k^2 p)$ --- Nearly 1st-order for $k \ll p$
- “Fast” local convergence rate (depending on “effective condition number”)

$$\text{(RLQN)} \quad \begin{cases} \mathbf{F}_n \mathbf{F}_n^T \approx \nabla^2 f(\boldsymbol{\theta}_n) + \tau_n \mathbf{I}_N & (\triangleright \text{Randomized low-rank Hessian approx.}) \\ \mathbf{H}_n \leftarrow \alpha_n (\mathbf{F}_n \mathbf{F}_n^T + \delta_n \mathbf{I}_N)^{-1} & (\triangleright \text{Preconditioning matrix}) \\ \boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}_n - \mathbf{H}_n \nabla f(\boldsymbol{\theta}_n) & (\triangleright \text{Parameter update}). \end{cases} \quad [\text{Duan, Lyu '24+}]$$

- Per-iteration computational cost = $O(k^2 p)$ --- Nearly 1st-order for $k \ll p$
- “Fast” local convergence rate (depending on “effective condition number”)
- Matching global convergence rate

$$\text{(RLQN)} \quad \begin{cases} \mathbf{F}_n \mathbf{F}_n^T \approx \nabla^2 f(\boldsymbol{\theta}_n) + \tau_n \mathbf{I}_N & (\triangleright \text{Randomized low-rank Hessian approx.}) \\ \mathbf{H}_n \leftarrow \alpha_n (\mathbf{F}_n \mathbf{F}_n^T + \delta_n \mathbf{I}_N)^{-1} & (\triangleright \text{Preconditioning matrix}) \\ \boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}_n - \mathbf{H}_n \nabla f(\boldsymbol{\theta}_n) & (\triangleright \text{Parameter update}). \end{cases} \quad [\text{Duan, Lyu '24+}]$$

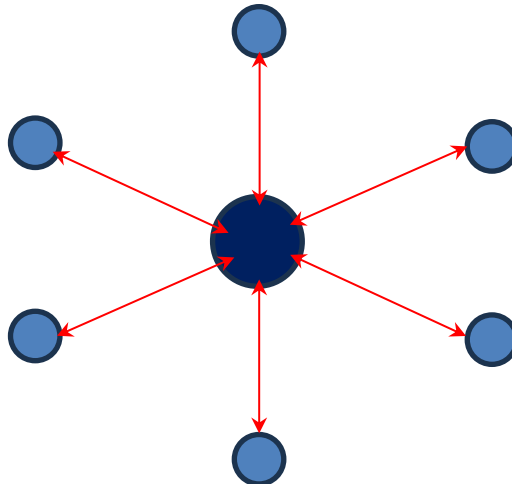
- Per-iteration computational cost = $O(k^2 p)$ --- Nearly 1st-order for $k \ll p$
- “Fast” local convergence rate (depending on “effective condition number”)
- Matching global convergence rate



A sampling approach:

Stochastic Aggregated Double-Averaging

Consensus, Consensus!



A finite-sum problem

$$\boldsymbol{\theta}^* \in \arg \min_{\boldsymbol{\theta} \in \Theta} \left\{ f(\boldsymbol{\theta}) := \sum_{v \in \mathcal{V}} f^v(\boldsymbol{\theta}) \right\}$$

A finite-sum problem

$$\boldsymbol{\theta}^* \in \arg \min_{\boldsymbol{\theta} \in \Theta} \left\{ f(\boldsymbol{\theta}) := \sum_{v \in \mathcal{V}} f^v(\boldsymbol{\theta}) \right\}$$

(PSGD)

$$v_n \sim \text{Uniform}(\mathcal{V}), \text{ independent}$$
$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} (\boldsymbol{\theta}_n - \alpha_n \nabla f^{v_n}(\boldsymbol{\theta}_n))$$

Memory

$$O(p)$$

Iteration Complexity

$$O(\epsilon^{-4})$$

A finite-sum problem

$$\boldsymbol{\theta}^* \in \arg \min_{\boldsymbol{\theta} \in \Theta} \left\{ f(\boldsymbol{\theta}) := \sum_{v \in \mathcal{V}} f^v(\boldsymbol{\theta}) \right\}$$

(PSGD)

$$v_n \sim \text{Uniform}(\mathcal{V}), \text{ independent}$$
$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} (\boldsymbol{\theta}_n - \alpha_n \nabla f^{v_n}(\boldsymbol{\theta}_n))$$

Memory

$$O(p)$$

Iteration Complexity

$$O(\epsilon^{-4})$$

(SAG)

$$v_n \sim \text{Uniform}(\mathcal{V}), \text{ independent}$$
$$\nabla^v \leftarrow \nabla f^v(\boldsymbol{\theta}_n), \text{ Keep other } \nabla^u \text{s}$$

$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} \left(\boldsymbol{\theta}_n - \alpha_n \frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \nabla^u \right)$$

$$O(pn)$$

$$O(\epsilon^{-2})$$

[Bottou '98]

A finite-sum problem

$$\boldsymbol{\theta}^* \in \arg \min_{\boldsymbol{\theta} \in \Theta} \left\{ f(\boldsymbol{\theta}) := \sum_{v \in \mathcal{V}} f^v(\boldsymbol{\theta}) \right\}$$

(PSGD)

$$v_n \sim \text{Uniform}(\mathcal{V}), \text{ independent}$$

$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} (\boldsymbol{\theta}_n - \alpha_n \nabla f^{v_n}(\boldsymbol{\theta}_n))$$

Memory

$$O(p)$$

Iteration Complexity

$$O(\epsilon^{-4})$$

(SAG)

$$v_n \sim \text{Uniform}(\mathcal{V}), \text{ independent}$$

$$\nabla^v \leftarrow \nabla f^v(\boldsymbol{\theta}_n), \text{ Keep other } \nabla^u \text{ s}$$

$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} \left(\boldsymbol{\theta}_n - \alpha_n \frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \nabla^u \right)$$

$$O(pn)$$

$$O(\epsilon^{-2})$$

[Bottou '98]

(MISO)

$$v_n \sim \text{Uniform}(\mathcal{V}), \text{ independent}$$

$$\boldsymbol{\theta}^{v_n} \leftarrow \boldsymbol{\theta}_n, \text{ Keep other } \boldsymbol{\theta}^u \text{ s}$$

$$\nabla^v \leftarrow \nabla f^v(\boldsymbol{\theta}_n), \text{ Keep other } \nabla^u \text{ s}$$

$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} \left(\frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \boldsymbol{\theta}^u - \alpha_n \frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \nabla^u \right)$$

$$O(pn)$$

$$O(\epsilon^{-2})$$

[Karimi et al. ICML'22]

[Mairal '13]

A finite-sum problem

$$\boldsymbol{\theta}^* \in \arg \min_{\boldsymbol{\theta} \in \Theta} \left\{ f(\boldsymbol{\theta}) := \sum_{v \in \mathcal{V}} f^v(\boldsymbol{\theta}) \right\}$$

(PSGD)

$$v_n \sim P(v_{n-1}, \cdot), \text{ Markovian}$$

$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} (\boldsymbol{\theta}_n - \alpha_n \nabla f^{v_n}(\boldsymbol{\theta}_n))$$

Memory

$$O(p)$$

Iteration Complexity

$$O(\epsilon^{-4})$$

(SAG)

$$v_n \sim P(v_{n-1}, \cdot), \text{ Markovian}$$

$$\nabla^v \leftarrow \nabla f^v(\boldsymbol{\theta}_n), \text{ Keep other } \nabla^u \text{s}$$

$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} \left(\boldsymbol{\theta}_n - \alpha_n \frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \nabla^u \right)$$

$$O(pn)$$

$$O(\epsilon^{-2})$$

[Even ICML '23]

(MISO)

$$v_n \sim P(v_{n-1}, \cdot), \text{ Markovian}$$

$$\boldsymbol{\theta}^{v_n} \leftarrow \boldsymbol{\theta}_n, \text{ Keep other } \boldsymbol{\theta}^u \text{s}$$

$$\nabla^v \leftarrow \nabla f^v(\boldsymbol{\theta}_n), \text{ Keep other } \nabla^u \text{s}$$

$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} \left(\frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \boldsymbol{\theta}^u - \alpha_n \frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \nabla^u \right)$$

$$O(pn)$$

??

A finite-sum problem

$$\boldsymbol{\theta}^* \in \arg \min_{\boldsymbol{\theta} \in \Theta} \left\{ f(\boldsymbol{\theta}) := \sum_{v \in \mathcal{V}} f^v(\boldsymbol{\theta}) \right\}$$

		Memory	Iteration Complexity
(PSGD)	$v_n \sim \text{Reshuffling}(\mathcal{V})$ $\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta}(\boldsymbol{\theta}_n - \alpha_n \nabla f^{v_n}(\boldsymbol{\theta}_n))$	$O(p)$	$O(\epsilon^{-3})$ [Mishchenko et al. '20] [Lu et al. '22]
(SAG)	$v_n \sim \text{Reshuffling}(\mathcal{V})$ $\nabla^v \leftarrow \nabla f^v(\boldsymbol{\theta}_n)$, Keep other ∇^u s $\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} \left(\boldsymbol{\theta}_n - \alpha_n \frac{1}{ \mathcal{V} } \sum_{u \in \mathcal{V}} \nabla^u \right)$	$O(pn)$??
(MISO)	$v_n \sim \text{Reshuffling}(\mathcal{V})$ $\boldsymbol{\theta}^{v_n} \leftarrow \boldsymbol{\theta}_n$, Keep other $\boldsymbol{\theta}^u$ s $\nabla^v \leftarrow \nabla f^v(\boldsymbol{\theta}_n)$, Keep other ∇^u s $\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} \left(\frac{1}{ \mathcal{V} } \sum_{u \in \mathcal{V}} \boldsymbol{\theta}^u - \alpha_n \frac{1}{ \mathcal{V} } \sum_{u \in \mathcal{V}} \nabla^u \right)$	$O(pn)$??

Minimization by Incremental Surrogate Optimization

$$\boldsymbol{\theta}^* \in \arg \min_{\boldsymbol{\theta} \in \Theta} \left\{ f(\boldsymbol{\theta}) := \sum_{v \in \mathcal{V}} f^v(\boldsymbol{\theta}) \right\}$$

(MISO)

$v_n \sim \text{Recurrent}(\mathcal{V})$

$\boldsymbol{\theta}^{v_n} \leftarrow \boldsymbol{\theta}_n$, Keep other $\boldsymbol{\theta}^u$ s

$\nabla^v \leftarrow \nabla f^v(\boldsymbol{\theta}_n)$, Keep other ∇^u s

$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} \left(\frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \boldsymbol{\theta}^u - \alpha_n \frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \nabla^u \right)$$

Iteration Complexity

$$O(\epsilon^{-2})$$

[Powell, Lyu ICML '24]

$$\boldsymbol{\theta}^* \in \arg \min_{\boldsymbol{\theta} \in \Theta} \left\{ f(\boldsymbol{\theta}) := \sum_{v \in \mathcal{V}} f^v(\boldsymbol{\theta}) \right\}$$

(MISO)

$v_n \sim \text{Recurrent}(\mathcal{V})$

$\boldsymbol{\theta}^{v_n} \leftarrow \boldsymbol{\theta}_n$, Keep other $\boldsymbol{\theta}^u$ s

$\nabla^v \leftarrow \nabla f^v(\boldsymbol{\theta}_n)$, Keep other ∇^u s

$$\boldsymbol{\theta}_{n+1} \leftarrow \Pi_{\Theta} \left(\frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \boldsymbol{\theta}^u - \alpha_n \frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \nabla^u \right)$$

Iteration Complexity

$$O(\epsilon^{-2})$$

[Powell, Lyu ICML '24]

- Recurrent sampling = Very general sampling scheme
e.g., i.i.d., Markovian, Recurrent, Cyclic, etc.

$$\boldsymbol{\theta}^* \in \arg \min_{\boldsymbol{\theta} \in \Theta} \left\{ f(\boldsymbol{\theta}) := \sum_{v \in \mathcal{V}} f^v(\boldsymbol{\theta}) \right\}$$

(MISO)

$$\begin{aligned} v_n &\sim \text{Recurrent}(\mathcal{V}) \\ \boldsymbol{\theta}^{v_n} &\leftarrow \boldsymbol{\theta}_n, \text{ Keep other } \boldsymbol{\theta}^u \text{ s} \\ \nabla^v &\leftarrow \nabla f^v(\boldsymbol{\theta}_n), \text{ Keep other } \nabla^u \text{ s} \\ \boldsymbol{\theta}_{n+1} &\leftarrow \Pi_{\Theta} \left(\frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \boldsymbol{\theta}^u - \alpha_n \frac{1}{|\mathcal{V}|} \sum_{u \in \mathcal{V}} \nabla^u \right) \end{aligned}$$

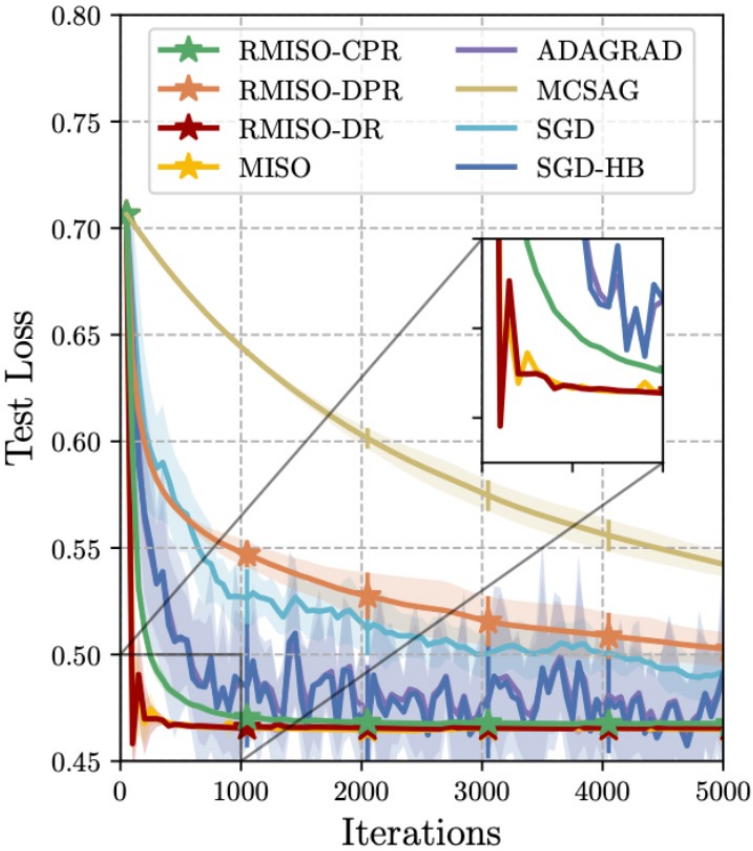
Iteration Complexity

$$O(\epsilon^{-2})$$

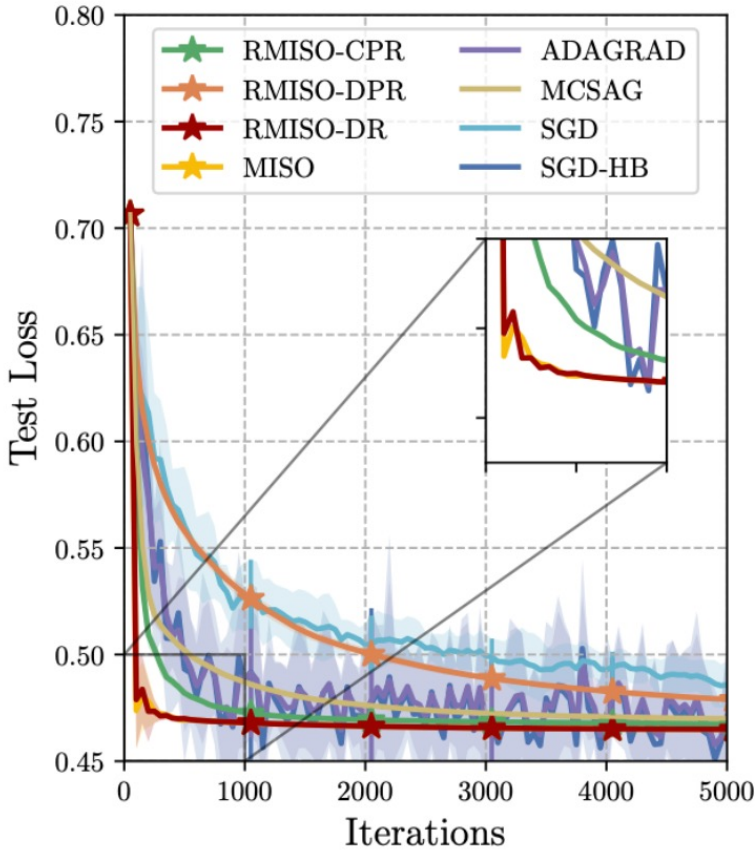
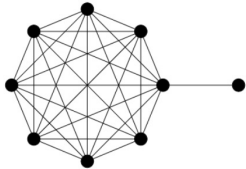
[Powell, Lyu ICML '24]

- Recurrent sampling = Very general sampling scheme
e.g., i.i.d., Markovian, Recurrent, Cyclic, etc.
- The implied constant in $O(\epsilon^{-2})$ is significantly improved

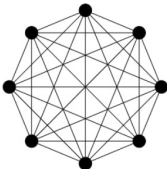
Minimization by Incremental Surrogate Optimization



(a) Lonely graph



(b) Complete graph



Takeaways

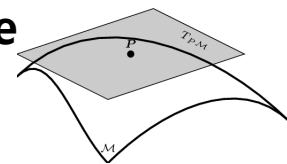
- Various optimization algorithms are based on simple mathematical principles

Takeaways

- Various optimization algorithms are based on simple mathematical principles
- Despite the extensive literature, there are many interesting challenges for improving these methods

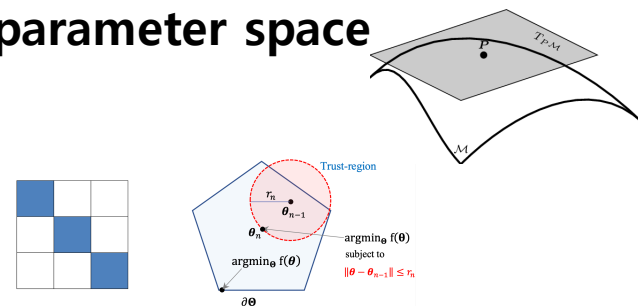
Takeaways

- Various optimization algorithms are based on simple mathematical principles
- Despite the extensive literature, there are many interesting challenges for improving these methods
 - **Leverage low-dim geometric structures in parameter space**
 - Riemannian Optimization



Takeaways

- Various optimization algorithms are based on simple mathematical principles
- Despite the extensive literature, there are many interesting challenges for improving these methods
 - **Leverage low-dim geometric structures in parameter space**
 - Riemannian Optimization
 - **Robustifying against stepsizes**
 - Block optimization and Trust-Region

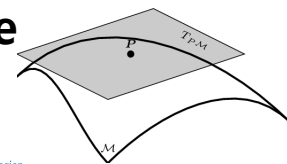


Takeaways

- Various optimization algorithms are based on simple mathematical principles
- Despite the extensive literature, there are many interesting challenges for improving these methods

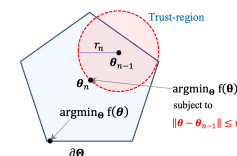
- **Leverage low-dim geometric structures in parameter space**

- Riemannian Optimization



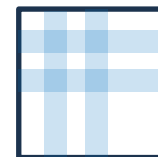
- **Robustifying against stepsizes**

- Block optimization and Trust-Region



- **Quasi-Newton "Lite"**

- Randomized low-rank Hessian approximation

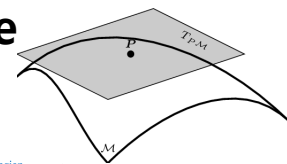


Takeaways

- Various optimization algorithms are based on simple mathematical principles
- Despite the extensive literature, there are many interesting challenges for improving these methods

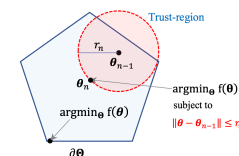
- **Leverage low-dim geometric structures in parameter space**

- Riemannian Optimization



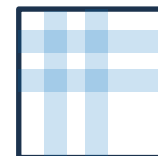
- **Robustifying against stepsizes**

- Block optimization and Trust-Region



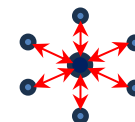
- **Quasi-Newton "Lite"**

- Randomized low-rank Hessian approximation



- **Finite-sum problem with arbitrary sampling**

- Aggregation with parameter and gradient averaging



Thank you very much!

References

1. Yifan Chen, Ethan N Epperly, Joel A Tropp, and Robert J Webber. "*Randomly pivoted cholesky: Practical approximation of a kernel matrix with few entry evaluations.*" ArXiv:2207.06503 (2022)
2. Joowon Lee, Hanbaek Lyu, and Weixin Yao, "*Exponentially Convergent Algorithms for Supervised Matrix Factorization*", NeurIPS 2023
3. Joowon Lee, Hanbaek Lyu, and Weixin Yao, "*Constrained Matrix Factorization: Local Landscape Analysis and Applications*" To appear in ICML 2024
4. Hanbaek Lyu, "*Stochastic regularized block majorization-minimization with weakly convex and multi-convex surrogates*" To appear in JMLR
5. Hanbaek Lyu, Christopher Strohmeier, and Deanna Needell, "*Online nonnegative tensor factorization and CP-Dictionary Learning for Markovian data*" JMLR 23(148):1–50, 2022
6. Hanbaek Lyu and Yuchen Li, "*Block majorization-minimization with diminishing radius for constrained nonconvex optimization*" Under Review in SIOPT
7. Yuchen Li, Laura Balzano, Deanna Needell, Hanbaek Lyu, "*Convergence and Complexity Guarantee for Inexact First-order Riemannian Optimization Algorithms.*" To appear in ICML 2024
8. William Powell and Hanbaek Lyu, "*Stochastic optimization with arbitrary recurrent data sampling.*" To appear in ICML 2024