

# Exponentially convergent algorithms for Supervised matrix factorization

Hanbaek Lyu

Department of Mathematics  
University of Wisconsin – Madison  
Institute for Foundations of Data Science

Joint work with  
Joowon Lee (UW Madison Statistics)  
Weixin Yao (UC Riverside Statistics)

SIAM MDS 2024

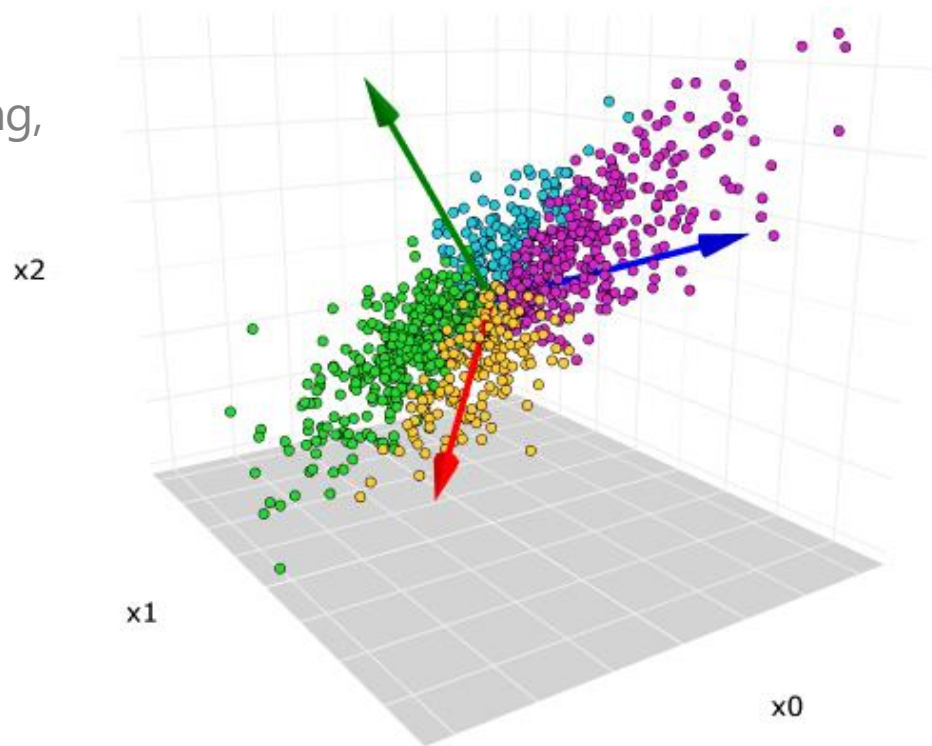
*Efficient and Robust Optimization Techniques for Structured Data Learning*

Oct. 22, 2024

# *Supervised Matrix Factorization*

## Dimension Reduction

- $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$  : high-dimensional data points, *presumably low-rank*
- How do we find the “best” low-dimensional representation?
- e.g., low-rank matrix factorization, dictionary learning, subspace learning, PCA
- PCA: Subspace that explains the most amount of data variance!

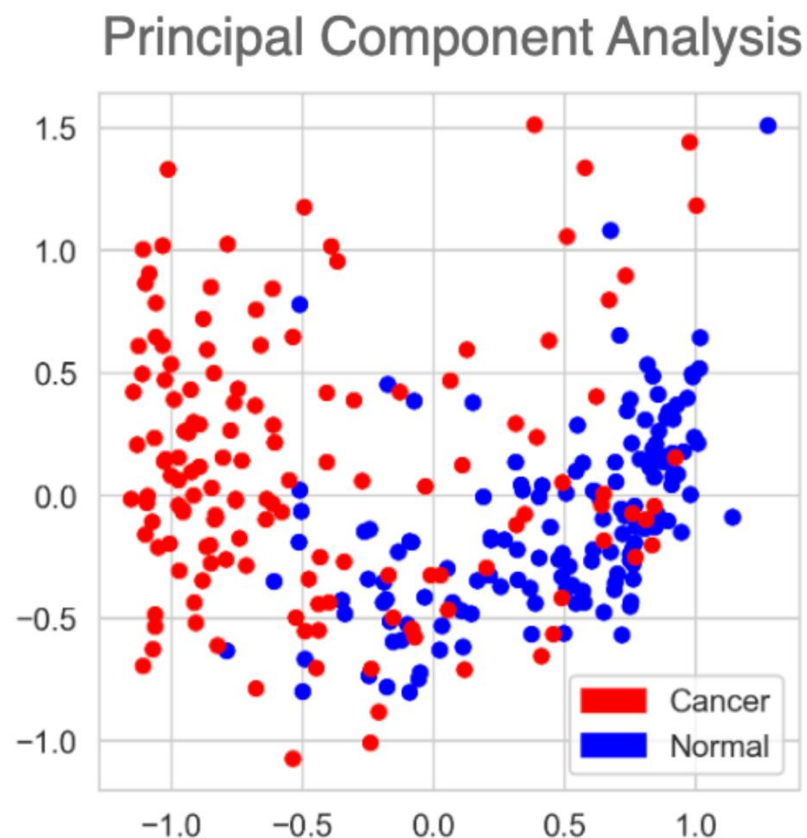


## Dimension Reduction

- $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$  : high-dimensional data points, *presumably low-rank*
- $y_1, \dots, y_n \in \{0,1\}$  : corresponding labels (e.g., cancer vs. normal)

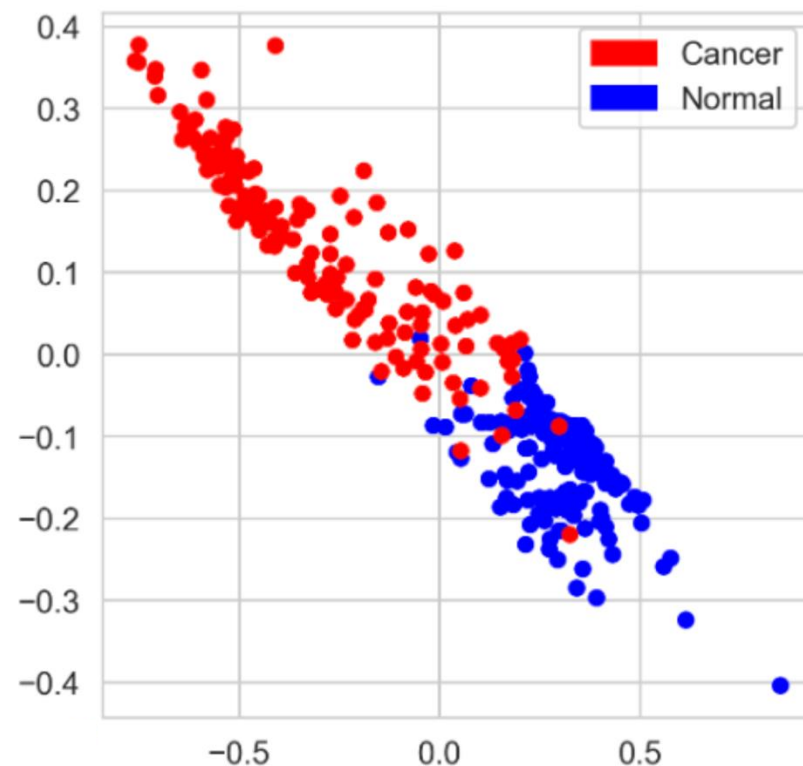
## Dimension Reduction

- $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$  : high-dimensional data points, *presumably low-rank*
- $y_1, \dots, y_n \in \{0,1\}$  : corresponding labels (e.g., cancer vs. normal)
- How do we find the “**best**” low-dimensional representation?
- “Explaining the most amount of data variance” — may not be the right criterion!

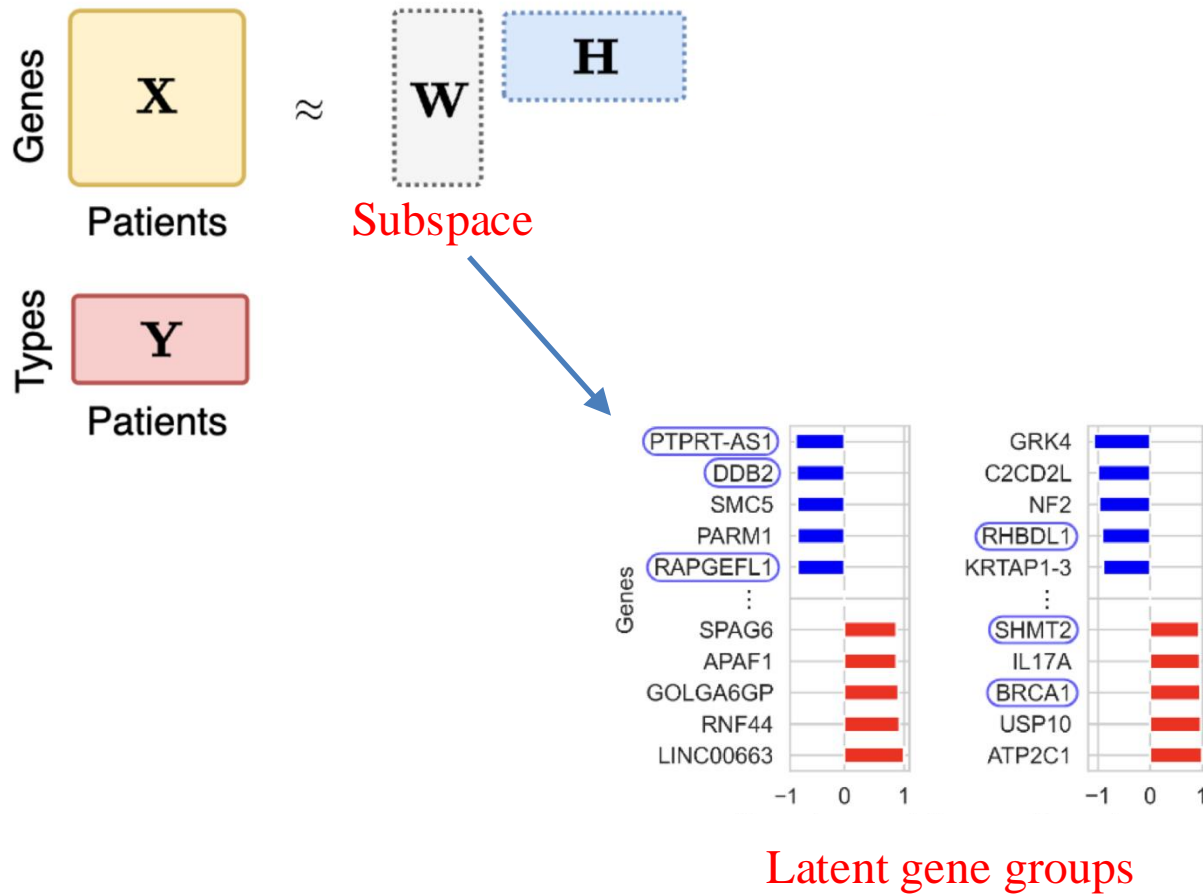


## Supervised Dimension Reduction

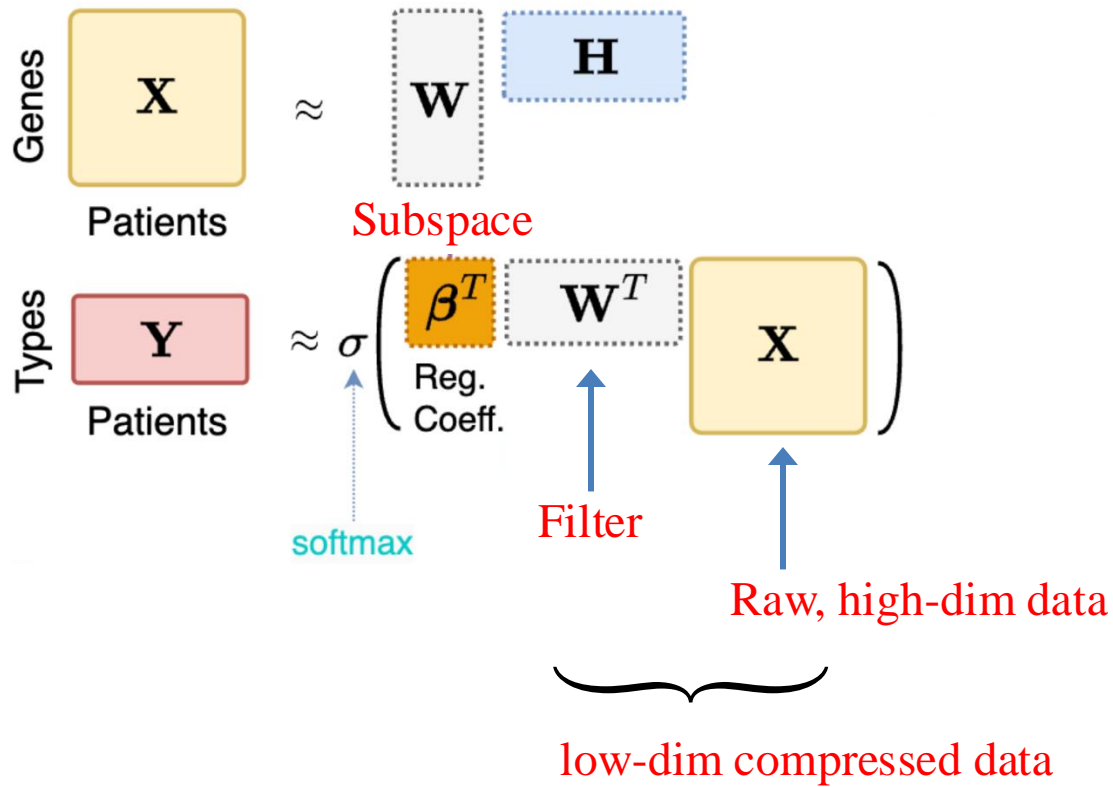
- $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$  : high-dimensional data points, *presumably low-rank*
- $y_1, \dots, y_n \in \{0,1\}$  : corresponding labels (e.g., cancer vs. normal)
- How do we find the “**best**” low-dimensional representation?
- Define a notion of “**label-aware low-rankness**”



## Supervised Matrix Factorization

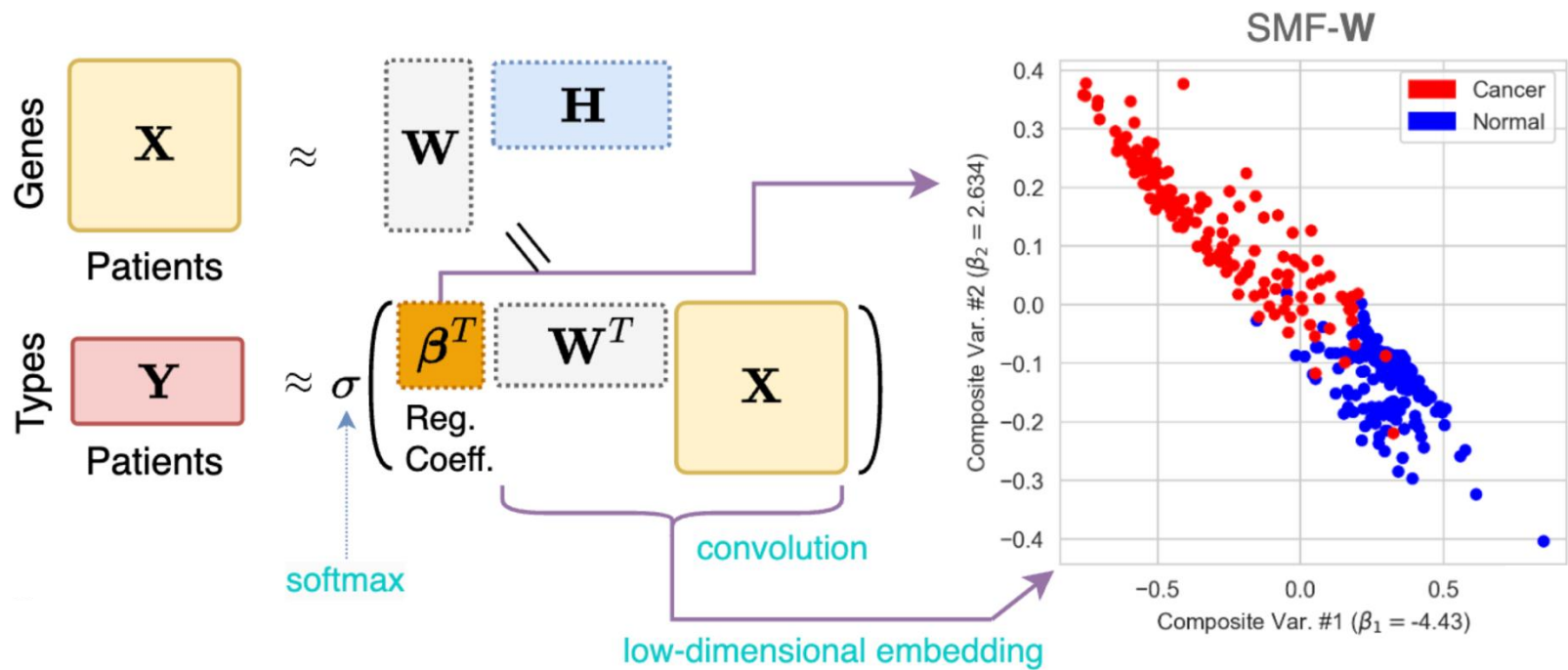


## Supervised Matrix Factorization



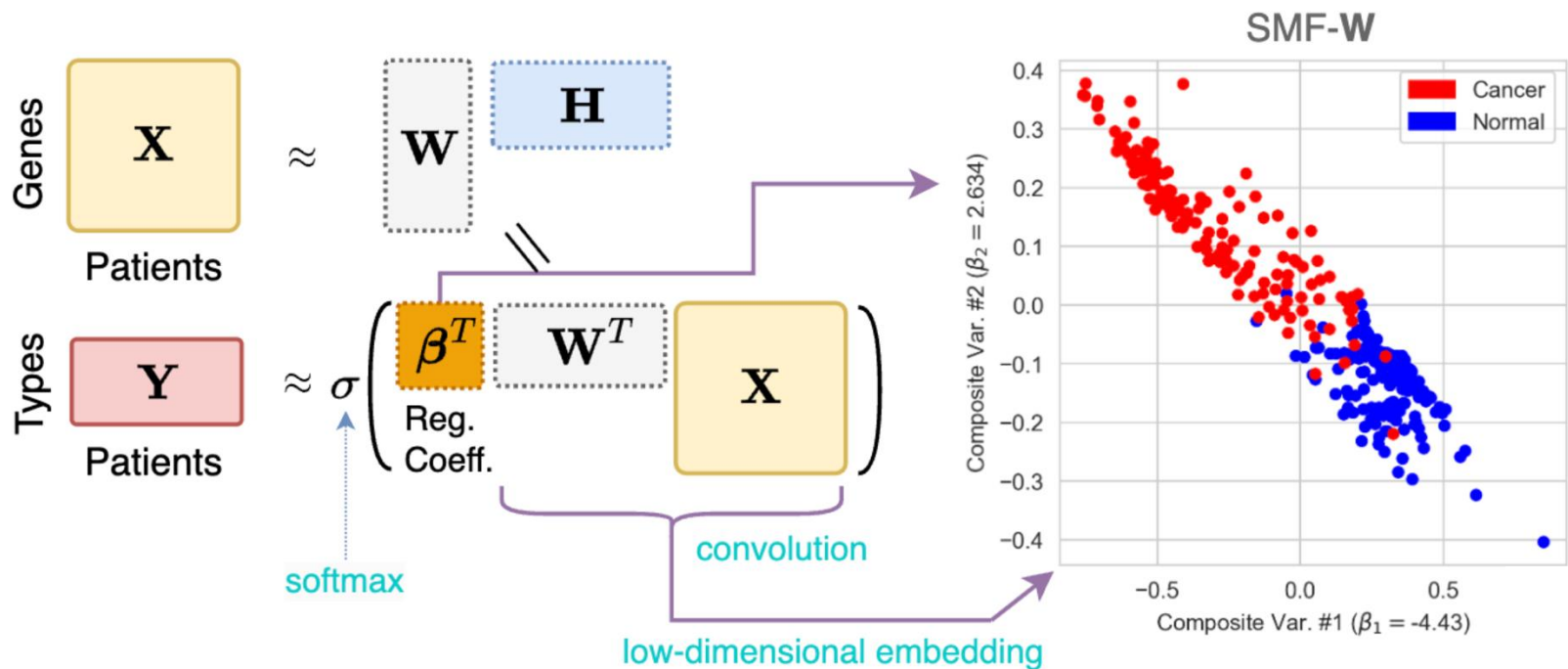


## Supervised Matrix Factorization



## Simultaneous dimension reduction and classification!

## Supervised Matrix Factorization



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

Negative log-likelihood under logistic model

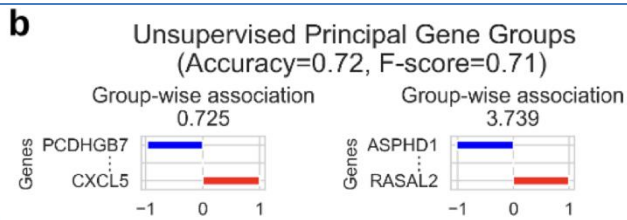
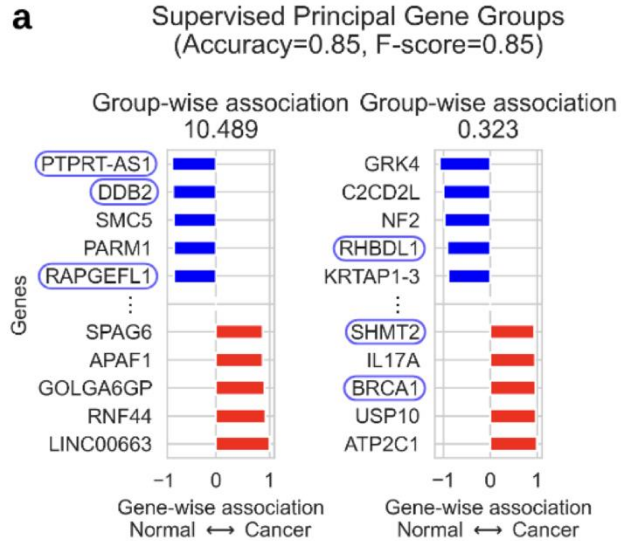
$$\ell(y, a) = \log(1 + \exp(a)) - \mathbf{1}_{\{y=1\}} a$$

Parameters

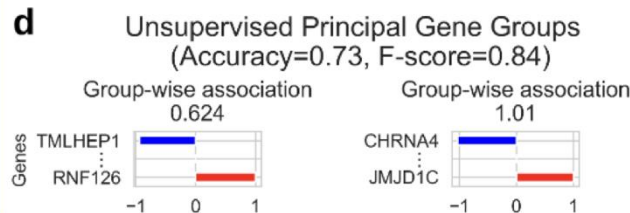
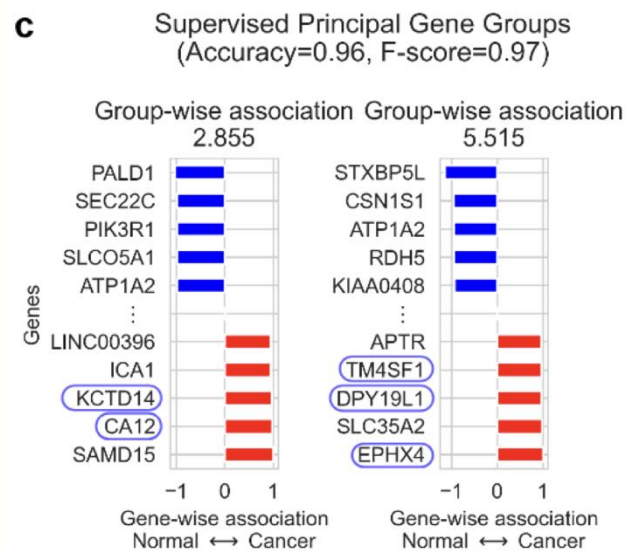
= Three factor matrices  $\mathbf{W}, \mathbf{H}, \beta$

# Supervised Matrix Factorization

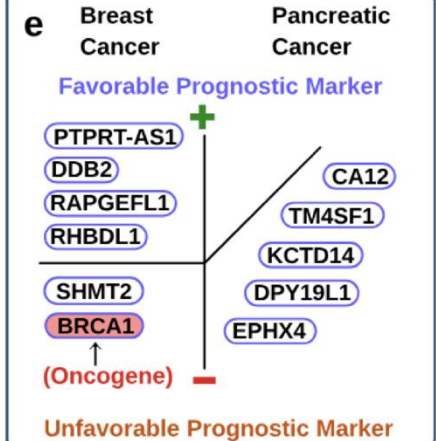
## Breast Cancer Detection



## Pancreatic Cancer Detection



Known cancer-associated genes detected by our method



**f**

Methods	Pancreatic	Breast
SMF-W (Ours)	0.869 (0.02)	<b>0.894 (0.02)</b>
SMF-H (Ours)	<b>0.885 (0.07)</b>	0.875 (0.01)
SMF-W (BCD)	0.785 (0.08)	0.753 (0.03)
SMF-H (BCD)	0.823 (0.06)	0.880 (0.02)
CNN	0.769 (0.07)	0.854 (0.06)
FFNN	0.816 (0.04)	0.890 (0.02)
MF-LR	0.747 (0.13)	0.454 (0.27)
NB	0.815 (0.07)	0.810 (0.02)
SVM	0.746 (0.09)	0.866 (0.02)
RF	0.815 (0.06)	0.844 (0.02)

SMF identifies latent gene groups associated with cancers

It even identifies **known oncogenes** and **prognostic markers**!

- SMF has been around from '08 (Mairal et al. )
- Has been used in lots of applications, but with ad hoc algorithms
- Not much rigorous results have been known before (both computationally and statistically)

- SMF has been around from '08 (Mairal et al. )
- Has been used in lots of applications, but with ad hoc algorithms
- Not much rigorous results have been known before (both computationally and statistically)

## Optimization for SMF

*1. Adaptive Block Coordinate Descent*

*2. Low-rank PGD*

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

- $\alpha_n$ : Stepsizes. How to choose them?

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

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

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

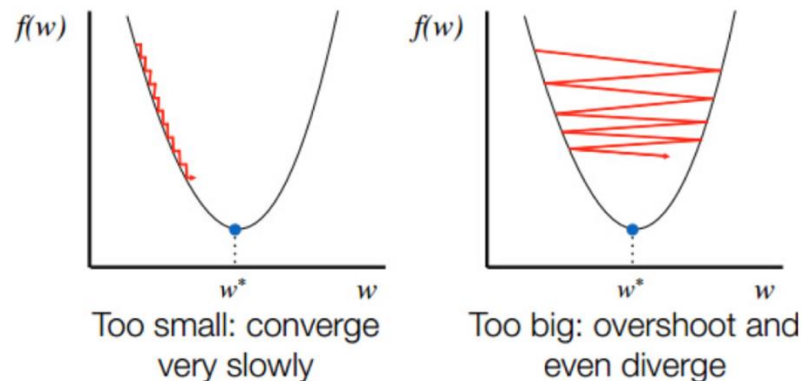
- $\alpha_n$ : Stepsizes. How to choose them?
  - "Small enough stepsize":  $\alpha_n \leq 1/L$ , where
    - $L =$  Lipschitz constant for  $\nabla f$  over  $\Theta$   
 $\approx$  **Largest absolute eigenvalue of  $\nabla^2 f$**  over  $\Theta$
  - But  $\alpha_n = 1/L$  is TOO SMALL!
    - In practice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



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

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

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



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

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

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

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

- Known to be much more robust against stepsize choices than PGD
  - e.g., low-rank matrix factorization, dictionary learning, tensor factorization, kernel learning

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

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

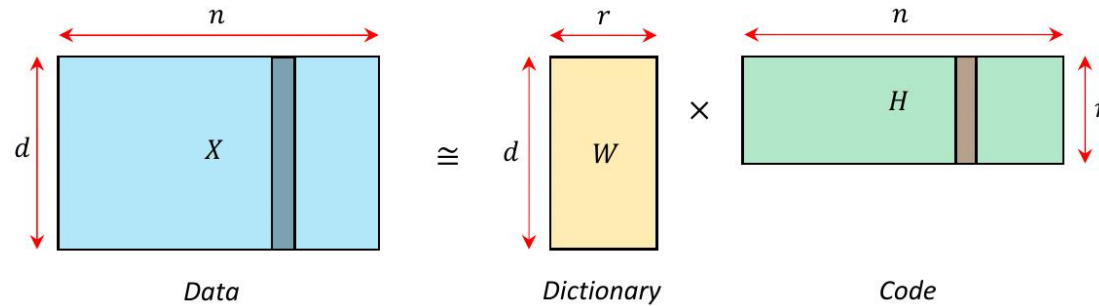
- Known to be much more robust against stepsize choices than PGD
  - e.g., low-rank matrix factorization, dictionary learning, tensor factorization, kernel learning
- [Large  $L \Leftrightarrow \nabla f$  changes wildly]
  - Sensitive dependence on stepsize

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

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

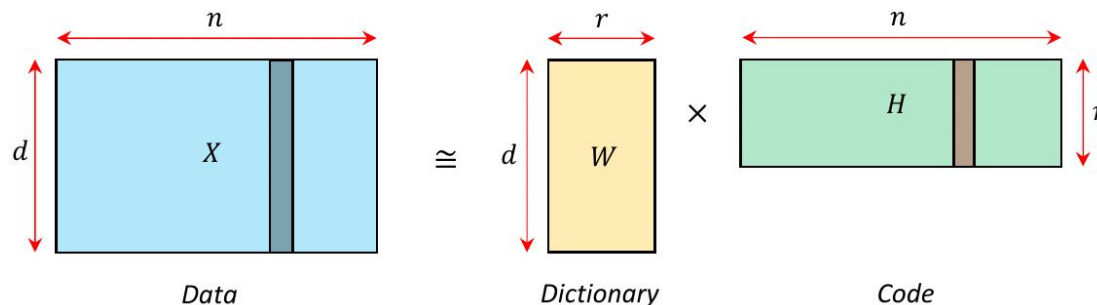
- Known to be much more robust against stepsize choices than PGD
  - e.g., low-rank matrix factorization, dictionary learning, tensor factorization, kernel learning
- [Large  $L \Leftrightarrow \nabla f$  changes wildly]
  - Sensitive dependence on stepsize
- Exploiting block structure → The “effective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

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



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

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



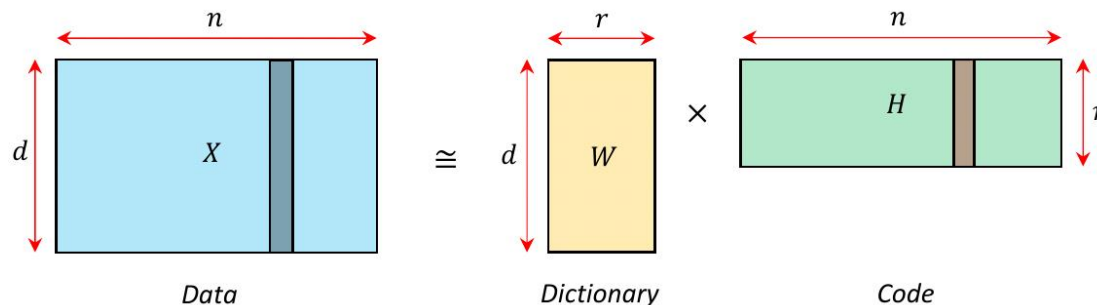
$$\begin{cases} \text{minimize} & f(W, H) = \|X - WH\|_F^2 \\ \text{subject to} & W \in \mathbb{R}_{\geq 0}^{p \times r}, H \in \mathbb{R}_{\geq 0}^{r \times n} \end{cases} \quad \begin{array}{l} \text{(Reconstruction error)} \\ \text{(Constraints)} \end{array}$$

**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)



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

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

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

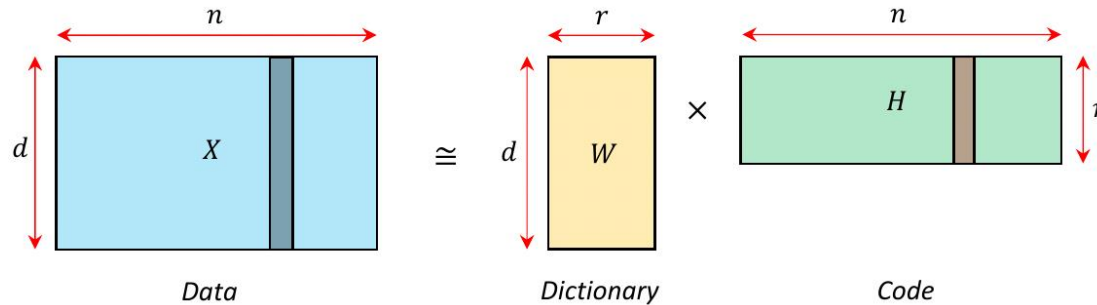
(Almost no one uses this 😊)

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

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



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



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

**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 ☺)

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

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

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

$= \text{Unbounded!}$

**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:**

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

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

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

**Adaptive BCD:**

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

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

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

**Largest Eval of *diagonal blocks* of the Hessian**

**<< Largest Eval of the entire Hessian**

## Supervised Matrix Factorization

---

### 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  $\mathbf{W} \in \mathcal{C}_1$ ,  $\mathbf{H} \in \mathcal{C}_2$ ,  $\boldsymbol{\beta} \in \mathcal{C}_3$ ,  $\boldsymbol{\Gamma} \in \mathcal{C}_4$
  - 5: **For**  $k = 1, 2, \dots, T$  **do:** ( $\triangleright$  For  $\alpha^+$  see 4.3.)
  - 6:   **(Update W)**
  - 7:   Update activation  $\mathbf{a}_1, \dots, \mathbf{a}_n$  and  $\mathbf{K}$
  - 8:    $\nabla_{\mathbf{W}} f(\mathbf{Z}) \leftarrow \mathbf{X}\mathbf{K}^T \boldsymbol{\beta}^T + 2\xi(\mathbf{W}\mathbf{H} - \mathbf{X})\mathbf{H}^T$
  - 9:   Choose  $\eta_{k,1}^{-1} > L_1 := \alpha^+ \|\boldsymbol{\beta}\|_2^2 \cdot \|\mathbf{X}\|_2^2 + 2\xi \|\mathbf{H}\|_2^2$
  - 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  $\boldsymbol{\beta}$ )**
  - 16:   Update activation  $\mathbf{a}_1, \dots, \mathbf{a}_n$  and  $\mathbf{K}$
  - 17:    $\nabla_{\boldsymbol{\beta}} f(\mathbf{Z}) \leftarrow \mathbf{W}^T \mathbf{X}\mathbf{K}^T$
  - 18:   Choose  $\eta_{k,3}^{-1} > L_3 := \alpha^+ \|\mathbf{W}\|_2^2 \cdot \|\mathbf{X}\|_2^2$
  - 19:    $\boldsymbol{\beta} \leftarrow \Pi_{\mathcal{C}_3}(\boldsymbol{\beta} - \eta_{k,3} \nabla_{\boldsymbol{\beta}} f(\mathbf{Z}))$
  - 20:   **(Update  $\boldsymbol{\Gamma}$ )**
  - 21:   Update activation  $\mathbf{a}_1, \dots, \mathbf{a}_n$  and  $\mathbf{K}$
  - 22:    $\nabla_{\boldsymbol{\Gamma}} f(\mathbf{Z}) \leftarrow \mathbf{X}_{\text{aux}} \mathbf{K}^T$
  - 23:   Choose  $\eta_{k,4}^{-1} > L_4 := \alpha^+ \|\mathbf{X}_{\text{aux}}\|_2^2$
  - 24:    $\boldsymbol{\Gamma} \leftarrow \Pi_{\mathcal{C}_4}(\boldsymbol{\Gamma} - \eta_{k,4} \nabla_{\boldsymbol{\Gamma}} f(\mathbf{Z}))$
  - 25: **End for**
  - 26: **Output:**  $\mathbf{Z} = (\mathbf{W}, \mathbf{H}, \boldsymbol{\beta}, \boldsymbol{\Gamma})$
- 

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

## Supervised Matrix Factorization

### 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  $\alpha^+$  see 4.3.)
- 6:   (Update  $\mathbf{W}$ )
- 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  $\mathbf{H}$ )
- 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  $\beta$ )
- 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  $\Gamma$ )
- 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:**  $\mathbf{Z} = (\mathbf{W}, \mathbf{H}, \beta, \Gamma)$

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

### Theorem. (Lee., L., Yao ICML '24 )

Adaptive BCD finds an  $\epsilon$ -stationary pt. of SMF-W within  $O(\epsilon^{-2})$  iterations

Depends only on the block-wise condition number

## Supervised Matrix Factorization

### 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  $\alpha^+$  see 4.3.)
- 6:   (Update  $\mathbf{W}$ )
- 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  $\mathbf{H}$ )
- 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  $\beta$ )
- 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  $\Gamma$ )
- 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:**  $\mathbf{Z} = (\mathbf{W}, \mathbf{H}, \beta, \Gamma)$

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

### Theorem. (Lee., L., Yao ICML '24 )

Adaptive BDC finds an  $\epsilon$ -stationary pt. of SMF-W within  $O(\epsilon^{-2})$  iterations

Depends only on the block-wise condition number  
Hessian analysis + Block Majorization-Minimization guarantees (w/ Yuchen Li ( $\rightarrow$  Job market))



## Supervised Matrix Factorization

### 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  $\alpha^+$  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  $\beta$ )**
- 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  $\Gamma$ )**
- 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:**  $\mathbf{Z} = (\mathbf{W}, \mathbf{H}, \beta, \Gamma)$

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

### Theorem. (Lee., L., Yao ICML '24 )

Adaptive BDC finds an  $\epsilon$ -stationary pt. of SMF-W within  $O(\epsilon^{-2})$  iterations

Depends only on the block-wise condition number  
Hessian analysis + Block Majorization-Minimization guarantees (w/ Yuchen Li)

### Theorem. (Lee., L., Yao ICML '24 )

With  $n$  i.i.d. samples from a generative model,  $\exists$  a local max. of the likelihood function within  $O(n^{-1/2})$  of the true parameter\*

Needs bound on the evals of the *entire* Hessian

\* Some true factors needs L2-reg.



## Exponentially fast optimization algorithm?

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

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

Negative log-likelihood under logistic model  
 $\ell(y, a) = \log(1 + \exp(a)) - \mathbf{1}_{\{y=1\}} a.$

- Highly non-convex, possibly constrained optimization problem
- Adaptive BCD gives standard  $O(\epsilon^{-2})$  polynomial convergence rate
- No hope for any **exponentially fast** algorithm? Yes there is!

## Exponentially fast optimization algorithm?

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

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

Negative log-likelihood under logistic model  
 $\ell(y, a) = \log(1 + \exp(a)) - \mathbf{1}_{\{y=1\}} a$

- Highly non-convex, possibly constrained optimization problem
- Adaptive BCD gives standard  $O(\epsilon^{-2})$  polynomial convergence rate
- No hope for any **exponentially fast** algorithm? Yes there is!

Reformulate it as a low-rank matrix estimation problem

$$\min_{\text{rk}(\mathbf{Z}) \leq r} F \left( \begin{pmatrix} \boxed{\mathbf{W}} & \boxed{\boldsymbol{\beta}} & \boxed{\mathbf{H}} \\ & & \mathbf{Z} \end{pmatrix} \right)$$

"lifted" objective: **strongly convex and smooth**

## Exponentially fast optimization algorithm?

$$\min_{\text{rk}(\mathbf{Z}) \leq r} F \left( \begin{array}{|c|c|c|} \hline \mathbf{W} & \boldsymbol{\beta} & \mathbf{H} \\ \hline & & \mathbf{Z} \\ \hline \end{array} \right)$$

"lifted" objective: **strongly convex and smooth**

$$\text{(Low-rk PGD)} \quad \mathbf{Z}_{k+1} \leftarrow \text{SVD}_r(\mathbf{Z}_k - \eta \nabla F(\mathbf{Z}_k))$$

## Exponentially fast optimization algorithm?

$$\min_{\text{rk}(\mathbf{Z}) \leq r} F \left( \begin{pmatrix} \boxed{\mathbf{W}} & \boxed{\boldsymbol{\beta}} & \boxed{\mathbf{H}} \\ & & \mathbf{Z} \end{pmatrix} \right)$$

"lifted" objective: **strongly convex and smooth**

$$\text{(Low-rk PGD)} \quad \mathbf{Z}_{k+1} \leftarrow \text{SVD}_r(\mathbf{Z}_k - \eta \nabla F(\mathbf{Z}_k))$$

NOT a Riemannian optimization alg.!

F is SC only in the Euclidean geometry!!

## Exponentially fast optimization algorithm?

$$\min_{\text{rk}(Z) \leq r} F \left( \begin{pmatrix} W & \beta & H \\ & & Z \end{pmatrix} \right)$$

"lifted" objective: **strongly convex and smooth**

$$(\text{Low-rk PGD}) \quad Z_{k+1} \leftarrow \text{SVD}_r(Z_k - \eta \nabla F(Z_k))$$

NOT a Riemannian optimization alg.!

F is SC only in the Euclidean geometry!!

**Theorem. (Lee, L., Yao NeurIPS '23 )**

If F is well-conditioned ( $\kappa < 3$ ) and  $\eta \ll 1$ ,  
LPGD converges exponentially fast to the global  
minimizer

## Exponentially fast optimization algorithm?

$$\min_{\text{rk}(Z) \leq r} F \left( \begin{pmatrix} W & \beta & H \\ & & Z \end{pmatrix} \right)$$

"lifted" objective: **strongly convex and smooth**

$$\text{(Low-rk PGD)} \quad Z_{k+1} \leftarrow \text{SVD}_r(Z_k - \eta \nabla F(Z_k))$$

NOT a Riemannian optimization alg.!

F is SC only in the Euclidean geometry!!

### Theorem. (Lee., L., Yao NeurIPS '23 )

If F is well-conditioned ( $\kappa < 3$ ) and  $\eta \ll 1$ ,  
LPGD converges exponentially fast to the global minimizer

F is well-conditioned,  $\text{SVD}_r \approx \text{Subspace Proj.}$

→ 2-Lipschitz, so stronger than  $\frac{1}{2}$  contraction suffices

## Exponentially fast optimization algorithm?

$$\min_{\text{rk}(\mathbf{Z}) \leq r} F \left( \begin{array}{|c|c|c|} \hline \mathbf{W} & \boldsymbol{\beta} & \mathbf{H} \\ \hline & & \mathbf{Z} \\ \hline \end{array} \right)$$

"lifted" objective: **strongly convex and smooth**

**(Low-rk PGD)**  $\mathbf{Z}_{k+1} \leftarrow \text{SVD}_r(\mathbf{Z}_k - \eta \nabla F(\mathbf{Z}_k))$

NOT a Riemannian optimization alg.!

F is SC only in the Euclidean geometry!!

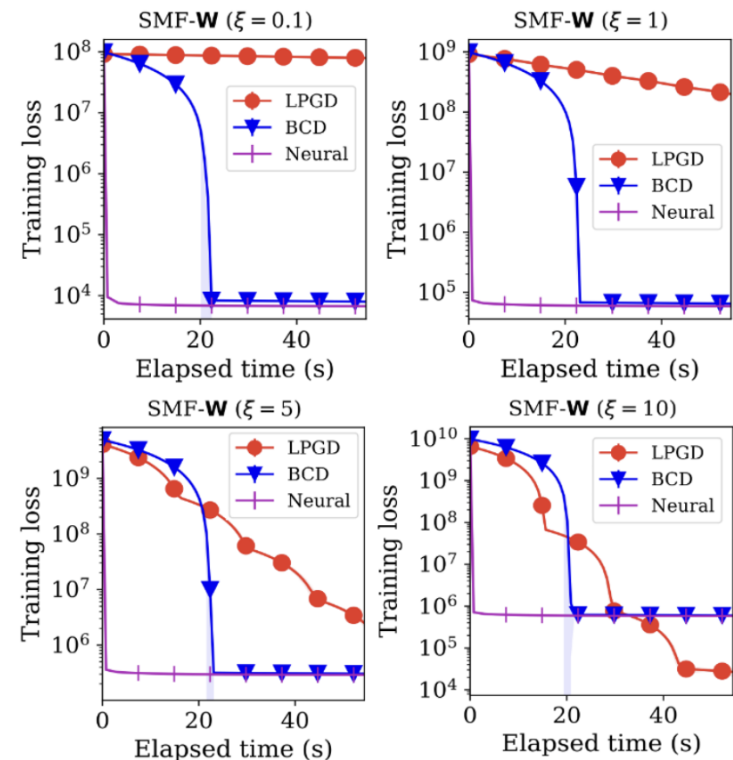
**Theorem. (Lee., L., Yao NeurIPS '23 )**

If F is well-conditioned ( $\kappa < 3$ ) and  $\eta \ll 1$ ,  
LPGD converges exponentially fast to the global minimizer

F is well-conditioned,  $\text{SVD}_r \approx \text{Subspace Proj.}$

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

Better  
Condition #



## Summary

SMF provides an interpretable low-dimensional, label-aware compression

- **Various applications where interpretability matters!**
  - Identifying:
    - cancer-related gene groups
    - survival-related gene groups (German Cancer Research Center)
    - driving patterns for auto insurance customers (AmFam)
    - climate patterns indicating future El Nino/La Nina (N. Chen)
- **Adaptive BCD**
  - Robust nonconvex constrained optimization
  - Depends only on the block-diagonal Hessian
  - Sublinear convergence
- **Low-rank PGD**
  - Lift non-convex problem to convex low-rank opt problem
  - Exponentially fast convergence with good condition number
- **Non-asymptotic Statistical Guarantees**



# Thank you very much!

## References

1. Joowon Lee, Hanbaek Lyu, and Weixin Yao, "*Exponentially Convergent Algorithms for Supervised Matrix Factorization*", NeurIPS 2023
2. Joowon Lee, Hanbaek Lyu, and Weixin Yao, "*Supervised Matrix Factorization: Local Landscape Analysis and Applications*" ICML 2024