

# Scaling the Hierarchical Topic Modeling Mountain

Neural NMF and Iterative Projection Methods

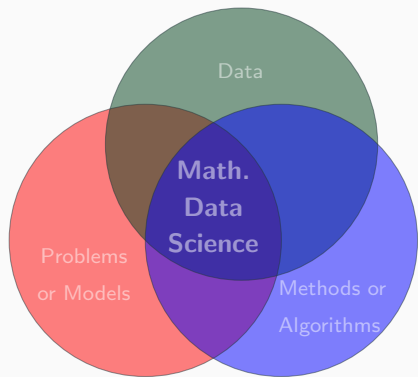
---

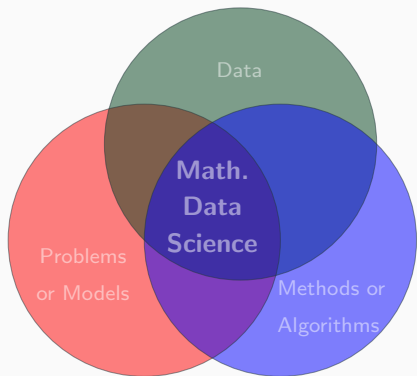
Jamie Haddock

Harvey Mudd College,  
January 28, 2020

Computational and Applied Mathematics  
UCLA

# Research Overview

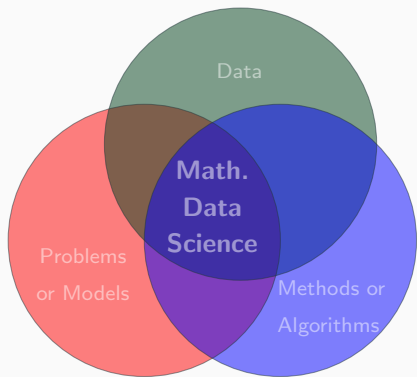




## Mathematical Tools:

- ▷ numerical analysis
- ▷ probability theory
- ▷ convex geometry/analysis
- ▷ combinatorics
- ▷ polyhedral theory
- ▷ ...

# Research Overview

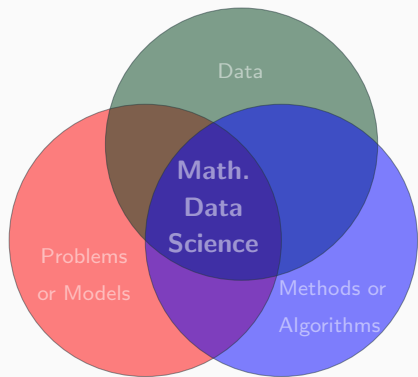


## Data:

- ▷ **MyLymeData surveys**
- ▷ 20newsgroup
- ▷ Netlib linear programs
- ▷ UCI repository
- ▷ computerized tomography
- ▷ NBA data
- ▷ ...



# Research Overview



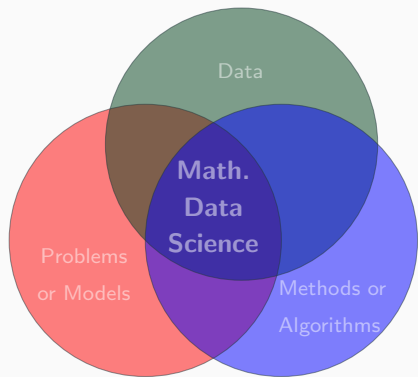
## Problems or Models:

- ▷ **linear least-squares**
- ▷ linear programs
- ▷ **nonnegative matrix factorization**
- ▷ neural networks
- ▷ compressed sensing
- ▷ ...

## Data:

- ▷ **MyLymeData surveys**
- ▷ 20newsgroup
- ▷ Netlib linear programs
- ▷ UCI repository
- ▷ computerized tomography
- ▷ NBA data
- ▷ ...

# Research Overview



## Problems or Models:

- ▷ **linear least-squares**
- ▷ linear programs
- ▷ **nonnegative matrix factorization**
- ▷ neural networks
- ▷ compressed sensing
- ▷ ...

## Data:

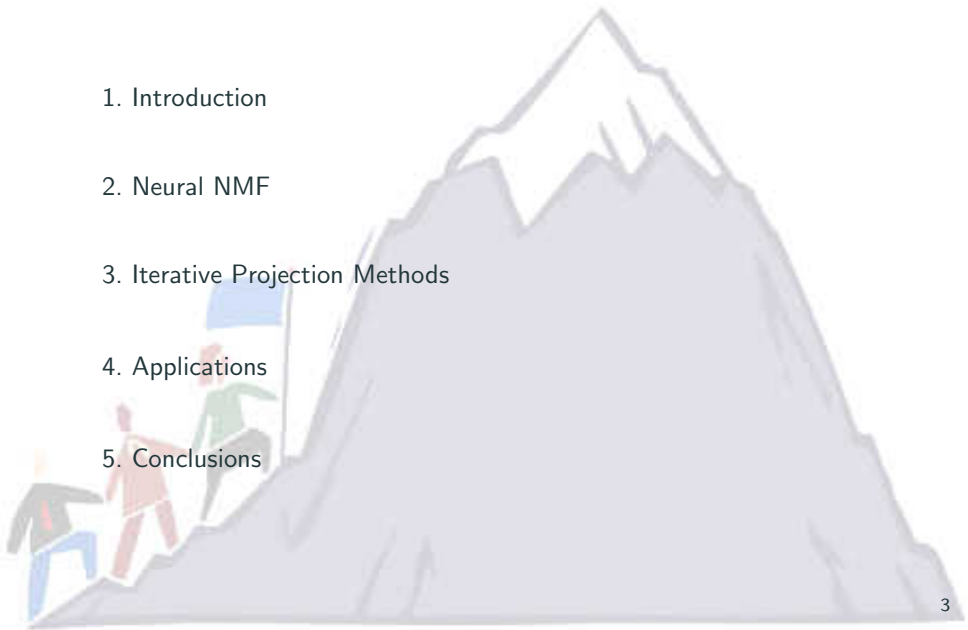
- ▷ **MyLymeData surveys**
- ▷ 20newsgroup
- ▷ Netlib linear programs
- ▷ UCI repository
- ▷ computerized tomography
- ▷ NBA data
- ▷ ...

## Methods or Algorithms:

- ▷ perceptron
- ▷ **iterative projections**
- ▷ Wolfe's method
- ▷ iterative hard thresholding
- ▷ backpropagation
- ▷ ...

# Talk Outline

1. Introduction
2. Neural NMF
3. Iterative Projection Methods
4. Applications
5. Conclusions



# Introduction

---

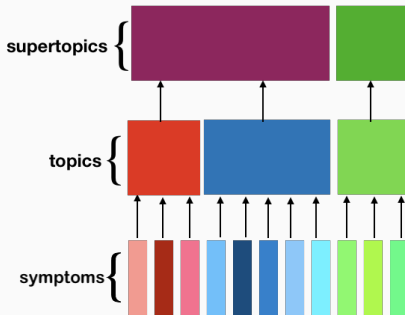
# Motivation

- ▷ MyLymeData: large collection of Lyme disease patient survey data collected by LymeDisease.org (~12,000 patients, 100s of questions)



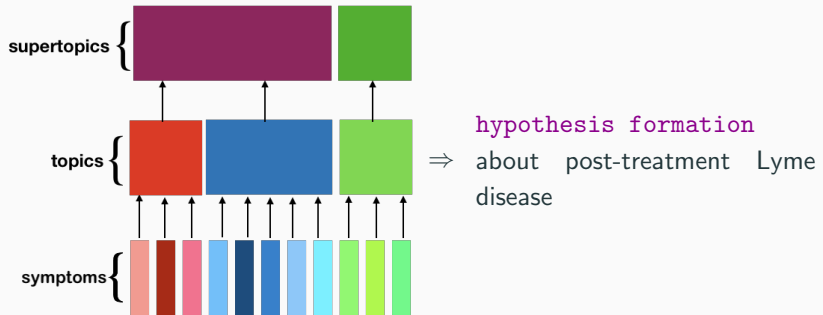
# Motivation

- ▷ MyLymeData: large collection of Lyme disease patient survey data collected by LymeDisease.org (~12,000 patients, 100s of questions)



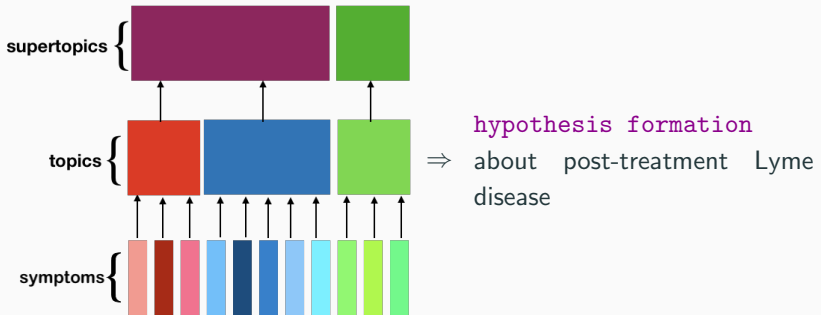
# Motivation

- ▷ MyLymeData: large collection of Lyme disease patient survey data collected by LymeDisease.org (~12,000 patients, 100s of questions)



# Motivation

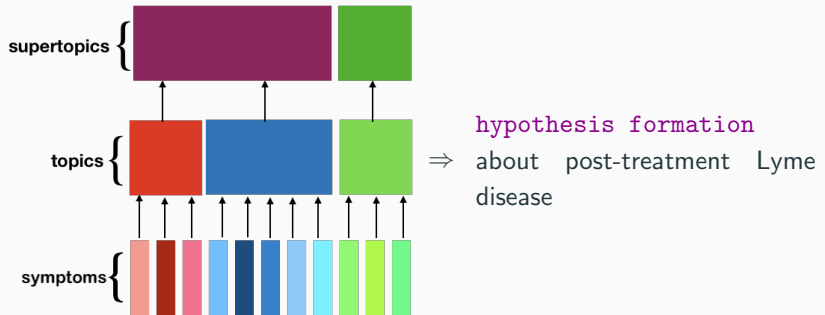
- ▷ MyLymeData: large collection of Lyme disease patient survey data collected by LymeDisease.org (~12,000 patients, 100s of questions)





# Motivation

- ▷ MyLymeData: large collection of Lyme disease patient survey data collected by LymeDisease.org (~12,000 patients, 100s of questions)



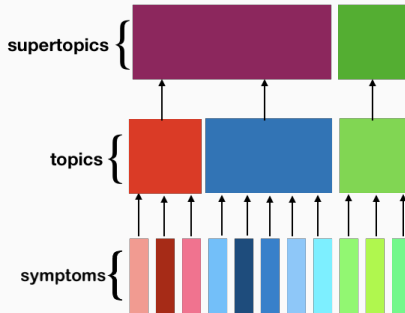
**Main question:** How can we identify the topic hierarchy of MyLymeData symptom questions?



# Motivation



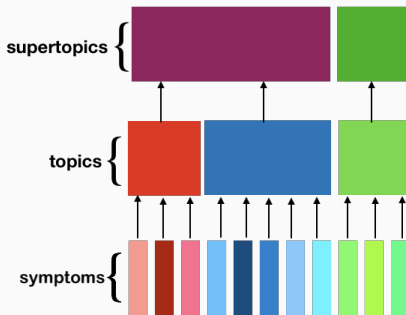
**Main question:** How can we identify the topic hierarchy of MyLymeData symptom questions?



# Motivation



**Main question:** How can we identify the topic hierarchy of MyLymeData symptom questions?



**Answer:** Neural Nonnegative Matrix Factorization

[Gao, H., Molitor, Needell, Sadovnik, Will, Zhang '19]

# Motivation

**Main question:** How can we identify the topic hierarchy of MyLymeData symptom questions?



**Answer:** Neural Nonnegative Matrix Factorization

[Gao, H., Molitor, Needell, Sadovnik, Will, Zhang '19]

**Main question:** How can we identify the topic hierarchy of MyLymeData symptom questions?



**Answer:** Neural Nonnegative Matrix Factorization

[Gao, H., Molitor, Needell, Sadovnik, Will, Zhang '19]

Sampling Kaczmarz-Motzkin Methods

[H., Ma '19], [De Loera, H., Needell '17]

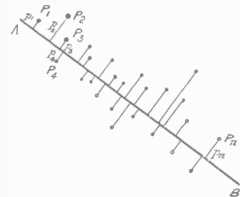


# Topic Modeling

▷ principal component analysis (PCA)

[Pearson 1901]

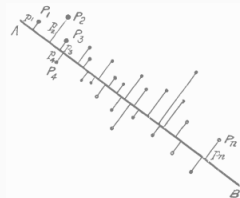
[Hotelling 1933]



Pearson, K. (1901) *On lines and planes of closest fit to systems of points in space.*

# Topic Modeling

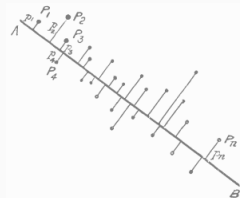
- ▷ principal component analysis (PCA)
  - [Pearson 1901]
  - [Hotelling 1933]
- ▷ latent dirichlet allocation (LDA)
  - [Pritchard, Stephens, Donnelly 2000]
  - [Blei, Ng, Jordan 2003]



Pearson, K. (1901) *On lines and planes of closest fit to systems of points in space.*

# Topic Modeling

- ▷ principal component analysis (PCA)
  - [Pearson 1901]
  - [Hotelling 1933]
- ▷ latent dirichlet allocation (LDA)
  - [Pritchard, Stephens, Donnelly 2000]
  - [Blei, Ng, Jordan 2003]
- ▷ clustering ( $k$ -means, Gaussian mixtures)
  - [Lloyd 1957]
  - [Pearson 1894]

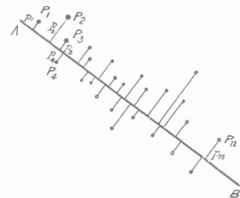


Pearson, K. (1901) *On lines and planes of closest fit to systems of points in space.*

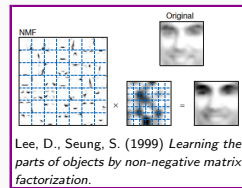


# Topic Modeling

- ▶ principal component analysis (PCA)
  - [Pearson 1901]
  - [Hotelling 1933]
- ▶ latent dirichlet allocation (LDA)
  - [Pritchard, Stephens, Donnelly 2000]
  - [Blei, Ng, Jordan 2003]
- ▶ clustering (*k*-means, Gaussian mixtures)
  - [Lloyd 1957]
  - [Pearson 1894]
- ▶ nonnegative matrix factorization (NMF)
  - [Paatero, Tapper 1994]
  - [Lee, Seung 1999]



Pearson, K. (1901) *On lines and planes of closest fit to systems of points in space.*

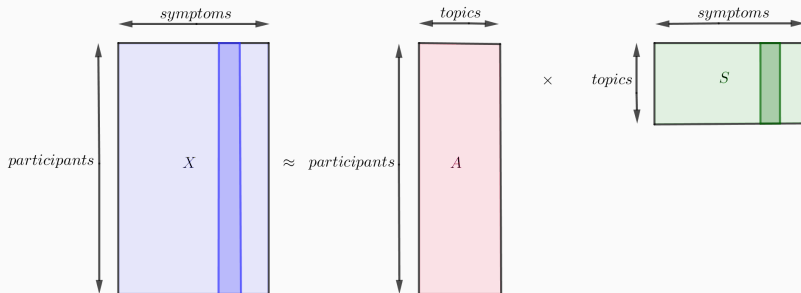


Lee, D., Seung, S. (1999) *Learning the parts of objects by non-negative matrix factorization.*

# Nonnegative Matrix Factorization (NMF)

**Model:** Given nonnegative data  $X$ , compute nonnegative  $A$  and  $S$  of lower rank so that

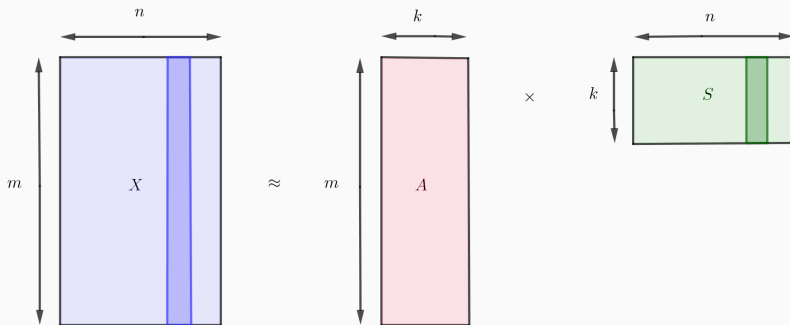
$$X \approx AS.$$



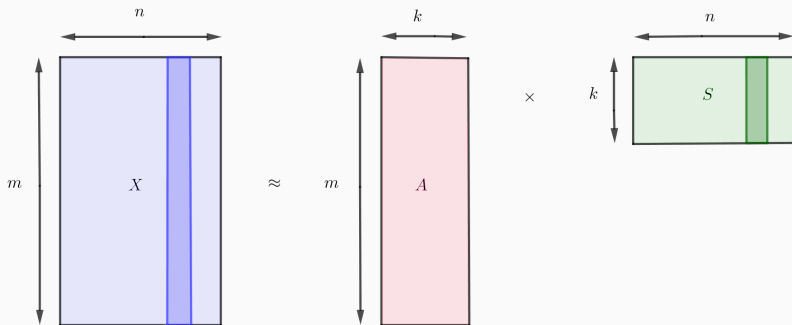
# Nonnegative Matrix Factorization (NMF)

**Model:** Given nonnegative data  $X$ , compute nonnegative  $A$  and  $S$  of lower rank so that

$$X \approx AS.$$



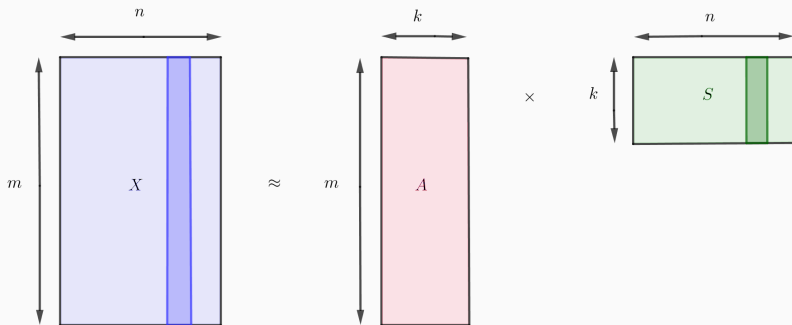
# Nonnegative Matrix Factorization (NMF)



▷ Often formulated as optimization problem

$$\min_{A \in \mathbb{R}_{\geq 0}^{m \times k}, S \in \mathbb{R}_{\geq 0}^{k \times n}} \|X - AS\|_F.$$

# Nonnegative Matrix Factorization (NMF)



- ▷ Often formulated as optimization problem

$$\min_{A \in \mathbb{R}_{\geq 0}^{m \times k}, S \in \mathbb{R}_{\geq 0}^{k \times n}} \|X - AS\|_F.$$

- ▷ Non-convex optimization problem, NP-hard to compute global optimum for fixed  $k$  [Vavasis 2008]

# Hierarchical NMF

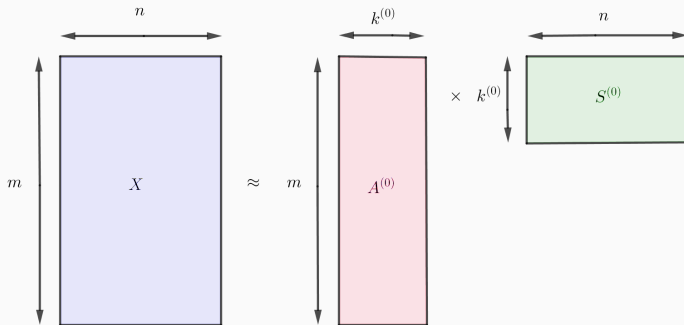
**Model:** Sequentially factorize

$$X \approx A^{(0)}S^{(0)}, S^{(0)} \approx A^{(1)}S^{(1)}, S^{(1)} \approx A^{(2)}S^{(2)}, \dots, S^{(\mathcal{L}-1)} \approx A^{(\mathcal{L})}S^{(\mathcal{L})}.$$

# Hierarchical NMF

**Model:** Sequentially factorize

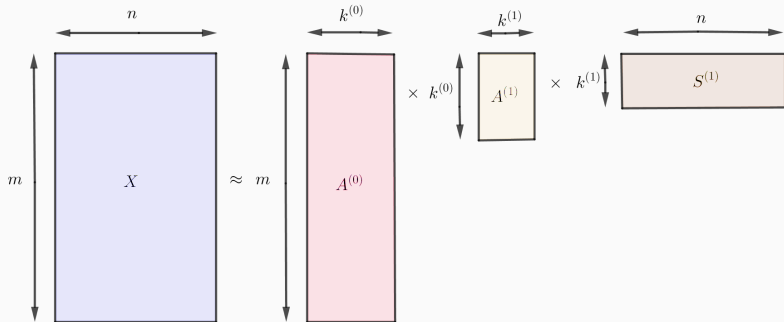
$$X \approx A^{(0)} S^{(0)}, S^{(0)} \approx A^{(1)} S^{(1)}, S^{(1)} \approx A^{(2)} S^{(2)}, \dots, S^{(\mathcal{L}-1)} \approx A^{(\mathcal{L})} S^{(\mathcal{L})}.$$



# Hierarchical NMF

**Model:** Sequentially factorize

$$X \approx A^{(0)} S^{(0)}, S^{(0)} \approx A^{(1)} S^{(1)}, S^{(1)} \approx A^{(2)} S^{(2)}, \dots, S^{(\mathcal{L}-1)} \approx A^{(\mathcal{L})} S^{(\mathcal{L})}.$$

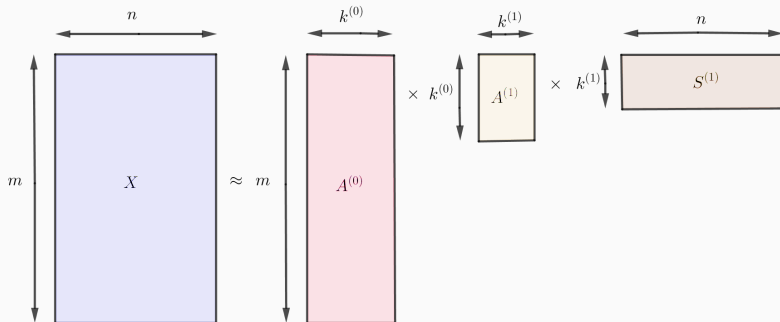




# Hierarchical NMF

**Model:** Sequentially factorize

$$X \approx A^{(0)} S^{(0)}, S^{(0)} \approx A^{(1)} S^{(1)}, S^{(1)} \approx A^{(2)} S^{(2)}, \dots, S^{(\mathcal{L}-1)} \approx A^{(\mathcal{L})} S^{(\mathcal{L})}.$$

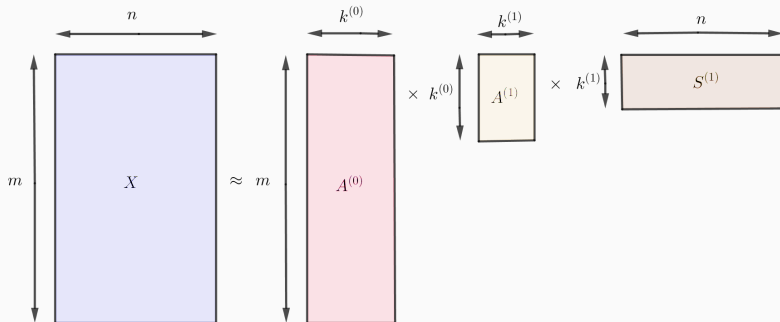


▷  $k^{(\ell)}$ : supertopics collecting  $k^{(\ell-1)}$  subtopics

# Hierarchical NMF

**Model:** Sequentially factorize

$$X \approx A^{(0)} S^{(0)}, S^{(0)} \approx A^{(1)} S^{(1)}, S^{(1)} \approx A^{(2)} S^{(2)}, \dots, S^{(\mathcal{L}-1)} \approx A^{(\mathcal{L})} S^{(\mathcal{L})}.$$



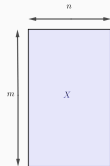
▷  $k^{(\ell)}$ : supertopics collecting  $k^{(\ell-1)}$  subtopics

▷ error propagates through layers

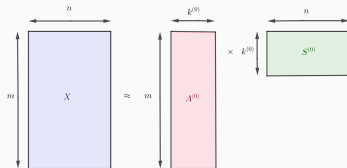
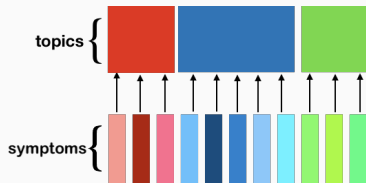
## Neural NMF

---

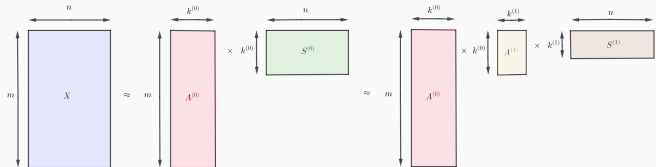
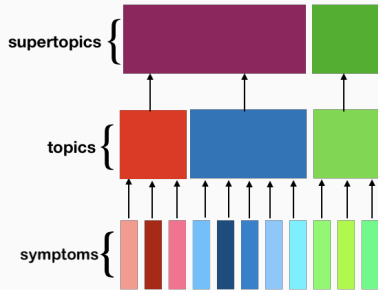
# Hierarchical NMF



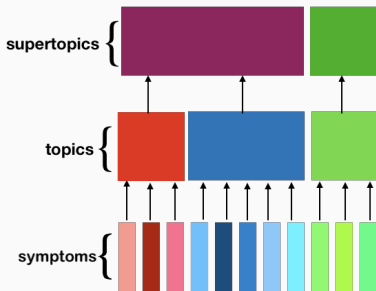
# Hierarchical NMF



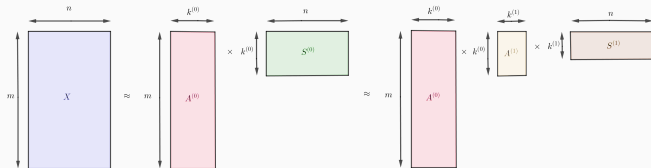
# Hierarchical NMF



# Hierarchical NMF



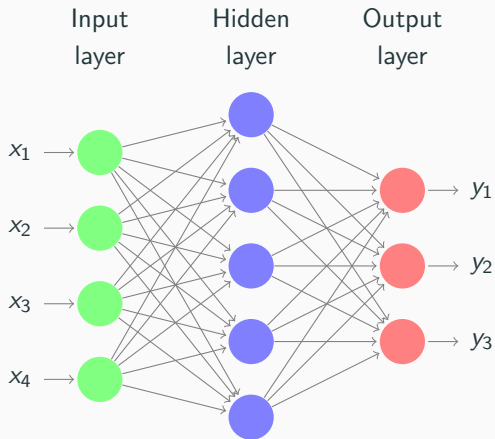
▷ hNMF can be implemented in a **feed-forward neural network** structure



# Feed-forward Neural Networks

**Goal:** Identify weights  $W_1, W_2, \dots, W_L$  to minimize model error

$$\sum_{n=1}^N E(\{W_i\}) = f(\mathbf{y}(\mathbf{x}_n, \{W_i\}), \mathbf{x}_n, \mathbf{t}_n).$$

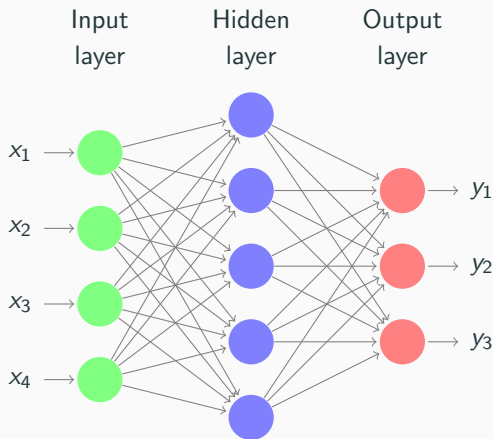




# Feed-forward Neural Networks

**Goal:** Identify weights  $W_1, W_2, \dots, W_L$  to minimize model error

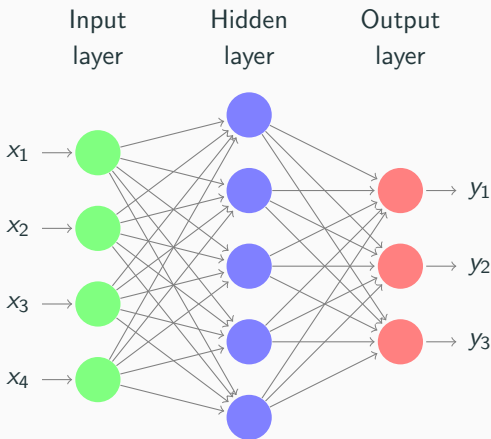
$$E(\{W_i\}) = \sum_{n=1}^N \|\mathbf{y}(\mathbf{x}_n, \{W_i\}) - \mathbf{t}_n\|_2^2.$$



# Feed-forward Neural Networks

**Goal:** Identify weights  $W_1, W_2, \dots, W_L$  to minimize model error

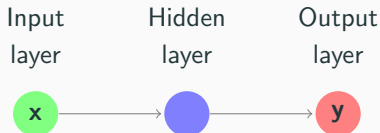
$$E(\{W_i\}) = \sum_{n=1}^N f(\mathbf{y}(\mathbf{x}_n, \{W_i\}), \mathbf{x}_n, \mathbf{t}_n).$$



# Feed-forward Neural Networks

**Goal:** Identify weights  $W_1, W_2, \dots, W_L$  to minimize model error

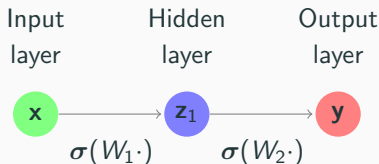
$$E(\{W_i\}) = \sum_{n=1}^N f(\mathbf{y}(\mathbf{x}_n, \{W_i\}), \mathbf{x}_n, \mathbf{t}_n).$$



# Feed-forward Neural Networks

**Goal:** Identify weights  $W_1, W_2, \dots, W_L$  to minimize model error

$$E(\{W_i\}) = \sum_{n=1}^N f(\mathbf{y}(\mathbf{x}_n, \{W_i\}), \mathbf{x}_n, \mathbf{t}_n).$$



## Training:

▷ forward

propagation:

$$\mathbf{z}_1 = \sigma(W_1 \mathbf{x}),$$

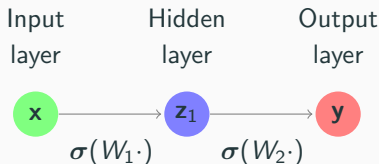
$$\mathbf{z}_2 = \sigma(W_2 \mathbf{z}_1), \dots,$$

$$\mathbf{y} = \sigma(W_L \mathbf{z}_{L-1})$$

# Feed-forward Neural Networks

**Goal:** Identify weights  $W_1, W_2, \dots, W_L$  to minimize model error

$$E(\{W_i\}) = \sum_{n=1}^N f(\mathbf{y}(\mathbf{x}_n, \{W_i\}), \mathbf{x}_n, \mathbf{t}_n).$$

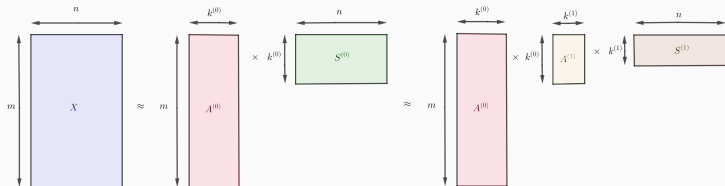


## Training:

- ▷ forward propagation:  
 $\mathbf{z}_1 = \sigma(W_1 \mathbf{x}),$   
 $\mathbf{z}_2 = \sigma(W_2 \mathbf{z}_1), \dots,$   
 $\mathbf{y} = \sigma(W_L \mathbf{z}_{L-1})$
- ▷ back propagation:  
update  $\{W_i\}$  with  $\nabla E(\{W_i\})$

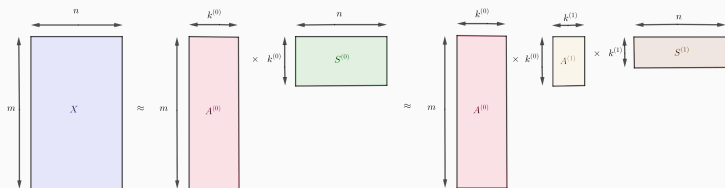
# Our method: Neural NMF

**Goal:** Develop true forward and back propagation algorithms for hNMF.



# Our method: Neural NMF

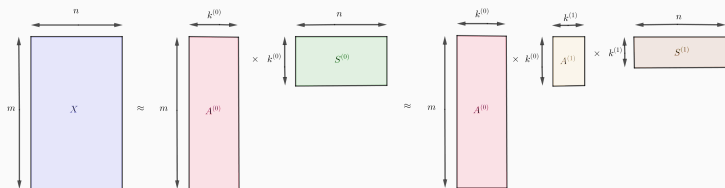
**Goal:** Develop true forward and back propagation algorithms for hNMF.



- ▷ Regard the  $A$  matrices as independent variables, determine the  $S$  matrices from the  $A$  matrices.

# Our method: Neural NMF

**Goal:** Develop true forward and back propagation algorithms for hNMF.

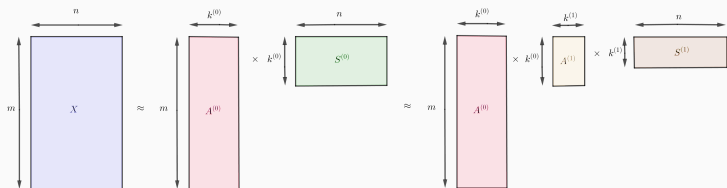


- ▷ Regard the  $A$  matrices as independent variables, determine the  $S$  matrices from the  $A$  matrices.
- ▷ Define  $q(X, A) := \operatorname{argmin}_{S \geq 0} \|X - AS\|_F^2$  (least-squares).



# Our method: Neural NMF

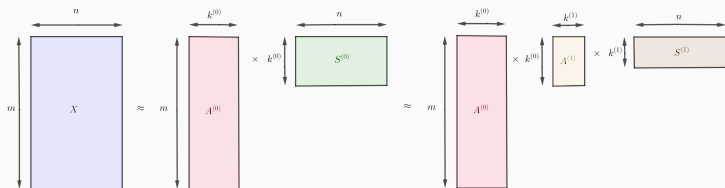
**Goal:** Develop true forward and back propagation algorithms for hNMF.



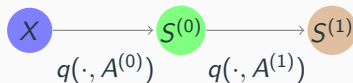
- ▷ Regard the  $A$  matrices as independent variables, determine the  $S$  matrices from the  $A$  matrices.
- ▷ Define  $q(X, A) := \operatorname{argmin}_{S \geq 0} \|X - AS\|_F^2$  (least-squares).
- ▷ Pin the values of  $S$  to those of  $A$  by recursively setting  $S^{(\ell)} := q(S^{(\ell-1)}, A^{(\ell)})$ .

# Our method: Neural NMF

**Goal:** Develop true forward and back propagation algorithms for hNMF.

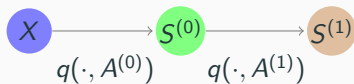
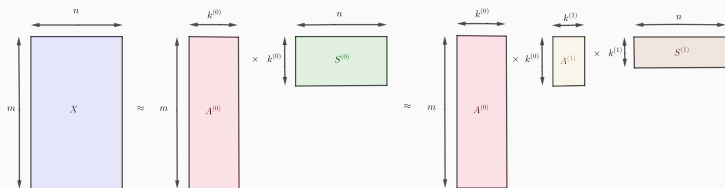


- ▷ Regard the  $A$  matrices as independent variables, determine the  $S$  matrices from the  $A$  matrices.
- ▷ Define  $q(X, A) := \operatorname{argmin}_{S \geq 0} \|X - AS\|_F^2$  (least-squares).
- ▷ Pin the values of  $S$  to those of  $A$  by recursively setting  $S^{(\ell)} := q(S^{(\ell-1)}, A^{(\ell)})$ .



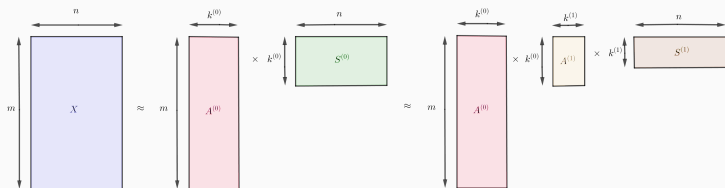
# Our method: Neural NMF

**Goal:** Develop true forward and back propagation algorithms for hNMF.

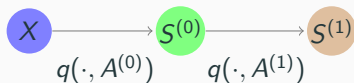


# Our method: Neural NMF

**Goal:** Develop true forward and back propagation algorithms for hNMF.

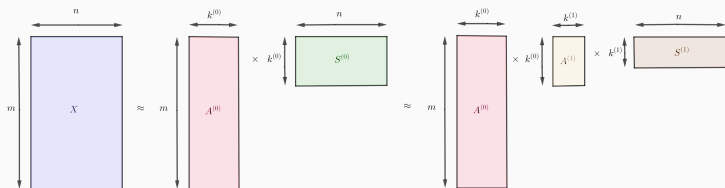


**Training:**



# Our method: Neural NMF

**Goal:** Develop true forward and back propagation algorithms for hNMF.

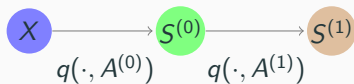


## Training:

▷ forward propagation:

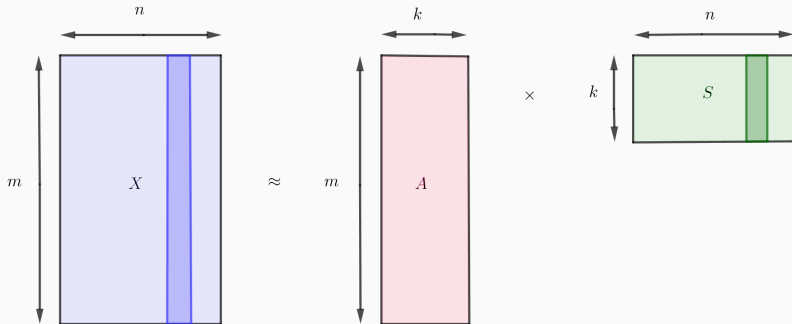
$$\begin{aligned} S^{(0)} &= q(X, A^{(0)}), \\ S^{(1)} &= q(S^{(0)}, A^{(1)}), \dots, \\ S^{(L)} &= q(S^{(L-1)}, A^{(L)}) \end{aligned}$$

▷ back propagation: update  $\{A^{(i)}\}$  with  $\nabla E(\{A^{(i)}\})$



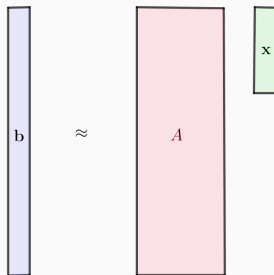
# Least-squares Subroutine

- ▷ least-squares is a fundamental subroutine in forward-propagation



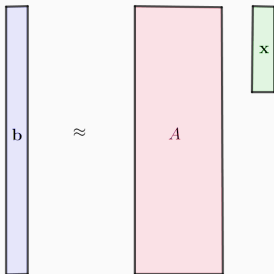
# Least-squares Subroutine

- ▷ least-squares is a fundamental subroutine in forward-propagation



# Least-squares Subroutine

- ▷ least-squares is a fundamental subroutine in forward-propagation



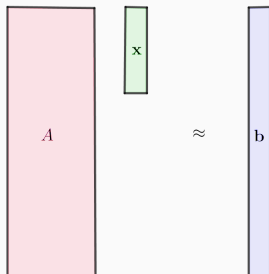
- ▷ **iterative projection methods** can solve these problems



# Iterative Projection Methods

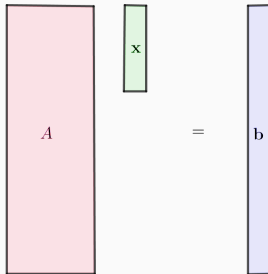
---

# General Setup



# General Setup

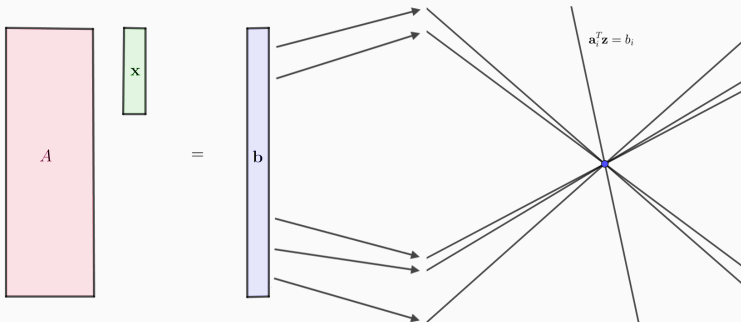
We are interested in solving **highly overdetermined systems of equations**,  $A\mathbf{x} = \mathbf{b}$ , where  $A \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$  and  $m \gg n$ . Rows are denoted  $\mathbf{a}_i^T$ .



A diagram illustrating the matrix equation  $A\mathbf{x} = \mathbf{b}$ . On the left, a tall pink rectangle labeled  $A$  represents the matrix. To its right is a short green rectangle labeled  $\mathbf{x}$ . An equals sign is placed between these two rectangles and a tall light blue rectangle labeled  $\mathbf{b}$  on the right. The vertical dimensions of the rectangles visually represent the dimensions of the vectors and matrix:  $A$  is  $m \times n$ ,  $\mathbf{x}$  is  $n \times 1$ , and  $\mathbf{b}$  is  $m \times 1$ .

# General Setup

We are interested in solving **highly overdetermined systems of equations**,  $Ax = \mathbf{b}$ , where  $A \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$  and  $m \gg n$ . Rows are denoted  $\mathbf{a}_i^T$ .



# Iterative Projection Methods

If  $\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}\}$  is nonempty, these methods construct an **approximation** to a solution:

1. Randomized Kaczmarz Method



Applications:

1. Tomography (Algebraic Reconstruction Technique)

# Iterative Projection Methods

If  $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} = \mathbf{b}\}$  is nonempty, these methods construct an **approximation** to a solution:

1. Randomized Kaczmarz Method
2. Motzkin's Method



Applications:

1. Tomography (Algebraic Reconstruction Technique)
2. Linear programming

# Iterative Projection Methods

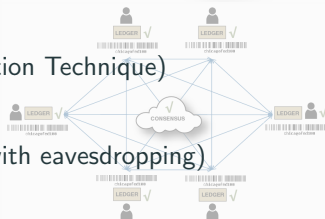
If  $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} = \mathbf{b}\}$  is nonempty, these methods construct an **approximation** to a solution:

1. Randomized Kaczmarz Method
2. Motzkin's Method
3. Sampling Kaczmarz-Motzkin Methods (SKM)



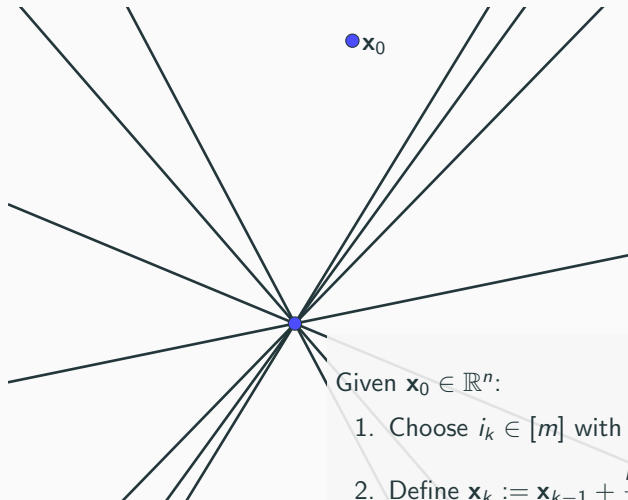
Applications:

1. Tomography (Algebraic Reconstruction Technique)
2. Linear programming
3. Average consensus (greedy gossip with eavesdropping)



When all node operators agree to the change and consensus is reached, the entire network will update their own ledgers. This ensures the immutability of records for network participants and end-users.

# Kaczmarz Method

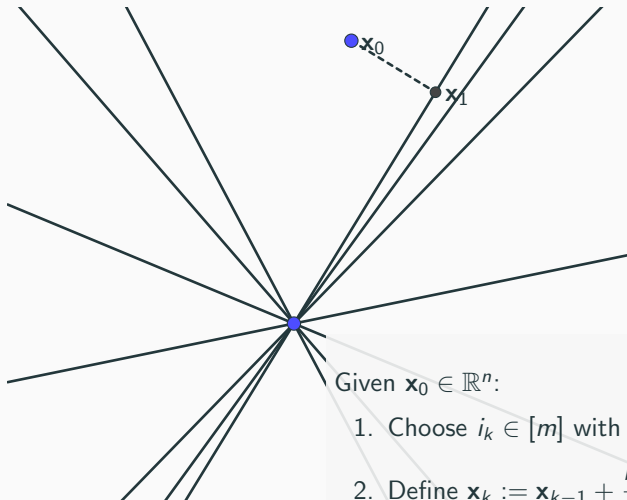


Given  $\mathbf{x}_0 \in \mathbb{R}^n$ :

1. Choose  $i_k \in [m]$  with probability  $\frac{\|\mathbf{a}_{i_k}\|^2}{\|A\|_F^2}$ .
2. Define  $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$ .
3. Repeat.



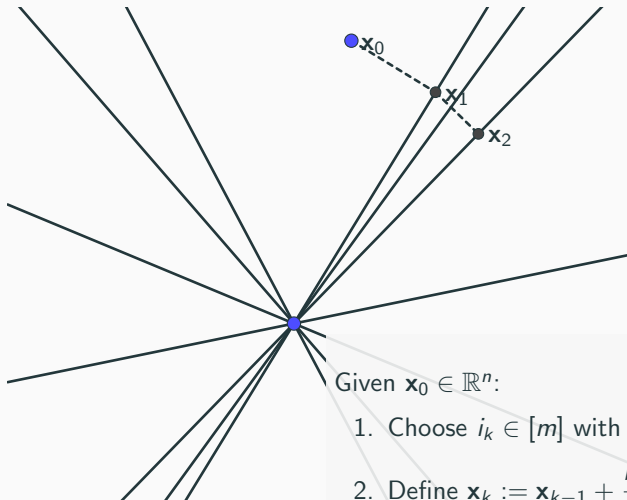
# Kaczmarz Method



Given  $\mathbf{x}_0 \in \mathbb{R}^n$ :

1. Choose  $i_k \in [m]$  with probability  $\frac{\|\mathbf{a}_{i_k}\|^2}{\|A\|_F^2}$ .
2. Define  $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$ .
3. Repeat.

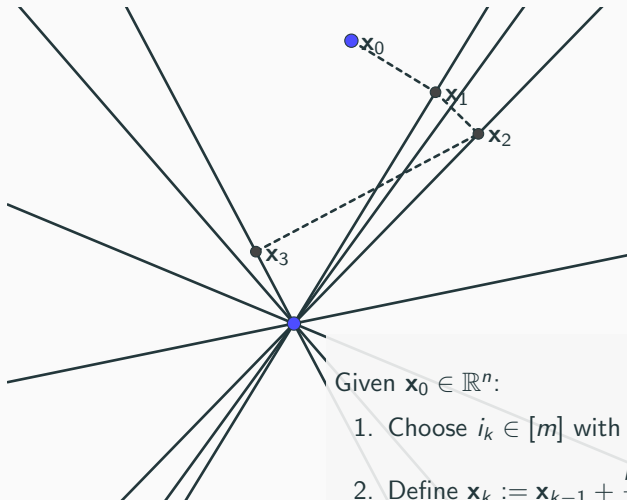
# Kaczmarz Method



Given  $\mathbf{x}_0 \in \mathbb{R}^n$ :

1. Choose  $i_k \in [m]$  with probability  $\frac{\|\mathbf{a}_{i_k}\|^2}{\|A\|_F^2}$ .
2. Define  $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$ .
3. Repeat.

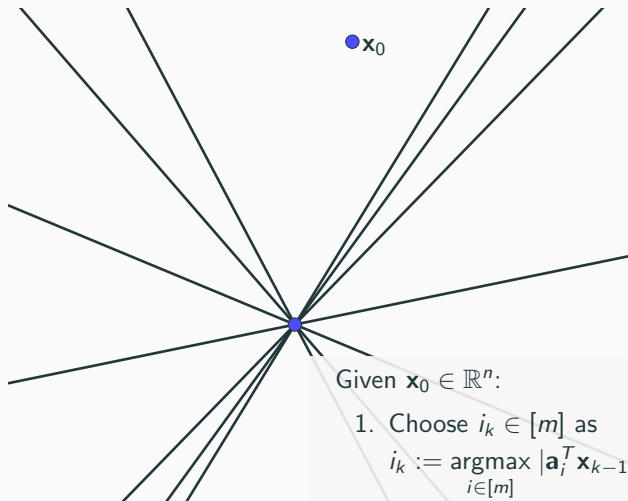
# Kaczmarz Method



Given  $\mathbf{x}_0 \in \mathbb{R}^n$ :

1. Choose  $i_k \in [m]$  with probability  $\frac{\|\mathbf{a}_{i_k}\|^2}{\|A\|_F^2}$ .
2. Define  $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$ .
3. Repeat.

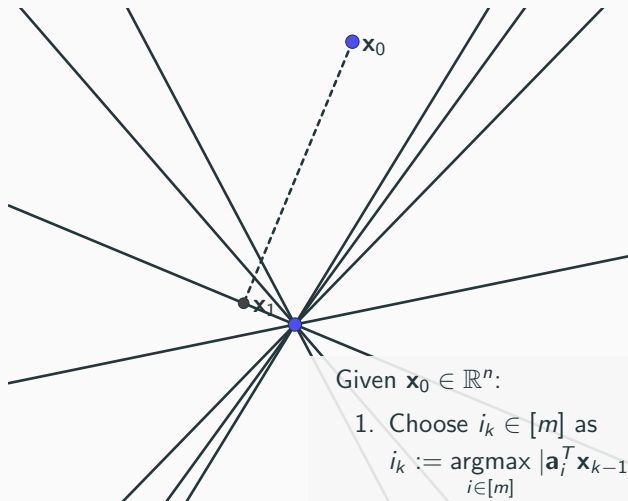
# Motzkin's Method



Given  $\mathbf{x}_0 \in \mathbb{R}^n$ :

1. Choose  $i_k \in [m]$  as
$$i_k := \operatorname{argmax}_{i \in [m]} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|.$$
2. Define  $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}.$
3. Repeat.

# Motzkin's Method



Given  $\mathbf{x}_0 \in \mathbb{R}^n$ :

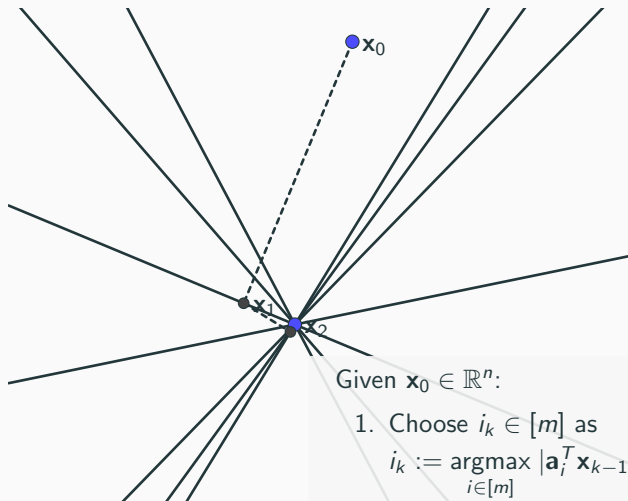
1. Choose  $i_k \in [m]$  as

$$i_k := \operatorname{argmax}_{i \in [m]} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|.$$

2. Define  $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$ .

3. Repeat.

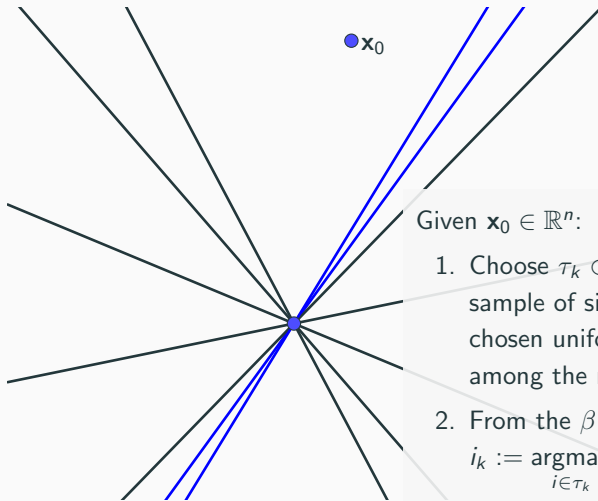
# Motzkin's Method



Given  $\mathbf{x}_0 \in \mathbb{R}^n$ :

1. Choose  $i_k \in [m]$  as
$$i_k := \operatorname{argmax}_{i \in [m]} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|.$$
2. Define  $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}.$
3. Repeat.

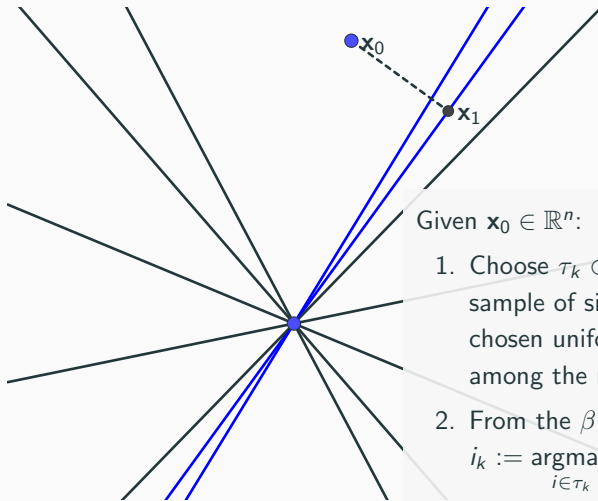
# Our Hybrid Method (SKM)



Given  $\mathbf{x}_0 \in \mathbb{R}^n$ :

1. Choose  $\tau_k \subset [m]$  to be a sample of size  $\beta$  constraints chosen uniformly at random among the rows of  $A$ .
2. From the  $\beta$  rows, choose  $i_k := \operatorname{argmax}_{i \in \tau_k} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|$ .
3. Define 
$$\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}.$$
4. Repeat.

# Our Hybrid Method (SKM)

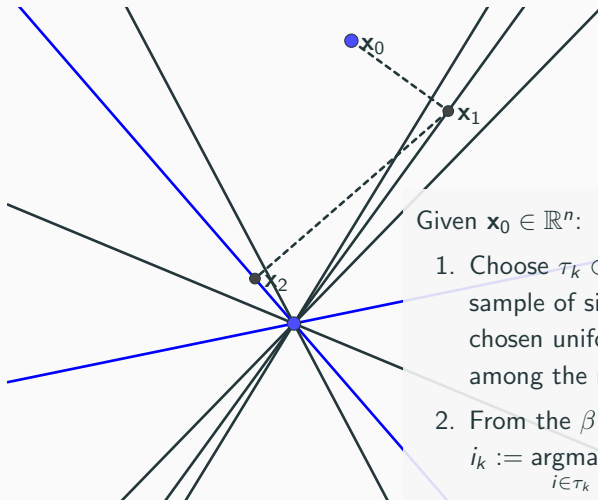


Given  $\mathbf{x}_0 \in \mathbb{R}^n$ :

1. Choose  $\tau_k \subset [m]$  to be a sample of size  $\beta$  constraints chosen uniformly at random among the rows of  $A$ .
2. From the  $\beta$  rows, choose  $i_k := \operatorname{argmax}_{i \in \tau_k} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|$ .
3. Define 
$$\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}.$$
4. Repeat.



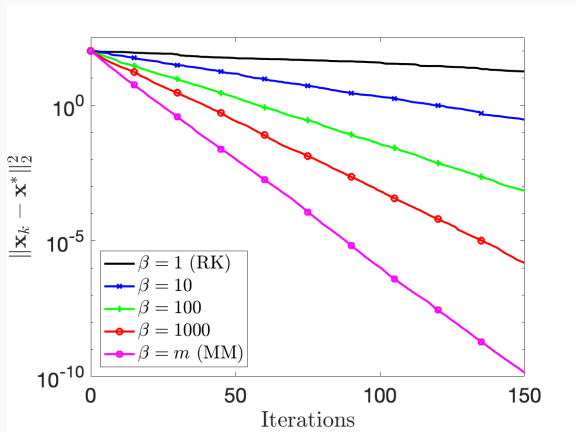
# Our Hybrid Method (SKM)



Given  $\mathbf{x}_0 \in \mathbb{R}^n$ :

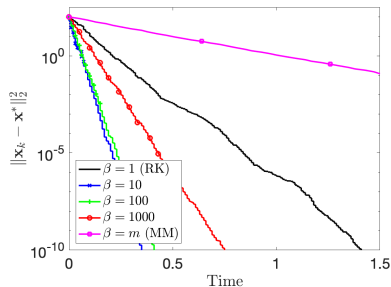
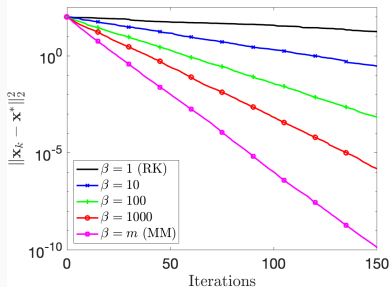
1. Choose  $\tau_k \subset [m]$  to be a sample of size  $\beta$  constraints chosen uniformly at random among the rows of  $A$ .
2. From the  $\beta$  rows, choose  $i_k := \operatorname{argmax}_{i \in \tau_k} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|$ .
3. Define 
$$\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}.$$
4. Repeat.

# Experimental Convergence



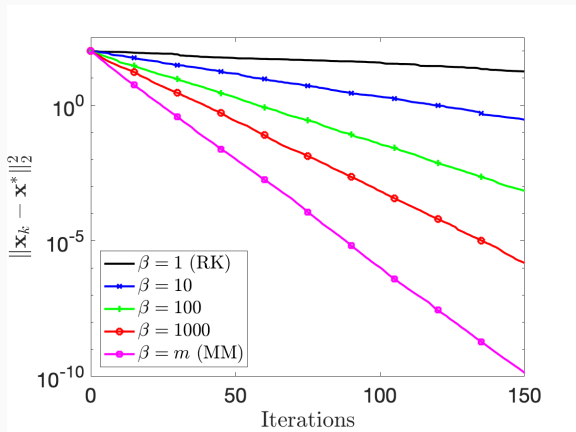
- ▷  $\beta$ : sample size
- ▷  $A$  is  $50000 \times 100$  Gaussian matrix, consistent system
- ▷ 'faster' convergence for larger sample size

# Experimental Convergence



- ▷  $\beta$ : sample size
- ▷  $A$  is  $50000 \times 100$  Gaussian matrix, consistent system
- ▷ 'faster' convergence for larger sample size

# Experimental Convergence



- ▷  $\beta$ : sample size
- ▷  $A$  is  $50000 \times 100$  Gaussian matrix, consistent system
- ▷ 'faster' convergence for larger sample size

# Convergence Rates

Below are the convergence rates for the methods on a system,  $A\mathbf{x} = \mathbf{b}$ , which is consistent with unique solution  $\mathbf{x}$ , whose rows have been normalized to have unit norm.

▷ RK (Strohmer, Vershynin '09):

$$\mathbb{E} \|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

# Convergence Rates

Below are the convergence rates for the methods on a system,  $A\mathbf{x} = \mathbf{b}$ , which is consistent with unique solution  $\mathbf{x}$ , whose rows have been normalized to have unit norm.

▷ RK (Strohmer, Vershynin '09):

$$\mathbb{E} \|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

▷ MM (Agmon '54):

$$\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

# Convergence Rates

Below are the convergence rates for the methods on a system,  $A\mathbf{x} = \mathbf{b}$ , which is consistent with unique solution  $\mathbf{x}$ , whose rows have been normalized to have unit norm.

- ▷ RK (Strohmer, Vershynin '09):

$$\mathbb{E} \|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

- ▷ MM (Agmon '54):

$$\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

- ▷ SKM (DeLoera, H., Needell '17):

$$\mathbb{E} \|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

# Convergence Rates

Below are the convergence rates for the methods on a system,  $A\mathbf{x} = \mathbf{b}$ , which is consistent with unique solution  $\mathbf{x}$ , whose rows have been normalized to have unit norm.

- ▷ RK (Strohmer, Vershynin '09):

$$\mathbb{E} \|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

- ▷ MM (Agmon '54):

$$\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

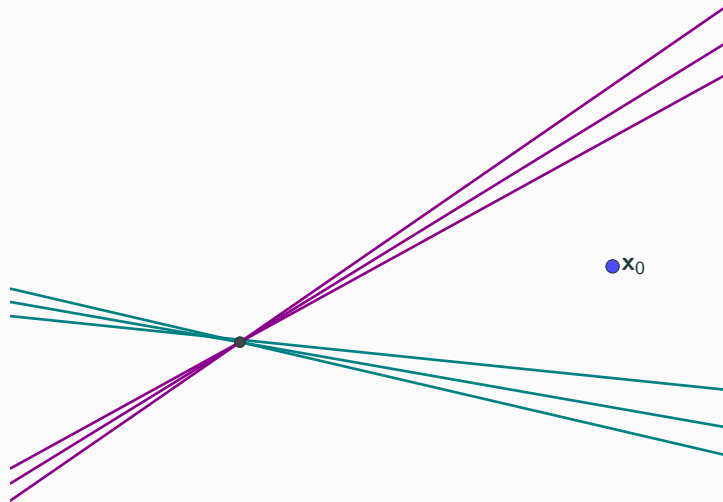
- ▷ SKM (DeLoera, H., Needell '17):

$$\mathbb{E} \|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(A)}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

Why are these all the same?



# A Pathological Example



# Structure of the Residual

Several works have used sparsity of the residual to improve the convergence rate of greedy methods.

[De Loera, [H.](#), Needell '17], [Bai, Wu '18], [Du, Gao '19]

# Structure of the Residual

Several works have used sparsity of the residual to improve the convergence rate of greedy methods.

[De Loera, [H.](#), Needell '17], [Bai, Wu '18], [Du, Gao '19]

However, not much sparsity can be expected in most cases. Instead, we'd like to use dynamic range of the residual to guarantee faster convergence.

$$\gamma_k := \frac{\sum_{\tau \in \binom{[m]}{\beta}} \|A_{\tau} \mathbf{x}_k - \mathbf{b}_{\tau}\|_2^2}{\sum_{\tau \in \binom{[m]}{\beta}} \|A_{\tau} \mathbf{x}_k - \mathbf{b}_{\tau}\|_{\infty}^2}$$

# Accelerated Convergence Rate

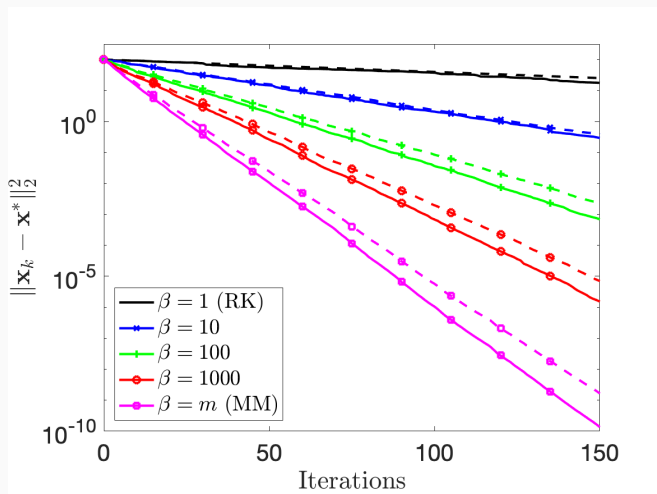
## Theorem (H. - Ma 2019)

Let  $A$  be normalized so  $\|\mathbf{a}_i\|_2 = 1$  for all rows  $i = 1, \dots, m$ . If the system  $A\mathbf{x} = \mathbf{b}$  is consistent with the unique solution  $\mathbf{x}^*$  then the SKM method converges at least linearly in expectation and the rate depends on the dynamic range of the random sample of rows of  $A$ ,  $\tau_j$ . Precisely, in the  $j + 1$ st iteration of SKM, we have

$$\mathbb{E}_{\tau_j} \|\mathbf{x}_{j+1} - \mathbf{x}^*\|_2^2 \leq \left(1 - \frac{\beta \sigma_{\min}^2(A)}{\gamma_j m}\right) \|\mathbf{x}_j - \mathbf{x}^*\|_2^2,$$

$$\text{where } \gamma_j := \frac{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_2^2}{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_\infty^2}.$$

# Accelerated Convergence Rate



- ▷  $A$  is  $50000 \times 100$  Gaussian matrix, consistent system
- ▷ bound uses dynamic range of sample of  $\beta$  rows

# What can we say about $\gamma_j$ ?

$$\text{Recall } \gamma_j := \frac{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_2^2}{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_\infty^2}.$$

$$1 \leq \gamma_j \leq \beta$$

# What can we say about $\gamma_j$ ?

$$\text{Recall } \gamma_j := \frac{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_2^2}{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_\infty^2}.$$

$$1 \leq \gamma_j \leq \beta$$

# What can we say about $\gamma_j$ ?

$$\text{Recall } \gamma_j := \frac{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_2^2}{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_\infty^2}.$$

$$1 \leq \gamma_j \leq \beta$$



# What can we say about $\gamma_j$ ?

$$\text{Recall } \gamma_j := \frac{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_2^2}{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_\infty^2}.$$

$$1 \leq \gamma_j \leq \beta$$

$$\mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \leq \alpha \|\mathbf{x}_{k-1} - \mathbf{x}^*\|_2^2$$

	Previous:
RK	$\alpha = 1 - \frac{\sigma_{\min}^2(A)}{m}$
SKM	$\alpha = 1 - \frac{\sigma_{\min}^2(A)}{m}$
MM	$1 - \frac{\sigma_{\min}^2(A)}{4} \leq \alpha \leq 1 - \frac{\sigma_{\min}^2(A)}{m}$

# What can we say about $\gamma_j$ ?

$$\text{Recall } \gamma_j := \frac{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_2^2}{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_\infty^2}.$$

$$1 \leq \gamma_j \leq \beta$$

$$\mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \leq \alpha \|\mathbf{x}_{k-1} - \mathbf{x}^*\|_2^2$$

	Previous:	Current:
RK	$\alpha = 1 - \frac{\sigma_{\min}^2(A)}{m}$	$\alpha = 1 - \frac{\sigma_{\min}^2(A)}{m}$
SKM	$\alpha = 1 - \frac{\sigma_{\min}^2(A)}{m}$	$1 - \frac{\beta \sigma_{\min}^2(A)}{m} \leq \alpha \leq 1 - \frac{\sigma_{\min}^2(A)}{m}$
MM	$1 - \frac{\sigma_{\min}^2(A)}{4} \leq \alpha \leq 1 - \frac{\sigma_{\min}^2(A)}{m}$	$1 - \sigma_{\min}^2(A) \leq \alpha \leq 1 - \frac{\sigma_{\min}^2(A)}{m}$

## What can we say about $\gamma_j$ ?

$$\text{Recall } \gamma_j := \frac{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_2^2}{\sum_{\tau \in \binom{[m]}{\beta}} \|A_\tau \mathbf{x}_j - \mathbf{b}_\tau\|_\infty^2}.$$

$$1 \leq \gamma_j \leq \beta$$

▷ nontrivial bounds on  $\gamma_k$  for Gaussian and average consensus systems

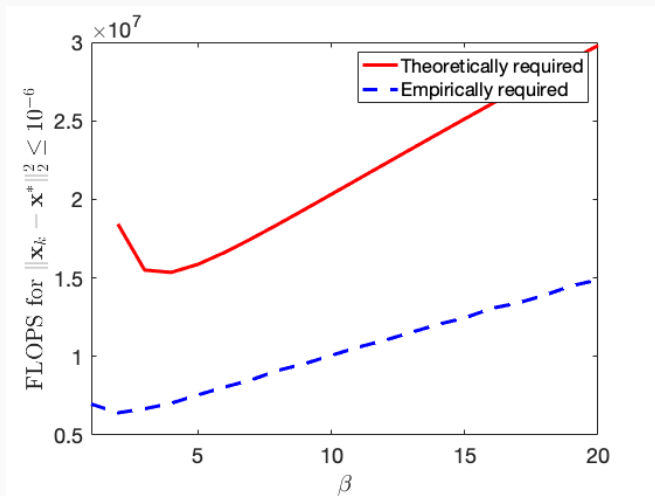
Now can we determine the optimal  $\beta$ ?

## Now can we determine the optimal $\beta$ ?

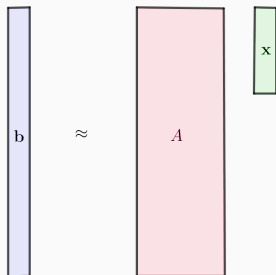
Roughly, if we know the value of  $\gamma_j$ , we can (just) do it.

## Now can we determine the optimal $\beta$ ?

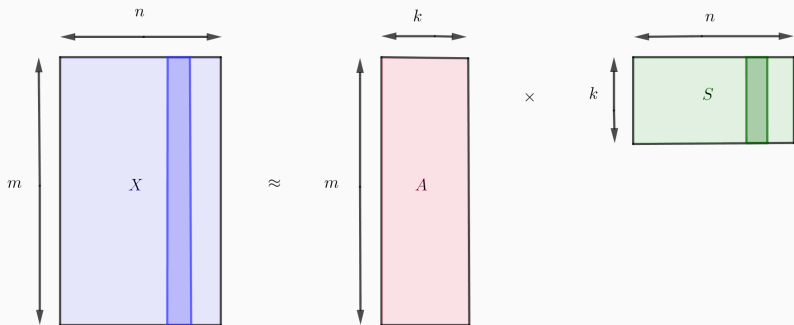
Roughly, if we know the value of  $\gamma_j$ , we can (just) do it.



## Back to Hierarchical NMF

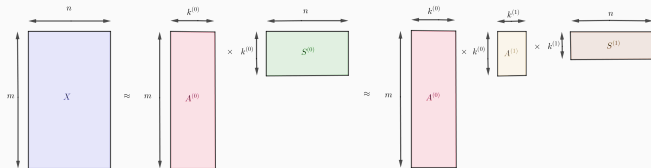
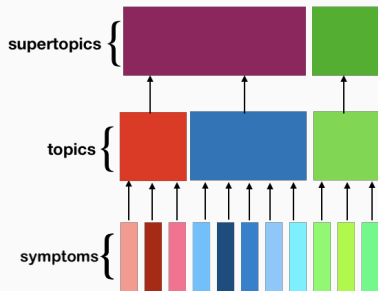


# Back to Hierarchical NMF

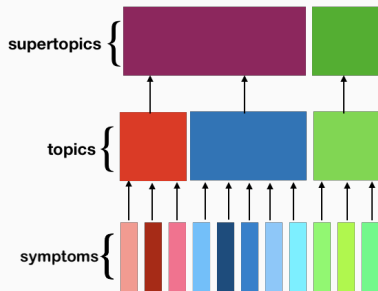




# Back to Hierarchical NMF

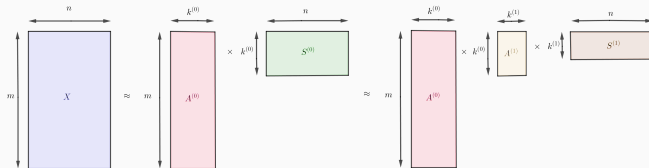


# Back to Hierarchical NMF

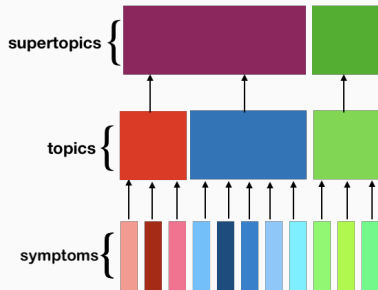


Compare:

▷ hNMF (sequential NMF)

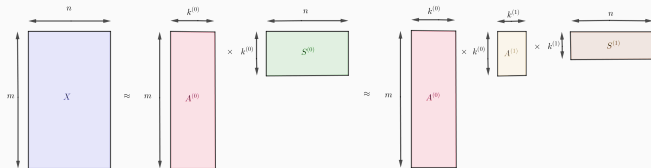


# Back to Hierarchical NMF

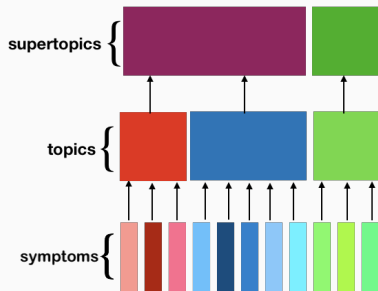


Compare:

- ▷ hNMF (sequential NMF)
- ▷ Deep NMF [Flenner, Hunter '18]

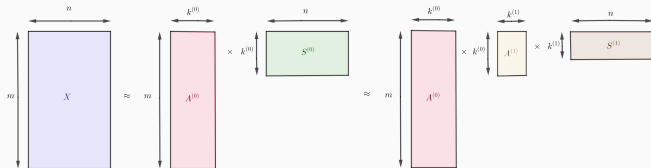


# Back to Hierarchical NMF



## Compare:

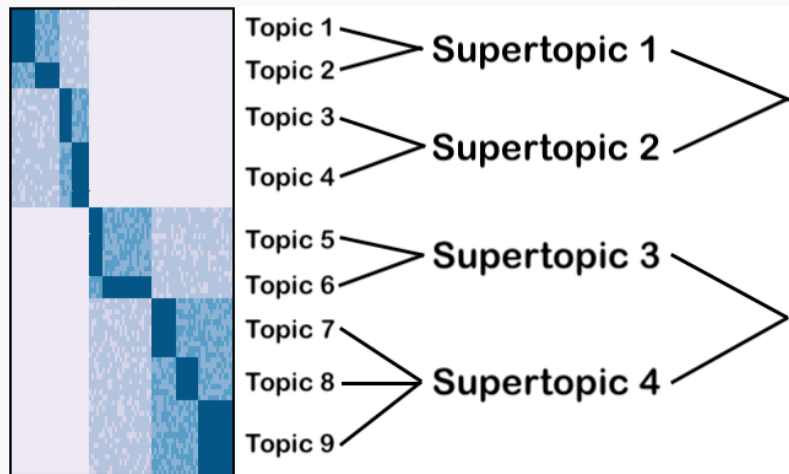
- ▷ hNMF (sequential NMF)
- ▷ Deep NMF [Flenner, Hunter '18]
- ▷ **Neural NMF**



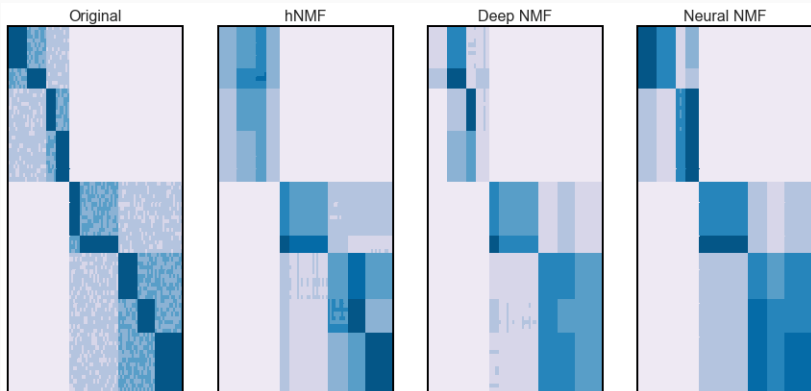
# Applications

---

## Experimental results: synthetic data

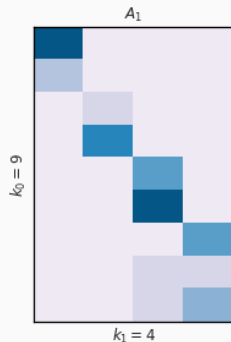
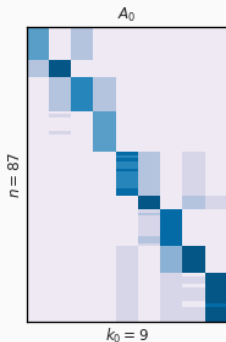
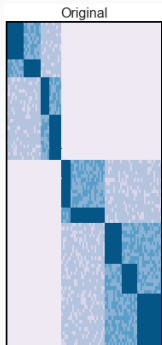


# Experimental results: synthetic data



- ▷ unsupervised reconstruction with two-layer structure  
( $k^{(0)} = 9, k^{(1)} = 4$ )

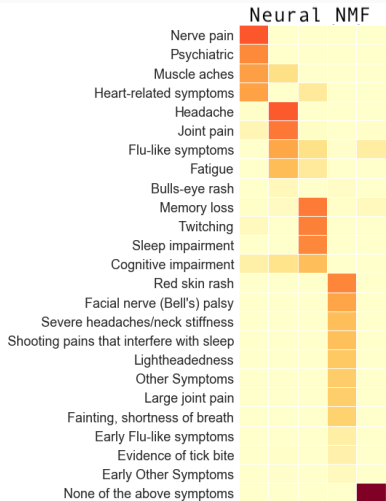
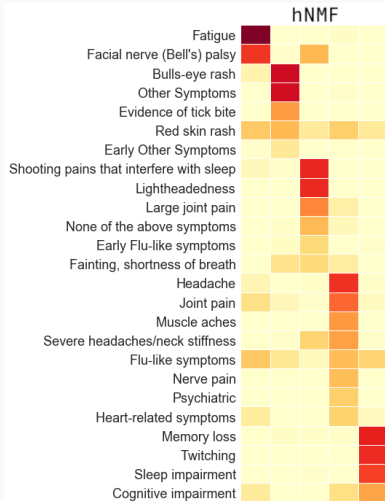
# Experimental results: synthetic data



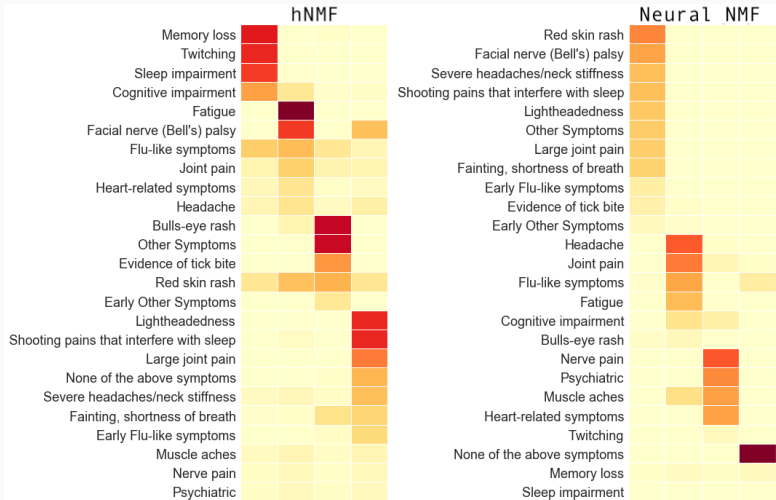
- ▷ unsupervised reconstruction with two-layer structure  
( $k^{(0)} = 9, k^{(1)} = 4$ )



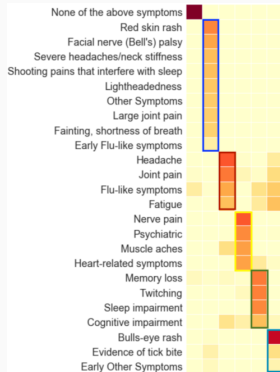
# Experimental results: MyLymeData



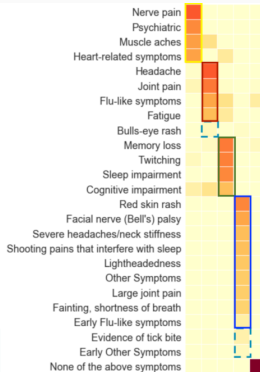
# Experimental results: MyLymeData



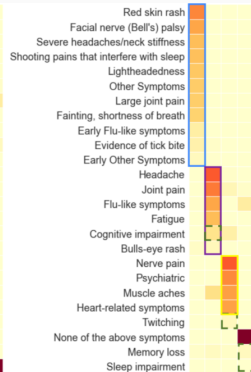
# Experimental results: MyLymeData



$$k^{(0)} = 6$$

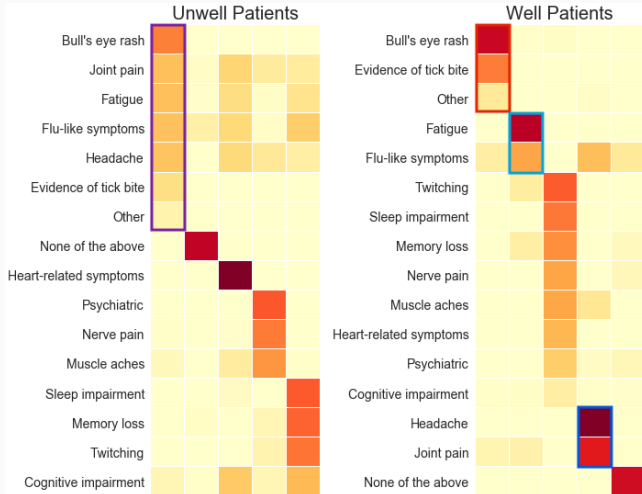


$$k^{(1)} = 5$$



$$k^{(2)} = 4$$

# Experimental results: MyLymeData





- ▷ bulls-eye rash (diagnosing symptoms) topic does not seem to persist for smaller number of topics



- ▷ bulls-eye rash (diagnosing symptoms) topic does not seem to persist for smaller number of topics
- ▷ unwell and well patients have very different presentation of bulls-eye rash symptom in topics



- ▷ bulls-eye rash (diagnosing symptoms) topic does not seem to persist for smaller number of topics
- ▷ unwell and well patients have very different presentation of bulls-eye rash symptom in topics
- ▷ patients unwell because lacking bulls-eye rash for diagnosis or indicative of different disease pathway?

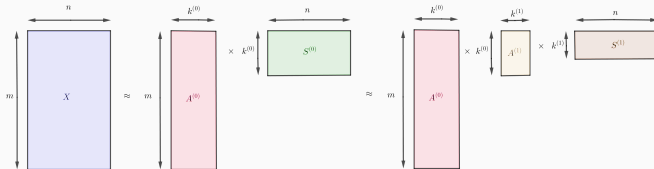
# Conclusions

---

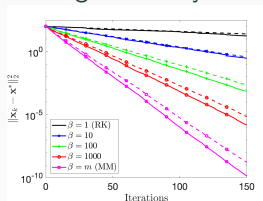


# Conclusions

- ▷ hNMF model can be implemented as a feed-forward neural network



- ▷ presented our method **Neural NMF**
- ▷ described family of algorithms which can solve fundamental least-squares subroutine
- ▷ presented accelerated convergence analysis for **SKM**



- ▷ applied Neural NMF to synthetic data and MyLymeData

## Nonnegative Tensor Decomposition (NTD):

- ▷ for dynamic topic modeling (stemming from WiSDM 2019)
- ▷ hierarchical NTD (joint with Needell, Vendrow\*)
- ▷ robustness of nonnegative CANDECOMP/PARAFAC decomposition (joint with Kassab•)
- ▷ Applications: NBA data (joint with Liu\*), temporal political data



## Iterative Projection Methods:

- ▷ dynamic SKM methods (joint with Ma)
- ▷ corruption robust methods (joint with Needell, Rebrova, Swartworth•)
- ▷ AutoML hyperparameter selection (joint with Heiner\*)
- ▷ Applications: linear network dynamics problems



\* denotes undergraduate collaborator, • denotes graduate collaborator

## **Combinatorial Methods:**

- ▷ Wolfe's method (joint with De Loera, Rademacher)
- ▷ Hansen-Lawson method
- ▷ Applications: metagenomic binning

## **Asynchronous Compressed Sensing:**

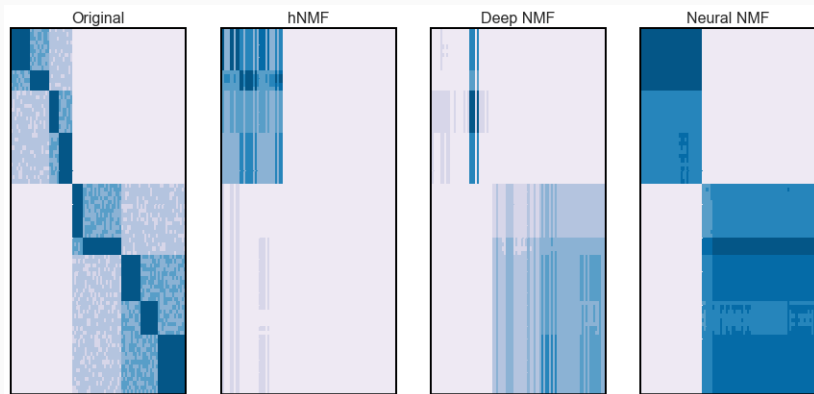
- ▷ Bayesian asynchronous methods (joint with Needell, Rahn timer, Zaeemzadeh)
- ▷ convergence analysis of IHT variants
- ▷ Sparse RK

# Thanks for listening!

## Questions?

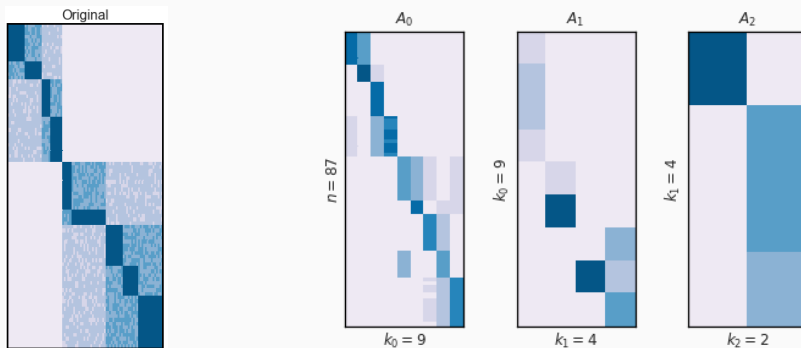
- [1] S. Agmon. **The relaxation method for linear inequalities.** Canadian J. Math., 6:382–392, 1954.
- [2] Z. Bai and W. Wu. **On greedy randomized Kaczmarz method for solving large sparse linear systems.** SIAM J. Sci. Comput., 40(1):A592–A606, 2018.
- [3] A. Cichocki and R. Zdunek. **Multilayer nonnegative matrix factorisation.** Electron. Lett., 42(16):947, 2006.
- [4] J. A. De Loera, J. Haddock, and D. Needell. **A sampling Kaczmarz-Motzkin algorithm for linear feasibility.** SIAM J. Sci. Comput., 39(5):S66–S87, 2017.
- [5] K. Du and H. Gao. **A new theoretical estimate for the convergence rate of the maximal weighted residual Kaczmarz algorithm.** Numer. Math. - Theory Me., 12(2):627–639, 2019.
- [6] M. Gao, J. Haddock, D. Molitor, D. Needell, E. Sadovnik, T. Will, and R. Zhang. **Neural nonnegative matrix factorization for hierarchical multilayer topic modeling.** In Proc. International Workshop on Computational Advances in Multi-Sensor Adaptive Processing, 2019.
- [7] J. Haddock and A. Ma. **Greedy works: An improved analysis of sampling Kaczmarz-Motzkin.** 2019. Submitted.
- [8] J. Haddock and D. Needell. **On Motzkins method for inconsistent linear systems.** BIT, 59(2):387–401, 2019.
- [9] S. Kaczmarz. **Angenäherte auflösung von systemen linearer gleichungen.** Bull. Int. Acad. Polon. Sci. Lett. Ser. A, pages 335–357, 1937.
- [10] D. D. Lee and H. S. Seung. **Learning the parts of objects by non-negative matrix factorization.** Nature, 401:788–791, 1999.
- [11] T. S. Motzkin and I. J. Schoenberg. **The relaxation method for linear inequalities.** Canadian J. Math., 6:393–404, 1954.
- [12] P. Paatero and U. Tapper. **Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values.** Environmetrics, 5(2):111–126, 1994.
- [13] T. Strohmer and R. Vershynin. **A randomized Kaczmarz algorithm with exponential convergence.** J. Fourier Anal. Appl., 15:262–278, 2009.

## Experimental results: synthetic data



- ▷ semisupervised reconstruction (40% labels) with three-layer structure ( $k^{(0)} = 9, k^{(1)} = 4, k^{(2)} = 2$ )

# Experimental results: synthetic data



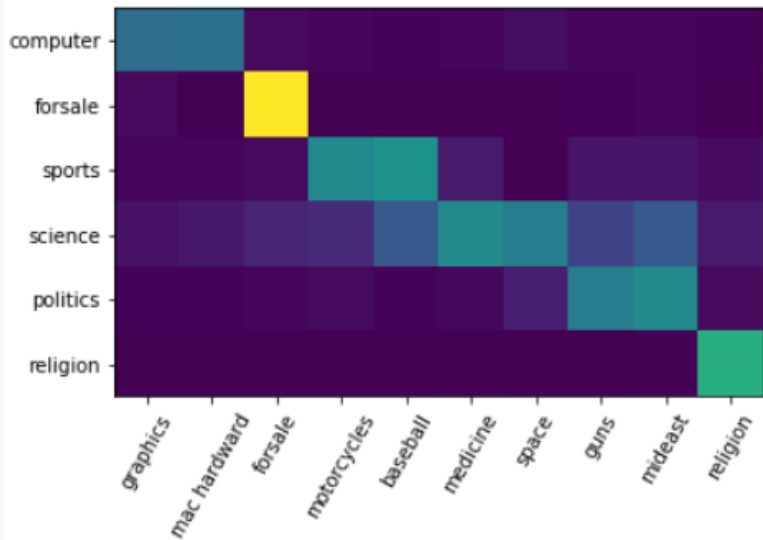
- ▷ semisupervised reconstruction (40% labels) with three-layer structure ( $k^{(0)} = 9, k^{(1)} = 4, k^{(2)} = 2$ )

# Experimental results: synthetic data

**Table 1:** Reconstruction error / classification accuracy

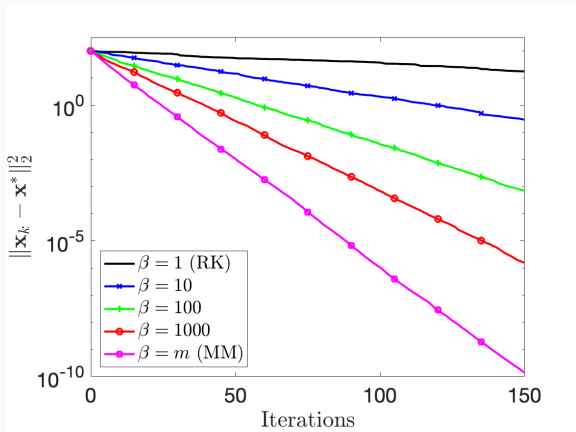
	Layers	Hier. NMF	Deep NMF	Neural NMF
Unsuper.	1	0.053	0.031	0.029
	2	0.399	0.414	<b>0.310</b>
	3	0.860	0.838	<b>0.492</b>
Semisuper.	1	0.049 / 0.933	0.031 / 0.947	0.042 / <b>1</b>
	2	0.374 / 0.926	0.394 / 0.911	<b>0.305</b> / <b>1</b>
	3	0.676 / 0.930	0.733 / 0.930	<b>0.496</b> / <b>0.990</b>
Supervised	1	0.052 / 0.960	0.042 / 0.962	0.042 / <b>1</b>
	2	0.311 / 0.984	0.310 / 0.984	0.307 / <b>1</b>
	3	0.495 / <b>1</b>	0.494 / <b>1</b>	0.498 / <b>1</b>

## Experimental results: 20 Newsgroups data



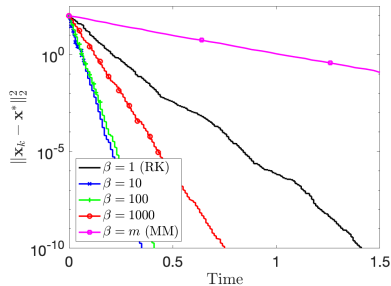
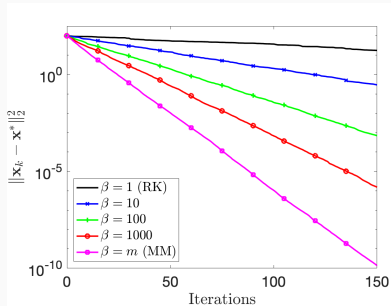


# Experimental Convergence



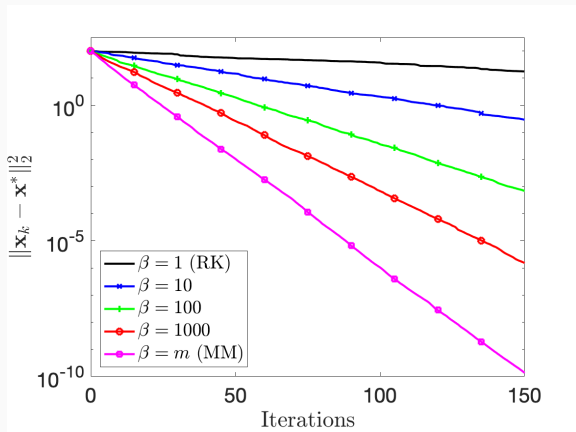
- ▷  $\beta$ : sample size
- ▷  $A$  is  $50000 \times 100$  Gaussian matrix, consistent system
- ▷ 'faster' convergence for larger sample size

# Experimental Convergence



- ▷  $\beta$ : sample size
- ▷  $A$  is  $50000 \times 100$  Gaussian matrix, consistent system
- ▷ 'faster' convergence for larger sample size

# Experimental Convergence



- ▷  $\beta$ : sample size
- ▷  $A$  is  $50000 \times 100$  Gaussian matrix, consistent system
- ▷ 'faster' convergence for larger sample size

**Goal:** Exploit similarities between neural networks and hierarchical NMF.

**Goal:** Exploit similarities between neural networks and hierarchical NMF.

▷ [Flenner, Hunter '18]

- introduces nonlinear pooling operator after each layer
- introduces multiplicative updates method meant to backpropagate

**Goal:** Exploit similarities between neural networks and hierarchical NMF.

- ▷ [Flenner, Hunter '18]
  - introduces nonlinear pooling operator after each layer
  - introduces multiplicative updates method meant to backpropagate
- ▷ [Trigeorgis, Bousmalis, Zafeiriou, Schuller '16]
  - relaxes some of nonnegativity constraints in hNMF

**Goal:** Exploit similarities between neural networks and hierarchical NMF.

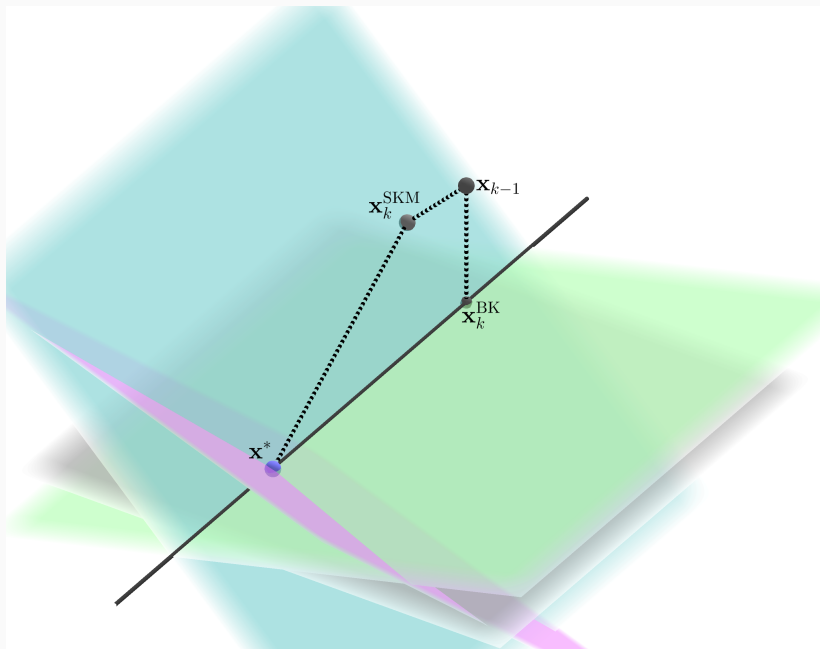
- ▷ [Flenner, Hunter '18]
  - introduces nonlinear pooling operator after each layer
  - introduces multiplicative updates method meant to backpropagate
- ▷ [Trigeorgis, Bousmalis, Zafeiriou, Schuller '16]
  - relaxes some of nonnegativity constraints in hNMF
- ▷ [Le Roux, Hershey, Weninger '15]
  - introduces NMF backpropagation algorithm with “unfolding” (no hierarchy)

**Goal:** Exploit similarities between neural networks and hierarchical NMF.

- ▷ [Flenner, Hunter '18]
  - introduces nonlinear pooling operator after each layer
  - introduces multiplicative updates method meant to backpropagate
- ▷ [Trigeorgis, Bousmalis, Zafeiriou, Schuller '16]
  - relaxes some of nonnegativity constraints in hNMF
- ▷ [Le Roux, Hershey, Weninger '15]
  - introduces NMF backpropagation algorithm with “unfolding” (no hierarchy)
- ▷ [Sun, Nasrabadi, Tran '17]
  - similar method lacking nonnegativity constraints



# Block Kaczmarz



## Bound on $\gamma_j$

$$\gamma_k \geq \frac{\beta}{m} \sigma_{\min}^2(A) \text{ when } A \text{ is row-normalized}$$