

# Matrix Decomposition: Applications and Algorithms

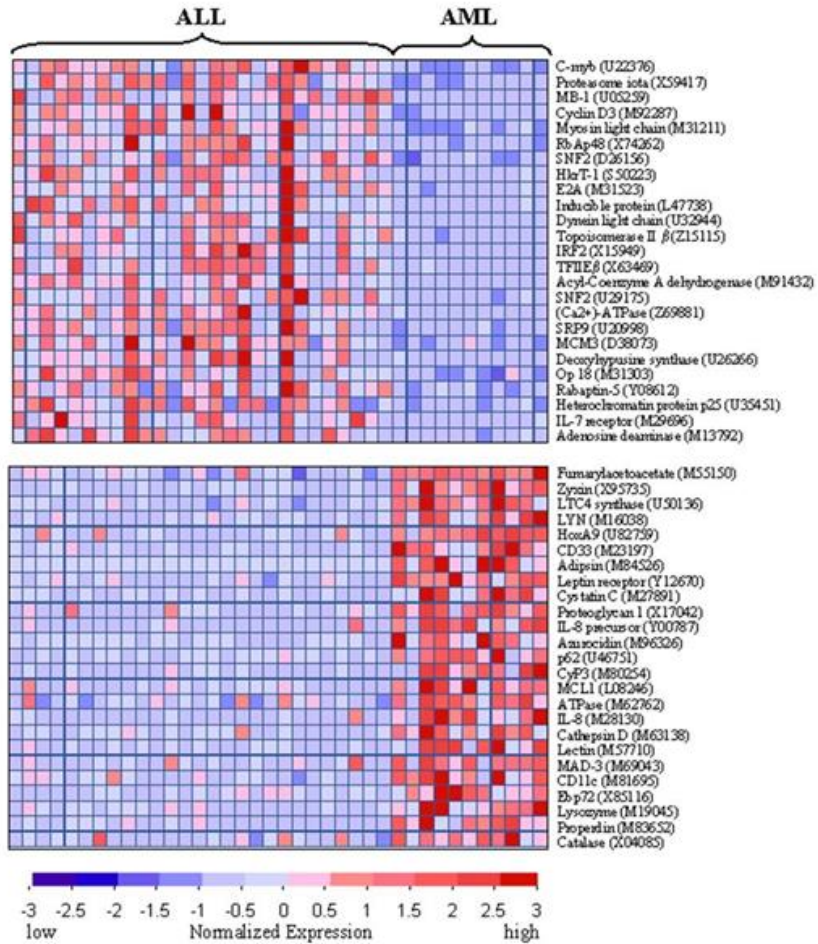
Praneeth Netrapalli

Microsoft Research India

# Extracting Information from Data Matrices

- Huge amounts of data available in a lot of domains: online search, social media, health records etc.
- Information *hidden* in the data; Very useful and has many applications
- Specific context: Data as matrices

# For example



Gene expression levels

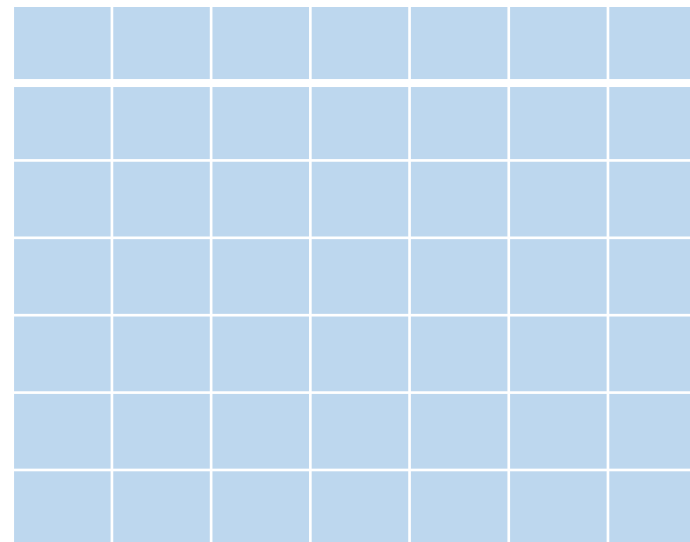
Image credit: Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring, T. R. Golub, et al., Science 286 (531) 1999

	Item 1	Item 2	Item 3	Item 4	Item 5
User 1	0	3	0	3	0
User 2	4	0	0	2	0
User 3	0	0	3	0	0
User 4	3	0	4	0	3
User 5	4	3	0	4	0

## User item ratings

Image credit:  
<https://katbailey.github.io/post/matrix-factorization-with-tensorflow/>

Frames



Frames in a video

# Matrix Decomposition

$$M \sim UV^T$$

$$s.t. U \in S_1, V \in S_2$$

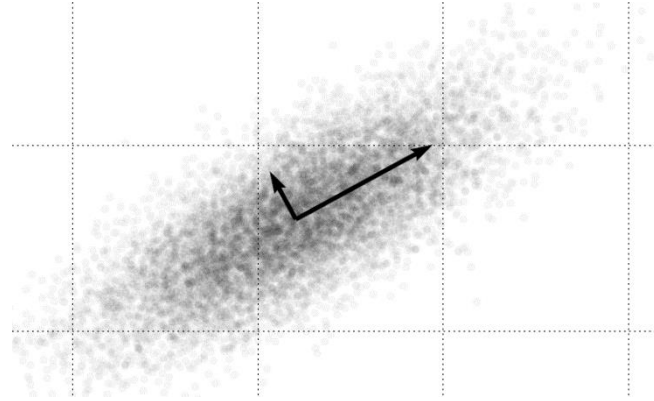
$$M \sim U + V$$

$$s.t. U \in S_1, V \in S_2$$

Helps extract information from matrices in several applications

Classical decompositions: Eigen, Schur, Cholesky, QR, Jordan,...

# Canonical Example: Principal Component Analysis



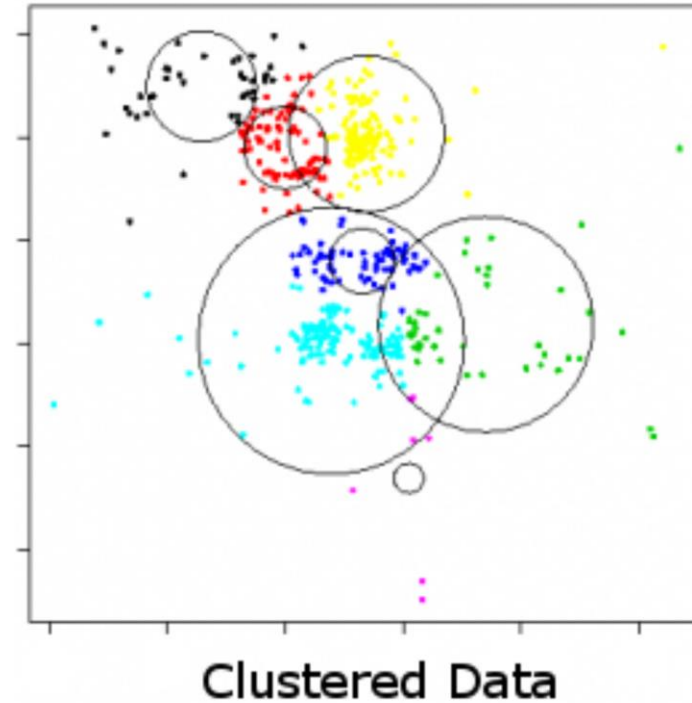
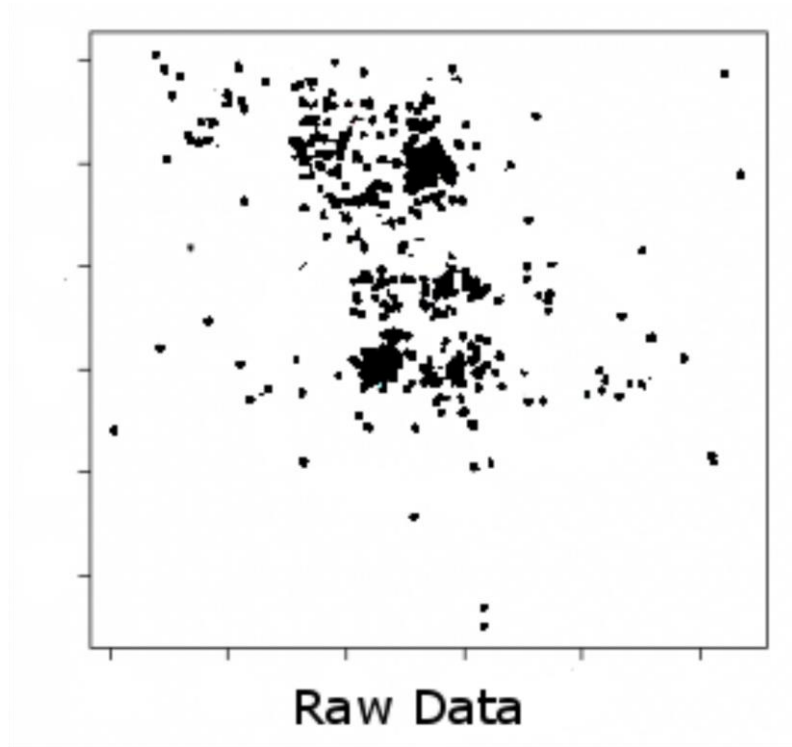
**Motivation:** Powerful tool for dimensionality reduction

# Canonical Example: Principal Component Analysis

$$\min_{U, V} \left\| \begin{array}{c} \square \\ M \\ \square \end{array} - \begin{array}{c} \square \\ U \\ \square \end{array} \begin{array}{c} \square \\ V^T \\ \square \end{array} \right\|_F^2$$

**Motivation:** Powerful tool for dimensionality reduction

# K-means Clustering



# K-means Clustering

$$\min_{U, V} \left\| M - UV^T \right\|_F^2$$

s.t. each column of  $V^T$  is an indicator vector



# In a Nutshell

- **Many more applications:** Matrix completion, Nonnegative PCA, Sparse coding, Subspace clustering ...
- **Several approaches**
  - Classical methods are either algebraic or based on convex relaxations
  - Not designed for modern data applications
    - Many passes over data not feasible
    - Noisy data

# Alternating Minimization (AltMin)

$$\min_{U \in S_1, V \in S_2} f(U, V)$$

1. Choose a random  $V$
2. Fix  $V$  and optimize over  $U$
3. Fix  $U$  and optimize over  $V$
4. Repeat Steps 2 & 3

In many cases, steps 2 and 3 efficiently solvable

**Empirically** widely used (e.g., k-means); has good performance

**Theoretically** significant progress in the last decade

# This talk

➤ Matrix completion

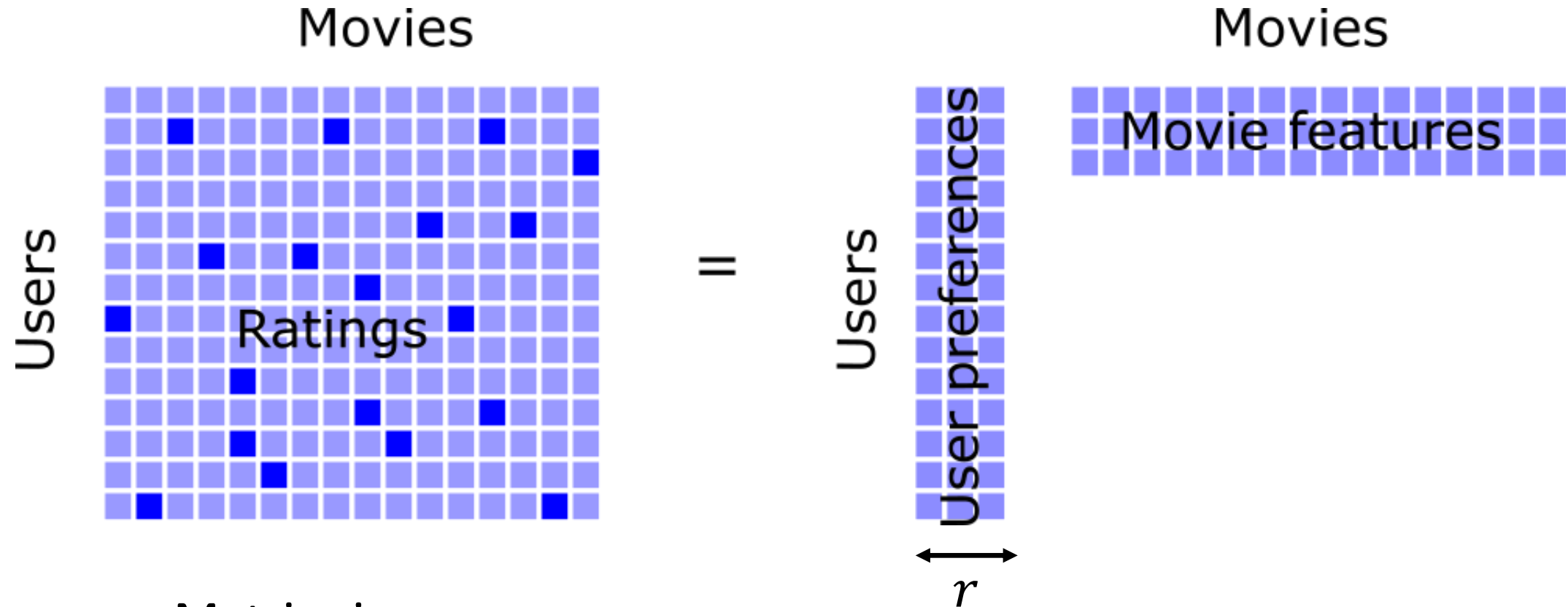
➤ Robust PCA

➤ Sparse coding

➤ Conclusion

# Part I: Matrix Completion

# Matrix Completion



- Matrix size:  $n \times m$
- $\text{Wlog } n \geq m$
- Rank =  $r \ll n, m$

**Goal:** Infer the entire matrix with the minimum number of observations.

# AltMin

$$\min_{U,V} \sum_{(i,j) \in \text{observed}} (M_{ij} - \underbrace{\langle U_i, V_j \rangle}_{\text{Predicted Rating}})^2$$

Diagram illustrating the AltMin optimization process:

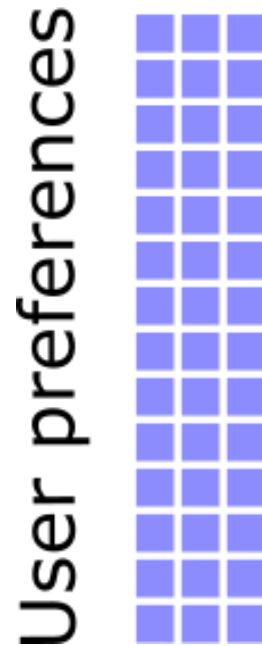
- The observed rating  $M_{ij}$  is compared against the predicted rating  $\langle U_i, V_j \rangle$ .
- The predicted rating is derived from the inner product of User  $i$ 's Preferences ( $U_i$ ) and Movie  $j$ 's Features ( $V_j$ ).
- The squared difference between the observed and predicted ratings is summed over all observed pairs  $(i,j)$ .

# AltMin

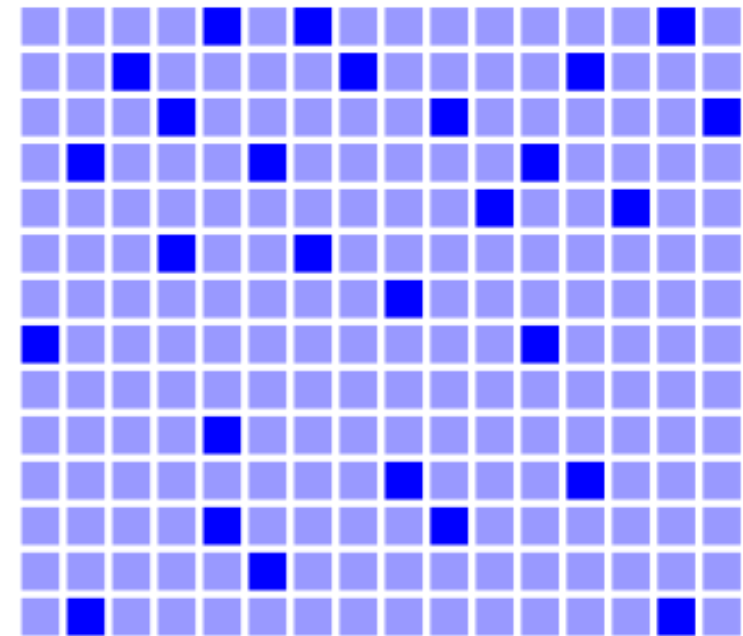
$$\min_{U,V} \sum_{(i,j) \in \text{observed}} (M_{ij} - \underbrace{\langle U_i, V_j \rangle}_{\text{Predicted Rating}})^2$$

Labels in diagram:  
-  $M_{ij}$ : Observed Rating  
-  $\langle U_i, V_j \rangle$ : Predicted Rating  
-  $U_i$ : User  $i$ 's Preferences  
-  $V_j$ : Movie  $j$ 's Features

Fix  $V$ , optimize over  $U$



Movie features



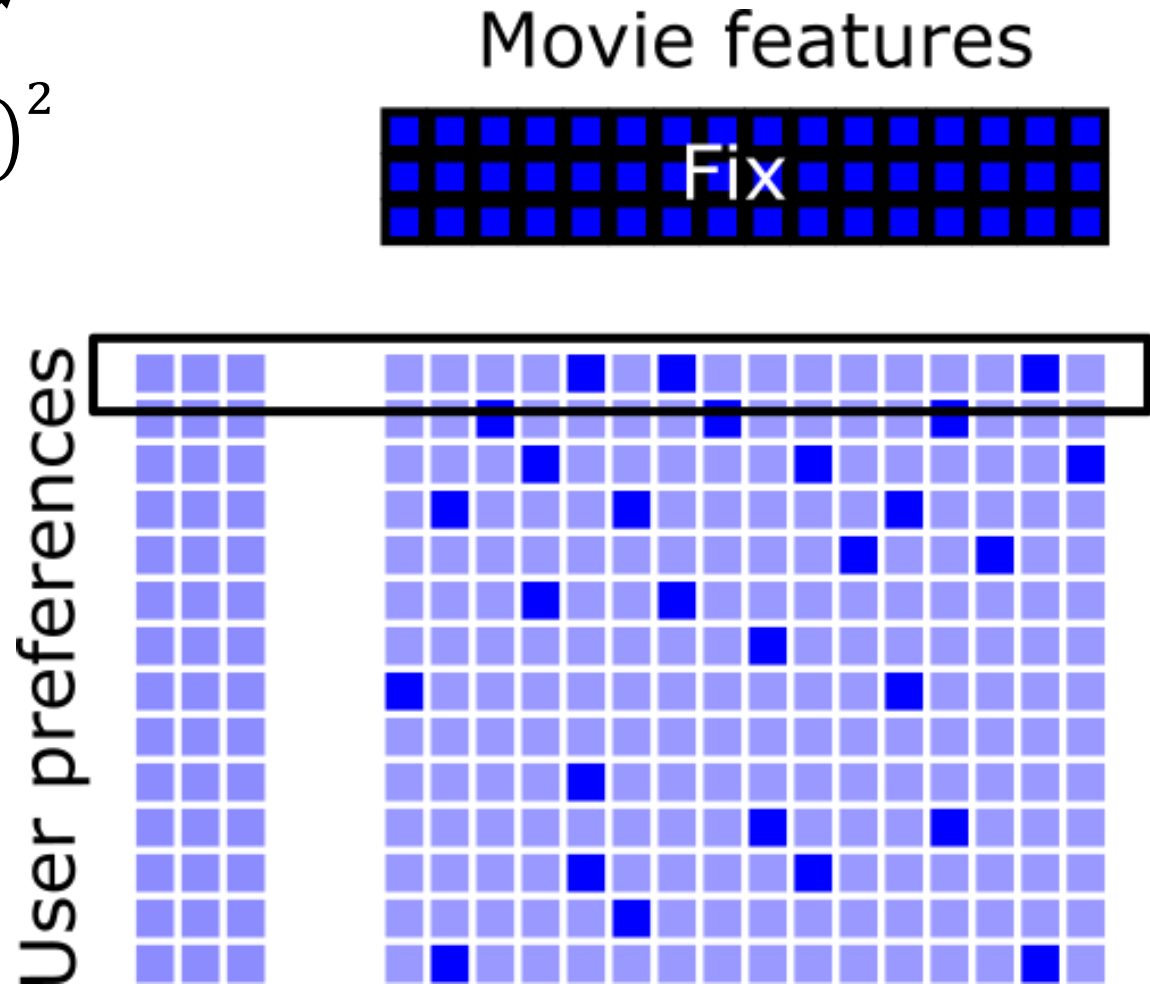
# AltMin

$$\min_{U, V} \sum_{(i, j) \in \text{observed}} (M_{ij} - \underbrace{\langle U_i, V_j \rangle}_{\text{Predicted Rating}})^2$$

Labels in diagram:  
-  $M_{ij}$ : Observed Rating  
-  $\langle U_i, V_j \rangle$ : Predicted Rating  
-  $U_i$ : User  $i$ 's Preferences  
-  $V_j$ : Movie  $j$ 's Features

Fix  $V$ , optimize over  $U$

Decouples into multiple least squares problems





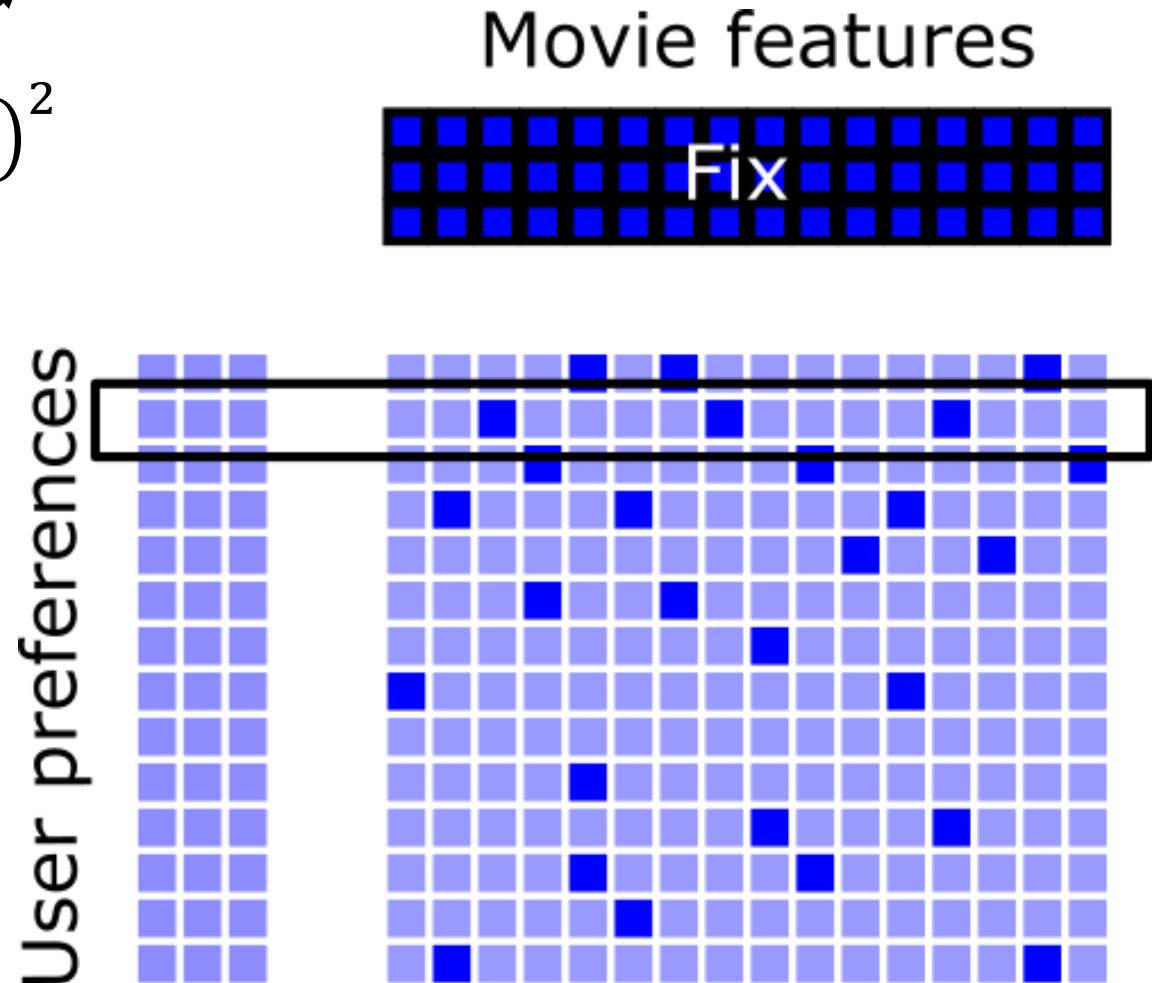
# AltMin

$$\min_{U, V} \sum_{(i, j) \in \text{observed}} (M_{ij} - \underbrace{\langle U_i, V_j \rangle}_{\text{Predicted Rating}})^2$$

Labels in diagram:  
-  $M_{ij}$ : Observed Rating  
-  $\langle U_i, V_j \rangle$ : Predicted Rating  
-  $U_i$ : User  $i$ 's Preferences  
-  $V_j$ : Movie  $j$ 's Features

Fix  $V$ , optimize over  $U$

Decouples into multiple least squares problems



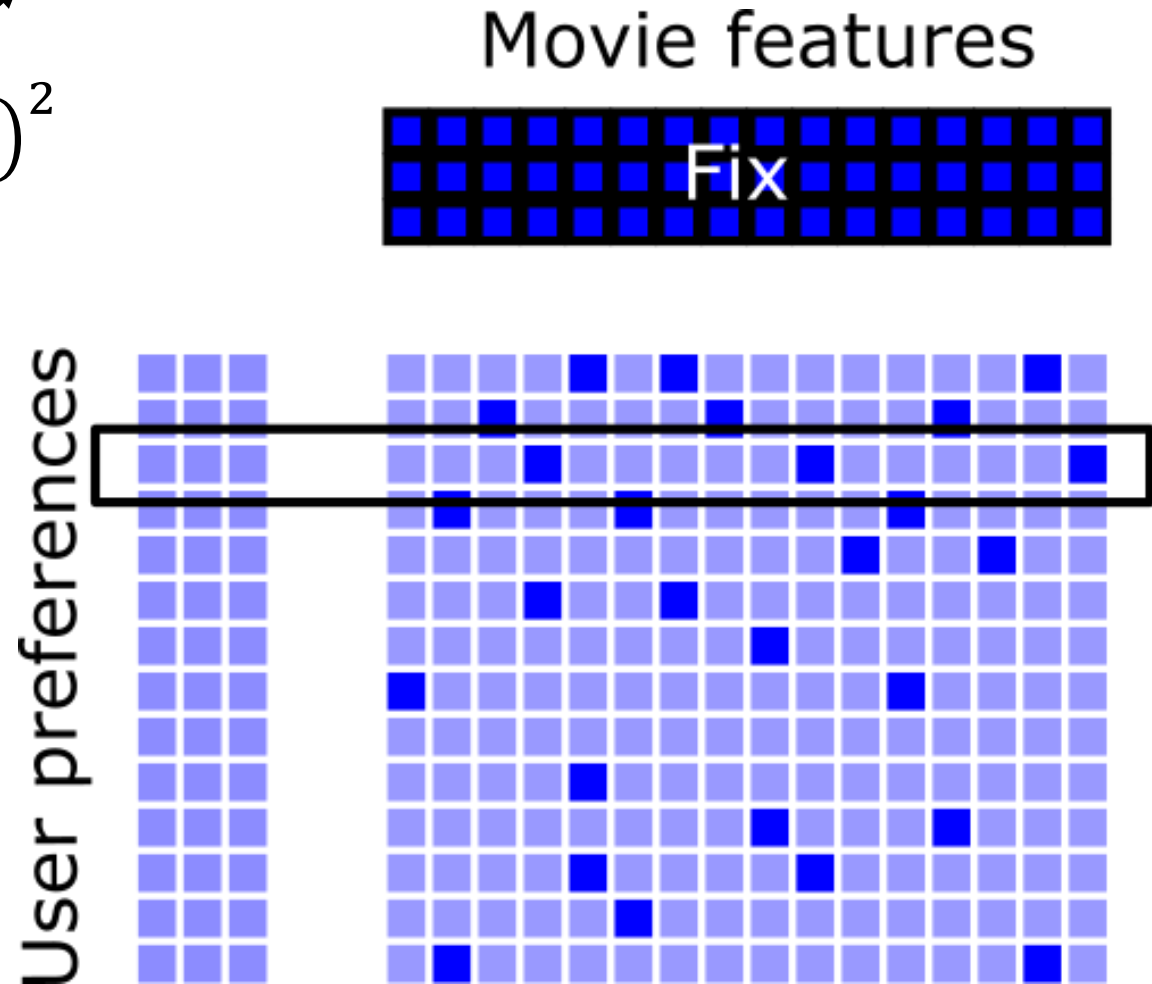
# AltMin

$$\min_{U,V} \sum_{(i,j) \in \text{observed}} (M_{ij} - \underbrace{\langle U_i, V_j \rangle}_{\text{Predicted Rating}})^2$$

Labels in diagram:  
-  $M_{ij}$ : Observed Rating  
-  $\langle U_i, V_j \rangle$ : Predicted Rating  
-  $U_i$ : User  $i$ 's Preferences  
-  $V_j$ : Movie  $j$ 's Features

Fix  $V$ , optimize over  $U$

Decouples into multiple least squares problems

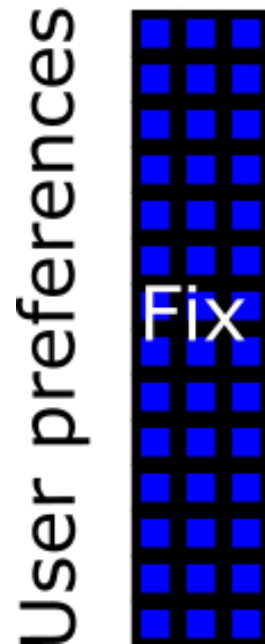


# AltMin

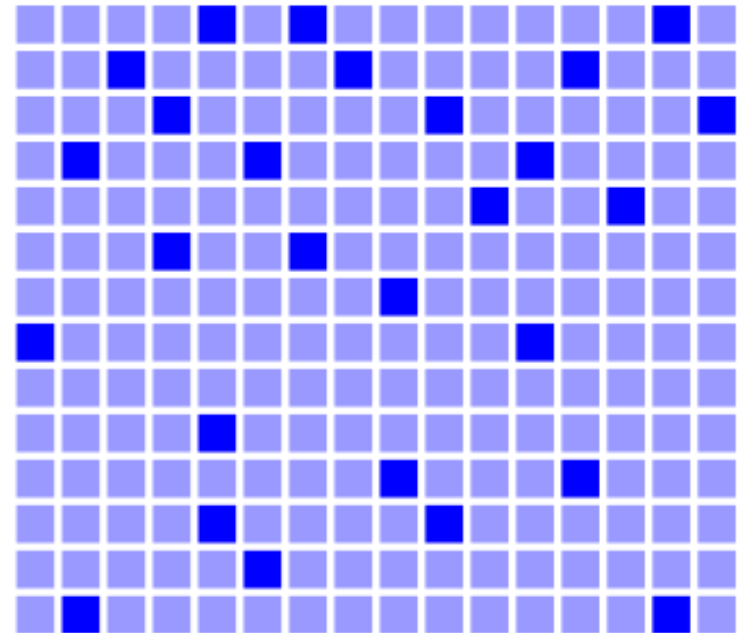
$$\min_{U,V} \sum_{(i,j) \in \text{observed}} (M_{ij} - \underbrace{\langle U_i, V_j \rangle}_{\text{Predicted Rating}})^2$$

Labels in diagram:  
User  $i$ 's Preferences (points to  $U_i$ )  
Movie  $j$ 's Features (points to  $V_j$ )  
Observed Rating (points to  $M_{ij}$ )  
Predicted Rating (points to  $\langle U_i, V_j \rangle$ )

Fix  $U$ , optimize over  $V$



Movie features



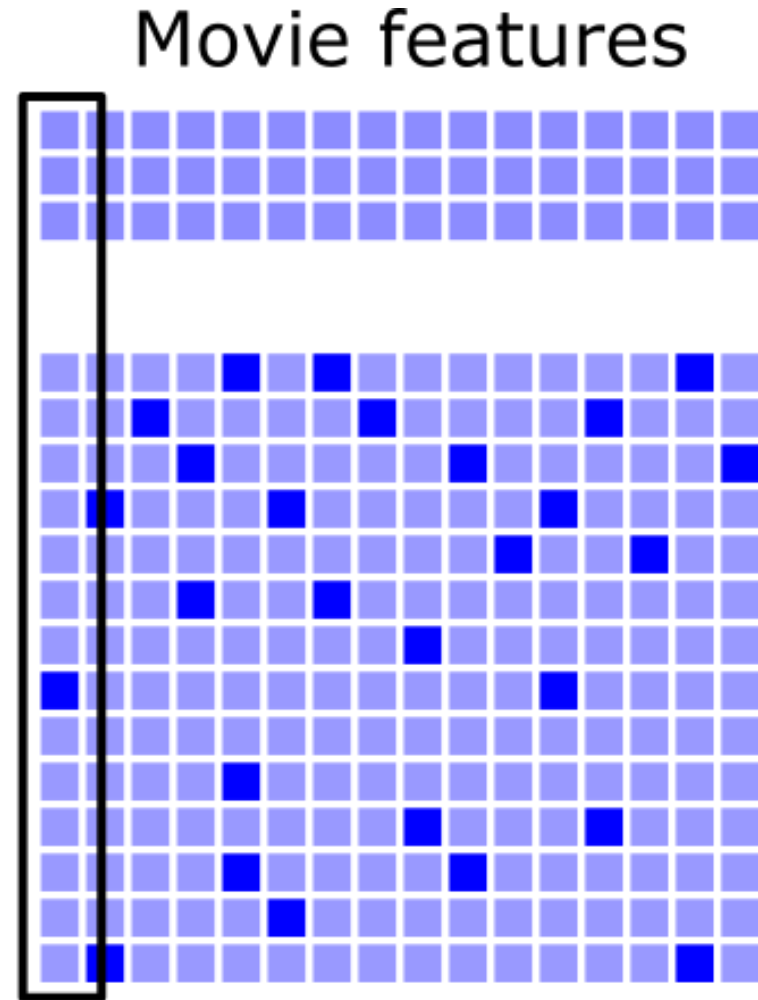
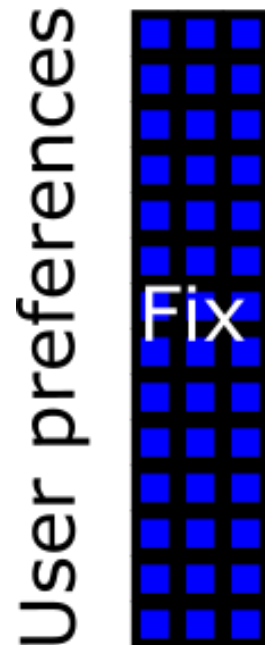
# AltMin

$$\min_{U,V} \sum_{(i,j) \in \text{observed}} (M_{ij} - \underbrace{\langle U_i, V_j \rangle}_{\text{Predicted Rating}})^2$$

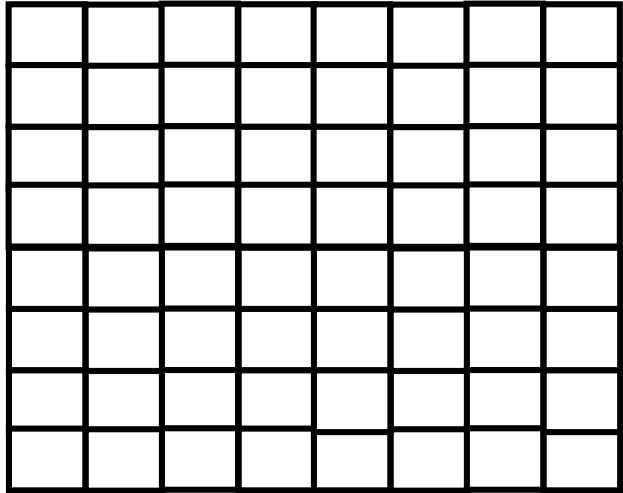
Labels in diagram:  
-  $M_{ij}$ : Observed Rating  
-  $\langle U_i, V_j \rangle$ : Predicted Rating  
-  $U_i$ : User  $i$ 's Preferences  
-  $V_j$ : Movie  $j$ 's Features

Fix  $U$ , optimize over  $V$

Again decouples into multiple least squares problems



# Need for Regularity Assumptions



# Need for Regularity Assumptions

	0				0		
			0				
0							
						0	
	0			0			
						0	
					0		
		0					0

# Need for Regularity Assumptions

?	0				0		
			0				
0							
						0	
	0			0			
						0	
					0		
		0					0

# Need for Regularity Assumptions

1	0				0		
			0				
0							
						0	
	0			0			
						0	
					0		
		0					0



# Need for Regularity Assumptions

2	0				0		
			0				
0							
						0	
	0			0			
						0	
					0		
		0					0

# Need for Regularity Assumptions

100	0				0		
			0				
0							
						0	
	0			0			
						0	
					0		
		0					0

- Prediction hard since user unique
- For good prediction, require no user to be unique

# Need for Regularity Assumptions

100	0				0		
			0				
0							
						0	
	0			0			
						0	
					0		
		0					0

- Prediction hard since user unique
- For good prediction, require no user to be unique

Non-uniqueness  $\cong$  Incoherence : if  $M = U\Sigma V^T$  SVD then

$$\max_i \left\{ \left\| \begin{array}{c} \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \square \end{array} \right\|_2 \right\} \rightarrow \|U_i\|_2 \leq \mu \sqrt{\frac{r}{n}} \quad \left( 1 \leq \mu \leq \sqrt{\frac{n}{r}} \right)$$

Similarly  $V$

Incoherence:  $\|U_i\|_2 \leq \mu \sqrt{\frac{r}{n}}$

$$r = 2$$

$\frac{n}{2}$  ↑

1	1	1	1	0	0	0	0
1	1	1	1	0	0	0	0
1	1	1	1	0	0	0	0
1	1	1	1	0	0	0	0
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1
0	0	0	0	1	1	1	1

Singular vectors →

$\sqrt{\frac{2}{n}}$  \*

1	0
1	0
1	0
1	0
0	1
0	1
0	1
0	1

- No user very unique
- $\mu = 1 \Leftrightarrow$  Very incoherent

Incoherence:  $\|U_i\|_2 \leq \mu \sqrt{\frac{r}{n}}$

$$r = 2$$

$\frac{n}{4}$   $\updownarrow$

1	1	0	0	0	0	0	0
1	1	0	0	0	0	0	0
0	0	1	1	1	1	1	1
0	0	1	1	1	1	1	1
0	0	1	1	1	1	1	1
0	0	1	1	1	1	1	1
0	0	1	1	1	1	1	1
0	0	1	1	1	1	1	1

Singular vectors  $\longrightarrow$

$\sqrt{\frac{2}{n}}$  \*

$\sqrt{2}$	0
$\sqrt{2}$	0
0	$\sqrt{2/3}$
0	$\sqrt{2/3}$
0	$\sqrt{2/3}$
0	$\sqrt{2/3}$
0	$\sqrt{2/3}$
0	$\sqrt{2/3}$

- Some users belong to a smaller group
- $\mu = \sqrt{2} \iff$  Still very incoherent, but a little less than earlier

Incoherence:  $\|U_i\|_2 \leq \mu \sqrt{\frac{r}{n}}$

$$r = 2$$

1  $\updownarrow$

1	0	0	0	0	0	0	0
0	1	1	1	1	1	1	1
0	1	1	1	1	1	1	1
0	1	1	1	1	1	1	1
0	1	1	1	1	1	1	1
0	1	1	1	1	1	1	1
0	1	1	1	1	1	1	1
0	1	1	1	1	1	1	1

Singular vectors  
 $\longrightarrow$

$\sqrt{\frac{2}{n}}$	$\sqrt{2}$	0
0	$\sqrt{1/2}$	
0	$\sqrt{1/2}$	
0	$\sqrt{1/2}$	
0	$\sqrt{1/2}$	
0	$\sqrt{1/2}$	
0	$\sqrt{1/2}$	
0	$\sqrt{1/2}$	

- User 1 very unique
- $\mu = \sqrt{\frac{n}{2}} \Leftrightarrow$  Not incoherent

# Regularity Assumptions

1. Low-rank matrix is **incoherent**: if  $M = U\Sigma V^T$  SVD then

$$\max_i \left\{ \begin{array}{c} U \\ \text{grid} \end{array} \right\} \rightarrow \|U_i\|_2 \leq \mu \sqrt{\frac{r}{n}} \quad \left( 1 \leq \mu \leq \sqrt{\frac{n}{r}} \right)$$

Similarly  $V$

2.  $d$  entries observed uniformly at random from all entries

# Result

## Theorem

AltMin recovers the full matrix if:

$$d \geq \text{poly}(\kappa, r) n \log n$$

Condition  
number

Rank

Dimension

[Keshavan 2012]

[Jain, Netrapalli, Sanghavi 2013]



# Result

## Theorem

AltMin recovers the full matrix if:

$$d \geq \text{poly}(\kappa, r) n \log n$$

Condition number      Rank      Dimension

[Keshavan 2012]

[Jain, Netrapalli, Sanghavi 2013]

Runtime	
AltMin	$\text{poly}(\kappa, r) n \log n$
Conv. Relaxation	$O(n^3)$

# Further Improvements

Resampling	}	<b>[Keshavan 2012]</b> <b>[Jain, Netrapalli, Sanghavi 2013]</b>	$\text{poly}(\kappa, r) n \log n$
		<b>[Hardt 2014]</b> <b>[Hardt, Wootters 2014]</b>	$\log \kappa \text{poly}(r) n \log n$
No resampling	}	<b>[Sun, Luo 2015]</b>	$\text{poly}(\kappa, r)n$
		<b>[Jain, Netrapalli 2015]</b>	$\text{poly}(r)n \log^3 n$

# Proof Outline (rank-1 case)

$$\min_{u,v} \sum_{(i,j) \text{ observed}} (M_{ij} - u_i v_j)^2 \xrightarrow{\text{AltMin}}$$

Approximate power  
method

# Proof Outline (rank-1 case)

$$\min_{u,v} \sum_{(i,j) \text{ observed}} (M_{ij} - u_i v_j)^2 \xrightarrow{\text{AltMin}}$$

Approximate power method

||| If all entries observed

$$\min_{u,v} \|M - uv^T\|_F^2$$

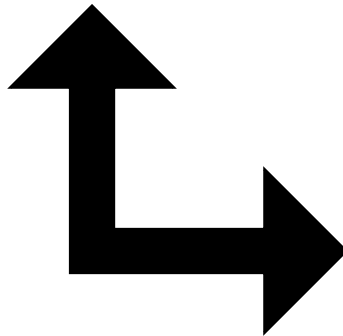
# Proof Outline (rank-1 case)

$$\min_{u,v} \sum_{(i,j) \text{ observed}} (M_{ij} - u_i v_j)^2 \xrightarrow{\text{AltMin}} \text{Approximate power method}$$

||| If all entries observed

$$\min_{u,v} \|M - uv^T\|_F^2$$

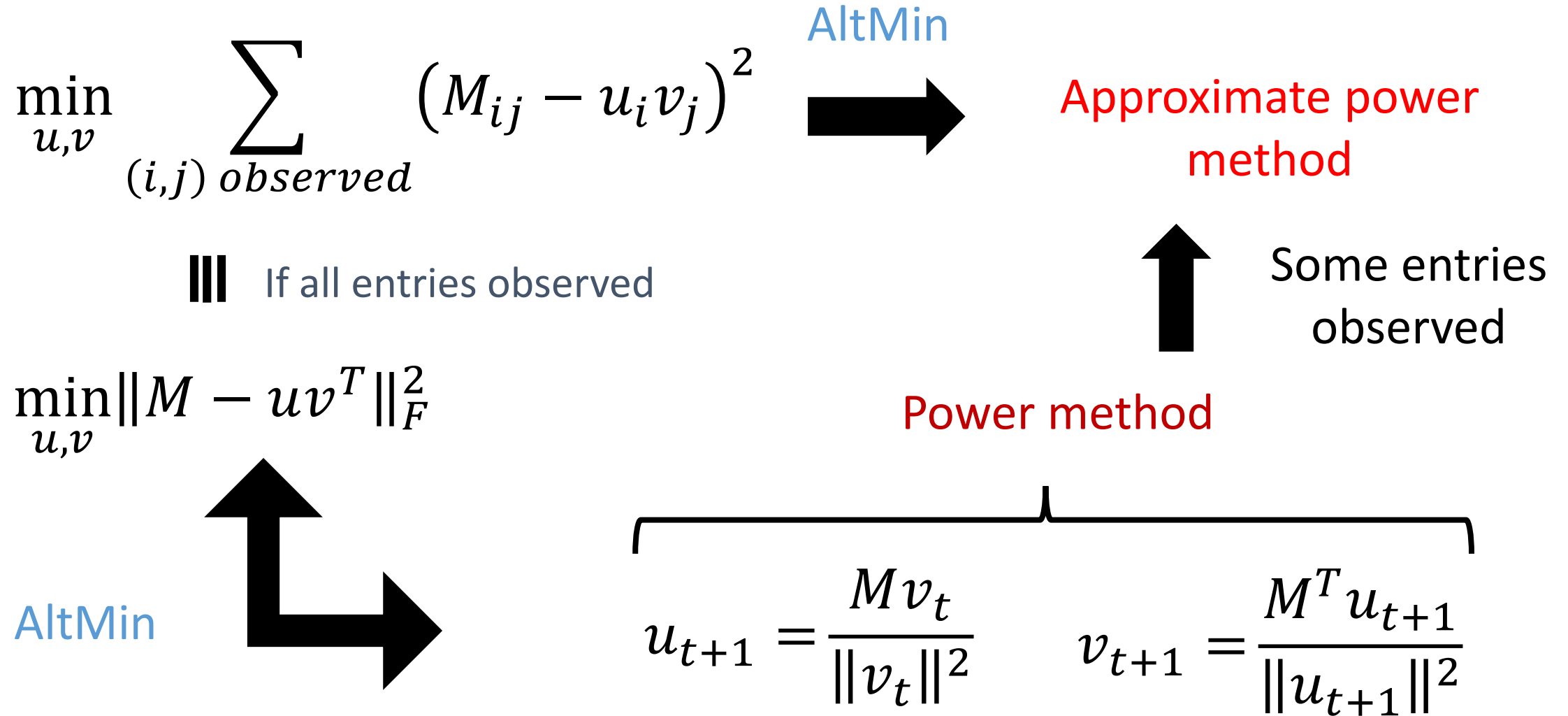
AltMin



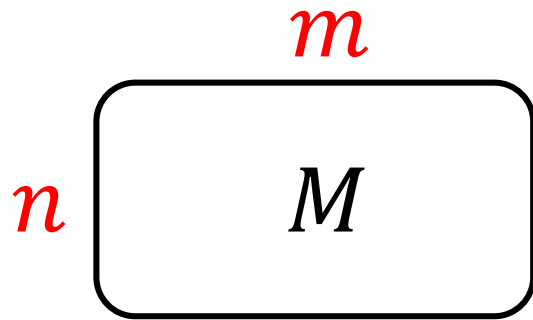
Power method

$$\left[ \begin{aligned} u_{t+1} &= \frac{M v_t}{\|v_t\|^2} & v_{t+1} &= \frac{M^T u_{t+1}}{\|u_{t+1}\|^2} \end{aligned} \right]$$

# Proof Outline (rank-1 case)

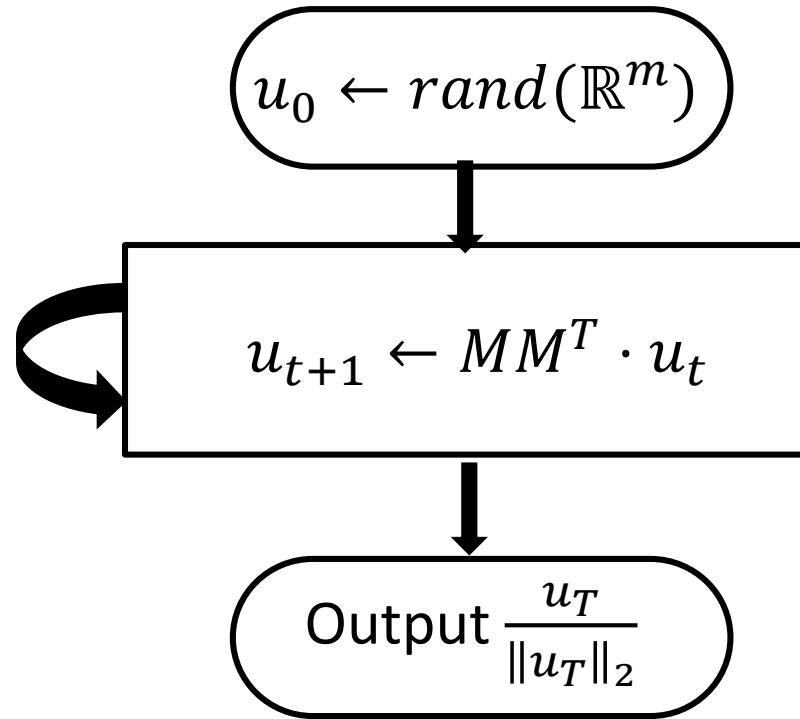


# Power Method for Computing Top Eigenvector



- $\text{Eigval}(MM^T): \lambda_1 \geq \lambda_2 \geq \dots$
- $\text{Eigvec}(MM^T): w_1, w_2, \dots$
- Eigengap:  $\delta \stackrel{\text{def}}{=} \frac{\lambda_1 - \lambda_2}{\lambda_1}$  (SNR)

Task  
Output  $u$  such that  
 $\|u - w_1\|_2 \leq \epsilon$



error  $\epsilon \Rightarrow O\left(\frac{1}{\delta} \log \frac{m}{\epsilon}\right)$  iterations

# Summary – Matrix Completion

- AltMin can provably solve low rank matrix completion
- AltMin performs approximate power method
- Improved understanding inspired the use of AltMin for many other problems

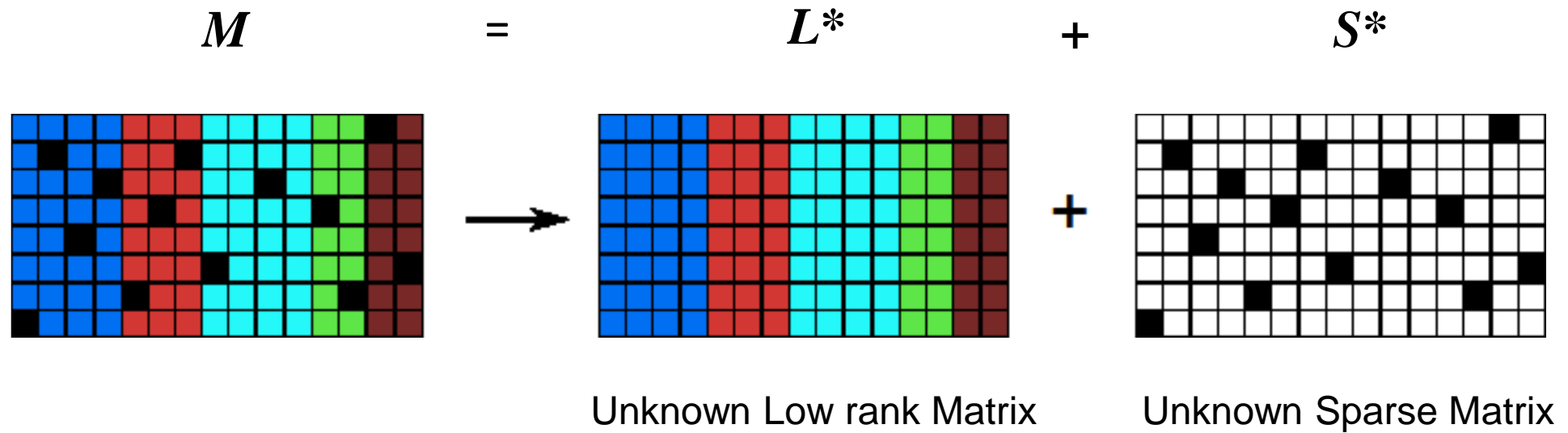


# Summary – Matrix Completion

Phase Retrieval	Netrapalli-Jain-Sanghavi 2013 Candes-Li-Soltanolkotabi 2014
Learning Sparsely Used Dictionaries	Arora-Ge-Moitra 2013 Agarwal-Anandkumar-Jain-Netrapalli 2013 Arora-Bhaskara-Ge-Ma 2014 Arora-Ge-Ma-Moitra 2014 Sun-Qu-Wright 2015
Mixed Linear Regression	Yi-Caramanis-Sanghavi 2014
Tensor decomposition	Anandkumar-Ge-JanZamin 2014 Jain-Oh 2014

## Part II: Robust PCA

# Robust PCA



**Task:** Given  $M$ , find  $L^*$  and  $S^*$

**Motivation:** PCA with missing/corrupted entries

# Foreground-Background Separation



Frames in a video



Background  
(Low rank)



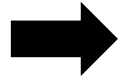
Foreground  
(Sparse)

# AltMin [Netrapalli et al. 2014]

$$\begin{aligned} \min_{L,S} & \|M - L - S\|_F^2 \\ \text{s.t.} & \text{rk}(L) = r \\ & S \text{ is sparse} \end{aligned}$$

# AltMin [Netrapalli et al. 2014]

$$\begin{aligned} \min_{L,S} & \|M - L - S\|_F^2 \\ \text{s.t.} & \text{rk}(L) = r \\ & S \text{ is sparse} \end{aligned}$$



Start:  $S = 0$

$$L \leftarrow \mathcal{P}_r(M - S)$$

Top- $r$  principal components

$$S \leftarrow \text{Threshold}(M - L)$$

Retain the largest elements

**Key Challenge:** Provable convergence for non-convex projections

# Regularity Assumptions

1	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$M$

=

1	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$L^*$

+

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$S^*$

# Regularity Assumptions

1	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$M$

=

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$L^*$

+

1	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$S^*$



# Regularity Assumptions

1	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$M$

=

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$L^*$

+

1	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$S^*$

**Issue:** Low rank matrix is also sparse

**Solution:** Ensure PCA directions not sparse – known as **incoherence**

# Incoherence

PCA  
directions

$$\sqrt{\frac{1}{n}} * \begin{array}{|c|} \hline 1 \\ \hline 1 \\ \hline 1 \\ \hline 1 \\ \hline 1 \\ \hline 1 \\ \hline 1 \\ \hline 1 \\ \hline \end{array}$$

Very incoherent

$$\sqrt{\frac{2}{n}} * \begin{array}{|c|} \hline 1 \\ \hline 1 \\ \hline 1 \\ \hline 1 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline \end{array}$$

Incoherent

$$\begin{array}{|c|} \hline 1 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline \end{array}$$

Not incoherent

# Regularity Assumptions

(deterministic)  $n \times n$  matrix  $M = L^* + S^*$  such that:

1. Rank- $r$  matrix is  **$\mu$ -incoherent**: if  $L^* = U\Sigma V^T$  then

$$\max_i \left\{ \begin{array}{c} U \\ \text{grid} \end{array} \right\} \rightarrow \|U_i\|_2 \leq \mu \sqrt{\frac{r}{n}} \quad \left( 1 \leq \mu \leq \sqrt{\frac{n}{r}} \right)$$

Similarly  $V$

2. Sparse matrix is **degree bounded**:

Any row or column of  $S^*$  has at most  $d$  non-zero elements

# Guarantee for AltMin

[Netrapalli et al. 2014]:  $L^*$  and  $S^*$  recovered by a non-convex algorithm if

$$d \leq \frac{n}{512\mu^2 r}$$

Information theoretically optimal  
up to constant factors

$$\text{Runtime: } O\left(r^2 n^2 \log \frac{1}{\epsilon}\right)$$

Factor  $r$  away from that of regular PCA

[Hsu, Kakade, Zhang 2011]:  $L^*, S^*$  recovered by a convex program if

$$d \leq c \frac{n}{\mu^2 r}$$

$$\text{Runtime: } O\left(\frac{n^3}{\epsilon^2}\right)$$

# Key Technical Result

$$L \leftarrow \mathcal{P}_r(M) = \mathcal{P}_r(L^* + S^*)$$

Need to show  $L \sim L^*$

$$\left. \begin{array}{l} \text{Incoherence of } L^* \\ \text{Sparsity of } S^* \end{array} \right\} \longrightarrow \|L - L^*\|_\infty \lesssim \frac{1}{10} \|S^*\|_\infty$$

**Standard Matrix Perturbation:**  $\|L - L^*\|_2 \leq \|S^*\|_2$

Does not use incoherence of  $L^*$ ; not sufficient

## Key Technique Contd.

$$L \leftarrow \mathcal{P}_r(M) = \mathcal{P}_r(L^* + S^*)$$

Rank 1:  $L^* = u^*u^{*T}$  and let  $L = uu^T$

## Key Technique Contd.

$$L \leftarrow \mathcal{P}_r(M) = \mathcal{P}_r(L^* + S^*)$$

Rank 1:  $L^* = u^*u^{*T}$  and let  $L = uu^T$   $\longrightarrow$   $(L^* + S^*)u = u$

Rearranging,  $u \approx (I - S^*)^{-1}L^*u \approx (I + S^* + S^{*2} + \dots)u^*$

# Key Technique Contd.

$$L \leftarrow \mathcal{P}_r(M) = \mathcal{P}_r(L^* + S^*)$$

Rank 1:  $L^* = u^*u^{*T}$  and let  $L = uu^T$   $\longrightarrow$   $(L^* + S^*)u = u$

Rearranging,  $u \approx (I - S^*)^{-1}L^*u \approx (I + S^* + S^{*2} + \dots)u^*$

Incoherence of  $L^*$   
Sparsity of  $S^*$   $\left. \vphantom{\begin{array}{l} \text{Incoherence of } L^* \\ \text{Sparsity of } S^* \end{array}} \right\} \longrightarrow \|u - u^*\|_\infty \lesssim d \|S^*\|_\infty \|u^*\|_\infty$



# Key Technique Contd.

$$L \leftarrow \mathcal{P}_r(M) = \mathcal{P}_r(L^* + S^*)$$

Rank 1:  $L^* = u^*u^{*T}$  and let  $L = uu^T$   $\longrightarrow$   $(L^* + S^*)u = u$

Rearranging,  $u \approx (I - S^*)^{-1}L^*u \approx (I + S^* + S^{*2} + \dots)u^*$

Incoherence of  $L^*$   
Sparsity of  $S^*$   $\left. \vphantom{\begin{matrix} \text{Incoherence of } L^* \\ \text{Sparsity of } S^* \end{matrix}} \right\} \begin{matrix} \longrightarrow \|u - u^*\|_\infty \lesssim d \|S^*\|_\infty \|u^*\|_\infty \\ \longrightarrow \|L - L^*\|_\infty \lesssim \frac{1}{10} \|S^*\|_\infty \end{matrix}$

# Key Technique Contd.

$$L \leftarrow \mathcal{P}_r(M) = \mathcal{P}_r(L^* + S^*)$$

Rank 1:  $L^* = u^*u^{*T}$  and let  $L = uu^T$   $\longrightarrow$   $(L^* + S^*)u = u$

Rearranging,  $u \approx (I - S^*)^{-1}L^*u \approx (I + S^* + S^{*2} + \dots)u^*$

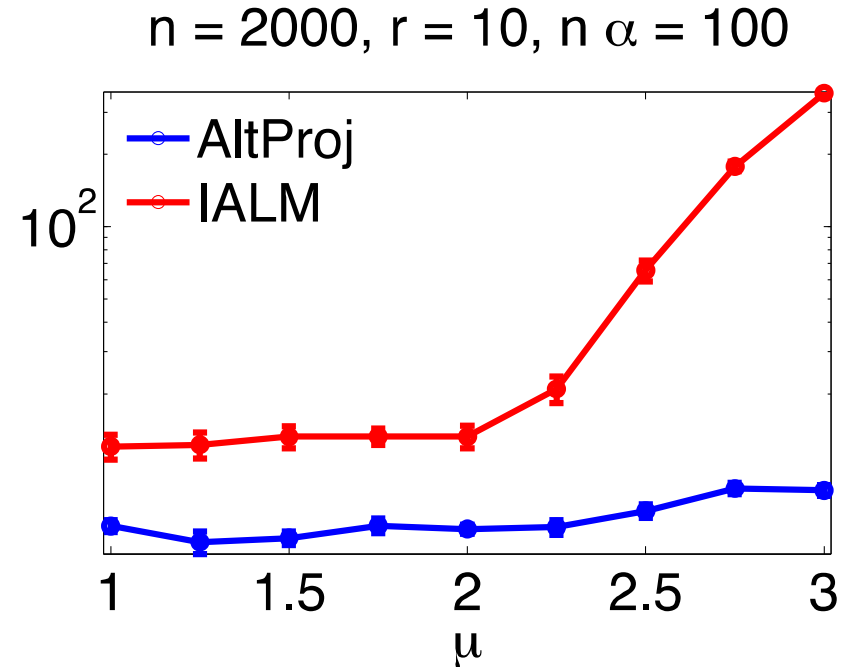
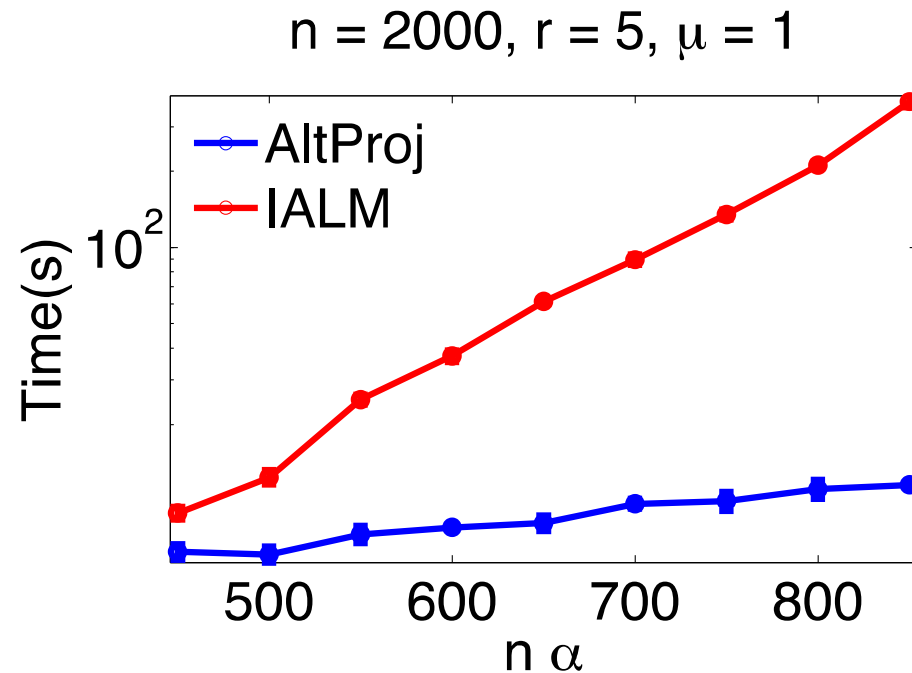


- Technique originally invented to study eigenvectors of random matrices: Erdos et al. 2011

# Comparison to Convex Solvers

IALM [Lin-Chen-Ma] : fastest / most popular special-purpose convex algorithm

Speed comparison, synthetic example. Single laptop.



# Empirical Results

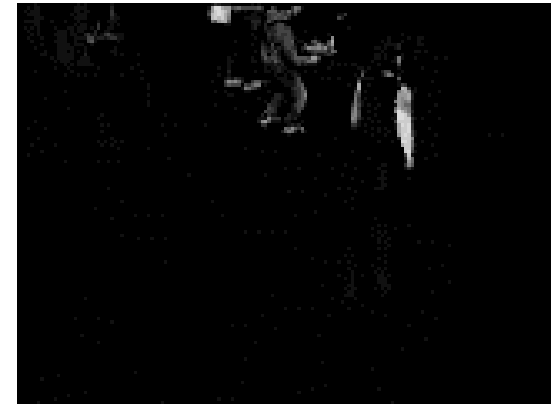
Convex Method. Runtime: 1700 sec



=



+



AltMin. Runtime: 70 sec



=



+



# Part III: Sparse Coding

# Sparse coding

 $Y$  $A^*$  $X^*$  $=$ 

Examples

Dictionary

Coefficients

Applications: Image compression, denoising, inpainting etc.

# Image reconstruction

50 % missing pixels



Learned reconstruction  
Average # coeffs: 4.0202  
MAE: 0.012977  
RMSE: 0.029204



Haar reconstruction OverComplete DCT reconstruction  
Average # coeffs: 4.7677 Average # coeffs: 4.7694  
MAE: 0.022833 MAE: 0.015719  
RMSE: 0.071107 RMSE: 0.037745



# Algorithm

$$\min_{A, X} \|Y - AX\|_F^2 \text{ s.t. } X \text{ is sparse}$$

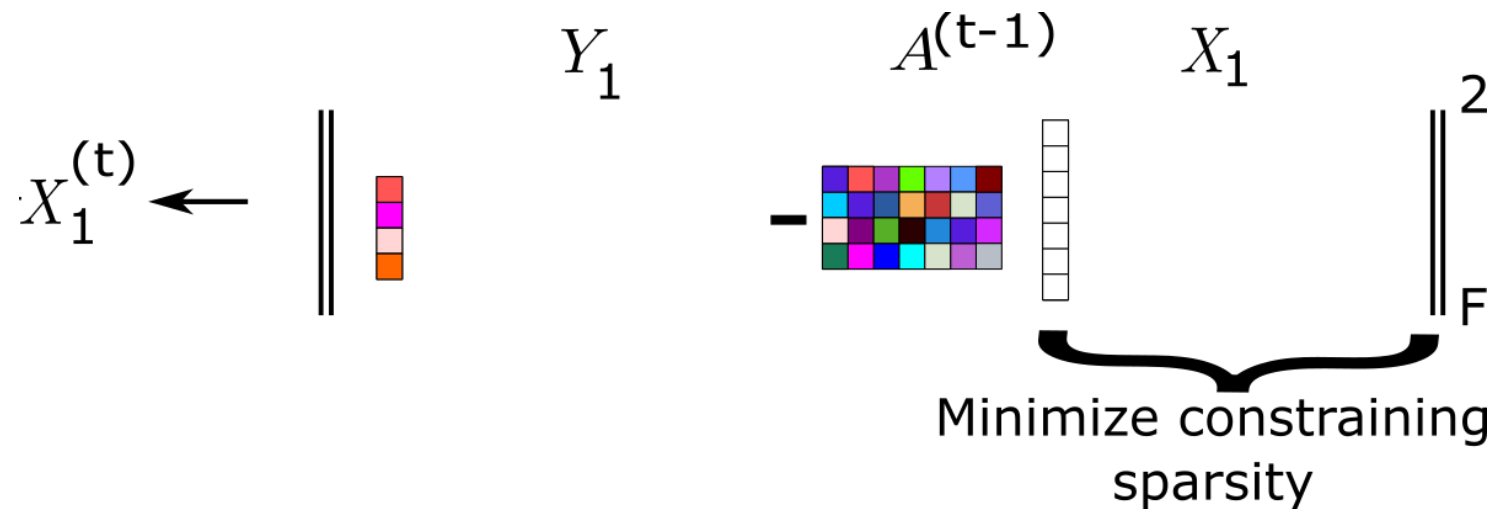
1. Choose an initial  $A^{(0)}$  using a **clustering** algorithm
2. In iteration  $t = 1, \dots$  do
  - Estimate coefficients  $X^{(t)}$  (Sparse recovery/compressed sensing)
  - Estimate dictionary  $A^{(t)}$  (Least squares)



# Algorithm

$$\min_{A, X} \|Y - AX\|_F^2 \text{ s.t. } X \text{ is sparse}$$

1. Choose an initial  $A^{(0)}$  using a **clustering** algorithm
2. In iteration  $t = 1, \dots$  do
  - Estimate coefficients  $X^{(t)}$  (Sparse recovery/compressed sensing)

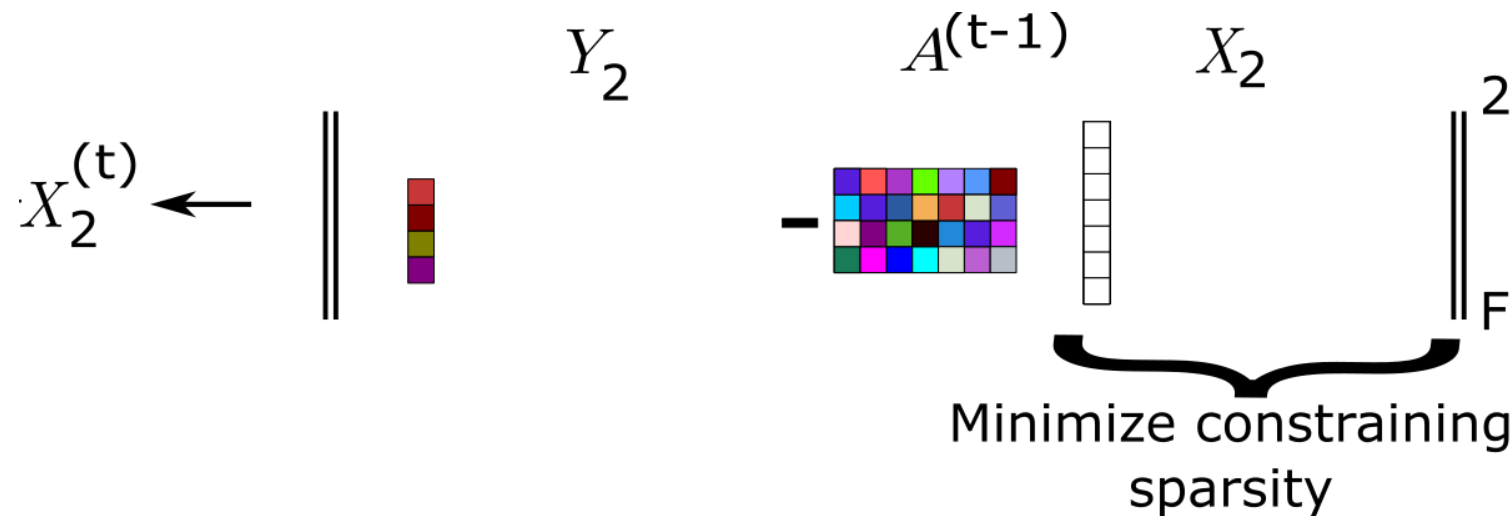


- Estimate dictionary  $A^{(t)}$  (Least squares)

# Algorithm

$$\min_{A, X} \|Y - AX\|_F^2 \text{ s.t. } X \text{ is sparse}$$

1. Choose an initial  $A^{(0)}$  using a **clustering** algorithm
2. In iteration  $t = 1, \dots$  do
  - Estimate coefficients  $X^{(t)}$  (Sparse recovery/compressed sensing)

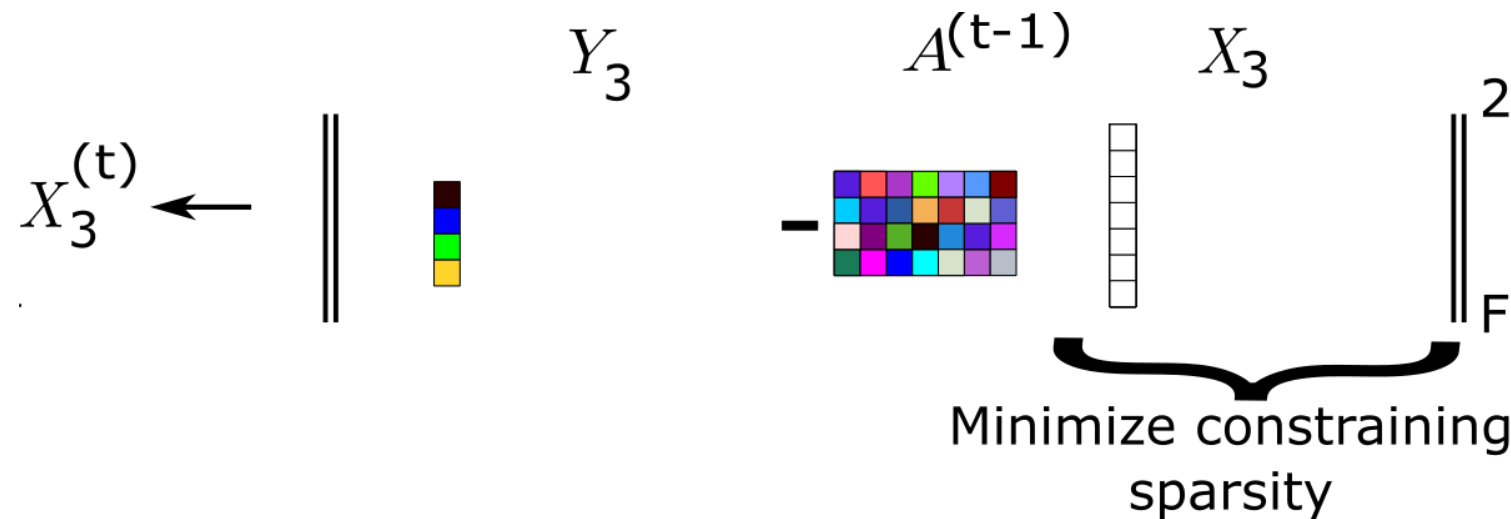


- Estimate dictionary  $A^{(t)}$  (Least squares)

# Algorithm

$$\min_{A, X} \|Y - AX\|_F^2 \text{ s.t. } X \text{ is sparse}$$

1. Choose an initial  $A^{(0)}$  using a **clustering** algorithm
2. In iteration  $t = 1, \dots$  do
  - Estimate coefficients  $X^{(t)}$  (Sparse recovery/compressed sensing)

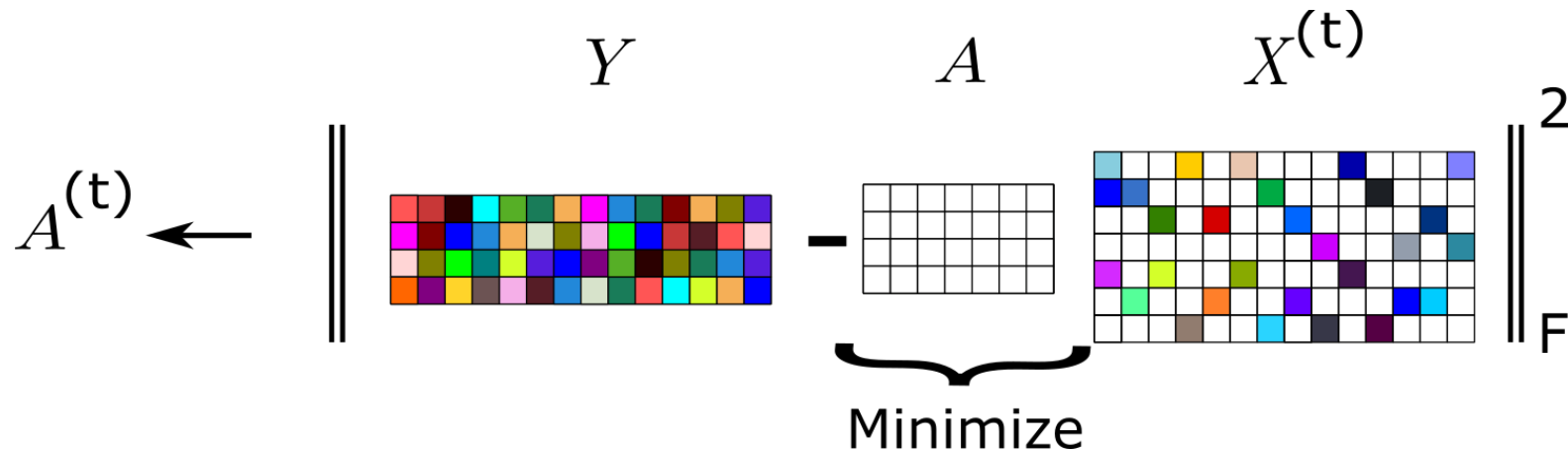


- Estimate dictionary  $A^{(t)}$  (Least squares)

# Algorithm

$$\min_{A, X} \|Y - AX\|_F^2 \text{ s.t. } X \text{ is sparse}$$

1. Choose an initial  $A^{(0)}$  using a **clustering** algorithm
2. In iteration  $t = 1, \dots$  do
  - Estimate coefficients  $X^{(t)}$  (Sparse recovery/compressed sensing)
  - Estimate dictionary  $A^{(t)}$  (Least squares)



# Result

Under some nondegeneracy conditions (incoherent dictionary and bounds on sparsity), recovers the decomposition.

[Arora, Ge, Moitra 2013], [Agarwal et al. 2013]

Many more results for simpler algorithms [Arora, Ge, Ma, Moitra 2015], noisy setting, nongenerative assumptions [Bhaskara, Tai 2019] etc.

# Conclusion

- Matrix decomposition is an extremely valuable tool for data analysis
- Finding the right structure is important
- AltMin is a very successful meta-algorithm
- Theoretical insights into AltMin help design new variants for new settings