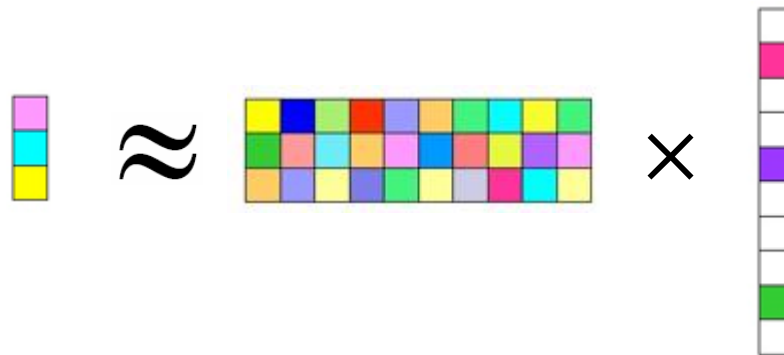


# Sparse Dictionary Learning: Algorithm Comparison and Challenges

$$b = Ax$$



Kerry Sun  
[sunx0486@umn.edu](mailto:sunx0486@umn.edu)  
05/06/2019



UNIVERSITY OF MINNESOTA

Driven to Discover<sup>SM</sup> 1/31

# Outline

- Background
  - Problem set-up
  - Motivation
  - Math
- Algorithms
- Experiments
- Limitations/challenges



# What is sparse Dictionary Learning?

Sparse vector      Input data vector

Dictionary

$$\mathbf{D}\mathbf{w} = \mathbf{x}, \mathbf{D} \in \mathbb{R}^{N \times K},$$

where  $K \gg N$  and  $\|\mathbf{w}\|_0 = k_0 \ll N$

$\mathbf{x} = \sum_{j=1}^K w_j \mathbf{d}_j$  where  $\mathbf{D} = \{\mathbf{d}_j\}$  of  $K$  vector,  $K > N$ .  $\mathbf{D}$  is an overcomplete set

1. Sparse Coding (obtaining a sparse  $\mathbf{w}$ )
2. Learning Dictionary (obtaining a matrix  $\mathbf{D}$ )

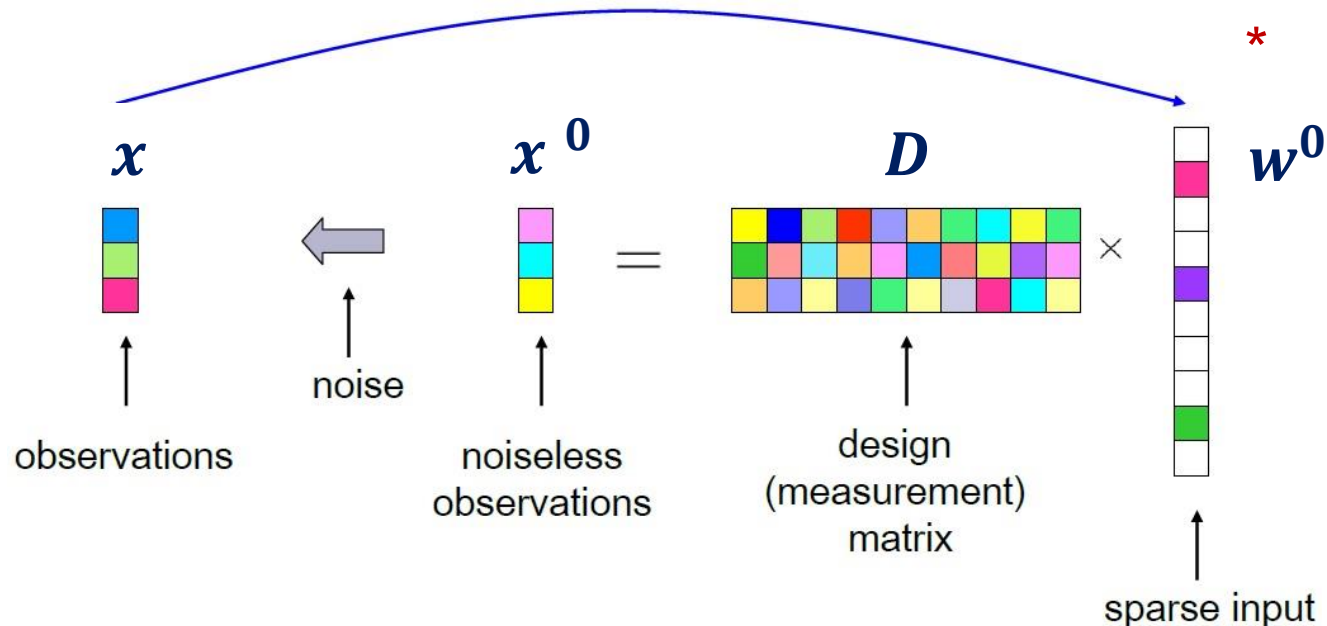


# What is sparse Dictionary Learning?

Sparse vector  
Data (observation) vector  
Dictionary

$$Dw = x, D \in \mathbb{R}^{N \times K},$$

where  $K \gg N$  and  $\|w\|_0 = k_0 \ll N$



# Why sparse Dictionary Learning?

Sparse vector      Input data vector

Dictionary →  $D\mathbf{w} = \mathbf{x}, D \in \mathbb{R}^{N \times K},$   
where  $K \gg N$  and  $\|\mathbf{w}\|_0 = k_0 \ll N$

1. Not always can find an optimal transform; want to represent input data using as few components as possible
2. Pre-constructed dictionaries (e.g. Fourier, wavelets, DCT\*) are usually **restricted** to signals/images of a certain type
3. Dictionary can be **learned** from the input data
4. Can be used in **many applications** such as Compressed sensing, signal recovery, de-noising.

# Applications

- Analysis

Given  $\mathbf{x}$ , can we determine the underlying vector  $\mathbf{w}_0$ ?

$$\|\mathbf{x} - \mathbf{D}\mathbf{w}_0\|_2 \leq \epsilon \quad \rightarrow \quad (P_0^\epsilon): \min_x \|\mathbf{w}\|_0 \text{ subject to } \|\mathbf{x} - \mathbf{D}\mathbf{w}\|_2 \leq \epsilon$$

- Compression

- De-noising

$$\tilde{\mathbf{x}} = \mathbf{x} + \mathbf{v}, \text{ where } \|\mathbf{v}\|_2 \leq \delta, \rightarrow (P_0^{\epsilon+\delta}): \min_x \|\mathbf{w}\|_0 \text{ subject to } \|\tilde{\mathbf{x}} - \mathbf{D}\mathbf{w}\|_2 \leq \epsilon + \delta$$

find approximation by  $\mathbf{D}\mathbf{w}_0^{\epsilon+\delta}$

- Compressed-Sensing

- Inverse Problems

- Morphological Component Analysis (MCA)

- Can we separate 2 sources?

$$\min_{\mathbf{w}_1, \mathbf{w}_2} \|\mathbf{w}_1\|_0 + \|\mathbf{w}_2\|_0 \text{ subject to } \|\mathbf{x} - \mathbf{D}_1\mathbf{w}_1 - \mathbf{D}_2\mathbf{w}_2\|_2^2 \leq \epsilon_1^2 + \epsilon_2^2$$



# Application-Denoising

Original clean image



Noisy image, 20.1721dB



Clean Image by Global Trained dictionary, 27.5672dB \*



\*



Inpainting result using the local-K-SVD

Left: the original image, Center: the degraded image with red text representing missing pixels, right: the recovered image (PSNR= 32.45dB)



# Application-Image Compression



Original

JPEG (15.81)

JPEG2000 (13.89)

PCA (10.66)

K-SVD (6.67)



Original

JPEG (14.67)

JPEG2000 (12.41)

PCA (9.44)

K-SVD (5.63)



Original

JPEG (15.3)

JPEG2000 (12.57)

PCA (10.27)

K-SVD (6.45)

\*

Values  
in () is RMSE



# Math Background

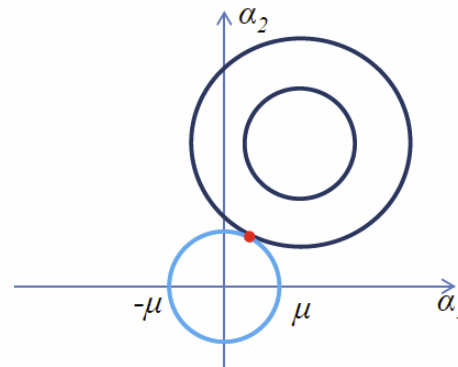
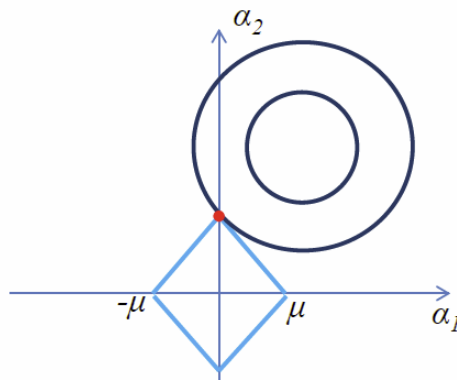
– the  $l_2$  norm.  $\| \alpha \|_2^2 \triangleq \sum_{i=1}^m \alpha_i^2$

– the  $l_0$  norm.  $\| \alpha \|_0 \triangleq \#\{i \mid \alpha_i \neq 0\}$

– the  $l_1$  norm.  $\| \alpha \|_1 \triangleq \sum_{i=1}^m |\alpha_i|$

} Sparsity inducing

$$\begin{aligned} \min_{\alpha \in \mathbb{R}} \frac{1}{2} \|x - \alpha\|_2^2 + \lambda \| \alpha \|_1 & \quad \min_{\alpha \in \mathbb{R}} \frac{1}{2} \|x - \alpha\|_2^2 + \lambda \| \alpha \|_2^2 \\ \Leftrightarrow \min_{\alpha \in \mathbb{R}} \frac{1}{2} \|x - \alpha\|_2^2 \text{ s.t. } \| \alpha \|_1 \leq \mu & \quad \Leftrightarrow \min_{\alpha \in \mathbb{R}} \frac{1}{2} \|x - \alpha\|_2^2 \text{ s.t. } \| \alpha \|_2 \leq \mu \end{aligned}$$



# Problem Definition (Setup)

**Given:** the input dataset  $\mathbf{X} = [x_1, \dots, x_L], x_l \in \mathbb{R}^N$

**Objective:** find

1) a dictionary  $\mathbf{D} \in \mathbb{R}^{N \times K}$

2) a coefficient set  $\mathbf{W} = [w_1, \dots, w_L], w_l \in \mathbb{R}^K$

Such that both  $\|\mathbf{X} - \mathbf{DW}\|_F^2$  is minimized\* and  $w_l$  are sparse ( $N \ll K$ )

$$\arg \min_{\mathbf{D} \in \mathcal{C}, \mathbf{w}_l \in \mathbb{R}^K} \sum_{l=1}^L \|\mathbf{x}_l - \mathbf{D}\mathbf{w}_l\|_2^2 + \lambda \|\mathbf{w}_l\|_p, \quad p \in [0, 1]$$

$$\mathcal{C} \equiv \{\mathbf{D} \in \mathbb{R}^{N \times K} : \|\mathbf{d}_i\|_2 \leq 1 \quad \forall i = 1, \dots, K\}$$

$$\lambda > 0$$

$$\|\mathbf{X} - \mathbf{DW}\|_F^2 = \sum_l \|\mathbf{x}_l - \mathbf{D}\mathbf{w}_l\|_2^2$$



# Algorithms

**Objective:** Find an optimal **sparse coding**  $W$  and a **dictionary**  $D$

**General strategy:** Split the problem into **2 parts**:

- 1) Keep  $D$  fixed, find  $W$  (sparse coding)
- 2) Keep  $W$  fixed, find  $D$  (dictionary learning)

## Popular Algorithms:

(combine **sparse coding** and **dictionary learning** together)

- MOD or ILS-DLA (1999)
- K-SVD – Compared to K-means (2006)
- RLS-DLA (2010)
- ODL (LASSO) (2010)



# Spare Coding

Keeping  $\mathbf{D}$  fixed and find  $\mathbf{W}$ :  $L$  independent problem

$$\arg \min_{\mathbf{D} \in \mathcal{C}, \mathbf{w}_l \in \mathbb{R}^K} \sum_{l=1}^L \|\mathbf{x}_l - \mathbf{D}\mathbf{w}_l\|_2^2 + \lambda \|\mathbf{w}_l\|_p, \quad p \in [0, 1]$$

**Strategy:** Vector selection algorithm

- **Pursuit algorithm** (e.g. Matching Pursuit (MP) and OMP based on  $l_0$  norm)
- LARS and LASSO algorithm (based on  $l_1$  norm)
- IRLS algorithm (based on  $l_p$  norm)
- Dantzig-Selector algorithm (based on  $l_\infty$  norm)



# Spare Coding - MP

What is matching pursuit (MP) algorithm?

it is a **sparse** approximation algorithm which finds the “best matching” projections of data onto of  $\mathbf{D}$

**Idea:** suppose optimal  $S = 1$ , so  $\mathbf{x} = c \cdot \mathbf{w}$

$$\epsilon(k) = \min_{w_k} \|\mathbf{d}_k w_k - \mathbf{x}\|_2^2, \quad w_k^* = \mathbf{d}_k^T \mathbf{x} / \|\mathbf{d}_k\|_2^2$$

A series of locally optimal single-term updates from  $\mathbf{w}^0 = \mathbf{0}$ , and iteratively construct a k-term  $\mathbf{w}^k$



# Spare Coding - MP

What is matching pursuit (MP) algorithm?

$$\min_{\mathbf{w} \in \mathbb{R}^K} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{w}\|_2^2 \quad s.t. \quad \|\mathbf{w}\|_0 \leq S$$

---

**Algorithm 1** Batch Matching Pursuit algorithm

---

**Input** : Signal  $\mathbf{x}$  and a normalized dictionary  $\mathbf{D}$

**Output:** List of coefficients  $(w_k)_{k=1}^K$  and indices for corresponding atoms  $(\gamma_k)_{k=1}^K$

**Initialization :**

$\mathbf{r} \leftarrow \mathbf{x}, \mathbf{w} \leftarrow \mathbf{0}$

**while** not Finished:

$\mathbf{c} \leftarrow \mathbf{D}^T \mathbf{r}$  (inner product)

find  $k : \operatorname{argmax}_k |c_k|$   Find the index of biggest  $c_k$

$w_k \leftarrow w_k + c_k$

$\mathbf{r} \leftarrow \mathbf{r} - c_k \cdot \mathbf{d}_k$

Finished:  $|\mathbf{r}| \leq \text{some limit}$

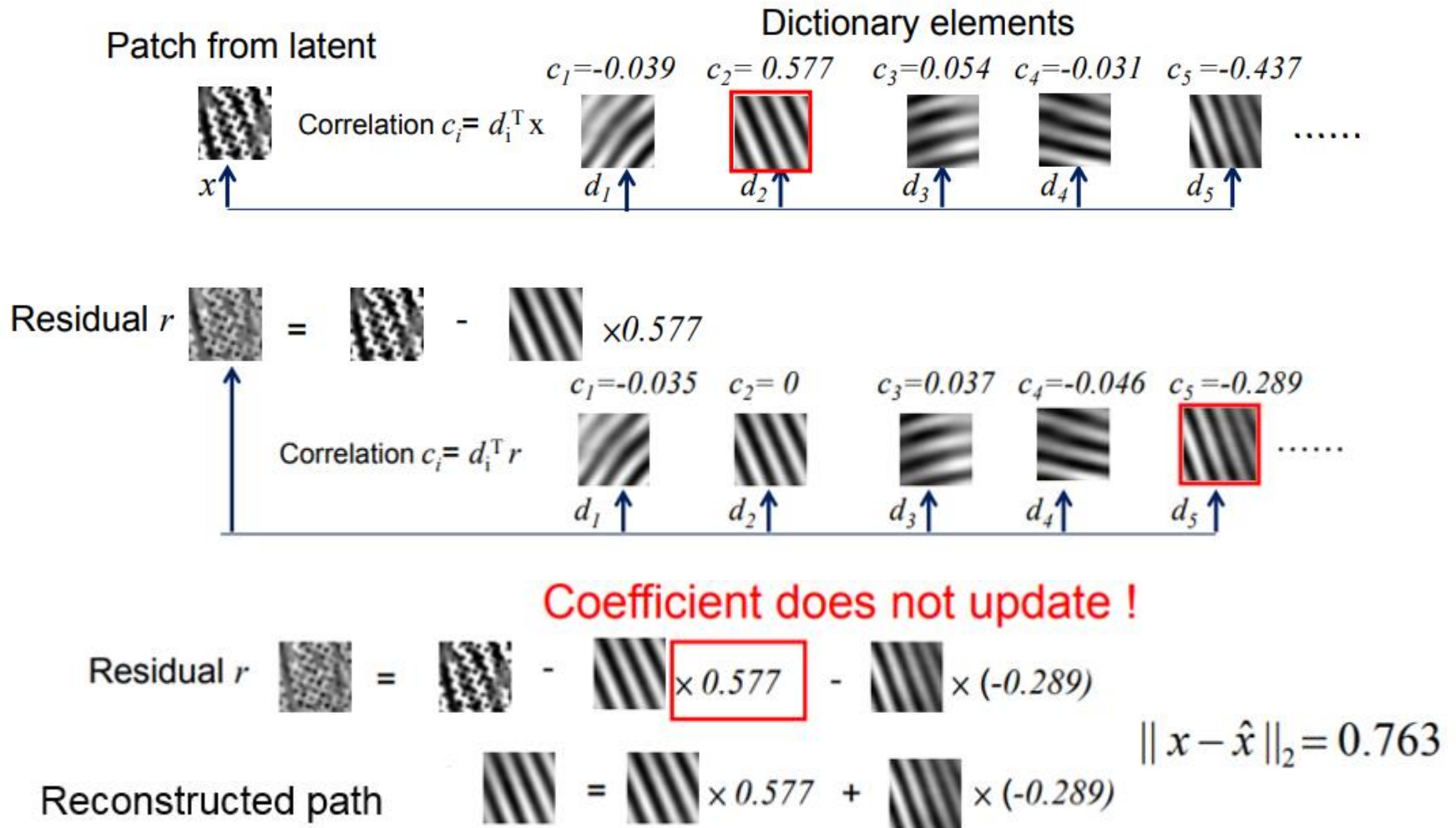
Finished:  $s$  non-zero entries in  $\mathbf{w}$

**end**

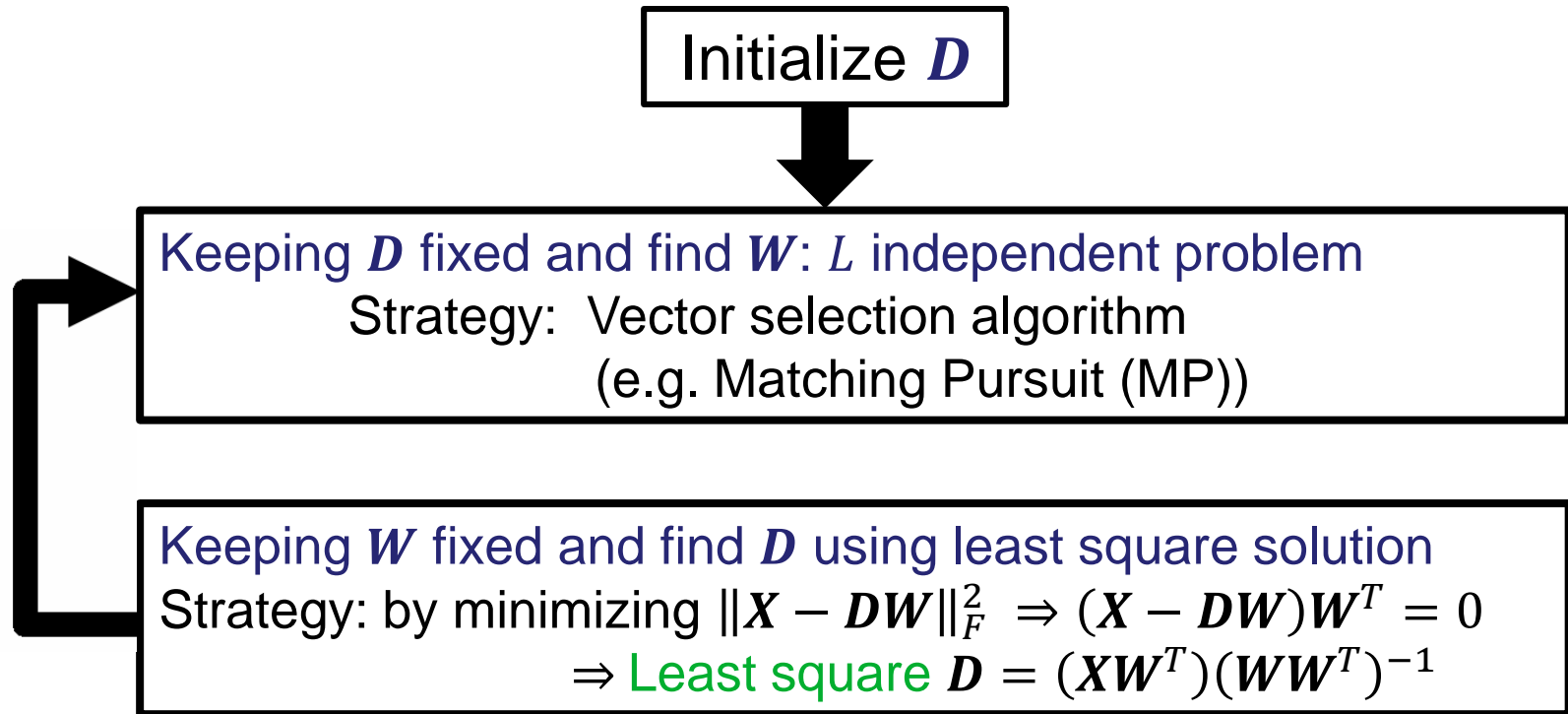
---



# Sparse Coding - MP



# Algorithm – Method of Optimal Direction (MOD)



## Issues:

- Can be slow when having large number of dictionary column because of the inversion
- Dictionary is updated before turning to re-evaluate the coefficients, which can inflict a server limitation on the training speed.



# K-SVD

$$\min_{\mathbf{D}, \mathbf{W}} \|\mathbf{X} - \mathbf{D}\mathbf{W}\|_F^2 \quad \text{subject to } \forall i, \|\mathbf{w}_i\|_0 \leq T_0$$

Initialize  $\mathbf{D}$

Keeping  $\mathbf{D}$  fixed and find  $\mathbf{W}$ :  $L$  independent problem

Strategy: Orthogonal Matching Pursuit (OMP)

$$\mathbf{D}_{\mathcal{S}^k}^T (\mathbf{D}_{\mathcal{S}^k} \mathbf{w}_{\mathcal{S}^k} - \mathbf{x}) = -\mathbf{D}_{\mathcal{S}^k}^T \mathbf{r}^k = 0$$

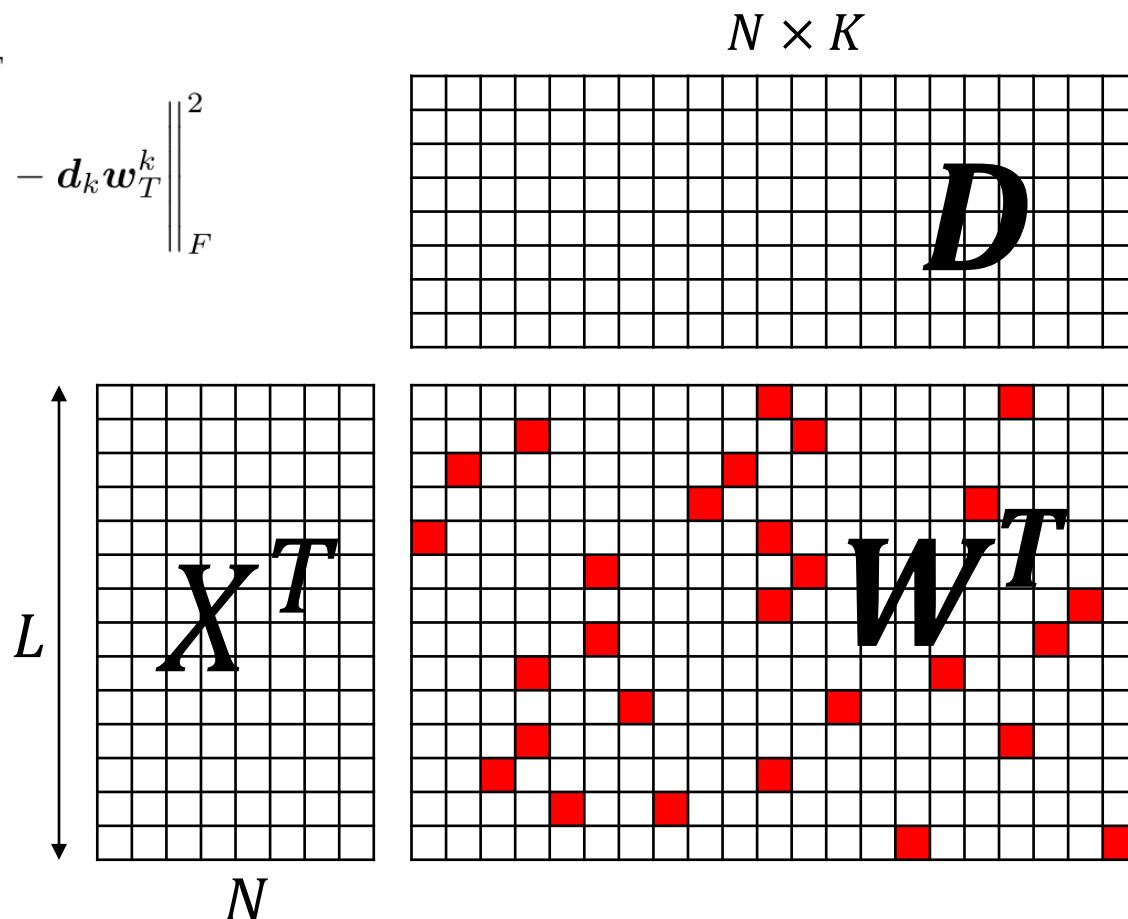
Keeping nonzero positions in  $\mathbf{W}$  fixed and find  $\mathbf{D}$  and  $\mathbf{W}$  using SVD decomposition.

Generalization of K-means clustering process



# K-SVD continued

$$\begin{aligned}\|X - DW\|_F^2 &= \left\| X - \sum_{j=1}^K d_j w_T^j \right\|_F^2 \\ &= \left\| \left( X - \sum_{j \neq k} d_j w_T^j \right) - d_k w_T^k \right\|_F^2 \\ &= \|E_k - d_k w_T^k\|_F^2\end{aligned}$$



# K-SVD continued

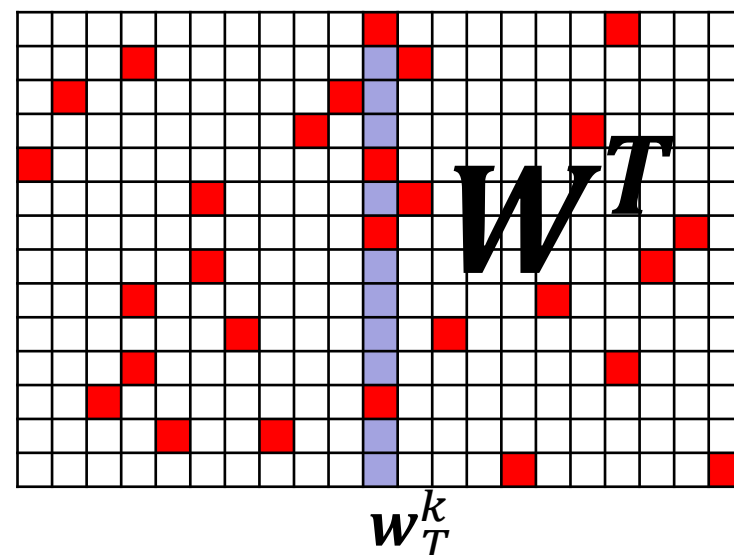
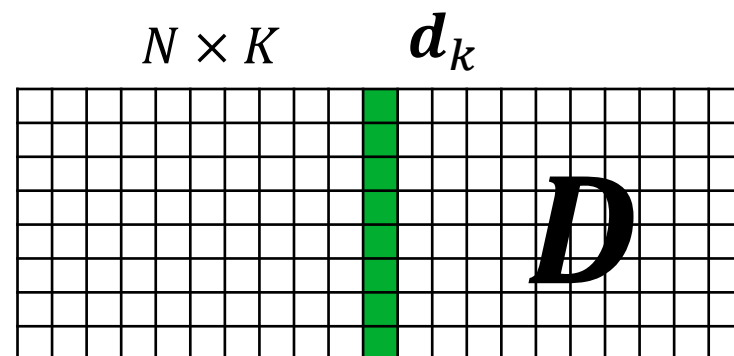
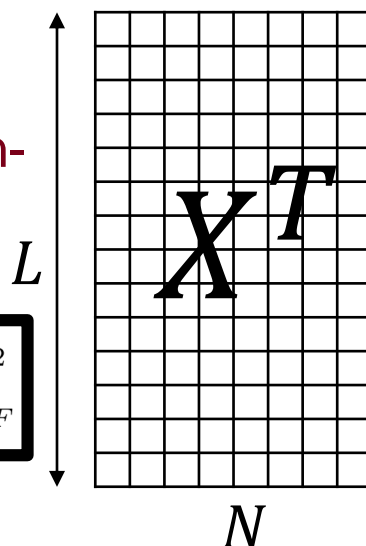
$$\begin{aligned}
 \|X - DW\|_F^2 &= \left\| X - \sum_{j=1}^K d_j w_T^j \right\|_F^2 \\
 &= \left\| \left( X - \sum_{j \neq k} d_j w_T^j \right) - d_k w_T^k \right\|_F^2 \\
 &= \|E_k - \underbrace{d_k w_T^k}_{\text{Rank-1}}\|_F^2
 \end{aligned}$$

Rank-1 approximation via SVD!

To keep cardinalities of  $W$ ,  
Define  $\Omega_k$  ( $L \times |\omega_l|$ ) such  
that it only pick out the non-  
zero entries.

$$\|E_k \Omega_K - d_k w_T^k \Omega_K\|_F^2 = \|E_k^R - d_k w_T^k\|_F^2$$

$$E_k^R = U \Delta V^T$$



# K-SVD continued

$$\begin{aligned}
 \|X - DW\|_F^2 &= \left\| X - \sum_{j=1}^K d_j w_T^j \right\|_F^2 \\
 &= \left\| \left( X - \sum_{j \neq k} d_j w_T^j \right) - d_k w_T^k \right\|_F^2 \\
 &= \|E_k - \underbrace{d_k w_T^k}_{\text{Rank-1}}\|_F^2
 \end{aligned}$$

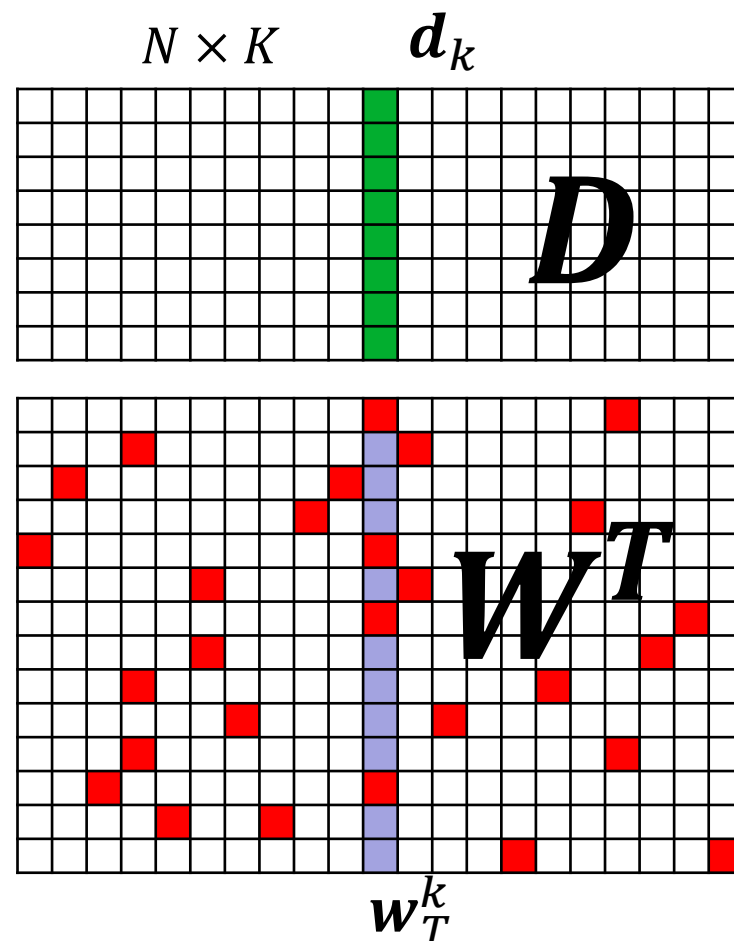
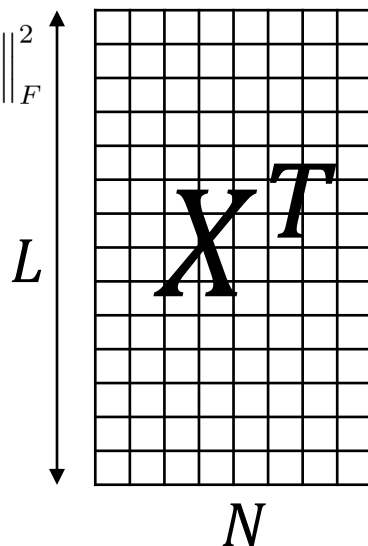
Rank-1 approximation via SVD!

$$\|E_k \Omega_K - d_k w_T^k \Omega_K\|_F^2 = \|E_k^R - d_k w_R^k\|_F^2$$

$$E_k^R = U \Delta V^T$$

$$d_k \rightarrow \tilde{d}_k = U(:, 1)$$

$$w_R^k \rightarrow V(:, 1) \Delta(1, 1)$$





# K-SVD continued

$$\begin{aligned}
 \|X - DW\|_F^2 &= \left\| X - \sum_{j=1}^K d_j w_T^j \right\|_F^2 \\
 &= \left\| \left( X - \sum_{j \neq k} d_j w_T^j \right) - d_k w_T^k \right\|_F^2 \\
 &= \|E_k - \underbrace{d_k w_T^k}_{\text{Rank-1}}\|_F^2
 \end{aligned}$$

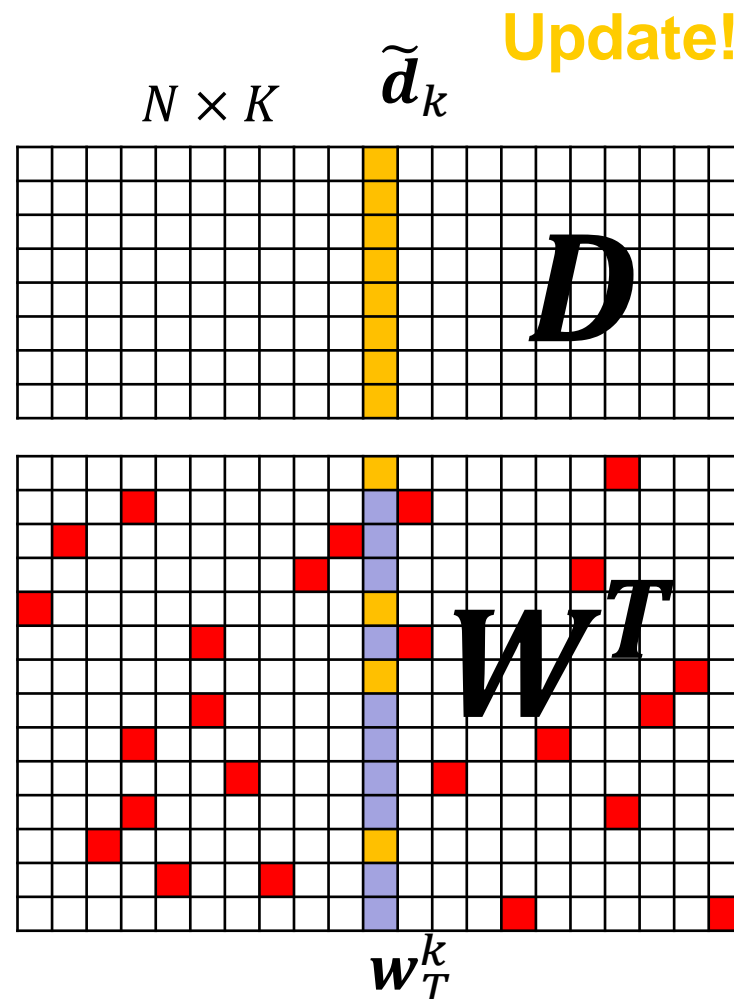
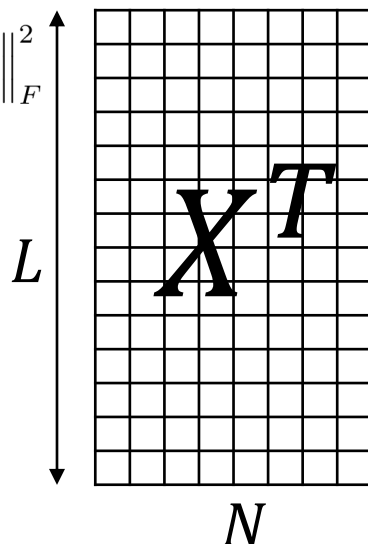
Rank-1 approximation via SVD!

$$\|E_k \Omega_K - d_k w_T^k \Omega_K\|_F^2 = \|E_k^R - d_k w_R^k\|_F^2$$

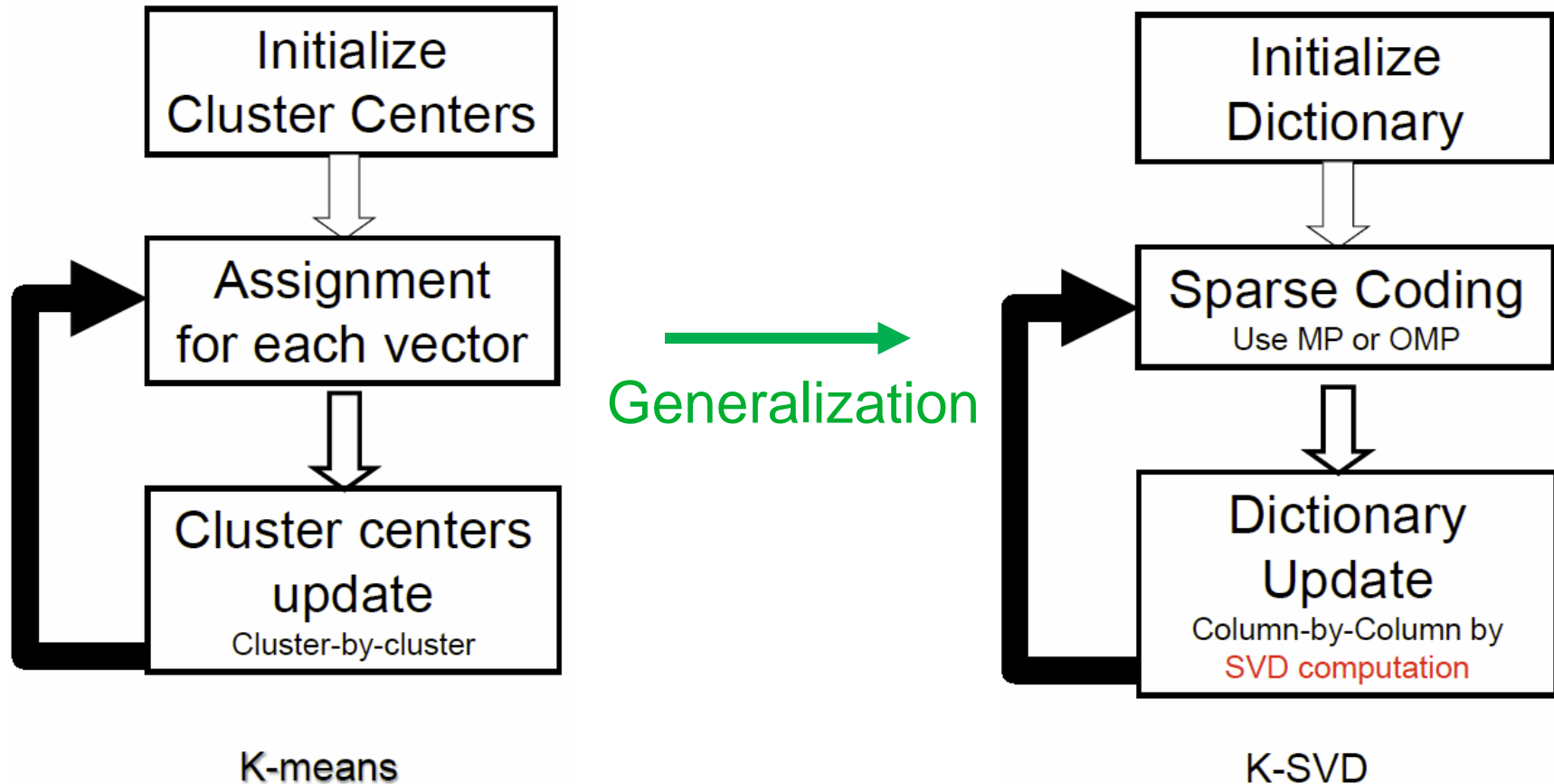
$$E_k^R = U \Delta V^T$$

$$d_k \rightarrow \tilde{d}_k = U(:, 1)$$

$$w_R^k \rightarrow V(:, 1) \Delta(1, 1)$$



# Compare K-SVD with K-means



$$\min_{C, X} \|Y - CX\|_F^2 \text{ subject to } \forall i, \mathbf{x}_i = \mathbf{e}_k \text{ for some } k \longrightarrow \min_{D, W} \|X - DW\|_F^2 \text{ subject to } \forall i, \|\mathbf{w}_i\|_0 \leq T_0$$

# Online Dictionary Learning (ODL)

## Algorithm 1 Online dictionary learning.

**Require:**  $\mathbf{x} \in \mathbb{R}^m \sim p(\mathbf{x})$  (random variable and an algorithm to draw i.i.d samples of  $p$ ),  $\lambda \in \mathbb{R}$  (regularization parameter),  $\mathbf{D}_0 \in \mathbb{R}^{m \times k}$  (initial dictionary),  $T$  (number of iterations).

- 1:  $\mathbf{A}_0 \leftarrow 0, \mathbf{B}_0 \leftarrow 0$  (reset the “past” information).
- 2: **for**  $t = 1$  to  $T$  **do**
- 3:   Draw  $\mathbf{x}_t$  from  $p(\mathbf{x})$ .
- 4:   Sparse coding: compute using LARS

$$\boldsymbol{\alpha}_t \triangleq \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^k} \frac{1}{2} \|\mathbf{x}_t - \mathbf{D}_{t-1} \boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1. \quad (8)$$

- 5:    $\mathbf{A}_t \leftarrow \mathbf{A}_{t-1} + \boldsymbol{\alpha}_t \boldsymbol{\alpha}_t^T$ .
- 6:    $\mathbf{B}_t \leftarrow \mathbf{B}_{t-1} + \mathbf{x}_t \boldsymbol{\alpha}_t^T$ .
- 7:   Compute  $\mathbf{D}_t$  using Algorithm 2, with  $\mathbf{D}_{t-1}$  as warm restart, so that

$$\begin{aligned} \mathbf{D}_t &\triangleq \arg \min_{\mathbf{D} \in \mathcal{C}} \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{x}_i - \mathbf{D} \boldsymbol{\alpha}_i\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1, \\ &= \arg \min_{\mathbf{D} \in \mathcal{C}} \frac{1}{t} \left( \frac{1}{2} \text{Tr}(\mathbf{D}^T \mathbf{D} \mathbf{A}_t) - \text{Tr}(\mathbf{D}^T \mathbf{B}_t) \right). \end{aligned} \quad (9)$$

- 8: **end for**
- 9: **Return**  $\mathbf{D}_T$  (learned dictionary).

## Algorithm 2 Dictionary Update.

**Require:**  $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_k] \in \mathbb{R}^{m \times k}$  (input dictionary),  
 $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_k] \in \mathbb{R}^{k \times k} = \sum_{i=1}^t \boldsymbol{\alpha}_i \boldsymbol{\alpha}_i^T$ ,  
 $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_k] \in \mathbb{R}^{m \times k} = \sum_{i=1}^t \mathbf{x}_i \boldsymbol{\alpha}_i^T$ .

- 1: **repeat**
- 2:   **for**  $j = 1$  to  $k$  **do**
- 3:     Update the  $j$ -th column to optimize for (9):

$$\begin{aligned} \mathbf{u}_j &\leftarrow \frac{1}{\mathbf{A}_{jj}} (\mathbf{b}_j - \mathbf{D} \mathbf{a}_j) + \mathbf{d}_j. \\ \mathbf{d}_j &\leftarrow \frac{1}{\max(\|\mathbf{u}_j\|_2, 1)} \mathbf{u}_j. \end{aligned} \quad (10)$$

- 4:   **end for**
- 5: **until convergence**
- 6: **Return**  $\mathbf{D}$  (updated dictionary).

Advantages of online learning:

1. Handle large and dynamic datasets,
2. Could be much faster than batch algorithms

# Dictionary properties

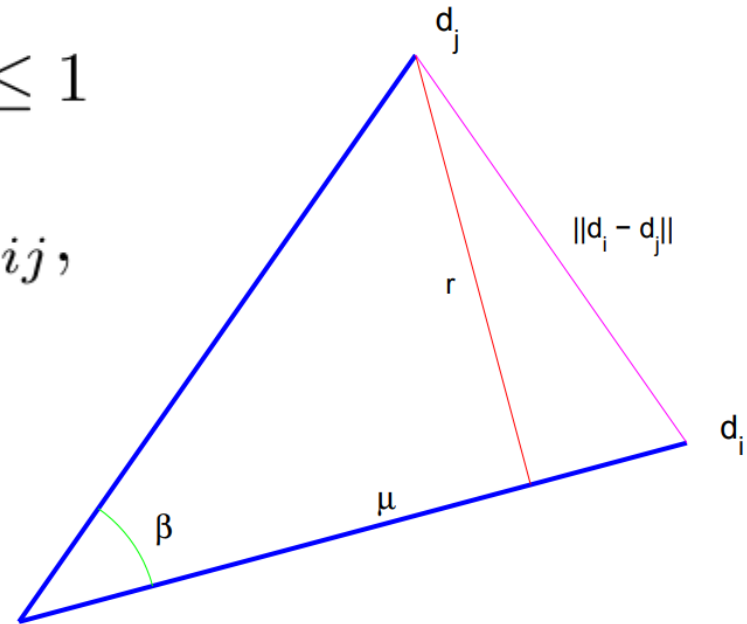
$$A \|x\|^2 \leq \sum_k |\langle x, \mathbf{d}_k \rangle|^2 \leq B \|x\|^2, \quad \text{for all } x \in H^*$$

$$\text{where } A = \|\mathbf{D}\|_2^2 = \sigma_N^2, \quad B = \|\mathbf{D}\|_2^2 = \sigma_1^2$$

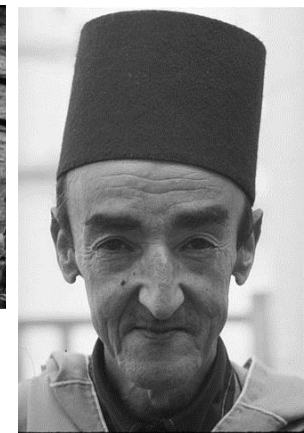
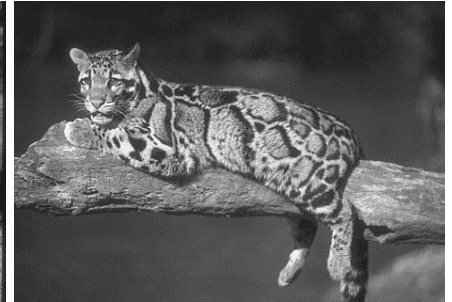
$$\mu_{ij} = |\mathbf{d}_i^T \mathbf{d}_j| = \cos \beta_{ij}, \quad 0 \leq \mu_{ij} \leq 1$$

$$\beta_{ij} = \arccos |\mathbf{d}_i^T \mathbf{d}_j| = \arccos \mu_{ij},$$

$$0 \leq \beta_{ij} \leq \pi/2$$



# Experiment: Image Compression



<https://www2.eecs.berkeley.edu/Research/Projects/CS/vision/grouping/segbench/>



UNIVERSITY OF MINNESOTA

Driven to Discover<sup>SM</sup> 25/ 31

# Experiment: Results

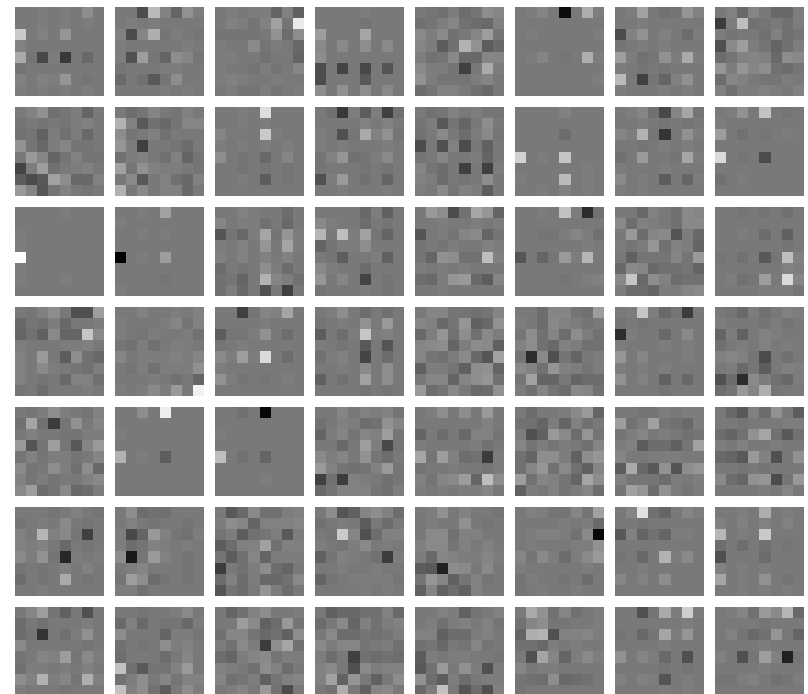
1. MOD (ILS-DLA)
2. RLS-DLA
3. K-SVD

Source Code:

<http://www.ux.uis.no/~karlsk/dle/>

The given matrix D has size 64x440

	MOD (ILS-DLA)	RLS-DLA	K-SVD
Iteration	200	200	200
tPSNR	36	36	36
Size	64 × 440	64 × 440	64 × 440
A	0.20	0.92	0.25
B	63.64	23.26	66.59
$\beta_{min}$	13.5	22.77	3.34
$\beta_{avg}$	47.09	57.58	46.29



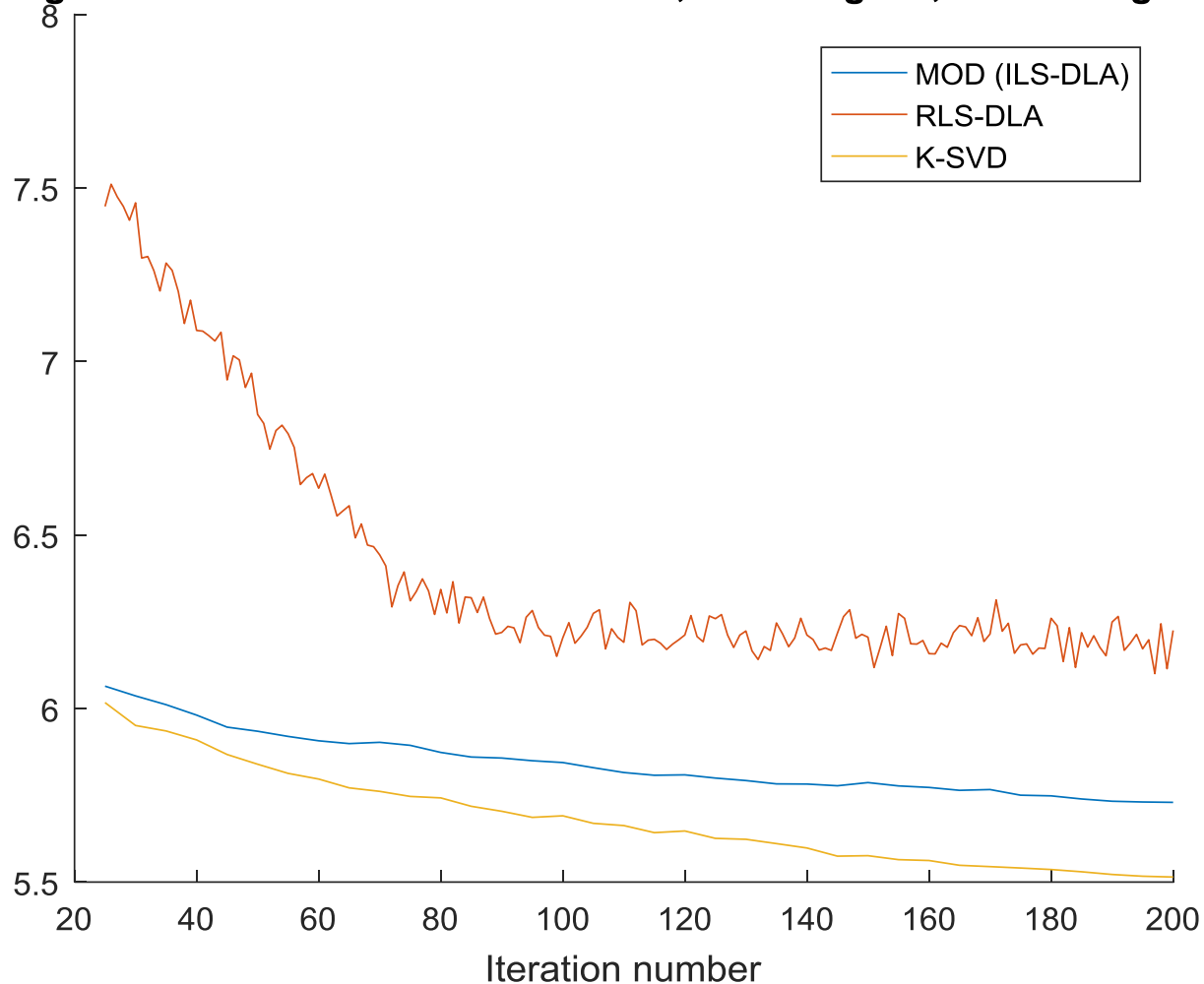
UNIVERSITY OF MINNESOTA

Driven to Discover<sup>SM</sup> 26/ 31



# Experiment: Results

Average number of non-zero coefficients, including DC, for training vectors.



# Experiment: Results

K-SVD

original Lena



512 x 512 (262144 pixels)

learned Lena



15051 nonzeros in w

$$\text{sparseness} = \frac{15051}{262144} = 0.0574$$



UNIVERSITY OF MINNESOTA

Driven to Discover<sup>SM</sup> 28/ 31

# Limitations and Challenges

- Restriction on low-dimension ( $n \leq 1000$ )
- Suitability of sparse dictionary learning model on general signal
  - “try and see” approach
- Ignores the existing dependencies between atoms
- Cannot synthesize reasonable signals from this model (e.g. natural image cannot be obtained by direct  $Dw$ )



# Thanks

## Any Questions?

### Key Reference:

K. Engan, S. O. Aase and J. Hakon Husoy, "Method of optimal directions for frame design," 1999 *IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings*

M. Aharon, M. Elad, and A. M. Bruckstein. The K-SVD: An algorithm for designing of overcomplete dictionaries for sparse representations. *IEEE Transactions on Signal Processing*, 54(11):4311-4322, November 2006.

Michael Elad. 2010. *Sparse and Redundant Representations: From Theory to Applications in Signal and Image Processing* (1st ed.). Springer Publishing Company, Incorporated.

J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research* 2010.

### Useful web:

<http://www.ux.uis.no/~karlsk/dle/>



# Backup slides-Image compression

