

# An Introduction to Sparse Coding and Dictionary Learning

Kai Cao

January 14, 2014

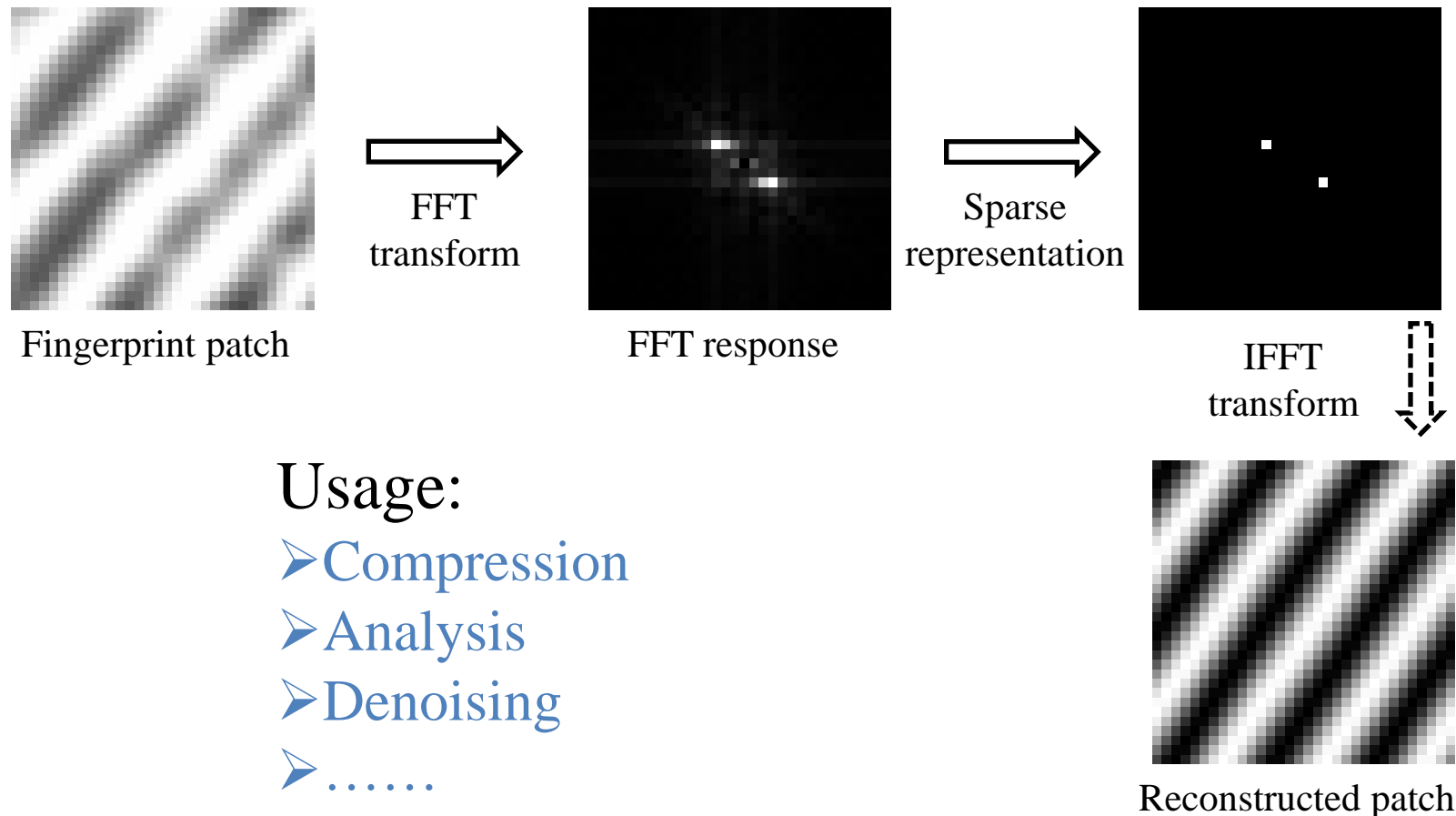
# Outline

- Introduction
- Mathematical foundation
- Sparse coding
- Dictionary learning
- Summary

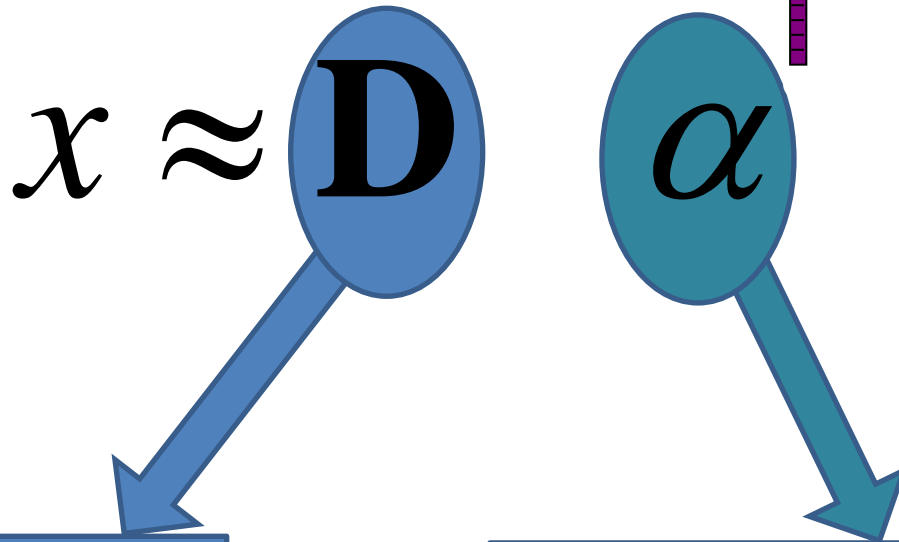
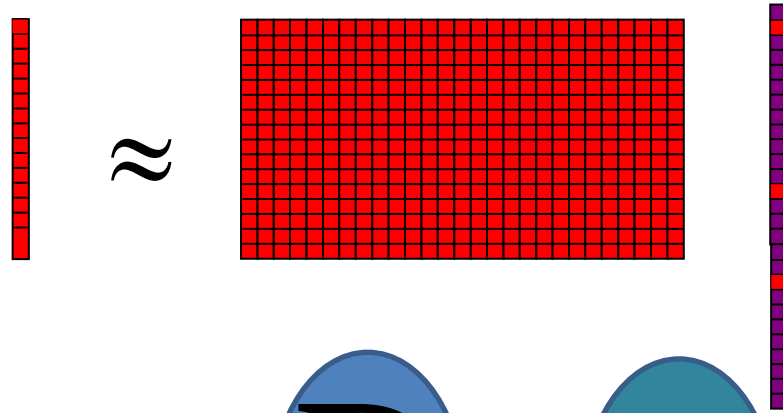
# Introduction

# What is sparsity?

- Sparsity implies many zeros in a vector or a matrix



# Sparse Representation



Dictionary Learning  
Problem

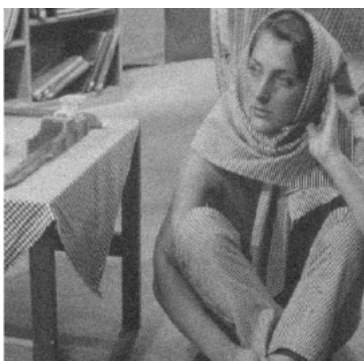
Sparse Coding  
Problem

# Application---Denoising

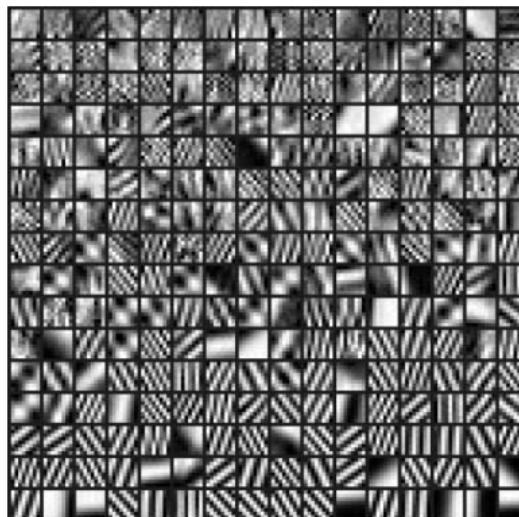
Source



Result 30.829dB



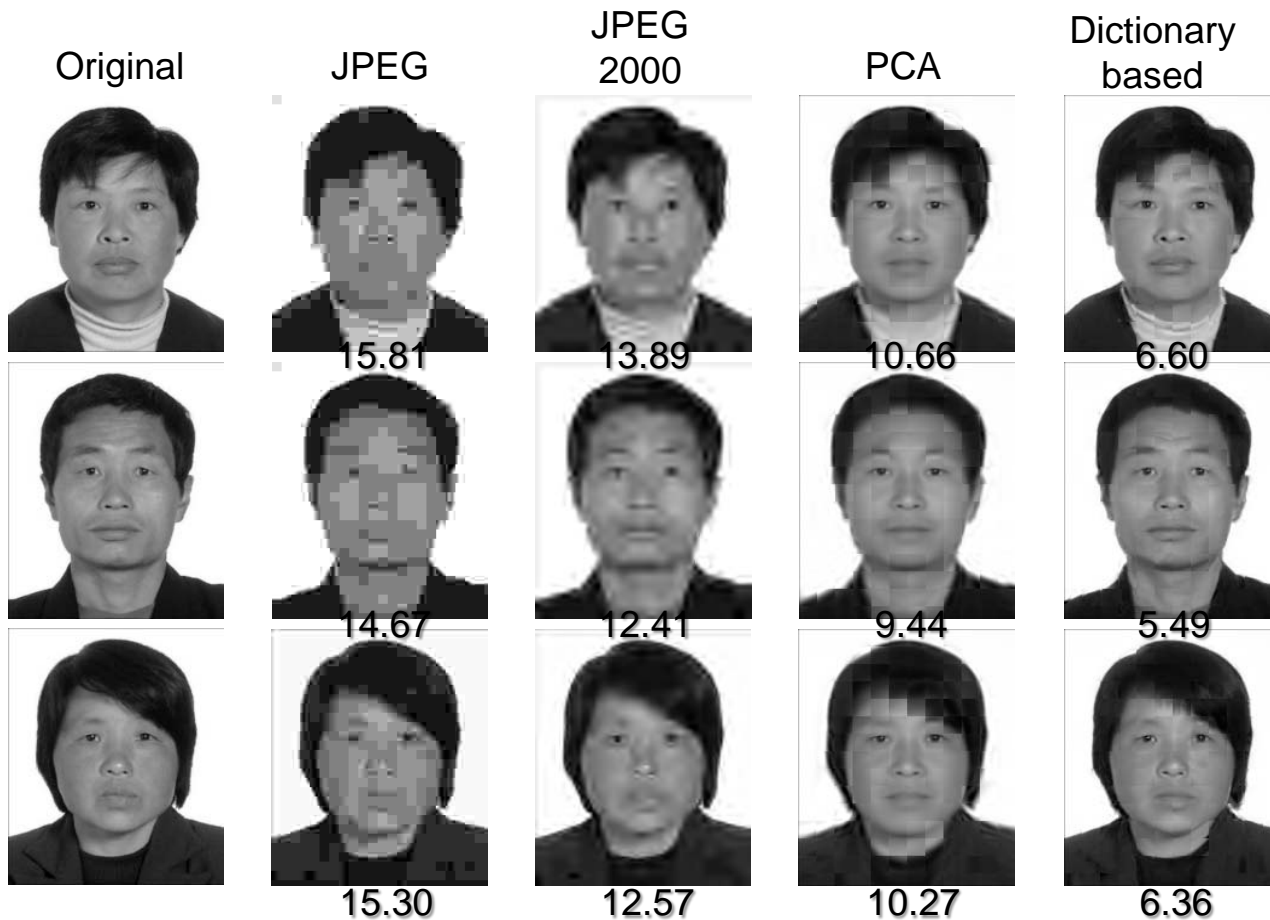
Noisy image



Dictionary

# Application---Compression

550 bytes per image



Bottom:  
RMSE values

# Mathematical foundation



# Derivatives of vectors

- First order

$$\frac{\partial a^T x}{x} = \frac{\partial x^T a}{x} = a$$

- Second order

$$\frac{\partial x^T B x}{\partial x} = (B + B^T) x$$

- Exercise

$$\min_{\alpha \in \mathbb{R}^m} \frac{1}{2} \|x - D\alpha\|_2^2 + \lambda \|\alpha\|_2^2, \quad x \in \mathbb{R}^n, D \in \mathbb{R}^{n \times m}$$



$$\alpha = (D^T D + \lambda I)^{-1} D^T x$$

# Trace of a Matrix

- Definition

$$\text{Tr}(A) = \sum_{i=1}^n a_{ii}, \quad A = (a_{ij}) \in R^{n \times n}$$

- Properties

$$\|A\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n a_{ij}^2 = \text{Tr}(A^T A),$$

$$\text{Tr}(A) = \text{Tr}(A^T),$$

$$\text{Tr}(A + B) = \text{Tr}(A) + \text{Tr}(B), \quad B \in R^{n \times n}$$

$$\text{Tr}(aA) = a\text{Tr}(A), \quad a \in R$$

$$\text{Tr}(AB) = \text{Tr}(BA), \quad B \in R^{n \times n}$$

$$\text{Tr}(ABC) = \text{Tr}(BCA) = \text{Tr}(CAB), \quad B, C \in R^{n \times n}$$

# Derivatives of traces

- First order

$$\frac{\partial}{\partial X} \text{Tr}(XA) = A^T$$

$$\frac{\partial}{\partial X} \text{Tr}(X^T A) = A$$

- Derivatives of traces

$$\frac{\partial}{\partial X} \text{Tr}(X^T XA) = XA^T + XA$$

$$\frac{\partial}{\partial X} \text{Tr}(X^T BX) = B^T X + BX$$

- Exercise

$$\min_{A \in R^{k \times m}} \|X - DA\|_F^2 + \lambda \|A\|_F^2, \quad X \in R^{n \times m}, D \in R^{n \times k}$$

$$\longrightarrow A = (D^T D + \lambda I)^{-1} D^T X$$

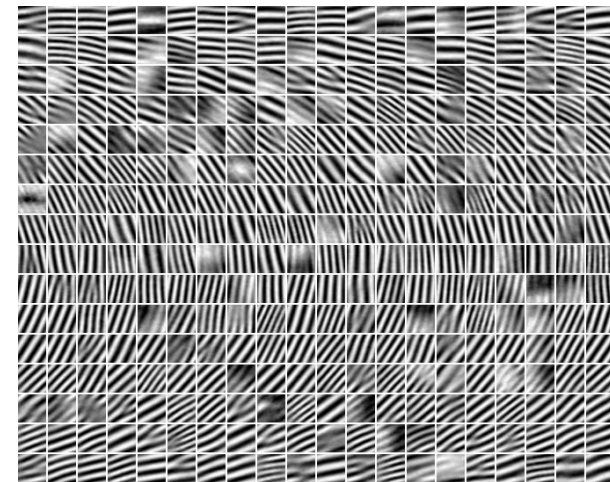
# Sparse coding

# Sparse linear model

- Let  $x \in R^n$  be a signal



- Let  $D = [d_1, d_2, \dots, d_m] \in R^{n \times m}$  be a set of normalized ( $d_i^T d_i = 1$ ) “basis vectors” (dictionary)



$D$

- Sparse representation is to find a sparse vector  $\alpha \in R^m$  such that  $x \approx D\alpha$ , where  $\alpha$  is regarded as **sparse code**

# The sparse coding model

- Objective function

$$\min_{\alpha \in \mathbb{R}^m} \underbrace{\frac{1}{2} \|x - D\alpha\|_2^2}_{\text{Data fitting term}} + \underbrace{\lambda \varphi(\alpha)}_{\text{Regularization term}}$$

- The regularization term  $\varphi$  can be

- 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

- ...

# Matching pursuit

$$\min_{\alpha \in \mathbb{R}^m} \frac{1}{2} \|x - D\alpha\|_2^2 \quad \text{s. t.} \quad \|\alpha\|_0 \leq L$$

1. Initialization:  $\alpha = 0$ , residual  $r = x$
2. while  $\|\alpha\|_0 < L$
3.     Select the element with maximum correlation with the residual

$$\hat{i} = \arg \max_{i=1, \dots, m} |d_i^T r|$$

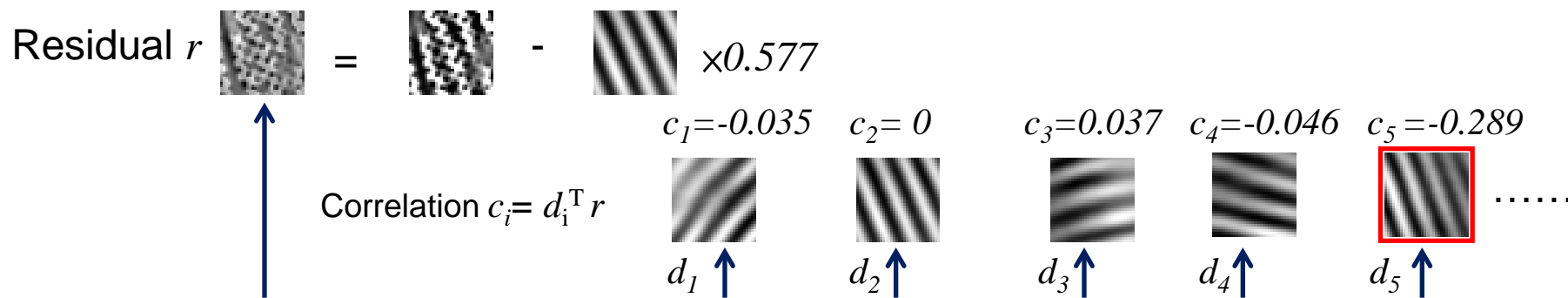
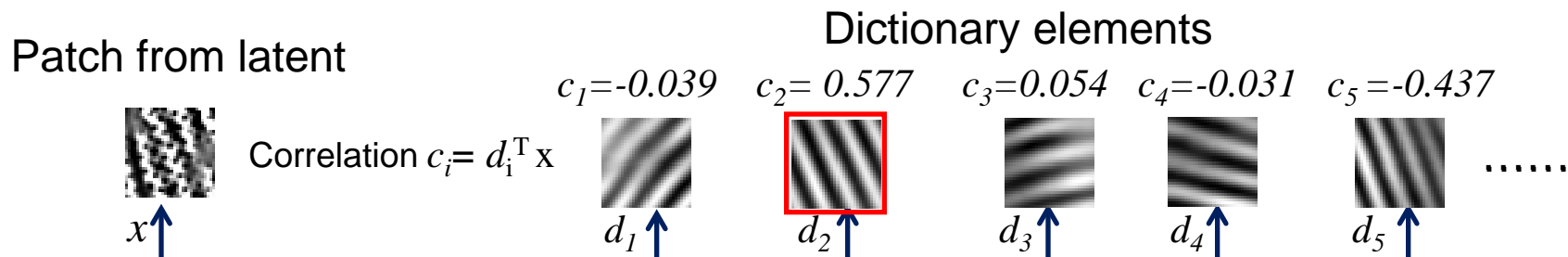
4.     Update the coefficients and residual

$$\alpha_{\hat{i}} = \alpha_{\hat{i}} + d_{\hat{i}}^T r$$

$$r = r - (d_{\hat{i}}^T r) d_{\hat{i}}$$

5. End while

# An example for matching pursuit



**Coefficient does not update !**

Residual  $r$

$r = x - d_2 \times 0.577 - d_5 \times (-0.289)$

Reconstructed patch  $\hat{x}$

$\hat{x} = d_2 \times 0.577 + d_5 \times (-0.289)$

$\|x - \hat{x}\|_2 = 0.763$



# Orthogonal matching pursuit

$$\min_{\alpha \in \mathbb{R}^m} \frac{1}{2} \|x - D\alpha\|_2^2 \quad \text{s. t.} \quad \|\alpha\|_0 \leq L$$

1. Initialization:  $\alpha = 0$ , residual  $r = x$ , active set  $\Omega = \emptyset$
2. while  $\|\alpha\|_0 < L$
3.     Select the element with maximum correlation with the residual

$$\hat{i} = \arg \max_{i=1, \dots, m} |d_i^T r|$$

4.     Update the active set, coefficients and residual

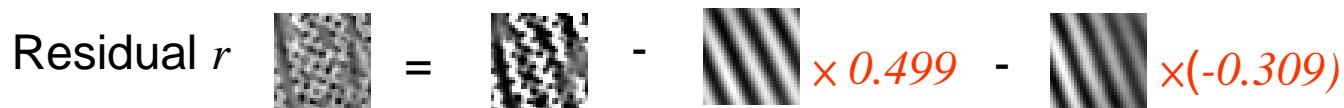
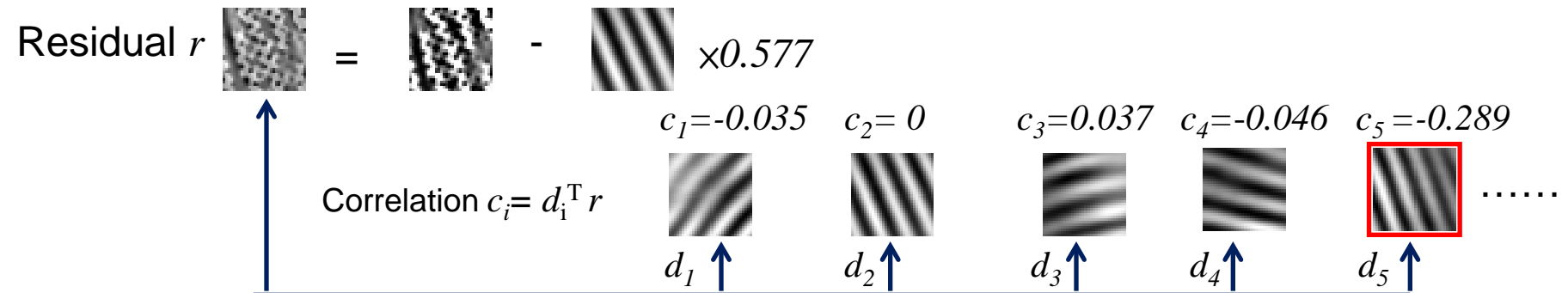
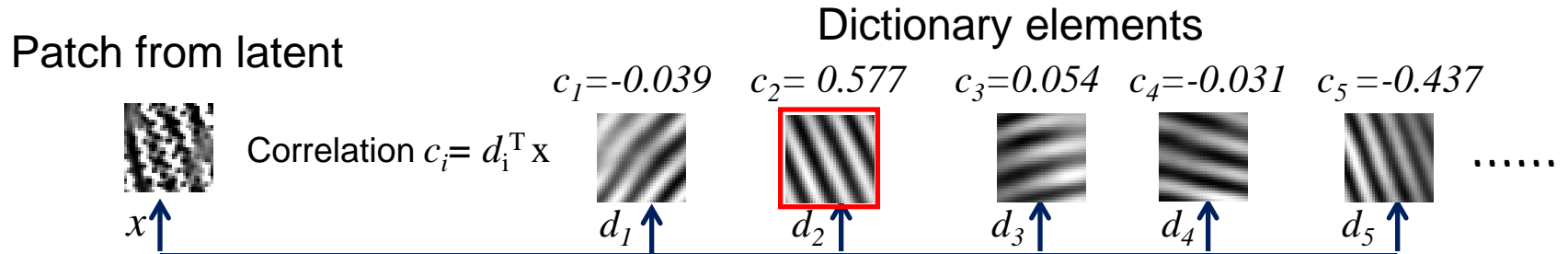
$$\Omega = \Omega \cup \hat{i}$$

$$\alpha_{\Omega} = (d_{\Omega}^T d_{\Omega})^{-1} d_{\Omega}^T r$$

$$r = x - d_{\Omega} \alpha_{\Omega}$$

5. End while

# An example for orthogonal matching pursuit



$$\|x - \hat{x}\|_2 = 0.759$$

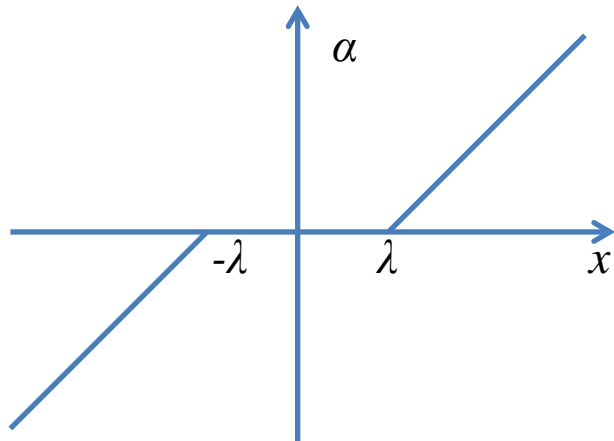


# Why does $l_1$ -norm induce sparsity?

- Analysis in 1D (comparison with  $l_2$ )

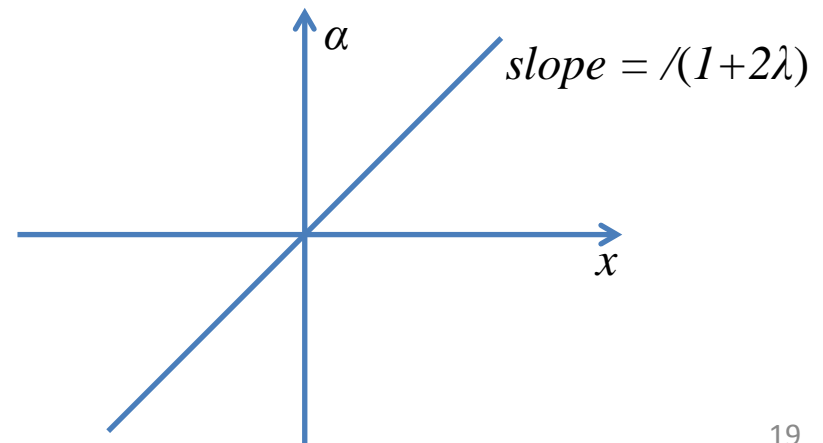
$$\min_{\alpha \in \mathbb{R}} \frac{1}{2}(x - \alpha)^2 + \lambda |\alpha|$$

$$\Rightarrow \begin{aligned} &\text{if } x \geq \lambda, \quad \alpha = x - \lambda \\ &\text{if } x \leq -\lambda, \quad \alpha = x + \lambda \\ &\text{else,} \quad \alpha = 0 \end{aligned}$$



$$\min_{\alpha \in \mathbb{R}} \frac{1}{2}(x - \alpha)^2 + \lambda \alpha^2$$

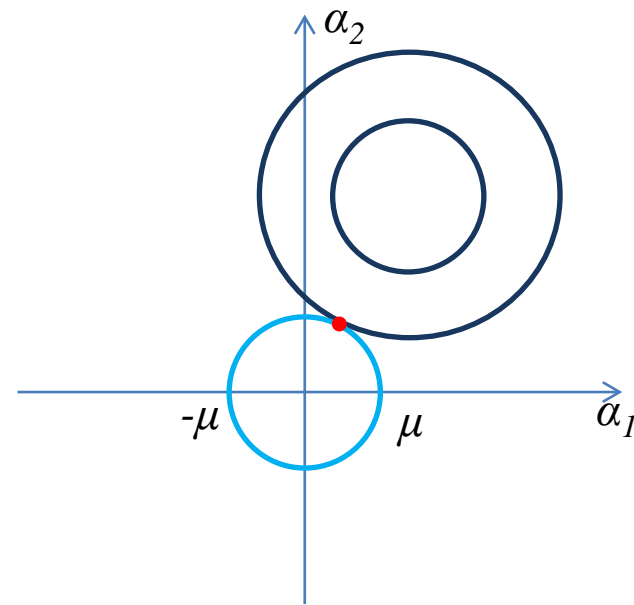
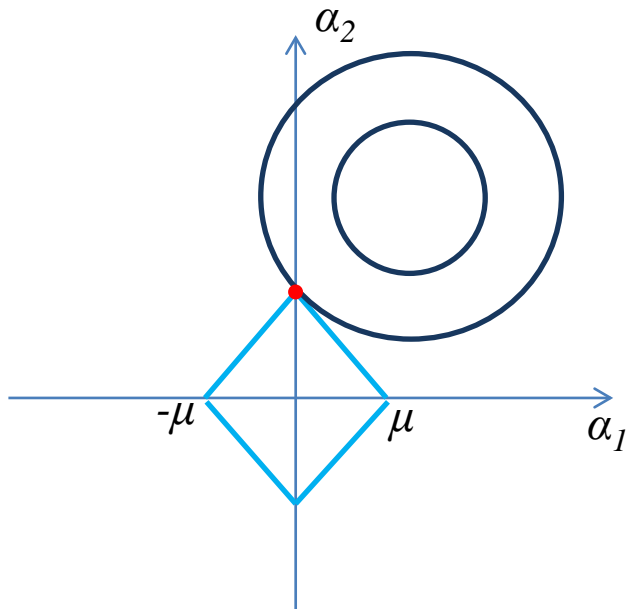
$$\Rightarrow \alpha = x/(1+2\lambda)$$



# Why does $l_1$ -norm induce sparsity?

- Analysis in 2D (comparison with  $l_2$ )

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



# Optimality condition for $l_1$ -norm regularization

$$\min_{\alpha \in \mathbb{R}^m} J(\alpha) = \frac{1}{2} \|x - D\alpha\|_2^2 + \lambda \|\alpha\|_1$$

- **Directional derivative** in the direction  $u$  at  $\alpha$

$$\nabla J(\alpha, u) = \lim_{t \rightarrow 0^+} \frac{J(\alpha + tu) - J(\alpha)}{t}$$

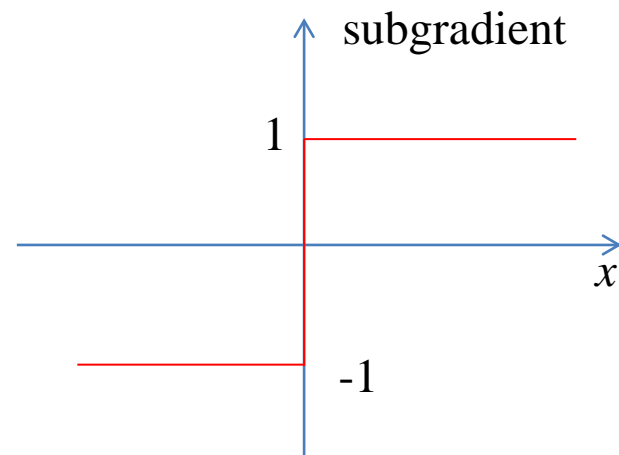
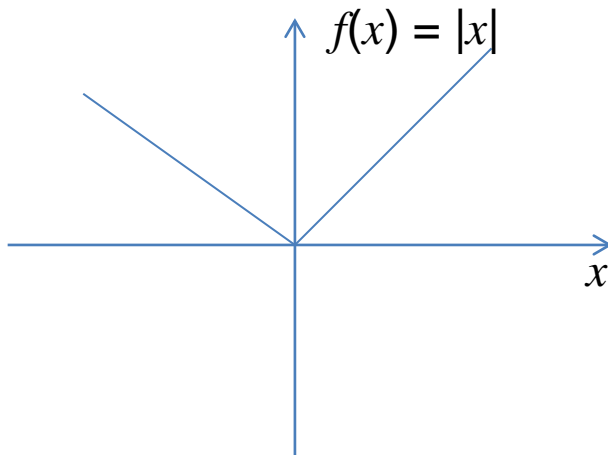
- $g$  is **subgradient** of  $J$  at  $\alpha$  if and only if

$$\forall t \in \mathbb{R}^m, J(t) \geq J(\alpha) + g^T (t - \alpha)$$

- **Proposition 1:**  $g$  is a subgradient  $\Leftrightarrow \forall u \in \mathbb{R}^m, g^T u \leq \nabla J(\alpha, u)$
- **Proposition 2:** if  $J$  is differentiable at  $\alpha$ ,  $\nabla J(\alpha, u) = \nabla J(\alpha)^T u$
- **Proposition 3:**  $\alpha$  is optimal if and only if for all  $u$ ,  $\nabla J(\alpha, u) \geq 0$

# Subgradient for $l_1$ -norm regularization

- Example:  $f(x) = |x|$



$$\nabla f(x, u) = \begin{cases} |u| & x = 0 \\ \text{sign}(x)u & x \neq 0 \end{cases}$$

# Subgradient for $l_1$ -norm regularization

$$\min_{\alpha \in \mathbb{R}^m} J(\alpha) = \frac{1}{2} \|x - D\alpha\|_2^2 + \lambda \|\alpha\|_1$$

➡  $\nabla J(\alpha, u) = -u^T D^T (x - D\alpha) + \lambda \sum_{i, a_i \neq 0} \text{sign}(a_i) u_i + \lambda \sum_{i, a_i = 0} |u_i|$

- $g$  is a subgradient at  $\alpha$  if and only if for all  $i$

$$|g_i - d_i^T (x - D\alpha)| \leq \lambda \quad \text{if} \quad a_i = 0$$

$$g_i = d_i^T (x - D\alpha) + \lambda \text{sign}(a_i) \quad \text{if} \quad a_i \neq 0$$

# First order method for convex optimization

- Differentiable objective
  - Gradient descent:  $\alpha_{t+1} = \alpha_t - \eta_t \nabla J(\alpha_t)$
  - With line search for a decent  $\eta_t$
  - Diminishing step size: e.g.,  $\eta_t = (t + t_0)^{-1}$
- Non differentiable objective
  - Subgradient decent:  $\alpha_{t+1} = \alpha_t - \eta_t g_t$ ,  $g_t$  is a subgradient
  - With line search
  - Diminishing step size



# Reformulation as quadratic program

$$\min_{\alpha \in \mathbb{R}^m} \frac{1}{2} \|x - D\alpha\|_2^2 + \lambda \|\alpha\|_1$$



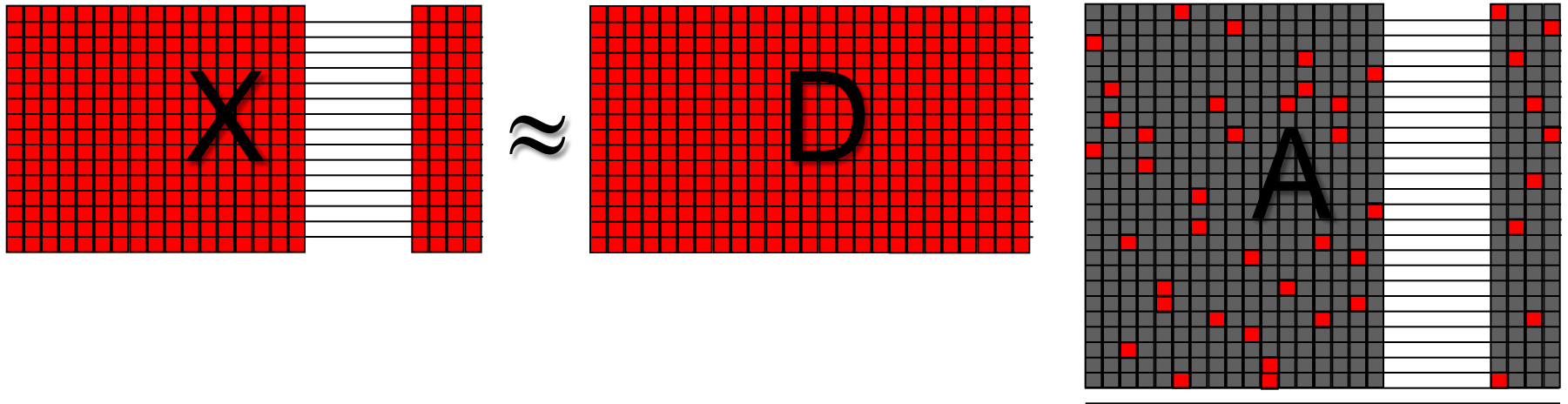
$$\min_{\alpha_+, \alpha_- \in \mathbb{R}_+^m} \frac{1}{2} \|x - D\alpha_+ + D\alpha_-\|_2^2 + \lambda(1^T \alpha_+ + 1^T \alpha_-)$$

# Dictionary Learning

# Dictionary selection

- Which D to use?
- A fixed set of basis:
  - Steerable wavelet
  - Contourlet
  - DCT Basis
  - .....
- Data adaptive dictionary – learn from data
  - K-SVD ( $l_0$ -norm)
  - On-line dictionary learning ( $l_1$ -norm)

# The objective function for K-SVD



$$\min_{D,A} \|X - DA\|_F^2$$



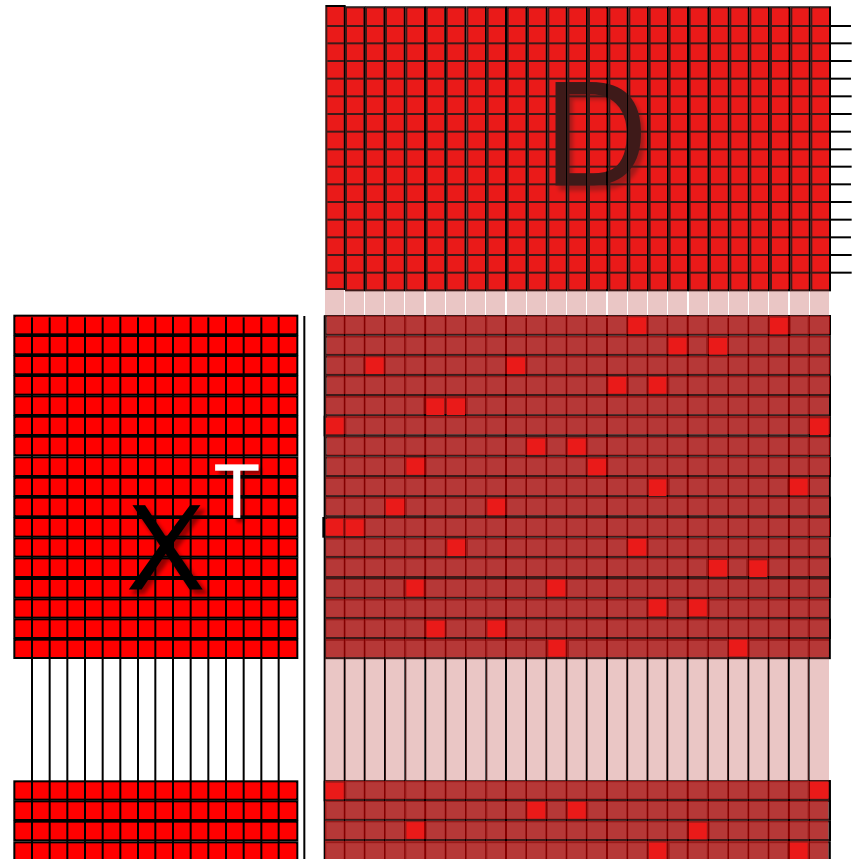
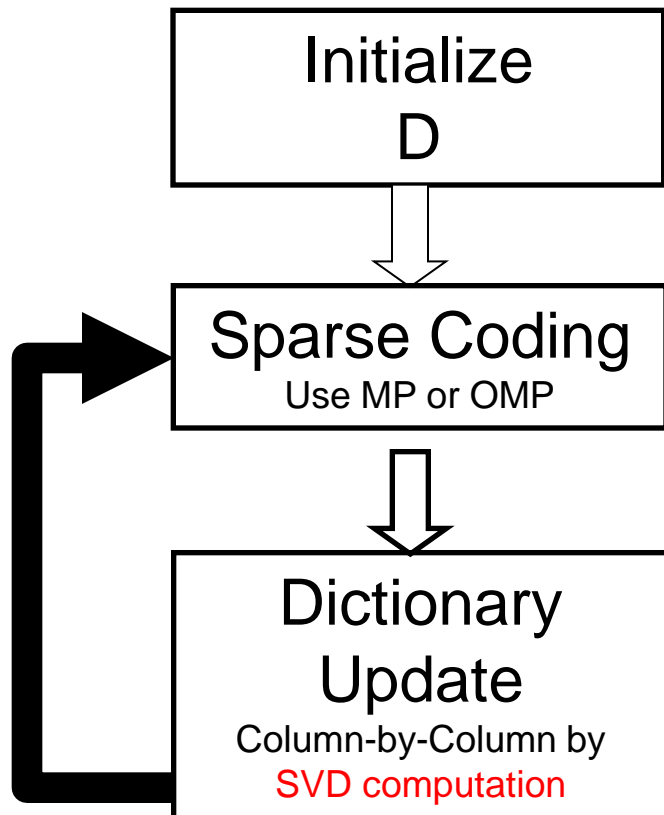
The examples are  
linear combinations  
of atoms from D

$$\forall j, \text{ s.t. } \|\alpha_j\|_0 \leq L$$



Each example has a  
sparse representation with  
no more than L atoms

# K-SVD – An Overview



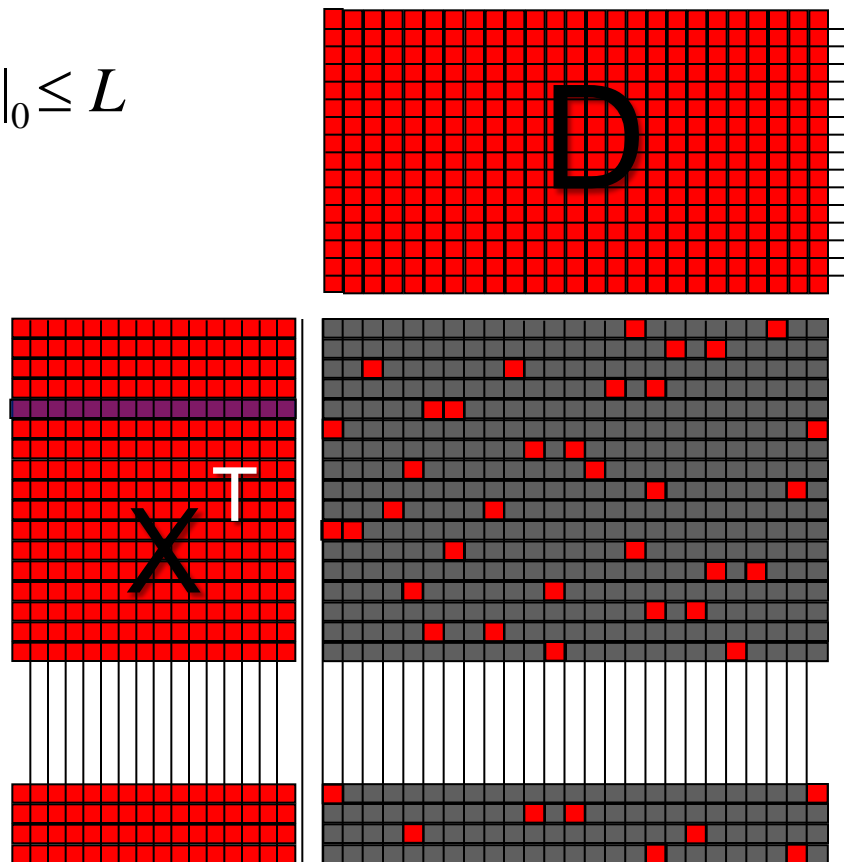
# K-SVD: Sparse Coding Stage

$$\min_A \|X - DA\|_F^2 \quad \forall j, \text{ s.t. } \|\alpha_j\|_0 \leq L$$

For the  $j^{\text{th}}$   
example  
we solve

$$\min_{\alpha} \|\mathbf{D}\alpha - x_j\|_2^2 \quad \text{s.t. } \|\alpha\|_0 \leq L$$

Ordinary Sparse Coding !



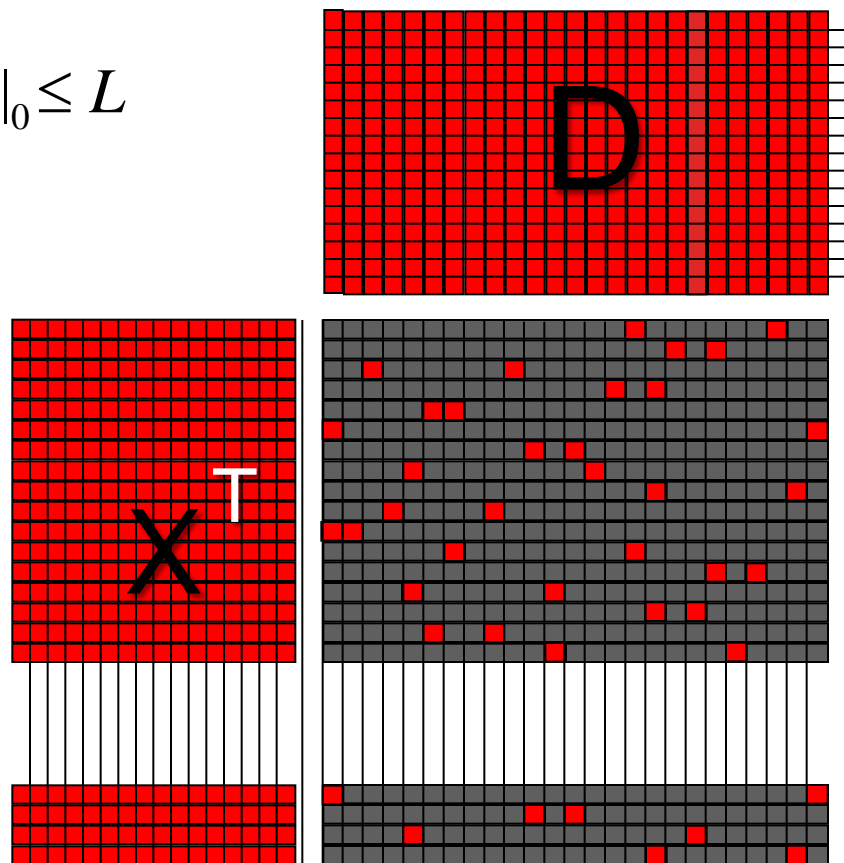
# K-SVD: Dictionary Update Stage

$$\min_D \|X - DA\|_F^2 \quad \forall j, \text{ s.t. } \|\alpha_j\|_0 \leq L$$

For the  $k^{\text{th}}$   
atom  
we solve

$$\min_{d_k} \|d_k \alpha_T^k - E_k\|_F^2$$

$$E_k = \sum_{i \neq k} d_i \alpha_T^i - X \quad (\text{the residual})$$

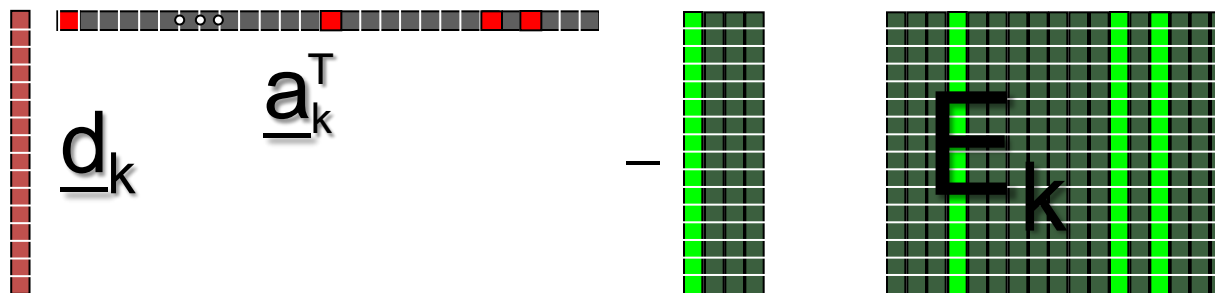


Solve with SVD

$$E_k = U \Lambda V^T \quad \Rightarrow \quad d_k = u_1$$

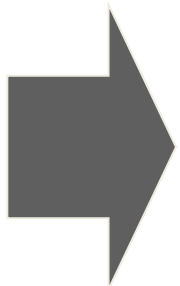
# K-SVD Dictionary Update Stage

We want to solve:

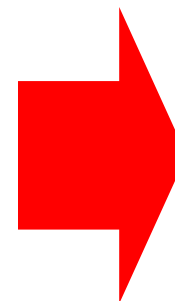


The diagram illustrates the equation  $\underline{d}_k^T \underline{a}_k^T - E_k$ . On the left,  $\underline{d}_k$  is a vertical column of red squares, and  $\underline{a}_k^T$  is a horizontal row of gray squares with several red squares. In the center is a minus sign. To the right of the minus sign is a vertical column of green squares, and further right is a square grid of green squares labeled  $E_k$ .

Only some of  
the examples  
use column  $\underline{d}_k$ !



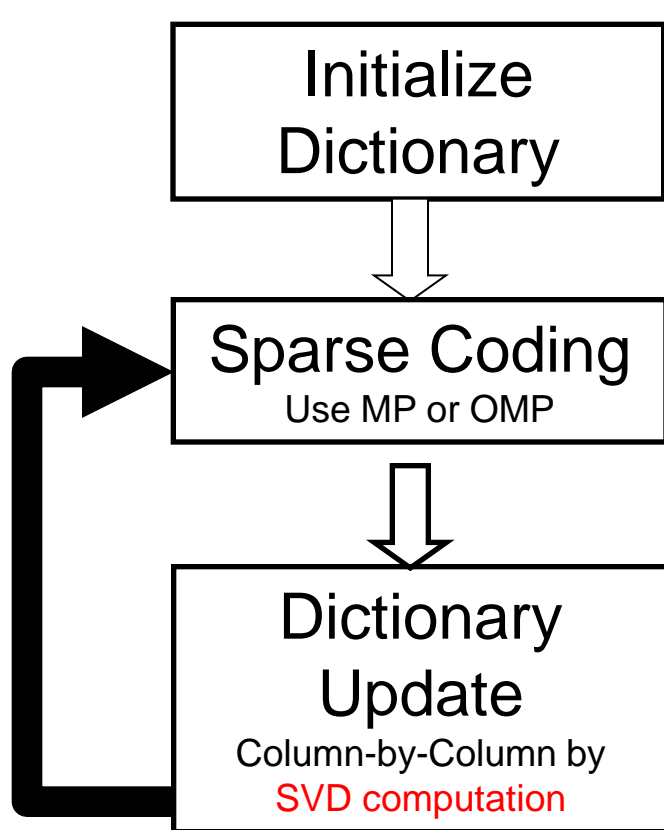
When updating  $\underline{a}_k$ ,  
only recompute  
the coefficients  
corresponding to  
those examples



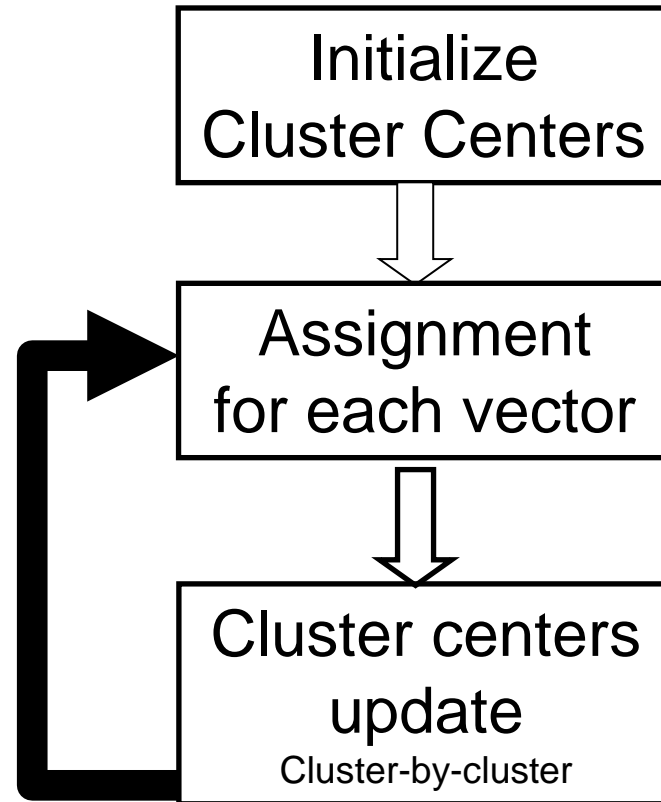
Solve with  
SVD!



# Compare K-SVD with K-means



K-SVD



K-means

# dictionary learning with $l_1$ -norm regularization

- Objective function for  $l_1$ -norm regularization

$$\min_D \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|x_i - D\alpha_i\|_2^2 + \lambda \|\alpha_i\|_1$$

where

$$\alpha_i \triangleq \arg \min_{\alpha \in R^m} \frac{1}{2} \|x_i - D\alpha\|_2^2 + \lambda \|\alpha\|_1$$

- Advantages of online learning:
  - Handle large and dynamic datasets,
  - Could be much faster than batch algorithms.

# dictionary learning with $l_1$ -norm regularization

$$\begin{aligned} F_t(D) &= \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|x_i - D\alpha_i\|_2^2 + \lambda \|\alpha_i\|_1 \\ &= \frac{1}{t} \left( \frac{1}{2} \text{Tr}(D^T D A_t) - \text{Tr}(D^T B_t) \right) + \lambda \sum_{i=1}^t \|\alpha_i\|_1 \end{aligned}$$

where

$$A_t = \sum_{i=1}^t \alpha_i \alpha_i^T, \quad B_t = \sum_{i=1}^t x_i \alpha_i^T$$



$$\frac{\partial F_t(D)}{\partial D} = \frac{1}{t} (D A_t - B_t)$$

For a new  $x_{t+1}$ ,  $A_{t+1} = A_t + \alpha_{t+1} \alpha_{t+1}^T$ ,  $B_{t+1} = B_t + x_{t+1} \alpha_{t+1}^T$

# On-line dictionary learning

- 1) Initialization:  $D_0 \in R^{n \times m}$ ;  $A_0=0$ ;  $B_0=0$ ;
- 2) For  $t=1, \dots, T$
- 3) Draw  $x_t$  from the training data set
- 4) Get sparse code

$$\alpha_t = \arg \min_{\alpha \in R^m} \frac{1}{2} \|x_t - D_{t-1} \alpha\|_2^2 + \lambda \|\alpha\|_1$$

- 5) Aggregate sufficient statistics

$$A_t = A_{t-1} + \alpha_t \alpha_t^T, \quad B_t = B_{t-1} + x_t \alpha_t^T,$$

- 6) Dictionary update

$$D_t = D_{t-1} - \rho \frac{\partial F_t(D)}{\partial D}$$

- 7) End for

# Toolbox - SPAMS

- SPArse Modeling Software:
  - Sparse coding
    - $l_0$ -norm regularization
    - $l_1$ -norm regularization
    - .....
  - Dictionary learning
    - K-SVD
    - Online dictionary learning
    - .....
- C++ implemented with Matlab interface
- <http://spams-devel.gforge.inria.fr/>

# Summary

---

- Sparsity and sparse representation
- Sparse coding with  $l_0$ - and  $l_1$ -norm regularization
  - Orthogonal matching pursuit/matching pursuit
  - Subgradient and optimal condition
- Dictionary learning with  $l_0$ - and  $l_1$ -norm regularization
  - K-SVD
  - Online dictionary learning
- Try to use it !!

# References

- T. T. Cai, Lie Wang, Orthogonal Matching Pursuit for Sparse Signal Recovery With Noise, *IEEE Transactions on Information Theory*, 57(7): 4680-4688, 2011
- Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Annals of statistics*, 32(2):407–499, 2004.
- 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.
- J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online dictionary learning for sparse coding. *In Proceedings of the International Conference on Machine Learning (ICML)*, 2009a.

Thank you for listening