# INTRODUCTORY APPLIED MACHINE LEARNING

Yan-Fu Kuo

Dept. of Bio-industrial Mechatronics Engineering

National Taiwan University

Today:

- Sparse coding

# Outline

- Goal of the lecture

- Denoising by sparse representations

- Sparsity and overcompleteness

- Theoretical and numerical foundations

- Dictionary learning and K-SVD algorithm
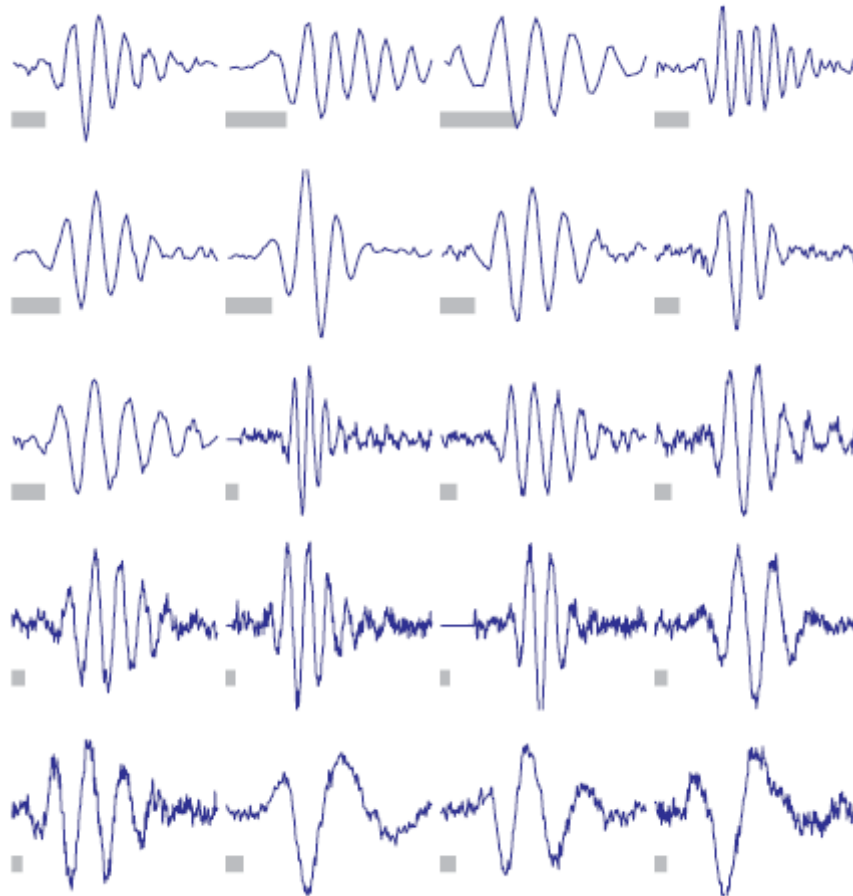
- Putting it all together

# Goals

- After this, you should be able to:

  - Understand the principles of sparse coding
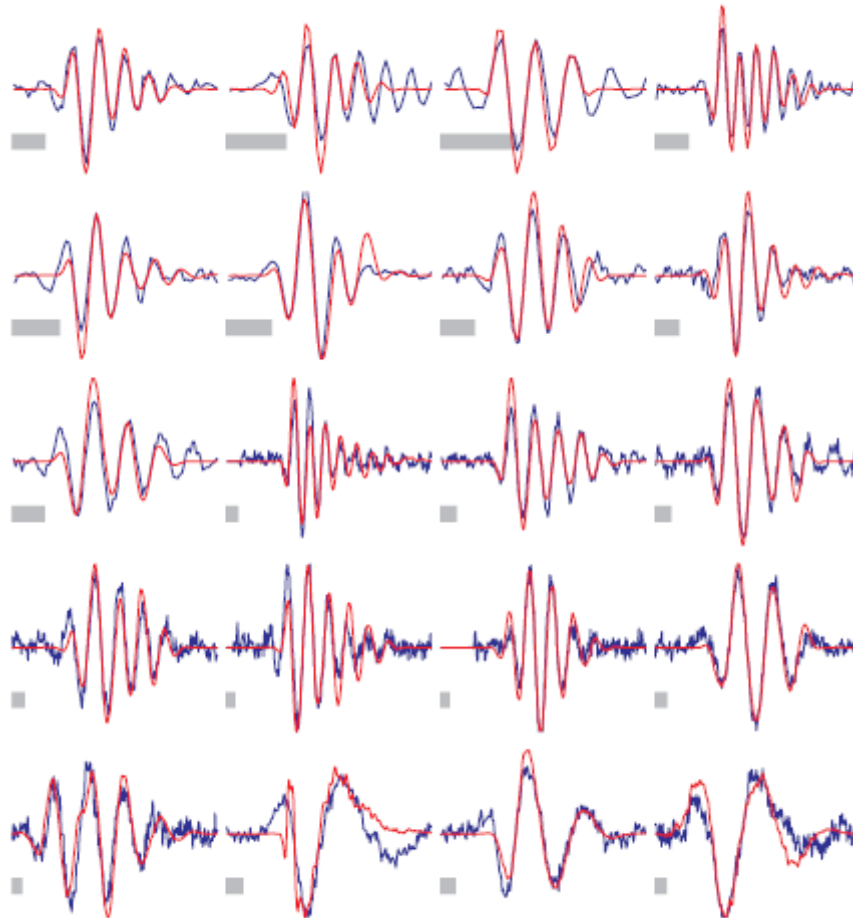
  - Apply sparse coding methods

# Noise Removal and Image Scaling Problems

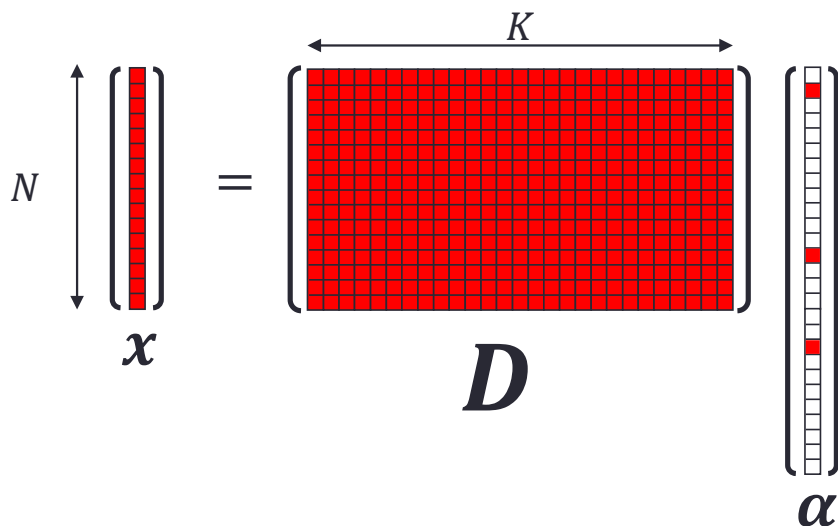# Audio Signal with Noise

# Noise Reduction

# Denoising by Energy Minimization

- Let $y \in \Re^N$ be measurements with noise, and let $x \in \Re^N$ be true signal (which is unknown) to be recovered

- Assume that $x$ can be calculated from a dictionary $D$, i.e., $x = D\alpha$

- Denoising is to minimize an energy function:

$$f(x) = \boxed{\frac{1}{2}\|x - y\|_2^2} + \boxed{\Pr(x) \;\; \lambda\|\alpha\|_0^0}$$

**Relation to measurements**     **Prior or regularization**

- For "sparse" representation, $\Pr(x) = \lambda\|\alpha\|_0^0$ for $x = D\alpha$
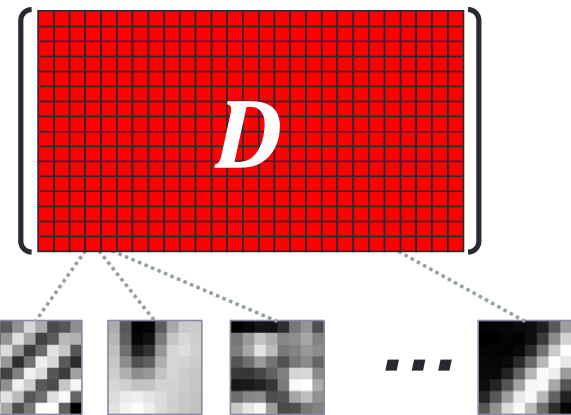
# True Signal in Sparse Representation



- The true signal $x$ is assumed to be a linear combination of some prototype signal (or atoms) from a dictionary $D \in \Re^{N \times K}$

- The coefficient vector $\alpha \in \Re^K$ is a vector with few (say $L$) non-zeros entries
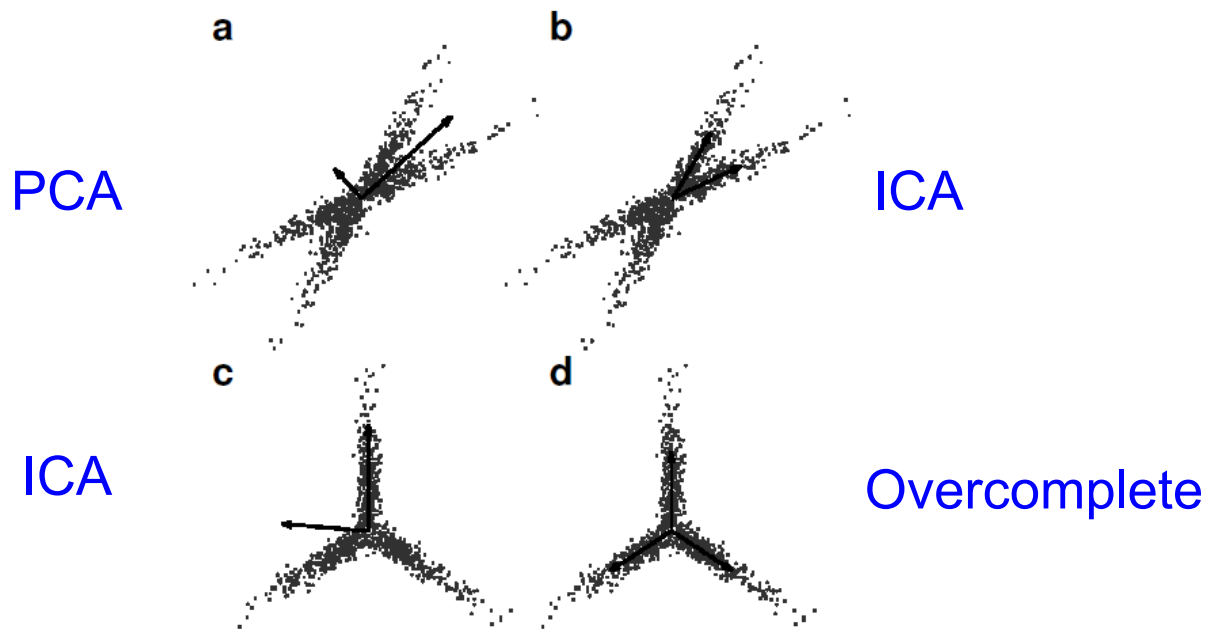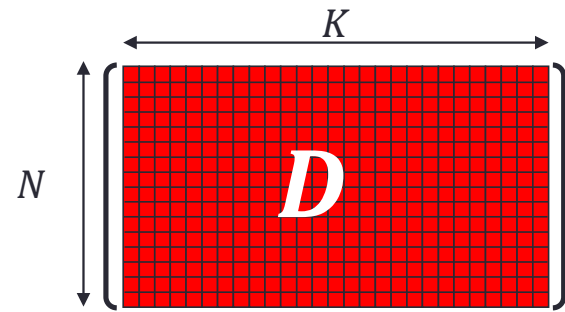
# What is A Dictionary *D*?

- A dictionary is a signal model that contains many basis (atoms)



- An image patch is a combination of these atoms



$$\approx 0.8 \times \qquad + 0.3 \times \qquad + 0.5 \times$$

# Why Overcomplete?

- The dictionary $D$ is usually overcomplete, i.e., $K > N$

- Why overcomplete?

$K$

$N$ $D$

PCA

ICA

ICA

Overcomplete

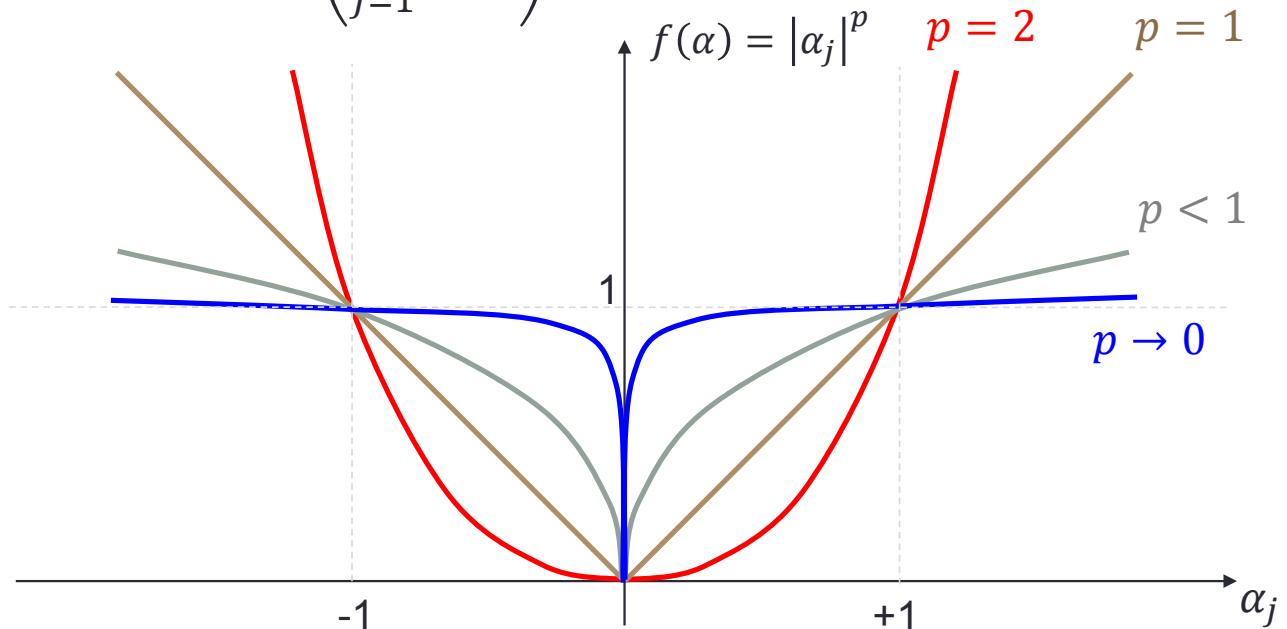# Why Sparse Representation?

- Simple: every signal $x$ is built as a linear combination of a <u>few</u> atoms from the dictionary $D$

- Rich: the obtained signals are a union of many low-dimensional spaces

- Effective: recent works adopt this model and successfully deploy it to applications

- Empirically established: neurological studies show similarity between this model and early vision processes

# How to Measure Sparsity? (Why $\|\boldsymbol{\alpha}\|_0^0$?)

- Need a measure of sparsity of $\boldsymbol{\alpha}$, i.e., $\|\boldsymbol{\alpha}\|_p^p = \#\{j: \alpha_j \neq 0\}$

- Note that $\|\boldsymbol{\alpha}\|_p^p \equiv \left( \sum_{j=1}^{K} |\alpha_j|^p \right)^{1/p}$



$f(\alpha) = |\alpha_j|^p$    $p = 2$    $p = 1$

$p < 1$

$p \to 0$

1

-1      +1      $\alpha_j$

# The Sparse Coding Problem

- Assume $D$ and $y$ are known

- What should $\alpha$ be?

$$\hat{\alpha} = \underset{\alpha}{\mathrm{argmin}} \ \frac{1}{2} \| \, x - y \, \|_2^2 \quad s.t. \ \|\alpha\|_0^0 \leq L$$

known

$$D\alpha - y = \qquad -$$

- Need to constrain number of non-zero entries in $\alpha$
- Since only a few ($L$ out of $K$) atoms can be merged to form the true signal, the noise cannot be fitted well

# Issues with the Formulation

- Numerical problem: how should we solve or approximate the solution of the problem?

$$\min_{\boldsymbol{\alpha}} \|\boldsymbol{D\alpha} - \boldsymbol{y}\|_2^2 \quad s.t. \quad \|\boldsymbol{\alpha}\|_0^0 \leq L$$

$$\min_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}\|_0^0 \quad s.t. \quad \|\boldsymbol{D\alpha} - \boldsymbol{y}\|_2^2 \leq \varepsilon^2$$

$$\min_{\boldsymbol{\alpha}} \|\boldsymbol{D\alpha} - \boldsymbol{y}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_0^0$$
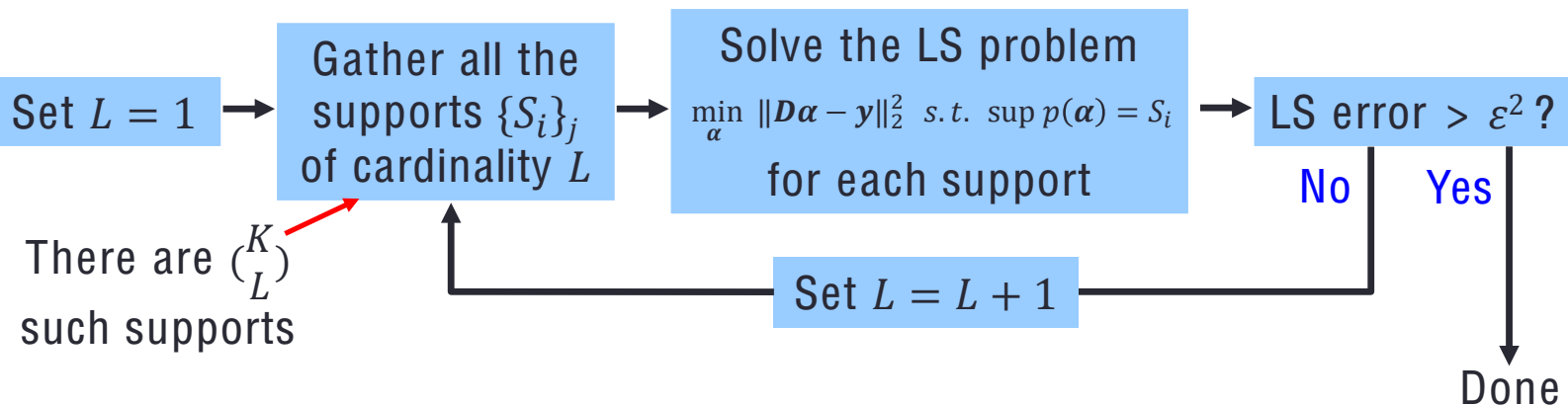
- Theoretical problem: is there a unique sparse representation?

- Practical problem: what dictionary $\boldsymbol{D}$ should we use, such that all this leads to effective denoising?

# Solving the Problem

$$\min_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}\|_0^0 \quad s.t. \quad \|\boldsymbol{D\alpha} - \boldsymbol{y}\|_2^2 \leq \varepsilon^2$$

This is a combinatorial problem, proven to be NP-Hard!

- Recipe for solving this problem:

Set $L = 1$ → Gather all the supports $\{S_i\}_j$ of cardinality $L$ → Solve the LS problem $\min_{\boldsymbol{\alpha}} \|\boldsymbol{D\alpha} - \boldsymbol{y}\|_2^2 \quad s.t. \quad \sup p(\boldsymbol{\alpha}) = S_i$ for each support → LS error > $\varepsilon^2$ ?

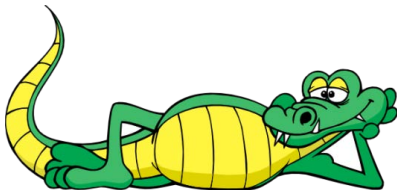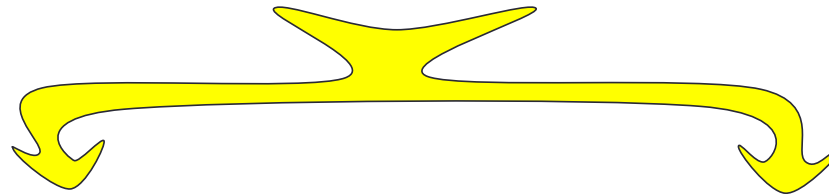There are $\binom{K}{L}$ such supports

Set $L = L + 1$

No    Yes

Done

Assume: $K = 2000$, $L = 10$, 1 nano-sec per each LS

We shall need $\sim 8 \times 10^9$ years to solve this problem !!!!!

# Approximation

$$\min_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}\|_0^0 \quad s.t. \quad \|\boldsymbol{D\alpha} - \boldsymbol{y}\|_2^2 \leq \varepsilon^2$$



**Relaxation methods**

Smooth the $L_0$ and use continuous optimization techniques

**Greedy methods**

Build the solution one non-zero element at a time

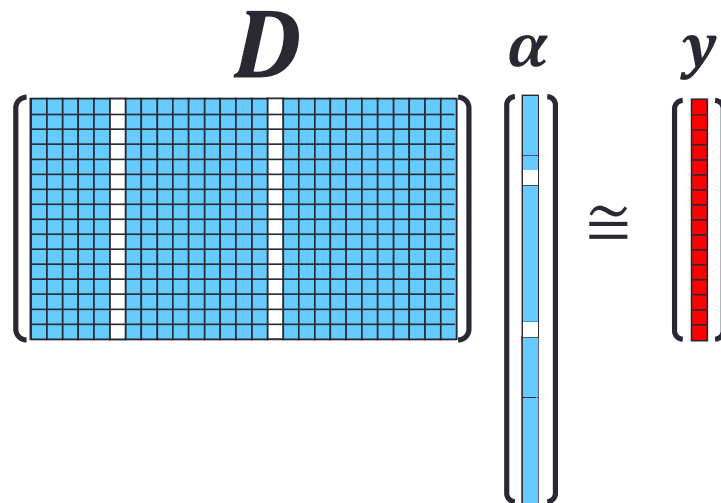# Relaxation Approach

**Solving this instead**

$$\min_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}\|_0^0 \quad s.t. \quad \|\boldsymbol{D\alpha} - \boldsymbol{y}\|_2^2 \leq \varepsilon^2$$

$$\min_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}\|_1^1 \quad s.t. \quad \|\boldsymbol{D\alpha} - \boldsymbol{y}\|_2^2 \leq \varepsilon^2$$

- Also known as basis pursuit

- The newly defined problem is convex and can be solved using quadratic programming techniques

- Very efficient solvers can be deployed

# Greedy Approach

- Also known as matching pursuit (MP)

- Finds one atom at a time

  - First step: find the one atom that best matches the signal

  - Next steps: given the previously found atoms, find the next one to <u>best fit</u> the residual



- The algorithm stops when $\|\boldsymbol{D\alpha} - \boldsymbol{y}\|_2^2 \leq \varepsilon^2$ is satisfied

# What Should the Dictionary *D* Be?

$$\min_{\alpha} \|\alpha\|_1^1 \quad s.t. \quad \|D\alpha - y\|_2^2 \leq \varepsilon^2$$

Assumption: good-behaved images have a sparse representation

*D* should be chosen such that it sparsifies the representations

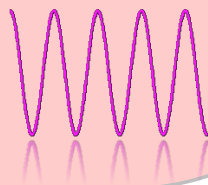Choosing *D* from a known set of transforms (Fourier, wavelet, consine, etc.)

Building *D* by training it, based on learning from image examples

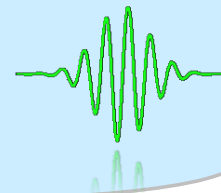# Some Analytic Dictionaries

**Fourier**

$$\phi_k(x) = e^{i2\pi kx}$$

**Smooth signals**

**Gabor**

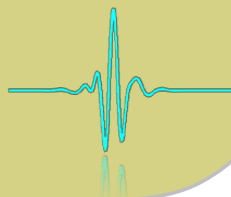$$\phi_{k,n}(x) = \boldsymbol{\omega}(x - \beta n)\, e^{i2\pi\alpha kx}$$

**Smooth signals**

**Wavelets**

$$\phi_{m,n}(x) = \alpha^{m/2} f(\alpha^m x - \beta n)$$
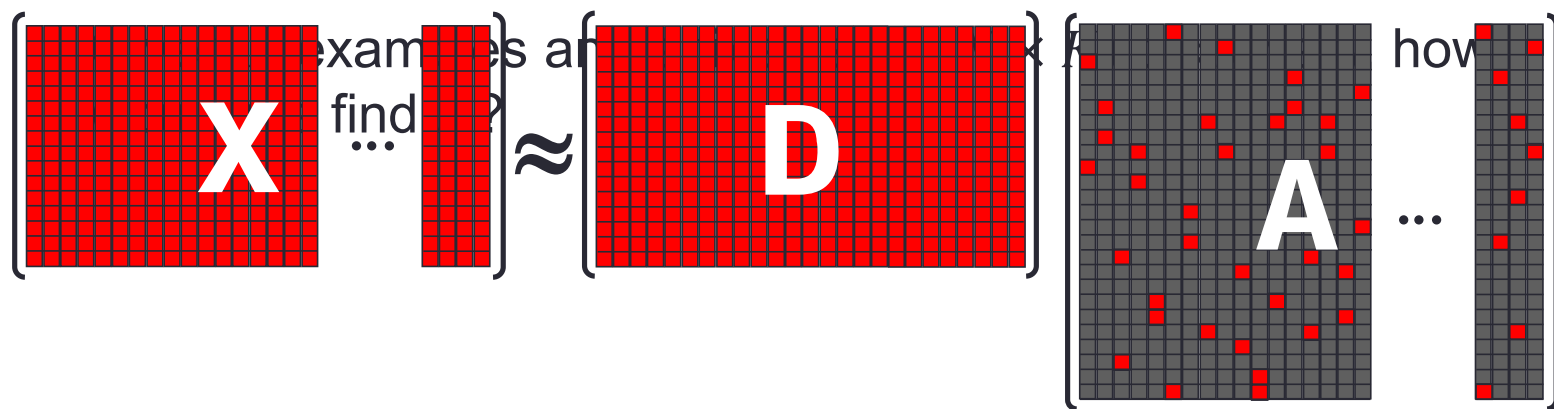
**Smooth + point singularities**

**Curvelets**

$$\phi_{m,n,\ell}(x) = \phi_m(R_{\Theta_\ell}(x - x_n^{m,\ell}))$$

**Smooth + curve singularities**

# The Dictionary Learning Problem

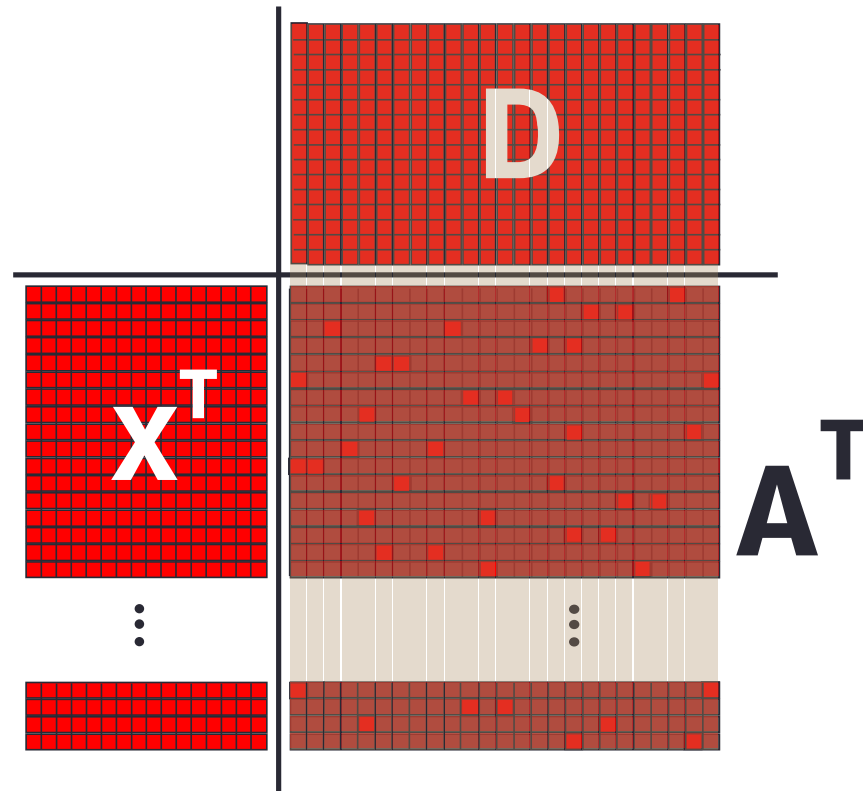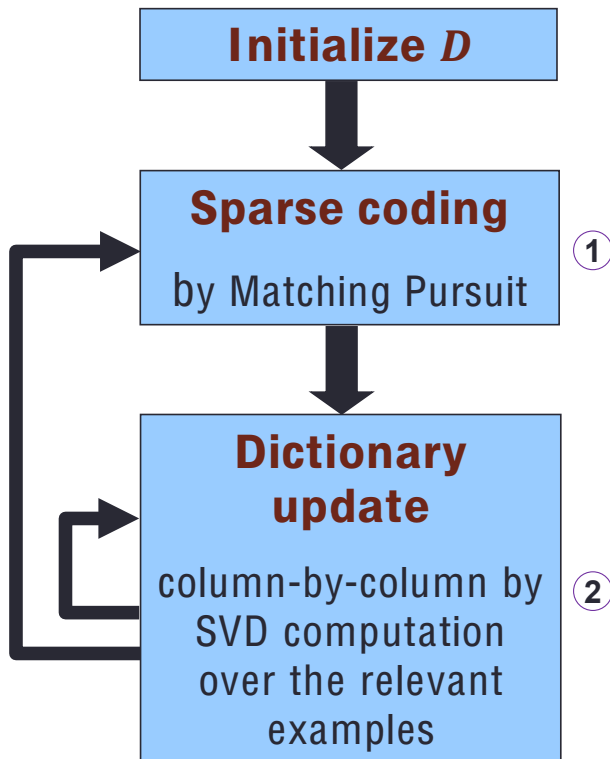...examples and... how find?

$$X \approx D \ A \ ...$$

$$\min_{D, A} \sum_{j=1}^{P} \| D\alpha_j - x_j \|_2^2 \quad s.t. \quad \| \alpha_j \|_0^0 \leq L \ \forall j$$

Each example is a linear combination of atoms from **D**

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

# K-SVD Algorithm – Overview

- Iterative process

# K-SVD Algorithm – Sparse Coding   ①

- Assume $D$ is known
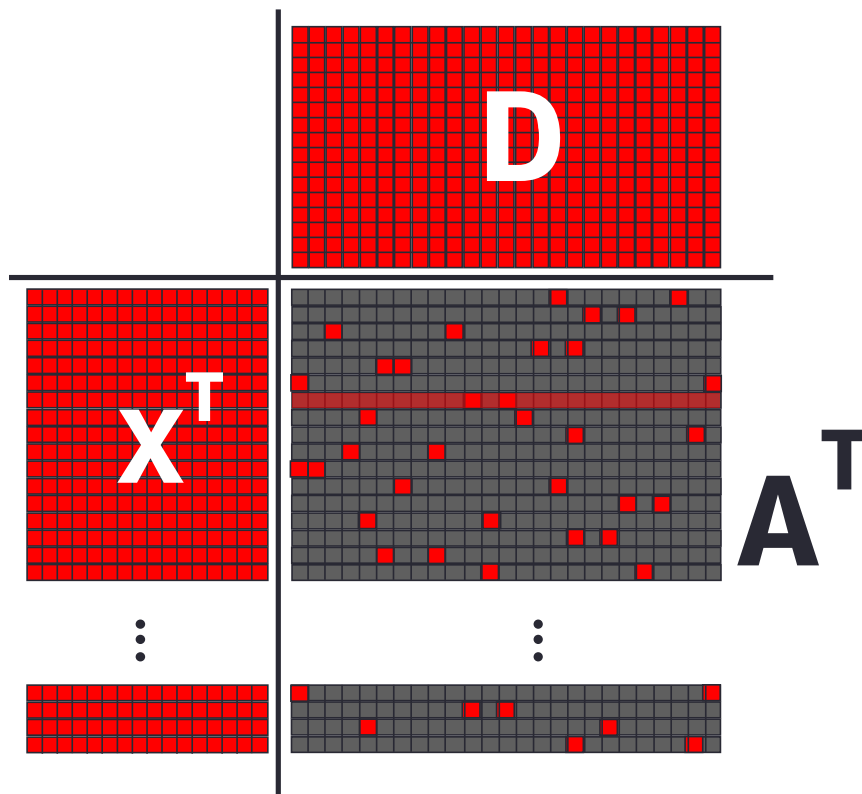
$$\min_{D, A} \sum_{j=1}^{P} \|D\alpha_j - x_j\|_2^2$$

$$s.t. \quad \|\alpha_j\|_0^0 \leq L \ \forall j$$

For the $k^{th}$ raw we solve
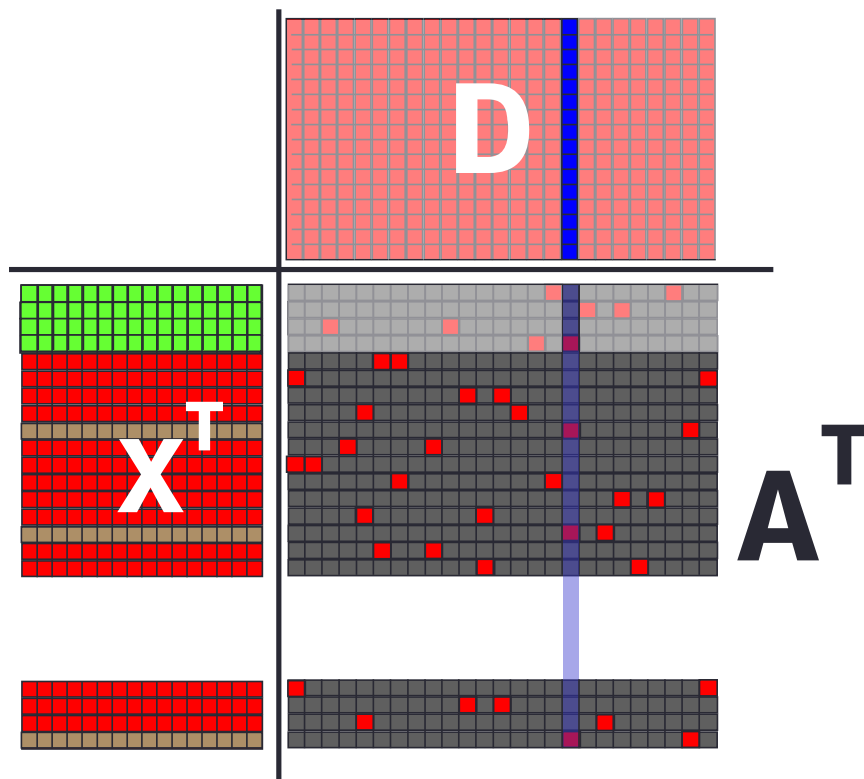
$$\min_{\alpha_k} \|D\alpha_k - x_k\|_2^2$$

$$s.t. \quad \|\alpha_k\|_p^p \leq L$$

Solved by Matching Pursuit

# K-SVD Algorithm – Dictionary Learning ②

- Update $D$ and $A$ simutaneously

- Target a column $d_k$
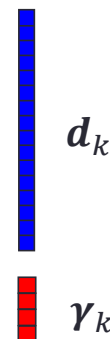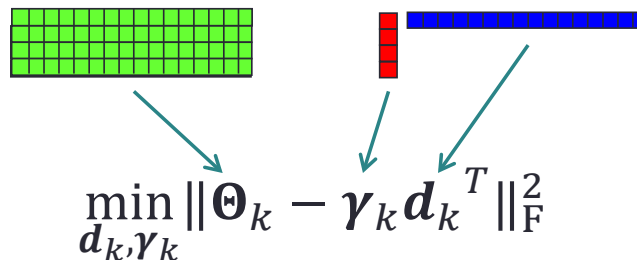
- Identify the examples that use the column $d_k$

# K-SVD Algorithm – Dictionary Update

- A good dictionary column $\boldsymbol{d}_k$ should well describe $\boldsymbol{\Theta}_k$

- Update $\boldsymbol{d}_k$ and $\boldsymbol{\gamma}_k$ by the cost function:

$$\min_{\boldsymbol{d}_k, \boldsymbol{\gamma}_k} \|\boldsymbol{\Theta}_k - \boldsymbol{\gamma}_k \boldsymbol{d}_k{}^{\mathrm{T}}\|_{\mathrm{F}}^2$$

by using SVD

$$\boldsymbol{\Theta}_k$$

$$\boldsymbol{d}_k$$

$$\boldsymbol{\gamma}_k$$

$$\min_{\boldsymbol{d}_k, \boldsymbol{\gamma}_k} \|\boldsymbol{\Theta}_k - \boldsymbol{\gamma}_k \boldsymbol{d}_k{}^{T}\|_{\mathrm{F}}^2$$

# K-SVD vs. K-means



**Initialize dictionary $D$**

**Sparse coding**

by Matching Pursuit

**Dictionary update**

column-by-column by SVD computation

**Initialize clustering centers**

**Assignment for each data point**

**Cluster centers update**

Cluster-by-cluster

# K-SVD vs. K-means

- Atom/centroid training
- Data point clustering

**K-SVD**

Atom $\phi_1$

Atom $\phi_2$

$\approx 0.3\phi_1 + 0.1\phi_2 + 0.6\phi_3$

Atom $\phi_3$
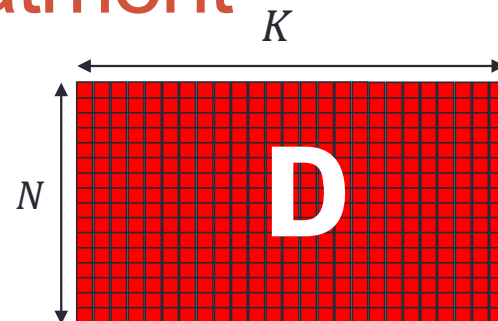
**K-means**

Centroid 1

Centroid 2

Centroid 3

# From Local to Global Treatment

$K$



$N$  **D**

- The K-SVD algorithm is reasonably fast for $N$ in the range of 10 to 400

- As $N$ grows, the complexity and the memory requirements of the K-SVD become prohibitive

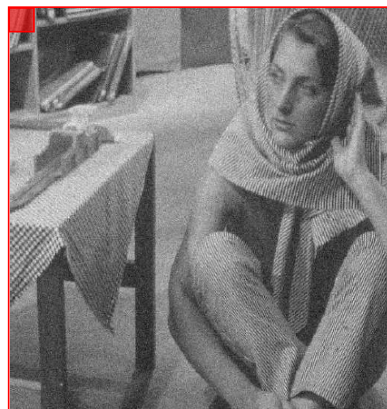- One solution: separate an image into patches of size $\sqrt{N}$-by-$\sqrt{N}$ in the image, including overlaps

$$\hat{x} = \underset{D, \{\alpha_{ij}\}_{ij}}{\mathrm{argmin}} \frac{1}{2} \|x - y\|_2^2 + \mu \sum_{ij} \|R_{ij}x - D\alpha_{ij}\|_2^2 \quad s.t. \quad \|\alpha_{ij}\|_0^0 \leq L$$

A binary matrix that extracts a patch in the $ij$ location

# What Data to Train On?

- Option 1:
  - Use a database of images
- Option 2:
  - Use the corrupted image itself!!
  - Simply sweep through all patches of size $\sqrt{N}$-by-$\sqrt{N}$ overlapping blocks

# Image Denoising Procedure

$$\widehat{x} = \underset{D, \{\alpha_{ij}\}_{ij}}{\operatorname{argmin}} \frac{1}{2} \|x - y\|_2^2 + \mu \sum_{ij} \|R_{ij}x - D\alpha_{ij}\|_2^2 \quad s.t. \quad \|\alpha_{ij}\|_0^0 \le L$$

$x = y$ and $D$ known

$x$ and $\alpha_{ij}$ known

$D$ and $\alpha_{ij}$ known

**Compute $\alpha_{ij}$ per patch**

$$\min_{D} \|R_{ij}x - D\alpha_{ij}\|_2^2$$
$$s.t. \quad \|\alpha_{ij}\|_0^0 \le L$$

**using the Matching Pursuit**

**Compute $D$**

$$\min_{D} \sum_{ij} \|R_{ij}x - D\alpha_{ij}\|_2^2$$

**using SVD, updating one column at a time**

**Compute $x$ by**

$$x = \left[I + \mu \sum_{ij} R_{ij}^T R_{ij}\right]^{-1}$$
$$\left[y + \mu \sum_{ij} R_{ij}^T D\alpha_{ij}\right]$$
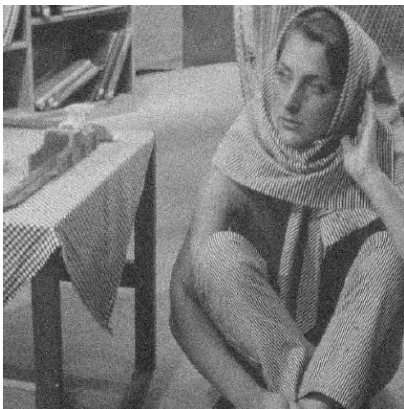
**by averaging of shifted patches**
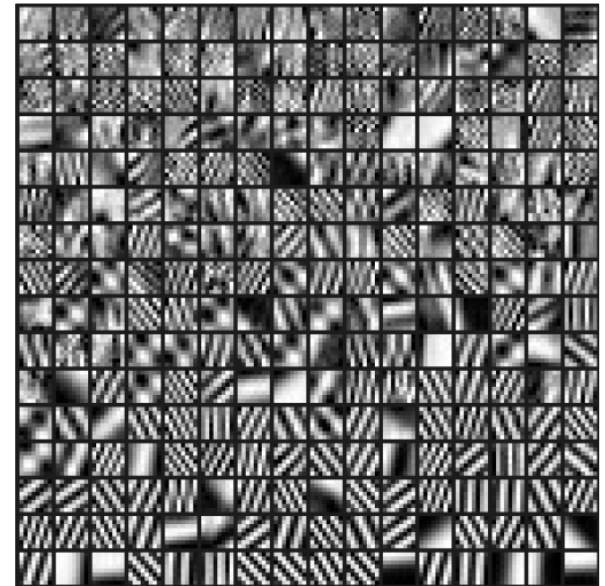
K-SVD

# Image Denoising

Source image



Noisy image

Result

The obtained dictionary after
10 iterations

# Summary

- Sparsity is an idea that can be used in designing tools to perform deonsing in signal/image processing

- The K-SVD algorithm is an efficient tool that can be applied to perform sparse coding and dictionary learning

# References

- R. Rubinstein, *Introduction to Sparse Representation and the K-SVD Algorithm*

- R. Rubinstein, *Sparsity-Based Signal Models and the Sparse K-SVD Algorithm*

- Andrew Ng, Image Classification using Sparse Coding, ECCV10 Tutorial