

K-SVD FOR DUMMIES

An Introduction to Sparse Representation and the K-SVD Algorithm



Ron Rubinstein

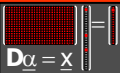
The CS Department
The Technion – Israel Institute of Technology
Haifa 32000, Israel

Friedrich-Alexander-Universität
Erlangen-Nürnberg



University of Erlangen - Nürnberg

April 2008

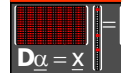


Noise Removal ?

Our story begins with image denoising ...



- Practical application
- A convenient platform (being the simplest inverse problem) for testing basic ideas in image processing.



An Introduction to
Sparse Representation
And the K-SVD Algorithm
Ron Rubinstein

Denoising By Energy Minimization

Many of the proposed denoising algorithms are related to the minimization of an energy function of the form

$$f(\underline{x}) = \frac{1}{2} \|\underline{x} - \underline{y}\|_2^2 + \Pr(\underline{x})$$

\underline{y} : Given measurements
 \underline{x} : Unknown to be recovered

Sanity (relation to measurements)

Prior or regularization

- This is in-fact a Bayesian point of view, adopting the Maximum-Aposteriori Probability (MAP) estimation.
- Clearly, the wisdom in such an approach is within the choice of the prior – **modeling the images** of interest.



Thomas Bayes
1702 - 1761

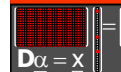


An Introduction to
Sparse Representation
And the K-SVD Algorithm
Ron Rubinstein

The Evolution Of $\Pr(\underline{x})$

During the past several decades we have made all sort of guesses about the prior $\Pr(\underline{x})$ for images:

$\Pr(\underline{x}) = \lambda \ \underline{x}\ _2^2$ Energy	$\Pr(\underline{x}) = \lambda \ \mathbf{L}\underline{x}\ _2^2$ Smoothness	$\Pr(\underline{x}) = \lambda \ \mathbf{L}\underline{x}\ _w^2$ Adapt+Smooth	$\Pr(\underline{x}) = \lambda p\{\mathbf{L}\underline{x}\}$ Robust Statistics
$\Pr(\underline{x}) = \lambda \ \nabla \underline{x}\ _1$ Total-Variation	$\Pr(\underline{x}) = \lambda \ \mathbf{W}\underline{x}\ _1$ Wavelet Sparsity	$\Pr(\underline{x}) = \lambda \ \mathbf{B}\underline{x}\ _1$ Bilateral Filter	$\Pr(\underline{x}) = \lambda \ \underline{\alpha}\ _0^0$ for $\underline{x} = \mathbf{D}\underline{\alpha}$ Sparse & Redundant



An Introduction to
Sparse Representation
And the K-SVD Algorithm
Ron Rubinstein

Agenda

1. A Visit to *Sparseland*

Introducing sparsity & overcompleteness

2. Transforms & Regularizations

How & why should this work?

3. What about the dictionary?

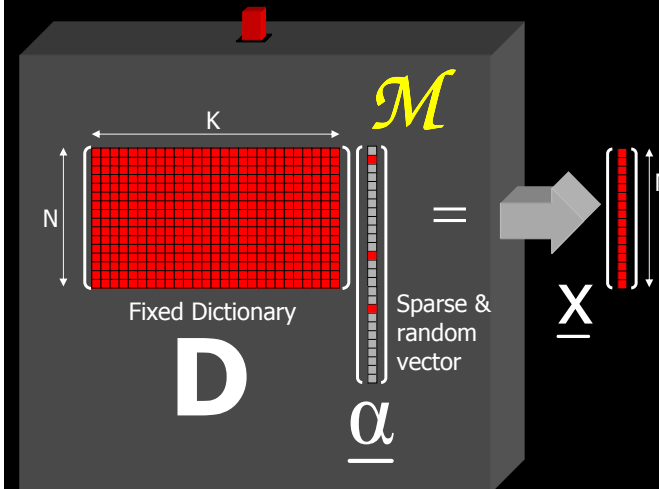
The quest for the origin of signals

4. Putting it all together

Image filling, denoising, compression, ...



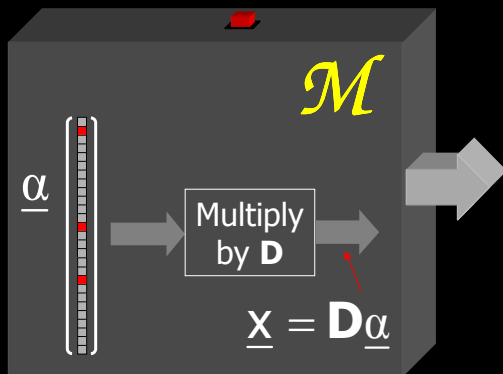
Generating Signals in *Sparseland*



- Every column in D (dictionary) is a prototype signal (atom).

- The vector α is generated randomly with few non-zeros in random locations and with random values.

Sparseland Signals Are Special



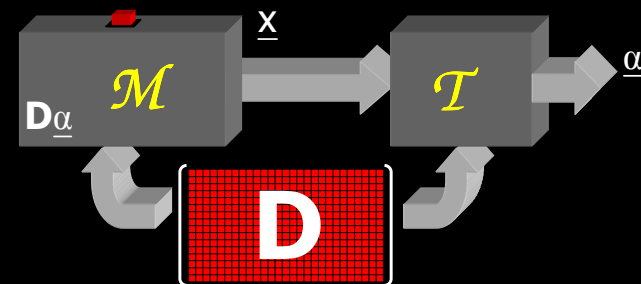
- **Simple:** Every signal is built as a linear combination of a few atoms from the dictionary D .

- **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. [Olshausen & Field ('96)]

Transforms in *Sparseland* ?

- Assume that x is known to emerge from M .
- How about "Given x , find the α that generated it in M " ?

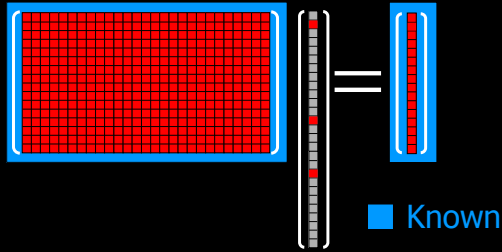


Problem Statement

We need to solve an under-determined linear system of equations:

$$\mathbf{D} \underline{\alpha} = \underline{x}$$

- Among all (infinitely many) possible solutions we want the **sparsest** !!
- We will measure sparsity using the L_0 "norm":



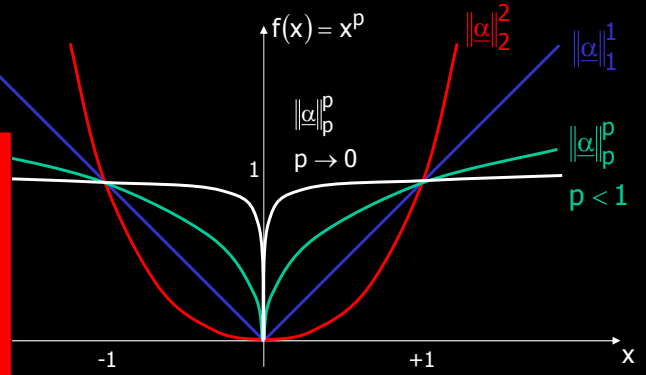
→ $\|\underline{\alpha}\|_0$

Measure of Sparsity?

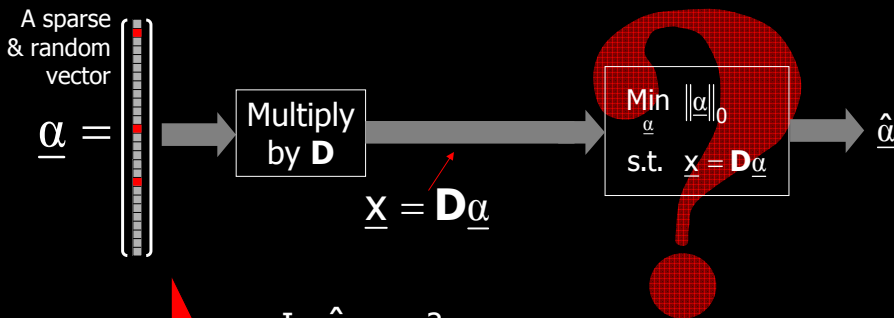
$$\|\underline{x}\|_p^p = \sum_{j=1}^k |x_j|^p$$

As $p \rightarrow 0$ we get a count of the non-zeros in the vector

→ $\|\underline{\alpha}\|_0$



Where We Are

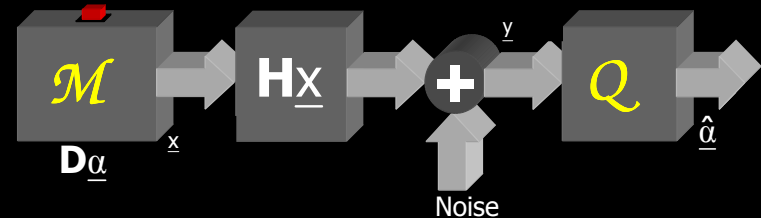


3 Major Questions

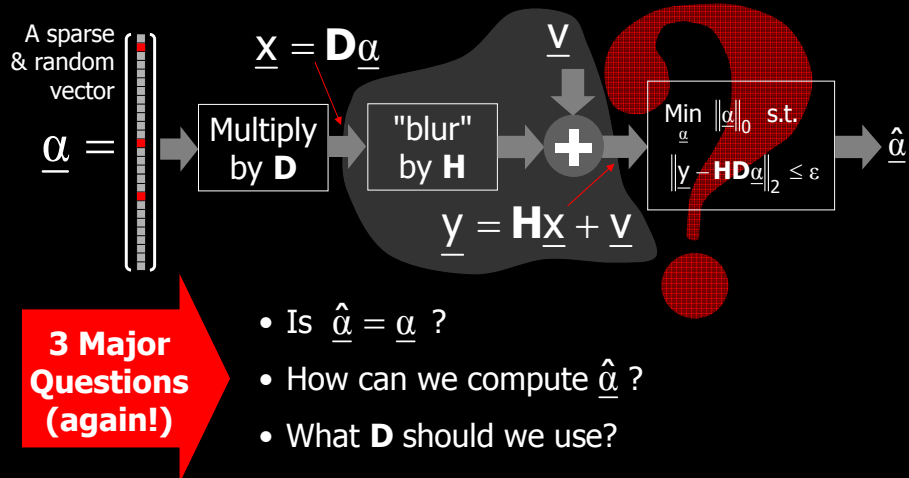
- Is $\hat{\alpha} = \alpha$?
- NP-hard: practical ways to get $\hat{\alpha}$?
- How do we know \mathbf{D} ?

Inverse Problems in *Sparseland* ?

- Assume that \underline{x} is known to emerge from \mathcal{M} .
- Suppose we observe $\underline{y} = \mathbf{H}\underline{x} + \underline{v}$, a degraded and noisy version of \underline{x} with $\|\underline{v}\|_2 \leq \epsilon$. How do we recover \underline{x} ?
- How about "find the $\underline{\alpha}$ that generated \underline{y} " ?

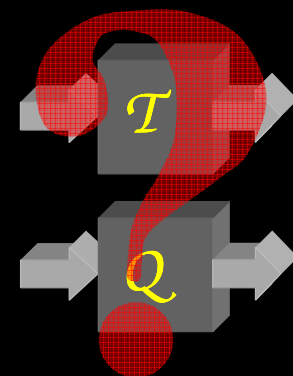


Inverse Problem Statement



Agenda

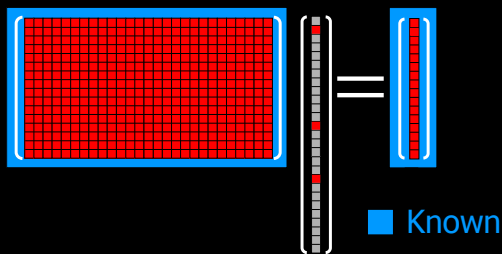
1. A Visit to *Sparseland*
Introducing sparsity & overcompleteness
2. **Transforms & Regularizations**
How & why should this work?
3. What about the dictionary?
The quest for the origin of signals
4. Putting it all together
Image filling, denoising, compression, ...



The Sparse Coding Problem

Our dream for now: Find the sparsest solution to

$$\mathbf{D}\underline{\alpha} = \underline{x}$$

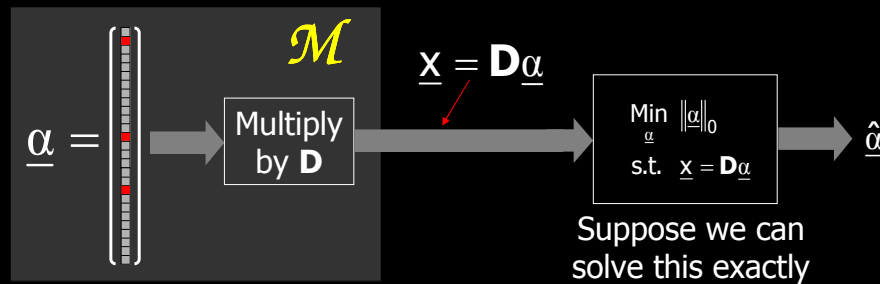


Put formally,

$$\text{Min}_{\underline{\alpha}} \|\underline{\alpha}\|_0 \text{ s.t. } \underline{x} = \mathbf{D}\underline{\alpha}$$

known

Question 1 – Uniqueness?



Why should we necessarily get $\hat{\underline{\alpha}} = \underline{\alpha}$?

It might happen that eventually $\|\hat{\underline{\alpha}}\|_0 < \|\underline{\alpha}\|_0$.

Matrix "Spark"

Definition: Given a matrix \mathbf{D} , $\sigma = \text{Spark}\{\mathbf{D}\}$ is the smallest number of columns that are linearly dependent.

Donoho & Elad ('02)

- By definition, if $\mathbf{D}\underline{v} = 0$ then $\|\underline{v}\|_0 \geq \sigma$
- Say I have $\underline{\alpha}_1$ and you have $\underline{\alpha}_2$, and the two are different representations of the same \underline{x} :

$$\underline{x} = \mathbf{D}\underline{\alpha}_1 = \mathbf{D}\underline{\alpha}_2 \Rightarrow \mathbf{D}(\underline{\alpha}_1 - \underline{\alpha}_2) = \underline{0}$$

$$\Rightarrow \|\underline{\alpha}_1 - \underline{\alpha}_2\|_0 \geq \sigma$$

Uniqueness Rule

- Now, what if my $\underline{\alpha}_1$ satisfies $\|\underline{\alpha}_1\|_0 < \frac{\sigma}{2}$?
- The rule $\|\underline{\alpha}_1 - \underline{\alpha}_2\|_0 \geq \sigma$ implies that $\|\underline{\alpha}_2\|_0 > \frac{\sigma}{2}$!

Uniqueness

If we have a representation that satisfies

$$\frac{\sigma}{2} > \|\underline{\alpha}\|_0$$

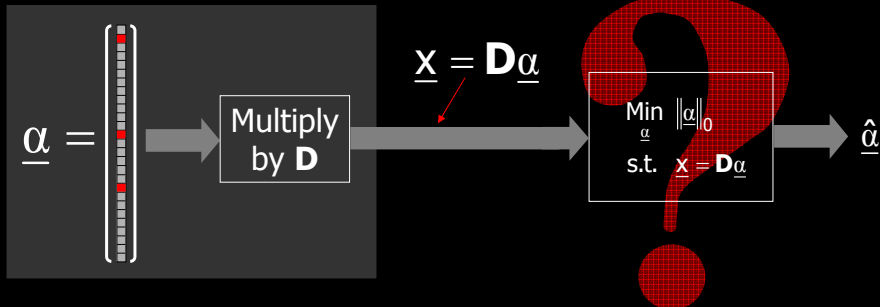
then necessarily it is the sparsest.

Donoho & Elad ('02)

So, if \mathcal{M} generates signals using "sparse enough" $\underline{\alpha}$, the solution of \rightarrow will find them exactly.

$$P_0 : \underset{\underline{\alpha}}{\text{Min}} \|\underline{\alpha}\|_0 \text{ s.t. } \underline{x} = \mathbf{D}\underline{\alpha}$$

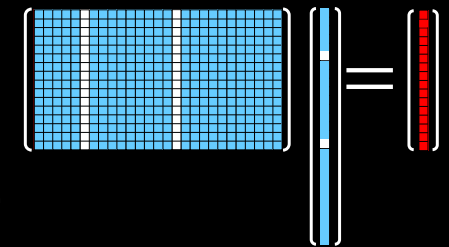
Question 2 – Practical P_0 Solver?



Are there reasonable ways to find $\hat{\underline{\alpha}}$?

Matching Pursuit (MP) Mallat & Zhang (1993)

- The MP is a greedy algorithm that finds one atom at a time.
- Step 1: find the one atom that **best matches** the signal.
- Next steps: given the previously found atoms, find the next **one** to **best fit** ...
- The Orthogonal MP (OMP) is an improved version that re-evaluates the coefficients after each round.



Basis Pursuit (BP) Chen, Donoho, & Saunders (1995)

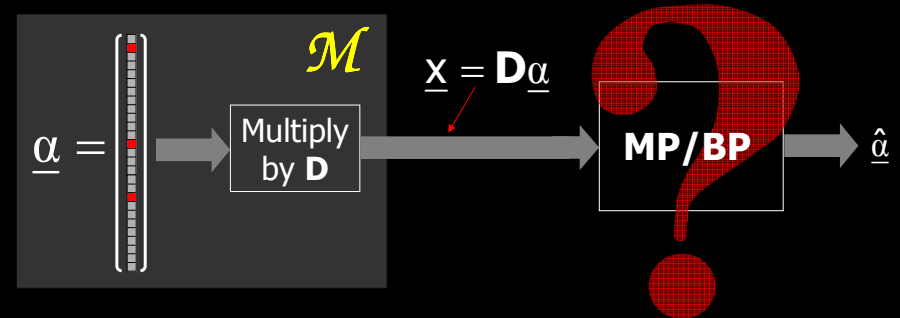
Instead of solving
 $\text{Min}_{\underline{\alpha}} \|\underline{\alpha}\|_0 \text{ s.t. } \underline{x} = \mathbf{D}\underline{\alpha}$



Solve this:
 $\text{Min}_{\underline{\alpha}} \|\underline{\alpha}\|_1 \text{ s.t. } \underline{x} = \mathbf{D}\underline{\alpha}$

- The newly defined problem is convex (linear programming).
- Very efficient solvers can be deployed:
 - Interior point methods [Chen, Donoho, & Saunders ('95)] ,
 - Iterated shrinkage [Figueroa & Nowak ('03), Daubechies, Defrise, & Demole ('04), Elad ('05), Elad, Matalon, & Zibulevsky ('06)].

Question 3 – Approx. Quality?



How effective are MP/BP ?

BP and MP Performance

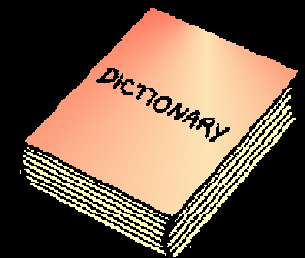
Donoho & Elad ('02)
 Gribonval & Nielsen ('03)
 Tropp ('03)
 Temlyakov ('03)

Given a signal \underline{x} with a representation $\underline{x} = \mathbf{D}\underline{\alpha}$, if $\|\underline{\alpha}\|_0 < (\text{some threshold})$ then BP and MP are **guaranteed** to find it.

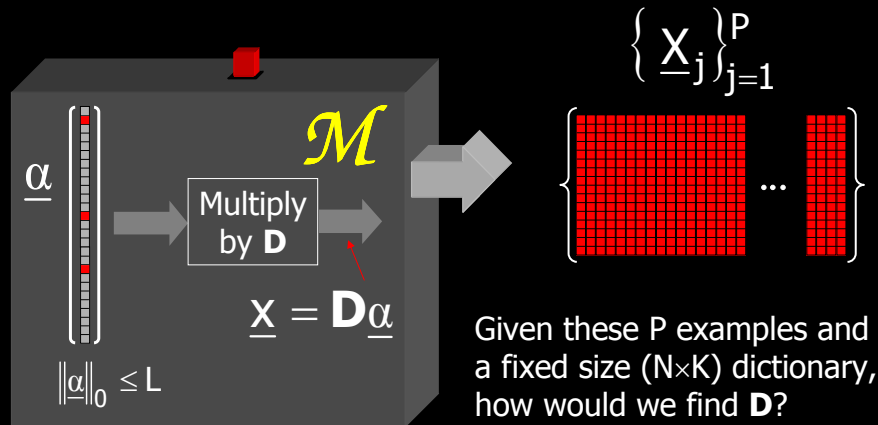
- MP and BP are different in general (hard to say which is better).
- The above results correspond to the worst-case.
- Average performance results available too, showing much better bounds [Donoho ('04), Candes et.al. ('04), Tanner et.al. ('05), Tropp et.al. ('06)].
- Similar results for general inverse problems [Donoho, Elad & Temlyakov ('04), Tropp ('04), Fuchs ('04), Gribonval et. al. ('05)].

Agenda

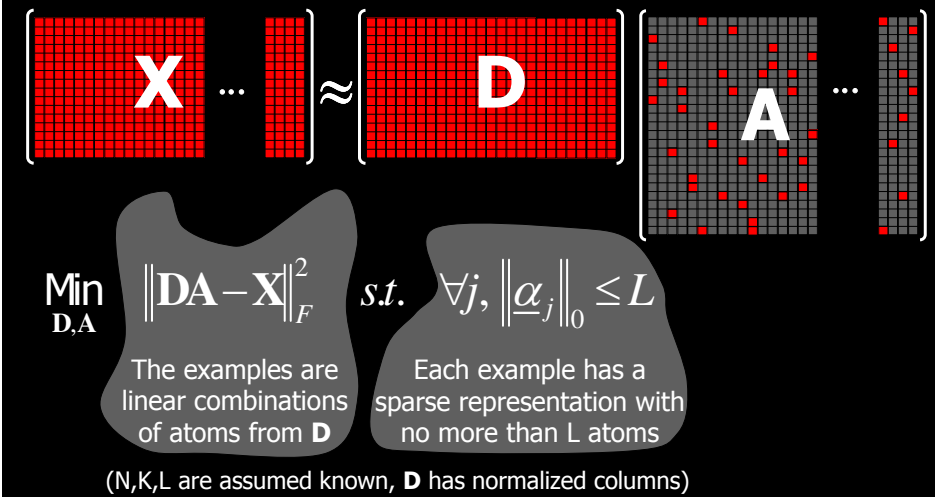
1. A Visit to *Sparseland*
 Introducing sparsity & overcompleteness
2. Transforms & Regularizations
 How & why should this work?
3. **What about the dictionary?**
The quest for the origin of signals
4. Putting it all together
 Image filling, denoising, compression, ...



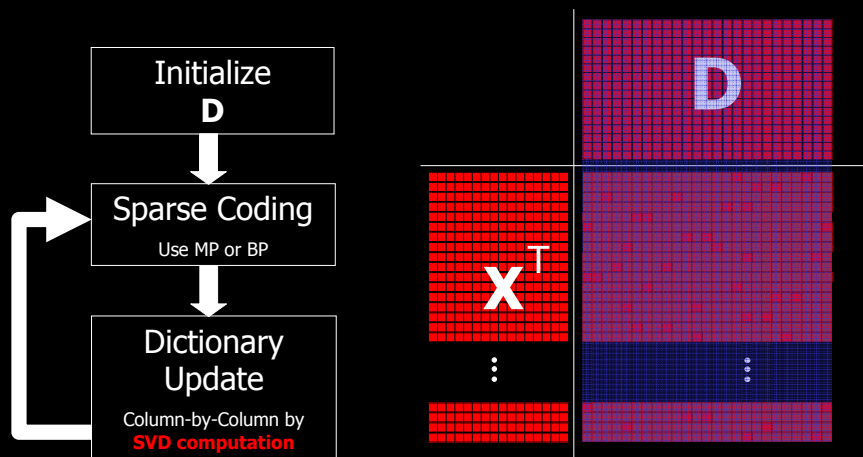
Problem Setting



The Objective Function

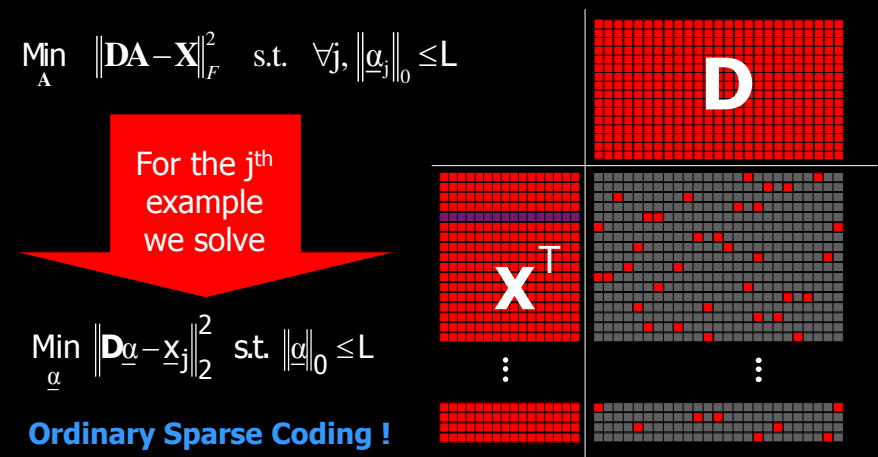


K-SVD – An Overview



Aharon, Elad & Bruckstein ('04)

K-SVD: Sparse Coding Stage



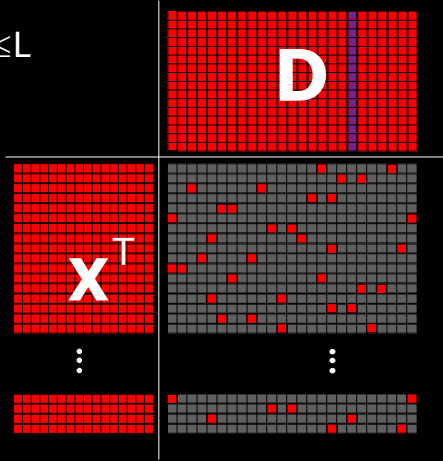
K-SVD: Dictionary Update Stage

$$\text{Min}_D \|\mathbf{D}\mathbf{A} - \mathbf{X}\|_F^2 \quad \text{s.t.} \quad \forall j, \|\mathbf{a}_j\|_0 \leq L$$

For the k^{th} atom we solve

$$\text{Min}_{\mathbf{d}_k} \|\mathbf{d}_k \mathbf{a}_k^T - \mathbf{E}_k\|_F^2$$

$$\mathbf{E}_k = \sum_{j \neq k} \mathbf{d}_j \mathbf{a}_j^T - \mathbf{X} \quad (\text{the residual})$$



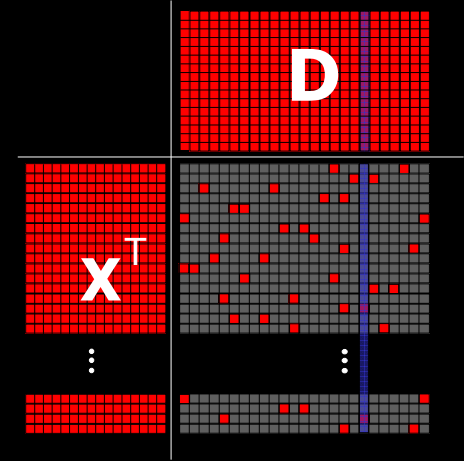
K-SVD Dictionary Update Stage

$$\text{Min}_{\mathbf{d}_k} \|\mathbf{d}_k \mathbf{a}_k^T - \mathbf{E}_k\|_F^2$$

We can do better than this

$$\text{Min}_{\mathbf{d}_k, \mathbf{a}_k} \|\mathbf{d}_k \mathbf{a}_k^T - \mathbf{E}_k\|_F^2$$

But wait! What about sparsity?



K-SVD Dictionary Update Stage

We want to solve:

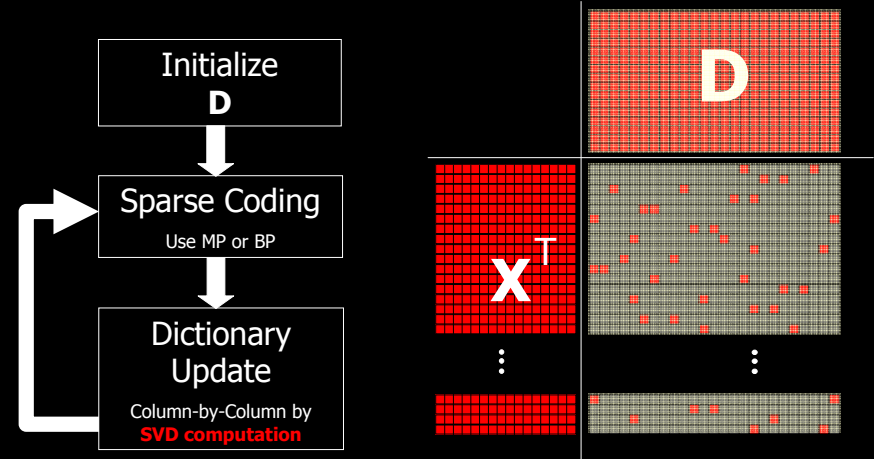
$$\text{Min}_{\mathbf{d}_k, \mathbf{a}_k} \left\| \begin{bmatrix} \mathbf{d}_k & \mathbf{a}_k^T \end{bmatrix} - \begin{bmatrix} \mathbf{E}_k \end{bmatrix} \right\|_F^2$$

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

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

Solve with **SVD!**

The K-SVD Algorithm – Summary



Agenda

1. A Visit to *Sparseland*
Introducing sparsity & overcompleteness
2. Transforms & Regularizations
How & why should this work?
3. What about the dictionary?
The quest for the origin of signals
4. **Putting it all together**
Image filling, denoising, compression, ...



Image Inpainting: Theory

- Assumption: the signal \underline{x} was created by $\underline{x} = D\alpha_0$ with a very sparse α_0 .
- Missing values in \underline{x} imply missing rows in this linear system.

$$D\alpha_0 = \underline{X}$$

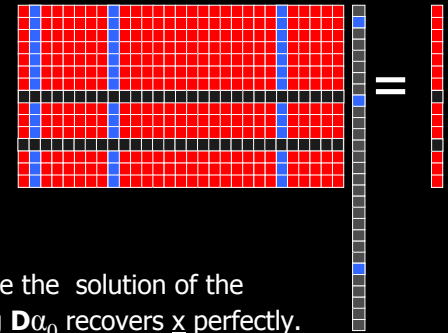
- By removing these rows, we get

$$\tilde{D}\alpha_0 = \tilde{\underline{X}}$$

- Now solve

$$\text{Min}_{\alpha} \|\alpha\|_0 \text{ s.t. } \tilde{\underline{X}} = \tilde{D}\alpha$$

- If α_0 was sparse enough, it will be the solution of the above problem! Thus, computing $D\alpha_0$ recovers \underline{x} perfectly.



Inpainting: The Practice

- We define a diagonal mask operator \mathbf{W} representing the lost samples, so that

$$\underline{y} = \mathbf{W}\underline{x} + \underline{v} \quad w_{i,i} \in \{0,1\}$$

- Given \underline{v} , we try to recover the representation of \underline{x} , by solving

$$\hat{\alpha} = \text{ArgMin}_{\alpha} \|\alpha\|_0 \text{ s.t. } \|\underline{y} - \mathbf{W}D\alpha\|_2 \leq \epsilon \quad \rightarrow \quad \hat{\underline{x}} = D\hat{\alpha}$$

- We use a dictionary that is the sum of two dictionaries, to get an effective representation of both texture and cartoon contents. This also leads to image separation [Elad, Starck, & Donoho ('05)]

Inpainting Results

Source

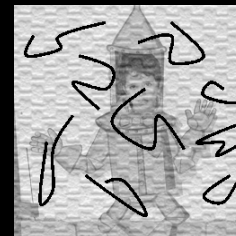


Image inpainting [2, 10, 20, 38] is the process of recovering data in a designated region of a still or video image. Applications range from removing objects from damaged paintings and photographs to produce a revised image in which the missing data is seamlessly merged into the image in a way that is undetectable by a typical viewer. Traditionally, inpainting has been done by professional artists. For photographs, inpainting is used to revert deterioration or to remove elements (e.g., removal of stamped text from photographs, the infamous "airbrushed out" enemies [20]). A current active area of research is the use of sparse representations for inpainting [20].

Dictionary:
Curvelet (cartoon)
+ Global DCT
(texture)

Outcome



Inpainting Results

Source



Image inpainting [2, 10, 20, 36] is the process of recovering data in a designated region of a still or motion picture. Applications range from removing objects from damaged paintings and photographs to producing a revised image in which the missing data is seamlessly merged into the image in a way that is not detectable by a typical viewer. Traditionally, inpainting has been done by professional artists. For photographs, inpainting is used to reverse deterioration, such as scratches or dust spots in film. It is also used to remove elements (e.g., removal of stamped numbers from photographs, the infamous "airbrushed" faces of the "airbrushed" enemies [20]). A current active area of research is the use of machine learning to improve inpainting results.

Dictionary:
Curvelet (cartoon)
+ Overlapped
DCT (texture)

Outcome



Inpainting Results



Denoising: Theory and Practice

Given a noisy image \underline{y} , we can clean it by solving

$$\hat{\alpha} = \underset{\alpha}{\text{ArgMin}} \|\alpha\|_0 \quad \text{s.t.} \quad \|\underline{y} - D\alpha\|_2 \leq \epsilon \quad \rightarrow \quad \hat{\underline{x}} = D\hat{\alpha}$$

Can we use the **K-SVD** dictionary?

With K-SVD, we cannot train a dictionary for an entire image. How do we go from **local** treatment of patches to a **global** prior?

Solution: force shift-invariant sparsity – for each $N \times N$ patch of the image, including overlaps.

From Local to Global Treatment

$$\hat{\alpha} = \underset{\alpha}{\text{ArgMin}} \|\alpha\|_0 \quad \text{s.t.} \quad \|\underline{y} - D\alpha\|_2 \leq \epsilon \quad \rightarrow \quad \hat{\underline{x}} = D\hat{\alpha}$$

For patches,
our MAP penalty
becomes

$$\hat{\underline{x}} = \underset{\underline{x}, \{\alpha_{ij}\}}{\text{ArgMin}} \frac{1}{2} \|\underline{x} - \underline{y}\|_2^2 + \mu \sum_{ij} \|\mathbf{R}_{ij} \underline{x} - D\alpha_{ij}\|_2^2 \quad \text{Extracts the } (i,j)^{\text{th}} \text{ patch}$$

s.t. $\|\alpha_{ij}\|_0 \leq L$ Our prior

What Data to Train On?

Option 1:

- ❑ Use a database of images: works quite well (~0.5-1dB below the state-of-the-art)



Option 2:

- ❑ Use the **corrupted image** itself !
- ❑ Simply sweep through all NxN patches (with overlaps) and use them to train
- ❑ Image of size 1000x1000 pixels $\Rightarrow \sim 10^6$ examples to use – more than enough.
- ❑ **This works much better!**

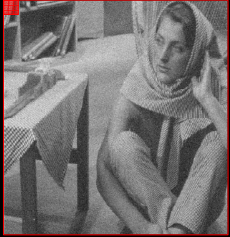


Image Denoising: The Algorithm

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

x and D known

x and α_{ij} known

D and α_{ij} known

Compute α_{ij} per patch

$$\alpha_{ij} = \underset{\alpha}{\text{Min}} \|R_{ij}x - D\alpha\|_2^2 \text{ s.t. } \|\alpha_{ij}\|_0 \leq L$$

using matching pursuit

Compute D to minimize

$$\underset{D}{\text{Min}} \sum_{ij} \|R_{ij}x - D\alpha_{ij}\|_2^2$$

using SVD, updating one column at a time

Compute x by

$$x = \left[I + \mu \sum_{ij} R_{ij}^T R_{ij} \right]^{-1} \left[y + \mu \sum_{ij} R_{ij}^T D \alpha_{ij} \right]$$

which is a simple averaging of shifted patches

K-SVD

Denoising Results

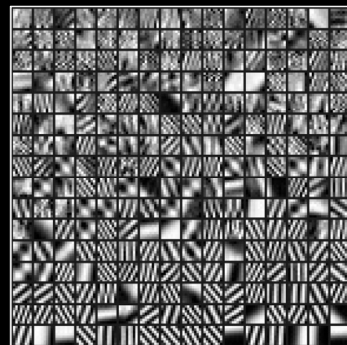


Source



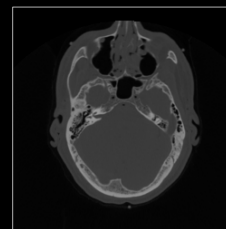
Result 30.829dB

Noisy image
PSNR = 22.1dB

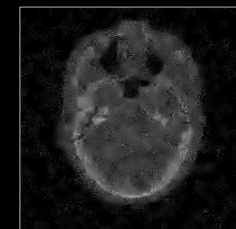


Obtained dictionary
after 10 iterations

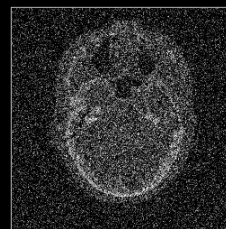
Denoising Results: 3D



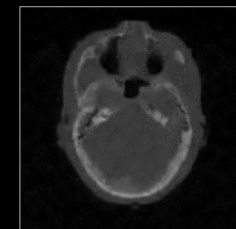
Source:
Vis. Male Head
(Slice #137)



2d-KSVD:
PSNR=27.3dB



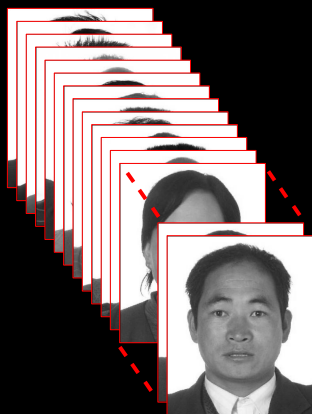
PSNR=12dB



3d-KSVD:
PSNR=32.4dB

Image Compression

- ❑ Problem: compressing **photo-ID images**.
- ❑ General purpose methods (JPEG, JPEG2000) do not take into account the specific family.
- ❑ By **adapting** to the image-content, better results can be obtained.



Compression: The Algorithm

Detect main features and align the images to a common reference (20 parameters)

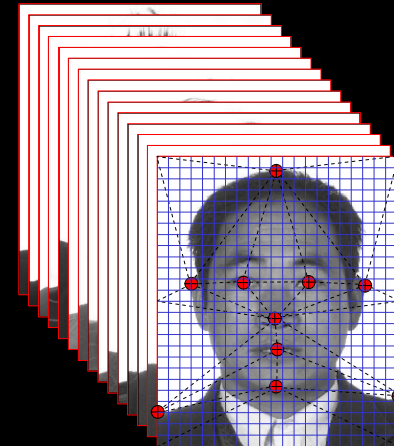
Divide each image to disjoint 15x15 patches, and for each compute a unique **dictionary**

Detect features and align

Divide to disjoint patches, and sparse-code each patch

Quantize and entropy-code

Training set (2500 images)



Training

Compression

Compression Results

Results for **820 bytes** per image

Original	JPEG	JPEG 2000	PCA	K-SVD
	11.99	10.49	8.81	5.56
	10.83	8.92	7.89	4.82
	10.93	8.71	8.61	5.58

Bottom: RMSE values

Compression Results

Results for **550 bytes** per image

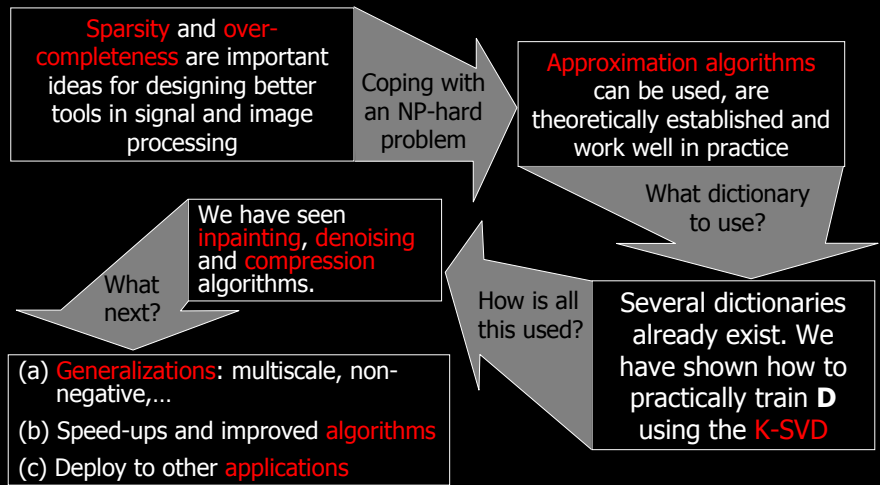
Original	JPEG	JPEG 2000	PCA	K-SVD
	15.81	13.89	10.66	6.60
	14.67	12.41	9.44	5.49
	15.30	12.57	10.27	6.36

Bottom: RMSE values

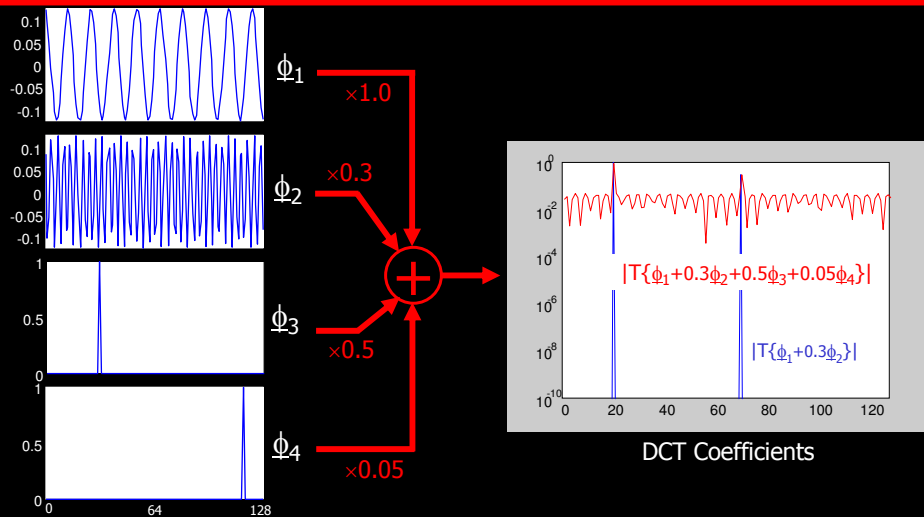
Today We Have Discussed

1. A Visit to *Sparseland*
Introducing sparsity & overcompleteness
2. Transforms & Regularizations
How & why should this work?
3. What about the dictionary?
The quest for the origin of signals
4. Putting it all together
Image filling, denoising, compression, ...

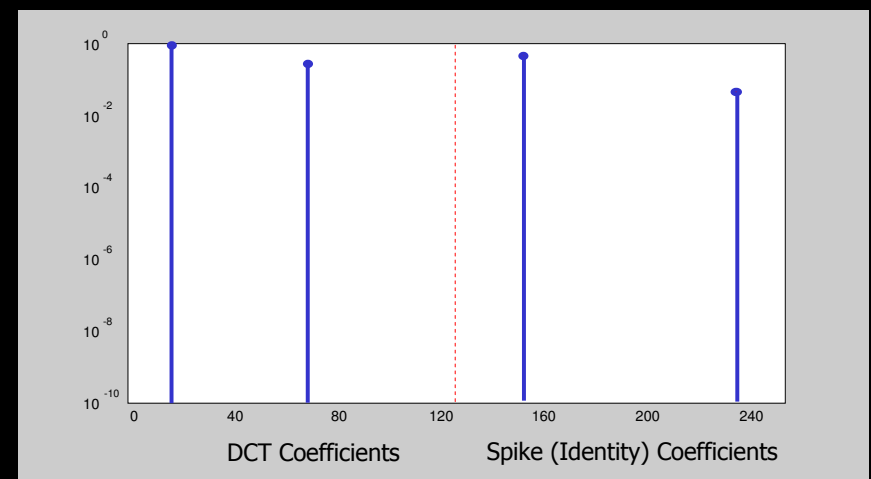
Summary



Why Over-Completeness?

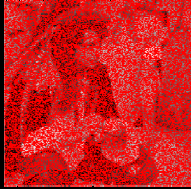


Desired Decomposition



Inpainting Results

70% Missing Samples



DCT (RMSE=0.04)



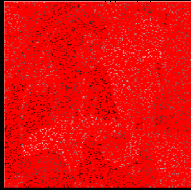
Haar (RMSE=0.045)



K-SVD (RMSE=0.03)



90% Missing Samples



DCT (RMSE=0.085_)



Haar (RMSE=0.07)



K-SVD (RMSE=0.06)

