# Sparse Modeling

## in

## Image Processing and Deep Learning

# Michael Elad

Computer Science Department
The Technion - Israel Institute of Technology
Haifa 32000, Israel

ICIP 2017
IEEE International Conference on Image Processing
September 17-20, 2017, Beijing, China

# This Lecture

| Sparseland<br>Sparse Representation Theory | → | CSC<br>Convolutional Sparse Coding | → | ML-CSC<br>Multi-Layered Convolutional Sparse Coding |
|:---:|:---:|:---:|:---:|:---:|

Sparsity-Inspired Models ⟶ Deep-Learning

Another underlying idea that will accompany us

▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪

Generative modeling of data sources enables
- A systematic algorithm development, &
- A theoretical analysis of their performance

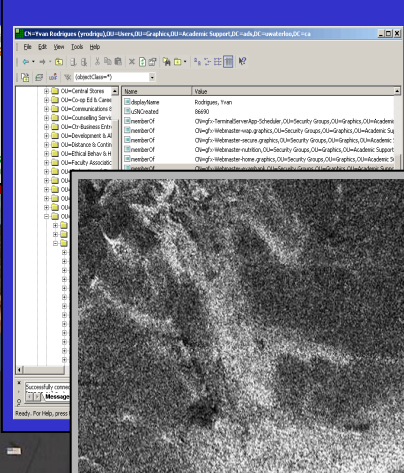# Multi-Layered Convolutional Sparse Modeling
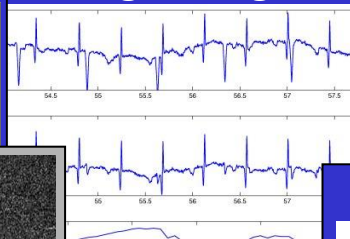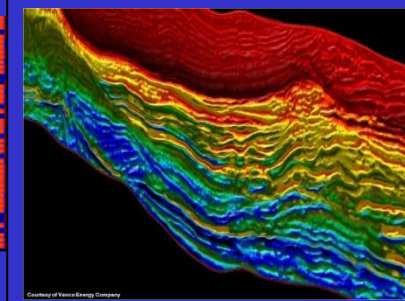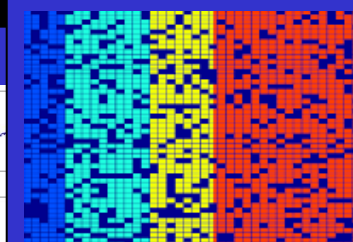
# Our Data is Structured

**Stock Market**

**Text Documents**

**Biological Signals**

**Matrix Data**

**Social Networks**

**Seismic Data**

**Still Images**

**Videos**

**Radar Imaging**

**Traffic info**

**Voice Signals**

**Medical Imaging**

**3D Objects**

o We are surrounded by various diverse sources of massive information

o Each of these sources have an internal structure, which can be exploited

o This structure, when identified, is the engine behind our ability to process this data

# Model?



**Fact 1**: This signal contains AWGN $\mathbb{N}(0,1)$

**Fact 2**: The clean signal is believed to be PWC

Effective removal of noise (and many other tasks) relies on an proper modeling of the signal

# Which Model to Choose?

- A model: a mathematical description of the underlying signal of interest, describing our beliefs regarding its structure

- The following is a partial list of commonly used models for images

- Good models should be simple while matching the signals

| Simplicity | ⟷ | Reliability |

- Models are almost always imperfect

Principal-Component-Analysis

Gaussian-Mixture

Markov Random Field

Laplacian Smoothness

DCT concentration

Wavelet Sparsity

Piece-Wise-Smoothness

C2-smoothness

Besov-Spaces

Total-Variation

Beltrami-Flow

# An Example: JPEG and DCT

**178KB – Raw data**

**24KB**

**20KB**

**12KB**

**8KB**

**4KB**

How & why does it works?

Discrete Cosine Trans.

The model assumption: after DCT, the top left coefficients to be dominant and the rest zeros

# Research in Signal/Image Processing

**Model**

**Problem (Application)**

**Signal**

**Numerical Scheme**

The fields of signal & image processing are essentially built of an evolution of models and ways to use them for various tasks

A New Research Work (and Paper) is Born

# What This Talk is all About?

## Data Models and Their Use

o Almost any task in data processing requires a model – true for denoising, deblurring, super-resolution, inpainting, compression, anomaly-detection, sampling, recognition, separation, and more

o Sparse and Redundant Representations offer a new and highly effective model – we call it

### *Sparseland*

o We shall describe this and descendant versions of it that lead all the way to … deep-learning

# Multi-Layered Convolutional

## Sparse Modeling

# A New Emerging Model

**Signal Processing**

Machine Learning

**Mathematics**

Wavelet Theory

Approximation Theory

Multi-Scale Analysis

*Sparseland*

Linear Algebra

Signal Transforms

Optimization Theory

Semi-Supervised Learning

Interpolation

Segmentation

Sensor-Fusion

Source-Separation

Classification

Compression

Inference (solving inverse problems)

Summarizing

Prediction

Denoising

Anomaly detection

Synthesis

Recognition

Clustering

Identification

Michael Elad
The Computer-Science Department
The Technion

# The *Sparseland* Model

o Task: model image patches of size 8×8 pixels

o We assume that a **dictionary** of such image patches is given, containing 256 **atom** images

o The *Sparseland* model assumption: every image patch can be described as a linear combination of **few** atoms
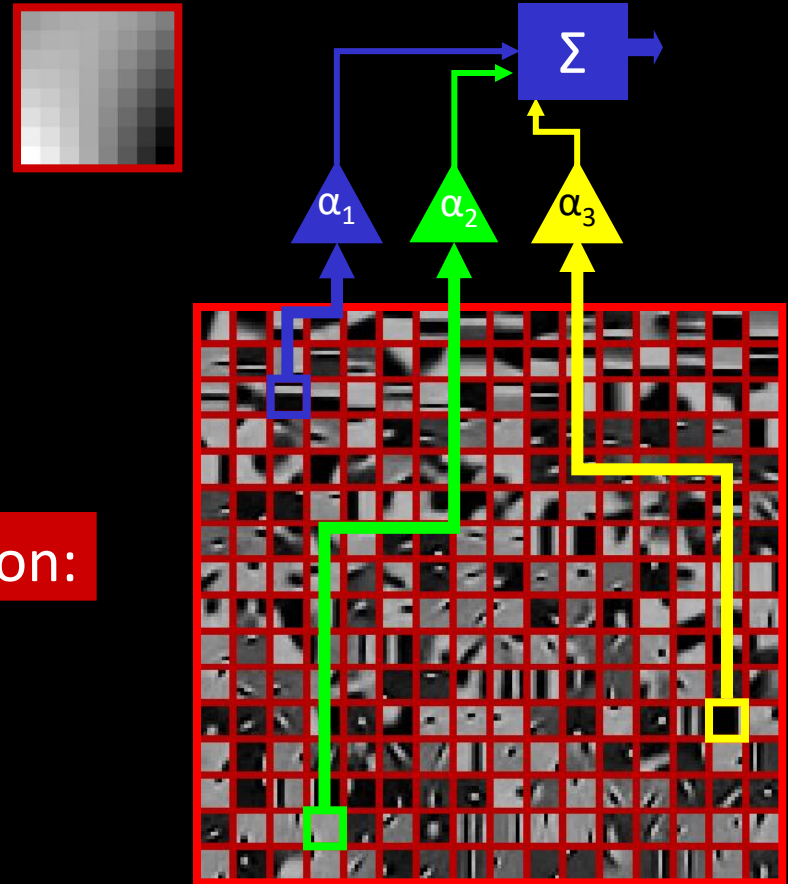
# The *Sparseland* Model

Properties of this model:

## Sparsity and Redundancy

○ We start with a 8-by-8 pixels patch and represent it using 256 numbers

   – This is a redundant representation

○ However, out of those 256 elements in the representation, only 3 are non-zeros

   – This is a sparse representation

○ Bottom line in this case: 64 numbers representing the patch are replaced by 6 (3 for the indices of the non-zeros, and 3 for their entries)

$\Sigma$

$\alpha_1$ $\alpha_2$ $\alpha_3$

# Chemistry of Data

We could refer to the *Sparseland*
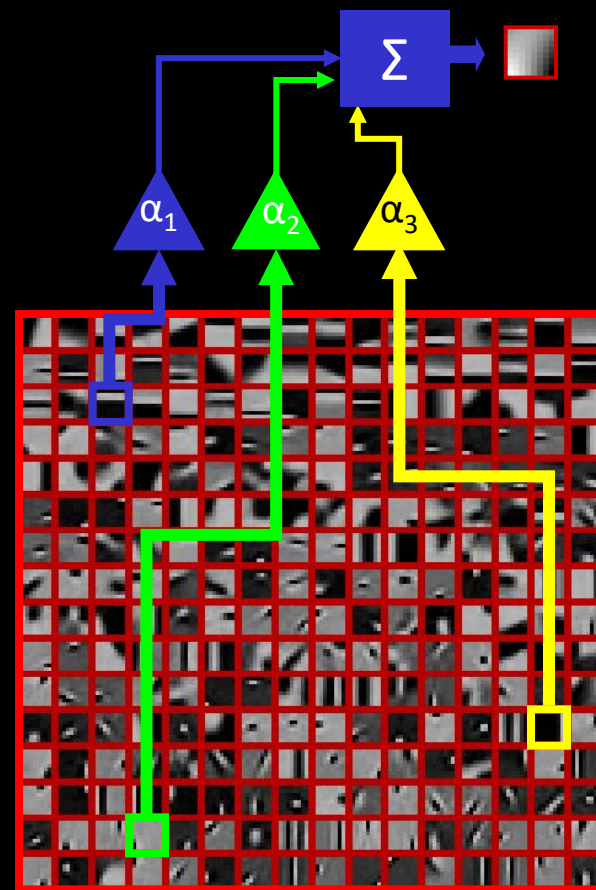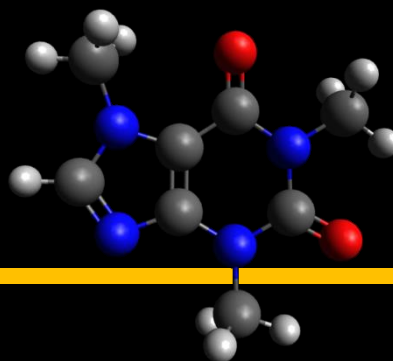model as the chemistry of information:

o Our dictionary stands for the Periodic Table
containing all the elements

o Our model follows a similar rationale:
Every molecule is built of **few** elements

# *Sparseland* : A Formal Description



$$\mathcal{M}$$

m

n

A Dictionary
**D**

α

A sparse vector

=

**X**

n

o Every column in **D** (dictionary) is a prototype signal (atom)

o The vector $\underline{\alpha}$ is generated with few non-zeros at arbitrary locations and values

o This is a generative model that describes how (we believe) signals are created

# Difficulties with *Sparseland*

o Problem 1: Given a signal, how can we find its atom decomposition?

o A simple example:

- There are 2000 atoms in the dictionary

- The signal is known to be built of 15 atoms

$$\binom{2000}{15} \approx 2.4e + 37 \quad \text{possibilities}$$

- If each of these takes 1nano-sec to test, will take ~7.5e20 years to finish !!!!!!

o So, are we stuck?

Michael Elad
The Computer-Science Department
The Technion

# Atom Decomposition Made Formal

$$\min_\alpha \|\alpha\|_0 \ \text{s.t.} \ x = \mathbf{D}\alpha$$

$$\min_\alpha \|\alpha\|_0 \ \text{s.t.} \ \|\mathbf{D}\alpha - y\|_2 \leq \varepsilon$$

Approximation Algorithms

**Relaxation methods**

Basis-Pursuit

**Greedy methods**

Thresholding/OMP

- $L_0$ – counting number of non-zeros in the vector

- This is a projection onto the *Sparseland* model

- These problems are known to be NP-Hard problem

Michael Elad
The Computer-Science Department
The Technion

# Pursuit Algorithms

$$\min_{\alpha} \|\alpha\|_0 \;\; \text{s.t.} \;\; \|\mathbf{D}\alpha - y\|_2 \leq \varepsilon$$

Approximation Algorithms

## Basis Pursuit        Matching Pursuit        Thresholding

Change the $L_0$ into $L_1$ and then the problem becomes convex and manageable

$$\min_{\alpha} \|\alpha\|_1$$
$$\text{s.t.}$$
$$\|\mathbf{D}\alpha - y\|_2 \leq \varepsilon$$

Find the support greedily, one element at a time

Multiply y by $\mathbf{D^T}$ and apply shrinkage:
$$\widehat{\alpha} = \mathcal{P}_{\beta}\{\mathbf{D^T}y\}$$

# Difficulties with *Sparseland*

o There are various pursuit algorithms

o Here is an example using the Basis Pursuit ($L_1$):



o Surprising fact: Many of these algorithms are often accompanied by theoretical guarantees for their success, if the unknown is sparse enough

# The Mutual Coherence

o Compute

$$\mathbf{D^T} \quad \mathbf{D} = \mathbf{D^T D}$$

Assume normalized columns

o The Mutual Coherence $\mu(\mathbf{D})$ is the largest off-diagonal entry in absolute value

o We will pose all the theoretical results in this talk using this property, due to its simplicity

o You may have heard of other ways to characterize the dictionary (Restricted Isometry Property - RIP, Exact Recovery Condition - ERC, Babel function, Spark, …)

# Basis-Pursuit Success

**Theorem:** **Given** a noisy signal $\mathrm{y} = \mathbf{D}\alpha + \mathrm{v}$ where $\|\mathrm{v}\|_2 \leq \varepsilon$ and $\alpha$ is sufficiently sparse,

$$\|\boldsymbol{\alpha}\|_0 < \frac{1}{4}\left(1 + \frac{1}{\mu}\right)$$

**then Basis-Pursuit:** $\min_\alpha \|\alpha\|_1 \quad \text{s.t.} \quad \|\mathbf{D}\alpha - \mathrm{y}\|_2 \leq \varepsilon$

**leads to a stable resul**t: $\|\widehat{\alpha} - \alpha\|_2^2 \leq \frac{4\varepsilon^2}{1 - \mu(4\|\alpha\|_0 - 1)}$

Donoho, Elad & Temlyakov ('06)

$\|\mathrm{x}$

$\mathbf{D}$

$\mathcal{M}$

$\alpha$

$+$

$\mathrm{y}$

$\|\mathrm{v}\|_2 \leq \varepsilon$

$\min_\alpha \|\alpha\|_0$
$\text{s.t.}$
$\|\mathbf{D}\alpha - \mathrm{y}\|_2 \leq \varepsilon$

$\widehat{\alpha}$

Comments:
- If $\varepsilon = 0 \rightarrow \widehat{\alpha} = \alpha$
- This is a worst-case analysis – better bounds exist
- Similar theorems exist for many other pursuit algorithms

# Difficulties with *Sparseland*

o Problem 2: Given a family of signals, how do we find the dictionary to represent it well?

o Solution: Learn! Gather a large set of signals (many thousands), and find the dictionary that sparsifies them

o Such algorithms were developed in the past 10 years (e.g., K-SVD), and their performance is surprisingly good

o We will not discuss this matter further in this talk due to lack of time

Michael Elad
The Computer-Science Department
The Technion

# Difficulties with *Sparseland*

o Problem 3: Why is this model suitable to describe various sources? e.g., Is it good for images? Audio? Stocks? …

o General answer: Yes, this model is extremely effective in representing various sources

- Theoretical answer: Clear connection to other models

- Empirical answer: In a large variety of signal and image processing (and later machine learning), this model has been shown to lead to state-of-the-art results

# Difficulties with *Sparseland* ?

o Problem 1: Given an image patch, how can we find its atom decomposition ?

o Problem 2: Given a family of signals, how do we find the dictionary to represent it well?

o Problem 3: Is this model flexible enough to describe various sources? E.g., Is it good for images? audio? …

**ALL ANSWERED POSITIVELY AND CONSTRUCTIVELY**

# This Field has been rapidly GROWING ...

o *Sparseland* has a great success in signal & image processing and machine learning tasks

o In the past 8-9 years, many books were published on this and closely related fields

# Coming Up: A Massive Open Online Course

# *Sparseland* for Image Processing

o When handling images, *Sparseland* is typically deployed on small overlapping patches due to the desire to train the model to fit the data better



o The model assumption is: each patch in the image is believed to have a sparse representation w.r.t. a common local dictionary

o What is the corresponding global model? This brings us to ... the Convolutional Sparse Coding (CSC)

# Multi-Layered Convolutional Sparse Modeling

### Joint work with



Yaniv Romano          Vardan Papyan          Jeremias Sulam

# Convolutional Sparse Coding (CSC)

$m$ filters convolved with their sparse representations

i-th feature-map: An image of the same size as $\mathbf{X}$ holding the sparse representation related to the i-filter

An image with $N$ pixels

$$[\mathbf{X}] = \sum_{i=1}^{m} \mathrm{d_i} * [\Gamma_\mathrm{i}]$$

The i-th filter of small size $n$

# CSC in Matrix Form

o Here is an alternative global sparsity-based model formulation

$$\mathbf{X} = \sum_{i=1}^{m} \mathbf{C}^i \mathbf{\Gamma}^i \ = \ [\mathbf{C}^1 \ \cdots \ \mathbf{C}^m] \begin{bmatrix} \mathbf{\Gamma}^1 \\ \vdots \\ \mathbf{\Gamma}^m \end{bmatrix} = \mathbf{D}\mathbf{\Gamma}$$

o $\mathbf{C}^i \in \mathbb{R}^{N \times N}$ is a banded and Circulant matrix containing a single atom with all of its shifts

$$\mathbf{C}^i =$$

o $\mathbf{\Gamma}^i \in \mathbb{R}^N$ are the corresponding coefficients ordered as column vectors

# The CSC Dictionary

$$[\mathbf{C}^1 \; \mathbf{C}^2 \; \mathbf{C}^3] =$$



$$\mathbf{D}_{\mathrm{L}}$$

$$\mathbf{D} =$$

# Why CSC?



$$\mathbf{X} = \mathbf{D\Gamma}$$

$$\mathbf{R_i X} = \mathbf{\Omega\gamma_i}$$

**stripe-dictionary**

**stripe vector**

Every patch has a sparse representation w.r.t. to the same local dictionary ($\mathbf{\Omega}$) just as assumed for images

# Why CSS?

$$R_{i+1}X$$

$$(2n-1)m$$

$$n$$

$$\boldsymbol{\gamma}_{i+1}$$

stripe-dictionary       stripe vector

$$X = D\Gamma$$

$$R_i X = \Omega \boldsymbol{\gamma}_i$$

$$R_{i+1} X = \Omega \boldsymbol{\gamma}_{i+1}$$

Every patch has a sparse representation w.r.t. to the same local dictionary ($\Omega$) just as assumed for images

# Classical Sparse Theory for CSC ?

$$\min_{\boldsymbol{\Gamma}} \quad \|\boldsymbol{\Gamma}\|_0 \quad \text{s.t.} \|\mathbf{Y} - \mathbf{D}\boldsymbol{\Gamma}\|_2 \leq \varepsilon$$

**Theorem**: **BP is guaranteed to "succeed" .... if** $\|\boldsymbol{\Gamma}\|_0 < \frac{1}{4}\left(1 + \frac{1}{\mu}\right)$

○ Assuming that $m = 2$ and $n = 64$ we have that [Welch, '74]

$$\mu \geq 0.063$$

○ Success of pursuits is guaranteed as long as

$$\|\boldsymbol{\Gamma}\|_0 < \frac{1}{4}\left(1 + \frac{1}{\mu(\mathbf{D})}\right) \leq \frac{1}{2}\left(1 + \frac{1}{0.063}\right) \approx 4.2$$

○ Only few (4) non-zeros GLOBALLY are allowed!!! This is a very pessimistic result!

○ The classic $Sparseland$ Theory does not cover well the CSC model

# Moving to Local Sparsity: Stripes

$\ell_{0,\infty}$ Norm: $\quad \|\mathbf{\Gamma}\|_{0,\infty}^{\mathrm{s}} = \max_i \ \|\mathbf{\gamma}_i\|_0$

$$\min_{\mathbf{\Gamma}} \ \ \|\mathbf{\Gamma}\|_{0,\infty}^{\mathrm{s}} \ \ \mathrm{s.\,t.} \ \ \|\mathbf{Y} - \mathbf{D\Gamma}\|_2 \leq \varepsilon$$

$\|\mathbf{\Gamma}\|_{0,\infty}^{\mathrm{s}}$ is low $\rightarrow$ all $\mathbf{\gamma}_i$ are sparse $\rightarrow$ every patch has a sparse representation over $\mathbf{\Omega}$

**The main question we aim to address is this:**

Can we generalize the vast theory of *Sparseland* to this new notion of local sparsity? For example, could we provide guarantees for success for pursuit algorithms?

$\mathbf{\gamma}_{i+1}$ $\mathbf{\gamma}_i$

$\mathbf{\Gamma}$

# Success of OMP

**Theorem:** If $\mathbf{Y} = \mathbf{D\Gamma} + \mathbf{E}$ where

$$\|\mathbf{\Gamma}\|_{0,\infty}^{s} < \frac{1}{2}\left(1 + \frac{1}{\mu}\right) - \frac{1}{\mu} \cdot \frac{\|\mathbf{E}\|_{2,\infty}^{p}}{|\mathbf{\Gamma}_{min}|}$$

then **OMP** run for $\|\mathbf{\Gamma}\|_{0}$ iterations

1. **Finds the correct support**

2. $\|\mathbf{\Gamma}_{\text{OMP}} - \mathbf{\Gamma}\|_{2}^{2} \leq \dfrac{\|\mathbf{E}\|_{2}^{2}}{1 - (\|\mathbf{\Gamma}\|_{0,\infty}^{s} - 1)\mu}$

Papyan, Sulam & Elad ('17)

This is a much better result – it allows
few non-zeros locally in each stripe, implying a
permitted $O(N)$ non-zeros globally

# Success of the Basis Pursuit

$$\Gamma_{BP} = \min_{\Gamma} \quad \frac{1}{2}\|Y - D\Gamma\|_2^2 + \lambda\|\Gamma\|_1$$

Recent works tackling the convolutional sparse coding problem via BP
[Bristow, Eriksson & Lucey '13]
[Wohlberg '14]
[Kong & Fowlkes '14]
[Bristow & Lucey '14]
[Heide, Heidrich & Wetzstein '15]
[Šorel & Šroubek '16]

Michael Elad
The Computer-Science Department
The Technion

# Success of the Basis Pursuit

$$\Gamma_{\mathrm{BP}} = \min_{\Gamma} \quad \frac{1}{2}\|Y - D\Gamma\|_2^2 + \lambda\|\Gamma\|_1$$

**Theorem:** For $Y = D\Gamma + E$, if $\lambda = 4\|E\|_{2,\infty}^p$ , **if**

$$\|\mathbf{\Gamma}\|_{\mathbf{0},\infty}^{\mathbf{s}} < \frac{\mathbf{1}}{\mathbf{3}}\left(\mathbf{1} + \frac{\mathbf{1}}{\mathbf{\mu(D)}}\right)$$

**then Basis Pursuit performs very-well:**

1. The support of $\Gamma_{\mathrm{BP}}$ is contained in that of $\Gamma$

2. $\|\Gamma_{\mathrm{BP}} - \Gamma\|_\infty \leq 7.5\|E\|_{2,\infty}^p$

3. Every entry greater than $7.5\|E\|_{2,\infty}^p$ is found

4. $\Gamma_{\mathrm{BP}}$ is unique

Papyan, Sulam & Elad ('17)

Michael Elad
The Computer-Science Department
The Technion

# Global Pursuit via Local Processing

o Could we suggest a solution of the global Basis Pursuit using only local (e.g. patch-based) operations ?

o The answer is positive !!

o We define image slices :

$$\mathbf{s_i} \equiv \mathbf{D_L}\alpha_i$$

$$\mathbf{\Gamma}_{\mathrm{BP}} = \min_{\mathbf{\Gamma}} \quad \frac{1}{2}\|\mathbf{Y} - \mathbf{D\Gamma}\|_2^2 + \lambda\|\mathbf{\Gamma}\|_1$$

$$\mathbf{X} = \mathbf{D\Gamma}$$

$$\mathbf{D_L}$$

$$\alpha_i$$

Michael Elad
The Computer-Science Department
The Technion

# Global Pursuit via Local Processing

$$(\mathbf{P}_1^\epsilon): \quad \mathbf{\Gamma}_{BP} = \min_{\mathbf{\Gamma}} \quad \frac{1}{2}\|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_2^2 + \lambda\|\mathbf{\Gamma}\|_1$$

These two are convex & equivalent

Redefine this problem using $s_i$ and $\alpha_i$

$$\min_{\mathbf{\alpha}_i, \mathbf{s}_i} \quad \frac{1}{2}\left\|\mathbf{Y} - \sum_i \mathbf{R}_i^T \mathbf{s}_i\right\|_2^2 + \lambda\sum_i \|\mathbf{\alpha}_i\|_1 \quad \text{s.t.} \quad \{\mathbf{s}_i = \mathbf{D}_L\mathbf{\alpha}_i\}_i$$

Update the $\mathbf{\alpha}_i$ by a local BP

Update the slices $\mathbf{s}_i$ by a simple LS & patch-averaging

If you apply the above two steps only once, you get a known patch-based denoising algorithm

# Global Pursuit via Local Processing

$$(\mathbf{P}_1^\epsilon): \quad \mathbf{\Gamma}_{\mathrm{BP}} = \min_{\mathbf{\Gamma}} \quad \frac{1}{2}\|\mathbf{Y} - \mathbf{D\Gamma}\|_2^2 + \lambda\|\mathbf{\Gamma}\|_1$$

These two are convex & equivalent

Redefine this problem using $s_i$ and $\alpha_i$

$$\min_{\mathbf{\alpha_i}, \mathbf{s_i}} \quad \frac{1}{2}\left\|\mathbf{Y} - \sum_i \mathbf{R}_i^T \mathbf{s_i}\right\|_2^2 + \lambda \sum_i \|\mathbf{\alpha_i}\|_1 \quad \text{s.t.} \quad \{\mathbf{s_i} = \mathbf{D}_L \mathbf{\alpha_i}\}_i$$

Update the $\alpha_i$ by a le...

Upd... ...lices $s_i$ ...LS & ...sing

If you ap... ...steps only once, you get a kn... ...tch-based denoising algorithm

This algorithm operates locally while guaranteeing to solve the global problem

# Two Comments About this Scheme

## We work with Slices and not Patches

Patches extracted from natural images, and their corresponding slices. Observe how the slices are far simpler, and contained by their corresponding patches

Patches
Slices

## The Proposed Scheme can be used for Dictionary ($\mathbf{D}_L$) Learning

Slice-based DL algorithm using standard patch-based tools, leading to a faster and simpler method, compared to existing methods

[Wohlberg, 2016]　　　　　　Ours

# Multi-Layered Convolutional Sparse Modeling

Michael Elad
The Computer-Science Department
The Technion

# CSC and CNN

o There is a rough analogy between CSC and CNN:

- Convolutional structure
- Data driven models
- ReLU is a sparsifying operator

o We shall now propose a principled way to analyze CNN

o But first, a brief review of CNN…

# CNN



**Y**     ReLU     ReLU

[LeCun, Bottou, Bengio and Haffner '98]
[Krizhevsky, Sutskever & Hinton '12]
[Simonyan & Zisserman '14]
[He, Zhang, Ren & Sun '15]

$$\text{ReLU}(z) = \max(\text{Thr}, z)$$

# CNN



$\mathbf{Y}$        $m_1$        $m_2$

$\mathbf{W_1}$        $\mathbf{W_2}$

[LeCun, Bottou, Bengio and Haffner '98]
[Krizhevsky, Sutskever & Hinton '12]
[Simonyan & Zisserman '14]
[He, Zhang, Ren & Sun '15]

$$\text{ReLU}(z) = \max(\text{Thr}, z)$$

# Mathematically...

$$f(\mathbf{Y}) = \text{ReLU}\left(\mathbf{b}_2 + \mathbf{W}_2^{\mathbf{T}}\, \text{ReLU}\left(\mathbf{b}_1 + \mathbf{W}_1^{\mathbf{T}}\mathbf{Y}\right)\right)$$

$\mathbf{Z}_2 \in \mathbb{R}^{Nm_2}$ $\quad \mathbf{b}_2 \in \mathbb{R}^{Nm_2}$ $\quad \mathbf{W}_2^{\mathbf{T}} \in \mathbb{R}^{Nm_2 \times Nm_1}$

$\mathbf{b}_1 \in \mathbb{R}^{Nm_1}$ $\qquad \mathbf{W}_1^{\mathbf{T}} \in \mathbb{R}^{Nm_1 \times N}$



$$\left\|\quad\right\| \boxminus \text{ReLU}\left\{ \left\|\quad\right\| \oplus \begin{bmatrix} n_1 m_1 \\ m_2 \\ \\ \\ \\ m_1 \\ \\ \\ \\ \end{bmatrix} \right\} \otimes \text{ReLU}\left\{ \left\|\quad\right\| \oplus \begin{bmatrix} n_0 \\ m_1 \\ \\ \\ \\ \end{bmatrix} \otimes \left\|\quad\right\| \right\}$$

$\mathbf{Y} \in \mathbb{R}^N$

# From CSC to Multi-Layered CSC

$\mathbf{X} \in \mathbb{R}^N$    $\mathbf{D}_1 \in \mathbb{R}^{N \times Nm_1}$    $\mathbf{\Gamma}_1 \in \mathbb{R}^{Nm_1}$

We propose to impose the same structure on the representations themselves

Convolutional sparsity (CSC) assumes an inherent structure is present in natural signals

$\mathbf{\Gamma}_1 \in \mathbb{R}^{Nm_1}$    $\mathbf{D}_2 \in \mathbb{R}^{Nm_1 \times Nm_2}$    $\mathbf{\Gamma}_2 \in \mathbb{R}^{Nm_2}$

**Multi-Layer CSC (ML-CSC)**

# Intuition: From Atoms to Molecules

$$\mathbf{X} \in \mathbb{R}^N \qquad \mathbf{D}_1 \in \mathbb{R}^{N \times Nm_1} \qquad \mathbf{\Gamma}_1 \in \mathbb{R}^{Nm_1} \quad \mathbf{D}_2 \in \mathbb{R}^{Nm_1 \times Nm_2} \qquad \mathbf{\Gamma}_2 \in \mathbb{R}^{Nm_2}$$

o We can chain the all the dictionaries into one effective dictionary
$$\mathbf{D}_{\text{eff}} = \mathbf{D}_1 \mathbf{D}_2 \mathbf{D}_3 \cdots \mathbf{D}_K \;\rightarrow\; \mathbf{x} = \mathbf{D}_{\text{eff}} \, \mathbf{\Gamma}_K$$

o This is a special $Sparseland$ (indeed, a CSC) model

o However:

$$\mathbf{\Gamma}_1 \in \mathbb{R}^{Nm_1}$$

- A key property in this model: sparsity of the intermediate representations

- The effective atoms: atoms $\rightarrow$ molecules $\rightarrow$ cells $\rightarrow$ tissue $\rightarrow$ body-parts ...

# A Small Taste: Model Training (MNIST)

MNIST Dictionary:
- $D_1$: 32 filters of [...] of 2 (dense)
- $D_2$: 128 fil[...] f 1 - 99.09 % sparse
- D3: 102[...] parse

$\mathbf{D}_1$ (7×7)

$\mathbf{D}_1\mathbf{D}_2$ (15×15)

$\mathbf{D}_1\mathbf{D}_2\mathbf{D}_3$ (28×28)

# A Small Taste: Model Training (CiFAR)



$\mathbf{D}_1$ (5×5×3)    $\mathbf{D}_1\mathbf{D}_2$ (13×13)    $\mathbf{D}_1\mathbf{D}_2\mathbf{D}_3$ (32×32)

The Technion

# ML-CSC: Pursuit

o **Deep–Coding Problem** ($\mathbf{DCP}_\lambda$) (dictionaries are known):

$$\left\{\begin{array}{ll} \mathbf{X} = \mathbf{D}_1 \boldsymbol{\Gamma}_1 & \|\boldsymbol{\Gamma}_1\|_{0,\infty}^{\mathrm{s}} \leq \lambda_1 \\ \boldsymbol{\Gamma}_1 = \mathbf{D}_2 \boldsymbol{\Gamma}_2 & \|\boldsymbol{\Gamma}_2\|_{0,\infty}^{\mathrm{s}} \leq \lambda_2 \\ \quad\vdots & \quad\vdots \\ \boldsymbol{\Gamma}_{K-1} = \mathbf{D}_K \boldsymbol{\Gamma}_K & \|\boldsymbol{\Gamma}_K\|_{0,\infty}^{\mathrm{s}} \leq \lambda_K \end{array}\right\}$$

o Or, more realistically for noisy signals,

$$\mathrm{Find} \ \ \{\boldsymbol{\Gamma}_j\}_{j=1}^{K} \quad s.t. \left\{\begin{array}{ll} \|\mathbf{Y} - \mathbf{D}_1 \boldsymbol{\Gamma}_1\|_2 \leq \mathcal{E} & \|\boldsymbol{\Gamma}_1\|_{0,\infty}^{\mathrm{s}} \leq \lambda_1 \\ \boldsymbol{\Gamma}_1 = \mathbf{D}_2 \boldsymbol{\Gamma}_2 & \|\boldsymbol{\Gamma}_2\|_{0,\infty}^{\mathrm{s}} \leq \lambda_2 \\ \quad\vdots & \quad\vdots \\ \boldsymbol{\Gamma}_{K-1} = \mathbf{D}_K \boldsymbol{\Gamma}_K & \|\boldsymbol{\Gamma}_K\|_{0,\infty}^{\mathrm{s}} \leq \lambda_K \end{array}\right\}$$

# A Small Taste: Pursuit



Y

X

$\Gamma_1$
94.51 % sparse
(213 nnz)

$\Gamma_2$
99.52% sparse
(30 nnz)

$\Gamma_3$
99.51% sparse
(5 nnz)

$$x = \mathbf{D}_1 \Gamma_1$$

$$x = \mathbf{D}_1 \mathbf{D}_2 \Gamma_2$$

$$x = \mathbf{D}_1 \mathbf{D}_2 \mathbf{D}_3 \Gamma_3$$

# ML-CSC: The Simplest Pursuit

*Keep it simple!*

The simplest pursuit algorithm (single-layer case) is the THR algorithm, which operates on a given input signal $\mathbf{Y}$ by:

$$\mathbf{Y} = \mathbf{D\Gamma} + \mathbf{E}$$

and $\mathbf{\Gamma}$ is sparse

$$\hat{\mathbf{\Gamma}} = \mathcal{P}_\beta(\mathbf{D}^{\mathrm{T}}\mathbf{Y})$$



Legend:
- $\mathcal{H}_\beta(z)$ - Hard
- $\mathcal{S}_\beta(z)$ - Soft
- $\mathcal{S}_\beta^+(z)$ - Soft Nonnegative

# Consider this for Solving the DCP

o Layered thresholding (LT):

Estimate $\boldsymbol{\Gamma}_1$ via the THR algorithm

$$\widehat{\boldsymbol{\Gamma}}_2 = \mathcal{P}_{\beta_2}\left(\mathbf{D}_2^{\mathrm{T}}\,\mathcal{P}_{\beta_1}\left(\mathbf{D}_1^{\mathrm{T}}\mathbf{Y}\right)\right)$$

Estimate $\boldsymbol{\Gamma}_2$ via the THR algorithm

$$\left(\mathbf{DCP}_\lambda^{\mathcal{E}}\right):\ \mathrm{Find}\ \ \{\boldsymbol{\Gamma}_j\}_{j=1}^{\mathrm{K}}\ \ \ s.t.$$

$$\left\{\begin{array}{ll} \|\mathbf{Y} - \mathbf{D}_1\boldsymbol{\Gamma}_1\|_2 \le \mathcal{E} & \|\boldsymbol{\Gamma}_1\|_{0,\infty}^{\mathrm{s}} \le \lambda_1 \\ \boldsymbol{\Gamma}_1 = \mathbf{D}_2\boldsymbol{\Gamma}_2 & \|\boldsymbol{\Gamma}_2\|_{0,\infty}^{\mathrm{s}} \le \lambda_2 \\ \vdots & \vdots \\ \boldsymbol{\Gamma}_{\mathrm{K}-1} = \mathbf{D}_{\mathrm{K}}\boldsymbol{\Gamma}_{\mathrm{K}} & \|\boldsymbol{\Gamma}_{\mathrm{K}}\|_{0,\infty}^{\mathrm{s}} \le \lambda_{\mathrm{K}} \end{array}\right\}$$

o Now let's take a look at how Conv. Neural Network operates:

$$f(\mathbf{Y}) = \mathrm{ReLU}\left(\mathbf{b}_2 + \mathbf{W}_2^{\mathrm{T}}\,\mathrm{ReLU}\left(\mathbf{b}_1 + \mathbf{W}_1^{\mathrm{T}}\mathbf{Y}\right)\right)$$

> The layered (soft nonnegative) thresholding and the CNN forward pass algorithm are the very same thing !!!

# Theoretical Path



$$X = D_1\Gamma_1$$
$$\Gamma_1 = D_2\Gamma_2$$
$$\vdots$$
$$\Gamma_{K-1} = D_K\Gamma_K$$

$\Gamma_i$ is $L_{0,\infty}$ sparse

$\mathcal{M}$

$X$

$Y$

$\mathcal{A}$

$(DCP_\lambda^\mathcal{E})$

Layered THR
(Forward Pass)

Maybe other?

$\{\hat{\Gamma}_i\}_{i=1}^K$

Armed with this view of a generative source model, we may ask new and daring questions

# Theoretical Path: Possible Questions

o Having established the importance of the ML-CSC model and its associated pursuit, the DCP problem, we now turn to its analysis

o The main questions we aim to address:

> I. Stability of the solution obtained via the hard layered THR algorithm (forward pass) ?
>
> II. Limitations of this (very simple) algorithm and alternative pursuit?

### … and here are questions we will not touch today:

> III. Algorithms for training the dictionaries $\{\mathbf{D}_i\}_{i=1}^{K}$ vs. CNN ?
>
> IV. New insights on how to operate on signals via CNN ?

# Success of the Layered-THR

**Theorem:** If $\|\mathbf{\Gamma}_i\|_{0,\infty}^{s} < \frac{1}{2}\left(1 + \frac{1}{\mu(\mathbf{D}_i)} \cdot \frac{|\mathbf{\Gamma}_i^{min}|}{|\mathbf{\Gamma}_i^{max}|}\right) - \frac{1}{\mu(\mathbf{D}_i)} \cdot \frac{\varepsilon_L^{i-1}}{|\mathbf{\Gamma}_i^{max}|}$

then the **Layered Hard THR** (with the proper thresholds)
**finds the correct supports** and $\left\|\mathbf{\Gamma}_i^{LT} - \mathbf{\Gamma}_i\right\|_{2,\infty}^{p} \leq \varepsilon_L^i,$ where

we have defined $\varepsilon_L^0 = \|\mathbf{E}\|_{2,\infty}^{p}$ and

$$\varepsilon_L^i = \sqrt{\|\mathbf{\Gamma}_i\|_{0,\infty}^{p} \cdot \left(\varepsilon_L^{i-1} + \mu(\mathbf{D}_i)\left(\|\mathbf{\Gamma}_i\|_{0,\infty}^{s} - 1\right)|\mathbf{\Gamma}_i^{max}|\right)}$$

Papyan, Romano & Elad ('17)

The stability of the forward pass is guaranteed if the underlying representations are **locally** sparse and the noise is **locally** bounded

Problems:
1. Contrast
2. Error growth
3. Error even if no noise

Michael Elad
The Computer-Science Department
The Technion

# Layered Basis Pursuit (BP)

o We chose the Thresholding algorithm due to its simplicity, but we do know that there are better pursuit methods – how about using them?

o Lets use the Basis Pursuit instead …

$$\left(\mathbf{DCP}_\lambda^\mathcal{E}\right): \text{ Find } \left\{\boldsymbol{\Gamma}_j\right\}_{j=1}^{K} \quad s.t.$$

$$\left\{ \begin{array}{ll} \|\mathbf{Y} - \mathbf{D}_1\boldsymbol{\Gamma}_1\|_2 \leq \mathcal{E} & \|\boldsymbol{\Gamma}_1\|_{0,\infty}^s \leq \lambda_1 \\ \boldsymbol{\Gamma}_1 = \mathbf{D}_2\boldsymbol{\Gamma}_2 & \|\boldsymbol{\Gamma}_2\|_{0,\infty}^s \leq \lambda_2 \\ \vdots & \vdots \\ \boldsymbol{\Gamma}_{K-1} = \mathbf{D}_K\boldsymbol{\Gamma}_K & \|\boldsymbol{\Gamma}_K\|_{0,\infty}^s \leq \lambda_K \end{array} \right\}$$

$$\boldsymbol{\Gamma}_1^{\text{LBP}} = \min_{\boldsymbol{\Gamma}_1} \frac{1}{2} \|\mathbf{Y} - \mathbf{D}_1\boldsymbol{\Gamma}_1\|_2^2 + \lambda_1 \|\boldsymbol{\Gamma}_1\|_1$$

$$\boldsymbol{\Gamma}_2^{\text{LBP}} = \min_{\boldsymbol{\Gamma}_2} \frac{1}{2} \left\|\boldsymbol{\Gamma}_1^{\text{LBP}} - \mathbf{D}_2\boldsymbol{\Gamma}_2\right\|_2^2 + \lambda_2 \|\boldsymbol{\Gamma}_2\|_1$$

> Deconvolutional networks
> [Zeiler, Krishnan, Taylor & Fergus '10]

# Success of the Layered BP

**Theorem:** **Assuming that** $\|\Gamma_i\|_{0,\infty}^s < \frac{1}{3}\left(1 + \frac{1}{\mu(D_i)}\right)$

then the Basis Pursuit performs very well:

1. The support of $\Gamma_i^{LBP}$ is contained in that of $\Gamma_i$

2. The error is bounded: $\left\|\Gamma_i^{LBP} - \Gamma_i\right\|_{2,\infty}^p \le \varepsilon_L^i$, where

$$\varepsilon_L^i = 7.5^i \|E\|_{2,\infty}^p \prod_{j=1}^i \sqrt{\|\Gamma_j\|_{0,\infty}^p}$$

3. Every entry in $\Gamma_i$ greater than

$$\varepsilon_L^i \Big/ \sqrt{\|\Gamma_i\|_{0,\infty}^p} \text{ will be found}$$

Papyan, Romano & Elad ('17)

Problems:
1. ~~Contrast~~
2. Error growth
3. ~~Error even if no noise~~

Michael Elad
The Computer-Science Department
The Technion

# Layered Iterative Thresholding

Layered BP:  $\mathbf{\Gamma}_j^{LBP} = \min_{\mathbf{\Gamma}_j} \frac{1}{2} \left\| \mathbf{\Gamma}_{j-1}^{LBP} - \mathbf{D}_j \mathbf{\Gamma}_j \right\|_2^2 + \xi_j \left\| \mathbf{\Gamma}_j \right\|_1$   $j$

Layered Iterative Soft-Thresholding:

$t$   $\mathbf{\Gamma}_j^t = \mathcal{S}_{\xi_j/c_j} \left( \mathbf{\Gamma}_j^{t-1} + \mathbf{D}_j^T (\widehat{\mathbf{\Gamma}}_{j-1} - \mathbf{D}_j \mathbf{\Gamma}_j^{t-1}) \right)$   $j$

Note that our suggestion implies that groups of layers share the same dictionaries

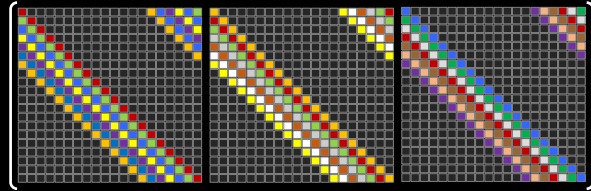Can be seen as a very deep recurrent neural network
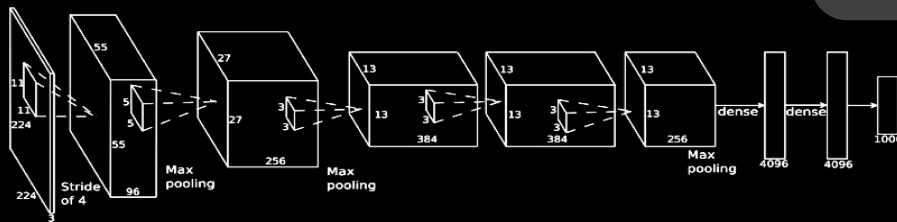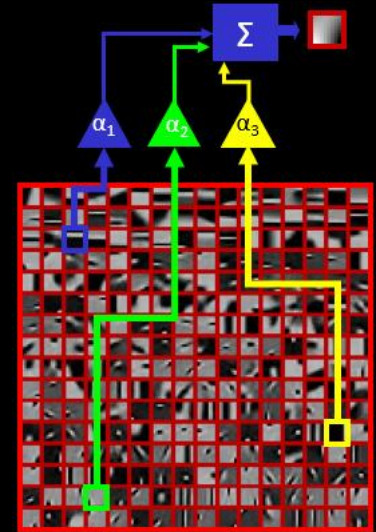[Gregor & LeCun '10]

# Time to Conclude

# This Talk

**The desire to model data**

*Sparseland*

Novel View of Convolutional Sparse Coding

Multi-Layer Convolutional Sparse Coding

A novel interpretation and theoretical understanding of CNN



$$\mathcal{M}$$

$$X = D_1 \Gamma_1$$
$$\Gamma_1 = D_2 \Gamma_2$$
$$\vdots$$
$$\Gamma_{K-1} = D_K \Gamma_K$$

$$\mathbf{X}$$

$\Gamma_i$ is $L_{0,\infty}$ sparse

Michael Elad
The Computer-Science Department
The Technion

# This Talk

**Take Home Message 1:**
Generative modeling of data sources enables algorithm development *along* with theoretically analyzing algorithms' performance
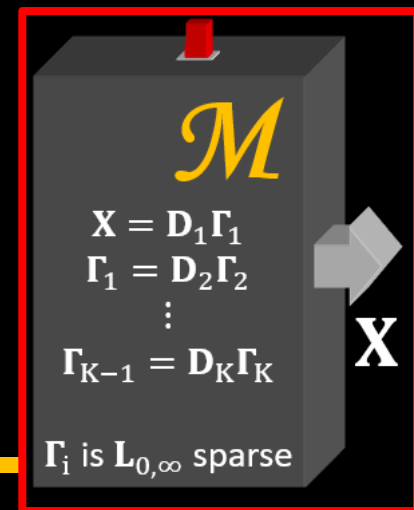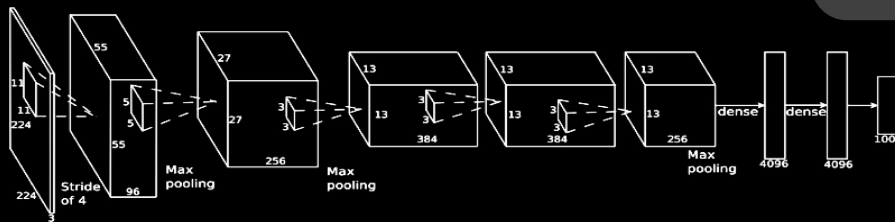
*Sparseland* ← The desire to model data

↓

Novel View of Convolutional Sparse Coding

↓

A novel interpretation and theoretical understanding of CNN ← Multi-Layer Convolutional Sparse Coding

# This Talk

**Sparseland** ← The desire to model data

Novel View of Convolutional Sparse Coding

Multi-Layer Convolutional Sparse Coding

A novel interpretation and theoretical understanding of CNN

**Take Home Message 2:** The Multi-Layer Convolutional Sparse Coding model could be a new platform for understanding and developing deep-learning solutions

More on these (including these slides and the relevant papers) can be found in http://www.cs.technion.ac.il/~elad

Michael Elad
The Computer-Science Department
The Technion