Welcome to Sparseland

Sparse & Redundant Representations and their Applications in Signal and Image Processing

Michael Elad

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



2017 Summer School on Signal Processing Meets Deep Learning September 4-8, 2017



The research leading to these results has been received funding from the European union's Seventh Framework Program (FP/2007-2013) ERC grant Agreement ERC-SPARSE- 320649



This Lecture







... is Built of Three Parts:

Part 1: A General Introduction to *Sparseland*: Presenting this Model and its Importance

Part 2: Diving In: Theory and Applications in Sparse Representations

Part 3: Relation to Deep-Learning: A Tale of three models: $Sparseland \rightarrow CSC \rightarrow CNN$



Welcome to Sparseland

Part 1: A General Introduction to Sparseland

Michael Elad

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



2017 Summer School on Signal Processing Meets Deep Learning September 4-8, 2017



The research leading to these results has been received funding from the European union's Seventh Framework Program (FP/2007-2013) ERC grant Agreement ERC-SPARSE- 320649



What This Field is All About ?



Michael Elad The Computer-Science Department The Technion

The Answer is Not Trivial

Depends whom you ask, as the researchers in this field come from various disciplines:

- Mathematics
- Applied Mathematics
- Statistics
- Signal & Image Processing: CS, EE, Bio-medical, ...
- Computer-Science Theory
- Machine-Learning
- Physics (optics)
- Geo-Physics
- Physics Astronomy
- Psychology (neuroscience)
- o ...



My Answer (For Now)

A New Transform for Signals

- We are all well-aware of the idea of transforming a signal and changing its representation
- We apply a transform to gain something efficiency, simplicity of the subsequent processing, speed, ...
- There are many known transforms: Fourier, DCT, Hadamard, Wavelets and its descendants, and more
- Our message: There is a new transform in town, based on sparse and redundant representations



Transforms – The General Picture





Michael Elad The Computer-Science Department The Technion

Redundancy?

- In a redundant transform, the representation vector is longer (m>n)
- This can still be done while preserving the linearity of the transform:





Sparse & Redundant Representation

n

m

zeros

- We shall keep the linearity of the inverse-transform
- As for the forward (computing $\underline{\alpha}$ from \underline{x}), there are infinitely many possible solutions
- We shall seek the sparsest of all soluti
 Sounds ... Boring !!!!
- This mal non-line Who cares about
- The field a new transform? ations is all about demonstrating electry this transform, solving various theoretical and numerical issues related to it, and showing how to use it in practice



α

What This Field is All About ?

Take 2



Michael Elad The Computer-Science Department The Technion

Lets Take a Wider Perspective Matrix Data





1600

1000

Model?



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





An Example: JPEG and DCT



How & why does it works?







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



Research in Signal/Image Processing





Again: What This Field is all About?

A Data Model and its Use

- Almost any task in data processing requires a model true for denoising, deblurring, super-resolution, inpainting, compression, anomaly-detection, sampling, recognition, ...
- There is a new model in town sparse and redundant representation – we will call it

Sparseland

This is a flexible model that can adjust to the signal



A New Emerging Model



A Closer Look at the *Sparseland* Model



Michael Elad The Computer-Science Department The Technion

The Sparseland Model

- Task: model image patches of size 8×8 pixels
- We assume that a dictionary of such image patches is given, containing 256 atom images
- 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)





Chemistry of Data

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

- Our dictionary stands for the Periodic Table containing all the elements
- Our model follows a similar rationale:
 Every molecule is built of few elements







Model vs. Transform ?

- The relation between the signal <u>x</u> and its representation
 <u>α</u> is the following linear system, just as described earlier
- We shall be interested in seeking sparse solutions to this system when deploying the sparse and redundant representation model
- \odot This is EXACTLY the transform we discussed earlier

Bottom Line: The transform and the model we described above are the same thing, and their impact on signal/image processing is profound and worth studying

n

m



- Problem 1: Given an image patch, how
 can we find its atom decomposition?
- A simple example:
 - There are 2000 atoms in the dictionary
 - The signal is known to be built of 15 atoms

$$\begin{pmatrix} 2000\\ 15 \end{pmatrix} \approx 2.4e + 37 \quad \text{possibilities}$$

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





- Various algorithms exist. Their theoretical analysis guarantees their success if the solution is sparse enough
- Here is an example the Iterative Reweighted LS:







- Problem 2: Given a family of signals, how do we find the dictionary to represent it well?
- Solution: Learn! Gather a large set of signals (many thousands), and find the dictionary that sparsifies them
- Many such algorithms were developed in the past 10 years (e.g., K-SVD), and their performance is surprisingly good
- This is only the beginning of a new era in signal processing ...





- Problem 3: Is this model flexible enough to describe various sources? e.g., Is it good for images? Audio? Stocks? ...
- General answer: Yes, this model is extremely effective in representing various sources
 - Theoretical answer: Relation to other known models
 - Empirical answer: we will see in this course, several image processing applications, where this model leads to the best known results (benchmark tests)





Problem 1: Given an image JOW \bigcirc can we find its at

ao we ANSNERED ent it while the dict and it while the dict and it while the second seco CONSTRUCTIVELY us sources? images? audio? ...





Probl

 \bigcirc

Who Works on This ?



Michael Elad The Computer-Science Department The Technion

Who is Working in This Field?



Donoho, Candes – Stanford



Tropp – CalTech



Baraniuk, W. Yin – Rice Texas



Gilbert, Vershynin, Plan– U-Michigan



Gribonval, Fuchs – INRIA France



Starck – CEA – France





Rao, Delgado – UC San-Diego





Davies – Edinburgh UK







Nowak, Willet - Wisconsin



Coifman – Yale



Romberg – GaTech



Lustig, Wainwright – Berkeley



Duke Sapiro, Daubachies – Duke



Friedlander – UBC Canada

Tarokh – Harvard



Elad, Zibulevsky, Bruckstein, Eldar, Segev, Mendelson – Technion



David L. Donoho

- An extremely talented mathematician and statistician from the Stanford Statistics Department
- He is among the few who founded this field sparse and redundant representations and its spin-off topic of compressed sensing
- In 2013 he won the
 Shaw prize ("the
 Nobel of the east")



Announcement and Citation

Announcement

The Shaw Prize in <u>Mathematical Sciences 2013</u>

is awarded to

David L Donoho

for his profound contributions to modern mathematical statistics and in particular the development of optimal algorithms for statistical estimation in the presence of noise and of efficient techniques for sparse representation and recovery in large data-sets.

28 May 2013 Hong Kong



This Field is rapidly GROWING ...

 Searching ISI-Web-of-Science (December 26th 2016): Topic= ((spars* and (represent* or approx* or solution or estimation) and (dictionary or pursuit or convex)) or (compres* and sens* and spars*))
 Ied to 6354 papers (it was ~4000 papers a year ago)

Here is how
 they spread
 over time
 (with ~146178
 citations):







Which Countries?

ield: Countries/Territories	Record Count	% of 6354	Bar Chart
PEOPLES R CHINA	2573	40.494 %	
USA	2171	34.167 %	
FRANCE	416	6.547 %	
ENGLAND	286	4.501 %	
CANADA	272	4.281 %	1
GERMANY	257	4.045 %	
ISRAEL	189	2.975 %	1.00
AUSTRALIA	170	2.675 %	1.00
SOUTH KOREA	160	2.518 %	1.00
ITALY	154	2.424 %	1.00
SWITZERLAND	144	2.266 %	1.00
JAPAN	141	2.219 %	1.00
SINGAPORE	137	2.156 %	1.00
INDIA	115	1.810 %	1.00
IRAN	115	1.810 %	1.00
SPAIN	84	1.322 %	1.00
BELGIUM	73	1.149 %	1.00
SCOTLAND	72	1.133 %	1.00
TURKEY	72	1.133 %	1.00
TAIWAN	71	1.117 %	1.00
SWEDEN	64	1.007 %	1
AUSTRIA	58	0.913 %	1
DENMARK	55	0.866 %	1
NETHERLANDS	55	0.866 %	1
FINLAND	51	0.803 %	1



Books in this field

- The following book was published in 2010, and it has served as the textbook for my advanced course
- Since then, it has been adopted by other courses worldwide (Stanford, Duke, Oxford, IISc, UCSD, ...)
- In the past 8-9 years, many other books were published in this and closely related fields





A Massive Open Online Course : Coming Up

Search:



Courses • Programs • Schools & Partners About •

Q

Sign In Register

Le Israel X

Sparse Representations in Signal and Image Processing

Learn the theory, tools and algorithms of sparse representations and their impact on signal and image processing.

Start the Professional Certificate Program



Courses in the Professional Certificate Program

Sparse Representations in Signal and Image Processing: Fundamentals Learn about the field of sparse representations by understanding its fundamental theoretical and algorithmic foundations. Learn more

Instructors

Sparse Representations in Image Processing: From Theory to Practice Learn about the deployment of the sparse representation model to signal and image processing. Learn more

Starts on October 25, 2017

Enroll Now

I would like to receive email from IsraelX and learn about other offerings related to Sparse Representations in Signal and Image Processing: Fundamentals.

Starts on February 28, 2018

Enroll Now

I would like to receive email from IsraelX and learn about other offerings related to Sparse Representations in Image Processing: From Theory to Practice.

Michael Elad The Computer-Science Department The Technion

Yaniv Romano

Michael Elad

Several Examples: Applications Leveraging this Model

Michael Elad The Computer-Science Department The Technion
Image Separation [Starck, Elad, & Donoho (`04)]





Inpainting [Starck, Elad, and Donoho ('05)]



Outcome

Source



Denoising [Mairal, Elad & Sapiro, ('06)]



Original



Noisy (20.43dB)



Result (30.75dB)



Poisson Denoising [Giryes & Elad ('14)]





Deblurring [Elad, Zibulevsky and Matalon, ('07)]



original

Measured

Restored (right): ISNR=7.0322 dB



Blind Deblurring [Shao and Elad ('14)]





Michael Elad The Computer-Science Department The Technion

Inpainting (Again!) [Mairal, Elad & Sapiro, ('06)]





To Summarize

An effective (yet simple) model for signals/images is key in getting better algorithms for various applications

Which model to choose? Sparse and redundant representations & trained dictionaries are drawing a considerable attention in recent years, due to the elegant theory and the impressive applications

Lets take a closer look at this model, clarify what we mean by the theory behind it, and show several of the algorithms it leads to for handling image processing tasks

So, what next?





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



Welcome to *Sparseland* Part 2: Diving In

Michael Elad

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



2017 Summer School on Signal Processing Meets Deep Learning September 4-8, 2017



The research leading to these results has been received funding from the European union's Seventh Framework Program (FP/2007-2013) ERC grant Agreement ERC-SPARSE- 320649



Agenda

Part I – Denoising by Sparse & Redundant Representations



Part II – Theoretical & Numerical Foundations

Part III – Dictionary Learning & The K-SVD Algorithm



Part V – Summary & Conclusions



Part IV – Back to Denoising ... and Beyond – handling stills and video denoising & inpainting, demosaicing, super-res., and compression

In this part we will show that

Sparsity and Redundancy are valuable and wellfounded tools for modeling data.

When used in image processing, they lead to state-of-the-art results.



Denoising by Sparse & Redundant Representations



Noise Removal?

Our story begins with image denoising ...



- Important: (i) Practical application; (ii) A convenient platform (being the simplest inverse problem) for testing basic ideas in image processing, and (iii) If you can denoise, you can do much more see Peyman's talk
- Many Considered Directions: Partial differential equations, Statistical estimators, Adaptive filters, Inverse problems & regularization, Wavelets, Example-based techniques, Sparse representations, ...



Denoising by Energy Minimization

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

$$f(\underline{x}) = \frac{1}{2} \|\underline{x} - \underline{y}\|_{2}^{2} + G(\underline{x})$$
rements
Relation to
Prior or regularization

 This is in-fact a Bayesian point of view, adopting the Maximum-A-posteriori Probability (MAP) estimation.

neasurements

 Clearly, the wisdom in such an approach is within the choice of the prior – modeling the images of interest.



Thomas Bayes 1702 - 1761



y : Given measu

<u>x</u> : Unknown to be recovered

The Evolution of G(<u>x</u>)

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

$$G(\underline{x}) = \lambda \|\underline{x}\|_{2}^{2} \quad G(\underline{x}) = \lambda \|\underline{L}\underline{x}\|_{2}^{2} \qquad G(\underline{x}) = \lambda \|\underline{L}\underline{x}\|_{w}^{2} \qquad G(\underline{x}) = \lambda \rho \{\underline{L}\underline{x}\}$$

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

$$G(\underline{x}) = \lambda \|\nabla\underline{x}\|_{1} \qquad G(\underline{x}) = \lambda \|W\underline{x}\|_{1}$$

$$G(\underline{x}) = \lambda \|W\underline{x}\|_{1} \qquad G(\underline{x}) = \lambda \|W\underline{x}\|_{1}$$

$$G(\underline{x}) = \lambda \|\nabla\underline{x}\|_{1} \qquad G(\underline{x}) = \lambda \|W\underline{x}\|_{1}$$

$$G(\underline{x}) = \lambda \|W\underline{x}\|_{1} \qquad G(\underline{x}) = \lambda \|W\underline{x}\|_{1}$$

$$G(\underline{x}) = \lambda \|G(\underline{x}) = \lambda \|W\underline{x}\|_{1}$$

$$G(\underline{x}) = \lambda \|G(\underline{x}) = \lambda \|G(\underline{x$$



Sparse Modeling of Signals





Sparseland Signals are Special



- Simple: Every generated signal is built as a linear combination of <u>few</u> atoms from our dictionary D
- Rich: A general model: the
 obtained signals are a union of many low-dimensional Gaussians
- Familiar: We have been using this model in other context for a while now (wavelet, JPEG2000, ...)



Sparse & Redundant Rep. Modeling?





Back to the MAP Energy Function

X

- $\circ \ L_0 \text{ norm is effectively} \\ \text{ counting the number of} \\ \text{ non-zeros in } \underline{\alpha}$
- The vector <u>α</u> is the representation (sparse/redundant) of the desired signal <u>x</u>
- The core idea: One cannot find a small set of atoms to represent the noise and thus it is "thrown" to the residual.
 → We obtain an effective projection of the noise onto a very low-dimensional space, thus getting denoising effect



Wait! There are Some Issues

Numerical Problems: How should we solve or approximate the solution of the problem

$$\min_{\underline{\alpha}} \left\| \mathbf{D}\underline{\alpha} - \underline{y} \right\|_{2}^{2} \text{ s.t. } \left\| \underline{\alpha} \right\|_{0} \leq L$$

$$\min_{\alpha} \|\underline{\alpha}\|_{0} \quad \text{s.t.} \quad \|\mathbf{D}\underline{\alpha} - \underline{y}\|_{2}^{2} \leq \varepsilon^{2}$$

$$\min_{\underline{\alpha}} \lambda \left\|\underline{\alpha}\right\|_{0} + \left\|\mathbf{D}\underline{\alpha} - \underline{y}\right\|_{2}^{2}$$

- Theoretical Problems: Is there a unique sparse representation? If we are to approximate the solution somehow, how close will we get?
- Practical Problems: What dictionary D should we use, such that all this leads to effective denoising? Will all this work in applications?



To Summarize So Far ...

Image denoising (and many other problems in image processing) requires a model for the desired image



We proposed a model for signals/images based on sparse and redundant representations

There are some issues:

- 1. Theoretical
- 2. How to approximate?
- 3. What about **D**?





Theoretical & Numerical Foundations



Lets Start with the Noiseless Problem

Suppose we build a signal by the relation

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

We aim to find the signal's representation:

$$\underline{\hat{\alpha}} = \operatorname{Arg\,Min}_{\underline{\alpha}} \|\underline{\alpha}\|_{0} \quad \text{s.t.} \quad \underline{\mathbf{X}} = \mathbf{D}\underline{\alpha}$$

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

It might happen that eventually $\|\hat{\alpha}\|_{0} < \|\alpha\|_{0}$.



Uniqueness

Known



Definition: Given a matrix **D**, σ=Spark{**D**} is the smallest number of columns that are linearly dependent

Donoho & E. ('02)

Example:

In tensor decomposition, Kruskal defined something similar already in 1989



Uniqueness Rule

Suppose this problem has been solved somehow

$$\underline{\hat{\alpha}} = \operatorname{ArgMin}_{\underline{\alpha}} \|\underline{\alpha}\|_{0} \quad \text{s.t.} \quad \underline{\mathbf{X}} = \mathbf{D}\underline{\alpha}$$

UniquenessIf we found a representation that satisfyDonoho & E. ('02) $\|\hat{\underline{\alpha}}\|_0 < \frac{\sigma}{2}$ Then necessarily it is unique (the sparsest).

This result implies that if \mathcal{M} generates signals using "sparse enough" $\underline{\alpha}$, the solution of the above will find it exactly



Our Goal





Lets Approximate

$$\min_{\underline{\alpha}} \|\underline{\alpha}\|_{0} \quad \text{s.t.} \quad \|\mathbf{D}\underline{\alpha} - \underline{y}\|_{2}^{2} \leq \varepsilon^{2}$$



Smooth the L₀ and use continuous optimization techniques

Greedy methods

Build the solution one non-zero element at a time



Relaxation – The Basis Pursuit (BP)



- This is known as the Basis-Pursuit (BP) [Chen, Donoho & Saunders ('95)]
- The newly defined problem is convex (quad. programming)
- Very efficient solvers can be deployed:
 - Interior point methods [Chen, Donoho, & Saunders ('95)] [Kim, Koh, Lustig, Boyd, & D. Gorinevsky (`07)]
 - Sequential shrinkage for union of ortho-bases [Bruce et.al. ('98)]
 - Iterative shrinkage [Figuerido & Nowak ('03)] [Daubechies, Defrise, & De-Mole ('04)]
 [E. ('05)] [E., Matalon, & Zibulevsky ('06)] [Beck & Teboulle (`09)] ...



Go Greedy: Matching Pursuit (MP)

- The MP is one of the greedy algorithms that finds one atom at a time [Mallat & Zhang ('93)]
- Step 1: find the one atom that
 best matches the signal



- Next steps: given the previously found atoms, find the next <u>one</u> to best fit the residual
- The algorithm stops when the error $\|\mathbf{p}_{\underline{\alpha}} \underline{y}\|_2$ is below the destination threshold.
- The Orthogonal MP (OMP) is an improved version that reevaluates the coefficients by Least-Squares after each round.



Pursuit Algorithms

$$\min_{\underline{\alpha}} \|\underline{\alpha}\|_{0} \quad \text{s.t.} \quad \|\mathbf{D}\underline{\alpha} - \underline{y}\|_{2}^{2} \leq \varepsilon^{2}$$

There are various algorithms designed for approximating the solution of this problem:

- Greedy Algorithms: Matching Pursuit, Orthogosal Matching Pursuit (OMP), Least-Squares-OMP, Weak Matching Nesuit, Block Matching Pursuit [1993-today]. Relaxation Algorithms: Basis Purpets.k.a. LASSO), Dnatzig Selector & numerical ways to har defilem [1995-today].
- Hybrid Algorithms: StOMP, CoSaMP, Subspace Pursuit, Iterative Hard-Thresholding [2007-today].



The Mutual-Coherence



- $\circ~$ The Mutual-Coherence μ is the largest off-diagonal entry in absolute value
- The Mutual-Coherence is a property of the dictionary (just like the "Spark"). In fact, the following relation can be shown:

$$\sigma \geq 1 + \frac{1}{\mu}$$



BP and MP Equivalence (No Noise)

Equivalence Donoho & E. ('02) Gribonval & Nielsen ('03) Tropp ('03) Temlyakov ('03) Given a signal <u>x</u> with a representation $\underline{X} = \underline{D}\underline{\alpha}$, assuming that $\|\underline{\alpha}\|_{0} < 0.5(1+1/\mu)$, BP and MP are guaranteed to find the sparsest solution.

- MP and $\hat{\mathbf{BP}}$ are different in general (hand to say which is better) $\underline{\alpha}$
- The above result corresponds to the worst-case, and as such, it is too pessimistic
- Average performance results are available too, showing much better bounds [Donoho (`04)] [Candes et.al. ('04)] [Tanner et.al. ('05)] [E. ('06)] [Tropp et.al. ('06)] ... [Candes et. al. ('09)]



BP Stability for the Noisy Case

Stability

Given a signal $y = \mathbf{D}\underline{\alpha} + \underline{v}$ with a representation satisfying $\|\underline{\alpha}\|_{0} < 1/3\mu$ and a white Gaussian Noise $\underline{\mathbf{v}} \sim \mathsf{N}(\mathbf{0}, \sigma^2 \mathbf{I})$

Ben-Haim, Eldar & E. ('09)

* With very high probability

- Ο
- For $\sigma=0$ we get a wea This result is the orac $\min_{\alpha} \lambda \|\underline{\alpha}\|_{1} + \|\mathbf{D}\underline{\alpha} \underline{y}\|_{2}^{2}$ Ο
- Similar results exist for other pursuit algorithms (Dantzig Selector, Orthogonal Ο Matching Pursuit, CoSaMP, Subspace Pursuit, ...)



To Summarize So Far ...

Image denoising (and many other problems in image processing) requires a model for the desired image



What

next?

We proposed a model for signals/images based on sparse and redundant representations

The Dictionary **D** should be found somehow !!! We have seen that there are approximation methods to find the sparsest solution, and there are theoretical results that guarantee their success.



Problems?

Dictionary Learning: The K-SVD Algorithm



What Should **D** Be?

$$\hat{\underline{\alpha}} = \underset{\underline{\alpha}}{\operatorname{argmin}} \|\underline{\alpha}\|_{0} \quad \text{s.t.} \quad \frac{1}{2} \| \mathbf{D}\underline{\alpha} - \underline{y} \|_{2}^{2} \leq \varepsilon^{2} \quad \Longrightarrow \quad \hat{\underline{x}} = \mathbf{D}\hat{\underline{\alpha}}$$

Our Assumption: Good-behaved Images have a sparse representation

D should be chosen such that it sparsifies the representations

One approach to choose **D** is from a known set of transforms (Steerable wavelet, Curvelet, Contourlets, Bandlets, Shearlets ...)

The approach we will take for building **D** is training it, based on Learning from Image Examples


Dictionary Learning: Problem Setting





Given these P examples and a fixed size [n×m] dictionary **D**:

1. Is **D** unique?

2. How would we find **D**?



Measure of Quality for D



۲X

[Aharon, E. & Bruckstein ('05)]

K–Means For Clustering





The K–SVD Algorithm – General





K–SVD: Sparse Coding Stage



K–SVD: Dictionary Update Stage



We should solve:



We refer only to the examples that use the column <u>d</u>_k

Fixing all **A** and **D** apart from the k^{th} column, and seek both \underline{d}_k and the k^{th} column in A to better fit the residual!



A Synthetic Experiment





Improved Dictionary Learning

$$\underset{\mathbf{D},\mathbf{A}}{\text{Min}} \sum_{j=1}^{P} \left\| \mathbf{D}\underline{\alpha}_{j} - \underline{x}_{j} \right\|_{2}^{2} \quad \text{s.t. } \forall j, \left\| \underline{\alpha}_{j} \right\|_{0} \leq L$$

MOD Algorithm

Fix **D** and update **A**

Fix **A** and update **D**

K-SVD Algorithm

Fix **D** and update **A**

for j=1:1:m

Fix A & D apart from the j-th atom its coefficients
Update d_i and its coef. in A

end



Improved Dictionary Learning

$$\begin{split} & \underset{\mathbf{D},\mathbf{A}}{\text{Min}} \sum_{j=1}^{P} \left\| \mathbf{D}\underline{\alpha}_{j} - \underline{\chi}_{j} \right\|_{2}^{2} \quad \text{s.t. } \forall j, \left\| \underline{\alpha}_{j} \right\|_{0} \leq \mathsf{L} \\ & \underset{\text{Improved Algorithm}}{\text{Ismith & E. 2013}} \\ & \text{Fix D and update A} \\ & \text{Fix D and update A} \\ & \text{First D and K-SVD can be considered added by the second of this methods} \\ & \text{First D and K-SVD can be done in two ways:} \\ & \text{First D and K-SVD can be done in two ways:} \\ & \text{His can be done in two ways:} \\ & \text{His can be done in the K-SVD, or} \\ & \text{Second added of the second of the seco$$



To Summarize So Far ...

Image denoising (and many other problems in image processing) requires a model for the desired image

What do we do?

What

next?

We proposed a model for signals/images based on sparse and redundant representations

Will it all work in applications? We have seen approximation methods that find the sparsest solution, and theoretical results that guarantee their success. We also saw a way to learn **D**



Problems?

BREAK

Back to Denoising ... and Beyond – Combining it All



Bringing Sparseland to Applications

• While the Sparseland model is clear and well-defined, there are various ways to bring it into an actual algorithm in applications



- The bad news: It is not obvious how to turn this model into a successful algorithm
- The good news: As we are about to see, there is a lot of room for ingenuity & originality in designing algorithms in image processing



From Local to Global Treatment

The K-SVD is reasonable for low-dimension \bigcirc signals (n in the range 10-400). As n grows, the complexity and the memory requirements of the K-SVD are prohibitive



- So, how should large images be handled? \bigcirc
- The solution: Force shift-invariant sparsity operate on patches \bigcirc of size n-by-n (n=8) in the image, including overlaps

$$\hat{\underline{x}} = \underset{\underline{x}, \{\underline{\alpha}_{ij}\}_{ij}}{\operatorname{ArgMin}} \quad \frac{1}{2} \left\| \underline{x} - \underline{y} \right\|_{2}^{2} + \mu \underset{ij}{\sum} \left\| \underline{R}_{ij} \underline{x} - \underline{D} \underline{\alpha}_{ij} \right\|_{2}^{2}$$

$$\text{Extracts a the (i,j) location of the (i,j) loca$$

Extracts a patch in the (i,j) location

What Data to Train On?

Option 1:

- Use a database of images
- We tried that, and it works fine (~0.5-1dB below the state-of-the-art)

Option 2:

- Use the corrupted image itself !!
- Simply sweep through all patches of size \sqrt{n} -by- \sqrt{n} (overlapping blocks)
- Image of size 1000^2 pixels $\implies \sim 10^6$ examples to use – more than enough
- This works much better!







K-SVD Image Denoising





Image Denoising (Gray) [E. & Aharon ('06)]





EPLL Improvement [Sulam and E. ('15)]

$$\hat{\underline{\mathbf{X}}} = \underset{\underline{\mathbf{X}}, \{\underline{\alpha}_{ij}\}_{ij}, \mathbf{D}}{\text{ArgMin}} \frac{1}{2} \left\| \underline{\mathbf{X}} - \underline{\mathbf{Y}} \right\|_{2}^{2} + \mu \sum_{ij} \left\| \mathbf{R}_{ij} \underline{\mathbf{X}} - \mathbf{D} \underline{\alpha}_{ij} \right\|_{2}^{2} \text{ s.t. } \left\| \underline{\alpha}_{ij} \right\|_{0} \leq L$$

- The algorithm we proposed updates <u>x</u> only once at the end
- Why not repeat the whole process several times?
- The rationale: The sparse representation model should be imposed on the patches of the FINAL image. After averaging, this is ruined





EPLL Improvement [Sulam and E. ('15)]

- Expected Patch Log Likelihood (EPLL) is an algorithm that came to fix this problem [Zoran and Weiss, ('11)], originally in the context of a GMM prior
- An extension of EPLL to *Sparsland* is proposed in [Sulam and E. ('15)]. The core idea is:
 - After the image has been computed, we proceed the iterative process, and apply several such overall rounds of updates
 - Sparse coding must be done with a new threshold, based on the remaining noise in the image. This is done by evaluating the noise level based on the linear projections (disregarding the support detection by the OMP)
 - This algorithm leads to state-of-the-art results, with 0.5-1dB improvement over the regular K-SVD algorithm shown before



EPLL Improvement [Sulam and E. ('15)]





Denoising (Color) [Mairal, E. & Sapiro ('08)]

U When turning to handle color images, the main





Denoising (Color) [Mairal, E. & Sapiro ('08)]

Our experiments lead to state-of-the-art denoising results, giving ~1dB better results compared to [Mcauley et. al. ('06)] which implements a learned MRF model (Field-of-Experts)



Original





Noisy (12.77dB)

Result (29.87dB)



Video Denoising [Protter & E. ('09)]



Denoised (PSNR=29.98)



Low-Dosage Tomography [Shtok, Zibulevsky & E. ('10)]

- In Computer-Tomography (CT) reconstruction, an image is recovered from a set of its projections
- In medicine, CT projections are obtained by X-ray, and it typically requires a high dosage of radiation in order to obtain a good quality reconstruction
- A lower-dosage projection implies a stronger noise (Poisson distributed) in data to work with
- Armed with sparse and redundant representation modeling, we can denoise the data and the final reconstruction ... enabling CT with lower dosage



Image Inpainting – The Basics

- Assume: the signal <u>x</u> has been created by $\underline{x}=D\underline{\alpha}_0$ with very sparse $\underline{\alpha}_0$
- Missing values in <u>x</u> imply missing rows in this linear system
- By removing these rows, we get $\tilde{\mathbf{D}}\underline{\alpha} = \underline{\tilde{\mathbf{X}}}$
- Now solve $\min_{\alpha} \|\underline{\alpha}\|_{0}$ s.t. $\underline{\tilde{\mathbf{X}}} = \mathbf{\tilde{D}}\underline{\alpha}$
- If $\underline{\alpha}_0$ was sparse enough, it will be the solution of the above problem! Thus, computing $\mathbf{D}\underline{\alpha}_0$ recovers <u>x</u> perfectly



Side Note: Compressed-Sensing

- Compressed Sensing is leaning on the very same principal, leading to alternative sampling theorems.
- Assume: the signal <u>x</u> has been created by $\underline{x} = \mathbf{D}\underline{\alpha}_0$ with very sparse $\underline{\alpha}_0$.
- Multiply this set of equations by the matrix **Q** which reduces the number of rows.
- The new, smaller, system of equations is
 QD<u>α</u> = QX → D<u>α</u> = X

 If <u>α</u>₀ was sparse enough, it will be the sparsest solution of the new system, thus, computing D<u>α</u>₀ recovers <u>x</u> perfectly.
- Compressed sensing focuses on conditions for this to happen, guaranteeing such recovery.



Inpainting Formulation [Mairal, E. & Sapiro ('08)]

$$\hat{\underline{\mathbf{x}}} = \underset{\underline{\mathbf{x}}, \{\underline{\alpha}_{ij}\}_{ij}, \mathbf{D}}{\text{ArgMin } \frac{1}{2} \left\| \underline{\mathsf{M}} \underline{\mathbf{x}} - \underline{\mathbf{y}} \right\|_{2}^{2} + \mu \sum_{ij} \left\| \mathbf{R}_{ij} \underline{\mathbf{x}} - \mathbf{D} \underline{\alpha}_{ij} \right\|_{2}^{2} \text{ s.t. } \left\| \underline{\alpha}_{ij} \right\|_{0} \leq L$$

The matrix M is a mask matrix, obtained by the identity matrix with some of its rows omitted, corresponding to the missing samples





Inpainting Formulation [Mairal, E. & Sapiro ('08)]

$$\hat{\underline{X}} = \underset{\underline{X}, \{\underline{\alpha}_{ij}\}_{ij}, \mathbf{D}}{\operatorname{Arg}} \frac{1}{2} \left\| \underbrace{M}\underline{X} - \underline{y} \right\|_{2}^{2} + \underset{ij}{\sum} \left\| \mathbf{R}_{ij}\underline{X} - \mathbf{D}\underline{\alpha}_{ij} \right\|_{2}^{2} \text{ s.t. } \left\| \underline{\alpha}_{ij} \right\|_{0} \leq \mathbf{L}$$

$$\underline{X} = \underline{Y} \text{ and } \mathbf{D} \text{ known} \qquad \underline{X} \text{ and } \alpha_{ij} \text{ known} \qquad \mathbf{D} \text{ and } \alpha_{ij} \text{ known}$$

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

$$\underline{\alpha}_{ij} = \underset{\alpha}{\operatorname{Min}} \left\| \underbrace{M}_{ij} \left(\mathbf{R}_{ij}\underline{X} - \mathbf{D}\underline{\alpha} \right) \right\|_{2}^{2}$$

$$\text{ s.t. } \left\| \underline{\alpha} \right\|_{0} \leq \mathbf{L}$$

$$using \text{ the matching pursuit}$$

$$Compute \mathbf{D} \text{ to minimize}$$

$$\underbrace{Min}_{\alpha} \sum_{ij} \left\| \underbrace{M}_{ij} \left(\mathbf{R}_{ij}\underline{X} - \mathbf{D}\underline{\alpha} \right) \right\|_{2}^{2}$$

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

$$\underbrace{Min}_{V} \underline{Y} + \mu \sum_{ij} \mathbf{R}_{ij}^{T} \mathbf{D}\underline{\alpha}_{ij} \right]$$

$$which is a again a simple averaging of patches$$



Inpainting [Mairal, E. & Sapiro ('08)]

For the Peppers image

Alg.	RMSE for 25% missing	RMSE for 50% missing
DCT: No-overlap	14.55	19.61
DCT: Overlap	9.00	11.55
K-SVD	8.1	10.05

This is a more challenging case, where the DCT is not a suitable dictionary.

- For Redundant DCT we get RMSE=16.13, and
- For K-SVD (15 iterations) we get RMSE=12.74



V

Inpainting [Mairal, E. & Sapiro ('08)]





Inpainting [Mairal, E. & Sapiro ('08)]

The same can be done for video, very much like the denoising treatment: (i) 3D patches, (ii) no need to compute the dictionary from scratch for each frame, and (iii) no need for explicit motion estimation



Original

80% missing

Result



Demosaicing [Mairal, E. & Sapiro ('08)]

 Our experiments lead to state-of-the-art demosaicing
 Today's cameras are sensing only one results, giving '0.2dB better results on a color per pixel, leaving the rest for compared to [chang & Chan ('06)]
 interpolated

 Generalizing the inpainting scheme to handle demosaicing is tricky because of the possibility to learn the mosaic pattern within the dictionary



 In order to avoid "over-fitting", we handle the demosaicing problem while forcing strong sparsity and applying only few iterations





Image Compression [Bryt and E. ('08)]

- The problem: Compressing photo-ID images
- General purpose methods (JPEG, JPEG2000)
 do not take into account the specific family
- By adapting to the image-content (PCA/K-SVD), better results could be obtained
- For these techniques to operate well, train locally (per patch) using a training set of images is required
- In PCA, only the (quantized) coefficients are stored, the K-SVD requires storage of the indices
- Geometric alignment of the image is very helpful should be done [Goldenberg, Kimmel, & E. ('05)]





Image Compression





Image Compression Results

Original JPEG JPEG-2000 Local-PCA QQ K-SVD **Results** for 820 Bytes per each file



Image Compression Results




Image Compression Results





Deblocking the Results [Bryt and E. (`09)]

550 bytes K-SVD results with and without deblocking



K-SVD (5.49)



Deblock (6.24)



Deblock (5.27)



K-SVD (6.45)





K-SVD (11.67)



Deblock (11.32)



Super-Resolution [Zeyde, Protter, & E. ('11)]

- Given a low-resolution image, we desire to enlarge it while producing a sharp looking result. This problem is referred to as "Single-Image Super-Resolution"
- Image scale-up using bicubic interpolation is far from being satisfactory for this task
- A brilliant and very different sparse and redundant representation technique was proposed [Yang, Wright, Huang, and Ma ('08)] for solving this problem, by training a coupleddictionaries for the low- and high res. images
- We extended and improved their algorithms and results



Super-Resolution – Results (1)

This book is about *convex optimization*, a special class of mathematical optimization problems, which includes least-squares and linear programming problems. It desis is well known that least-squares and linear programming problems have a fairly mia complete theory, arise in a variety of applications, and can be solved numerically Ind very efficiently. The basic point of this book is that the same can be said for the It is larger class of convex optimization problems. Inol

While the mathematics of convex optimization has been studied for about a century, several related recent developments have stimulated new interest in the topic. The first is the recognition that interior-point methods, developed in the 1980s to solve linear programming problems, can be used to solve convex optimiza tion problems as well. These new methods allow us to solve certain new classes of convex optimization problems, such as semidefinite programs and second-order cone programs, almost as easily as linear programs.

The second development is the discovery that convex optimization problems cam be (beyond least-squares and linear programs) are more prevalent in practice than was previously thought. Since 1990 many applications have been discovered in areas such as automatic control systems, estimation and signal processing, comon hos munications and networks, electronic circuit design, data analysis and modeling iti réne statistics, and finance. Convex optimization has also found wide application in combinatorial optimization and global optimization, where it is used to find bounds or the optimal value, as well as approximate solutions. We believe that many other solutions. apply d applications of convex optimization are still waiting to be discovered.

There are great advantages to recognizing or formulating a problem as a convex optimization problem. The most basic advantage is that the problem can then be solved, very reliably and efficiently, using interior-point methods or other special designer methods for convex optimization. These solution methods are reliable enough to be modifies embedded in a computer-aided design or analysis tool, or even a real-time reactive iston ma or automatic control system. There are also theoretical or conceptual advantages of formulating a problem as a convex optimization problem. The associated dual

FJINN-14.000D

Ideal Image

The trainir

×717

An amazing variety of practical proble design, analysis, and operation) can be mization problem, or some variation such Indeed, mathematical optimization has b It is widely used in engineering, in elect trol systems, and optimal design probler and aerospace engineering. Optimization design and operation, finance, supply ch other areas. The list of applications is st

For most of these applications, mathe a human decision maker, system designer process, checks the results, and modifies when necessary. This human decision ma by the optimization problem, e.g., buyin iding aportiolio.

<u>89 training pates</u>

Given Image



and

desis

othe

a hu

proc

when

by t

port

Super-Resolution – Results (2)



Given image



Scaled-Up (factor 2:1) using the proposed algorithm, PSNR=29.32dB (3.32dB improvement over bicubic)



Super-Resolution – Results (2)



The Original

Bicubic Interpolation

SR result



Super-Resolution – Results (2)



The Original

Bicubic Interpolation

SR result



Poisson Denoising







$$peak = 0.1$$



$$P(y \mid x) = \frac{x^{y}}{y!}e^{-x}$$

peak $\triangleq \max_{i,j} \{x_{i,j}\}$



Poisson Denoising [Salmon et. al., 2011] [Giryes et. al., 2013]

- Anscombe transform converts Poisson distributed noise into an approximately Gaussian one, with variance 1 using the following formula [Anscombe, 1948]: $f_{Anscombe} \left(y \right) = 2\sqrt{y + \frac{3}{8}}$
- However, this is of reasonable accuracy only if peak>4.
- For lower peaks (poor illumination), we use the patch-based approach with dictionary learning, BUT ... in the exponent domain:

$$\left\{ \underbrace{\mathbf{x}}_{\mathbf{x}} = \mathbf{D}\underline{\alpha} \\ \text{where } \|\underline{\alpha}\|_{0} \leq \mathsf{L} \right\} \longrightarrow \begin{cases} \underline{\mathbf{x}}_{\mathbf{x}} = \exp\{\mathbf{D}\underline{\alpha}\} \\ \text{where } \|\underline{\alpha}\|_{0} \leq \mathsf{L} \end{cases}$$



Poisson Denoising – Results (1)







Original

Noisy (peak=1)

Result (PSNR=22.59dB)

Dictionary learned atoms:





Poisson Denoising – Results (2)





Other Applications?

- Poisson Denoising & Inpainting
- Blind deblurring
- Audio inpainting
- Dynamic MRI reconstruction
- Clutter reduction in ultrasound
- □ Single image interpolation
- Anomaly detection





Summary and Conclusion



In this Part we Have Seen that ...

Sparsity, Redundancy, and the use of examples are important ideas that can be used in designing better tools in signal/image processing

What do we do?

In our work on we cover theoretical, numerical, and applicative issues related to this model and its use in practice

Many of the results we got focused on patchbased methods – it is time to understand better this choice and its limitations, with the hope to lead to new insights

Next we focus on ... We keep working on:

- Improving the model
- Improve the dictionaries
- Demonstrate on other applications

X

What

next

Thank You

All this Work is Made Possible Due to

my teachers and mentors



colleagues & friends collaborating with me





and my students

G. Sapiro J.L. Starck

I. Yavneh M. Zibulevsky P. Milanfar







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



Welcome to Sparseland Part 3: A Tale of Three Models Sparseland→CSC→CNN

Michael Elad

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





Michael Elad The Computer-Science Department The Technion The research leading to these results has been received funding from the European union's Seventh Framework Program (FP/2007-2013) ERC grant Agreement ERC-SPARSE- 320649



In This Talk



The Underlying Idea

Generative Modeling

of data sources enable

- A systematic algorithm development, &
- A theoretical analysis of their performance

* Only CNN? What about other architectures ?



The Results Presented

... are the fruit of a joint work with







Yaniv Romano Vardan Papyan Jeremias Sulam



Michael Elad The Computer-Science Department The Technion

Part I Motivation and Background



Michael Elad The Computer-Science Department The Technion

Our Starting Point: Image Denoising



Many (thousands) image denoising algorithms have been proposed over the years, some of which are extremely effective









Michael Elad The Computer-Science Department The Technion Published Items in Each Year



Leading Image Denoising Methods...

are built upon powerful patch-based local models:









Popular local models: GMM

Sparse-Representation Example-based Low-rank Field-of-Experts & Neural networks



Patch-Based Image Denoising

- K-SVD: sparse representation modeling of image patches
 [Elad & Aharon, '06]
- BM3D: combines sparsity and self-similarity
 [Dabov, Foi, Katkovnik & Egiazarian '07]
- EPLL: uses GMM of the image patches
 [Zoran & Weiss '11]
- MLP: multi-layer perceptron
 [Burger, Schuler & Harmeling '12]
- NCSR: non-local sparsity with centralized coefficients
 [Dong, Zhang, Shi & Li '13]
- WNNM: weighted nuclear norm of image patches
 [Gu, Zhang, Zuo & Feng '14]
- SSC–GSM: nonlocal sparsity with a GSM coefficient model
 [Dong, Shi, Ma & Li '15]





Recall K-SVD Denoising [Elad & Aharon, '06]



• Despite its simplicity, this is a very well-performing algorithm

- $\,\circ\,$ Its origins can be traced back to Guleryuz's local DCT recovery
- A small modification of this method leads to state-of-the-art results [Mairal, Bach, Ponce, Spairo, Zisserman, `09]



What is Missing?

 Over the years, many kept revisiting this algorithm and its line of thinking, with a clear feeling that key features are still lacking



- What is missing? Here is what **WE** thought of...
 - A multi-scale treatment [Ophir, Lustig & Elad '11] [Sulam, Ophir & Elad '14] [Papyan & Elad '15]
 - Exploiting self-similarities [Ram & Elad '13] [Romano, Protter & Elad '14]
 - Pushing to better agreement on the overlaps [Romano & Elad '13] [Romano & Elad '15]
 - Enforcing the local model on the final patches (EPLL) [Sulam & Elad '15]
- $\circ\,$ Eventually, we realized that the key part that is missing is

A Theoretical Backbone



Missing Theoretical Backbone?

 \circ The core global-local model assumption on $\mathbf{X} \in \mathbb{R}^N$:

 $\forall i \quad \mathbf{R}_i \mathbf{X} = \mathbf{\Omega} \mathbf{\gamma}_i \quad \text{where} \quad \|\mathbf{\gamma}_i\|_0 \leq k$

Every patch in the unknown signal is expected to have a sparse representation w.r.t. the same dictionary Ω

Questions to consider:

- Who are the signals belonging to this model? Do they exist?
- How should we project a signal on this model (pursuit)?
- Could we offer theoretical guarantees for this model/algorithms?
- Could we offer a global pursuit algorithm that operates locally?
- How should we learn Ω if this is indeed the model?

 As we will see, all these questions are very relevant to recent developments in signal processing and machine learning



Coming Up



Limitations of patch averaging



Convolutional Sparse Coding (CSC) model

Multi-Layer Convolutional Sparse Coding (ML-CSC)



Convolutional neural networks (CNN)



Fresh view of CNN through the eyes of sparsity



Part II Convolutional Sparse Coding

(IEEE-TSP)

Working Locally Thinking Globally: Theoretical Guarantees for Convolutional Sparse Coding Vardan Papyan, Jeremias Sulam and Michael Elad

(ICCV 2017)

Convolutional Dictionary Learning via Local Processing Vardan Papyan, Yaniv Romano, Jeremias Sulam, and Michael Elad



Convolutional Sparse Coding (CSC)



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





Intuitively ...





Michael Elad The Computer-Science Department The Technion ○ Here is an alternative global sparsity-based model formulation

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

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

Cⁱ

 $\mathbf{r}^{i} \in \mathbb{R}^{N}$ are the corresponding coefficients





Two Interpretations





Michael Elad The Computer-Science Department The Technion

Why CSC?



stripe-dictionary ______stripe vector -

Every patch has a sparse representation w.r.t. to the same local dictionary Ω , just as we have assumed



X

R

= DI

CSC Relation to Our Story

- \circ A clear global model: every patch has a sparse representation w.r.t. to the same local dictionary Ω , just as we have assumed
- \odot No notion of disagreement on the patch overlaps
- \odot Related to the current common practice of patch averaging (\boldsymbol{R}_{i}^{T})
 - put the patch $\Omega \gamma_i$ back in the i-th location of the global vector)

$$\mathbf{X} = \mathbf{D}\mathbf{\Gamma} = \frac{1}{n} \sum_{i} \mathbf{R}_{i}^{\mathrm{T}} \mathbf{\Omega} \mathbf{\gamma}_{i}$$

- What about the Pursuit?
 - "Patch averaging": independent sparse coding for each patch
 - CSC: should seek all the representations together
- \odot Is there a bridge between the two? We'll come back to this later ...





 This model has been used in the past [Lewicki & Sejnowski '99] [Hashimoto & Kurata, '00]

 Most works have focused on solving *efficiently* its associated pursuit, called **convolutional sparse coding**, using the BP algorithm

 $(\mathbf{P}_1^{\epsilon}): \min_{\mathbf{\Gamma}} \|\mathbf{\Gamma}\|_1 + \lambda \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_2^2$

Convolutional dictionary

Several applications were demonstrated:

- Pattern detection in images and the analysis of instruments in music signals [Mørup, Schmidt & Hansen '08]
- Inpainting [Heide, Heidrich & Wetzstein '15]
- Super-resolution [Gu, Zuo, Xie, Meng, Feng & Zhang '15]

However, little is known regrading its theoretical aspects. Why?
 Perhaps because the regular SparsLand theory is sufficient?



Classical Sparse Theory (Noiseless)

$$(\mathbf{P}_0): \min_{\mathbf{\Gamma}} \|\mathbf{\Gamma}\|_0 \text{ s.t. } \mathbf{X} = \mathbf{D}\mathbf{\Gamma}$$

Definition: Mutual-Coherence: $\mu(\mathbf{D}) = \max_{i \neq j} |d_i^T d_j|$

[Donoho & Elad '03]



then this solution is necessarily the sparsest

[Donoho & Elad '03]

Theorem: The OMP and BP are guaranteed to recover the true sparse code assuming that $\|\Gamma\|_0 < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D})}\right)$

[Tropp '04], [Donoho & Elad '03]


The Need for a Theoretical Study

 \odot Assuming that m=2 and n=64 we have that [Welch, '74] $\mu(\mathbf{D}) \geq 0.063$

 As a result, uniqueness and success of pursuits is guaranteed as long as

$$\|\mathbf{\Gamma}\|_{0} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D})}\right) \le \frac{1}{2} \left(1 + \frac{1}{0.063}\right) \approx 8$$

Less than 8 non-zeros GLOBALLY are allowed!!!
 This is a very pessimistic result!

- Repeating the above for the noisy case leads to even worse performance predictions
- Bottom line: Classic SparseLand Theory cannot provide good explanations for the CSC model





Moving to Local Sparsity: Stripes

 $\ell_{0,\infty}$ Norm: $\|\boldsymbol{\Gamma}\|_{0,\infty}^{s} = \max_{i} \|\boldsymbol{\gamma}_{i}\|_{0}$

 $(\mathbf{P}_{0,\infty})$: min $\|\mathbf{\Gamma}\|_{0,\infty}^{s}$ s.t. $\mathbf{X} = \mathbf{D}\mathbf{\Gamma}$

 $\|\Gamma\|_{0,\infty}^{s}$ is low \rightarrow all γ_{i} are sparse \rightarrow every patch has a sparse representation over Ω

The Main Questions we Aim to Address:

- I. Is the solution to this problem unique ?
- II. Can we recover the solution via a global OMP/BP ?



m = 2

Stripe-Spark and Uniqueness

$$\begin{pmatrix} \mathbf{P}_{0,\infty} \end{pmatrix}: \quad \min_{\Gamma} \ \|\Gamma\|_{0,\infty}^{s} \text{ s.t. } \mathbf{X} = \mathbf{D}\Gamma$$
efinition: Stripe Spark $\eta_{\infty}(\mathbf{D}) = \min_{\Delta} \ \|\Delta\|_{0,\infty}^{s} \text{ s.t. } \left\{ \begin{matrix} \mathbf{D}\Delta = 0 \\ \Delta \neq 0 \end{matrix} \right\}$
Theorem: If a solution Γ is found for $(\mathbf{P}_{0,\infty})$ such that:
$$\|\Gamma\|_{0,\infty}^{s} < \frac{1}{2}\eta_{\infty}$$
then it is necessarily the optimal solution to this problem
Theorem: The relation between the
Stripe-Spark and the Mutual Coherence is:
$$\eta_{\infty}(\mathbf{D}) \ge 1 + \frac{1}{\mu(\mathbf{D})}$$



D

Uniqueness via Mutual Coherence

$$(\mathbf{P}_{0,\infty})$$
: min $\|\mathbf{\Gamma}\|_{0,\infty}^{s}$ s.t. $\mathbf{X} = \mathbf{D}\mathbf{\Gamma}$

Theorem: If a solution Γ is found for $(\mathbf{P}_{0,\infty})$ such that: $\|\Gamma\|_{0,\infty}^{s} < \frac{1}{2}\left(1 + \frac{1}{\mu(\mathbf{D})}\right)$ then this is necessarily the unique optimal solution to

this problem

This result is exciting: This and later results pose a local constraint for a global guarantee, and as such, they are far more optimistic compared to the global guarantees For k non-zeros per stripe, and filters of length n, we get $\|\mathbf{\Gamma}\|_0 \cong \frac{k}{2n-1} \cdot N$ non-zeros globally



Recovery Guarantees

$$\begin{pmatrix} \mathbf{P}_{0,\infty} \end{pmatrix}: \quad \min_{\Gamma} \| \| \Gamma \|_{0,\infty}^{s} \text{ s.t. } \mathbf{X} = \mathbf{D}\Gamma$$
Lets solve this problem via OMP or BP^{*}, applied globally
Theorem: If a solution Γ of $(\mathbf{P}_{0,\infty})$ satisfies:
$$\| \Gamma \|_{0,\infty}^{s} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D})} \right)$$
then global OMP and BP are guaranteed to find it

Both OMP and BP do not assume **local** sparsity but still guaranteed to succeed. One could propose algorithms that rely on this assumption

* How about variants that would exploit the local sparsity?



From Ideal to Noisy Signals

 \odot So far, we have assumed an ideal signal $X=D\Gamma$

- \odot However, in practice we usually have $Y=D\Gamma+E$ where E is due to noise or model deviations
- \odot To handle this, we redefine our problem as:

$$(\mathbf{P}_{0,\infty}^{\epsilon})$$
: min $\|\mathbf{\Gamma}\|_{0,\infty}^{s}$ s.t. $\|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_{2} \leq \epsilon$

\odot The Main Questions We Aim to Address:

- I. Stability of the solution to this problem ?
- II. Stability of the solution obtained via global OMP/BP ?
- III. Could the same recovery be done via local (patch) operations ?



Stability of via Stripe-RIP

If you carefully review this result, you should be disappointed, as we see a noise magnification !! Is this true ?

Answer: No!!! This is a worst-case (with adversarial noise) analysis



Local Noise Assumption

- \circ Thus far, our analysis relied on the local sparsity of the underlying solution Γ , which was enforced through the $\ell_{0,\infty}$ norm
- \odot In what follows, we present stability guarantees for both OMP and BP that will also depend on the local energy in the noise vector E
- \circ This will be enforced via the $\ell_{2,\infty}$ norm, defined as:

 $\|\mathbf{E}\|_{2,\infty}^{p} = \max_{i} \|\mathbf{R}_{i}\mathbf{E}\|_{2}$



Stability of OMP

Theorem: If $\mathbf{Y} = \mathbf{D}\mathbf{\Gamma} + \mathbf{E}$ where $\|\mathbf{\Gamma}\|_{0,\infty}^{s} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D})}\right) - \frac{1}{\mu(\mathbf{D})} \cdot \frac{\|\mathbf{E}\|_{2,\infty}^{p}}{|\Gamma_{\min}|}$ then OMP run for $\|\mathbf{\Gamma}\|_{0}$ iterations will 1. Find the correct support 2. $\|\mathbf{\Gamma}_{OMP} - \mathbf{\Gamma}\|_{2}^{2} \le \frac{\|\mathbf{E}\|_{2}^{2}}{1 - (\|\mathbf{\Gamma}\|_{0,\infty}^{s} - 1)\mu(\mathbf{D})}$



Michael Elad The Computer-Science Department The Technion

Stability of Lagrangian BP

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

Theorem: For $\mathbf{Y} = \mathbf{D}\mathbf{\Gamma} + \mathbf{E}$, if $\lambda = 4 \|\mathbf{E}\|^p$ and

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

Then we are guaranteed that

- 1. The support of $\Gamma_{\rm BP}$ is contained
- 2. $\|\mathbf{\Gamma}_{\mathrm{BP}} \mathbf{\Gamma}\|_{\infty} \le 7.5 \|\mathbf{E}\|_{2,\infty}^{\mathrm{p}}$
- 3. Every entry greater than 7.5
- 4. $\Gamma_{\rm BP}$ is unique

Theoretical foundation for 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

Global Pursuit via Local Processing

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

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

 α_i

- While CSC is a global model, its theoretical guarantees rely on local properties
- We aim to show that this global-local relation can also be exploited for solving the global BP problem using only local operations



Global Pursuit via Local Processing (2)

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





Global Pursuit via Local Processing (2)

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

Turning to the local form and using the Augmented Lagrangian $\min_{\boldsymbol{\alpha}_{i},\boldsymbol{s}_{i}} \frac{1}{2} \left\| \mathbf{Y} - \sum_{i} \mathbf{R}_{i}^{T} \mathbf{s}_{i} \right\|_{2}^{2} + \sum_{i} \left(\lambda \|\boldsymbol{\alpha}_{i}\|_{1} + \frac{\rho}{2} \|\mathbf{s}_{i} - \mathbf{D}_{L} \boldsymbol{\alpha}_{i} + \mathbf{u}_{i}\|_{2}^{2} \right)$

- \circ These two problems are equivalent, and convex w.r.t their variables
- \odot The new formulation targets the local slices, and their sparse representations
- \odot The vectors \mathbf{u}_i are the Lagrange multipliers for the constraints $\boldsymbol{s}_i {=} \boldsymbol{D}_L \boldsymbol{\alpha}_i$



Global Pursuit via Local Processing (2)

$$\min_{\boldsymbol{\alpha}_{i},\boldsymbol{s}_{i}} \frac{1}{2} \left\| \mathbf{Y} - \sum_{i} \mathbf{R}_{i}^{\mathrm{T}} \mathbf{s}_{i} \right\|_{2}^{2} + \sum_{i} \left(\lambda \|\boldsymbol{\alpha}_{i}\|_{1} + \frac{\rho}{2} \|\mathbf{s}_{i} - \mathbf{D}_{\mathrm{L}} \boldsymbol{\alpha}_{i} + \mathbf{u}_{i}\|_{2}^{2} \right)$$

$$ADMM$$



Comment: One iteration of this procedure amounts to ... the very same patch-averaging algorithm we started with



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



The Proposed Scheme can be used for Dictionary (D_L) Learning

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





Partial Summary of CSC

 What we have seen so far is a new way to analyze the global CSC model using local sparsity constraints. We proved:



Uniqueness of the solution for the noiseless problem



Stability of the solution for the noisy problem



Guarantee of success and stability of both OMP and BP



We obtained guarantees and algorithms that operate locally while claiming global optimality



We mentioned briefly the mater of learning the model (i.e. dictionary learning for CSC), and presented our competitive approach which is based on simple local steps



Part III Going Deeper

(JMLR 2017)

Convolutional Neural Networks Analyzed via Convolutional Sparse Coding Vardan Papyan, Yaniv Romano and Michael Elad



Michael Elad The Computer-Science Department The Technion

CSC and CNN

 \odot There seems to be a connection between CSC and CNN:

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

 \odot We propose a principled way to analyze CNN

 \odot But first, a short review of CNN...





CNN



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



Michael Elad The Computer-Science Department The Technion ReLU(z) = max(Thr, z)

CNN



Notice that we do not include a pooling stage:

- Can be replaced by a convolutional layer with increased stride without loss in performance [Springenberg, Dosovitskiy, Brox & Riedmiller '14]
- The current state-of-the-art in image recognition does not use it [He, Zhang, Ren & Sun '15]



Mathematically...

 $f(\mathbf{Y}, {\mathbf{W}_{i}}, {\mathbf{b}_{i}}) = \text{ReLU}(\mathbf{b}_{2} + \mathbf{W}_{2}^{T} \text{ReLU}(\mathbf{b}_{1} + \mathbf{W}_{1}^{T}\mathbf{X}))$





Training Stage of CNN

 \odot Consider the task of classification

 \circ Given a set of signals $\{Y_j\}_j$ and their corresponding labels $\{h(Y_j)\}_i$, the CNN learns an end-to-end mapping





Back to CSC



We propose to impose the same structure on the representations themselves $\mathbf{D}_2 \in \mathbb{R}^{Nm_1 \times Nm_2}$ $\Gamma_2 \in \mathbb{R}^{Nm_2}$ m_{2} m_1

Multi-Layer CSC (ML-CSC)



Intuition: From Atoms to Molecules



- $\begin{array}{c} \circ \mbox{ We can chain the all the dictionaries} \\ \mbox{ into one effective dictionary} \\ \mathbf{D}_{eff} = \mathbf{D}_1 \mathbf{D}_2 \mathbf{D}_3 \cdots \mathbf{D}_K \rightarrow \mathbf{x} = \mathbf{D}_{eff} \, \mathbf{\Gamma}_K \end{array} \end{array}$
- This is a special *Sparseland* (indeed, a CSC) model, however:
 - A key property in this model:
 sparsity of intermediate representations

 $\Gamma_1 \in \mathbb{R}^{Nm_1}$

• The effective atoms are combinations of the original atoms \rightarrow molecules \rightarrow cell \rightarrow tissue \rightarrow body-part ...



 Nm_2

 $\overline{\Gamma}_2 \in \mathbb{R}$

A Small Taste: Model Training (MNIST)





A Small Taste: Model Training (CFAR)

 $D_1 (5 \times 5 \times 3)$ $D_1 D_2 (13 \times 13)$

CIFAR Dictionary:

- D₁: 64 filters of size 5x5x3, stride of 2 dense
- D₂: 256 filters of size 5x5x64, stride of 2 82.99 % sparse
- D₃: 1024 filters of size 5x5x256
 90.66 % sparse





ML-CSC: Pursuit

 $\circ \text{ Deep-Coding Problem (DCP}_{\lambda}) \text{ (dictionaries are known):} \\ \left(\begin{array}{c} \mathbf{X} = \mathbf{D}_{1}\mathbf{\Gamma}_{1} \\ \mathbf{\Gamma}_{1} \|_{0,\infty}^{s} \leq \lambda_{1} \end{array} \right)$

Find
$$\{\mathbf{\Gamma}_{j}\}_{j=1}^{K}$$
 s.t. $\begin{cases} \mathbf{\Gamma}_{1} = \mathbf{D}_{2}\mathbf{\Gamma}_{2} & \|\mathbf{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \mathbf{\Gamma}_{K-1} = \mathbf{D}_{K}\mathbf{\Gamma}_{K} & \|\mathbf{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{cases}$

• Or, more realistically for noisy signals,

Find
$$\{\mathbf{\Gamma}_{j}\}_{j=1}^{K}$$
 s.t.
$$\begin{cases} \|\mathbf{Y} - \mathbf{D}_{1}\mathbf{\Gamma}_{1}\|_{2} \leq \mathcal{E} & \|\mathbf{\Gamma}_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{1} = \mathbf{D}_{2}\mathbf{\Gamma}_{2} & \|\mathbf{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \mathbf{\Gamma}_{K-1} = \mathbf{D}_{K}\mathbf{\Gamma}_{K} & \|\mathbf{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{cases}$$



A Small Taste: Pursuit





Michael Elad The Computer-Science Department The Technion

ML-CSC: Dictionary Learning

• Deep-Learning Problem (**DLP** $_{\lambda}$):

Find
$$\{\mathbf{D}_{i}\}_{i=1}^{K}$$
 s.t.
$$\begin{cases} \left\| \mathbf{Y}_{j} - \mathbf{D}_{1}\mathbf{\Gamma}_{1}^{j} \right\|_{2}^{2} \leq \mathcal{E} \quad \left\| \mathbf{\Gamma}_{1}^{j} \right\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{2}^{j} = \mathbf{D}_{2}\mathbf{\Gamma}_{1}^{2} \quad \left\| \mathbf{\Gamma}_{2}^{j} \right\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \mathbf{\Gamma}_{K}^{j} = \mathbf{D}_{K}\mathbf{\Gamma}_{K}^{j} \quad \left\| \mathbf{\Gamma}_{K}^{j} \right\|_{0,\infty}^{s} \leq \lambda_{K} \end{cases}_{j=1}$$

 $\min_{\{\mathbf{D}_i\}_{i=1}^{K}, \mathbf{U}} \sum_{i} \ell\left(h(\mathbf{Y}_i), \mathbf{U}, \mathbf{D}\mathbf{C}\mathbf{P}^{\star}(\mathbf{Y}_i, \{\mathbf{D}_i\})\right)$

The deepest representation Γ_{K}

obtained by solving the DCP

 While the above is an unsupervised DL, a supervised version can be envisioned

[Mairal, Bach & Ponce '12]



ML-CSC: The Simplest Pursuit

Keep it simple! \circ The simplest pursuit algorithm (single-layer case) is the THR algorithm, which operates on a given input signal Y by:

10 $\mathbf{Y} = \mathbf{D}\mathbf{\Gamma} + \mathbf{E}$ $\mathcal{H}_{\beta}(z)$ - Hard 8 $\mathcal{S}_{\beta}(z)$ - Soft and Γ is sparse 6 $\mathcal{S}^+_{\beta}(z)$ - Soft Nonnegative 4 $\mathbf{2}$ 0 $\widehat{\mathbf{\Gamma}} = \mathcal{P}_{\beta}(\mathbf{D}^{\mathrm{T}}\mathbf{Y})$ -2-4-6-8

-10

-8

-6 -4

-2

2



Michael Elad The Computer-Science Department The Technion

8

10

Consider this for Solving the DCP

 \circ Layered thresholding (LT):

Estimate Γ_1 via the THR algorithm

$$\widehat{\boldsymbol{\Gamma}}_{2} = \mathcal{P}_{\beta_{2}}\left(\boldsymbol{D}_{2}^{\mathrm{T}}\mathcal{P}_{\beta_{1}}\left(\boldsymbol{D}_{1}^{\mathrm{T}}\boldsymbol{Y}\right)\right)$$

Estimate $\Gamma_{\!2}$ via the THR algorithm

○ Forward pass of CNN:

 $f(\mathbf{X}) = \text{ReLU}(\mathbf{b}_2 + \mathbf{W}_2^{\mathrm{T}} \text{ReLU}(\mathbf{b}_1 + \mathbf{W}_1^{\mathrm{T}} \mathbf{Y}))$

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



$$\begin{pmatrix} \mathbf{D}\mathbf{C}\mathbf{P}_{\lambda}^{\mathcal{E}} \end{pmatrix}: \text{ Find } \left\{ \mathbf{\Gamma}_{j} \right\}_{j=1}^{K} \quad s. t. \\ \begin{cases} \|\mathbf{Y} - \mathbf{D}_{1}\mathbf{\Gamma}_{1}\|_{2} \leq \mathcal{E} & \|\mathbf{\Gamma}_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{1} = \mathbf{D}_{2}\mathbf{\Gamma}_{2} & \|\mathbf{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \mathbf{\Gamma}_{K-1} = \mathbf{D}_{K}\mathbf{\Gamma}_{K} & \|\mathbf{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{pmatrix}$$

Consider this for Solving the DLP

 $O \text{ DLP (supervised}^*): \\ \min_{\{\mathbf{D}_i\}_{i=1}^{K}, \mathbf{U}} \sum_{j} \ell\left(h(\mathbf{Y}_j), \mathbf{U}, \frac{\mathbf{D}\mathbf{CP}^*(\mathbf{Y}_j, \{\mathbf{D}_i\})}{\mathbf{U}_j}\right)$

The thresholds for the DCP should also learned

Estimate via the layered THR algorithm

 \odot CNN training:

$$\min_{\{\mathbf{W}_i\},\{\mathbf{b}_i\},U}\sum_{j}\ell\left(h(\mathbf{Y}_j),\mathbf{U},f(\mathbf{Y},\{\mathbf{W}_i\},\{\mathbf{b}_i\})\right)$$

The problem solved by the training stage of CNN and the DLP are equivalent as well, assuming that the DCP is approximated via the layered thresholding algorithm

 Recall that for the ML-CSC, there exists an unsupervised avenue for training the dictionaries that has no simple parallel in CNN



Theoretical Path



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



Theoretical Path: Possible Questions

- Having established the importance of the ML-CSC model and its associated pursuit, the DCP problem, we now turn to its analysis
- \odot The main questions we aim to address:
 - I. Uniqueness of the solution (set of representations) to the (DCP_{λ}) ?
 - II. Stability of the solution to the $(\mathbf{DCP}_{\lambda}^{\mathcal{E}})$ problem ?
 - III. Stability of the solution obtained via the hard and soft layered THR algorithms (forward pass) ?
 - IV. Limitations of this (very simple) algorithm and alternative pursuit?
 - V. Algorithms for training the dictionaries $\{\mathbf{D}_i\}_{i=1}^{K}$ vs. CNN ? VI. New insights on how to operate on signals via CNN ?



Uniqueness of (DCP_{λ})



The feature maps CNN aims to recover are unique



Stability of $(\mathbf{D}\mathbf{C}\mathbf{P}_{\lambda}^{\mathcal{E}})$

 \circ The problem we aim to solve is this

$$\begin{split} \left(\mathbf{D}\mathbf{C}\mathbf{P}_{\lambda}^{\mathcal{E}} \right) &: \text{Find a set of representations satisfying} \\ \|\mathbf{Y} - \mathbf{D}_{1}\mathbf{\Gamma}_{1}\|_{2} \leq \mathcal{E} \quad \|\mathbf{\Gamma}_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{1} &= \mathbf{D}_{2}\mathbf{\Gamma}_{2} \quad \|\mathbf{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ &\vdots &\vdots \\ \mathbf{\Gamma}_{K-1} &= \mathbf{D}_{K}\mathbf{\Gamma}_{K} \quad \|\mathbf{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{split}$$





• The question we pose is How close is $\widehat{\Gamma}_i$ to Γ_i ?


Stability of $(\mathbf{D}\mathbf{C}\mathbf{P}_{\lambda}^{\mathcal{E}})$

Theorem: If the true representations $\{\Gamma_i\}_{i=1}^K$ satisfy $\|\boldsymbol{\Gamma}_{i}\|_{0,\infty}^{s} \leq \lambda_{i} < \frac{1}{2} \left(1 + \frac{1}{\mu(\boldsymbol{D}_{i})}\right)$ then the set of solutions $\left\{ \widehat{\mathbf{\Gamma}}_{\mathbf{i}}
ight\}_{\mathbf{i=1}}^{\mathrm{K}}$ obtained by solving this problem (somehow) must obey $\left\|\widehat{\Gamma}_{i} - \Gamma_{i}\right\|_{2}^{2} \leq \mathcal{E}_{i}^{2}$ for $\mathcal{E}_{0}^{2} = 4\mathcal{E}^{2}, \qquad \mathcal{E}_{i}^{2} = \frac{\mathcal{E}_{i-1}^{2}}{1 - (2\lambda_{i} - 1)\mu(\mathbf{D}_{i})}$

The problem CNN aims to solve is stable under certain conditions

Observe this annoying effect of error magnification as we dive into the model



Stability of Layered-THR

$$\begin{split} \text{Theorem: If } \|\boldsymbol{\Gamma}_{i}\|_{0,\infty}^{s} &< \frac{1}{2} \left(1 + \frac{1}{\mu(\boldsymbol{D}_{i})} \cdot \frac{\left|\boldsymbol{\Gamma}_{i}^{min}\right|}{\left|\boldsymbol{\Gamma}_{i}^{max}\right|}\right) - \frac{1}{\mu(\boldsymbol{D}_{i})} \cdot \frac{\boldsymbol{\epsilon}_{L}^{i-1}}{\left|\boldsymbol{\Gamma}_{i}^{max}\right|} \\ \text{then the layered hard THR (with the proper thresholds) will } \\ \text{find the correct supports and} \\ \left\|\boldsymbol{\Gamma}_{i}^{LT} - \boldsymbol{\Gamma}_{i}\right\|_{2,\infty}^{p} \leq \boldsymbol{\epsilon}_{L}^{i} \\ \text{where we have defined } \boldsymbol{\epsilon}_{L}^{0} = \|\boldsymbol{E}\|_{2,\infty}^{p} \text{ and} \\ \boldsymbol{\epsilon}_{L}^{i} = \sqrt{\|\boldsymbol{\Gamma}_{i}\|_{0,\infty}^{p}} \cdot \left(\boldsymbol{\epsilon}_{L}^{i-1} + \mu(\boldsymbol{D}_{i})\left(\|\boldsymbol{\Gamma}_{i}\|_{0,\infty}^{s} - 1\right)|\boldsymbol{\Gamma}_{i}^{max}|\right) \end{split}$$

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



 \circ (**DCP**_{λ}) Noiseless: Find a set of representations satisfying

$$\begin{split} \mathbf{X} &= \mathbf{D}_{1}\mathbf{\Gamma}_{1} & \|\mathbf{\Gamma}_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{1} &= \mathbf{D}_{2}\mathbf{\Gamma}_{2} & \|\mathbf{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \mathbf{\Gamma}_{K-1} &= \mathbf{D}_{K}\mathbf{\Gamma}_{K} & \|\mathbf{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{split}$$

 \odot So far we proposed the Layered THR:

$$\widehat{\mathbf{\Gamma}}_{K} = \mathcal{P}_{\beta_{K}} \left(\mathbf{D}_{K}^{\mathrm{T}} \dots \mathcal{P}_{\beta_{2}} \left(\mathbf{D}_{2}^{\mathrm{T}} \mathcal{P}_{\beta_{1}} \left(\mathbf{D}_{1}^{\mathrm{T}} \mathbf{X} \right) \right) \right)$$

 \odot The motivation is clear – getting close to what CNN use

 However, this is the simplest and weakest pursuit known in the field of sparsity – Can we offer something better?



Layered Basis Pursuit (Noiseless)

 \circ Our Goal: (**DCP**_{λ}): Find a set of representations satisfying

$$\begin{split} \mathbf{X} &= \mathbf{D}_{1}\mathbf{\Gamma}_{1} & \|\mathbf{\Gamma}_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{1} &= \mathbf{D}_{2}\mathbf{\Gamma}_{2} & \|\mathbf{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \mathbf{\Gamma}_{K-1} &= \mathbf{D}_{K}\mathbf{\Gamma}_{K} & \|\mathbf{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{split}$$

• We can propose a Layered Basis Pursuit Algorithm:

$$\Gamma_{1}^{\text{LBP}} = \min_{\Gamma_{1}} \|\Gamma_{1}\|_{1} \text{ s.t. } \mathbf{X} = \mathbf{D}_{1}\Gamma_{1}$$
$$\Gamma_{2}^{\text{LBP}} = \min_{\Gamma_{2}} \|\Gamma_{2}\|_{1} \text{ s.t. } \Gamma_{1}^{\text{LBP}} = \mathbf{D}_{2}\Gamma_{2}$$

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



Guarantee for Success of Layered BP

 As opposed to prior work in CNN, we can do far more than just proposing an algorithm – we can analyze its terms for success:



Theorem: If a set of representations $\{\Gamma_i\}_{i=1}^K$ of the Multi-Layered CSC model satisfy $\|\Gamma_i\|_{0,\infty}^s \leq \lambda_i < \frac{1}{2} \left(1 + \frac{1}{\mu(D_i)}\right)$ then the Layered BP is guaranteed to find them

• Consequences:

- The layered BP can retrieve the underlying representations in the noiseless case, a task in which the forward pass fails to provide
- The Layered-BP's success does not depend on the ratio $|\Gamma_i^{min}|/|\Gamma_i^{max}|$



Layered Basis Pursuit (Noisy)

$$\Gamma_{1}^{\text{LBP}} = \min_{\Gamma_{1}} \frac{1}{2} \|\mathbf{Y} - \mathbf{D}_{1}\Gamma_{1}\|_{2}^{2} + \lambda_{1}\|\Gamma_{1}\|_{1}$$
$$\Gamma_{2}^{\text{LBP}} = \min_{\Gamma_{2}} \frac{1}{2} \|\Gamma_{1}^{\text{LBP}} - \mathbf{D}_{2}\Gamma_{2}\|_{2}^{2} + \lambda_{2}\|\Gamma_{2}\|_{1}$$

We can invoke a result we have seen already, referring to the BP for the CSC model:

RECALL For
$$\mathbf{Y} = \mathbf{D}\mathbf{\Gamma} + \mathbf{E}$$
, if
 $\|\mathbf{\Gamma}\|_{0,\infty}^{s} < \frac{1}{3} \left(1 + \frac{1}{\mu(\mathbf{D})}\right)$
then we are guaranteed that
 $\|\mathbf{\Delta}\|_{2,\infty}^{p} \le 7.5 \ \varepsilon_{L}^{0} \sqrt{\|\mathbf{\Gamma}\|_{0,\infty}^{p}}$



Stability of Layered BP

Theorem: Assuming that $\|\Gamma_i\|_{0,\infty}^s < \frac{1}{3}\left(1 + \frac{1}{\mu(D_i)}\right)$ then for correctly chosen $\{\lambda_i\}_{i=1}^K$ we are guaranteed that

- 1. The support of $\Gamma_i^{
 m LBP}$ is contained in that of Γ_i
- 2. The error is bounded: $\|\boldsymbol{\Gamma}_{i}^{\text{LBP}} \boldsymbol{\Gamma}_{i}\|_{2,\infty}^{p} \leq \varepsilon_{L}^{i}$, where

$$\varepsilon_{\mathrm{L}}^{\mathrm{i}} = 7.5^{\mathrm{i}} \|\mathbf{E}\|_{2,\infty}^{\mathrm{p}} \prod_{\mathrm{j}=1}^{\mathrm{r}} \sqrt{\|\mathbf{\Gamma}_{\mathrm{j}}\|_{0,\infty}^{\mathrm{p}}}$$

3. Every entry in Γ_i greater than $\epsilon_L^i / \sqrt{\|\Gamma_i\|_{0,\infty}^p}$ will be found

Problems:

- .. Contrast
- 2. Error growth
- 3. Error even if no noise



Layered Iterative Thresholding

Layered BP:
$$\Gamma_{j}^{\text{LBP}} = \min_{\Gamma_{j}} \frac{1}{2} \left\| \Gamma_{j-1}^{\text{LBP}} - \mathbf{D}_{j} \Gamma_{j} \right\|_{2}^{2} + \xi_{j} \left\| \Gamma_{j} \right\|_{1}$$

Layered Iterative Soft-Thresholding:

$$\mathbf{\Gamma}_{j}^{t} = S_{\xi_{j}/c_{j}} \left(\mathbf{\Gamma}_{j}^{t-1} + \frac{1}{c_{j}} \mathbf{D}_{j}^{T} (\widehat{\mathbf{\Gamma}}_{j-1} - \mathbf{D}_{j} \mathbf{\Gamma}_{j}^{t-1}) \right) \mathbf{j}$$

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

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



Michael Elad The Computer-Science Department The Technion

*
$$c_i > 0.5 \lambda_{max} (\mathbf{D}_i^T \mathbf{D}_i)$$

Time to Conclude



Michael Elad The Computer-Science Department The Technion

This Talk





A Massive Open Online Course: Coming Up

Search:



Courses - Programs - Schools & Partners About -

Q

Sign In Register

/ Israel X

Sparse Representations in Signal and Image Processing

Learn the theory, tools and algorithms of sparse representations and their impact on signal and image processing.

Start the Professional Certificate Program





Sparse Representations in Signal and Image Processing: Fundamentals Learn about the field of sparse representations by understanding its fundamental heoretical and algorithmic foundations. earn more

Sparse Representations in Image Processing: From Theory to Practice Learn about the deployment of the sparse representation model to signal and image processing.

Starts on October 25, 2017

Enroll Now

I would like to receive email from IsraelX and learn about other offerings related to Sparse Representations in Signal and Image Processing: Fundamentals.

Starts on February 28, 2018

Enroll Now

I would like to receive email from IsraelX and learn about other offerings related to Sparse Representations in Image Processing: From Theory to Practice.



Instructors

Learn more



Michael Elad The Computer-Science Departme Computer The Technion



Yaniv Romano



Michael Flad



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



Michael Elad The Computer-Science Department The Technion