Signal Modeling: From Convolutional Sparse Coding to Convolutional Neural Networks

Vardan Papyan The Computer Science Department Technion – Israel Institute of technology Haifa 32000, Israel



Joint work with





Jeremias Sulam Yaniv Romano Prof. Michael Elad





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

Part I

Motivation and Background





Our Starting Point: Image Denoising



Many image denoising algorithms can be cast as the minimization of an energy function of the form

$$\frac{1}{2} \frac{\|\mathbf{X} - \mathbf{Y}\|_{2}^{2}}{\text{Relation to}} + \begin{array}{c} \mathbf{G}(\mathbf{X}) \\ \mathbf{Prior or} \\ \text{regularization} \end{array}$$





Leading Image Denoising Methods...

are built upon powerful patch-based local models:

- 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]
- CSR: clustering and sparsity on patches [Dong, Li, Lei & Shi '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]







The Sparse-Land Model

• Assumes every patch a linear combination of a few columns, called atoms, from a matrix that is termed a dictionary.





* \boldsymbol{R}_i for 1D signals

Patch Denoising



Consider this Algorithm [Elad & Aharon, '06]



measurements

Prior or regularization





What is Missing?

- Over the years, many researchers kept revisiting this algorithm and the line of thinking behind it, with a clear feeling that the final word has not been said, and that key features are still lacking.
- What is missing? Here is what our group 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]

Beyond all these, a key part that is missing is a **theoretical** backbone for the local model as a way to characterize the **global** unknown image







Missing Theoretical Backbone?

• The core global-local model assumption on X:

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

Every patch in the unknown signal is expected to have a sparse representation w.r.t. the same dictionary $\pmb{\Omega}$

- Questions to consider:
 - I. Who are those signals belonging to this model? Do they exist?
 - II. Under which conditions on Ω would this model be feasible?
 - III. How does one sample from this model?
 - IV. How should we perform pursuit properly (and locally) under this model?
 - V. How should we learn Ω if this is indeed the model?





In this Talk



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

Working Locally Thinking Globally-Part I: Theoretical Guarantees for Convolutional Sparse Coding

Working Locally Thinking Globally-Part II: Stability and Algorithms for Convolutional Sparse Coding

Vardan Papyan, Jeremias Sulam and Michael Elad







Convolutional Sparsity Assumes...

<u></u>





Convolutional Sparse Representation

• Formally, consider the following global sparsity-based model

$$\mathbf{X} = \sum_{i=1}^{m} \mathbf{C}^{i} \mathbf{\Gamma}^{i} = \mathbf{D} \mathbf{\Gamma}$$

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

Cⁱ

• $\Gamma^{i} \in \mathbb{R}^{N}$ are its corresponding coefficients.



Two Interpretations







- 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.

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

Convolutional dictionary

- Several applications were demonstrated:
 - Inpainting [Heide, Heidrich & Wetzstein '15]
 - Super-resolution [Gu, Zuo, Xie, Meng, Feng & Zhang '15]
 - Pattern detection in images and the analysis of instruments in music signals [Mørup, Schmidt & Hansen '08]
- However, little is known regrading its theoretical aspects.



Classical Sparse Theory (Noiseless)

<u>Mutual Coherence</u>: $\mu(\mathbf{D}) = \max_{i \neq i} |\mathbf{d}_i^T \mathbf{d}_j|$

[Donoho & Elad '03]

For a signal
$$\mathbf{X} = \mathbf{D}\mathbf{\Gamma}$$
, if $\|\mathbf{\Gamma}\|_0 < \frac{1}{2}\left(1 + \frac{1}{\mu(\mathbf{D})}\right)$ then this solution is necessarily the sparsest.

[Donoho & Elad '03]

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

- Assuming that m = 2 and n = 64 we have that $\mu(\mathbf{D}) \ge 0.063$.
- As a result, uniqueness and success of pursuits is guaranteed as long as

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

• This is a very pessimistic result!





The Local Representation



the patches of X



The Local Representation



Inherent Positive Properties

✓ A clear global model with a shift invariant local prior.

- Every patch has a sparse representation w.r.t. to a local dictionary Ω .
- ✓ No disagreement on the patch overlaps.

✓ Related to the current common practice of patch averaging.

• The signal can be written as

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

- \mathbf{R}_{i}^{T} puts the patch $\mathbf{\Omega} \mathbf{\gamma}_{i}$ in the i-th location in the *N*-dimensional vector.
- The patch averaging scheme solves the sparse coding problem independently for every patch while convolutional sparse coding seeks for the representations of all the patches together.



The $\ell_{0,\infty}$ Norm and the $\mathbf{P}_{0,\infty}$ Problem

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

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

A global sparse vector is likely if it can represent every patch in the signal sparsely.

The Main Questions We Aim to Address:

- I. Uniqueness of the solution to this problem ?
- II. Guaranteed recovery of the solution via global OMP/BP?



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 globally optimal solution to this problem.

We should be excited about this result and later ones because they pose a local constraint for a global guarantee, and as such, they are far more optimistic compared to the global guarantees





8 non-zeros per stripe can result in $0.06 \cdot N$ non-zeros globally

Recovery Guarantees







From Ideal to Noisy Signals

- So far, we have assumed an ideal signal $X = D\Gamma$.
- However, in practice we usually have $\mathbf{Y} = \mathbf{D}\mathbf{\Gamma} + \mathbf{E}$ where \mathbf{E} is due to noise or model deviations.
- To handle this, we redefine our problem as:

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

- 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. The same recovery done via local operations ?

[Candes & Tao '05]

Stability of $(\mathbf{P}_{0,\infty}^{\epsilon})$ via Stripe-RIP

$$\begin{split} (\mathbf{P}_{0,\infty}^{e}) &: & \min_{\Gamma} \|\|\Gamma\|_{0,\infty}^{s} \text{ s.t. } \|\|\mathbf{Y} - \mathbf{D}\Gamma\|_{2} \leq \epsilon & \Gamma \\ \hline \mathbf{Definition:} \ \mathbf{D} \text{ is said to satisfy Stripe-RIP with constant } \delta_{k} \text{ if:} \\ & (1 - \delta_{k}) \|\Delta\|_{2}^{2} \leq \|\mathbf{D}\Delta\|_{2}^{2} \leq (1 + \delta_{k}) \|\Delta\|_{2}^{2} \\ \hline \text{for any vector } \Delta \text{ with } \|\Delta\|_{0,\infty}^{s} = k. \end{split}$$

Theorem: If the true representation Γ satisfies $\|\Gamma\|_{0,\infty}^{s} = k < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D})}\right)$ Then a solution $\widehat{\Gamma}$ for $(\mathbf{P}_{0,\infty}^{\epsilon})$ must be close to it $\|\widehat{\Gamma} - \Gamma\|_{2}^{2} \le \frac{4\epsilon^{2}}{1 - \delta_{2k}} \le \frac{4\epsilon^{2}}{1 - (2k - 1)\mu(\mathbf{D})}$

 $\delta_k \leq (k-1)\mu(\mathbf{D})$



Local Noise Assumption

- Thus far, our analysis relied on the local sparsity of the underlying solution Γ, which was enforced through the ℓ_{0,∞} norm.
- 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.
- 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 $\|\mathbf{\Gamma}\|_{0,\infty}^{s} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D})}\right) - \frac{1}{\mu(\mathbf{D})} \cdot \frac{\|\mathbf{E}\|_{2,\infty}^{p}}{|\mathbf{\Gamma}_{\min}|}$ Then OMP run for $\|\mathbf{\Gamma}\|_{0}$ iterations will 1. Find the correct support 2. $\|\mathbf{\Gamma}_{OMP} - \mathbf{\Gamma}\|_{2}^{2} \le \frac{\|\mathbf{E}\|_{2}^{2}}{1 - (\|\mathbf{\Gamma}\|_{0,\infty}^{s} - 1)\mu(\mathbf{D})}$





Stability of Lagrangian BP

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

<u>Theorem</u>: For $\mathbf{Y} = \mathbf{D}\mathbf{\Gamma} + \mathbf{E}$, if $\xi = 4 \|\mathbf{E}\|_{2,\infty}^p$ and $\|\mathbf{\Gamma}\|_{0,\infty}^s < \frac{1}{3} \left(1 + \frac{1}{\mu(\mathbf{D})}\right)$

Then we are guaranteed that

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





Stability of Lagrangian BP

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

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

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

Then we are guaranteed that

- 1. The support of $\Gamma_{\rm BP}$ is contain
- 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.

and 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]





Proof relies on the work of [Tropp '06] 28

Global Pursuit via Local Processing

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

- Thus far, we have seen that while the CSC is a global model, its theoretical guarantees rely on local properties. Yet this global-local relation can also be exploited for practical purposes. Next, we show how one can solve the global BP problem using only local operations.
- Iterative Soft Thresholding [Blumensath & Davies '08]:

Projection $\Gamma^{t} = S_{\xi/c} \left(\Gamma^{t-1} + \frac{1}{c} D^{T} (Y - D\Gamma^{t-1}) \right)$ onto L_{1} ball Gradient step global aggregation

• This can be equally written as:

local sparse code

 $\forall i \quad \boldsymbol{\alpha}_{i}^{t} = \mathcal{S}_{\xi/c} \left(\boldsymbol{\alpha}_{i}^{t-1} + \boldsymbol{D}_{L}^{T} \boldsymbol{R}_{i} (\boldsymbol{Y} - \boldsymbol{D}\boldsymbol{\Gamma}^{t-1}) \right)$

local dictionary local residual

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

Technion Israel Institute of Technology α_i

Simulation

Details:

- Signal length: N = 300
- Patch size: n = 25
- Unique atoms: p = 5
- Global sparsity: k = 40
- Number of iterations: 400
- Lagrangian: $\xi = 4 \|\mathbf{E}\|_{2,\infty}^{p}$



True Sparse Code
 Iterative Soft Thresholding



Partial Summary of CSC

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

Uniqueness of the solution to the noiseless problem.



Stability of the solution to 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.



Part III Going Deeper

Convolutional Neural Networks Analyzed via Convolutional Sparse Coding

Vardan Papyan, Yaniv Romano and Michael Elad







CSC and CNN

- There seems to be a relation between CSC and CNN:
 - Both have a convolutional structure.
 - Both use a data driven approach for training their model.
 - The most popular non-linearity employed in CNN, called ReLU, is known to be connected to sparsity.
- In this part, we aim to make this connection clear in order to provide a theoretical understanding of CNN through the eyes of sparsity.
- 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]





ReLU(z) = max(0, z)

CNN



No 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...

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

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





Training Stage of CNN

- Consider the task of classification, for example.
- Given a set of signals $\{X_j\}_j$ and their corresponding labels $\{h(X_j)\}_j$, the CNN learns an end-to-end mapping.



Back to CSC

 $\mathbf{X} \in \mathbb{R}^{N} \qquad \mathbf{D}_{1} \in \mathbb{R}^{N \times Nm_{1}} \quad \mathbf{\Gamma}_{1} \in \mathbb{R}^{Nm_{1}}$ $\begin{bmatrix} & n_{0} & & \\ & & & & \\ & & & \\ & & & & \\ & & & \\ & & & &$

Convolutional sparsity assumes an inherent structure is present in natural signals. Similarly, the representations themselves could also be assumed to have such a structure.

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



Multi-Layer CSC (ML-CSC)





 $\Gamma_2 \in \mathbb{R}^{Nm_2}$

Deep Coding and Learning Problems





Deep Coding and Learning Problems

大致

$$\begin{array}{l} \left(DCP_{\lambda}^{\mathcal{E}} \right) : \text{Find a set of representations satisfying} \\ \| Y - D_1 \Gamma_1 \|_2 \leq \mathcal{E}_0 & \| \Gamma_1 \|_{0,\infty}^s \leq \lambda_1 \\ \| \Gamma_1 - D_2 \Gamma_2 \|_2 \leq \mathcal{E}_1 & \| \Gamma_2 \|_{0,\infty}^s \leq \lambda_2 \\ \vdots & \vdots \\ \| \Gamma_{K-1} - D_K \Gamma_K \|_2 \leq \mathcal{E}_{K-1} & \| \Gamma_K \|_{0,\infty}^s \leq \lambda_K \\ \end{array}$$

$$\begin{array}{l} \left(DLP_{\lambda}^{\mathcal{E}} \right) : & \min_{\{D_i\}_{i=1}^K, U} \sum_j \ell \left(h(X_j), U, DCP^*(Y_j, \{D_i\}) \right) \\ & & & & & \\ \end{array}$$

$$True \ \text{label} \ \ \text{Classifier} \quad \text{Deepest representation solution by solving the DCP} \\ \end{array}$$



Keep it simple!

The simplest pursuit in the sparse representation is the thresholding algorithm. Given an input signal **X**, this operates by:





Consider this for Solving the DCP

• Layered thresholding (LT):

Estimate $\Gamma_{\!1}$ via the thresholding algorithm

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

Estimate Γ_2 via the thresholding algorithm

• Forward pass of CNN:

 $\begin{array}{ll} (\mathbf{D}\mathbf{C}\mathbf{P}_{\lambda}) \text{: Find a set of} \\ \text{representations satisfying} \\ \mathbf{X} = \mathbf{D}_{1}\mathbf{\Gamma}_{1} & \|\mathbf{\Gamma}_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \mathbf{\Gamma}_{1} = \mathbf{D}_{2}\mathbf{\Gamma}_{2} & \|\mathbf{\Gamma}_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ & \vdots & \vdots \\ \mathbf{\Gamma}_{K-1} = \mathbf{D}_{K}\mathbf{\Gamma}_{K} & \|\mathbf{\Gamma}_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{array}$

 $f(\mathbf{X}) = \operatorname{ReLU}\left(\mathbf{b}_{2} + \mathbf{W}_{2}^{\mathrm{T}}\operatorname{ReLU}\left(\mathbf{b}_{1} + \mathbf{W}_{1}^{\mathrm{T}}\mathbf{X}\right)\right)$

The layered (soft nonnegative) thresholding and the forward pass algorithm are equal !!!





Consider this for Solving the DLP

• DLP:

$$\min_{\{\mathbf{D}_i\}_{i=1}^{K}, \mathbf{U}} \sum_{j} \ell\left(h(\mathbf{X}_j), \mathbf{U}, \frac{\mathbf{D}\mathbf{C}\mathbf{P}^{\star}(\mathbf{X}_j, \{\mathbf{D}_i\})}{\sqrt{2}}\right)$$

Estimate via the layered thresholding algorithm

• CNN training:

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

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





Theoretical Questions





Uniqueness of (DCP_{λ})

 (\mathbf{DCP}_{λ}) : Find a set of representations satisfying $\mathbf{X} = \mathbf{D}_1 \mathbf{\Gamma}_1 \qquad \|\mathbf{\Gamma}_1\|_{0,\infty}^{\mathsf{s}} \le \lambda_1$ ls this set $\Gamma_1 = \mathbf{D}_2 \Gamma_2 \qquad \|\Gamma_2\|_{0,\infty}^{s} \leq \overline{\lambda_2}$ unique? $\Gamma_{K-1} = \mathbf{D}_{K}\Gamma_{K} \quad \|\Gamma_{K}\|_{0,\infty}^{s} \leq \lambda_{K}$ **<u>Theorem</u>**: If a set of solutions $\{\Gamma_i\}_{i=1}^K$ is found for (**DCP** $_{\lambda}$) such that: $\|\boldsymbol{\Gamma}_{i}\|_{0,\infty}^{s} \leq \lambda_{i} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D}_{i})}\right)$ Then these are necessarily the unique solution to this problem.

The feature maps CNN aims to recover are unique





[Papyan, Sulam & Elad '16]

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

Theorem: If the true representations $\{\Gamma_i\}_{i=1}^K$ satisfy $\|\mathbf{\Gamma}_{\mathbf{i}}\|_{0,\infty}^{\mathbf{s}} \leq \lambda_{\mathbf{i}} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D}_{\mathbf{i}})}\right)$ And the error thresholds for $(\mathbf{DCP}_{\lambda}^{\mathcal{E}})$ are $\mathcal{E}_{0}^{2} = \|\mathbf{E}\|_{2}^{2}, \qquad \mathcal{E}_{i}^{2} = \frac{4\mathcal{E}_{i-1}^{2}}{1 - (2\|\boldsymbol{\Gamma}_{i}\|_{0,\infty}^{s} - 1)\mu(\mathbf{D}_{i})}$ Then the set of solutions $\{\widehat{\Gamma}_i\}_{i=1}^K$ obtained by solving this problem must be close to the true ones $\|\widehat{\Gamma}_{i} - \Gamma_{i}\|_{2}^{2} \leq \mathcal{E}_{i}^{2}$

The problem CNN aims to solve is stable under certain conditions





[Papyan, Sulam & Elad '16] 49

Stability of LT

$$\begin{split} \underline{\text{Theorem}} &: \text{If } \|\Gamma_{i}\|_{0,\infty}^{s} < \frac{1}{2} \left(1 + \frac{1}{\mu(D_{i})} \cdot \frac{\left|\Gamma_{i}^{min}\right|}{\left|\Gamma_{i}^{max}\right|}\right) - \frac{1}{\mu(D_{i})} \cdot \frac{\epsilon_{L}^{i-1}}{\left|\Gamma_{i}^{max}\right|} \\ \text{Then the layered soft thresholding will*} \\ 1. & \text{Find the correct supports} \\ 2. & \left\|\Gamma_{i}^{LT} - \Gamma_{i}\right\|_{2,\infty}^{p} \leq \epsilon_{L}^{i} \\ \text{We have defined } \epsilon_{L}^{0} = \|\mathbf{E}\|_{2,\infty}^{p} \text{ and} \\ \epsilon_{L}^{i} = \sqrt{\|\Gamma_{i}\|_{0,\infty}^{p}} \cdot \left(\epsilon_{L}^{i-1} + \mu(\mathbf{D}_{i})\left(\|\Gamma_{i}\|_{0,\infty}^{s} - 1\right)|\Gamma_{i}^{max}| \end{split}$$

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



* For correctly chosen thresholds

Limitations of the Forward Pass

 The stability analysis reveals several inherent limitations of the forward pass algorithm:



Even in the noiseless case, it is incapable of recovering the solution of the DCP problem.



Its success depends on the ratio $|\Gamma_i^{\min}|/|\Gamma_i^{\max}|$. This is a direct

consequence of the forward pass algorithm relying on the simple thresholding operator.



The distance between the true sparse vector and the estimated one increases exponentially as function of the layer depth.

 In the next and final part we propose a new algorithm attempting to solve some of these problems.



Part IV What Next?

Convolutional Neural Networks Analyzed via Convolutional Sparse Coding

Vardan Papyan, Yaniv Romano and Michael Elad







Layered Basis Pursuit (Noiseless)

- Our Goal: $\begin{array}{ll} (DCP_{\lambda}) : \text{Find a set of representations satisfying} \\ X = D_{1}\Gamma_{1} & \|\Gamma_{1}\|_{0,\infty}^{s} \leq \lambda_{1} \\ \Gamma_{1} = D_{2}\Gamma_{2} & \|\Gamma_{2}\|_{0,\infty}^{s} \leq \lambda_{2} \\ \vdots & \vdots \\ \Gamma_{K-1} = D_{K}\Gamma_{K} & \|\Gamma_{K}\|_{0,\infty}^{s} \leq \lambda_{K} \end{array}$
- Layered thresholding: $\widehat{\Gamma}_2 = \mathcal{P}_{\beta_2} \left(\mathbf{D}_2^T \, \mathcal{P}_{\beta_1} (\mathbf{D}_1^T \mathbf{X}) \right)$
- Thresholding is the simplest pursuit known in the field of sparsity.

• Layered BP: $\Gamma_1^{\text{LBP}} = \min_{\Gamma_1} \|\Gamma_1\|_1$ s.t. $\mathbf{X} = \mathbf{D}_1\Gamma_1$ $\Gamma_2^{\text{LBP}} = \min_{\Gamma_1} \|\Gamma_2\|_1$ s.t. $\Gamma_1^{\text{LBP}} = \mathbf{D}_2\Gamma_2$ Deconvolution [Zeiler, Krish

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



Guarantee for Success of Layered BP

 (\mathbf{DCP}_{λ}) : Find a set of representations satisfying • Our Goal: $\mathbf{X} = \mathbf{D}_1 \mathbf{\Gamma}_1 \qquad \|\mathbf{\Gamma}_1\|_{0,\infty}^{\mathsf{s}} \le \lambda_1$ $\Gamma_1 = \mathbf{D}_2 \Gamma_2 \qquad \|\Gamma_2\|_{0,\infty}^{s} \le \lambda_2$ $\Gamma_{K-1} = \mathbf{D}_{K}\Gamma_{K} \quad \|\Gamma_{K}\|_{0,\infty}^{s} \leq \lambda_{K}$ **<u>Theorem</u>**: If a set of solutions $\{\Gamma_i\}_{i=1}^K$ of (DCP_{λ}) satisfy $\|\mathbf{\Gamma}_{\mathbf{i}}\|_{0,\infty}^{\mathbf{s}} \leq \lambda_{\mathbf{i}} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{D}_{\mathbf{i}})}\right)$ Then the layered BP is guaranteed to find them.

✓ The layered BP can retrieve the underlying representations in the noiseless case, a task in which the forward pass fails.

✓ Its success does not depend on the ratio $|\Gamma_i^{\min}| / |\Gamma_i^{\max}|$.



Stability of Layered BP

Theorem: Assuming that

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

Then we are guaranteed that*

- 1. The support of Γ_i^{LBP} is contained in that of Γ_i
- 2. $\left\| \boldsymbol{\Gamma}_{i}^{LBP} \boldsymbol{\Gamma}_{i} \right\|_{2,\infty} \leq \varepsilon_{L}^{i}$
- 3. Every entry in Γ_i greater than $\varepsilon_L^i / (\|\Gamma_i\|_{0,\infty}^p)$ will be found

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

* For correctly chosen $\{\xi_i\}_{i=1}^K$

[Papyan, Sulam & Elad '16] 58

Layered Iterative Thresholding

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

$$\Gamma_{2}^{\text{LBP}} = \min_{\Gamma_{2}} \frac{1}{2} \|\Gamma_{1}^{\text{LBP}} - \mathbf{D}_{2}\Gamma_{2}\|_{2}^{2} + \xi_{2}\|\Gamma_{2}\|_{1}$$
Can be seen as a recurrent neural network
[Gregor & LeCun '10]
$$\Gamma_{1}^{\text{t}} = S_{\Gamma_{1}} \left(\Gamma_{1}^{\text{t}-1} + \frac{1}{2}\mathbf{D}_{1}^{\text{T}}(\mathbf{Y} - \mathbf{D}_{1}\Gamma_{1}^{\text{t}-1})\right)$$

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

IT

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

$$\Gamma_{2}^{t} = S_{\xi_{2}/c_{2}} \left(\Gamma_{2}^{t-1} + \frac{1}{c_{2}} \mathbf{D}_{2}^{T} (\widehat{\Gamma}_{1} - \mathbf{D}_{2} \Gamma_{2}^{t-1}) \right)$$



RA

Conclusion



We described the limitations of patch based processing as a motivation for the CSC model.



We then presented a theoretical study of this model both in a noiseless and a noisy setting.



A multi-layer extension for it, tightly connected to CNN, was proposed and similarly analyzed.



Finally, an alternative to the forward pass algorithm was presented.

Future Work: leveraging theoretical insights into practical implications

THEORY INTO PRACTICE





Questions?



Æ