The Dichotomy Between Global Processing & Local Modeling *

Michael Elad The Computer Science Department The Technion – Israel Institute of technology Haifa 32000, Israel

December 7-11, 2015 Technical University of Berlin

2. International Matheon Conference on Compressed Sensing and its Applications



Global Locca

* Joint work with





Dima Batenkov

Jere. Sulam





Vardan Papyan

Yaniv Romano



The research leading to these results has received funding from the European Research Council under European Union's Seventh Framework Program, ERC Grant agreement no. 320649, and by the Intel Collaborative Research Institute for Computational Intelligence



Part I

Motivating this Discussion



Our Starting Point: Image Denoising



Many of the proposed image denoising algorithms are cast as the minimization of an energy function of the form

$$f(\underline{X}) = \frac{1}{2} \|\underline{X} - \underline{Y}\|_{2}^{2} + \begin{array}{c} G(\underline{X}) \\ \text{Relation to} \\ \text{measurements} \end{array} + \begin{array}{c} G(\underline{X}) \\ \text{Prior or} \\ \text{regularization} \end{array} + \begin{array}{c} G(\underline{X}) \\ \underline{Y}: \text{Given measurements} \\ \underline{X}: \text{Unknown to be recovered} \end{array}$$



Leading Image Denoising Methods

are built upon powerful patch-based (local) image models:

- K-SVD: sparse representation modeling of image patches [Elad & Aharon, '06]
- BM3D: combines a sparsity and self-similarity [Dabov, Foi, Katkovnik & Egiazarian '07]
- EPLL: uses GMM of the image patches [Yu, Sapiro & Mallat '10] [Zoran & Weiss '11]
- CSR: clustering and sparsity on patches [Dong, Li, Lei & Shi '11]

IN THIS TALK

- ❑ We aim to dive into this strange choice of modeling signals/images locally for regularizing global inverse problems. What is the rationale behind this? Is this enough? Can it be improved?
- □ We shall start with a practical view and gradually move to a theory.

□ We will have more open question than answers ...



Consider this Algorithm [Elad & Aharon, '06]



- □ This method is very effective for image denoising. Many variants of it were developed over the years in order to extend/improve it.
- □ The above is only one in the large family **of patch-based algorithms** that denoise an image by decomposing it into patches, processing them separately and then merging them back by **plain averaging** to form the final outcome.



This Algorithm's Origin

location

$$\begin{split} \hat{\underline{X}} &= \underset{\underline{X}, \{\underline{\alpha}_{k}\}_{k}, \mathbf{D}}{\operatorname{ArgMin}} \frac{1}{2} \|\underline{X} - \underline{Y}\|_{2}^{2} + \mu \underset{k}{\sum} \|\mathbf{R}_{k}\underline{X} - \mathbf{D}\underline{\alpha}_{k}\|_{2}^{2} \text{ s.t. } \|\underline{\alpha}_{k}\|_{0} \leq L \\ \Box \text{ The expression } \mathbf{R}_{k}\underline{X} \text{ extracts a patch of size n from } \underline{X} \in \mathbb{R}^{N} : \\ \mathbf{R}_{k} &= \left[\underbrace{k}_{k} \operatorname{th} \underbrace{k} \operatorname{th} \underbrace{k}_{k} \operatorname{th} \underbrace{k}_{k} \operatorname{th} \underbrace{k} \operatorname{th} \operatorname{th} \underbrace{k} \operatorname{th} \underbrace{k} \operatorname{th} \underbrace{k} \operatorname{th} \underbrace{k} \operatorname{th} \operatorname{th} \underbrace{k} \operatorname{th} \operatorname{th} \underbrace{k} \operatorname{th} \operatorname{th} \operatorname{th} \operatorname{th} \underbrace{k} \operatorname{th} \operatorname{th$$

 $\Box \mathbf{R}_k^T \underline{z}$ puts the patch \underline{z} in the kth location in the N-dim. vector.

 \Box Considering cyclic patch-extraction, we have: $\frac{1}{n}\sum \mathbf{R}_{k}^{\mathsf{T}}\mathbf{R}_{k} = \mathbf{I}$



Κ

This Algorithm's Origin

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

The Proposed Regularization / Model:

Every patch in the unknown image is expected to have a sparse representation w.r.t. the dictionary **D**

- This algorithm seeks the "most appropriate" dictionary to fulfil this expectation, along with the sparse representations.
- □ The denoising itself is obtained by projecting each of the patches to this model (via OMP).
- \Box When optimizing over <u>X</u>, this amounts to plain patch-averaging.



So, What is Missing?

- □ Over the past several 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 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]²
 - 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 unknown image.



Theoretical Backbone?

The core model Assumption on \underline{X} :

$$\forall \mathbf{k} \quad \mathbf{R}_{k} \underline{X} = \mathbf{D} \underline{\alpha}_{k} \text{ where } \|\underline{\alpha}_{k}\|_{0} \leq \mathbf{L}$$

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

Questions to consider:

- □ Who are those signals belonging to this model?
- □ Under which conditions on **D** would this model be feasible?
- □ How does one sample from this model?
- □ How should we perform pursuit properly (& locally) under this model?
- □ How should we learn **D** if this is indeed the model?





Why Sampling ?

Surely, you are familiar with this line of work ...

Generate a random vector $\underline{\alpha}$ of length m and with L<<m non-zeros

Multiply $\underline{\alpha}$ by the dictionary **D** of size n-by-m, and obtain $\underline{x} = \mathbf{D}\underline{\alpha}$

Add noise (WAGN) with STD= σ to the signal and get $\underline{y} = \underline{x} + \underline{v}$

Questions:

n

Given \underline{y} (& L, σ , and some properties of **D**),

m

(i) Can we recover the true support of <u>α</u>?
(ii) How efficiently can we denoise <u>y</u>?



Why Sampling ?

Generate ignal X
such the earlier fits
patches has parse
representation
$$\forall k \ \mathbf{R}_k \underline{X} = \mathbf{D}\underline{\alpha}_k$$

where $\|\underline{\alpha}_k\|_0 \leq L$

Add noise (WAGN) with STD= σ to the signal and get $\underline{Y}=\underline{X}+\underline{V}$



Questions:

Given <u>Y</u> (& L, σ , and some properties of **D**),

- (i) Can we recover the true supports of $\underline{\alpha}_{k}$?
- (ii) How efficiently can we denoise <u>Y</u>?What can the oracle do?

(iii) Can we do all this by local processing?



This Talk



For such signals, we address fundamental questions

that are of great relevance to image processing.



Part II

Toy Problem: The Gaussian Case



Problem Definition

□ We are given a Gaussian signal $\underline{X} \in \mathbb{R}^N$ of known statistics, $\underline{X} \sim N\{\underline{0}, \Sigma\}$ contaminated by WAGN:

$$\underline{\mathbf{Y}} = \underline{\mathbf{X}} + \underline{\mathbf{V}}, \quad \underline{\mathbf{V}} \sim \mathbf{N}\left\{\underline{\mathbf{0}}, \mathbf{I}\right\}$$

\Box Our goal – denoise <u>Y</u> to get as close as possible to <u>X</u>.

* For simplicity we assume zero mean signal X and a unit variance noise.





Sounds Easy – Wiener Filtering !

□ Use the Wiener Filter for the denoising, as it gives the Minimum Mean Square Error (MMSE) result:

$$\min_{\underline{X}} \frac{1}{2} \|\underline{Y} - \underline{X}\|_{2}^{2} + \frac{1}{2} \underline{X}^{\mathsf{T}} \Sigma^{-1} \underline{X}$$
This is -log P(X|Y)
$$\hat{\underline{X}}_{\text{Global}} = \Sigma \left(\mathbf{I} + \Sigma\right)^{-1} \underline{Y}$$

□ This is the best we could perform! However,

- It calls for knowing Σ , and
- Inverting a matrix of size N-by-N.



Sounds Easy – Wiener Filtering !

For example, using the following Σ (N=1000), the obtained linear filter is this:

We are choosing Σ to be a circulant matrix for reasons that will be clear shortly





Logarithmic Scale

□ Are there simpler alternatives?



Local Processing by Patch-Averaging

□ A patch of length n<<N extracted from <u>X</u> by $\mathbf{R}_k \underline{X}$ is also Gaussian-distributed, obtained as the marginal distribution:

$$\underline{\mathbf{X}}_{k} = \mathbf{R}_{k} \underline{\mathbf{X}} \sim \mathbf{N} \Big\{ \underline{\mathbf{0}}, \mathbf{R}_{k} \boldsymbol{\Sigma} \mathbf{R}_{k}^{\mathsf{T}} \Big\}$$

We could apply the Wiener filter to each of these patches (overlapped) and then seek a global solution that best fits these local results:

$$\min_{\underline{X}} \sum_{k} \left\| \mathbf{R}_{k} \underline{X} - \underline{\hat{\mathbf{X}}}_{k} \right\|_{2}^{2}$$



Local Processing by Patch-Averaging

□ This process amounts to a simple aggregation by averaging:

$$\hat{\underline{X}}_{k} = \mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{\mathsf{T}} \left(\mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{\mathsf{T}} + \mathbf{I} \right)^{-1} \mathbf{R}_{k} \underline{Y}$$

$$\hat{\underline{X}}_{LPA} = \frac{1}{n} \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \cdot \hat{\underline{X}}_{k}$$

$$= \frac{1}{n} \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \cdot \mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{\mathsf{T}} \left(\mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{\mathsf{T}} + \mathbf{I} \right)^{-1} \mathbf{R}_{k} \underline{Y}$$



Local Patch-Averaging (LPA)

- □ By assuming a circulant Σ , all local models are the same, which is the typical case practiced in image processing of assuming the same model for all patches.
- □ The obtained linear filter now is this:

$$\hat{\mathbf{X}}_{\mathsf{LPA}} = \frac{1}{n} \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{\mathsf{T}} \left(\mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{\mathsf{T}} + \mathbf{I} \right)^{-1} \mathbf{R}_{k} \mathbf{Y}$$

□ Naturally, this matrix is circulant and banded, with width 2n-1 (n=40).



Could We Do Better ? Yes – EPLL !

□ We could refer to the global signal as the unknown, but impose the local Gaussian models. This is the rational of the EPLL (Expected Patch Log Likelihood) [Zoran & Weiss `10].

□ The MMSE in this case reads

$$\min_{\underline{X}} \frac{1}{2} \|\underline{Y} - \underline{X}\|_{2}^{2} + \frac{c}{2n} \sum_{k} \underline{X}^{\mathsf{T}} \mathbf{R}_{k}^{\mathsf{T}} \left(\mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{\mathsf{T}}\right)^{-1} \mathbf{R}_{k} \underline{X}$$
$$\hat{\underline{X}}_{\mathsf{EPLL}} = \left[I + \frac{c}{n} \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \left(\mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{\mathsf{T}}\right)^{-1} \mathbf{R}_{k}^{\mathsf{T}}\right]^{-1} \mathbf{Y}$$



EPLL via Local Processing (ADMM)

□ Wait! This looks no simpler than the global Wiener, so where is the benefit in thinking locally? $\hat{\underline{X}}_{EPLL} = \left[I + \frac{C}{n} \sum_{k} \mathbf{R}_{k}^{T} \left(\mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{T}\right)^{-1} \mathbf{R}_{k}\right]^{-1} \underline{Y}$

☐ Answer: ADMM

$$\begin{split} &\left\{ \underline{\hat{z}}_{k} = \left(\mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{\mathsf{T}} \right) \left(\mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{\mathsf{T}} + \frac{c}{n\lambda} I \right)^{-1} \left(\mathbf{R}_{k} \underline{\hat{X}} - \underline{\hat{u}}_{k} \right) \right\}_{k=1,2,\dots} \text{A Wiener filter on each patch} \\ & \underline{\hat{X}} = \left(I + \lambda \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \mathbf{R}_{k} \right)^{-1} \left(\underline{Y} + \lambda \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \left(\underline{\hat{z}}_{k} + \underline{\hat{u}}_{k} \right) \right) \quad \text{: Patch averaging} \\ & \underline{\hat{u}}_{k} = \underline{\hat{u}}_{k} - \mathbf{R}_{k} \underline{\hat{X}} + \underline{\hat{z}}_{k} \end{split}$$

 \Box Notice that the first iteration (with <u>u</u>=<u>0</u>) coincides with LPA.



EPLL – Back to the Example





Going Multi-Scale ...

One could work with fixed-sized patches, related to different scales of the signal. This may enable better proximity to the global filter. Scale invariance may prove valuable ...

$$\Box \text{ For example } \dots \quad \min_{\underline{X}} \quad \frac{1}{2} \|\underline{Y} - \underline{X}\|_{2}^{2} + \frac{c_{1}}{2n} \sum_{k} \underline{X}^{\mathsf{T}} \mathbf{R}_{k}^{\mathsf{T}} \left(\mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{\mathsf{T}}\right)^{-1} \mathbf{R}_{k} \underline{X} \\ + \frac{c_{2}}{2n} \sum_{k} \underline{X}^{\mathsf{T}} \mathbf{Q}_{k}^{\mathsf{T}} \left(\mathbf{Q}_{k} \Sigma \mathbf{Q}_{k}^{\mathsf{T}}\right)^{-1} \mathbf{Q}_{k} \underline{X}$$

- **R**_k extracts a patch of length n from the signal
- Q_k extracts a patch of size 2n from the signal and then reduce it to length n by filtering followed by 2:1 decimation

□ Naturally, this can be repeated in several scales ...



Going Multi-Scale ...

$$\hat{\underline{X}}_{\text{Global}} = \underline{\Sigma \left(\mathbf{I} + \Sigma \right)^{-1} \underline{Y}}$$



$$\hat{\underline{X}}_{EPLL} = \left[I + \frac{c}{n} \sum_{k} \mathbf{R}_{k}^{T} \left(\mathbf{R}_{k} \Sigma \mathbf{R}_{k}^{T} \right)^{-1} \mathbf{R}_{k} \right]^{-1} \underline{Y}$$



 $\underline{\mathsf{Y}}$

Going Multi-Scale ...

These are the absolute error matrices of the filters versus the global Wiener

True Scale





Average Denoising Performance





The Gaussian Case: A Summary

What can be learned from this toy problem?

- □ The signal \underline{X} has a clear global statistical model, which induces the ultimate denoising approach the Wiener filtering.
- Nevertheless, we may handle the denoising task by operating on overlapped patches, using only local marginals, and by operating locally we can get to near-ideal performance.
- □ Under some conditions on Σ , the local and the global methods may lead to the same performance.
- □ In image processing, identifying/formulating a global model is simply impossible, while learning a local model is within reach.
- □ We are using local marginals to "characterize" (but not to reconstruct) the global distribution.



Part III

Back to Sparse Representations: Who Are These Signals?



The Model to Explore

The model Assumption on \underline{X} :

$$\forall \mathbf{k} \quad \mathbf{R}_{k} \underline{X} = \mathbf{D} \underline{\alpha}_{k} \text{ where } \|\underline{\alpha}_{k}\|_{0} \leq L$$

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

Questions to consider:

- □ Who are those signals belonging to this model?
- □ Under which conditions on **D** would this model be feasible?
- □ How does one sample from this model?
- □ How should we perform pursuit properly (& locally) under this model?
- \square How should we learn **D** if this is indeed the model?





Globalizing the Model (1)

□ We start with the model: $\forall k \quad \mathbf{R}_{k} \underline{X} = \mathbf{D}\underline{\alpha}_{k}, \quad \|\underline{\alpha}_{k}\|_{0} \leq \mathbf{L}$

 $\Box \mathbf{R}_k$ extracts a patch of size n from $\underline{X} \in \mathbb{R}^N$:



 $\Box \text{ Exploiting } \frac{1}{n} \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \mathbf{R}_{k} = \mathbf{I} \not \rightarrow \underbrace{X} = \frac{1}{n} \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \mathbf{R}_{k} \underline{X}$ i.e., the global signal is built of averaged local pieces.



Globalizing the Model (2)

$$\Box \text{ Defining } \underline{\Gamma} = \begin{bmatrix} \underline{\alpha}_{1} \\ \underline{\alpha}_{2} \\ \vdots \\ \underline{\alpha}_{N} \end{bmatrix} \xrightarrow{X} = \frac{1}{n} \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \mathbf{R}_{k} \underline{X} = \sum_{k} \frac{1}{n} \mathbf{R}_{k}^{\mathsf{T}} \mathbf{D} \underline{\alpha}_{k}$$
leads to
$$\begin{bmatrix} \mathbf{I} & \mathbf{I} \\ \frac{1}{n} \mathbf{R}_{1}^{\mathsf{T}} \mathbf{D} & \frac{1}{n} \mathbf{R}_{2}^{\mathsf{T}} \mathbf{D} & \cdots & \frac{1}{n} \mathbf{R}_{N}^{\mathsf{T}} \mathbf{D} \\ \mathbf{I} & \mathbf{I} & \mathbf{I} \end{bmatrix} \underline{\Gamma} = \underline{X} = \mathbf{D}_{G} \underline{\Gamma}$$

□ This suggests the existence of a global sparsity-based model of the kind we are familiar with ...



Globalizing the Model (3)





Globalizing the Model (4)

 \Box However, the vector $\underline{\Gamma}$ is structured ...

$$\begin{split} \underline{X} &= \sum_{k} {}^{1}_{\overline{n}} \mathbf{R}_{k}^{\mathsf{T}} \mathbf{D} \underline{\alpha}_{k} = \mathbf{D}_{G} \underline{\Gamma} & \underline{\Gamma} = \\ \forall j \quad \mathbf{R}_{j} \underline{X} &= \sum_{k} {}^{1}_{\overline{n}} \mathbf{R}_{j} \mathbf{R}_{k}^{\mathsf{T}} \mathbf{D} \underline{\alpha}_{k} = \mathbf{R}_{j} \mathbf{D}_{G} \underline{\Gamma} = \mathbf{D} \underline{\alpha}_{j} \\ &= \begin{bmatrix} 0 & \cdots & \mathbf{D} & \cdots & 0 \end{bmatrix} \underline{\Gamma} \end{split}$$





Globalizing the Model (5)

$$\forall j \quad \mathbf{R}_{j}\mathbf{D}_{G}\underline{\Gamma} = \begin{bmatrix} 0 & \cdots & \mathbf{D} & \cdots & 0 \end{bmatrix}\underline{\Gamma}$$

□ This defines N sets of n equations each, which can be expressed as ... $M\Gamma = 0$

where \mathbf{M} is of size $nN \times mN$

□ M is some sort of "Laplacian" in nature, capturing the fact that the averaged patches over the overlaps should have a complete agreement on the content.



Globalizing the Model (6)

□ Another limitation on $\underline{\Gamma}$ is the fact that $\forall k \|\underline{\alpha}_k\|_0 \leq L$, which we will denote as

$$\left\|\underline{\Gamma}\right\|_{0,\infty} = \mathbf{L}$$

A word about our notation: this definition works on blocks of m elements. Later on we will see a similar definition with a sliding window ...

$$\left\|\underline{\Gamma}\right\|_{0,\infty}^{\mathsf{m}} = \mathbf{L} \longleftrightarrow \left\|\underline{\Gamma}\right\|_{0,\infty}^{\mathsf{m-Sliding}} = \mathbf{L}$$



Globalizing the Model – Summary



implying a useless model !

□ Even if this model is feasible, how can we sample from it?

□ All this might seem hopeless, but ... do not despair ...



Special Case 1: PWC Signals



Piece-Wise-Constant signals obey the local model: If every patch of length n contains (up to) L-1 steps, then these patches can be described as

$$\forall \mathbf{k} \quad \mathbf{R}_{\mathbf{k}} \underline{X} = \mathbf{D} \underline{\alpha}_{\mathbf{k}}, \quad \left\| \underline{\alpha}_{\mathbf{k}} \right\|_{0} \leq \mathbf{L}, \quad \mathbf{D} = \mathbf{D} \mathbf{L}$$





Special Case 1: PWC Signals

$$\underline{X} = \mathbf{D}_{G}\underline{\Gamma}, \quad \mathbf{M}\underline{\Gamma} = \underline{0}, \quad \|\underline{\Gamma}\|_{0,\infty} = \mathbf{L}$$

□ In the PWC case, **M** has a null space of dimension N. We denote this null-space by **Z** (of size nN×N),

$$\mathbf{Z}\underline{\Phi} = \underline{\Gamma}$$

□ A proper choice for describing Z leads to the regular global sparse model with the

$$\underline{X} = \mathbf{D}_{G} \underline{\Gamma} = \mathbf{D}_{G} \mathbf{Z} \underline{\Phi}$$
Global Heaviside
dictionary

 \Box The representation $\underline{\Phi}$ satisfies $\left\|\underline{\Phi}\right\|_{0,\infty}^{n-Sliding} \leq L-1$.



Special Case 2: Convolutional Sparsity

□ Consider a global sparsity-based model for signals $\underline{X} = \mathbf{H} \underline{\Phi}$, in which the dictionary is **circulant** and **banded** in the following way :





Special Case 2: Convolutional Sparsity

□ In this case, every patch of length n≤b will have a sparse representation w.r.t. the same local dictionary:





Special Case 2: Convolutional Sparsity

$$\underline{X} = \mathbf{D}_{G}\underline{\Gamma}, \ \mathbf{M}\underline{\Gamma} = \underline{0}, \ \|\underline{\Gamma}\|_{0,\infty} = \mathbf{L}$$

□ In terms of the above, we get that **M** has a null space of dimension N, and the situation is similar to the one encountered in the PWC case.

□ The same holds true for a global signal emerging from a sparsity-based model with a union of circulant and banded matrices

$$\underline{\mathbf{X}} = \sum_{j} \mathbf{H}_{j} \underline{\Phi}_{j}$$

we will come back to this later on ...



Our Global-Local Model: A Summary

Local Model:
$$\forall k \quad \mathbf{R}_{k} \underline{X} = \mathbf{D}_{\underline{\alpha}_{k}}, \quad \left\| \underline{\alpha}_{k} \right\|_{0} \leq \mathbf{L}$$

Globalized Model: $\underline{X} = \mathbf{D}_{G} \underline{\Gamma}, \quad \mathbf{M} \underline{\Gamma} = \underline{0}, \quad \left\| \underline{\Gamma} \right\|_{0,\infty}^{n} = \mathbf{I}$



We have identified special cases of local dictionaries **D** that enable this model.

Furthermore, in these cases we can easily synthesize signals.



Our Global-Local Model: A Summary

Clarification:

The migration from the local to the globalized model is not constructive. So why bother? Because this will serve us later on when we define the pursuit algorithms.

Open Questions :

Are there other local dictionaries D that could enable the local-global model? We do know of other cases ...
 Is it guaranteed that any such permissible D will have a global model?



Part IV Pursuit & Denoising for These Signals



The Model to Explore

The model Assumption on \underline{X} :

$$\forall \mathbf{k} \quad \mathbf{R}_{k} \underline{X} = \mathbf{D} \underline{\alpha}_{k} \text{ where } \|\underline{\alpha}_{k}\|_{0} \leq L$$

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

Questions to consider:

- □ Who are those signals belonging to this model?
- □ Under which conditions on **D** would this model be feasible?
- □ How does one sample from this model?

□ How should we perform pursuit properly (& locally) under this model?

□ How should we learn **D** if this is indeed the model?





\Box The signal <u>X</u> obeys our local model

 $\forall k \quad \mathbf{R}_{k} \underline{X} = \mathbf{D} \underline{\alpha}_{k} \text{ where } \|\underline{\alpha}_{k}\|_{0} \leq L$

or the globalized one,

$$\underline{X} = \mathbf{D}_{\mathbf{G}}\underline{\Gamma}, \quad \mathbf{M}\underline{\Gamma} = \underline{0}, \quad \left\|\underline{\Gamma}\right\|_{0,\infty} = \mathbf{L}$$

□ Given $\underline{Y} = \underline{X} + \underline{V}$, (~ $N\{\underline{0}, \sigma^2 I\}$), our goal is to recover \underline{X} .

□ The specific questions we aim to address are:

- How would the oracle perform in this case?
- How can we perform this denoising using local pursuit and how well will this work?



The Pursuit Goal

The signal is believed to belong to the following model $\underline{X} = \mathbf{D}_{G}\underline{\Gamma}, \ \mathbf{M}\underline{\Gamma} = \underline{0}, \ \|\underline{\Gamma}\|_{0,\infty} = \mathbf{L}$

Pursuit is simply a projection $\min_{\underline{\Gamma}} \| \mathbf{D}_{G} \underline{\Gamma} - \underline{Y} \|_{2}^{2}$ s.t. $\mathbf{M} \underline{\Gamma} = \underline{0}, \| \underline{\Gamma} \|_{0,\infty} = \mathbf{L}$ onto the model

□ This is very similar to the regular pursuit we are accustomed to, and it is just as complex (NP-Hard).

□ Remember: while we write the pursuit in terms of D_G and M, our true goal is to break this into local pursuit steps while being equivalent to this formulation.



The Oracle (1)

□ The oracle knows the locations of the $|S| \le LN$ non-zeros in $\underline{\Gamma}$. The operator extracting these non-zeros is **P** (size $|S| \times mN$). Thus, $\mathbf{P}^{\mathsf{T}} \mathbf{P} \underline{\Gamma} = \mathbf{P}^{\mathsf{T}} \underline{\Gamma}_{s} = \underline{\Gamma}$

 \Box Our pursuit task simplifies to the search of the optimal $\underline{\Gamma}_s$:

$$\min_{\underline{\Gamma}_{s}} \left\| \mathbf{D}_{G} \mathbf{P}^{\mathsf{T}} \underline{\Gamma}_{s} - \underline{\mathbf{Y}} \right\|_{2}^{2} \quad \text{s.t.} \quad \mathbf{M} \mathbf{P}^{\mathsf{T}} \underline{\Gamma}_{s} = \underline{\mathbf{0}}$$

□ Denoting by **B** (size nN×d (d≤|S|)), the null-space of **MP**^T, we have $\underline{\Gamma}_s = \mathbf{B}\underline{\theta}$, and the pursuit becomes

$$\min_{\underline{\theta}} \left\| \mathbf{D}_{\mathbf{G}} \mathbf{P}^{\mathsf{T}} \mathbf{B} \underline{\theta} - \underline{\mathbf{Y}} \right\|_{2}^{2}$$





$$\min_{\underline{\theta}} \left\| \mathbf{D}_{\mathsf{G}} \mathbf{P}^{\mathsf{T}} \mathbf{B} \, \underline{\theta} - \underline{\mathbf{Y}} \right\|_{2}^{2}$$

\Box This matrix is of size N×d, and it is full-rank.

 \Box ... We started with a noise power of N σ^2 in the signal <u>Y</u>, and the oracle ends with noise energy of d σ^2 .

□ Special cases of interest:

□ In the PWC case, d=number of flat regions, while |S|≈d[•]n.
 □ In the convolutional sparsity case (with one circulant matrix),

$$\mathbf{d} = \left\|\underline{\Phi}\right\|_{0}, \ \left|\mathbf{S}\right| = \mathbf{d} \cdot (2\mathbf{n} - 1)$$



Oracle via LPA

□ Under the regime of an oracle (the supports of $\underline{\alpha}_k$ are known), lets apply the Local Patch Averaging (LPA):

$$\min_{\underline{X}} \sum_{k} \left\| \mathbf{R}_{k} \underline{X} - \underbrace{\mathbf{D}_{S_{k}} \mathbf{D}_{S_{k}}^{\dagger} \mathbf{R}_{k} \underline{Y}}_{\frac{\hat{X}_{k}}{\underline{X}_{k}}} \right\|_{2}^{-} =$$

$$\underline{X} = \frac{1}{n} \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \hat{\underline{\mathbf{X}}}_{k} = \frac{1}{n} \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \mathbf{D}_{\mathbf{S}_{k}} \mathbf{D}_{\mathbf{S}_{k}}^{\mathsf{T}} \mathbf{R}_{k} \underline{\mathbf{Y}}$$

Theorem: Applying the oracle's LPA algorithm iteratively leads to the optimal solution.



Oracle via EPLL

□ Under the regime of an oracle (the supports of $\underline{\alpha}_k$ are known), lets apply the EPLL:

$$\min_{\underline{X}} \quad \frac{1}{2} \left\| \underline{Y} - \underline{X} \right\|_{2}^{2} + \frac{\lambda}{2} \sum_{k} \left\| \underbrace{\left(\mathbf{I} - \mathbf{D}_{S_{k}} \mathbf{D}_{S_{k}}^{\dagger} \right) \mathbf{R}_{k} \underline{X}}_{\mathbf{P}_{k}} \right\|_{2}^{2}$$
$$\underbrace{\hat{X}}_{\text{EPLL}} = \left[\mathbf{I} + \lambda \sum_{k} \mathbf{R}_{k}^{\mathsf{T}} \mathbf{P}_{k}^{\mathsf{T}} \mathbf{P}_{k}^{\mathsf{T}} \mathbf{R}_{k} \right]^{-1} \underline{Y}$$

Theorem: Applying the oracle's EPLL **iteratively** with arbitrary $\lambda > 0$ leads to the optimal solution.



The Oracle: A Summary

□ As expected, the "complexity" of the signal (the value d) governs the denoising performance.

The above two results regarding LPA and EPLL are misleading, as they may suggest that it is enough to apply local pursuit (projection onto an appropriate subspace) and average in order to lead to the ideal solution.

While this is indeed the case if the supports are known, as we depart from the oracle regime, we will see that the local pursuits must "communicate" in order to lead to a more successful global recovery.



Patch-Based Pursuit

Recall our pursuit talks:





Patch-Based Pursuit - Demo

Details:

- Signal length: N = 700
- Patch size: n = 25
- Dictionary = PWC
- Global sparsity: k = 25
- Noise with $\sigma = 3$

Results:

- Noise (per sample): σ^2
- ADMM 0.078σ²
- LPA 0.106σ²
- Iter-LPA $0.113\sigma^2$





Part V

A Closer Look at the Convolutional Sparsity (The Noiseless Case)



The Convolutional Sparsity Case

□ Let us return to the special case of convolutional sparsity.

□ We will describe this not via the local model or its globalized form, but rather through the global model.

Consider a global sparsity-based model of the form $\underbrace{X}_{j=1}^{p} \mathbf{H}_{j} \underbrace{\Phi}_{j}$ where \mathbf{H}_{j} are banded and circulant.

\Box For simplicity we assume that b=n.

The **Global** Dictionary





The Local Representations $\underline{\alpha}_k$





The Representation Problem

$$\begin{pmatrix} \mathbf{P}_{0,\infty} \end{pmatrix} \min_{\underline{\Phi}} \left\| \underline{\Phi} \right\|_{0,\infty} \quad \text{s.t. } \underline{X} = \mathbf{H}_{\mathsf{T}} \underline{\Phi}$$

where we have defined $\left\|\underline{\Phi}\right\|_{0,\infty} = \max_{k} \left\|\underline{\alpha}_{k}\right\|_{0}$

The Main Questions We Aim to Address:

- Uniqueness of the solution to this problem ?
- Guaranteed Recovery of the solution via global OMP/BP ?
- The same recovery done via local operations ?



Stripe Spark and Uniqueness

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

Theorem: If a solution is found for $(\mathbf{P}_{0,\infty})$ such that $\left\|\underline{\Phi}\right\|_{0,\infty} < \frac{1}{2}\sigma_{S}\left(\mathbf{H}_{T}\right)$

Then this is necessarily the globally optimal solution to this problem.



Stripe Spark vs. Mutual Coherence

 $\mu(\mathbf{H}_{\mathsf{T}}) \triangleq \max_{i \neq i} \left| \mathbf{H}_{\mathsf{T}}^{\mathsf{T}} \mathbf{H}_{\mathsf{T}} \right|$

This is the classic coherence, defined over the global dictionary, assuming normalized columns

Theorem: The relation between the stripe spark and the global coherence is

 $\sigma_{S}\left(\boldsymbol{H}_{T}\right) \geq 1 + \frac{1}{\mu\left(\boldsymbol{H}_{T}\right)}$

Thus uniqueness of the solution is guaranteed if

$$\left\|\underline{\Phi}\right\|_{0,\infty} < \frac{1}{2} \left(1 + \frac{1}{\mu \left(\mathbf{H}_{T}\right)}\right)$$



Recovery Guarantees

Lets solve this problem via OMP or BP, applied globally

$$\left(\mathbf{P}_{0,\infty} \right) \quad \min_{\underline{\Phi}} \ \left\| \underline{\Phi} \right\|_{0,\infty} \quad \text{s.t. } \underline{X} = \mathbf{H}_{\mathsf{T}} \underline{\Phi}$$

Theorem: If a solution of $(\mathbf{P}_{0,\infty})$ satisfies $\left\|\underline{\Phi}\right\|_{0,\infty} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{H}_{T})}\right)$

Then global OMP and BP are guaranteed to find it.



Local Recovery Guarantee via ADMM

$$\left(\mathbf{P}_{1}^{\varepsilon} \right) \quad \min_{\underline{\Phi}} \lambda \left\| \underline{\Phi} \right\|_{1} + \frac{1}{2} \left\| \underline{X} - \mathbf{H}_{T} \underline{\Phi} \right\|_{2}^{2}$$

.... can be converted to the following equivalent format

$$\min_{\underline{\Phi}, \{\underline{\Phi}_k\}, \{\underline{\alpha}_k\}} \sum_{k} \lambda \left\| \underline{\Phi}_k \right\|_1 + \frac{1}{2} \left\| \mathbf{R}_k \underline{X} - \mathbf{D}\underline{\alpha}_k \right\|_2^2 \quad \text{s.t.} \quad \begin{cases} \underline{\alpha}_k = \mathbf{Q}_k \underline{\Phi} \\ \underline{\Phi}_k = \mathbf{S}_k \underline{\alpha}_k \end{cases}$$

which can be solved iteratively by

- Updating $\underline{\phi}_k$ via simple soft shrinkage.
- Updating $\underline{\alpha}_k$ via simple multiplication by a matrix, and
- Updating $\underline{\Phi}$ via patch-averaging.



Local Recovery Guarantee via ADMM

$$\begin{pmatrix} \mathbf{P}_{1}^{\varepsilon} \end{pmatrix} \min_{\underline{\Phi}} \lambda \|\underline{\Phi}\|_{1} + \frac{1}{2} \|\underline{X} - \mathbf{H}_{T} \underline{\Phi}\|_{2}^{2}$$

Theorem: If the solution of $(\mathbf{P}_{0,\infty})$ satisfies $\left\|\underline{\Phi}\right\|_{0,\infty} < \frac{1}{2} \left(1 + \frac{1}{\mu(\mathbf{H}_{T})}\right)$

Then the ADMM solver of the global BP, with $\lambda \rightarrow 0$, is guaranteed to find it.



Local Pursuit in Action

Details:

- Signal length: N = 300
- Patch size: n = 25
- Dictionary = Convolutional
- Unique atoms: p = 5
- Global sparsity: k = 50
- Number Iterations: 1000
- λ= 1/50

The graph shows $\underline{\Phi}$ - a vector of length 1500, with 50 non-zeros





Part VI **Concluding Remarks** Global **Local**



Local Model for a Global Signal

- □ This work suggests an interesting extension of sparse approximation theory to these new breed of models.
- □ Key questions to address:
 - □ Who are those signals obeying this model ?
 - □ Who are the appropriate local dictionaries ?
 - □ How should these signals be processed via local operations ?
 - □ Can we derive performance bounds for such algorithms ?
 - □ How should the dictionary be learned ?
- □ As we have seen today, some of these questions were answered. Much work remains in order to fully map this field ...
- □ Stay tuned we are working on this.



Thank You for Your Time & Patience ... and ... Thanks to The Organizers:

Holger Boche - Technical University Munich Giuseppe Caire - Technical University Berlin Robert Calderbank - Duke University Gitta Kutyniok - Technical University Berlin Rudolf Mathar - RWTH Aachen University

Questions?



Image Denoising is a Popular Problem

Probably the most studied problem in image processing ...





- □ There are ~22,000 journal papers on image denoising.
- Searching "image and Gaussian and noise and (denois* or remov* or filter* or clean) " in ISI WoS leads to ~1800 papers

Why is it so popular? Here are few possible explanations:
(i) It does come up in many applications
(ii) It is the simplest inverse problem, platform for new ideas
(iii) Many other problems can be recast as an iterated denoising, and ...
(iv) It is misleadingly simple



Restricted Isometry Property ...

$$\min_{\underline{\Gamma}} \left\| \mathbf{D}_{\mathsf{G}} \underline{\Gamma} - \underline{\mathsf{Y}} \right\|_{2}^{2} \quad \text{s.t.} \quad \mathbf{M} \underline{\Gamma} = \underline{\mathsf{0}}, \ \left\| \underline{\Gamma} \right\|_{\mathbf{0},\infty} = \mathsf{L}$$

Definition: The globalized model, as characterized by D_G , M and k, is said to have the generalized RIP with δ_k if:

$$\left(1-\delta_{k}\right)\left\|\underline{V}\right\|_{2}^{2} \leq \left\|\mathbf{D}_{G}\underline{V}\right\|_{2}^{2} \leq \left(1+\delta_{k}\right)\left\|\underline{V}\right\|_{2}^{2}$$

for any vector obeying $\mathbf{M}\underline{V} = \underline{0}, \|\underline{V}\|_{0\infty} = \mathbf{k}$

Armed with this definition, we can derive a stability result for the above pursuit problem, show its near-oracle performance, derive a similar result for L₁ replacing the L_{0,∞}, and then via ADMM, show that this can be achieved by local operations

