Image Denoising and Beyond via Learned Dictionaries and Sparse Representations*

Michael Elad

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

* Joint work with







M. Aharon A.M. Bruckstein O. Bryt D.L. Donoho



D.L. Donono



Tel-Aviv University Applied Math Department Approximation Seminar June 26, 2008



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 then generalizing to more complex problems.
- Many Considered Directions: Partial differential equations, Statistical estimators, Adaptive filters, Inverse problems & regularization, Wavelets, Example-based techniques, Sparse representations, ...



Agenda





Part I Denoising by Sparse & Redundant Representations



Denoising By Energy Minimization

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

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

 \underline{y} : Given measurements

 $\underline{\mathbf{x}}$: Unknown to be recovered

Relation to measurements

Prior or regularization

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



Thomas Bayes 1702 - 1761



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

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





Sparse Modeling of Images



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

The vector <u>α</u> is generated randomly with few (say L) non-zeros at random locations and with random values.



Back to Our MAP Energy Function

- □ We L_0 norm is effectively counting the number of non-zeros in α .
- □ The vector $\underline{\alpha}$ is the representation (sparse/redundant).

The core idea: while few (L out of K) atoms can be merged to form the true signal, the noise cannot be fitted well. Thus, we obtain an effective projection of the noise onto a very low-dimensional space, thus getting denoising effect.



There are Some Issues ...

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

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

□ Theoretical Problem: If we are to approximate the solution somehow, how close will we get?

Practical Problem: What dictionary D should we use, such that all this leads to effective denoising?

These are the topics of the next 2 parts.



PartTheoretical &Numerical Foundations





Dictionaries and Sparse representations

By: Michael Elad

Lets Approximate

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



Smooth the L₀ and use continuous optimization techniques



Greedy methods

Build the solution one non-zero element at a time





- □ 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)].
 - Sequential shrinkage for union of ortho-bases [Bruce et.al. ('98)].
 - Iterated shrinkage [Figuerido & Nowak ('03)] [Daubechies, Defrise, & Demole ('04)] [Elad ('05)] [Elad, Matalon, & Zibulevsky ('06)].



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 rsidual.



- □ The algorithm stops when the error $\|D\underline{\alpha} \underline{y}\|_2$ is below the destination threshold.
- □ The Orthogonal MP (OMP) is an improved version that re-evaluates the coefficients by Least-Squares after each round.



Equivalence For Min $\|\underline{\alpha}\|_{0}$ s.t. $\underline{D}\underline{\alpha} = y$

Equivalence

Given a signal y with a representation $y = D\alpha$, assuming that $\|\underline{\alpha}\|_{0} < \text{Thr}\{D\}$, BP and OMP are [Donoho & Elad ('02)] [Gribonval & Nielsen ('03)] **Guaranteed** to find the sparsest solution. [Tropp (`03)] [Temlyakov ('03)]

□ MP and BP are different in general (hard to say which is better).

□ 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)] [Elad ('06)] [Tropp et.al. ('06)].



BP Stability for $Min \|\underline{\alpha}\|_0$ s.t. $\|\underline{D}\underline{\alpha} - \underline{y}\|_2 \le \varepsilon$

Stability 1

Given a signal $\underline{y} = \mathbf{D}\underline{\alpha} + \underline{v}$ with a representation satisfying $\|\underline{\alpha}\|_0 < 0.5 \text{Thr}\{\mathbf{D}\}$ and bounded noise $\|\underline{v}\|_2 \le \varepsilon$, BP will give stability, i.e., $\|\underline{\hat{\alpha}}_{BP} - \underline{\alpha}\|_2^2 < \text{Const}_1\{\mathbf{D}\} \cdot \varepsilon^2$

[Donoho, Elad & Temlyakov ('06)] [Tropp ('06)] [Donoho & Elad ('07)]

 \Box For ε =0 we get a weaker version of the previous result.

□ Surprising - the error is independent of the SNR.

□ This result is useless for assessing denoising performance.

□ Worst case versus average performance [Candes et. al. ('07)] [Donoho ('07)].



OMP Stability for $Min \|\underline{\alpha}\|_0$ s.t. $\|\underline{D}\underline{\alpha} - \underline{y}\|_2 \le \varepsilon$

Stability 2 Given a signal $\underline{Y} = \mathbf{D}\underline{\alpha} + \underline{Y}$ with bounded noise $\|\underline{Y}\|_2 \le \varepsilon_r$ and a sparse representation, $\|\underline{\alpha}\|_0 < \mathrm{Thr}\{\mathbf{D}\} - \frac{\varepsilon \cdot \mathrm{Const}_2\{\mathbf{D}\}}{\min_k\{|\alpha(k)|\}}$ OMP will give stability, i.e., $\|\underline{\hat{\alpha}}_{\mathsf{MP}} - \underline{\alpha}\|_2^2 < \mathrm{Const}_3\{\mathbf{D}\} \cdot \varepsilon^2$

[Donoho, Elad & Temlyakov ('06)] [Tropp ('06)]

 \Box For ε =0 we get the results shown already.

□ Here the error is dependent of the SNR, and

□ There are additional results on the sparsity pattern recovery.



Part II Dictionary Learning: The K-SVD Algorithm



What Should D Be?

$$\underline{\hat{\alpha}} = \underset{\alpha}{\operatorname{arg\,min}} \|\underline{\alpha}\|_{0}^{0} \quad \text{s.t.} \ \frac{1}{2} \| \mathbf{D}\underline{\alpha} - \underline{y} \|_{2}^{2} \le \varepsilon^{2} \longrightarrow \hat{\underline{x}} = \mathbf{D}\underline{\hat{\alpha}}$$

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, ...)

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



Measure of Quality for D





The K–SVD Algorithm – General [Aharon, Elad & Bruckstein ('04,'05)] Initialize D Sparse Coding Use Matching Pursuit Dictionary Update Column-by-Column by SVD computation over the relevant examples



Image Denoising & Beyond Via Learned Dictionaries and Sparse representations By: Michael Elad

21

K–SVD: Sparse Coding Stage



Image Denoising & Beyond Via Learned Dictionaries and Sparse representations By: Michael Elad

K–SVD: Dictionary Update Stage



We should solve:



We refer only to the examples that use the column \underline{d}_k

Fixing all **A** and **D** apart from the kth column, and seek both <u>d</u>_k and the kth column in **A** to better fit the **residual**!



Part IV Back to Denoising ... and Beyond – Combining It All



From Local to Global Treatment

 The K-SVD algorithm is reasonable for lowdimension signals (N in the range 10-400).
 As N grows, the complexity and the memory requirements of the K-SVD become prohibitive.



- □ So, how should large images be handled?
- □ The solution: Force shift-invariant sparsity on each patch of size N-by-N (N=8) in the image, including overlaps [Roth & Black ('05)].

$$\hat{\underline{x}} = \operatorname{ArgMin}_{\underline{x}, \{\underline{\alpha}_{ij}\}_{ij}} \frac{1}{2} \left\| \underline{x} - \underline{y} \right\|_{2}^{2} + \mu \sum_{ij} \left\| \mathbf{R}_{ij} \underline{x} - \mathbf{D} \underline{\alpha}_{ij} \right\|_{2}^{2}$$
Extracts a patch in the ij location
$$s.t. \quad \left\| \underline{\alpha}_{ij} \right\|_{0}^{0} \leq L$$
Our prior



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 N-by-N (overlapping blocks),
- □ Image of size 1000^2 pixels → $\sim 10^6$ examples to use more than enough.
- □ This works much better!











Image Denoising (Gray) [Elad & Aharon ('06)]





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

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) direct generalization (working with R+G+B patches) leads to color artifacts. The solution was found to be a bias in the pursuit towards the color content.

Original

Noisy (12.77dB) Result (29.87dB)



Image Denoising & Beyond Via Learned Dictionaries and Sparse representations By: Michael Elad

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

Our experiments lead to state-of-the-art inpainting results.





Image Denoising & Beyond Via Learned Dictionaries and Sparse representations By: Michael Elad

Video Denoising [Protter & Elad ('06)]

When turning to handle video, one could improve over the previous scheme in two important ways:

Our experiments lead to state-of-the-art video denoising results, giving ~0.5dB better results on average compared to [Boades, Coll & Morel ('05)] and comparable to [Rusanovskyy, Dabov, & Egiazarian ('06)]

> compensation can and should be Noisy (0=15) Denoised (PSNR=29.98) avoided [Buades, Col, and Morel ('06)].

3. Motion estimation and

implicitly.

Original



Image Denoising & Beyond Via Learned Dictionaries and Sparse representations By: Michael Elad .62)

Image Compression [Bryt and Elad ('06)]

- □ 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 dictionaries locally (per patch) using a training set of images is required.
- In PCA, only the (quantized) coefficients are stored, whereas the K-SVD requires storage of the indices as well.
- Geometric alignment of the image is very helpful and should be done [Goldenberg, Kimmel, & Elad ('05)].





Image Compression Results





Image Compression Results





Part V Interesting & Recent Advancement



What If ...

Consider the denoising problem

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

and suppose that we can find a group of J candidate solutions

$$\{ \underline{\alpha}_j \}_{j=1}^{J}$$

such that

$$\forall j \quad \begin{cases} \left\|\underline{\alpha}_{j}\right\|_{0}^{0} << N \\ \left\|\underline{D}\underline{\alpha}_{j} - \underline{y}\right\|_{2}^{2} \le \varepsilon^{2} \end{cases}$$

Basic Questions:

- ❑ What could we do with such a set of competing solutions in order to better denoise <u>y</u>?
- □ Why should this work?
- □ How shall we practically find such a set of solutions?

These questions were studied and answered recently [Elad and Yavneh ('08)]



Motivation

Why bother with such a set?

- Because of the intriguing relation to example-based techniques, where several nearest-neighbors for the signal are used jointly.
- Because each representation conveys a different story about the desired signal.
- Because pursuit algorithms are often wrong in finding the sparsest representation, and then relying on their solution becomes too sensitive.
- □ ... Maybe there are "deeper" reasons?





Generating Many Representations

Our Answer: Randomizing the OMP





Lets Try

Proposed Experiment :

- □ Form a random **D**.
- \Box Multiply by a sparse vector $\underline{\alpha}_0$ ($\|\underline{\alpha}_0\|_0^0 = 10$).
- □ Add Gaussian iid noise (σ =1) and obtain <u>y</u>.
- □ Solve the problem
 - $\min_{\underline{\alpha}} \|\underline{\alpha}\|_{0}^{0} \text{ s.t. } \|\mathbf{D}\underline{\alpha} \underline{y}\|_{2}^{2} \leq 100$ using OMP, and obtain $\underline{\alpha}^{OMP}$.
- Use RandOMP and obtain $\left\{ \begin{array}{c} \alpha_{j}^{\text{RandOMP}} \end{array} \right\}_{i=1}^{1000}$
- □ Lets look at the obtained representations ...





Some Observations



We see that

- The OMP gives the sparsest solution
- Nevertheless, it is not the most effective for denoising.
- The cardinality of a representation does not reveal its efficiency.



Image Denoising & Beyond Via Learned Dictionaries and Sparse representations By: Michael Elad

The Surprise ...







Is It Consistent? Yes!

Here are the results of 1000 trials with the same parameters ...





Image Denoising & Beyond Via Learned Dictionaries and Sparse representations By: Michael Elad

The Explanation – Our Signal Model



Signal Model Assumed

D is fixed and known

 \Box <u>a</u> is built by:

- Choosing the support S w.p. P(S) of all the 2^κ possibilities Ω,
- Choosing the coefficients using iid Gaussian entries* N(0,σ_x): P(<u>α</u>|S).

□ The ideal signal is $\underline{x} = D\underline{\alpha}$.

The p.d.f. of the signal P(
$$\underline{x}$$
) is: P(\underline{x}) = $\sum_{S \in \Omega} P(\underline{x}|S)P(S)$

* Not exactly, but this does not change our analysis.



The Explanation – Adding Noise



Noise Assumed:

The noise \underline{v} is additive white Gaussian vector with probability $P_v(\underline{v})$

$$\mathsf{P}(\underline{\mathsf{y}}|\underline{\mathsf{x}}) = \mathsf{C} \cdot \mathsf{exp}\left\{-\frac{\left\|\underline{\mathsf{x}} - \underline{\mathsf{y}}\right\|^{2}}{2\sigma^{2}}\right\}$$

The p.d.f. of the noisy signal $P(\underline{y})$, and the conditionals $P(\underline{y}|S)$ and $P(S|\underline{y})$ are clear and well-defined (although nasty).



Image Denoising & Beyond Via Learned Dictionaries and Sparse representations By: Michael Elad

Some Algebra Leads To

$$\underline{\hat{x}}^{MMSE} = E \left\{ \underline{x} \middle| \underline{y} \right\}$$

$$P(S|\underline{y}) \propto exp\left\{\frac{\sigma_{X}^{2}}{\sigma_{X}^{2} + \sigma^{2}} \cdot \frac{\left\|\underline{y}_{S}\right\|^{2}}{2\sigma^{2}}\right\}$$

Projection of the signal <u>y</u> onto the support S $\underline{y}_{S} = \mathbf{D} \cdot \left\{ \underset{\underline{\alpha}}{\operatorname{ArgMin}} \| \underline{y} - \underline{D}\underline{\alpha} \| \text{ s.t. sup } \underline{x} \right\} = S$

Implications:

The best estimator (in terms of L₂ error) is a weighted average of many sparse representations!!!



As It Turns Out ...

- □ The MMSE estimation we got requires a sweep through all 2^{κ} supports (i.e. combinatorial search) impractical.
- □ Similarly, an explicit expression for $P(\underline{x}/\underline{y})$ can be derived and maximized this is the MAP estimation, and it also requires a sweep through all possible supports impractical too.

□ The OMP is a (good) approximation for the MAP estimate.

□ The RandOMP is a (good) approximation of the Minimum-Mean-Squared-Error (MMSE) estimate. It is close to the Gibbs sampler of the probability $P(S|\underline{y})P(S)$ from which we should draw the weights.

Back to the beginning: Why Use Several Representations? Because their average leads to provable better noise suppression.



Example

The following results correspond to a small dictionary (20×30), where the combinatorial formulas can be evaluated as well.

Parameters:

- N=20, K=30
- True support=3
- $\sigma_x = 1$
- J=10
- Averaged over 1000
 experiments





Image Denoising & Beyond Via Learned Dictionaries and Sparse representations By: Michael Elad

Part VI Summary and Conclusion



Today We Have Seen that ...



More on these (including the slides, the papers, and a Matlab toolbox) in http://www.cs.technion.ac.il/~elad





All this Work is Made Possible Due to



my teachers and mentors

colleagues & friends collaborating with me



G. Sapiro J.L. Starck I. Yavneh M. Zibulevsky

and my students



M. Aharon O. Bryt

J. Mairal

M. Protter R. Rubinstein J. Shtok



Image Denoising & Beyond Via Learned Dictionaries and Sparse representations By: Michael Elad

Sparseland Signals Are Special



Interesting Model:

- 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 low-dimensional Gaussians (or Laplacians).
- □ Familiar: We have been using this model in other context for a while now (wavelet, JPEG, ...).



K–Means For Clustering





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

- □Orodexpectimentaslead sensible of typenart demosaicing colesytes, gixelygleavi2gEthettest results on interpolated compared to [Chang & Chan ('06)]
- Generalizing the previous 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 have handled the demosaicing problem while forcing strong sparsity and only few iterations.
- □ The same concept can be deployed to inpainting.



Image Compression Results





Image Compression



