

SRINKAGE FOR REDUNDANT REPRESENTATIONS ?

Michael Elad

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



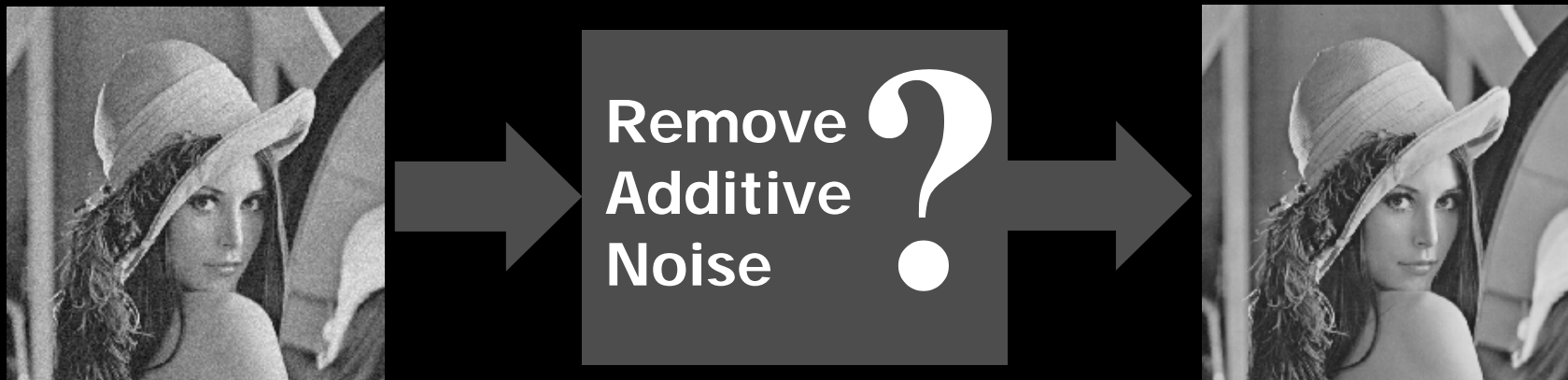
SPARS'05

Signal Processing with Adaptive Sparse Structured Representations
November 16-18, 2005 – Rennes, France



Noise Removal

Our story begins with signal/image denoising ...

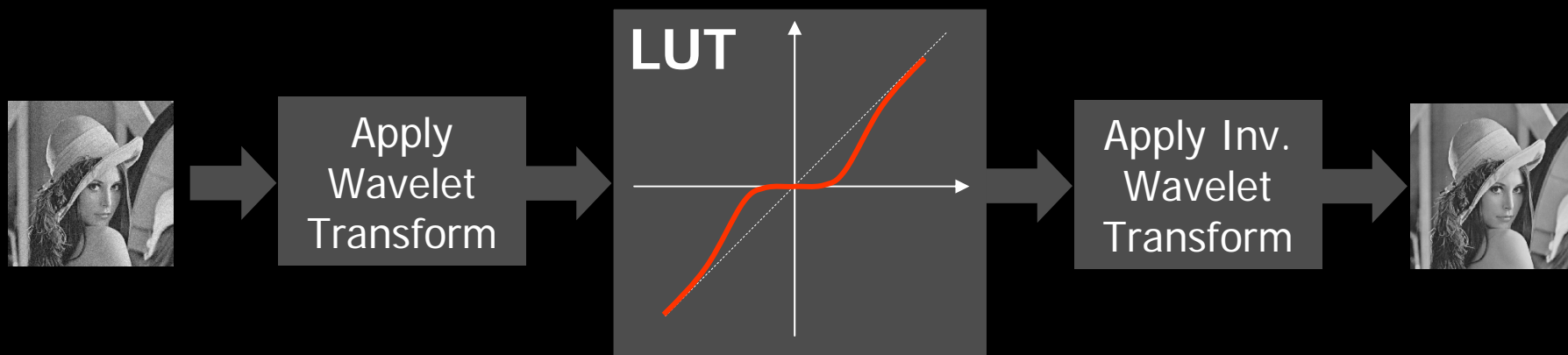


- ❑ 100 years of activity – numerous algorithms.
- ❑ Considered Directions include: PDE, statistical estimators, adaptive filters, inverse problems & regularization, sparse representations, ...



Shrinkage For Denoising

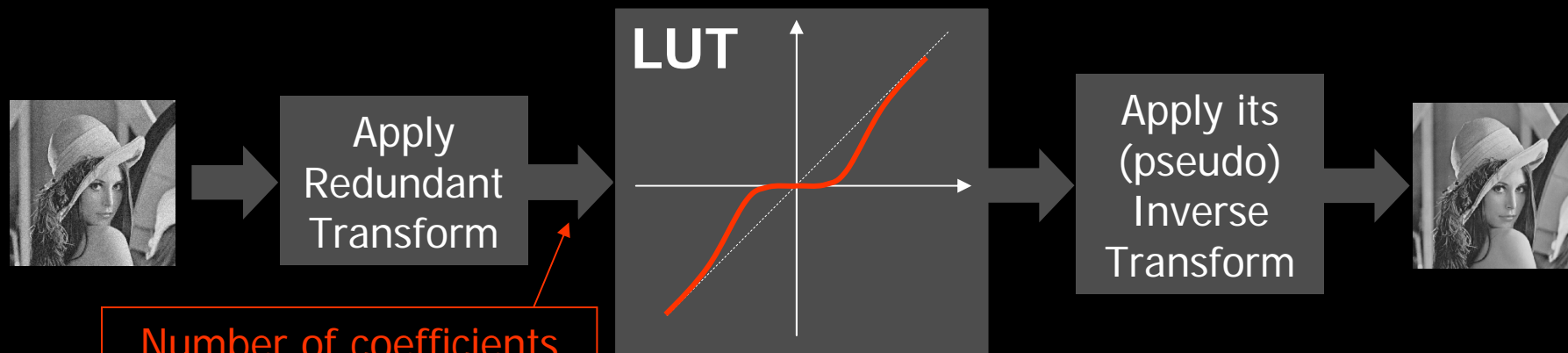
- ❑ Shrinkage is a simple yet effective denoising algorithm [Donoho & Johnstone, 1993].
- ❑ Justification 1: minimax near-optimal over the Besov (smoothness) signal space (complicated!!!!).



- ❑ Justification 2: Bayesian (MAP) optimal [Simoncelli & Adelson 1996, Moulin & Liu 1999].
- ❑ In both justifications, an additive **Gaussian white** noise and a **unitary** transform are crucial assumptions for the optimality claims.



Redundant Transforms?



- ❑ This scheme is still applicable, and it works fine (tested with curvelet, contourlet, undecimated wavelet, and more).
- ❑ However, it is no longer the optimal solution for the MAP criterion.

TODAY'S FOCUS:

IS SHRINKAGE STILL RELEVANT WHEN HANDLING REDUNDANT (OR NON-UNITARY) TRANSFORMS?
HOW? WHY?



Agenda

1. **Bayesian Point of View – a Unitary Transform**
Optimality of shrinkage
2. What About Redundant Representation?
Is shrinkage is relevant? Why? How?
3. Conclusions



Thomas Bayes
1702 - 1761



The MAP Approach

Minimize the following function with respect to \underline{x} :

$$f(\underline{x}) = \frac{1}{2} \|\underline{x} - \underline{y}\|_2^2 + \lambda \cdot \text{Pr}(\underline{x})$$

Log-Likelihood term

Prior or regularization

Unknown to be recovered

Given measurements



Image Prior?

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

$$\Pr(\underline{x}) = \lambda \|\underline{x}\|_2^2$$



Energy

$$\Pr(\underline{x}) = \lambda \|\mathbf{L}\underline{x}\|_2^2$$



Smoothness

$$\Pr(\underline{x}) = \lambda \|\mathbf{L}\underline{x}\|_{\mathbf{W}}^2$$



Adapt +
Smooth

$$\Pr(\underline{x}) = \lambda \rho\{\mathbf{L}\underline{x}\}$$



Robust
Statistics

$$\Pr(\underline{x}) = \lambda \|\|\nabla \underline{x}\|\|_1$$



Total-
Variation

$$\Pr(\underline{x}) = \lambda \|\mathbf{W}\underline{x}\|_1 \quad \Pr(\underline{x}) = \lambda \|\mathbf{T}\underline{x}\|_1$$

Today's Focus

- Mumford & Shah formulation,
- Compression algorithms as priors,
- ...



(Unitary) Wavelet Sparsity

L_2 is unitarily invariant

$$f(\underline{x}) = \frac{1}{2} \|\underline{x} - \underline{y}\|_2^2 + \lambda \cdot \|\mathbf{W}\underline{x}\|_1$$

$$f(\hat{\underline{x}}) = \frac{1}{2} \|\mathbf{W}^H(\hat{\underline{x}} - \hat{\underline{y}})\|_2^2 + \lambda \cdot \|\hat{\underline{x}}\|_1$$

$$f(\underline{x}) = \frac{1}{2} \|\hat{\underline{x}} - \hat{\underline{y}}\|_2^2 + \lambda \cdot \|\hat{\underline{x}}\|_1$$

$$= \sum_k \left[\frac{1}{2} (\hat{x}_k - \hat{y}_k)^2 + \lambda |\hat{x}_k| \right]$$

Define
 $\hat{\underline{x}} = \mathbf{W}\underline{x}$
 $\underline{x} = \mathbf{W}^H\hat{\underline{x}}$

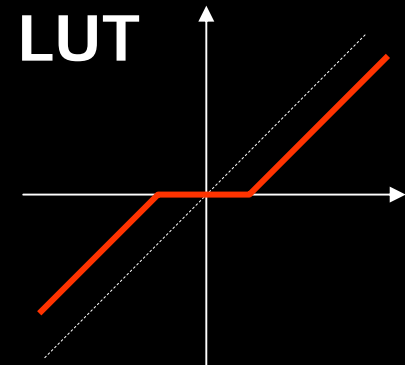
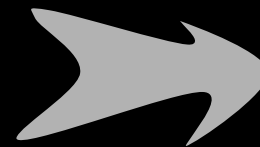
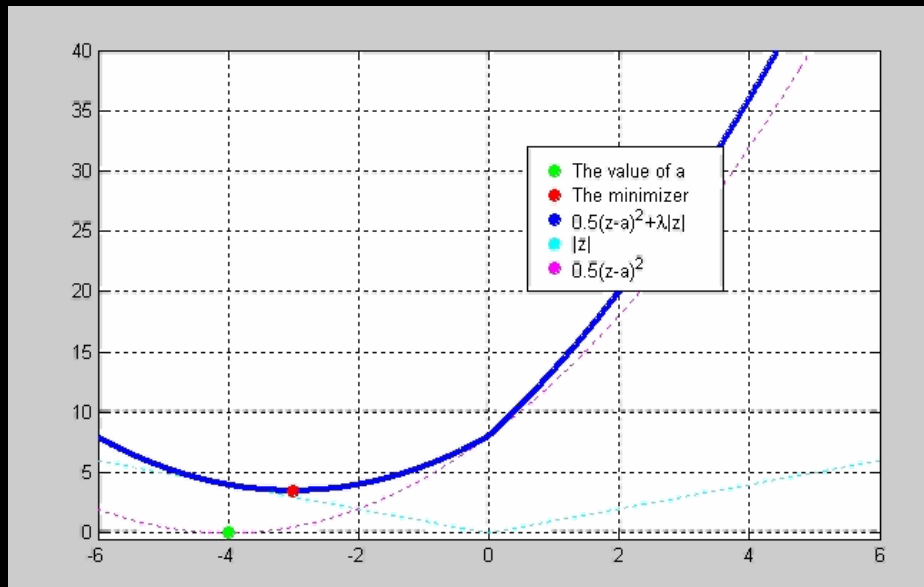
We got a separable set of 1D optimization problems



Why Shrinkage?

Want to minimize this 1-D function with respect to z

$$f(z) = \frac{1}{2}(z - a)^2 + \lambda|z|$$



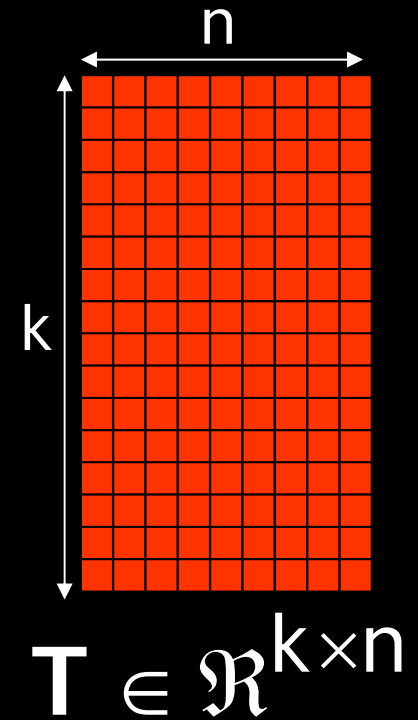
$$z_{\text{opt}} = \begin{cases} a - \lambda & a \geq \lambda \\ 0 & |a| < \lambda \\ a + \lambda & a \leq -\lambda \end{cases}$$

A LUT can be built for any other robust function (replacing the $|z|$), including non-convex ones (e.g., L_0 norm)!!



Agenda

1. Bayesian Point of View – a Unitary Transform
Optimality of shrinkage
2. **What About Redundant Representation?**
Is shrinkage is relevant? Why? How?
3. Conclusions



An Overcomplete Transform

$$\mathbf{T} \underline{x} = \begin{matrix} \text{10x10 grid} \\ \text{1x10 column} \end{matrix} = \begin{matrix} \text{1x10 column} \end{matrix} = \underline{\alpha}$$

$$f(\underline{x}) = \frac{1}{2} \|\underline{x} - \underline{y}\|_2^2 + \lambda \cdot \|\mathbf{T} \underline{x}\|_1$$

Redundant transforms are important because they can

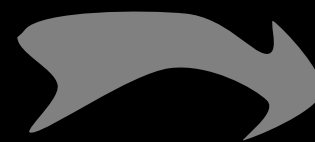
- (i) Lead to a shift-invariance property,
- (ii) Represent images better (because of orientation/scale analysis),
- (iii) Enable deeper sparsity (and thus give more structured prior).



Analysis versus Synthesis

Analysis
Prior:

$$f(\underline{x}) = \frac{1}{2} \|\underline{x} - \underline{y}\|_2^2 + \lambda \cdot \|\mathbf{T}\underline{x}\|_1$$



Define

$$\underline{\alpha} = \mathbf{T}\underline{x}$$



$$\underline{x} = \mathbf{T}^+ \underline{\alpha}$$

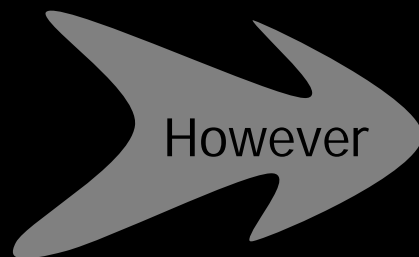
$$\tilde{f}(\underline{\alpha}) = \frac{1}{2} \|\mathbf{T}^+ \underline{\alpha} - \underline{y}\|_2^2 + \lambda \cdot \|\underline{\alpha}\|_1$$



Synthesis
Prior:

$$\tilde{f}(\underline{\alpha}) = \frac{1}{2} \|\mathbf{D}\underline{\alpha} - \underline{y}\|_2^2 + \lambda \cdot \|\underline{\alpha}\|_1$$

Basis Pursuit



However

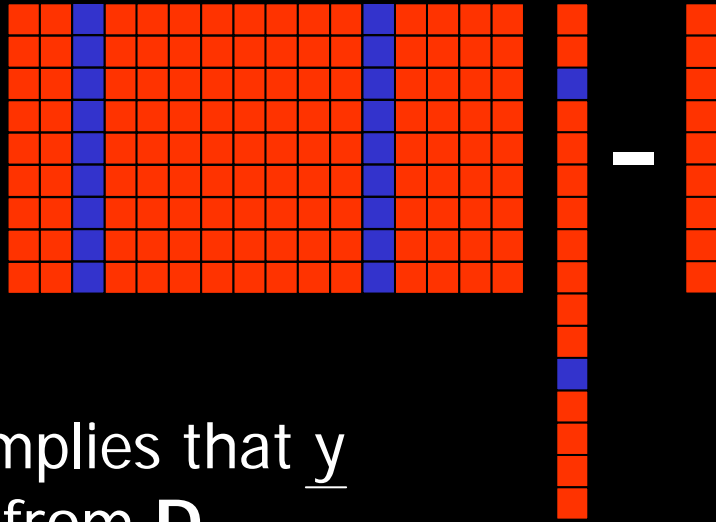
$$\mathbf{D} \cdot \underset{\underline{\alpha} = \mathbf{T}\mathbf{T}^+ \underline{\alpha}}{\text{Arg min}} \tilde{f}(\underline{\alpha}) = \underset{\underline{x}}{\text{Arg min}} f(\underline{x})$$



Basis Pursuit As Objective

Our Objective: $\tilde{f}(\underline{\alpha}) = \frac{1}{2} \|\underline{D}\underline{\alpha} - \underline{y}\|_2^2 + \lambda \cdot \|\underline{\alpha}\|_1$



$\underline{D}\underline{\alpha} - \underline{y} =$ 

Getting a sparse solution implies that \underline{y} is composed of few atoms from \underline{D}



Sequential Coordinate Descent

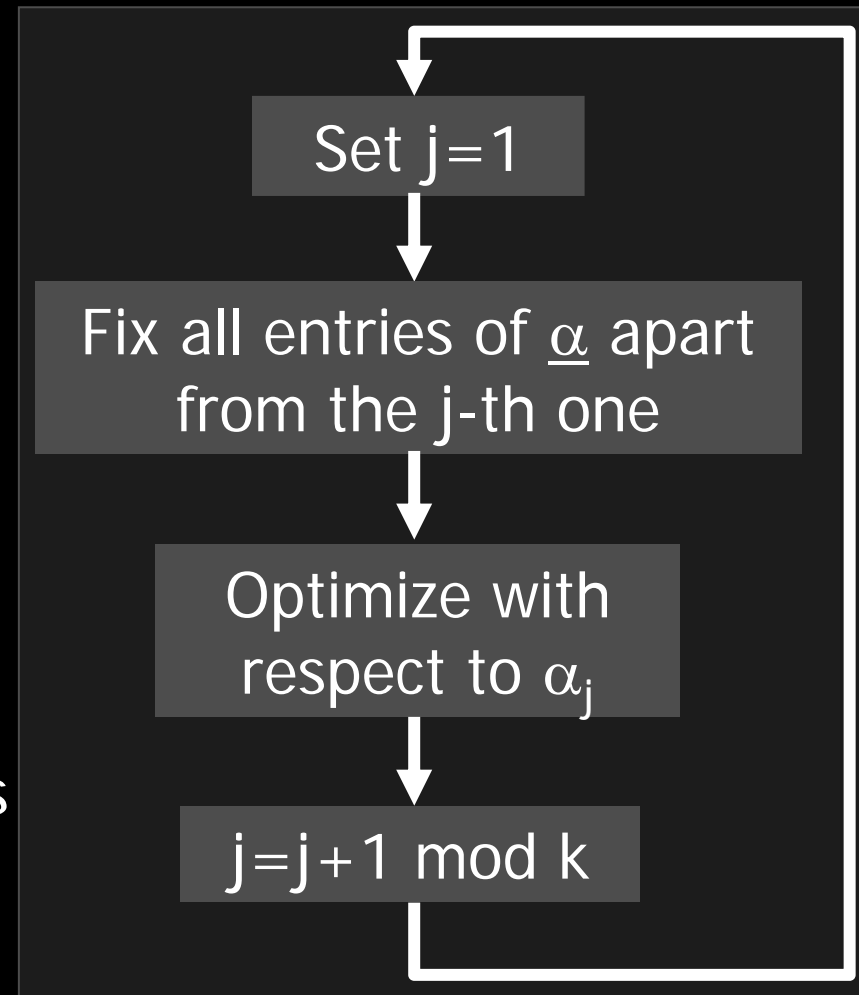
Our objective

$$\tilde{f}(\underline{\alpha}) = \frac{1}{2} \left\| \mathbf{D}\underline{\alpha} - \underline{y} \right\|_2^2 + \lambda \cdot \|\underline{\alpha}\|_1$$



- The unknown, $\underline{\alpha}$, has k entries.
- How about optimizing w.r.t. each of them sequentially?
- The objective per each becomes

$$\tilde{f}(z) = \frac{1}{2} \left\| z\underline{d}_j - \tilde{\underline{y}} \right\|_2^2 + \lambda|z|$$



We Get Sequential Shrinkage

BEFORE: We had this 1-D function to minimize

$$f(z) = \frac{1}{2}(z - a)^2 + \lambda|z|$$

and the solution was $z_{\text{opt}} = \mathbf{S}\{a, \lambda\} = \begin{cases} a - \lambda & a \geq \lambda \\ 0 & |a| < \lambda \\ a + \lambda & a \leq -\lambda \end{cases}$

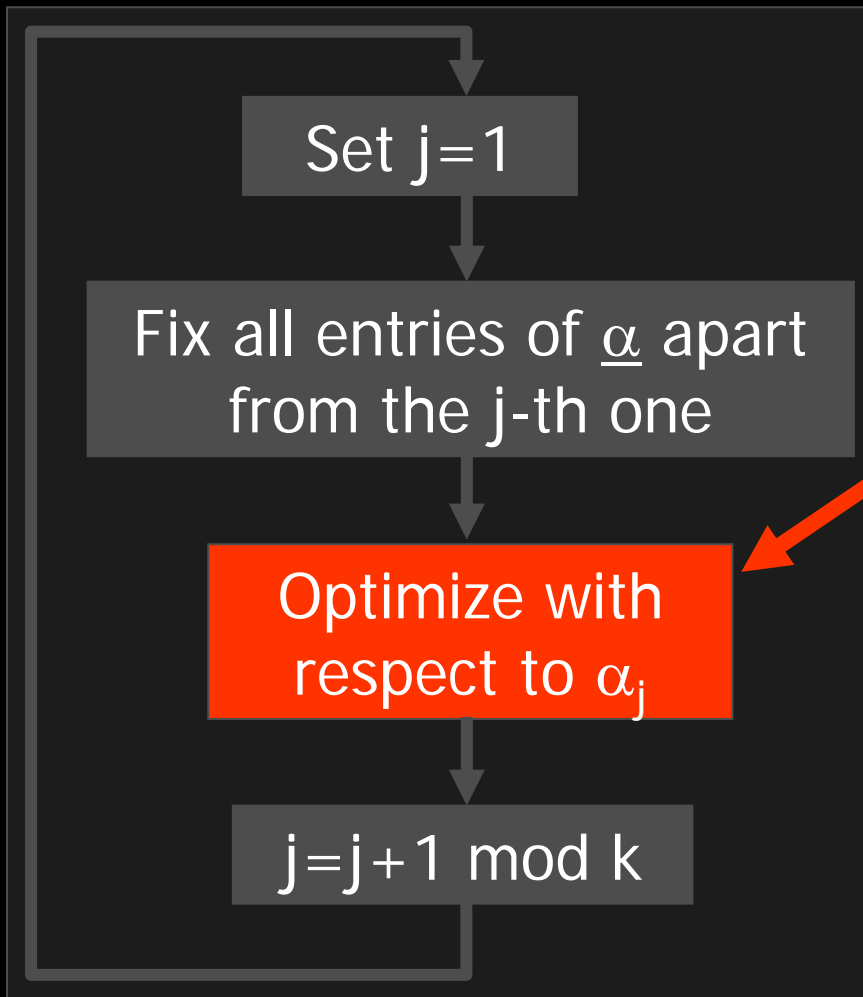
NOW: Our 1-D objective is

$$\tilde{f}(z) = \frac{1}{2} \|z \underline{d}_j - \tilde{\underline{y}}\|_2^2 + \lambda|z|$$

and the solution now is $z_{\text{opt}} = \mathbf{S}\left\{ \frac{\underline{d}_j^H \tilde{\underline{y}}}{\|\underline{d}_j\|_2^2}, \frac{\lambda}{\|\underline{d}_j\|_2^2} \right\}$



Sequential? Not Good!!



$$\tilde{\underline{y}} = \underline{y} - \mathbf{D}\underline{\alpha} + \underline{d}_j\alpha_j \quad \text{and}$$


$$\alpha_j^{\text{opt}} = \mathbf{S} \left\{ \begin{array}{c} \underline{d}_j^H \tilde{\underline{y}} \\ \|\underline{d}_j\|_2^2 \end{array} , \begin{array}{c} \lambda \\ \|\underline{d}_j\|_2^2 \end{array} \right\}$$

- ❑ This method requires drawing one column at a time from \mathbf{D} .
- ❑ In most transforms this is not comfortable at all !!!

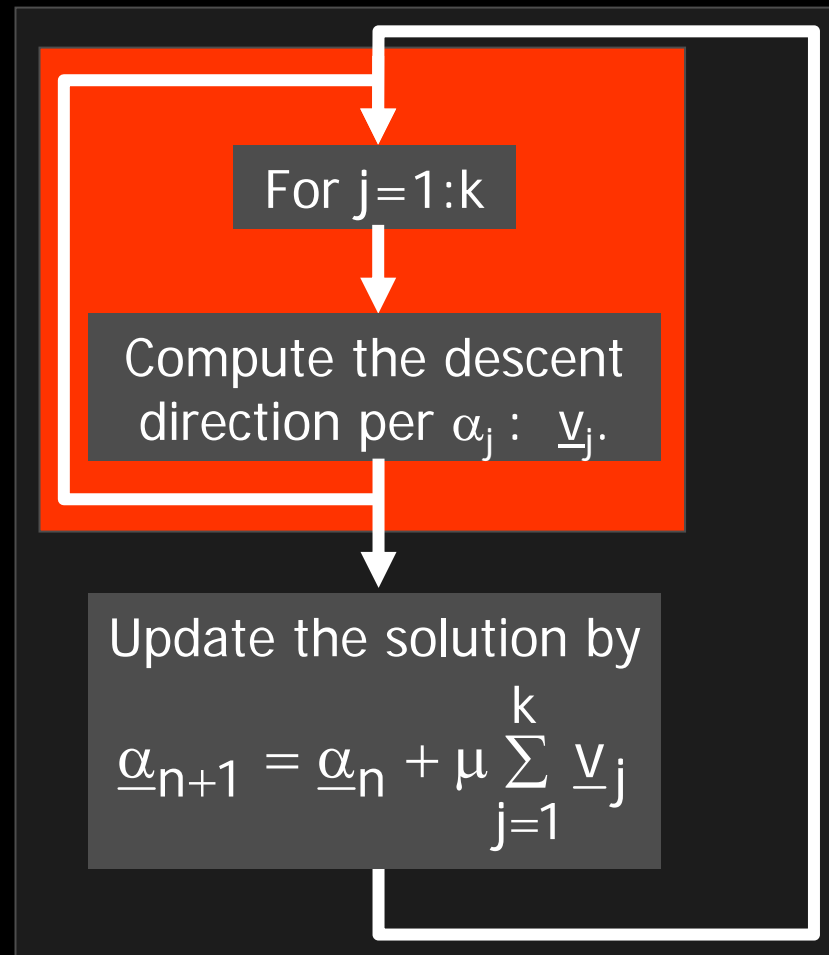


How About Parallel Shrinkage?

Our objective

$$\tilde{f}(\underline{\alpha}) = \frac{1}{2} \left\| \mathbf{D}\underline{\alpha} - \underline{y} \right\|_2^2 + \lambda \cdot \left\| \underline{\alpha} \right\|_1$$


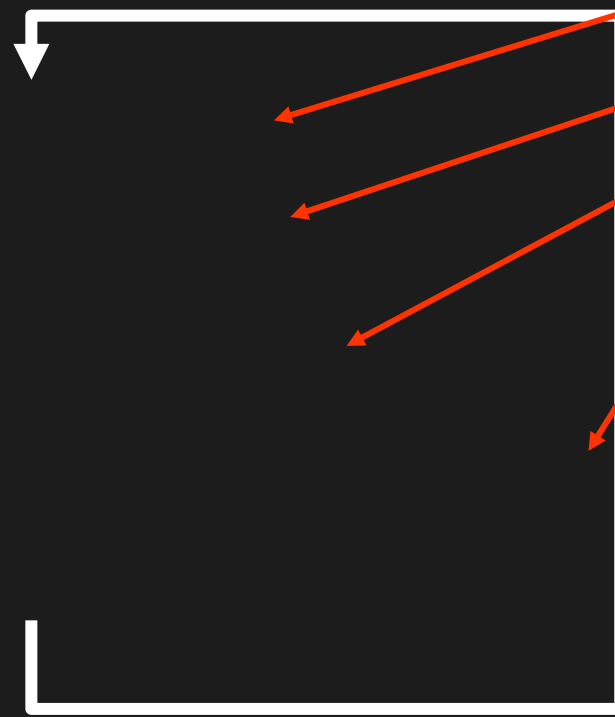
- Assume a current solution $\underline{\alpha}_n$.
- Using the previous method, we have k descent directions obtained by a simple shrinkage.
- How about taking all of them at once, with a proper relaxation?
- Little bit of math lead to ...



The Proposed Algorithm

Initialize $\underline{\alpha}_0 = \underline{0}$ & $k = 0$.

Compute



The synthesis error

Back-projection to the signal domain

Shrinkage operation

Update by line-search

At all stages, the dictionary is applied as a whole, either directly, or via its adjoint

$$(*) \mathbf{W} = \text{diag}^{-1} \{ \mathbf{D}^H \mathbf{D} \}$$



The First Iteration – A Closer Look

Initialize $\underline{\alpha}_0 = \underline{0}$ & $k = 0$.

Compute

1. $\underline{e} = \underline{y} - \mathbf{D}\underline{\alpha}_k$

2. $\underline{e}_T = \mathbf{W}\mathbf{D}^H \underline{e}$ (*)

3. $\underline{e}_T^S = \mathbf{S}\{\underline{e}_T + \underline{\alpha}_k, \lambda \mathbf{W} \cdot \underline{1}\}$

4. $\underline{\alpha}_{k+1} = \underline{\alpha}_k + \mu(\underline{e}_T^S - \underline{\alpha}_k)$

5. $k = k + 1$

(*) $\mathbf{W} = \text{diag}^{-1}\{\mathbf{D}^H \mathbf{D}\}$

$$\underline{\alpha}_0 = \underline{0}$$

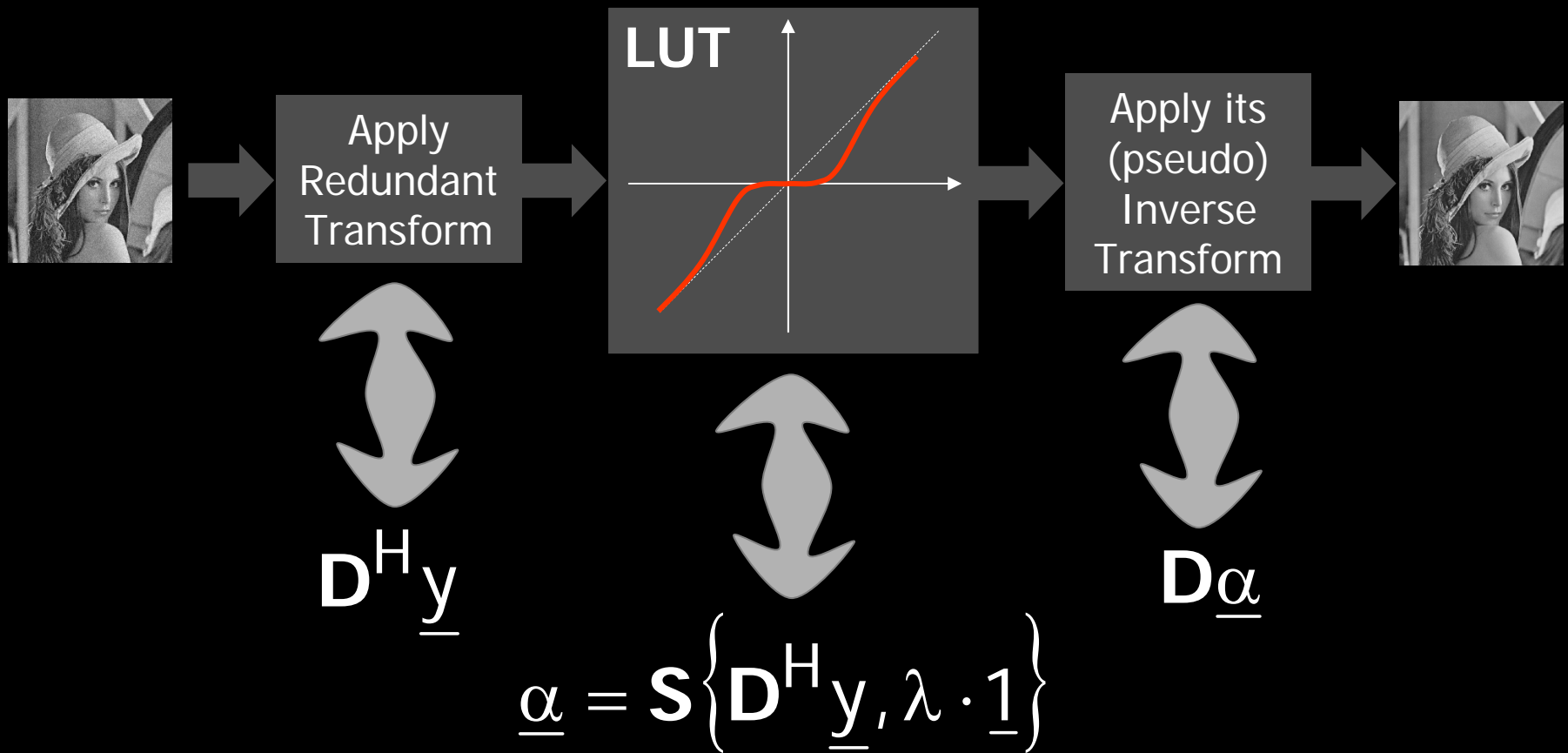
$$\underline{\alpha}_1 = \mu \mathbf{S}\{\mathbf{W}\mathbf{D}^H \underline{y}, \lambda \mathbf{W} \cdot \underline{1}\}$$

For Example:
Tight ($\mathbf{D}\mathbf{D}^H = c \cdot \mathbf{I}$)
and normalized
($\mathbf{W} = \mathbf{I}$) frame

$$\hat{\underline{x}}_1 = \mathbf{D} \cdot \mathbf{S}\{\mathbf{D}^H \underline{y}, \lambda \cdot \underline{1}\}$$



Relation to Simple Shrinkage



A Simple Example

Minimize

$$\tilde{f}(\underline{\alpha}) = \frac{1}{2} \left\| \mathbf{D}\underline{\alpha} - \underline{y} \right\|_2^2 + \lambda \cdot \|\underline{\alpha}\|_1$$

\mathbf{D} : a 100×1000 , union of 10
random unitary matrices,
 \underline{y} : $\mathbf{D}\underline{\alpha}$, with $\underline{\alpha}$ having 15 non-
zeros in random locations,

$\lambda=1$, $\underline{\alpha}_0 = \underline{0}$,

Line-Search: Armijo

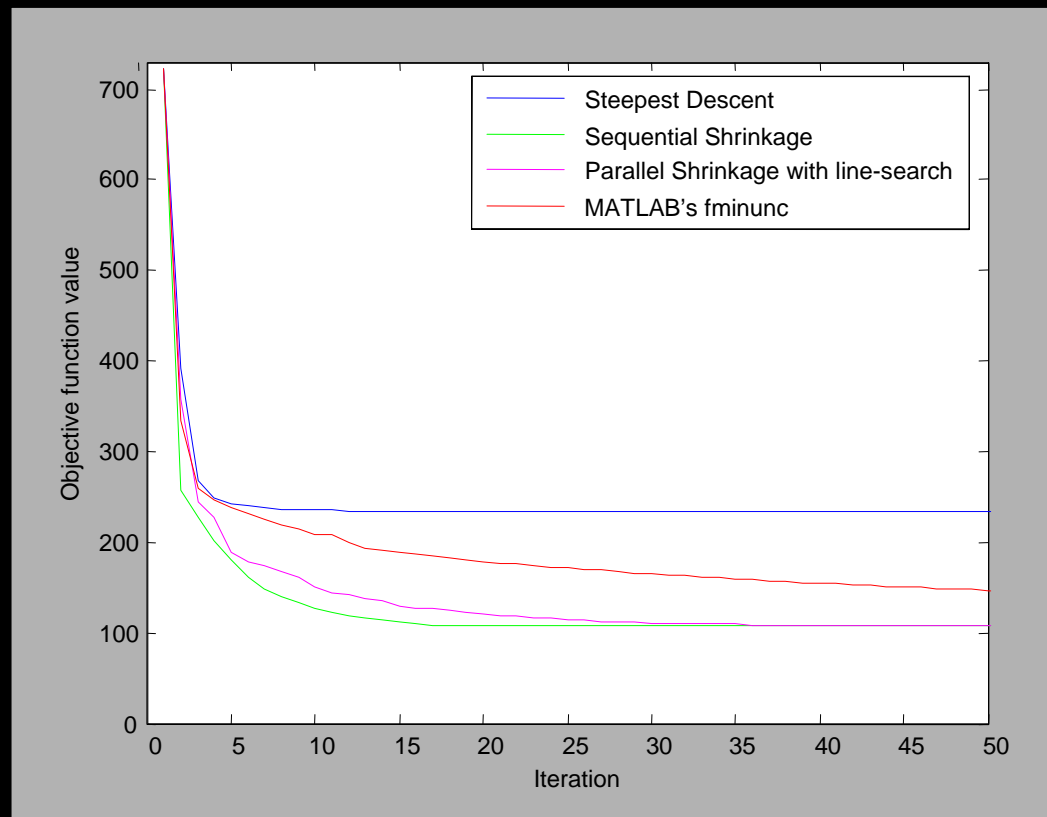


Image Denoising

Minimize

$$\tilde{f}(\underline{\alpha}) = \frac{1}{2} \|\mathbf{D}\underline{\alpha} - \underline{y}\|_2^2 + \lambda \cdot \|\mathbf{W}\underline{\alpha}\|_1$$

- The Matrix \mathbf{W} gives a variance per each coefficient, learned from the corrupted image.
- \mathbf{D} is the contourlet transform (recent version).
- The length of $\underline{\alpha}$: $\sim 1e+6$.
- The Seq. Shrinkage algorithm can no longer be simulated

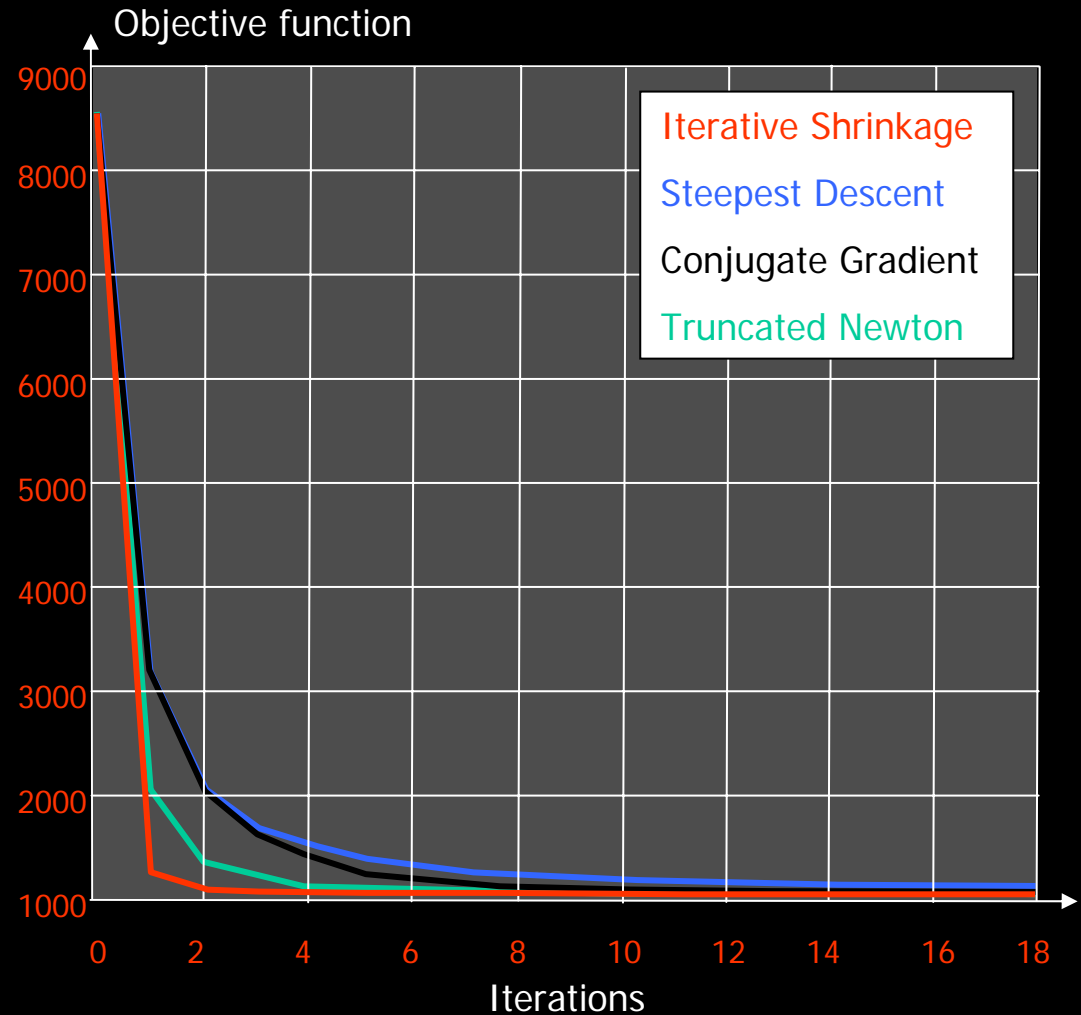


Image Denoising

Minimize

$$\tilde{f}(\underline{\alpha}) = \frac{1}{2} \|\mathbf{D}\underline{\alpha} - \underline{y}\|_2^2 + \lambda \cdot \|\mathbf{W}\underline{\alpha}\|_1$$



Evaluate $\|\mathbf{D}\hat{\underline{\alpha}} - \underline{x}_{\text{True}}\|_2^2$



Even though one iteration of our algorithm is equivalent in complexity to that of the SD, the performance is much better

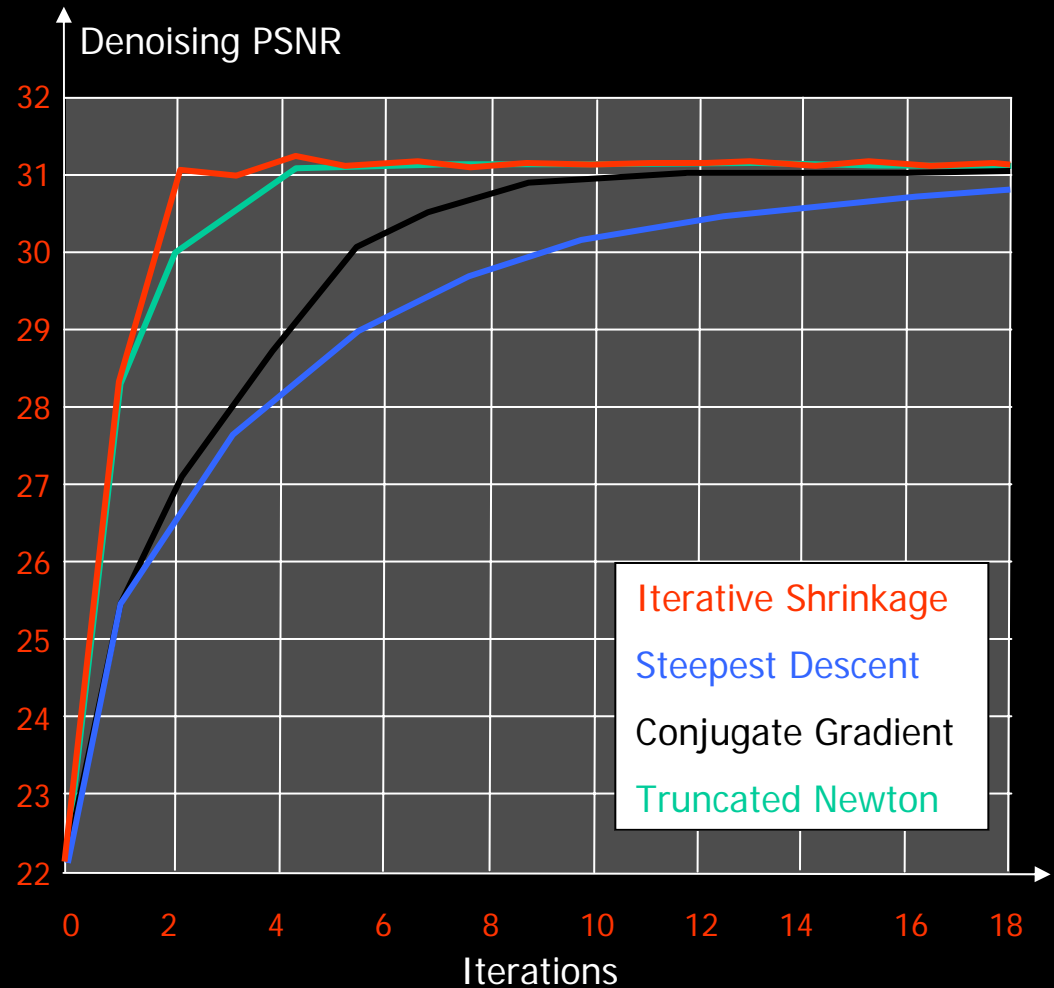
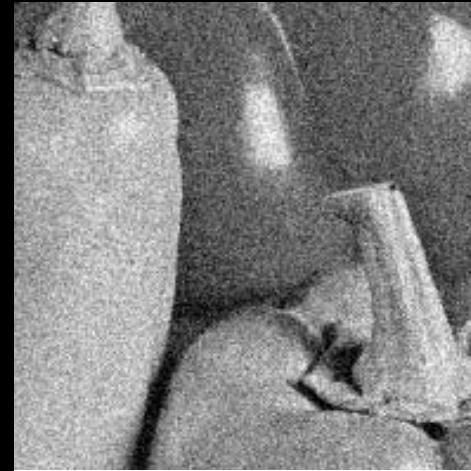


Image Denoising

Original
Image



Noisy
Image with
 $\sigma=20$



Iterated
Shrinkage –
First Iteration
PSNR=28.30dB



Iterated
Shrinkage –
second iteration
PSNR=31.05dB



Closely Related Work

- The “same” algorithm was derived in several other works:
 - Sparse representation over curvelet [Starck, Candes, Donoho, 2003].
 - E-M algorithm for image restoration [Figueiredo & Nowak 2003].
 - Iterated Shrinkage for problems of the form Minimize $\| \mathbf{K} \underline{x} - \underline{y} \|_2^2 + \lambda \cdot \| \mathbf{W} \underline{x} \|_1$ [Daubechies, Defrise, & De-Mol, 2004].
- The work proposed here is different in several ways:
 - **Motivation:** Shrinkage for redundant representation, rather than general inverse problems.
 - **Derivation:** We used a parallel CD-algorithm, and others used the EM or a sequence of surrogate functions.
 - **Algorithm:** We obtain a slightly different algorithm, where the norms of the atoms are used differently, different thresholds are used, the choice of μ is different.



Agenda

1. Bayesian Point of View – a Unitary Transform
Optimality of shrinkage
2. What About Redundant Representation?
Is shrinkage is relevant? Why? How?
3. **Conclusions**



Conclusion

Shrinkage is an appealing signal denoising technique

When optimal?

For additive Gaussian noise and unitary transforms

What if the transform is redundant?

Compute all the CD directions, and use the average

How?

Go Parallel

How to avoid the need to extract atoms?

Option 1: apply sequential coordinate descent which leads to a sequential shrinkage algorithm

Getting what?

We obtain a very easy to implement parallel shrinkage algorithm that requires forward transform, scalar shrinkage, and inverse transform.

