Fast L0-based Sparse Signal Recovery

Nick Kingsbury and Yingsong Zhang



Signal Processing Group Department of Engineering

ACVT seminar, Adelaide - May 2011

(日) (四) (三) (三) (三)

Outline

Introduction

- 2 L0 reweighted-L2 Minimization
 - IRLS
 - Basic Model
 - Basic Algorithm
 - \bullet Fast Algorithm ${\rm L}_0{\rm RL}_2$
- 3 Continuation Strategy
 - Geometry of the penalty function
 - $\bullet \ {\rm L}_0 {\rm R} {\rm L}_2$ with continuation
 - 4 Numerical Results
- 5 Super-resolution
- 6 Experimental results
- 7 Conclusion

Introduction

The L_0 minimisation problem

Find **x** which gives min $\|\mathbf{x}\|_0$ subject to $\|\mathbf{y} - \Phi \mathbf{x}\|_2 \le \xi$ (1) where Φ is known and of size $M \times N$ with M < N

- $\mathsf{L}_0:$ NP-hard combinatorial search for exact solution
- L₁: convex problem, but poor computation speed for large problems; greedy algorithms can be used, e.g. CoSaMP, NESTA.
- L_p : iterative reweighted techniques, e.g. IRL1, IRLS
- $L_0:$ greedy algorithms usually used, e.g. IHT, $\ell_0\text{-}AP,\,SL0$

We propose a highly efficient form of IRLS which approximates L₀ minimisation, allows $N \sim 10^7$ or 10^8 , and has good reconstruction performance.

イロト 不得下 イヨト イヨト

Iterative reweighted least-squares minimization (IRLS)

min
$$\mathbf{x}^T \mathbf{S} \mathbf{x}$$
 subject to $\|\mathbf{y} - \Phi \mathbf{x}\|_2 = 0$

 ${\boldsymbol{\mathsf{S}}}$ is a diagonal weight matrix. The solution of each step is the solution of

min
$$\|\mathbf{v}\|^2$$
 subject to $\|\mathbf{y} - \Phi \mathbf{x}\|_2 = 0$ and $\mathbf{x} = \mathbf{S}^{-\frac{1}{2}} \mathbf{v}$

which is unique and the reweighted iterative rule is:

$$\mathbf{x} = \mathbf{S}^{-\frac{1}{2}} (\Phi \mathbf{S}^{-\frac{1}{2}})^{\dagger} \mathbf{y}, \tag{2}$$

where [†] denotes the pseudo inverse (i.e. $\mathbf{H}^{\dagger} = (\mathbf{H}^{T}\mathbf{H})^{-1}\mathbf{H}^{T}$). The inverse of $(\mathbf{H}^{T}\mathbf{H})$ is difficult for large problems.

Iterative reweighted least squares minimizations (IRLS)

- Gorodnitsky & Rao (1997) $S_j = \frac{1}{(|x_j|^2)^{(1-\frac{p}{2})}}$, for $0 \le p \le 1$. Chartrand & Yin (2008) $S_j = \frac{1}{(|x_j|^2 + \epsilon^2)^{(1-\frac{p}{2})}}$, for $0 \le p \le 1$
 - Theoretical analysis shows that local convergence rate of the algorithm is superlinear; smaller p values result in faster convergence rate.
 - Experimentally shows that L_p norm with $p \rightarrow 0$ can help to achieve a higher success rate in exact signal recovery.

However, IRLS in this form strictly only models noise-free observations.

Basic Model – Gaussian measurement noise

Consider the noisy system

$$\mathbf{y} = \Phi \mathbf{x} + \mathbf{n}$$

where we assume :

 Φ is a known $M \times N$ matrix, $\mathbf{n} \sim \mathcal{N}(0, \nu^2)$ is noise and the prior of \mathbf{x} is a zero-mean scaled adaptive Gaussian model such that

$$p(\mathbf{x}) \propto \sqrt{|\mathbf{S}|} \exp(-\frac{1}{2}\mathbf{x}^T \mathbf{S} \mathbf{x}),$$

where **S** is a diagonal matrix, whose j^{th} diagonal entry $S_i = 1/\sigma_i^2$.

This is well suited for modeling wavelet coefs.

Basic Model – Prior for σ_i

We further assume an independent prior $\propto \exp(-\frac{1}{2}\epsilon^2/\sigma_i^2)$ for each σ_j .



Figure: Illustration of the prior for σ_j .

This prior can be regarded as an approximation to the lower bounded uniform prior $U(\epsilon, +\infty)$. It is approximately uniformly distributed in the region $(3\epsilon, +\infty)$. Meanwhile it tends to 0 as σ_j approaches 0 so as to prevent σ^2 getting too small and hence avoid numerical instability in the model.

Log MAP function and basic algorithm

The prior on σ_j gives the following negative log MAP function:

Negative log MAP function

$$J(\mathbf{x}, \mathbf{S}) = \nu^2 \left(\mathbf{x}^T \mathbf{S} \mathbf{x} - \ln |\mathbf{S}| + \epsilon^2 \sum_j S_j \right) + \|\mathbf{y} - \mathbf{\Phi} \mathbf{x}\|^2$$
(4)

Minimizing $J(\mathbf{x}, \mathbf{S})$ results in the following iteration rules:

Basic algorithm

$$\mathbf{x} = \underset{\mathbf{x}}{\arg\min} J(\mathbf{x}, \mathbf{S}) = (\Phi^T \Phi + \nu^2 \mathbf{S})^{-1} \Phi^T \mathbf{y}$$

$$\sigma_j^2 = \underset{\sigma_j^2}{\arg\min} J(\mathbf{x}, \mathbf{S}) = |x_j|^2 + \epsilon^2$$

$$S_j = \frac{1}{\sigma_j^2} = \frac{1}{|x_j|^2 + \epsilon^2} \quad \forall j$$
(5)

Fast Algorithm L_0RL_2 – to avoid $(\Phi^T \Phi + \nu^2 S)^{-1}$

For wavelet-like signal spaces, we introduce the vector $\boldsymbol{\alpha} = [\alpha_1 \dots \alpha_K]$ and the diagonal operator Λ_{α} that multiplies the k^{th} subspace/subband by α_k :

$$(\Lambda_{\alpha}\mathbf{x})_k = \alpha_k \mathbf{x}_k \quad \text{for } k = 1 \cdots K$$

where \mathbf{x}_k is a masked version of \mathbf{x} with non-zeros only in the subband k, and where α_k may be optimized independently for each subband to be the minimum α_k such that

$$\alpha_k \mathbf{x}_k^T \mathbf{x}_k \geq \| \Phi \mathbf{x}_k \|^2$$
 for any k and \mathbf{x} .

Then using majorisation minimization (MM) [3], we have the following

new auxiliary function:

$$J_{lpha}(\mathbf{x},\mathbf{S},\mathbf{z}) = J(\mathbf{x},\mathbf{S}) + (\mathbf{x}-\mathbf{z})^{T}\Lambda_{lpha}(\mathbf{x}-\mathbf{z}) - \|\Phi(\mathbf{x}-\mathbf{z})\|^{2}$$
 (

(日) (周) (三) (三)

/ 28

э

6)

Fast Algorithm $\mathrm{L}_0\mathrm{RL}_2$

The new auxiliary function eliminates the difficult $\mathbf{x}^T \Phi^T \Phi \mathbf{x}$ term from J(x, S) and results in the following iteration rules:

$$L_{0}RL_{2}$$

$$\mathbf{x}_{n+1} = (\Lambda_{\alpha} + \nu^{2}\mathbf{S}_{n})^{-1} \left[(\Lambda_{\alpha} - \Phi^{T}\Phi)\mathbf{z}_{n} + \Phi^{T}\mathbf{y} \right]$$

$$\mathbf{z}_{n+1} = \arg\min_{\mathbf{z}} J_{\alpha}(\mathbf{x}_{n+1}, \mathbf{S}_{n}, \mathbf{z}) = \mathbf{x}_{n+1}$$

$$\mathbf{S}_{n+1} = \operatorname{diag}(\left[|x_{j,n+1}|^{2} + \epsilon^{2} \right]_{j=1,\cdots,N}^{-1})$$
(7a)
(7b)

Note that $(\Lambda_{\alpha} + \nu^2 \mathbf{S}_n)$ is now a diagonal matrix and hence is easy to invert!

$$J(\mathbf{x}, \mathbf{S}) = \nu^2 \left(\mathbf{x}^T \mathbf{S} \mathbf{x} - \ln |\mathbf{S}| + \epsilon^2 \sum_j S_j \right) + \|\mathbf{y} - \mathbf{\Phi} \mathbf{x}\|^2 \qquad (4)$$

$$\mathbf{x}_{n+1} = (\Lambda_{\alpha} + \nu^{2} \mathbf{S}_{n})^{-1} \left[(\Lambda_{\alpha} - \Phi^{T} \Phi) \mathbf{z}_{n} + \Phi^{T} \mathbf{y} \right]$$

$$\mathbf{z}_{n+1} = \arg\min_{\mathbf{z}} J_{\alpha} (\mathbf{x}_{n+1}, \mathbf{S}_{n}, \mathbf{z}) = \mathbf{x}_{n+1}$$

$$\mathbf{S}_{n+1} = \operatorname{diag} (\left[|x_{j,n+1}|^{2} + \epsilon^{2} \right]_{j=1,\dots,N}^{-1})$$
(7a)
(7b)

Substituting eq(7b) into log MAP eq(4) gives

cost function

$$J(\mathbf{x}, \mathbf{S}) = \nu^2 \left(\text{const} + \sum_j \ln(x_j^2 + \epsilon^2) / \epsilon \right) + \|\mathbf{y} - \mathbf{\Phi}\mathbf{x}\|^2$$
(9)

► < Ξ ►</p>

3

$$J(\mathbf{x}, \mathbf{S}) = \nu^2 \left(\mathbf{x}^T \mathbf{S} \mathbf{x} - \ln |\mathbf{S}| + \epsilon^2 \sum_j S_j \right) + \|\mathbf{y} - \mathbf{\Phi} \mathbf{x}\|^2 \qquad (4)$$

$$\mathbf{x}_{n+1} = (\Lambda_{\alpha} + \nu^{2} \mathbf{S}_{n})^{-1} \left[(\Lambda_{\alpha} - \Phi^{T} \Phi) \mathbf{z}_{n} + \Phi^{T} \mathbf{y} \right]$$

$$\mathbf{z}_{n+1} = \arg\min_{\mathbf{z}} J_{\alpha} (\mathbf{x}_{n+1}, \mathbf{S}_{n}, \mathbf{z}) = \mathbf{x}_{n+1}$$

$$\mathbf{S}_{n+1} = \operatorname{diag} (\left[|x_{j,n+1}|^{2} + \epsilon^{2} \right]_{j=1,\dots,N}^{-1})$$
(7a)
(7b)

Substituting eq(7b) into log MAP eq(4) gives

cost function

$$J(\mathbf{x}, \mathbf{S}) = \nu^2 \left(\operatorname{const} + \sum_j \ln(x_j^2 + \epsilon^2) / \epsilon \right) + \|\mathbf{y} - \Phi \mathbf{x}\|^2$$
(9)

→ Ξ →

3

Geometry of the log-sum penalty function $f_{\ln,\epsilon}$

$$f_{\ln,\epsilon} = C \ln \frac{x^2 + \epsilon^2}{\epsilon^2} \text{ with } \epsilon = 0.1, \text{ compared to } \|x\|_1 \text{ and thresholded } \|x\|_0.$$



Figure: Compared to the L1 norm, the geometry of the log-sum penalty function lends itself well to detecting sparsity as $\epsilon \rightarrow 0$.

Geometry: unit ball of $f_{\ln,\epsilon}(x)$



Figure: The effect of ϵ on the geometry of the unit ball of $f_{\ln,\epsilon}(x)$. When ϵ is large, the geometry of $f_{\ln,\epsilon}$ approximates the L2 ball; when ϵ becomes small, the geometry of $f_{\ln,\epsilon}$ approaches that of the L0 norm.

Acceleration and parameter selection

Now we consider the minimization of our problem:

$$F_{\epsilon,\nu}(\mathbf{x}) = \nu^2 \sum_j \ln \frac{|x_j|^2 + \epsilon^2}{\epsilon} + \|\mathbf{y} - \Phi \mathbf{x}\|^2,$$
(10)

where there are two parameters which decide the solution path. ν balances the fidelity and sparsity, and ϵ decides the geometry of the penalty function. We formulate a sequence of such problems $F_{\epsilon_n,\nu_n}(\mathbf{x})$:

$$F_{\epsilon_n,\nu_n}(\mathbf{x}) = \nu_n^2 \sum_j \ln \frac{|x_j|^2 + \epsilon_n^2}{\epsilon_n} + \|\mathbf{y} - \Phi \mathbf{x}\|^2$$
(11)

starting from large ν_0 and ϵ_0 , then simultaneously reducing ν_n and ϵ_n until $\nu_n = \nu$ and $\epsilon_n = \epsilon$. It is easy to see that $F_{\epsilon_n,\nu_n}(\mathbf{x})$ continuously deforms to $F_{\epsilon,\nu}(\mathbf{x})$. We try to ensure that the path of the global minima of $F_{\epsilon_n,\nu_n}(\mathbf{x})$ leads to the global minimum of $F_{\epsilon,\nu}(\mathbf{x})$.

Kingsbury & Zhang (University of Cambridge)

ACVT – Ma

11 14

| / 28

Geometry: changes of the multidimensional cost function



(a) $J(\mathbf{x}, \epsilon)$. The region where $\frac{\partial J(\mathbf{x}, \epsilon)}{\partial \epsilon} > 0$ is (b) Projection of $J(\mathbf{x}, \epsilon)$ to plane $\epsilon = 0$ coloured red.

Figure: The geometry changes of $J(\mathbf{x}, \epsilon) = \sum_{j=1}^{3} \sum_{j} \ln(x_j^2 + \epsilon^2)/\epsilon$ on a toy example: $\Phi = [3 \ 2 \ 3 \ ; 3 \ 3 \ 1.5]$ and $\mathbf{y} = [6; 6]$. The bold line connects all global minima as ϵ reduces towards zero. Sparsest solution is $\mathbf{x} = [2; 0; 0]$.

Algorithm with continuation

Let $r_L(\mathbf{x})$ denote the *L*th largest amplitude element of vector \mathbf{x} and L_{\max} be the maximum number of nonzero terms we expect in \mathbf{x} .

Set the initial $\epsilon = \| \Phi^T \mathbf{y} \|_{\infty}$, and then reduce ϵ gradually.

 ΔL controls the convergence rate and is chosen according to the size of the worst-case sparsity $L_{\rm max}$ and the desired tradeoff between convergence rate and probability of reaching the global minimum.

1-D random signal



(a) noise free

(b) noise s.d. = 0.05, SNR \approx 42dB

Figure: Plots showing how L_0RL_2 gradually selects the non-zero components of x as ϵ is reduced. The sparse input signal has 1500 elements, among which there are 45 non-zeros. This signal is similar to that used by Daubechies et al in [2]. In the example (b) with noisy observations, ν converges to 0.0576, whereas the true s.dev of the added noise is 0.05.

Kingsbury & Zhang (University of Cambridge)

ACVT – May 201

1-D random signal: noise free example

Table: Time comparison of different algorithms

	RMSE	Time (s)	RMSE	Time (s)
LORL2	8.80E-3	0.11	8.80E-7	0.19
IRLS			9.60E-5	32.50
SL0	5.20E-3	0.61	2.29E-5	0.84
IRL1			2.80E-6	0.59
L1-SpaRSA	8.90E-3	0.18	2.50E-4	0.67
L1-SPGL1	9.10E-3	0.31	1.20E-5	1.30

$$RMSE = \|\mathbf{x} - \mathbf{s}\|_2 / \|\mathbf{s}\|_2, \qquad (12)$$

where \mathbf{x} is the final estimation and \mathbf{s} is the true input signal.

1-D Heavisine signal



(a) Original signal

(b) L1-optimization

Figure: The heavisine signal has sparse representation on wavelet basis, which is well structured. Utilizing such information in signal reconstruction often results in improved performance. We use dual-tree complex wavelets (DT-CWT) for optimal performance (see our paper [4] for results with structural constraints).

1-D Heavisine signal



(c) IRL1, RMSE = 0.042

(d) $L_0 RL_2$, RMSE = 0.026

Figure: The heavisine signal has sparse representation on wavelet basis, which is well structured. Utilizing such information in signal reconstruction often results in improved performance. We use dual-tree complex wavelets (DT-CWT) for optimal performance (see our paper [4] for results with structural constraints).

Super-resolution

A real problem:

In some situations the point-spread function (psf) \mathbf{H} at a higher sampling rate is available, although \mathbf{y} at the same sampling rate is not, e.g. low inter-slice sampling rate in 3D microscope or medical scan data. Assume the observation is only available at a lower sampling rate, which we model as:

$$\mathbf{D}(\mathbf{y}) = \mathbf{D}(\mathbf{H}\mathbf{x}) + \mathbf{D}(\mathbf{n}) \tag{13}$$

where \mathbf{D} is a matrix that represents the subsampling operation. For simplicity, we denote this as:



Embedding tree structure into the algorithm

To further enhance the algorithm we can fairly easily build structural tree dependencies into the reweighting matrix \mathbf{S} , which serves as a regularization constraint to confine the support of the large coefficients. This is especially effective with the DT-CWT because of its shift-invariance, its complex coefs, and its directionally selective filter properties.

Tree-model options

- Enforce the tree structure in the model to generate **S**
- Use bivariate shrinkage to denoise $Wx^{(n+1)}$, and then use the denoised coefficient energies to update **S**.

We find the second method gives the best results so far ...

Experimental results

Super-resolving a simple ring in 2D:



(part of) original

Alternate columns are omitted from the original

Difference from original:

cubic spline

 $\mathrm{L}_0\mathrm{RL}_2$



Kingsbury & Zhang (University of Cambridge)

ACVT – May 2011

22 / 28

Results on Cameramen Gaussian filter, [16 \times 16], $\sigma_f = 1$ pel, BSNR = 40dB.



Results on Cameramen Gaussian filter, $[16 \times 16]$, $\sigma_f = 1$ pel, BSNR = 40dB.



(c) Cubic Spline interpolation of (d) (d) deconvolved result from (b) (4) E (4) E (4)

Kingsbury & Zhang (University of Cambridge)

< <>></>

Results on Cameramen Gaussian filter, [16 \times 16], $\sigma_f = 1$ pel, BSNR = 40dB.



(e) restored to higher resolution



(f) odd columns of (e)

Results on Cameramen: quantitative results

We use ISNR to quantify the improvement:

$$\begin{split} \text{ISNR}(\mathbf{x}^{(n)}) &= 10 \log_{10}(\frac{\|\mathbf{y} - \mathbf{x}\|^2}{\|\mathbf{x}^{(n)} - \mathbf{x}\|^2})\\ \text{ISNR}(\mathbf{D}(\mathbf{x}^{(n)})) &= 10 \log_{10}(\frac{\|\bar{\mathbf{y}} - \mathbf{D}\mathbf{x}\|^2}{\|\mathbf{D}\mathbf{x}^{(n)} - \mathbf{D}\mathbf{x}\|^2}) \end{split}$$

and get the following ISNR(dB) measures:

	(c)	(d)	(e)	(f)
Gaussian filter $16 imes16$	3.1734	3.2720	4.2188	4.0900
Uniform blur filter 8×8	4.3641	4.5664	5.5276	5.5179

3D microscopy dataset

- Wide-field fluorescence 3D datasets are often not equally sampled in all directions to save time and cost. (Similarly for many medical 3D datasets.)
- The 3D DT CWT representation of the dataset gives 28 planar-wave oriented subbands in each scale. These emphasize planar and linear object features, while noise and aliasing effects in other directions are suppressed by the sparsity-inducing regularization process.
- The coefficients in different subbands are denoised by bivariate shrinkage [Sendur & Selesnick 2002] to improve the quality of **S**, and thus the final result.

Original image



horizontal

Carlow Carlow

vertical

Kingsbury & Zhang (University of Cambridge)

ヨト・イヨト

æ

Spline interpolated image



Kingsbury & Zhang (University of Cambridge)

ACVT – May 2

Experimental results

Our recovery - more results available on our website



Kingsbury & Zhang (University of Cambridge)

ACVT – May 2

∃ ►

Conclusion

The L_0RL_2 algorithm

- is suitable for large-scale problems, e.g. 3D datasets, because it only requires matrix-vector multiplications and element-wise operations in each iteration. Convergence is fast for many problems.
- relies (in its basic form) only on one preset parameter, the sparsity level $L_{\rm max}$. A loose estimate of $L_{\rm max}$ is enough to achieve good results.
- is easy to implement and very flexible to allow the integration of prior knowledge of signal structure (see [4] for explanation and results).
- is well-suited to handling noisy measurement data (because of its L₂ heritage), and provides good noise suppression in the final results.

Our model turns out to be similar to hierarchical (sparse) Bayesian modelling, often solved by the Relevance Vector Machine [1].

References

- C.M. Bishop and M.E. Tipping.

Variational relevance vector machines.

In *Proceedings of the 16th Conference on Uncertainty in Artificial Intelligence*, pages 46–53. Citeseer, 2000.

- I. Daubechies, R. DeVore, M. Fornasier, and S. Gunturk. Iteratively re-weighted least squares minimization for sparse recovery. *Communications on Pure and Applied Mathematics*, 2009.

Cédric Vonesch and Michael Unser.

A fast thresholded landweber algorithm for wavelet-regularized multidimensional deconvolution.

IEEE Transactions on Image Processing, 17(4):539–549, 2008.

Y. Zhang and N. Kingsbury.

Fast I0-based sparse signal recovery.

In IEEE International Workshop on Machine Learning for Signal Processing (MLSP 2010),, Kittila, Finland., Aug 29 - Sept 1 2010.

A B M A B M