# One-Bit Compressed Sensing using the Soft-Thresholded Landweber Iteration

Andrew E. Yagle

Department of EECS, The University of Michigan, Ann Arbor, MI 48109-2122

*Abstract*— **One-bit compressed sensing reconstructs a sparse signal from only signs of linear combinations of the signal values. The problem formulates easily as a linear programming problem, but simpler methods are desirable. The binary iterative hard thresholded Landweber (BIHT) algorithm requires an oracle to know the exact sparsity of the signal. We reformulate the one-bit problem with noisy observation as a non-negative least-squares problem and solve the latter using the Landweber iteration with a non-negativity constraint, which becomes a soft-thresholded Landweber iteration with both shrinkage and thresholding. We apply this to image reconstruction from two bits of Fourer phase without an image support constraint.**

Phone: 734-763-9810. Fax: 734-763-1503.
Email: aey@eecs.umich.edu.

## I. INTRODUCTION

### A. Problem Statement

We are given as data only the *signs* $\{s_i\}$

$$s_i = \text{sign}[y_i] = \begin{cases} +1 & \text{if } y_i > 0 \\ -1 & \text{if } y_i < 0 \end{cases} \quad (1)$$

of the $M$ linear combinations $y$ of an unknown $x$

$$y = Ax \quad (2)$$

- $x$ is an unknown sparse (mostly zero) $N$-vector;
- $A$ is a known $M \times N$ full-rank matrix with $M \geq N$;
- Each component $x_i$ of $x$ is a random variable with

$$f_{x_i}(X) = C \underbrace{e^{-\lambda|X|}}_{\text{SPARSE}} \underbrace{e^{-\mu X^2}}_{\text{SMALL}}. \quad (3)$$

The Gaussian prior penalizes too-large values of $x$. The Laplacian (two-sided exponential) sparsifies $x$. The goal is to reconstruct $x$ from the signs $s_i$ of $y_i$, i.e., from only one (sign) bit of each value $y_i$ of $y$.

### B. Problem Reformulation

We assume that the $\{s_i\}$ are independent and equally likely to be $\pm 1$, and let matrix $S=\text{diag}[s_i]$. Premultiplying (2) by $S$ gives (here $|y|$ means $[|y_i|]$)

$$z = |y| = Sy = (SA)x. \quad (4)$$

- $z_i=|y_i|$. The given data is precisely that $z \geq 0$;
- $SA$ is still a known full-rank matrix ($S$ is known).

The cost functional to be minimized is then

$$\mathcal{E} = \underbrace{||\mathbf{01} + z - (SA)x||_2^2/2}_{\text{EQUALITY}} + \underbrace{\lambda||x||_1}_{\text{SPARSE}} + \underbrace{\mu||x||_2^2}_{\text{RIDGE}} \quad (5)$$

with the additional constraint $z \geq 0$ (which is the data) and where we have defined the usual two norms

$$||x||_1 = \sum_{i=1}^{N} |x_i|; \quad ||x||_2^2 = \sum_{i=1}^{N} x_i^2; \quad \mathbf{01} = \begin{bmatrix} \mathbf{0}\}M-1 \\ 1\} \quad 1 \end{bmatrix}.$$

$\mathcal{E}$ is the one-bit compressed sensing version of the elastic net (LASSO plus ridge) cost functional. The ridge provides regularization when $A$ is near-singular.

The signs $s_i$ of $y_i$ only determine $x$ to a positive scale factor. To keep $x$ from being driven to zero, the last row of the first term of (5) is $1+z_M-s_M a_M^T x$, where $a_M^T$ is the last row of $A$ and $z_M \geq 0$ makes $s_M a_M^T x \geq 1$ and keeps $x \neq 0$. This has been done previously by constraining $||x||=1$, but this is a non-convex constraint, which complicates optimization.

The $\ell_1$-norm penalty term $||x||_1$ penalizes small deviations of the elements of $x$ from zero. A considerable amount of research since 2000 has proven what the geophysical community has observed since the 1960s: The $\ell_1$ norm produces sparse solutions.

The $\ell_2$-norm penalty term $||x||_2^2$ does not penalize small deviations, since its slope is zero at zero. But it does penalize large deviations more heavily than the $\ell_1$ norm, and thus stabilizes the solution in the presence of noise. Use of an $\ell_2$ norm penalty term for this purpose is called Tikhonov regularization.

Minimization of $\mathcal{E}$ in the limit $\lambda \to 0$ with $\mu=0$ becomes minimization of $||x||_1$ subject to the equality constraint $\mathbf{01}+z=(SA)x$ and the inequality $z \geq 0$. This can be formulated as a linear programming problem and solved using any number of algorithms. Alternatively, any number of convex optimization algorithms can be applied to the minimization of $\mathcal{E}$.

However, in many applications $A$ is not represented as a matrix, but as a sequence of operations, such as wavelet or fast Fourier transforms. In this case, $Ax$ and $A^T y$ can be computed much more quickly than a typical matrix-vector multiplication, e.g., $N \log N$. This motivates use of the Landweber iteration below.

## II. Reformulation of $\mathcal{E}$ Minimization as a Non-Negative Least Squares Problem

First, we use the usual procedure of defining the positive $x^+$ and negative $x^-$ parts of $x$ as

$$x_i^+ = \begin{cases} +x_i & \text{if } x_i \geq 0 \\ 0 & \text{if } x_i \leq 0 \end{cases} \quad x_i^- = \begin{cases} -x_i & \text{if } x_i \leq 0 \\ 0 & \text{if } x_i \geq 0 \end{cases} \geq 0. \tag{6}$$

Then $x_i^+$ and $x_i^-$ are related to $x_i$ and $|x_i|$ by

$$x = x^+ - x^- \tag{7}$$

$$\|x\|_1 = \sum_{i=1}^{N}(x_i^+ + x_i^-)$$

$$\|x\|_2^2 = \|x^+\|_2^2 + \|x^-\|_2^2$$

Now consider the $(N+M) \times (2N+M)$ problem

$$\underbrace{\begin{bmatrix} \mathbf{01} \\ \frac{-\lambda}{\sqrt{2\mu}}\mathbf{1} \end{bmatrix}}_{\tilde{y}} = \underbrace{\begin{bmatrix} SA & -SA & -I \\ \sqrt{2\mu}I & \sqrt{2\mu}I & 0 \end{bmatrix}}_{\tilde{A}} \underbrace{\begin{bmatrix} x^+ \\ x^- \\ z \end{bmatrix}}_{\tilde{x}} \tag{8}$$

where we have defined the two column vectors

$$\mathbf{01} = [\underbrace{0 \ldots 0}_{M-1}\underbrace{1}_{1}]^T; \quad \mathbf{1} = [\underbrace{1,1\ldots 1}_{N}]^T \tag{9}$$

and we impose constraints $x_i^+ \geq 0$; $x_i^- \geq 0$; $z \geq 0$.

The squared $\ell_2$ error is then (note that $x_i^+ x_i^- = 0$)

$$\begin{aligned}
\|\tilde{y} - \tilde{A}\tilde{x}\|_2^2 &= \|(SA)(x^+ - x^-) - z - \mathbf{01}\|_2^2 \\
&+ \|\sqrt{2\mu}(x^+ + x^-) + \frac{\lambda}{\sqrt{2\mu}}\|_2^2 \\
&= \|(SA)x - z - \mathbf{01}\|_2^2 \\
&+ 2\mu\|x\|_2^2 + 2\lambda\|x\|_1 + \frac{N\lambda^2}{2\mu} \\
&= 2\mathcal{E} + \frac{N\lambda^2}{2\mu}. \tag{10}
\end{aligned}$$

The final term in (10) does not affect the argmax, so computing the *non-negative least-squares* solution to (8) minimizes the cost function $\mathcal{E}$. Also note that

$$(SA)^T(SA) = A^T(S^T S)A = A^T A \tag{11}$$

independent of the diagonal matrix $S$ of signs of $y_i$.

Minimization of $\mathcal{E}$ is a non-negative least-squares problem. It can be solved using Matlab code like:

```
clear;M=30;N=10;L=.0001;E=.0000005;
e=sqrt(2*E);randn('state',1);X(3)=1;X(7)=-2;
A=randn(M,N);X(N)=0;X=X';Y=A*X;
S=diag(sign(Y));AA=[S*A -S*A -eye(M)];
AA=[AA;e*eye(N) e*eye(N) zeros(N,M)];
YY=[zeros(M-1,1);1;-L/e*ones(N,1)];
Z=lsnonneg(AA,YY);[X Z(1:N)-Z(N+1:2*N)]
```

## III. Solution of the Non-Negative Least-Squares Using Landweber

### A. Review of Landweber Iteration

The basic Landweber iteration is

$$x^{k+1} = x^k + A^T(y - Ax), \quad x^0 = 0 \tag{12}$$

where $x^k$ is the estimate of $x$ at the $k^{th}$ iteration. The Landweber iteration can be viewed as a steepest descent algorithm for minimizing the cost function

$$\begin{aligned}
f(x) &= (1/2)\|y - Ax\|_2^2 \\
\nabla f(x) &= -A^T(y - Ax). \tag{13}
\end{aligned}$$

The basic steepest descent algorithm is

$$\begin{aligned}
x^{k+1} &= x^k - \nabla f(x) \\
&= x^k + A^T(y - Ax). \tag{14}
\end{aligned}$$

Here, we use the basic Landweber iteration

$$\tilde{x}^{k+1} = \tilde{x}^k + \tilde{A}^T(\tilde{y} - \tilde{A}\tilde{x}^k) \tag{15}$$

with a non-negativity constraint at each iteration

$$\tilde{x}_i^{k+1} = \max[\tilde{x}_i^{k+1}, 0]. \tag{16}$$

We call this the *non-negative* Landweber iteration.

Since the cost functional $f(x)$ and the non-negativity constraints $x_i^+ \geq 0$ and $x_i^- \geq 0$ are all convex, the non-negative Landweber iteration will converge if the maximum eigenvalue of $\tilde{A}^T \tilde{A} \leq 2$.

### B. Application to Present Problem

Substitution of (8) in (15) gives after some algebra

$$\begin{aligned}
(x^+)^{k+1} &= (x^+)^k - A^T A[(x^+)^k - (x^-)^k] + A^T S\mathbf{01} \\
&+ A^T Sz^k - \lambda\mathbf{1} - 2\mu[(x^+)^k + (x^-)^k] \\
(x^-)^{k+1} &= (x^-)^k + A^T A[(x^+)^k - (x^-)^k] - A^T S\mathbf{01} \\
&- A^T Sz^k - \lambda\mathbf{1} - 2\mu[(x^+)^k + (x^-)^k] \\
z^{k+1} &= SA[(x^+)^k - (x^-)^k] - \mathbf{01} \tag{17}
\end{aligned}$$

followed by the three non-negativity constraints

$$\begin{aligned}
(x^+)_i^{k+1} &= \max[(x^+)_i^{k+1}, 0] \\
(x^-)_i^{k+1} &= \max[(x^-)_i^{k+1}, 0] \\
z_i^{k+1} &= \max[z_i^{k+1}, 0] \tag{18}
\end{aligned}$$

so most of the computations in the Landweber iteration are the matrix-vector multiplications $Ax$ and $A^T Sz$, which can often be implemented using a fast algorithm such as the fast Fourier transform, the fast wavelet algorithm or a sparse matrix-times-vector.

## C. Convergence of Non-Negative Landweber

$\tilde{A}$ is underdetermined. The nonzero eigenvalues of $\tilde{A}^T\tilde{A}$ are the eigenvalues of $\tilde{A}\tilde{A}^T$, which from

$$
\begin{aligned}
\tilde{A}\tilde{A}^T &= \begin{bmatrix} SA & -SA & -I \\ \sqrt{2\mu}I & \sqrt{2\mu}I & 0 \end{bmatrix} \begin{bmatrix} A^TS & \sqrt{2\mu}I \\ -A^TS & \sqrt{2\mu}I \\ -I & 0 \end{bmatrix} \\
&= \begin{bmatrix} 2(SA)(A^TS)+I & 0 \\ 0 & 4\mu I \end{bmatrix} \qquad (19)
\end{aligned}
$$

are double those of $(SA)(A^TS)$ plus one, and $4\mu$. The non-negative Landweber iteration will converge if $\mu$ and all eigenvalues of $(SA)(A^TS)$ are less than $\frac{1}{2}$.

Since $SA$ is overdetermined, the nonzero eigenvalues of $(SA)(SA)^T$ in (19) are the eigenvalues of $(A^TS)(SA)=A^TA$, independent of data $S$.

## D. Derivation of Thresholded Landweber

Let $\mu=0$ and recall $x^k=(x^+)^k-(x^-)^k$ and $-(x^-)^k$ are the negative values of $x^k$. Then (17) becomes

$$
\begin{aligned}
x^{k+1} &= x^k - A^TAx^k + A^TS(\mathbf{01}+z^k) \\
&= x^k + A^TS[(\mathbf{01}+z^k)-SAx^k] \\
z^{k+1} &= SAx^k - \mathbf{01} \qquad (20)
\end{aligned}
$$

followed by shrinkage of $|x^k|$ by $\lambda$, thresholding values of $|x^k| < \lambda$ to 0, and imposing non-negative $z^{k+1}$:

$$
\begin{aligned}
x_i^{k+1} &= \begin{cases} x_i^{k+1}-\lambda & \text{if } x_i^{k+1} > +\lambda \\ x_i^{k+1}+\lambda & \text{if } x_i^{k+1} < -\lambda \\ 0 & \text{if } |x_i^{k+1}| < \lambda \end{cases} \\
z^{k+1} &= \begin{cases} z^{k+1} & \text{if } z^{k+1} \geq 0 \\ 0 & \text{if } z^{k+1} \leq 0 \end{cases} \qquad (21)
\end{aligned}
$$

This is similar to soft-thresholded Landweber with

- "Observations" $(\mathbf{01}+z^k)$ instead of $y$, where the
- Signs of $y_i$ (the data) are included in $SA$ and
- $z^{k+1} \geq 0$ forces $(SA)x^k \geq 0$, which is the data.

## E. Orthonormal Columns

Let overdetermined $A$ have orthonormal columns $A^TA = I$. Then the above algorithm simplifies to

$$
\begin{aligned}
z^{k+1} &= (SA)x^k - \mathbf{01} \\
x^{k+1} &= (A^TS)(z^k + \mathbf{01}) \qquad (22)
\end{aligned}
$$

along with thresholding and shrinkage steps (21).

If $z^k$ is replaced with $z^{k+1}$ in the second equation, this is similar to an alternating-projections algorithm with a sign constraint in one domain and thresholding and shrinkage (for sparsity) in the other domain.

Multiplication by $SA$ allows the sign constraint to be implemented as a non-negativity constraint. Multiplication by $A^TS$ undoes multiplication by $S$ in $SA$.

## IV. IMAGE RECONSTRUCTION FROM TWO BITS OF DISCRETE FOURIER TRANSFORM PHASE

### A. Background

The problem of reconstructing a signal or image from two bits of the phase of its Fourier transform arises in phase-shifting digital holography. This reduced bit-depth allows real-time reconstruction by a spatial light modulator, which modulates phase.

Previous algorithms on reconstruction from two bits of Fourier phase have been alternating projection algorithms that impose a support constraint in the spatial domain and a sign constraint on the phase in the Fourier domain. If the support is chosen to be a half-plane, then one bit of phase is sufficient, since the even part of the image with such a constraint determines the entire image. This seems somewhat contrived. Knowledge of one bit of phase amounts to approximate knowledge of the zero crossings of the real part of the Fourier transform, which under mild assumptions determines the real part of the image to an overall scale factor. Two bits of phase amounts to approximate knowledge of the zero crossings of the real and imaginary parts of the Fourier transform, which determine the even and odd parts of the image, respectively, to separate scale factors for each.

Here we replace the support constraint (which is usually not available) with a sparsity constraint, which consists of soft thresholding. This arises directly from the above algorithm. The matrix $A$ is the concatenation of a 2-D discrete Fourier transform and the inverse of a sparsifying operator, such as a wavelet transform. Both of these can be implemented directly as algorithms, so that the matrix $A$ need never be formed or stored. $A^H$ is then an inverse 2-D discrete Fourier transform followed by a wavelet transform, again implemented as algorithms.

### B. Problem Formulation

The $N \times N$-point 2-D discrete Fourier transform of the $M \times M$ image $x[m,n]$ is defined as

$$
X_{k_1,k_2} = \sum_{m=0}^{M-1}\sum_{n=0}^{M-1} x[m,n]e^{-j\frac{2\pi}{N}(mk_1+nk_2)}. \qquad (23)
$$

The inverse 2-D discrete Fourier transform is

$$
x[m,n] = \frac{1}{N^2}\sum_{k_1=0}^{N-1}\sum_{k_2=0}^{N-1} X_{k_1,k_2}e^{j\frac{2\pi}{N}(mk_1+nk_2)}. \qquad (24)
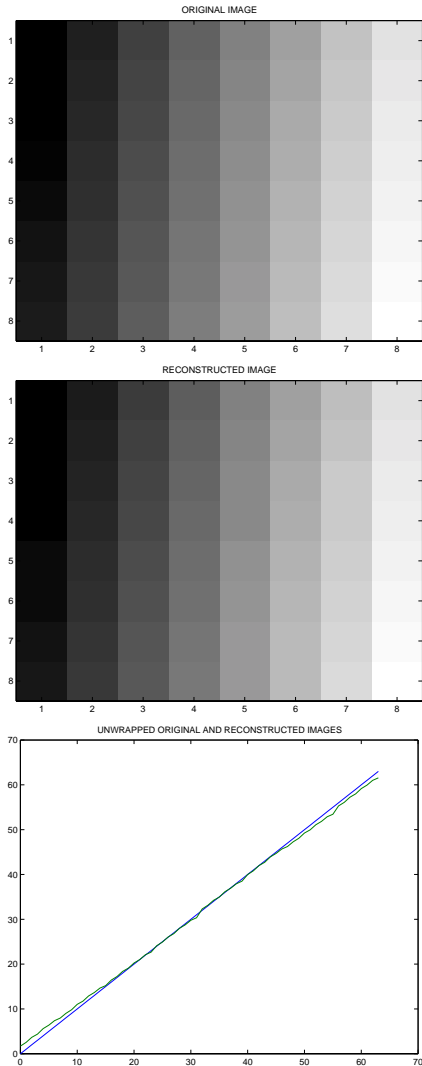$$

The 2-D DFT is a unitary linear transformation, except for a factor of $\frac{1}{N}$ (these are collected in (24)).

"Two bits of phase" means the signs of $\text{Re}[X_{k_1,k_2}]$ and $\text{Im}[X_{k_1,k_2}]$. This specifies the phase as lying in one of four quadrants, hence "two bit" quantization.