
Bayesian Compressive Sensing and Projection Optimization

Shihao Ji
Lawrence Carin

SHJI@ECE.DUKE.EDU
LCARIN@ECE.DUKE.EDU

Department of Electrical and Computer Engineering, Duke University, Durham, NC 27708 USA

Abstract

This paper introduces a new problem for which machine-learning tools may make an impact. The problem considered is termed “compressive sensing”, in which a real signal of dimension N is measured accurately based on $K \ll N$ real measurements. This is achieved under the assumption that the underlying signal has a sparse representation in some basis (e.g., wavelets). In this paper we demonstrate how techniques developed in machine learning, specifically sparse Bayesian regression and active learning, may be leveraged to this new problem. We also point out future research directions in compressive sensing of interest to the machine-learning community.

1. Introduction

Over the last two decades there have been significant advances in the development of orthonormal bases for compact representation of a wide class of discrete signals. An important example of this is the wavelet transform [Mallat, 1998], with which general signals are represented in terms of atomic elements localized in time and frequency, yielding highly compact representations of many natural signals. Let the $N \times N$ matrix \mathbf{B} represent a wavelet basis, with basis functions defined by associated columns; a general signal $\mathbf{f} \in \mathbb{R}^N$ may be represented as $\mathbf{f} = \mathbf{B}\mathbf{w}$, where $\mathbf{w} \in \mathbb{R}^N$ represents the wavelet and scaling function coefficients [Mallat, 1998]. For most natural signals \mathbf{f} , most components of the vector \mathbf{w} have negligible amplitude. Therefore, if $\hat{\mathbf{w}}$ represents the weights \mathbf{w} with the smallest $N-M$ coefficients set to zero, and $\hat{\mathbf{f}} = \mathbf{B}\hat{\mathbf{w}}$, then the relative error $\|\mathbf{f} - \hat{\mathbf{f}}\|_2 / \|\mathbf{f}\|_2$ is often negligibly small for $M \ll N$. This property has led to the develop-

ment of state-of-the-art compression algorithms based on wavelet-based transform coding [Pearlman et al., 2004].

In conventional applications one first measures the N -dimensional signal \mathbf{f} , \mathbf{f} is then transformed into a new basis where the signal is sparse, and the new basis coefficients \mathbf{w} are then quantized [Pearlman et al., 2004]. This invites the following question: If the underlying signal is ultimately compressible, is it possible to perform a compact (“compressive”) set of measurements directly, thereby offering the potential to simplify the sensing system (reduce the number of required measurements)? This question has recently been answered in the affirmative [Candès et al., 2006, Donoho, 2006], introducing the field of compressive sensing (CS).

In its earliest form the relationship between the underlying signal \mathbf{f} and the CS measurements \mathbf{g} has been constituted through random projections [Candès et al., 2006, Donoho, 2006]. Specifically, assume that the signal \mathbf{f} is compressible in some basis \mathbf{B} (not necessarily a wavelet basis), the k th CS measurement g_k (k th component of \mathbf{g}) is constituted by projecting \mathbf{f} onto a “random” basis that is constituted with a “random” linear combination of the basis functions in \mathbf{B} , i.e., $g_k = \mathbf{f}^T(\mathbf{B}\mathbf{r}_k)$, where $\mathbf{r}_k \in \mathbb{R}^N$ is a column vector with each element an i.i.d. draw of a random variable, with arbitrary alphabet (e.g., real or binary) [Tsaig & Donoho, 2006].

Based on the above discussion, the CS measurements may be represented as $\mathbf{g} = \mathbf{\Phi}\mathbf{B}^T\mathbf{f} = \mathbf{\Phi}\mathbf{w}$, where $\mathbf{\Phi} = [\mathbf{r}_1 \dots \mathbf{r}_K]^T$ is an $K \times N$ matrix, assuming K random measurements are made. Since typically $K < N$ we have fewer measurements than degrees of freedom for the signal \mathbf{f} . Therefore, inversion for the weights \mathbf{w} (and hence \mathbf{f}) is ill-posed. However, if one exploits the fact that \mathbf{w} is sparse with respect to a known orthonormal basis \mathbf{B} , then one may approximate \mathbf{w} accurately via an ℓ_1 -regularized formulation [Donoho, 2006]

$$\tilde{\mathbf{w}} = \arg \min_{\mathbf{w}} \{ \|\mathbf{g} - \mathbf{\Phi}\mathbf{w}\|_2^2 + \rho \|\mathbf{w}\|_1 \}, \quad (1)$$

where the scalar ρ controls the relative importance ap-

Appearing in *Proceedings of the 24th International Conference on Machine Learning*, Corvallis, OR, 2007. Copyright 2007 by the author(s)/owner(s).

plied to the Euclidian error and the sparseness term. This basic framework has been the starting point for several recent CS inversion algorithms, including linear programming [Chen et al., 1999] and greedy algorithms [Tropp & Gilbert, 2005, Donoho et al., 2006], for a point estimate of the weights \mathbf{w} .

In the discussion that follows we demonstrate that the solution to this problem may exploit many of the tools developed recently in the machine-learning community, specifically sparse Bayesian regression [Tipping, 2001, Wipf et al., 2004] and active learning [Fedorov, 1972, MacKay, 1992]. Moreover, the results of this machine-learning-based CS analysis significantly advance the state of the art in CS. The encouraging nature of this initial analysis is meant to introduce the machine-learning community to a new problem for which it is poised to make an important contribution.

2. Compressive-Sensing Inversion From Bayesian Viewpoint

2.1. Compressive Sensing as Linear Regression

It was assumed at the start that \mathbf{f} is compressible in the basis \mathbf{B} . Therefore, let \mathbf{w}_s represent an N -dimensional vector that is identical to the vector \mathbf{w} for the M elements in \mathbf{w} with largest magnitude; the remaining $N-M$ elements in \mathbf{w}_s are set to zero. Similarly, we introduce a vector \mathbf{w}_e that is identical to \mathbf{w} for the smallest $N-M$ elements in \mathbf{w} , with all remaining elements of \mathbf{w}_e set to zero. We therefore have $\mathbf{w} = \mathbf{w}_s + \mathbf{w}_e$, and

$$\mathbf{g} = \Phi \mathbf{w} = \Phi \mathbf{w}_s + \Phi \mathbf{w}_e = \Phi \mathbf{w}_s + \mathbf{n}_e, \quad (2)$$

where $\mathbf{n}_e = \Phi \mathbf{w}_e$. Since it was assumed at the start that Φ is constituted through random samples, the components of \mathbf{n}_e may be approximated as a zero-mean Gaussian noise as a consequence of Central-Limit Theorem [Papoulis & Pillai, 2002] for large $N-M$. We also note that the CS measurements may be noisy, with the measurement noise, denoted by \mathbf{n}_m , represented by a zero-mean Gaussian distribution, and therefore

$$\mathbf{g} = \Phi \mathbf{w}_s + \mathbf{n}_e + \mathbf{n}_m = \Phi \mathbf{w}_s + \mathbf{n}, \quad (3)$$

where the components of \mathbf{n} are approximated as a zero-mean Gaussian noise with unknown variance σ^2 .

The above analysis has converted the CS problem of inverting for the sparse weights \mathbf{w}_s into a linear-regression problem with a constraint (prior) that \mathbf{w}_s is sparse, or more relevantly, sparse Bayesian regression [Tipping, 2001, Wipf et al., 2004]. Assuming knowledge of Φ , the quantities to be estimated based

on the CS measurements \mathbf{g} are the sparse weights \mathbf{w}_s and the noise variance σ^2 . In a Bayesian analysis we seek a full posterior density function for \mathbf{w}_s and σ^2 .

2.2. Sparseness Prior and MAP approximation

In a Bayesian formulation our understanding of the fact that \mathbf{w}_s is sparse is formalized by placing a sparseness-promoting prior on \mathbf{w}_s . A widely used sparseness prior is the Laplace density function [Figueiredo, 2002]:

$$p(\mathbf{w}|\lambda) = (\lambda/2)^N \exp(-\lambda \sum_{i=1}^N |w_i|), \quad (4)$$

where in (4) and henceforth we drop the subscript s on \mathbf{w} , recognizing that we are always interested in a sparse solution for the weights. Given the CS measurements \mathbf{g} , and assuming the model in (3), it is straightforward to demonstrate that the solution in (1) corresponds to a maximum *a posteriori* (MAP) estimate for \mathbf{w} using the prior in (4).

3. Estimate of Full Posterior for Sparse Weights

3.1. Hierarchical Sparseness Prior

The above discussion demonstrated that conventional CS inversion for the weights \mathbf{w} corresponds to a MAP approximation to a Bayesian linear-regression analysis, with a Laplace sparseness prior on \mathbf{w} . This then raises the question of whether the Bayesian analysis may be carried further, to realize an estimate of the full posterior on \mathbf{w} and σ^2 .

Rather than imposing a Laplace prior on \mathbf{w} , we develop a hierarchical prior [Tipping, 2001, Figueiredo, 2002] that has similar properties but that allows convenient conjugate-exponential analysis. Specifically, we introduce the prior

$$p(\mathbf{w}|\boldsymbol{\alpha}) = \prod_{i=1}^N \mathcal{N}(w_i|0, \alpha_i^{-1}), \quad (5)$$

where $\mathcal{N}(w_i|0, \alpha_i^{-1})$ is a zero-mean Gaussian density function with precision (inverse-variance) α_i . We further place the following Gamma prior on $\boldsymbol{\alpha}$

$$p(\boldsymbol{\alpha}|a, b) = \prod_{i=1}^N \Gamma(\alpha_i|a, b). \quad (6)$$

The overall prior on \mathbf{w} is evaluated as

$$p(\mathbf{w}|a, b) = \prod_{i=1}^N \int_0^\infty \mathcal{N}(w_i|0, \alpha_i^{-1}) \Gamma(\alpha_i|a, b) d\alpha_i. \quad (7)$$

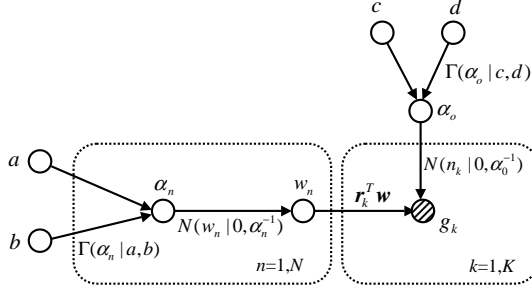


Figure 1. Graphical model of the Bayesian CS formulation.

Density function $\Gamma(\alpha_i|a, b)$ is the conjugate prior for α_i , when w_i plays the role of observed data and $\mathcal{N}(w_i|0, \alpha_i^{-1})$ is a likelihood function; consequently the integral $\int_0^\infty \mathcal{N}(w_i|0, \alpha_i^{-1})\Gamma(\alpha_i|a, b)d\alpha_i$ can be evaluated analytically, and it corresponds to the Student- t distribution [Tipping, 2001]. With appropriate choice of the hyperparameters a and b , the Student- t distribution is strongly peaked about $w_i=0$, and therefore the prior in (7) favors most w_i being zero (i.e., it is a sparseness prior). Similarly, we introduce the Gamma prior $\Gamma(\alpha_0|c, d)$ on the inverse of the noise variance $\alpha_0=1/\sigma^2$.

To see the advantage of the sparseness prior in (7), consider the graphical structure of the model as reflected in Fig. 1, for generation of the observed data \mathbf{g} . Following consecutive blocks in Fig. 1 (following the direction of the arrows), let p_k represent the parameter associated with block k , and p_{k+1} represents the next parameter in the sequence. For all steps in Fig. 1, the density function for p_k is the conjugate prior for the likelihood defined in terms of the density function for p_{k+1} , assuming that all parameters except p_k are held constant (i.e., all parameters other than p_k temporarily play the role of fixed data). This structural form is very convenient for implementing iterative algorithms for evaluation of the posterior density function for \mathbf{w} and α_0 . For example, one may conveniently implement a Markov Chain Monte Carlo (MCMC) [Gilks et al., 1996] or, more efficiently and approximately, a variational Bayesian (VB) analysis [Bishop & Tipping, 2000]. While the VB analysis is efficient relative to MCMC, we here consider a type-II maximum-likelihood (ML) analysis, with the objective of achieving highly efficient computations while still preserving accurate results [Tipping, 2001].

3.2. The Bayesian CS algorithm

As shown by Tipping [2001], in the context of the relevance vector machine (RVM), if \mathbf{g} , $\boldsymbol{\alpha}$ and $\alpha_0=1/\sigma^2$ are known, then the posterior for \mathbf{w} can be expressed

analytically as a multivariate Gaussian distribution with mean and covariance:

$$\boldsymbol{\mu} = \alpha_0 \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{g}, \quad (8)$$

$$\boldsymbol{\Sigma} = (\mathbf{A} + \alpha_0 \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1}, \quad (9)$$

where $\mathbf{A} = \text{diag}(\alpha_1, \dots, \alpha_N)$. Further, the marginal likelihood for $\boldsymbol{\alpha}$ and α_0 , or equivalently, its logarithm $\mathcal{L}(\boldsymbol{\alpha}, \alpha_0)$ can be expressed analytically by integrating out the weights \mathbf{w} , to yield

$$\begin{aligned} \mathcal{L}(\boldsymbol{\alpha}, \alpha_0) &= \log p(\mathbf{g}|\boldsymbol{\alpha}, \alpha_0) = \log \int p(\mathbf{g}|\mathbf{w}, \alpha_0) p(\mathbf{w}|\boldsymbol{\alpha}) d\mathbf{w} \\ &= -\frac{1}{2} [K \log 2\pi + \log |\mathbf{C}| + \mathbf{g}^T \mathbf{C}^{-1} \mathbf{g}], \quad (10) \end{aligned}$$

with $\mathbf{C} = \sigma^2 \mathbf{I} + \boldsymbol{\Phi} \mathbf{A}^{-1} \boldsymbol{\Phi}^T$. Then a type-II ML solution employs the point estimates for $\boldsymbol{\alpha}$ and α_0 that maximize (10). This can be implemented readily via the EM algorithm or direct differentiation [Tipping, 2001], to yield:

$$\alpha_i^{new} = \gamma_i / \mu_i^2, \quad i \in \{1, 2, \dots, N\}, \quad (11)$$

with $\gamma_i \triangleq 1 - \alpha_i \boldsymbol{\Sigma}_{ii}$, where $\boldsymbol{\Sigma}_{ii}$ is the i th diagonal element from $\boldsymbol{\Sigma}$ in (9), and

$$1/\alpha_0^{new} = \frac{\|\mathbf{g} - \boldsymbol{\Phi} \boldsymbol{\mu}\|_2^2}{N - \sum_i \gamma_i}. \quad (12)$$

Note that $\boldsymbol{\alpha}^{new}$ and α_0^{new} are a function of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, while $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are a function of $\boldsymbol{\alpha}$ and α_0 ; this suggests an iterative algorithm, where one iterates between (8)-(9) and (11)-(12), and in this process α_i becomes very large for those w_i that have insignificant amplitudes for representation of $\mathbf{g} = \boldsymbol{\Phi} \mathbf{w}$. Only a relatively small set of w_i , for which the corresponding α_i remains relatively small, contribute for representation of \mathbf{g} , and the level of sparseness (size of M) is determined automatically [Wipf et al., 2004]. It is also important to note that, as a result of the type-II ML solution, the point estimates (rather than the posterior densities) of $\boldsymbol{\alpha}$ and α_0 are sought. Therefore, there is no need to set a , b , c and d on the Gamma hyperpriors.

While it is useful to have a measure of uncertainty in the weights \mathbf{w} , the quantity of most interest is the signal $\mathbf{f} = \mathbf{B} \mathbf{w}$. Since \mathbf{w} is drawn from a multivariate Gaussian distribution with mean and covariance defined in (8)-(9), then \mathbf{f} is also drawn from a multivariate Gaussian distribution, with mean and covariance

$$E(\mathbf{f}) = \mathbf{B} \boldsymbol{\mu}, \quad (13)$$

$$\text{Cov}(\mathbf{f}) = \mathbf{B} \boldsymbol{\Sigma} \mathbf{B}^T. \quad (14)$$

The diagonal elements of the covariance matrix in (14) provide ‘‘error bars’’ on the accuracy of the inversion of \mathbf{f} , as represented in terms of its mean.

While the iterative algorithm described above has been demonstrated to yield a highly accurate sparse linear-regression representation [Tipping, 2001], we note the following practical limitation. When evaluating (9) one must invert matrices of size $N \times N$, which has complexity $\mathcal{O}(N^3)$, thereby making this approach relatively slow for data \mathbf{f} of large dimension N . This motivates development of a fast Bayesian algorithm with the objective of achieving highly efficient computations comparable to existing CS algorithms, such as OMP [Tropp & Gilbert, 2005] and StOMP [Donoho et al., 2006].

Fortunately, this fast Bayesian algorithm has been developed in [Faul & Tipping, 2002] by analyzing the properties of the marginal likelihood function in (10). This enables a principled and efficient sequential addition and deletion of candidate basis function (columns of Φ) to monotonically maximize the marginal likelihood. Because of limited space, we omit the detailed discussion of this fast algorithm, and refer the reader to [Tipping & Faul, 2003] for more information. We here only briefly summarize some of its key properties. Compared with the iterative algorithm presented above, the fast algorithm operates in a constructive manner, i.e., sequentially adds terms to the model until all M nonzero weights have been added. So, the complexity of the algorithm is more related to M than N . Further, by using the matrix inverse identity, the inverse operation in (9) has been implemented by an iterative update formula with reduced complexity. Detailed analysis shows that this fast algorithm has complexity $\mathcal{O}(NM^2)$, which is more efficient than the original RVM, especially when the underlying signal is truly sparse ($M \ll N$). Unlike other related CS algorithms (e.g., OMP and StOMP), the fast algorithm has the operation of deleting a basis function from the model (i.e., setting $\alpha_i = \infty$). This deletion operation is the likely explanation for the improvement in sparsity of this fast algorithm demonstrated in the experiments (see Sec. 5). More importantly, while conventional CS algorithms yield a point estimate for \mathbf{f} , the Bayesian analysis considered here also yields the error bars defined in (14).

4. Adaptive Compressive Sensing

4.1. Selecting Projections to Reduce Signal Uncertainty

In the original CS construction, the projections represented by $\Phi = [\mathbf{r}_1 \dots \mathbf{r}_K]^T$ were constituted via i.i.d. realizations of an underlying random variable [Tsaig & Donoho, 2006]. In addition, previous CS algorithms [Donoho, 2006, Tropp & Gilbert, 2005, Donoho

et al., 2006] focused on estimating \mathbf{w} have employed a point estimate like that in (1); such approaches do not provide a measure of uncertainty in \mathbf{f} , and therefore adaptive design of Φ was not feasible. The Bayesian CS (BCS) algorithm introduced in Sec. 3.2 allows efficient computation of \mathbf{f} and associated error bars, and therefore one may consider the possibility of adaptively selecting projections \mathbf{r}_k with the goal of reducing uncertainty. Such a framework has been studied previously in machine learning under the name of experiment design or active learning [Fedorov, 1972, MacKay, 1992]. Further, the error bars also give a way to determine how many measurements are enough for faithful CS reconstruction, i.e., when the change in the uncertainty is not significant, it may be assumed that one is simply reconstructing the noise \mathbf{n} in (3), and therefore the adaptive sensing may be stopped.

As discussed above, the estimated posterior on the signal \mathbf{f} is a multivariate Gaussian distribution, with mean $E(\mathbf{f}) = \mathbf{B}\boldsymbol{\mu}$ and covariance $Cov(\mathbf{f}) = \mathbf{B}\boldsymbol{\Sigma}\mathbf{B}^T$. The differential entropy [Cover & Thomas, 1991] for \mathbf{f} therefore satisfies:

$$\begin{aligned} h(\mathbf{f}) &= \frac{1}{2} \log |\mathbf{B}\boldsymbol{\Sigma}\mathbf{B}^T| + c = \frac{1}{2} \log |\boldsymbol{\Sigma}| + c \\ &= -\frac{1}{2} \log |\mathbf{A} + \alpha_0 \Phi^T \Phi| + c, \end{aligned} \quad (15)$$

where c is a constant, independent of Φ . Recall that $\mathbf{A} = \text{diag}(\alpha_1, \dots, \alpha_N)$, and therefore the dependence of the differential entropy on the observed CS measurements \mathbf{g} is defined by the point estimates of $\boldsymbol{\alpha}$ and α_0 (from the type-II ML estimates discussed in Sec. 3.2).

We may now ask which new projection \mathbf{r}_{K+1} would be optimal for minimizing the differential entropy in (15). Toward this end, we augment Φ by adding a $(K+1)$ th row represented by \mathbf{r}_{K+1}^T . If we let $h_{new}(\mathbf{f})$ represent the new differential entropy as a consequence of adding this new projection measurement, via the matrix determinant identity we have

$$h_{new}(\mathbf{f}) = h(\mathbf{f}) - \frac{1}{2} \log [1 + \alpha_0 \mathbf{r}_{K+1}^T \boldsymbol{\Sigma} \mathbf{r}_{K+1}], \quad (16)$$

where α_0 and $\boldsymbol{\Sigma}$ are based on estimates found using the previous K projections. In order to minimize (16) the next projection \mathbf{r}_{K+1} must be designed to maximize the variance of the *expected* measurement g_{K+1} since

$$\mathbf{r}_{K+1}^T \boldsymbol{\Sigma} \mathbf{r}_{K+1} = \mathbf{r}_{K+1}^T Cov(\mathbf{w}) \mathbf{r}_{K+1} = \text{Var}(g_{K+1}). \quad (17)$$

In other words, \mathbf{r}_{K+1} must be selected to constitute the measurement g_{K+1} for which the data is most uncertain (and hence access to the associated measurement would be most informative).

There are multiple ways this adaptive algorithm may be utilized in practice. If it is possible to design new projections \mathbf{r}_{K+1} adaptively “on the fly”, then one might perform an eigen-decomposition of the matrix $\mathbf{\Sigma}$, and select for representation of \mathbf{r}_{K+1} the eigenvector with largest eigenvalue. Alternatively, if from a hardware standpoint such flexibility in design of \mathbf{r}_{K+1} is not feasible, then one might *a priori* design a library \mathbf{L} of possible next projections, with \mathbf{r}_{K+1} selected from \mathbf{L} with the goal of maximizing (17).

An additional issue needs to be clarified if the eigenvector of $\mathbf{\Sigma}$ is used for the next projection \mathbf{r}_{K+1} . Because of the sparse Bayesian solution, $\mathbf{\Sigma}$ only employs elements corresponding to the associated nonzero components of \mathbf{w} found based on the fast algorithm (i.e., $\mathbf{\Sigma}$ is reduced in general to a small matrix). So when constructing the next projection based on the eigenvector, some entries of \mathbf{r}_{K+1} will be empty. If we impute all those empty entries with zeros, we are under the risk of being wrong. The initial estimate of \mathbf{w} can be inaccurate; if we impute with zeros, the estimate will be always biased and has no chance to be corrected since the corresponding contributions from underlying true \mathbf{w} are always ignored. To mitigate this problem, we impute those empty entries with random samples drawn i.i.d. from a Gaussian distribution $\mathcal{N}(0, 1)$. After the imputation, we re-scale the magnitude of the imputed entries to 0.01. In this way, we utilize the optimized projection, and at the same time allow some contributions from the empty entries. Overall, the final projection \mathbf{r}_{K+1} has the magnitude $\|\mathbf{r}_{K+1}\|_2 = 1.01$.

4.2. Approximate Adaptive CS

The error bars on the estimate of \mathbf{f} play a critical role in implementing the above adaptive CS scheme, with this a direct product from the Bayesian analysis. Since there are established CS algorithms based on a point estimate of \mathbf{w} , one may ask whether these algorithms may be modified, utilizing insights from the Bayesian analysis. The advantage of such an approach is that, if possible, one would access some of the advantages of the Bayesian analysis, in an approximate sense, while being able to retain the advantages of existing CS algorithms.

The adaptive algorithm in (16) relies on computation of the covariance matrix $\mathbf{\Sigma} = (\mathbf{A} + \alpha_0 \mathbf{\Phi}^T \mathbf{\Phi})^{-1}$; since $\mathbf{\Phi}$ is assumed known, this indicates that what is needed are estimates for α_0 and α , the latter required for the diagonal matrix \mathbf{A} . From (12) we have $\sigma^2 = 1/\alpha_0 = \|\mathbf{g} - \mathbf{\Phi}\boldsymbol{\mu}\|_2^2 / (N - \sum_i \gamma_i)$, where the denominator $N - \sum_i \gamma_i$ may be viewed as an estimate for the number of components of the weight vector \mathbf{w} that

have negligible amplitude. Consequently, assume that a CS algorithm such as OMP or StOMP is used to yield a point estimate of the weights \mathbf{w} , denoted \mathbf{w}_p , and assume that there are M_0 non-zero elements in \mathbf{w}_p ; then one may approximate the “noise” variance as $\sigma^2 = \|\mathbf{g} - \mathbf{\Phi}\mathbf{w}_p\|_2^2 / (N - M_0)$.

Concerning the diagonal matrix \mathbf{A} , it may be viewed as a regularization of the matrix $(\alpha_0 \mathbf{\Phi}^T \mathbf{\Phi})$, to assure that the matrix inversion is well posed. While the Bayesian analysis in Sec. 3.2 indicates that the loading represented by \mathbf{A} should be non-uniform, we may simply make \mathbf{A} diagonalized uniformly, with value corresponding to a small fraction of the average value of the diagonal elements of $(\alpha_0 \mathbf{\Phi}^T \mathbf{\Phi})$. In Sec. 5, when presenting example results, we make comparisons between the rigorous implementation discussed in Sec. 4.1 and the approximate scheme discussed here (as applied to the BCS algorithm). However, similar modifications may be made to other related algorithms, such as OMP and StOMP.

5. Example Results

We test the performance of BCS on several example problems considered widely in the CS literature, with comparisons made to Basis Pursuit (BP) [Chen et al., 1999] and StOMP [Donoho et al., 2006]. While BP is a relatively computationally expensive algorithm that involves linear programming, StOMP may be one of the state-of-the-art fast CS algorithms. In the experiments we evaluate the reconstruction error as $\|\mathbf{f}_{method} - \mathbf{f}\|_2 / \|\mathbf{f}\|_2$. All the computations presented here were performed using Matlab run on a 3.4GHz Pentium machine.

5.1. BCS and Projection Optimization

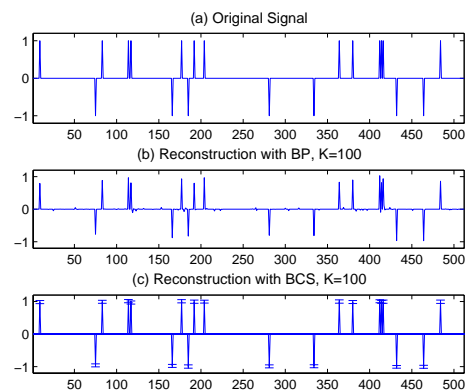


Figure 2. Reconstruction of Spikes. (a) Original signal; (b) Reconstruction with BP, $err_{BP}=0.158$, $t_{BP}=1.56$ secs; (c) Reconstruction with BCS, $err_{BCS}=0.015$, $t_{BCS}=0.63$ secs.

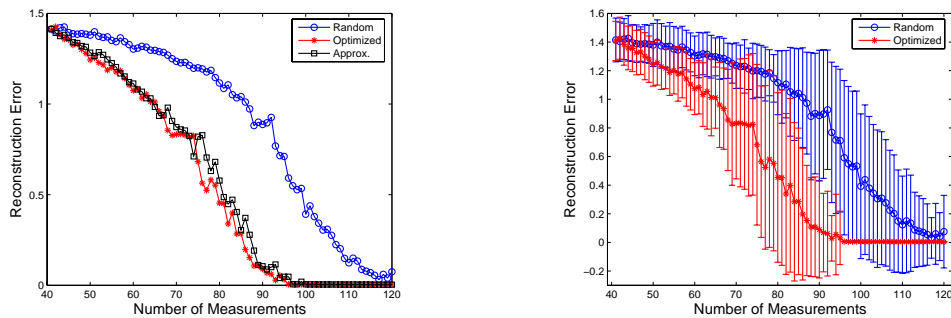


Figure 3. Comparison of adaptive and random projections, with the first 40 projections performed randomly. (a) Reconstruction error of BCS with random projections, optimized projections (Sec. 4.1) and approximated projections (Sec. 4.2); the results are averaged over 100 runs; (b) the variances of the reconstruction error of BCS with random projections and optimized projections (Sec. 4.1); the variance for the approximate scheme in Sec. 4.2 is very similar to that for Sec. 4.1, and thus is omitted to improve visibility.

In the first example we consider a length $N=512$ signal that contains $M=20$ spikes created by choosing 20 locations at random and then putting ± 1 at these points (Fig. 2(a)). The projection matrix Φ is constructed by first creating a $K \times N$ matrix with i.i.d. draws of a Gaussian distribution $\mathcal{N}(0, 1)$, and then the rows of Φ are normalized to unit amplitude. To simulate measurement noise, zero-mean Gaussian noise with standard deviation $\sigma_m = 0.005$ is added to each of the K measurements that define the data \mathbf{g} . In the experiment $K = 100$, and the reconstructions are implemented by BP and BCS. For the BP implementation, we used the ℓ_1 -magic package available online at <http://www.acm.caltech.edu/l1magic/>.

Figures 2(b-c) demonstrate the reconstruction results with BP and BCS, respectively. Because it is a noisy reconstruction problem, BP cannot recover the underlying sparse signal exactly. Consequently, the BCS reconstruction is much cleaner than BP, as $M=20$ spikes are correctly recovered with (about 10 times) smaller reconstruction error relative to BP. In addition, BCS yields “error-bars” for the estimated signal, indicating the confidence for the current estimation. Regarding the computation time, BCS also outperforms BP.

As discussed in Sec. 4, the Bayesian analysis also allows designing projection matrix Φ for adaptive CS. In the second experiment, we use the same dataset as in Fig. 2 and study the performance of BCS for projection design. The initial 40 measurements are conducted by using the random projections as in Fig. 2, except that the rows of Φ are normalized to 1.01 for the reasons discussed in Sec. 4.1. The remaining 80 measurements are sequentially conducted by optimized projections, with this compared to using random projections. In the experiment, after each projec-

tion vector \mathbf{r}_{K+1} is determined, the associated reconstruction error is also computed. For the optimized projection, \mathbf{r}_{K+1} is constructed by using the eigenvector of Σ that has the largest eigenvalue. When examining the approximate scheme discussed in Sec. 4.2, we used 10% of the average value of the diagonal elements of $(\alpha_0 \Phi^T \Phi)$ for diagonal loading. Because of the randomness in the experiment (i.e., the initial 40 random projections and the empty-entries imputation for \mathbf{r}_{K+1} , etc.), we execute the experiment 100 times with the average performance reported in Fig. 3.

It is demonstrated in Fig. 3 that the reconstruction error of the optimized projection is much smaller than that of the random projection, indicating the superior performance of this optimization. Further, the approximate scheme in Sec. 4.2 yields results very comparable to the more-rigorous analysis in Sec. 4.1. This suggests that existing CS software may be readily modified to implement the optimization procedure, and yield results comparable to that of the full BCS solution.

5.2. BCS vs. BP and StOMP

In the following set of experiments, the performance of BCS is compared to BP and StOMP (equipped with CFDR and CFAR thresholding) on two example problems included in the *Sparselab* package that is available online at <http://sparselab.stanford.edu/>. Following the experiment setting in the package, all the projection matrix Φ here are drawn from a uniform spherical distribution [Tsaig & Donoho, 2006].

5.2.1. RANDOM-BARS

Figure 4 shows the reconstruction results for *Random-Bars* that has been used in [Tsaig & Donoho, 2006].

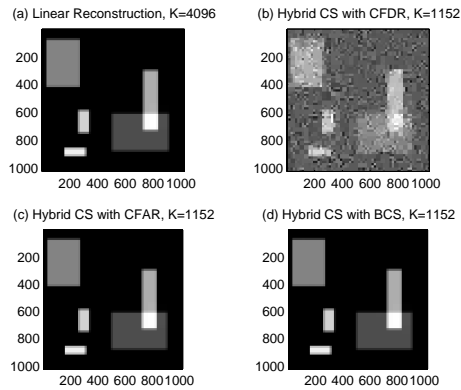


Figure 4. Reconstruction of *Random-Bars* with hybrid CS. (a) Linear reconstruction from $K=4096$ samples, $err_{LIN}=0.2271$; (b) Reconstruction with CFDR, $err_{CFDR}=0.5619$, $t_{CFDR}=4$ secs; (c) Reconstruction with CFAR, $err_{CFAR}=0.2271$, $t_{CFAR}=5$ secs; (d) Reconstruction with BCS, $err_{BCS}=0.2271$, $t_{BCS}=15$ secs. BP took 108 secs with the reconstruction error 0.2279, which is not shown here.

We used the Haar wavelet expansion, which is naturally suited to images of this type, with a coarsest scale $j_0 = 3$, and a finest scale $j_1 = 6$. Figure 4(a) shows the result of linear reconstruction with $K = 4096$ samples, which represents the best performance that could be achieved by all the CS implementations used, whereas Figs. 4(b-d) have results for the hybrid CS scheme [Tsaig & Donoho, 2006] with $K = 1152$ hybrid compressed samples. It is demonstrated that BCS and StOMP with CFAR yield the near optimal reconstruction error (0.2271); among all the CS algorithms considered StOMP is the fastest one. However, as we have noted, the performance of StOMP strongly relies on the thresholding parameters selected. For the *Random-Bars* problem considered, the performance of StOMP with CFDR is very sensitive to its parameter-setting, with one typical example result shown in Fig. 4(b).

5.2.2. MONDRIAN

Figure 5 displays a photograph of a painting by Piet Mondrian, the Dutch neo-plasticist. Despite being a simple geometric example, this image still presents a challenge, as its wavelet expansion is not as sparse as the examples considered above. We used a multiscale CS scheme [Tsaig & Donoho, 2006] for image reconstruction, with a coarsest scale $j_0 = 4$, and a finest scale $j_1 = 6$ on the “symmlet8” wavelet. Figure 5(a) shows the result of linear reconstruction with $K = 4096$ samples, which represents the best performance that could be achieved by all the CS implementations used, whereas Figs. 5(b-d) have results for the multiscale CS

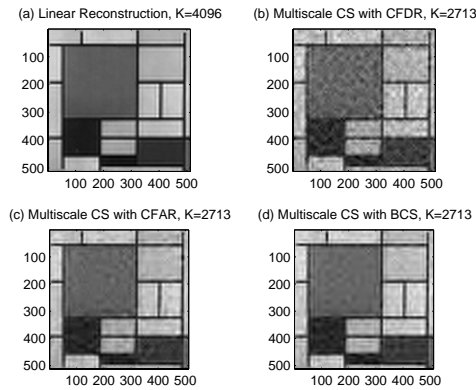


Figure 5. Reconstruction of *Mondrian* with multiscale CS. (a) Linear reconstruction from $K=4096$ samples, $err_{LIN}=0.1333$; (b) Reconstruction with CFDR, $err_{CFDR}=0.1826$, $t_{CFDR}=10$ secs; (c) Reconstruction with CFAR, $err_{CFAR}=0.1508$, $t_{CFAR}=28$ secs; (d) Reconstruction with BCS, $err_{BCS}=0.1503$, $t_{BCS}=18$ secs. BP took 162 secs with the reconstruction error 0.1416, which is not shown here.

scheme with $K=2713$ multiscale compressed samples. In the example results in Figs. 5(b-c), we used the same parameters-setting for StOMP as those used in the *Sparselab* package. It is demonstrated that all the CS implementations yielded a faithful reconstruction to the original image, while BCS produced the second smallest reconstruction error (0.1503) using the second smallest computation time (18 secs).

Table 1. Summary of the performances of BP, StOMP and BCS on *Mondrian*.

	BP	CFDR	CFAR	BCS
# Nonzeros	3840	1766	926	615
Time (secs)	162	10	28	18
Reconst. Error	0.1416	0.1826	0.1508	0.1503

To understand why BCS is more efficient than StOMP on this problem, we checked the number of nonzero weights recovered by BCS and StOMP, with the results reported in Table 1. Evidently, BCS found the sparsest solution (with 615 nonzeros) relative to the two StOMP implementations, but yielded the second smallest reconstruction error (0.1503). This indicates that although each iteration of StOMP allows multiple nonzero weights to be added into the “active set” [Donoho et al., 2006], this process may be a too generous usage of weights without reducing the reconstruction error. The sparser solution of BCS is the likely explanation of its relative higher speed compared to StOMP in this example.

6. Conclusions

Compressive sensing has been considered from a Bayesian perspective. It has been demonstrated that by utilizing the previously derived algorithms (relevance vector machine) from the sparse Bayesian literature, problems in CS can be solved more effectively. In practice we have found that the results from this Bayesian analysis are often sparser than existing CS solutions [Chen et al., 1999, Donoho et al., 2006]. On the examples considered from the literature, the BCS solution typically has computation time comparable to state-of-the-art algorithms such as StOMP [Donoho et al., 2006]; in some cases BCS is even faster as a consequence of the improved sparsity. We have also considered adaptive CS by optimizing the projection matrix Φ . Experiments on synthetic data demonstrate a significantly accelerated rate of convergence compared to the original CS construction. Finally, a simple approximate scheme has been introduced, which allows existing CS algorithms to benefit immediately from this Bayesian analysis.

There is a clear connection between CS and regression shrinkage and selection via the Lasso [Tibshirani, 1996, Efron et al., 2004] as both focus on solving the objective function (1). Research on Lasso has produced algorithms that might have some relevance to CS and BCS. In addition to this, other possible areas of future research may include simultaneous inversion of multiple data sets, borrowing ideas from multi-task learning [Caruana, 1997], and a theoretical analysis of adaptive CS, which can be an important complement to the existing analysis for the conventional CS formulation [Haupt & Nowak, 2006].

Acknowledgments

The authors wish to thank Ya Xue for the initial discussions and suggestions on compressive sensing, and the anonymous reviewers for the constructive suggestions. The authors also thank E. Candès, J. Romberg and D. Donoho *et al.* for sharing the ℓ_1 -Magic and SparseLab online. Their generous distribution of the code made the experimental comparisons in this paper very convenient.

References

Bishop, C. M., & Tipping, M. E. (2000). Variational relevance vector machines. *UAI 16* (pp. 46–53).

Candès, E., Romberg, J., & Tao, T. (2006). Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Information Theory*, *52*, 489–509.

Caruana, R. (1997). Multitask learning. *Machine Learning*, *28*, 41–75.

Chen, S., Donoho, D. L., & Saunders, M. A. (1999). Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, *20*, 33–61.

Cover, T. M., & Thomas, J. A. (1991). *Elements of information theory*. New York, NY: Wiley.

Donoho, D. L. (2006). Compressed sensing. *IEEE Trans. Information Theory*, *52*, 1289–1306.

Donoho, D. L., Tsai, Y., Drori, I., & Starck, J.-C. (2006). Sparse solution of underdetermined linear equations by stagewise orthogonal matching pursuit. Preprint.

Efron, B., Hastie, T., Johnstone, I., & Tibshirani, R. (2004). Least angle regression. *The Annals of Statistics*, *32*, 407–499.

Faul, A. C., & Tipping, M. E. (2002). Analysis of sparse Bayesian learning. *NIPS 14*.

Fedorov, V. V. (1972). *Theory of optimal experiments*. Academic Press.

Figueiredo, M. (2002). Adaptive sparseness using Jeffreys prior. *NIPS 14*.

Gilks, W. R., Richardson, S., & Spiegelhalter, D. J. (1996). *Markov Chain Monte Carlo in Practice*. Chapman & Hall.

Haupt, J., & Nowak, R. (2006). Signal reconstruction from noisy random projections. *IEEE Trans. Information Theory*, *52*, 4036–4048.

MacKay, D. (1992). Information-based objective functions for active data selection. *Neural Computation*, *4*, 590–604.

Mallat, S. (1998). *A wavelet tour of signal processing*. Academic Press. 2nd edition.

Papoulis, A., & Pillai, S. U. (2002). *Probability, random variables and stochastic processes*. McGraw-Hill. 4th edition.

Pearlman, W. A., Islam, A., Nagaraj, N., & Said, A. (2004). Efficient, low-complexity image coding with a set-partitioning embedded block coder. *IEEE Trans. Circuits Systems Video Technology*, *14*, 1219–1235.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *J. Royal. Statist. Soc B.*, *58*, 267–288.

Tipping, M. E. (2001). Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, *1*, 211–244.

Tipping, M. E., & Faul, A. C. (2003). Fast marginal likelihood maximisation for sparse Bayesian models. *Proc. of the 9th International Workshop on AISTats*.

Tropp, J. A., & Gilbert, A. C. (2005). Signal recovery from partial information via orthogonal matching pursuit. Preprint.

Tsai, Y., & Donoho, D. L. (2006). Extensions of compressed sensing. *Signal Processing*, *86*, 549–571.

Wipf, D., Palmer, J., & Rao, B. (2004). Perspectives on sparse Bayesian learning. *NIPS 16*.