

is a continuous and unimodal function of σ^2 , with the unique maximum achieved at $\sigma^2 = (\sigma^2)^{(p+1)}$, see also (9a). We conclude that $(\sigma^2)^{(p)} - (\sigma^2)^{(p+1)}$ must go to zero. The second claim of Theorem 1 follows.

REFERENCES

- [1] *IEEE Signal Process. Mag. (Special Issue on Compressive Sampling)*, vol. 25, no. 2, Mar. 2008.
- [2] E. J. Candès and T. Tao, "Decoding by linear programming," *IEEE Trans. Inf. Theory*, vol. 51, pp. 4203–4215, Dec. 2005.
- [3] E. J. Candès and T. Tao, "Near-optimal signal recovery from random projections: Universal encoding strategies?," *IEEE Trans. Inf. Theory*, vol. 52, no. 12, pp. 5406–5425, 2006.
- [4] D. L. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, pp. 1289–1306, Apr. 2006.
- [5] J. A. Tropp and S. J. Wright, "Computational methods for sparse solution of linear inverse problems," *Proc. IEEE*, vol. 98, no. 6, pp. 948–958, 2010.
- [6] P. Boufounos and R. Baraniuk, "1-bit compressive sensing," in *Proc. 42nd Annu. Conf. Inf. Sci. Syst.*, Princeton, NJ, Mar. 2008, pp. 16–21.
- [7] W. Dai, H. V. Pham, and O. Milenkovic, "Distortion-rate functions for quantized compressive sensing," in *IEEE Inf. Theory Workshop Netw. Inf. Theory*, Volos, Greece, Jun. 2009, pp. 171–175.
- [8] L. Jacques, D. Hammond, and M. Fadili, "Dequantizing compressed sensing: When oversampling and non-Gaussian constraints combine," *IEEE Trans. Inf. Theory*, vol. 57, no. 1, pp. 559–571, 2011.
- [9] A. Zymnis, S. Boyd, and E. Candès, "Compressed sensing with quantized measurements," *IEEE Signal Process. Lett.*, vol. 17, pp. 149–152, Feb. 2010.
- [10] K. Qiu and A. Dogandžić, "A GEM hard thresholding method for reconstructing sparse signals from quantized noisy measurements," in *Proc. 4th IEEE Int. Workshop Comput. Adv. Multi-Sensor Adapt. Process. (CAMSAP)*, San Juan, Puerto Rico, Dec. 2011, pp. 349–352.
- [11] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Englewood Cliffs, NJ: Prentice-Hall, 1993.
- [12] N. L. Johnson and S. Kotz, *Continuous Univariate Distributions*, 2nd ed. ed. New York: Wiley, 1994, vol. 1.
- [13] T. Blumensath and M. E. Davies, "Iterative hard thresholding for compressed sensing," *Appl. Comput. Harmon. Anal.*, vol. 27, no. 3, pp. 265–274, 2009.
- [14] K. Qiu and A. Dogandžić, "Double overrelaxation thresholding methods for sparse signal reconstruction," in *Proc. 44th Annu. Conf. Inf. Sci. Syst.*, Princeton, NJ, 2010, pp. 1–6.
- [15] K. Qiu and A. Dogandžić, "ECME thresholding methods for sparse signal reconstruction," Iowa State Univ., Ames, Tech. Rep., Apr. 2010 [Online]. Available: <http://arxiv.org/abs/1004.4880>
- [16] K. Qiu and A. Dogandžić, "Variance-component based sparse signal reconstruction and model selection," *IEEE Trans. Signal Process.*, vol. 58, pp. 2935–2952, Jun. 2010.
- [17] T. T. Do, T. D. Tran, and L. Gan, "Fast compressive sampling with structurally random matrices," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, Las Vegas, NV, Apr. 2008, pp. 3369–3372.
- [18] J. Starck, F. Murtagh, and J. Fadili, *Sparse Image and Signal Processing: Wavelets, Curvelets, Morphological Diversity*. New York: Cambridge Univ. Press, 2010.
- [19] D. A. Harville, *Matrix Algebra From a Statistician's Perspective*. New York: Springer-Verlag, 1997.
- [20] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. Roy. Statist. Soc. Ser. B*, vol. 39, no. 1, pp. 1–38, 1977, with discussion.

Exact Reconstruction Conditions for Regularized Modified Basis Pursuit

Wei Lu and Namrata Vaswani

Abstract—In this work, we obtain sufficient conditions for exact recovery of regularized modified basis pursuit (reg-mod-BP) and discuss when the obtained conditions are weaker than those for modified compressive sensing or for basis pursuit (BP). The discussion is also supported by simulation comparisons. Reg-mod-BP provides a solution to the sparse recovery problem when both an erroneous estimate of the signal's support, denoted by T , and an erroneous estimate of the signal values on T are available.

Index Terms—Compressive sensing, modified-CS, partially known support, sparse reconstruction.

I. INTRODUCTION

In this work, we obtain sufficient conditions for exact recovery of regularized modified basis pursuit (reg-mod-BP) and discuss when the obtained conditions are weaker than those for modified compressive sensing [2] or for basis pursuit (BP) [3], [4]. Reg-mod-BP was briefly introduced in our earlier work [2] as a solution to the sparse recovery problem when both an erroneous estimate of the signal's support, denoted by T , and an erroneous estimate of the signal values on T , denoted by $(\hat{\mu})_T$, are available. The problem is precisely defined in Section I-A. Reg-mod-BP, given in (11), tries to find a vector that is sparsest outside the set T among all solutions that are close enough to $(\hat{\mu})_T$ on T and satisfy the data constraint. In practical applications, T and $(\hat{\mu})_T$ may be available from prior knowledge, or in recursive reconstruction applications, e.g., recursive dynamic MRI [2], [5], recursive compressive sensing (CS) based video compression [6], [7], or recursive projected CS (ReProCS) [8], [9] based video layering, one can use the support and signal estimate from the previous time instant for this purpose.

Basis pursuit (BP) was introduced in [3] as a practical (polynomial complexity) solution to the problem of reconstructing a sparse $m \times 1$ vector, x , with support denoted by N , from an $n \times 1$ measurements' vector, $y := Ax$, when $n < m$. BP solves the following convex (actually linear) program:

$$\min_{\beta} \|\beta\|_1 \text{ subject to } y = A\beta. \quad (1)$$

The recent CS literature has provided strong exact recovery results for BP that are either based on the restricted isometry property (RIP) [4], [10] or that use the geometry of convex polytopes to obtain "exact recovery thresholds" on the n needed for exact recovery with high probability [11], [12]. BP is often just referred to as CS in recent works and our work also occasionally does this.

Manuscript received April 25, 2011; revised July 25, 2011 and October 20, 2011; accepted January 09, 2012. Date of publication February 03, 2012; date of current version April 13, 2012. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Jerome Idier. This work was supported in part by NSF Grant CCF-0917015. A part of this work was presented at the Forty-Fourth Asilomar Conference of Signals, Systems and Computing, 2010 [1].

The authors are with the Department of Electrical and Computer Engineering, Iowa State University, Ames IA 50010 USA (e-mail: luwei@iastate.edu; namrata@iastate.edu).

Digital Object Identifier 10.1109/TSP.2012.2186445

In recent work [2], we introduced the problem of sparse reconstruction with partial and partly erroneous support knowledge, denoted by T , and proposed a solution called modified compressive sensing (mod-CS). We obtained exact reconstruction conditions for mod-CS and showed when they are weaker than those for BP. Mod-CS tries to find the solution that is sparsest outside the set T among all solutions of $y = A\beta$, i.e., it solves

$$\min_{\beta} \|\beta_{T^c}\|_1 \text{ subject to } y = A\beta. \quad (2)$$

Ideally, the above should be referred to as mod-BP, but since we used the term mod-CS when we introduced it, we will retain it here. Similar problems were also studied in parallel work by von Borries *et al.* [13] and Khajehnejad *et al.* [14]. In [14], the authors assumed a probabilistic prior on the support, solved the following weighted ℓ_1 problem, and obtained exact recovery thresholds similar to those in [12]:

$$\min_{\beta} \|\beta_{T^c}\|_1 + \gamma \|\beta_T\|_1 \text{ subject to } y = A\beta. \quad (3)$$

In another related work [15], Wang *et al.* showed how to iteratively improve recovery of a *single* signal by solving BP in the first iteration, obtaining a support estimate, solving (2) with this support estimate and repeating this. They also obtained exact recovery guarantees for a single iteration.

Another related idea is *CS-diff* or *CS-residual*, which recovers the residual signal $x - \hat{\mu}$ by solving (1) with y replaced by $y - A\hat{\mu}$. This is related to our earlier least squares CS-residual (LS-CS) and Kalman filtered CS (KF-CS) ideas [5], [16]. However, as explained in [2], the residual signals using all these methods have a support size that is equal to or slightly larger than that of x (except if $(\hat{\mu})_T = x_T$). As a result, these do not achieve exact recovery with fewer measurements. The limitations of some other variants of this are also discussed in detail in [17]. Reg-mod-BP may also be interpreted as a Bayesian or a model-based CS approach. Recent work in this area include [18]–[20].

This paper is organized as follows. We introduce reg-mod-BP in Section II. In Section III, we obtain the exact reconstruction result, discuss its implications and give the key lemmas leading to its proof. Simulation comparisons are given in Section IV and conclusions in Section V.

Notation and Problem Definition

For a set T , $T^c = \{i \in [1, \dots, m], i \notin T\}$. \emptyset is the empty set. We use $|\cdot|$ to denote the cardinality of a set. The same notation is also used for the absolute value of a scalar. The meaning is clear from context.

For a vector b , $(b)_T$, or just b_T , denotes a subvector containing the elements of b with indices in T . $\|b\|_k$ means the ℓ_k norm of the vector b . The notation $b \succeq 0$ ($b \succ 0$) means that each element of the vector b is greater than or equal to (strictly greater than) zero. Similarly $b \preceq 0$ ($b \prec 0$) means each element is less than or equal to (strictly less than) zero. We define the sign pattern, $\text{sgn}(b)$ as

$$[\text{sgn}(b)]_i = \begin{cases} \frac{b_i}{|b_i|} & \text{if } b_i \neq 0 \\ 0 & \text{if } b_i = 0. \end{cases} \quad (4)$$

We use $'$ for matrix transpose. For a matrix A , A_T denotes the submatrix containing the columns of A with indices in T . Also, $\|A\| := \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2}$ is the induced 2 norm.

Our goal is to solve the sparse reconstruction problem, i.e., reconstruct an m -length sparse vector, x , with support, N , from an $n < m$ length measurement vector,

$$y := Ax \quad (5)$$

when an erroneous estimate of the signal's support, denoted by T ; and an erroneous estimate of the signal values on T , denoted by $(\hat{\mu})_T$, are available. The support estimate, T , can be rewritten as

$$T = N \cup \Delta_e \setminus \Delta, \text{ where } \Delta := N \setminus T \text{ and } \Delta_e := T \setminus N \quad (6)$$

are the errors (Δ contains the misses while Δ_e contains the extras) in the support estimate.

The signal value estimate is assumed to be zero along T^c , i.e.,

$$\hat{\mu} = \begin{bmatrix} (\hat{\mu})_T \\ \mathbf{0}_{T^c} \end{bmatrix}$$

and it satisfies

$$(\hat{\mu})_T = (x)_T + \nu, \text{ with } \|\nu\|_{\infty} \leq \rho. \quad (7)$$

The restricted isometry constant (RIC) [4], δ_s , for A , is defined as the smallest positive real number satisfying $(1 - \delta_s)\|c\|_2^2 \leq \|A_S c\|_2^2 \leq (1 + \delta_s)\|c\|_2^2$ for all subsets S of cardinality $|S| \leq s$ and all real vectors c of length $|S|$. The restricted orthogonality constant (ROC) [4], θ_{s_1, s_2} , is defined as the smallest positive real number satisfying $|c_1' A_{T_1}' A_{T_2} c_2| \leq \theta_{s_1, s_2} \|c_1\|_2 \|c_2\|_2$ for all disjoint sets T_1, T_2 with $|T_1| \leq s_1, |T_2| \leq s_2$ and $s_1 + s_2 \leq m$, and for all vectors c_1, c_2 of length $|T_1|, |T_2|$ respectively. Both δ_s and θ_{s_1, s_2} are nondecreasing functions of s and of s_1, s_2 , respectively [4].

We will frequently use the following functions of the RIC and ROC of A in Section III:

$$a_k(s, \tilde{s}) := \frac{\theta_{s, s} + \frac{\theta_{s, k} \theta_{s, k}}{1 - \delta_k}}{1 - \delta_s - \frac{\theta_{s, k}^2}{1 - \delta_k}} \quad (8)$$

$$K_k(u) := \frac{\sqrt{1 + \delta_u}}{1 - \delta_u - \frac{\theta_{u, k}^2}{1 - \delta_k}}. \quad (9)$$

For the matrix A , and for any set S for which $A_S' A_S$ is full rank, we define the matrix $M(S)$ as

$$M(S) := I - A_S (A_S' A_S)^{-1} A_S'. \quad (10)$$

II. REGULARIZED MODIFIED BASIS PURSUIT

Mod-CS given in (2) puts no cost on β_T and no explicit constraint except $y = A\beta$. Thus, when very few measurements are available, β_T can become larger than required in order to satisfy $y = A\beta$ with the smallest $\|\beta_{T^c}\|_1$. A similar, though less, bias will also occur with (3) when $\gamma < 1$. However, if a signal value estimate on T , $(\hat{\mu})_T$, is also available, one can use that to constrain β_T . One way to do this, as suggested in [2], is to add $\lambda \|\beta_T - \hat{\mu}_T\|_2^2$ to the mod-CS cost. However, as we saw from simulations, while this does achieve lower reconstruction error, it cannot achieve exact recovery with fewer measurements (smaller n) than mod-CS [2]. The reason is it puts a cost on the entire ℓ_2 distance from $(\hat{\mu})_T$ and so encourages elements on the extras set, Δ_e , to be closer to $(\hat{\mu})_{\Delta_e}$ which is nonzero.

On the other hand, if we instead use the ℓ_{∞} distance from $(\hat{\mu})_T$, and add it as a constraint, then, at least in certain situations, we can achieve exact recovery with a smaller n than mod-CS. Thus, we study

$$\min_{\beta} \|\beta_{T^c}\|_1, \text{ subject to } y = A\beta \text{ and } \|\beta_T - \hat{\mu}_T\|_{\infty} \leq \rho \quad (11)$$

and call it *reg-mod-BP*. We see from simulations, that whenever one or more of the inequality constraints are active at x , i.e., $|x_i - \hat{\mu}_i| = \rho$

for some $i \in T$, (11) does achieve exact recovery with fewer measurements than mod-CS. We use this observation to derive a better exact recovery result below.¹

III. EXACT RECONSTRUCTION CONDITIONS

In this section, we obtain exact reconstruction conditions for reg-mod-BP by exploiting the above fact. We give the result and discuss its implications below in Section III-A. The key lemmas leading to its proof are given in Section III-B and the proof outline in Section III-C.

A. Exact Reconstruction Result

Let us begin by defining the two types of active sets (set of indices for which the inequality constraint is active), T_{a+} and T_{a-} , and the inactive set, T_{in} , as follows:

$$\begin{aligned} T_{a+} &:= \{i \in T : x_i - \hat{\mu}_i = \rho\} \\ T_{a-} &:= \{i \in T : x_i - \hat{\mu}_i = -\rho\} \\ T_{in} &:= \{i \in T : |x_i - \hat{\mu}_i| < \rho\}. \end{aligned} \quad (12)$$

In the result below, we try to find the sets $T_{a+g} \subseteq T_{a+}$ and $T_{a-g} \subseteq T_{a-}$ so that $|T_{a+g}| + |T_{a-g}|$ is maximized while T_{a+g} and T_{a-g} satisfy certain constraints. We call these the “good” sets. We define the “bad” subset of T , as $T_b := T \setminus (T_{a+g} \cup T_{a-g})$. As we will see, the smaller the size of this bad set, the weaker are our exact recovery conditions.

Theorem 1 (Exact Recovery Conditions): Consider recovering a sparse vector, x , with support N , from $y := Ax$ by solving (11). The support estimate, T , and the misses and extras in it, Δ , Δ_e , satisfy (6). The signal estimate, $\hat{\mu}$, satisfies (7), i.e., $\|x_T - \hat{\mu}_T\|_\infty \leq \rho$. Define the sizes of the sets T and Δ as

$$k := |T|, \quad u := |\Delta|. \quad (13)$$

The true x is the unique minimizer of (11) if

- 1) $\delta_{k+u} < 1$, $\delta_{2u} + \delta_k + \theta_{k,2u}^2 < 1$, and
- 2) $a_k(2u, u) + a_{k_b}(u, u) < 1$ where

$$\begin{aligned} T_b &:= T \setminus (T_{a+g} \cup T_{a-g}), \text{ and} \\ k_b &:= |T_b| \\ \{T_{a+g}, T_{a-g}\} &= \arg \max_{\tilde{T}_{a+g}, \tilde{T}_{a-g}} (|\tilde{T}_{a+g}| + |\tilde{T}_{a-g}|) \text{ subject to} \end{aligned}$$

$$\begin{aligned} \tilde{T}_{a+g} &\subseteq T_{a+}, \quad \tilde{T}_{a-g} \subseteq T_{a-}, \\ A_i' w &> 0 \quad \forall i \in \tilde{T}_{a+g}, \text{ and} \\ A_i' w &< 0 \quad \forall i \in \tilde{T}_{a-g}, \end{aligned}$$

where

$$\begin{aligned} w &:= M(\tilde{T}_b) A_\Delta (A_\Delta' M(\tilde{T}_b) A_\Delta)^{-1} \text{sgn}(x_\Delta) \\ \tilde{T}_b &:= T \setminus (\tilde{T}_{a+g} \cup \tilde{T}_{a-g}), \end{aligned} \quad (14)$$

$M(S)$ is specified in (10), $a_k(s, \bar{s})$ is defined in (8), and the sets T_{a+} , T_{a-} are defined in (12). ■

Notice that $a_k(s, \bar{s})$ is a nondecreasing function of k . Since $k_b = k - |T_{a+g}| - |T_{a-g}|$, thus, finding the largest possible sets T_{a+g} and T_{a-g} ensures that the condition $a_k(2u, u) + a_{k_b}(u, u) < 1$ is the

¹One can also try to constrain the ℓ_2 distance instead of the ℓ_∞ distance. When the ℓ_2 constraint is active, one should again need a smaller \mathbf{n} for exact recovery. When we check this via simulations, this does happen, but since it is at most one active constraint, the reduction in \mathbf{n} required is small compared to what is achieved by (11) and hence we do not study this further.

weakest. The reason for defining T_{a+g} and T_{a-g} in the above fashion will become clear in the proof of Lemma 2.

Notice also that the first condition of the above result ensures that $\delta_k < 1$. Since $|\tilde{T}_b| \leq k$, thus, $A_{\tilde{T}_b}' A_{\tilde{T}_b}$ is positive definite and thus invertible. Thus $M(\tilde{T}_b)$ is always well defined. The first condition also ensures that $a_k(2u, u) > 0$. Since $k_b \leq k$, and since δ_s and θ_{s_1, s_2} are nondecreasing functions of s , s_1 , s_2 , it also ensures that $a_{k_b}(u, u) > 0$.

Remark 1 (Applicability): A practical case where some of the inequality constraints will be active with nonzero probability is when dealing with quantized signals and quantized signal estimates. If the range of values that the signal estimate can take given the signal (or vice versa) is known, the smallest choice of ρ is easily computed. We show some examples in Section IV. In general, even if just the range of values both can take is known, we can compute ρ . The fewer the number values that $x_i - \hat{\mu}_i$ can take, the larger will be the expected size of the active set, $T_a := T_{a+} \cup T_{a-}$. Also, the condition (14) will hold for nonempty $T_g := T_{a+g} \cup T_{a-g}$ with nonzero probability. Some real applications where quantized signals and signal estimates occur are recursive CS based video compression [6], [7] (the original video itself is quantized) or in recursive projected CS (ReProCS) [8], [9] based moving or deforming foreground objects' extraction (e.g., a person moving towards a camera) from very large but correlated noise (e.g., very similar looking but slowly changing backgrounds), particularly when the videos are coarsely quantized (low bit rate). A common example where low bit rate videos occur is mobile telephony applications. In any of these applications, if we know a bound on the maximum change of the sparse signal's value from one time instant to the next, that can serve as ρ .

Remark 2 (Comparison With BP, Mod-CS, Other Results): The worst case for Theorem 1 is when both the sets T_{a+g} and T_{a-g} are empty either because no constraint is active (T_{a+} and T_{a-} are both empty) or because (14) does not hold for any pair of subsets of T_{a+} and T_{a-} . In this case, we have $k_b = k$ and so the required sufficient conditions are the same as those of mod-CS [2, Theorem 1]. A small extra requirement is that x satisfies (7). Thus, in the worst case, Theorem 1 holds under the same conditions on A (needs the same number of measurements) as mod-CS [2]. In [2], we have already argued that the mod-CS result holds under weaker conditions than the results for BP [4], [10] as long as the size of the support errors, $|\Delta|$, $|\Delta_e|$, are small compared to the support size, $|N|$, and hence the same can be said about Theorem 1. For example, we argued that when $|\Delta| = |\Delta_e| = 0.02|N|$ (numbers taken from a recursive dynamic MRI application), the mod-CS conditions are weaker than those of BP. Small $|\Delta|$, $|\Delta_e|$ is a valid assumption in recursive recovery applications like recursive dynamic MRI, recursive CS based video compression, or ReProCS based foreground extraction from large but correlated background noise.

Moreover, if some inequality constraints are active and (14) holds, as in case of quantized signals and signal estimates, Theorem 1 holds under weaker conditions on A than the mod-CS result.

As noted by an anonymous reviewer, our exact recovery conditions require knowledge of x . However this is an issue with many results in sparse recovery, e.g., [21], and especially those that use more prior knowledge, e.g., [18].

Remark 3 (Small Reconstruction Error): The reconstruction error of reg-mod-BP is significantly smaller than that of mod-CS, weighted ℓ_1 or BP, even when none of the constraints is active, as long as ρ is small (see Table III). On the other hand, the exact recovery conditions do not depend on the value of ρ , but only on the size of the good subsets of the active sets. This is also observed in our simulations. In Table III, we show results for $\rho = 0.1$. Even when we tried $\rho = 0.5$, the exact reconstruction probability or the smallest n needed for exact reconstruction remained the same, but the reconstruction error increased.

Remark 4 (Computation Complexity): Finding the best T_{a+g} and T_{a-g} requires that one check all possible subsets of T_{a+} and T_{a-} and find the pair with the largest sum of sizes that satisfies (14). To do this, one would start with $\tilde{T}_{a+g} = T_{a+}$, $\tilde{T}_{a-g} = T_{a-}$; compute \tilde{T}_b and w and check if (14) holds; if it does not, remove one element from \tilde{T}_{a+g} and then check (14); then remove an element from \tilde{T}_{a-g} and check (14); keep doing this until one finds a pair for which (14) holds. In the worst case, one will need to check (14) $2^{|T_{a+}|+|T_{a-}|}$ times. However, the complexity of computing the RIC $\delta_{|T|}$ or any of the ROC's is anyway exponential in $|T|$ and $|T| \geq |T_{a+}| + |T_{a-}|$. In summary, computing the conditions of Theorem 1 has complexity that is exponential in the support size, but the same is true for all sparse recovery results that use the RIC. We should mention though that, for certain random matrices, e.g., random Gaussian, there are results that upper bound the RIC values with high probability, e.g., see [4]. However, the resulting bounds are usually quite loose.

B. Proof of Theorem 1: Key Lemmas

Our overall proof strategy is similar to that of [4] for BP and of [2] for mod-CS. We first find a set of sufficient conditions on an $n \times 1$ vector, w , that help ensure that x is the unique minimizer of (11). This is done in Lemma 1. Next, we find sufficient conditions that the measurement matrix A should satisfy so that one such w can be found. This is done in an iterative fashion in the theorem's proof. The proof uses Lemma 2 at the zeroth iteration, followed by applications of Lemma 3 at later iterations.

To obtain the sufficient conditions on w , as suggested in [4], we first write out the Karush–Kuhn–Tucker (KKT) conditions for x to be a minimizer of (11)[22, Ch. 5]. By strengthening these a little, we get a set of *sufficient* conditions for w to be the *unique* minimizer. The necessary conditions for x to be a minimizer are: there exists an $n \times 1$, vector w (Lagrange multiplier for the constraints in $y = Ax$), a $|T_{a+}| \times 1$ vector, λ_1 , and a $|T_{a-}| \times 1$ vector, λ_2 , such that (s.t.)

- 1) every element of λ_1 and λ_2 is nonnegative, i.e., $\lambda_1 \succeq 0$ and $\lambda_2 \succeq 0$;
- 2) $A_{T_{in}}'w = 0$, $A_{T_{a+}}'w = \lambda_1$, $A_{T_{a-}}'w = -\lambda_2$, $A_{\Delta}'w = \text{sgn}(x_{\Delta})$, and $\|A_{(T \cup \Delta)}'w\|_{\infty} \leq 1$.

As we will see in the proof of Lemma 1, strengthening $\|A_{(T \cup \Delta)}'w\|_{\infty} \leq 1$ to $\|A_{(T \cup \Delta)}'w\|_{\infty} < 1$, keeping the other conditions the same, and requiring that $\delta_{k+u} < 1$ gives us a set of *sufficient* conditions.

Lemma 1: Let x be as defined in Theorem 1. x is the unique minimizer of (11) if $\delta_{k+u} < 1$ and if we can find an $n \times 1$ vector, w , s.t.

- 1) $A_{T_{in}}'w = 0$, $A_{T_{a+}}'w \succeq 0$, $A_{T_{a-}}'w \preceq 0$;
- 2) $A_{\Delta}'w = \text{sgn}(x_{\Delta})$;
- 3) $|A_j'w| < 1$ for all $j \notin T \cup \Delta$.

Recall that T_{a+} , T_{a-} and T_{in} are defined in (12) and k, u in Theorem 1. ■

Proof: The proof is given in Appendix A.

Notice that the first condition is weaker than that of Lemma 1 of mod-CS [2] (which requires $A_T'w = 0$), while the other two are the same. Next, we try to obtain sufficient conditions on the measurement matrix, A (on its RIC's and ROC's) to ensure that such a w can be found. This is done by using Lemmas 2 and 3 given below. Lemma 2 helps ensure that the first two conditions of Lemma 1 hold and provides the starting point for ensuring that the third condition also holds. Then, Lemma 3 applied iteratively helps ensure that the third condition also holds.

Lemma 2: Assume that $k+u \leq m$. Let \tilde{s} be such that $k+u+\tilde{s} \leq m$. If $\delta_u + \delta_{k_b} + \theta_{k_b, u}^2 < 1$, then there exists an $n \times 1$ vector \tilde{w} and an “exceptional” set, E , disjoint with $T \cup \Delta$, s.t.

- 1) $A_{T_b}'\tilde{w} = 0$, $A_{T_{a+g}}'\tilde{w} > 0$, $A_{T_{a-g}}'\tilde{w} < 0$;
- 2) $A_{\Delta}'\tilde{w} = \text{sgn}(x_{\Delta})$;
- 3) $|E| < \tilde{s}$, $\|A_E'\tilde{w}\|_2 \leq a_{k_b}(u, \tilde{s})\sqrt{u}$, $|A_j'\tilde{w}| \leq \frac{a_{k_b}(u, \tilde{s})}{\sqrt{\tilde{s}}}\sqrt{u} \forall j \notin T \cup \Delta \cup E$;

$$4) \|\tilde{w}\|_2 \leq K_{k_b}(u)\sqrt{u}.$$

Recall that $a_k(s, \tilde{s})$, $K_k(s)$ are defined in (8), (9) and T_{a+g} , T_{a-g} , T_b , k_b , k and u in Theorem 1. ■

Notice that because we have assumed that $\delta_u + \delta_{k_b} + \theta_{k_b, u}^2 < 1$, $a_{k_b}(u, \tilde{s})$ and $K_{k_b}(u)$ are positive. We call the set E an “exceptional” set, because except on the set $E \subseteq (T \cup \Delta)^c$, everywhere else on $(T \cup \Delta)^c$, $|A_j'\tilde{w}|$ is bounded. This notion is taken from [4]. Notice that the first two conditions of the above lemma are one way to satisfy the first two conditions of Lemma 1 since $T_b = T_{in} \cup (T_{a+} \setminus T_{a+g}) \cup (T_{a-} \setminus T_{a-g})$.

Proof: The proof is given in Appendix B. We let $\tilde{w} = M(T_b)A_{\Delta}(A_{\Delta}'M(T_b)A_{\Delta})^{-1}\text{sgn}(x_{\Delta})$. Since the good sets T_{a+g} , T_{a-g} are appropriately defined [see (14)], the first two conditions hold. The rest of the proof bounds $\|\tilde{w}\|_2$, and finds the set $E \subseteq (T \cup \Delta)^c$ of size $|E| < \tilde{s}$ so that $|A_j'\tilde{w}|$ is bounded for all $i \notin T \cup \Delta \cup E$ and also $\|A_E'\tilde{w}\|_2$ is bounded.

Lemma 3 [2, Lemma 2]: Assume that $k \leq m$. Let s, \tilde{s} be such that $k+s+\tilde{s} \leq m$. Assume that $\delta_s + \delta_k + \theta_{k, s}^2 < 1$. Let T_d be a set that is disjoint with T , of size $|T_d| \leq s$ and let c be a $|T_d| \times 1$ vector. Then there exists an $n \times 1$ vector, \tilde{w} , and a set, E , disjoint with $T \cup T_d$, s.t. i) $A_T'\tilde{w} = 0$, ii) $A_{T_d}'\tilde{w} = c$, iii) $|E| < \tilde{s}$, $\|A_E'\tilde{w}\|_2 \leq a_k(s, \tilde{s})\|c\|_2$, $|A_j'\tilde{w}| \leq \frac{a_k(s, \tilde{s})}{\sqrt{\tilde{s}}}\|c\|_2$, $\forall j \notin T \cup T_d \cup E$, and iv) $\|\tilde{w}\|_2 \leq K_k(s)\|c\|_2$.

Recall that $a_k(s, \tilde{s})$, $K_k(s)$ are defined in (8), (9), and k, u in Theorem 1. ■

Proof: The proof of Lemma 3 is given in [2] and also in [23, App. C].

Notice that because we have assumed that $\delta_s + \delta_k + \theta_{k, s}^2 < 1$, $a_k(s, \tilde{s})$ and $K_k(s)$ are positive.

C. Proof Outline of Theorem 1

The proof is very similar to that of [2]. Hence we give only the outline here. The complete proof is in [23]. At iteration zero, we apply Lemma 2 with $\tilde{s} \equiv u$, to get a w_1 and an exceptional set $T_{d,1}$, disjoint with $T \cup \Delta$, of size less than u . Lemma 2 can be applied because $k_b \leq k$ and condition 1 of the theorem holds. At iteration $r > 0$, we apply Lemma 3 with $T_d \equiv \Delta \cup T_{d,r}$ (so that $s \equiv 2u$), $c_{\Delta} \equiv 0$, $c_{T_d} \equiv A_{T_d}'w_r$ and $\tilde{s} \equiv u$ to get a w_{r+1} and an exceptional set $T_{d,r+1}$ disjoint with $T \cup \Delta \cup T_{d,r}$ of size less than u . Lemma 3 can be applied because condition 1 of the theorem holds. Define $w := \sum_{r=1}^{\infty} (-1)^{r-1} w_r$. We then argue that if condition 2 of the theorem holds, w is well-defined and satisfies the conditions of Lemma 1. Applying Lemma 1, the result follows.

IV. NUMERICAL EXPERIMENTS

In this section, we show two types of numerical experiments. The first simulates quantized signals and signal estimates. This is the case where some constraints are active with nonzero probability. The good set, $T_g = T_{a+g} \cup T_{a-g}$ is also non empty with nonzero probability. Hence, for a given small enough n , reg-mod-BP has significantly higher exact reconstruction probability, $p_{\text{exact}}(n)$, as compared to both mod-CS [2] and weighted ℓ_1 [14] and much higher than that of BP [3], [4]. Alternatively, it also requires a significantly reduced n for exact reconstruction with probability one, $n_{\text{exact}}(1)$. In computing $p_{\text{exact}}(n)$ we average over the distribution of x, T and $\hat{\mu}$, as also in [2], [4]. All numbers are computed based on 100 Monte Carlo simulations. To compute $n_{\text{exact}}(1)$, we tried various values of n for each algorithm and computed the smallest n required for exact recovery always (in all 100 simulations).

We also do a second simulation where signal estimates are not quantized.

In the following steps, the notation

$$z \sim \text{discrete-uniform}(a_1, a_2, \dots, a_n)$$

TABLE I

QUANTIZED SIGNALS AND SIGNAL ESTIMATES. RECALL THAT $k = |T| = 26$. FOR $2K = 4$, THE EXPECTED SIZES OF T_a , T_g AND T_b ARE $\mathbb{E}[|T_a|] = 10.01$, $\mathbb{E}[|T_g|] = 5.27$ AND $\mathbb{E}[|T_b|] = 20.73$. FOR $2K = 10$, $\mathbb{E}[|T_a|] = 4.28$, $\mathbb{E}[|T_g|] = 2.3$ AND $\mathbb{E}[|T_b|] = 23.7$

	$2K$	BP	mod-CS	weighted ℓ_1	Reg-mod-BP
$p_{\text{exact}}(0.15m)$	4	0	0.18	0.16	0.64
N-RMSE(0.15m)	4	1.011	0.059	0.060	0.029
$n_{\text{exact}}(1)$	4	0.39m	0.21m	0.21m	0.18m
$p_{\text{exact}}(0.15m)$	10	0	0.18	0.16	0.39
N-RMSE(0.15m)	10	1.011	0.059	0.060	0.032
$n_{\text{exact}}(1)$	10	0.4m	0.21m	0.21m	0.20m

means that z is equally likely to be equal to a_1, a_2, \dots or a_n . We use $\pm a$ as short for $+a, -a$. Also, $z \sim \text{uniform}(a, b)$ generates a scalar uniform random variable in the range $[a, b]$. The notation $x_i \stackrel{\text{iid}}{\sim} P$ for all $i \in S$ means that, for all $i \in S$, each x_i is identically distributed according to P and is independent of all the others.

For the quantized case, x was an $m = 256$ length sparse vector with support size $|N| = 0.1m = 26$ and support estimate error sizes $u = |\Delta| = |\Delta_c| = 0.1|N| = 3$. We generated the matrix A once as an $n \times m$ random Gaussian matrix (generate an $n \times m$ matrix with i.i.d zero mean Gaussian entries and normalize each column to unit ℓ_2 norm). The following steps were repeated 100 times.

- 1) The support set, N , of size $|N|$, was generated uniformly at random from $[1, m]$. The support misses set, Δ , of size u , was generated uniformly at random from the elements of N . The support extras set, Δ_c , also of size u , was generated uniformly at random from the elements of N^c . The support estimate, $T = N \cup \Delta_c \setminus \Delta$ and thus $|T| = |N| = 26$.
- 2) We generated $x_i \stackrel{\text{iid}}{\sim} \text{discrete-uniform}(\pm 1)$ for $i \in N \cap T$; $x_i \stackrel{\text{iid}}{\sim} \text{discrete-uniform}(\pm 0.1)$ for $i \in \Delta$, and $x_i = 0$ for $i \in N^c$. $x_{N \cap T}$ and x_Δ are also independent of each other. We generated $\hat{\mu}_T = x_T + \nu$ where $\nu_i \stackrel{\text{iid}}{\sim} \text{discrete-uniform}(0, \pm \frac{\rho}{K}, \pm 2\frac{\rho}{K}, \dots, \pm \rho)$ for $i \in T \cap N$ and $\nu_i \stackrel{\text{iid}}{\sim} \text{discrete-uniform}(\pm \frac{\rho}{K}, \pm 2\frac{\rho}{K}, \dots, \pm \rho)$ for $i \in \Delta_c$. We used $\rho = 0.1$ and tried two choices of K . Notice that, for a given K , the number of equally likely values that $x_i - \hat{\mu}_i$ for $i \in T$ can take are roughly $2K + 1$ ($2K$ when $i \in \Delta_c$). The constraint is active when $x_i - \hat{\mu}_i$ is equal to $\pm \rho$. Thus, the expected size of the active set is roughly $\frac{2}{2K+1}|T|$.
- 3) We generated $y = Ax$. We solved reg-mod-BP given in (11) with $\rho = 0.1$; BP given in (1); mod-CS given in (2); and weighted ℓ_1 given in (3) with various choices of γ : [0.1 0.05 0.01 0.001]. We used the CVX optimization package, <http://www.stanford.edu/boyd/cvx/>, which uses primal-dual interior point method for solving the minimization problem.

We computed $p_{\text{exact}}(n)$ as the number of times \hat{x} was equal to x ("equal" was defined as $\frac{\|\hat{x} - x\|_2}{\|x\|_2} < 10^{-5}$) divided by 100. For weighted ℓ_1 , we computed $p_{\text{exact}}(n)$ for each choice of γ and recorded the largest one. This corresponded to $\gamma = 0.1$. We tabulate results in Table I. In the first row, we record $p_{\text{exact}}(0.15m)$ for all the methods, when using $K = 2$. We also record the Monte Carlo average of the sizes of the active set $|T_a| = |T_{a+} \cup T_{a-}|$; of the good set, $|T_g| = |T_{a+g} \cup T_{a-g}|$ and of the bad set $|T_b| = k - |T_g|$. In the second row, we record the normalized root mean-square error (N-RMSE). In the third row, we record $n_{\text{exact}}(1)$. In the next three rows, we repeat the same things with $K = 5$.

As can be seen, $|T_g|$ is about half the size of the active set, $|T_a|$. As K is increased, $|T_a|$ and hence $|T_g|$ reduces ($|T_b|$ increases) and thus $p_{\text{exact}}(0.15m)$ decreases and $n_{\text{exact}}(1)$ increases. Also, for mod-CS and weighted ℓ_1 , $p_{\text{exact}}(0.15m)$ is significantly smaller than for reg-mod-BP, while $n_{\text{exact}}(1)$ is larger.

TABLE II

QUANTIZED SIGNALS AND SIGNAL ESTIMATES: CASE 2. RECALL THAT $k = |T| = 26$. THE EXPECTED SIZES OF T_a , T_g AND T_b ARE $\mathbb{E}[|T_a|] = 9.02$, $\mathbb{E}[|T_g|] = 4.58$ AND $\mathbb{E}[|T_b|] = 21.42$

	BP	mod-CS	weighted ℓ_1	Reg-mod-BP
$p_{\text{exact}}(0.15m)$	0	0.26	0.26	0.57
N-RMSE(0.15m)	0.967	0.152	0.152	0.082
$n_{\text{exact}}(1)$	0.4m	0.21m	0.21m	0.20m

TABLE III
THE NON-QUANTIZED CASE

	BP	mod-CS	weighted ℓ_1	Reg-mod-BP
$p_{\text{exact}}(0.18m)$	0	0.87	0.87	0.87
N-RMSE(0.18m)	0.961	0.0175	0.0177	0.0123
N-RMSE(0.11m)	1.05	0.179	0.175	0.0635
$n_{\text{exact}}(1)$	0.39m	0.21m	0.21m	0.21m

Next, we simulated a more realistic scenario—the case of 3-bit quantized images (both x and $\hat{\mu}$ take integer values between 0 to 7). Here again $m = 256$, $|N| = 0.1m = 26$, and $u = |\Delta| = |\Delta_c| = 0.1|N| = 3$. The sets N , Δ , Δ_c and T were generated as before. We generated $x_i \stackrel{\text{iid}}{\sim} \text{discrete-uniform}(3, 4, \dots, 7)$ for $i \in N \cap T$; $x_i \sim \text{discrete-uniform}(1, 2)$ for $i \in \Delta$; and $x_i = 0$ for $i \in N^c$. Also, $\hat{\mu}_T = \text{clip}(x_T + \nu)$ where $\nu_i \sim \text{discrete-uniform}(-2, -1, 0, 1, 2)$ for $i \in T \cap N$; and $\nu_i \sim \text{discrete-uniform}(-2, -1, 1, 2)$ for $i \in \Delta_c$. Also $\text{clip}(z)$ clips any value more than 7 to 7 and any value less than zero to zero. Clearly, in this case $\rho = 2$. We record our results in Table II. Similar conclusions as before can be drawn.

Finally, we simulated the nonquantized case. We used $m = 256$, $|N| = 0.1m = 26$, and $u = |\Delta| = |\Delta_c| = 0.1|N| = 3$. We generated $x_i \stackrel{\text{iid}}{\sim} \text{discrete-uniform}(\pm 1)$ for $i \in N \cap T$; $x_i \stackrel{\text{iid}}{\sim} \text{discrete-uniform}(\pm 0.1)$ for $i \in \Delta$, and $x_i = 0$ for $i \in N^c$. The signal estimate, $\hat{\mu}_T = x_T + \nu$ where $\nu_i \stackrel{\text{iid}}{\sim} \text{uniform}(-\rho, \rho)$ with $\rho = 0.1$. We tabulate our results in Table III.

Since ν is a real vector (not quantized), the probability of any constraint being active is zero. Thus, as expected, p_{exact} and n_{exact} are the same for reg-mod-BP and mod-CS and weighted ℓ_1 , though significantly better than BP. However, the N-RMSE for reg-mod-BP is significantly lower than that for mod-CS and weighted ℓ_1 also, particularly when $n = 0.11m$.

V. CONCLUSION

In this work, we obtained sufficient exact recovery conditions for reg-mod-BP, (11), and discussed their implications. Our main conclusion is that if some of the inequality constraints are active and if even a subset of the set of active constraints satisfies certain conditions (given in (14)), then reg-mod-BP achieves exact recovery under weaker conditions than what mod-CS needs. A practical situation where this would happen is when both the signal and its estimate are quantized. In other cases, the conditions are only as weak as those for mod-CS. In either case they are much weaker than those for BP as long as T is a good support estimate. From simulations, we see that even without any active constraints, the reg-mod-BP reconstruction error is much lower than that of mod-CS or weighted ℓ_1 .

APPENDIX

A. Proof of Lemma 1

Denote a minimizer of (11) by β . Since $y = Ax$ and x satisfies (7), x is feasible for (11). Thus,

$$\|\beta_{T^c}\|_1 \leq \|x_{T^c}\|_1 = \|x_\Delta\|_1. \quad (15)$$

Next, we use the conditions on w given in Lemma 1 and the fact that x is supported on $N \subseteq T \cup \Delta$ to show that $\|\beta_{T^c}\|_1 \geq \|x_{T^c}\|_1$ and hence $\|x_{T^c}\|_1 = \|\beta_{T^c}\|_1$. Notice that

$$\begin{aligned} \|\beta_{T^c}\|_1 &= \sum_{j \in \Delta} |x_j + \beta_j - x_j| + \sum_{j \notin T \cup \Delta} |\beta_j| \\ &\geq \sum_{j \in \Delta} |x_j + \beta_j - x_j| + \sum_{j \notin T \cup \Delta} w' A_j \beta_j \end{aligned} \quad (16)$$

$$\begin{aligned} &\geq \sum_{j \in \Delta} \text{sgn}(x_j)(x_j + (\beta_j - x_j)) \\ &\quad + \sum_{j \notin T \cup \Delta} w' A_j (\beta_j - x_j) \end{aligned} \quad (17)$$

$$\begin{aligned} &= \|x_\Delta\|_1 + \sum_{j \notin T} w' A_j (\beta_j - x_j) \\ &= \|x_\Delta\|_1 + w'(A\beta - Ax) \\ &\quad - \sum_{j \in T} w' A_j (\beta_j - x_j) \end{aligned} \quad (18)$$

$$= \|x_\Delta\|_1 - \sum_{j \in T} w' A_j (\beta_j - \hat{\mu}_j + \hat{\mu}_j - x_j) \quad (19)$$

$$\begin{aligned} &= \|x_\Delta\|_1 - \sum_{j \in T_{a+}} w' A_j (\beta_j - \hat{\mu}_j - \rho) \\ &\quad - \sum_{j \in T_{a-}} w' A_j (\beta_j - \hat{\mu}_j + \rho) \end{aligned} \quad (20)$$

$$\geq \|x_\Delta\|_1 = \|x_{T^c}\|_1. \quad (21)$$

In the above, the inequality in (16) follows because $w' A_j \leq |w' A_j| < 1$ for $j \notin T \cup \Delta$ and because $|\beta_j| \geq \beta_j$. Inequality (17) uses the fact that $|z| \geq \text{sgn}(b)z$ for any two scalars z and b and that $x_j = 0$ for $j \notin T \cup \Delta$. In (18), the first equality uses $\text{sgn}(x_j)x_j = |x_j|$ and $w' A_j = \text{sgn}(x_j)$ for $j \in \Delta$. The second equality just rewrites the second term in a different form. In (19), we use the fact that $A\beta = Ax = y$ (since both β and x are feasible) to eliminate $w'(A\beta - Ax)$. Equation (20) uses $w' A_j = 0$ for $j \in T_{in}$ and the definitions of T_{a+} and T_{a-} given in (12). Finally, (21) follows because $-\sum_{j \in T_{a+}} w' A_j (\beta_j - \hat{\mu}_j - \rho) - \sum_{j \in T_{a-}} w' A_j (\beta_j - \hat{\mu}_j + \rho) \geq 0$. This holds since $-\rho \leq \beta_j - \hat{\mu}_j \leq \rho$ for all $j \in T$; $w' A_j \geq 0$ for $j \in T_{a+}$; and $w' A_j \leq 0$ for $j \in T_{a-}$.

Both inequalities (15) and (16)–(21) can hold only when $\|\beta_{T^c}\|_1 = \|x_{T^c}\|_1$, i.e., all the inequalities in (16)–(21) hold with equality. Consider the inequality in (16). Since $|w' A_j| < 1$ for $j \notin T \cup \Delta$, this holds with equality only if $\beta_j = 0$ for all $j \notin T \cup \Delta$. Since $A\beta = y = Ax$ and since both β and x are supported on $T \cup \Delta$ (or on its subset), $A_{T \cup \Delta}(\beta_{T \cup \Delta} - x_{T \cup \Delta}) = 0$. Since $\delta_{k+u} < 1$, $A_{T \cup \Delta}$ has full rank. Therefore, this means that $\beta_{T \cup \Delta} = x_{T \cup \Delta}$. Thus, we can conclude that $\beta = x$, i.e., x is the unique minimizer.

B. Proof of Lemma 2

This proof uses the following simple facts. Let $\lambda_{\min}(M)$, $\lambda_{\max}(M)$ denote the minimum and maximum eigenvalues of a matrix M . i) For positive semi-definite matrices, M , Q , $\|M\| = \lambda_{\max}(M)$; $\|MQ\| \leq \|M\|\|Q\|$; $\lambda_{\min}(M - Q) \geq \lambda_{\min}(M) - \lambda_{\max}(Q)$; and for a positive definite matrix, M , $\|M^{-1}\| = \frac{1}{\lambda_{\min}(M)}$; ii) for any matrices, B , C , $\|B - C\| \leq \|B\| + \|C\|$; iii) for disjoint sets T_1 , T_2 , $\|A_{T_1} A_{T_2}'\| \leq \theta_{|T_1|, |T_2|}$ [2, eq. (3)]; iv) $1 - \delta_{|T_1|} \leq \lambda_{\min}(A_{T_1}' A_{T_1}) \leq \lambda_{\max}(A_{T_1}' A_{T_1}) \leq 1 + \delta_{|T_1|}$ [4]; v) $M(T_b)$ is a projection matrix and so $M(T_b)M(T_b)' = M(T_b)$ and $\|M(T_b)\| = 1$; and vi) $\|\text{sgn}(x_\Delta)\|_2 = \sqrt{u}$.

The lemma assumes that $\delta_u + \delta_{k_b} + \theta_{k_b, u}^2 < 1$. This implies that a) $\delta_u < 1$ and so $A_\Delta' A_\Delta$ is positive definite and so $u \leq n$; b) $\delta_{k_b} < 1$ and so $A_{T_b}' A_{T_b}$ is positive definite and $M(T_b)$ is well-defined; and c) as we show next,

$A_\Delta' M(T_b) A_\Delta$ is positive definite and hence full rank. Since $A_\Delta' M(T_b) A_\Delta = A_\Delta' A_\Delta - A_\Delta' A_{T_b} (A_{T_b}' A_{T_b})^{-1} A_{T_b}' A_\Delta$ is a difference of two positive semi-definite matrices, thus,

$$\begin{aligned} &\lambda_{\min}(A_\Delta' M(T_b) A_\Delta) \\ &\geq \lambda_{\min}(A_\Delta' A_\Delta) - \lambda_{\max}(A_\Delta' A_{T_b} (A_{T_b}' A_{T_b})^{-1} A_{T_b}' A_\Delta) \\ &\geq (1 - \delta_u) - \frac{\theta_{k_b, u}^2}{1 - \delta_{k_b}} > 0. \end{aligned} \quad (22)$$

Thus, $A_\Delta' M(T_b) A_\Delta$ is positive definite. The first inequality in (22) follows from fact i). The second one follows because

$$\begin{aligned} &\lambda_{\min}(A_\Delta' A_\Delta) \geq (1 - \delta_u) \text{ (using fact iv)}; \\ &\lambda_{\max}(A_\Delta' A_{T_b} (A_{T_b}' A_{T_b})^{-1} A_{T_b}' A_\Delta) \\ &= \|A_\Delta' A_{T_b} (A_{T_b}' A_{T_b})^{-1} A_{T_b}' A_\Delta\| \\ &\leq \|A_\Delta' A_{T_b}\| \|(A_{T_b}' A_{T_b})^{-1}\| \|A_{T_b}' A_\Delta\| \text{ (using fact i)}; \\ &\|A_\Delta' A_{T_b}\| = \|A_{T_b}' A_\Delta\| \leq \theta_{k_b, u} \text{ (using fact iii)}; \text{ and} \\ &\|(A_{T_b}' A_{T_b})^{-1}\| = \frac{1}{\lambda_{\min}(A_{T_b}' A_{T_b})} \leq \frac{1}{1 - \delta_{k_b}} \text{ (since } A_{T_b}' A_{T_b} \\ &\text{ is positive definite, this follows using fact i) and fact iv)}. \end{aligned}$$

The third inequality of (22) follows because $(1 - \delta_u) - \frac{\theta_{k_b, u}^2}{1 - \delta_{k_b}} = \frac{1 - \delta_u - \delta_{k_b} + \delta_u \delta_{k_b} - \theta_{k_b, u}^2}{1 - \delta_{k_b}} > 0$. Both the numerator and the denominator are positive because we have assumed that $\delta_u + \delta_{k_b} + \theta_{k_b, u}^2 < 1$.

Using fact v), $A_\Delta' M(T_b) A_\Delta = A_\Delta' M(T_b) M(T_b)' A_\Delta$. Thus, using the above, $A_\Delta' M(T_b) M(T_b)' A_\Delta$ is positive definite and hence has full rank u . Thus, the $u \times n$ fat matrix, $A_\Delta' M(T_b)$ has full rank, u .

To prove the lemma, we first try to construct an $n \times 1$ vector, \tilde{w} , that satisfies the first two conditions of the lemma. Then, we show that we can find an exceptional set E so that the constructed \tilde{w} and E satisfy all the required conditions. Any \tilde{w} that satisfies $A_{T_b}' \tilde{w} = 0$ lies in the null space of A_{T_b}' and hence is of the form $\tilde{w} = M(T_b)\gamma$. To satisfy the second condition, we need a γ that satisfies $A_\Delta' M(T_b)\gamma = \text{sgn}(x_\Delta)$. As shown above, $A_\Delta' M(T_b)$ is full rank and so this system of equations has a solution (in fact has infinitely many solutions). We can compute the minimum ℓ_2 norm solution in closed form as $\gamma = M(T_b)' A_\Delta (A_\Delta' M(T_b) M(T_b)' A_\Delta)^{-1} \text{sgn}(x_\Delta)$. Since $M(T_b) M(T_b)' = M(T_b)$, $\tilde{w} = M(T_b)\gamma$ can be rewritten as

$$\tilde{w} = M(T_b) A_\Delta (A_\Delta' M(T_b) A_\Delta)^{-1} \text{sgn}(x_\Delta). \quad (23)$$

Using the definition of T_{a+g} , T_{a-g} given in (14) in Theorem 1, we can see that \tilde{w} satisfies the first two conditions of the lemma. Recall that $A_i' w > 0$ for all $i \in T_{a+g}$ is equivalent to $A_{T_{a+g}}' w \succ 0$, and similarly, $A_i' w < 0$ for all $i \in T_{a-g}$ is equivalent to $A_{T_{a-g}}' w \prec 0$.

The rest of the proof is similar to that of [2, Lemma 2]. Consider any set \tilde{T}_d disjoint with $T \cup \Delta$ of size $|\tilde{T}_d| \leq \tilde{s}$. Then,

$$\begin{aligned} \|A_{\tilde{T}_d}' \tilde{w}\|_2 &\leq \|A_{\tilde{T}_d}' M(T_b) A_\Delta\| \|(A_\Delta' M(T_b) A_\Delta)^{-1}\| \|\text{sgn}(x_\Delta)\|_2 \\ &\leq \left(\theta_{\tilde{s}, u} + \frac{\theta_{\tilde{s}, k_b} \theta_{u, k_b}}{1 - \delta_{k_b}} \right) \frac{1}{1 - \delta_u - \frac{\theta_{u, k_b}^2}{1 - \delta_{k_b}}} \sqrt{u} \\ &= a_{k_b}(u, \tilde{s}) \sqrt{u}. \end{aligned} \quad (24)$$

Notice that $a_{k_b}(u, \tilde{s})$ is positive because we have assumed that $\delta_u + \delta_{k_b} + \theta_{k_b, u}^2 < 1$. The bound in (24) follows using the simple facts given in the beginning. We obtain (24) as follows. Consider the first term $\|A_{\tilde{T}_d}' M(T_b) A_\Delta\|$. Using the definition of $M(T_b)$ and fact ii), $\|A_{\tilde{T}_d}' M(T_b) A_\Delta\| \leq \|A_{\tilde{T}_d}' A_\Delta\| + \|A_{\tilde{T}_d}' A_{T_b} (A_{T_b}' A_{T_b})^{-1} A_{T_b}' A_\Delta\|$. Using fact iii), $\|A_{\tilde{T}_d}' A_\Delta\| \leq \theta_{\tilde{s}, u}$, $\|A_{\tilde{T}_d}' A_{T_b}\| \leq \theta_{\tilde{s}, k_b}$ and $\|A_{T_b}' A_\Delta\| \leq \theta_{u, k_b}$. Since $A_{T_b}' A_{T_b}$ is positive definite, using fact i) and fact iv),

$\|(A_{T_b}' A_{T_b})^{-1}\| = \frac{1}{\lambda_{\min}(A_{T_b}' A_{T_b})} \leq \frac{1}{1-\delta_{k_b}}$. Thus, we get $\|A_{T_d}' M(T_b) A_{\Delta}\| \leq (\theta_{s,u} + \frac{\theta_{s,k_b} \theta_{u,k_b}}{1-\delta_{k_b}})$. Consider the second term $\|(A_{\Delta}' M(T_b) A_{\Delta})^{-1}\|$. Since $A_{\Delta}' M(T_b) A_{\Delta}$ is positive definite, using fact i) and (22), $\|(A_{\Delta}' M(T_b) A_{\Delta})^{-1}\| = \frac{1}{\lambda_{\min}(A_{\Delta}' M(T_b) A_{\Delta})} \leq \frac{1}{(1-\delta_u) - \frac{\theta_{u,k_b}}{1-\delta_{k_b}}}$. Using fact vi), the third term, $\|\text{sgn}(x_{\Delta})\|_2 = \sqrt{u}$.

Define the set, E , as $E := \{j \in (T \cup \Delta)^c : |A_j' \tilde{w}| > \frac{a_{k_b}(u, \tilde{s})\sqrt{u}}{\sqrt{s}}\}$. Notice that $|E|$ must obey $|E| < \tilde{s}$ since otherwise we can contradict (24) by taking $\tilde{T}_d \subseteq E$. Since $|E| < \tilde{s}$ and E is disjoint with $T \cup \Delta$, (24) holds for $\tilde{T}_d \equiv E$, i.e., $\|A_E' \tilde{w}\|_2 \leq a_{k_b}(u, \tilde{s})\sqrt{u}$. Also, by definition of E , $|A_j' \tilde{w}| \leq \frac{a_{k_b}(u, \tilde{s})\sqrt{u}}{\sqrt{s}}$, for all $j \notin T \cup \Delta \cup E$. Thus, \tilde{w} satisfies the third condition of the lemma.

Finally, $\|\tilde{w}\|_2 \leq \|M(T_b)\| \|A_{\Delta}\| \|(A_{\Delta}' M(T_b) A_{\Delta})^{-1}\| \sqrt{u} \leq K_{k_b}(u)\sqrt{u}$. This follows using fact v); $\|A_{\Delta}\| \leq \sqrt{1+\delta_u}$; and fact i) and (22). Thus, we have found a \tilde{w} and E that satisfy all required conditions.

REFERENCES

- [1] W. Lu and N. Vaswani, "Exact reconstruction conditions and error bounds for regularized modified basis pursuit (reg-modified-BP)," in *Proc. 44th Asilomar Conf. Signals, Syst., Comput.*, Pacific Grove, CA, Nov. 7–10, 2010, pp. 763–767.
- [2] N. Vaswani and W. Lu, "Modified-CS: Modifying compressive sensing for problems with partially known support," *IEEE Trans. Signal Process.*, vol. 58, no. 9, pp. 4595–4607, Sep. 2010.
- [3] S. Chen, D. Donoho, and M. Saunders, "Atomic decomposition by basis pursuit," *SIAM J. Sci. Comput.*, vol. 20, pp. 33–61, 1998.
- [4] E. Candes and T. Tao, "Decoding by linear programming," *IEEE Trans. Inf. Theory*, vol. 51, no. 12, pp. 4203–4215, 2005.
- [5] N. Vaswani, "LS-CS-residual (LS-CS): compressive sensing on the least squares residual," *IEEE Trans. Signal Process.*, vol. 58, no. 8, pp. 4108–4120, Aug. 2010.
- [6] V. Stankovic, L. Stankovic, and S. Cheng, "Compressive video sampling," presented at the Eur. Signal Process. Conf. (EUSIPCO), Lausanne, Switzerland, Aug. 25–29, 2008.
- [7] J. Y. Park and M. B. Wakin, "A multiscale framework for compressive sensing of video," presented at the Picture Coding Symp. (PCS), Chicago, IL, May 6–8, 2009.
- [8] C. Qiu and N. Vaswani, "ReProCS: A missing link between recursive robust PCA and recursive sparse recovery in large but correlated noise," arXiv: 1106.3286, 2011.
- [9] C. Qiu and N. Vaswani, "Support-predicted modified-CS for principal components' pursuit," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, St. Petersburg, Russia, 2011, pp. 668–672.
- [10] E. Candes, "The restricted isometry property and its implications for compressed sensing," *Compte Rendus de l'Academie des Sciences, Paris, Series I*, no. 346, pp. 589–592, 2008.
- [11] D. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, Apr. 2006.
- [12] D. Donoho and J. Tanner, "High-dimensional centrally symmetric polytopes with neighborliness proportional to dimension," *Discr. Comput. Geom.*, vol. 35, no. 4, pp. 617–652, 2006.
- [13] R. Von Borries, C. J. Miosso, and C. Potes, "Compressive sensing reconstruction with prior information by iteratively reweighted least-squares," *IEEE Trans. Signal Process.*, vol. 57, no. 6, pp. 2424–2431, Jun. 2009.
- [14] A. Khajehnejad, W. Xu, A. Avestimehr, and B. Hassibi, "Analyzing weighted l_1 minimization for sparse recovery with nonuniform sparse models," *IEEE Trans. Signal Process.*, vol. 59, no. 5, pp. 1985–2001, May 2011.
- [15] Y. Wang and W. Yin, "Sparse signal reconstruction via iterative support detection," *SIAM J. Imag. Sci.*, vol. 3, pp. 462–491, 2010.
- [16] N. Vaswani, "Kalman filtered compressed sensing," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, San Diego, CA, Oct. 12–15, 2008, pp. 893–896.

- [17] W. Lu and N. Vaswani, "Regularized modified BPDN for noisy sparse reconstruction with partial erroneous support and signal value knowledge," *IEEE Trans. Signal Process.*, vol. 60, no. 1, pp. 182–196, Jan. 2012.
- [18] R. Baraniuk, V. Cevher, M. Duarte, and C. Hegde, "Model-based compressive sensing," *IEEE Trans. Inf. Theory*, vol. 56, pp. 1982–2001, Apr. 2010.
- [19] P. Schniter, L. Potter, and J. Ziniel, "Fast Bayesian matching pursuit: Model uncertainty and parameter estimation for sparse linear models," in *Proc. Inf. Theory Appl. (ITA)*, La Jolla, VA, Jan. 27–Feb. 1, 2008.
- [20] S. Som, L. C. Potter, and P. Schniter, "Compressive imaging using approximate message passing and a Markov-tree prior," in *Proc. 44th Asilomar Conf. Signals, Syst., Comput.*, Pacific Grove, CA, Nov. 7–10, 2010, pp. 243–247.
- [21] J. A. Tropp, "Just relax: Convex programming methods for identifying sparse signals in noise," *IEEE Trans. Inf. Theory*, vol. 52, no. 3, pp. 1030–1051, Mar. 2006.
- [22] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [23] W. Lu and N. Vaswani, "Exact reconstruction conditions for regularized modified basis pursuit (reg-modified-BP)," [Online]. Available: <http://arxiv.org/pdf/1108.3350.pdf>, arXiv:1108.3350v1, 2011

Generalized New Mersenne Number Transforms

Said Boussakta, Monir T. Hamood, and Nick Rutter

Abstract—Two new number theoretic transforms named as odd and odd-squared new Mersenne number transforms are introduced for incorporation into a generalized new Mersenne number transforms (GNMNTs) suite, which are defined in finite fields modulo Mersenne primes where arithmetic operations and residue reductions are simple to implement. This suite is categorized by type, with detailed instructions regarding their derivations. An example is given which shows their suitability for the calculation of different types of convolutions, along with an analysis of their arithmetic complexities for radix-2 and split radix algorithms. This in turn shows that these new transforms are suitable for fast error free calculation of convolutions/correlations for signal processing and other applications.

Index Terms—New Mersenne number transform (NMNT), number theoretic transforms (NTTs), odd new Mersenne number transform (ONMNT), odd-squared new Mersenne number transform (O²NMNT).

I. INTRODUCTION

The use of number theoretic transforms (NTTs) have been firmly established within the field of signal processing [1]. This is owing to their contributing ability to perform error-free calculations over a field or a ring of integers whilst maintaining the Cyclic Convolution Property (CCP). In contrast to other methods of calculation, such as the fast Fourier transform (FFT) which involves complex arithmetic with rounding and/or truncation errors in its calculations; errors also arise in the multiplication with cosine and sine functions which are irrational, preventing exact representation in a finite precision machine [2]. Additionally, the use of NTTs have been proven to provide such

Manuscript received February 15, 2011; revised June 29, 2011 and October 17, 2011; accepted January 02, 2012. Date of publication January 26, 2012; date of current version April 13, 2012. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Sofia C. Olhede.

The authors are with the School of Electrical, Electronic and Computer Engineering, Newcastle University, NE1 7RU Newcastle Upon Tyne, U.K. (e-mail: s.boussakta@ncl.ac.uk; m.t.hamood@ncl.ac.uk; nick.rutter@ncl.ac.uk).

Digital Object Identifier 10.1109/TSP.2012.2186131