

## Claremont Colleges Scholarship @ Claremont

---

CMC Faculty Publications and Research

CMC Faculty Scholarship

---

6-10-2016

# Optimizing quantization for Lasso recovery

Xiaoyi Gu

*University of California, Los Angeles*

Shenyinying Tu

*University of California, Los Angeles*

Hao-Jun Michael Shi

*University of California, Los Angeles*

Mindy Case

*University of California, Los Angeles*

Deanna Needell

*Claremont McKenna College*

*See next page for additional authors*

---

### Recommended Citation

X. Gu, S. Tu, H-J.M. Shi, M. Case, D. Needell, and Y. Plan. "Optimizing quantization for Lasso recovery." 2016.

This Article - preprint is brought to you for free and open access by the CMC Faculty Scholarship at Scholarship @ Claremont. It has been accepted for inclusion in CMC Faculty Publications and Research by an authorized administrator of Scholarship @ Claremont. For more information, please contact [scholarship@cuc.claremont.edu](mailto:scholarship@cuc.claremont.edu).

---

**Authors**

Xiaoyi Gu, Shenyinying Tu, Hao-Jun Michael Shi, Mindy Case, Deanna Needell, and Yaniv Plan

# Optimizing quantization for Lasso recovery

Xiaoyi Gu, Shenyinying Tu, Hao-Jun Michael Shi, Mindy Case, Deanna Needell, and Yaniv Plan\*

June 10, 2016

## Abstract

This letter is focused on quantized Compressed Sensing, assuming that Lasso is used for signal estimation. Leveraging recent work, we provide a framework to optimize the quantization function and show that the recovered signal converges to the actual signal at a quadratic rate as a function of the quantization level. We show that when the number of observations is high, this method of quantization gives a significantly better recovery rate than standard Lloyd-Max quantization. We support our theoretical analysis with numerical simulations.

## 1 Introduction

We consider the structured linear model  $\mathbf{y} = \Phi \mathbf{x} \in \mathbb{R}^M$ ,  $\mathbf{x} \in K \subset \mathbb{R}^N$ ,  $\Phi \in \mathbb{R}^{N \times M}$ . Given  $\mathbf{y}, \Phi, K$ , a goal of interest is to recover  $\mathbf{x}$ . Here,  $K$  can encode sparsity as in the *Compressed Sensing* (CS) setting, or more generally, small total variation, low-rank matrices, sparsity in a dictionary, etc.

In the noisy model ( $\mathbf{y} \approx \Phi \mathbf{x}$ ), this problem can be solved by minimizing the  $\ell_2$  loss function, called  $K$ -Lasso in [1]:

$$\text{minimize } \|\Phi \mathbf{x}' - \mathbf{y}\|_2 \quad \text{subject to } \mathbf{x}' \in K. \quad (1)$$

In practice,  $\Phi \mathbf{x}$  must be quantized to fit in the digital domain, inducing inaccuracy in  $\mathbf{y}$ . Work on quantization in CS includes worst case or average distortion [2, 3], reconstruction error bounds for Basis Pursuit [4], the CoSaMP method [3, 5, 6], and relaxed Brief Propagation [7]. The uniform quantizers [8, 9, 10] and standard Lloyd-Max quantizers [3] are some typical examples of quantization schemes.

Recently, Plan and Vershynin [1] analyzed the non-linear model

$$\mathbf{y}_i = f((\Phi \mathbf{x})_i), \quad \forall i = 1, \dots, M, \quad \mathbf{x} \in K, \quad (2)$$

for some non-linear function  $f$ , giving tight recovery error bounds for  $K$ -Lasso (1). In this letter, we specialize and synthesize these bounds when the non-linearity encodes quantization. We propose to choose the quantization function which optimizes the tight error bound. This method has also been recently proposed in [11], based on a different analysis and without the focus of quantization. We carefully simplify and bound the error rate with the optimized quantization function, and show that often it significantly outperforms conventional quantization methods.

---

\*X.G., S.T., H.M.S., and M.C. are with Univ. of California, Los Angeles CA 90095, USA. D.N. is with Claremont McKenna College, Claremont CA 91711, USA. Y.P. is with Univ. of British Columbia, Vancouver BC V6T 1Z2, Canada.

## 2 Preliminaries and Model

**Preliminaries** We use the same notation as in [1]. We say an event has “high probability” if it has probability at least 0.99. The notation  $\lesssim$  hides an absolute constant, and  $g$  ( $\mathbf{g}$ ) is a standard normal variable (vector). Throughout the paper, we assume that  $\Phi$  has independent standard normal entries,  $\mathbf{y}$  follows the model (2), and  $M \gtrsim d(K)$ . We utilize the *mean width* and *tangent cone* to give an effective measure of the dimension of  $K$ . For more background on *local mean width*, see [1].

**Definition 2.1.** Let  $B_2 \subset \mathbb{R}^N$  be the unit Euclidean ball. The local mean width of a subset  $K \subset \mathbb{R}^N$  is defined as  $w(K) = \mathbb{E} \sup_{\mathbf{x} \in K \cap B_2} \langle \mathbf{x}, \mathbf{g} \rangle$ . The tangent cone of set  $K$  at  $\mathbf{x}$  is  $D(K, \mathbf{x}) := \{\tau \mathbf{h} : \tau \geq 0, \mathbf{h} \in K - \mathbf{x}\}$ .

**Theorem 2.2.** [1] Suppose  $\mathbf{x} \in S^{N-1}$ , and  $\mu \mathbf{x} \in K$ . Let  $d(K) := w(D(K, \mu \mathbf{x}))^2$ . Then with high probability, the solution  $\mathbf{x}'$  to problem (1) satisfies

$$\|\mathbf{x}' - \mu \mathbf{x}\|_2 \lesssim \left( \sqrt{d(K)} \sigma + \eta \right) / \sqrt{M}, \quad (3)$$

where  $\mu = \mathbb{E}[f(g) \cdot g]$ ,  $\sigma^2 = \mathbb{E}[(f(g) - \mu g)^2]$ , and  $\eta^2 = \mathbb{E}[(f(g) - \mu g)^2 g^2]$ .

The quantity  $d(K)$  plays the role of the dimension of  $K$  locally near the point  $\mu \mathbf{x}$ . For example, if  $K$  is an  $n$ -dimensional subspace then  $d(K) \approx n$ .

**Model** We assume the model (2) with  $f = f_Q$  for a quantization function  $f_Q$ . We assume that  $\mathbf{x} \in S^{N-1}$ , or equivalently that  $\|\mathbf{x}\|_2$  is known (and can thus be scaled); in Section 3.3 we relax this assumption.

We denote the *quantizer* or *quantization function* as

$$f_Q(x) = m_i \text{ for } x \in [\tau_i, \tau_{i+1}] \quad (4)$$

where  $\{\tau_1 = -\infty, \tau_2, \dots, \tau_{Q+1} = \infty\}$  partitions the real line and  $\{m_1, m_2, \dots, m_Q\}$  are the associated quantization values.

A conventional Lloyd-Max quantizer is derived from the following idea. The distribution of  $(\Phi \mathbf{x})_i$  is standard normal and thus known. It is natural to choose the quantizer to minimize the Mean Squared Error (MSE) function defined as

$$MSE := \mathbb{E}[(f_Q(g) - g)^2] = \sum_{i=1}^Q \int_{\tau_i}^{\tau_{i+1}} (m_i - \tau)^2 p_g(\tau) d\tau, \quad (5)$$

where  $p_g$  is the probability density function of the standard normal  $g$ . Such a minimization problem can be solved iteratively by using the Lloyd Max algorithm [12].

Note that the MSE only helps control the right-hand side of (3), but not the scaling factor  $\mu$ . This is typically not one, and may thus lead to sub-optimal recovery error. In this letter, we investigate the behavior of  $K$ -Lasso with non-linear measurements if  $f = f_Q$  is obtained from (i) minimizing the MSE (5) (conventional Lloyd-Max quantization) and (ii) minimizing the MSE (5) with restriction to  $\mu = 1$  (our proposed quantization function). We find that enforcing the latter restriction can significantly improve the error rate. That is the main result of our paper. We also note that both methods of quantization do not require knowledge of the structure,  $K$ , that is used for the  $K$ -Lasso. Thus, the results of this paper can be useful both for the various signal structures associated with CS, and also when  $x$  belongs to a linear subspace, and vanilla least-squares estimation is used.

### 3 Theoretical Results

#### 3.1 Unit norm signals and MSE without restriction

When  $f_Q$  is taken to minimize the MSE (5), we have the following corollary.

**Corollary 3.1.** *Suppose  $\mathbf{x} \in S^{N-1}$  and the quantizer of the form (4) is the minimizer:*

$$f_Q = \operatorname{argmin}_{m_i, \tau_i} \mathbb{E}[(f_Q(g) - g)^2]. \quad (6)$$

Then with high probability, the solution  $\mathbf{x}'$  to problem (1) satisfies

$$\left| Q^{-2} - \frac{\sqrt{d(K)}(Q^{-2} - Q^{-4}) + 2 - Q^{-2}}{\sqrt{M}} \right| \lesssim \|\mathbf{x}' - \mathbf{x}\|, \quad (7)$$

and

$$\|\mathbf{x}' - \mathbf{x}\| \lesssim \frac{\sqrt{d(K)}(Q^{-2} - Q^{-4}) + 2 - Q^{-2}}{\sqrt{M}} + Q^{-2}.$$

*Proof.* We use the notation of Theorem 2.2 and explicitly compute  $\mu$ ,  $\sigma$  and  $\eta$ . Since  $f_Q$  minimizes the MSE, we differentiate (5) with respect to  $\tau_i$  and  $m_i$  to get

$$\int_{\tau_i}^{\tau_{i+1}} (m_i - \tau) p_g(\tau) d\tau = 0, \quad \tau_i = \frac{m_i + m_{i-1}}{2}, \quad \forall i. \quad (8)$$

By the definition of  $\mu$  (3) and (8)

$$\mu = \sum_{i=1}^Q \int_{\tau_i}^{\tau_{i+1}} m_i^2 p_g(\tau) d\tau = \mathbb{E}[f_Q(g)^2]. \quad (9)$$

Next, set

$$e_Q := \inf_{m_i, \tau_i} \text{MSE} = \mathbb{E}[f_Q(g)^2] - 2\mathbb{E}[f_Q(g)g] + \mathbb{E}[g^2] = -\mu + 1.$$

Then,  $\sigma^2 = \mathbb{E}[f_Q(g)^2] - \mu^2 = \mu - \mu^2 = e_Q - e_Q^2$ . By the Minkowski and Hölder inequalities,

$$\begin{aligned} \eta^2 &\leq (\mathbb{E}[f_Q(g)^2 g^2]^{1/2} + \mathbb{E}[\mu^2 g^4]^{1/2})^2 \\ &\leq (\mathbb{E}[f_Q(g)^4]^{1/4} \mathbb{E}[g^4]^{1/4} + \sqrt{3}\mu)^2 \\ &= (3^{1/4} \mathbb{E}[f_Q(g)^4]^{1/4} + \sqrt{3}\mu)^2. \end{aligned}$$

In addition,

$$\begin{aligned} \mathbb{E}[g^4] - \mathbb{E}[f_Q(g)^4] &= \sum_{i=1}^Q \int_{\tau_i}^{\tau_{i+1}} (\tau^4 - m_i^4) p_g(\tau) d\tau \\ &\geq C \sum_{i=1}^Q \int_{\tau_i}^{\tau_{i+1}} (\tau^2 - m_i^2) p_g(\tau) d\tau \\ &= C e_Q \geq 0, \end{aligned} \quad (10)$$

for some positive constant  $C$ , so we have  $\mathbb{E}[f_Q(g)^4] \leq \mathbb{E}[g^4] = 3$ . Summarizing,  $\eta^2 \leq 3(\mu + 1)^2 \leq 3(2 - e_Q)^2$ . As in e.g. [13], define  $R(f_Q) := \log_2 Q$ , which represents the rate of quantizer coding. Then we have

$$e_Q = \inf_{f_Q: R_{f_Q} \leq R} \text{MSE} \cong \frac{1}{12} \left( \int p_g(\tau)^{1/3} d\tau \right)^3 2^{-2R} \cong \frac{1}{12} 6\pi\sqrt{3} Q^{-2}. \quad (11)$$

Substituting (11) into the above bounds for  $\mu$ ,  $\sigma^2$ ,  $\eta^2$  completes the proof.

The lower bound of  $\|\mathbf{x}' - \mathbf{x}\|$  follows from the fact that  $\|\mathbf{x}' - \mathbf{x}\|_2 \geq \|\mu\mathbf{x} - \mathbf{x}\|_2 - \|\mathbf{x}' - \mu\mathbf{x}\|_2 = |\mu - 1| - \|\mathbf{x}' - \mu\mathbf{x}\|_2$ .  $\square$

We see that as  $Q$  increases,  $\mathbf{x}'$  converges quadratically to  $\mathbf{x}$ . The constant 2 in the numerator eventually fades out as the number of measurements increases.

### 3.2 Optimal Quantization with restriction $\mu = 1$

We next minimize (5) while enforcing  $\mu = 1$  to obtain a bound for  $\|\mathbf{x}' - \mathbf{x}\|_2$  directly, and compare the results to the previous section.

**Corollary 3.2.** *Suppose  $\mathbf{x} \in S^{N-1}$ , and the quantizer of the form (4) is the minimizer:*

$$f_Q = \underset{m_i, \tau_i}{\text{argmin}} \mathbb{E}[(f_Q(g) - g)^2] \quad \text{s.t. } \mu = 1, \quad (12)$$

with  $\mu$  defined in (3). Then with high probability, the solution  $\mathbf{x}'$  to problem (1) satisfies

$$\|\mathbf{x}' - \mathbf{x}\| \lesssim \frac{\sqrt{d(K)}(Q^{-2} + Q^{-4}) + 2 + (Q^{-2} + Q^{-4})}{\sqrt{M}}. \quad (13)$$

*Proof.* To solve this optimization problem, we use Lagrange Multipliers to solve  $\nabla \text{MSE} = \lambda \nabla \mu$  with constraint  $\mu = 1$ . Equivalently,

$$\int_{\tau_i}^{\tau_{i+1}} ((\lambda + 2)\tau - 2m_i) p_g(\tau) d\tau = 0, \quad \tau_i = \frac{m_i + m_{i-1}}{\lambda + 2}, \quad \forall i. \quad (14)$$

Then similar to the computations in Section 3.1, we have  $\sigma^2 = \mathbb{E}[f_Q(g)^2] - \mu^2$  and

$$1 = \mu = \sum_{i=1}^Q \int_{\tau_i}^{\tau_{i+1}} \frac{2}{2 + \lambda} m_i^2 p_g(\tau) d\tau = \frac{2}{2 + \lambda} \mathbb{E}[f_Q(g)^2]. \quad (15)$$

Define  $e'_Q := \min_{m_i, \tau_i} \text{MSE}$  under the constraint  $\mu = 1$ , then  $e'_Q = \mathbb{E}[f_Q(g)^2] - 2\mathbb{E}[f_Q(g)g] + 1 = \lambda/2 = \sigma^2$  and  $\eta^2 \leq 3(e'_Q + 2)^2$  by similar computations as in Section 3.1.

Next we analyze the relationship between  $e_Q$  and  $e'_Q$ . Let  $\mathbf{m} = (m_1, \dots, m_Q)$  and let the optimal quantization levels be  $\{m_i^*, \tau_i^*\}$ . Treat the  $\text{MSE} = \text{MSE}(m_i, \tau_i^*)$  and  $\mu = \mu(m_i, \tau_i^*)$  as functions of  $\{m_i\}$  evaluated at  $\{\tau_i^*\}$ . It suffices to find the local Lipschitz constants for  $\mu$  and the MSE near  $\{m_i^*\}$ , then

$$\left\| \frac{d\mu}{d\mathbf{m}} \right\|_1 = \sum_{i=1}^Q \left| \frac{d\mu}{dm_i} \right| = \sum_{i=1}^Q \frac{1}{2\sqrt{2\pi}} \left| \int_{\tau_i^*}^{\tau_{i+1}^*} \tau e^{-\tau^2/2} d\tau \right| \gtrsim \tilde{C}, \quad (16)$$

for some  $\tilde{C} > 0$ . Define  $\delta := e_Q / \|\frac{d\mu}{d\mathbf{m}}\|_1 \leq e_Q / \tilde{C}$ . Then by the continuity of  $\mu$ , there exists  $\{m'_i\}$  that lies inside the  $\delta$ -ball of  $\{m_i^*\}$  such that  $\mu(m'_i, \tau_i^*) = 1$ . Next,

$$\begin{aligned} \left\| \frac{d\text{MSE}}{d\mathbf{m}} \right\|_1 &\leq 2 \sum_{i=1}^Q \left| \int_{\tau_i^*}^{\tau_{i+1}^*} (m_i^* - \tau) p_g d\tau \right| + 2 \sum_{i=1}^Q \left| \int_{\tau_i^*}^{\tau_{i+1}^*} \delta p_g d\tau \right| \\ &\leq 2\delta, \end{aligned}$$

gives  $|e'_Q - e_Q| \leq 2\delta^2 = 2e_Q^2 / \tilde{C}^2$ . Applying (11), we have  $e'_Q \lesssim Q^{-2} + Q^{-4}$ . Substituting this expression for  $e_Q$  into the above bounds on  $\mu$ ,  $\sigma^2$ , and  $\eta^2$  gives the desired result.  $\square$

**Remark 3.3** (Comparison of proposed method to standard Lloyd-Max quantization). *Comparing the error bound from our proposed method (13) and those of Lloyd-Max quantization (7), we see that the former is proportional to  $1/\sqrt{M}$  whereas the latter has a term which does not decrease with  $M$ . Thus, as the number of observations  $M$  increases, the proposed method gives much more accurate recovery.*

### 3.3 Quantization robustness

We next consider the case where  $\|\mathbf{x}\|_2$  is approximately 1 and study the robustness of the Lloyd-Max quantizer to the unit norm assumption. The following is a simple corollary of Theorem 2.2 which removes the assumption that  $\|\mathbf{x}\|_2 = 1$  by rescaling.

**Corollary 3.4.** *Assume that  $\mu\mathbf{x} \in K$  and let  $d(K) := w(D(K, \mu\mathbf{x}))^2$ . Then with high probability, the solution  $\mathbf{x}'$  to problem (1) satisfies*

$$\left\| \mathbf{x}' - \mu_p \frac{\mathbf{x}}{\|\mathbf{x}\|_2} \right\|_2 \lesssim \frac{\sqrt{d(K)}\sigma_p + \eta_p}{\sqrt{M}} \quad (17)$$

where  $\mu_p := \mathbb{E}[f_Q(\|\mathbf{x}\|_2 g)g]$ ,  $\sigma_p^2 := \mathbb{E}[(f_Q(\|\mathbf{x}\|_2 g) - \mu_p g)^2]$ ,  $\eta_p^2 := \mathbb{E}[(f_Q(\|\mathbf{x}\|_2 g) - \mu_p g)^2 g^2]$ .

Assume the model  $1 - \delta \leq \|\mathbf{x}\|_2 \leq 1 + \delta$  for  $\delta$  small. As before, let  $f_{Q(1)}$  be the optimal quantizer obtained from minimizing the MSE  $\mathbb{E}[(f_Q(g) - g)^2]$ . Our result shows that the recovery rate is linearly proportional to the perturbation and is inversely proportional to the quantization level with quadratic rate.

**Corollary 3.5.** *Suppose  $1 - \delta \leq \|\mathbf{x}\|_2 \leq 1 + \delta$ , and the quantizer of the form (4) is the minimizer:*

$$f_{Q(1)} = \underset{m_i, \tau_i}{\operatorname{argmin}} \mathbb{E}[(f_Q(g) - g)^2]. \quad (18)$$

Then with high probability, the solution  $\mathbf{x}'$  to problem (1) satisfies

$$\|\mathbf{x}' - \mathbf{x}\|_2 \lesssim \frac{\sqrt{d(K)}(Q^{-2} + Q^{-4} + \delta) + 2 - Q^{-2} + \delta}{\sqrt{M}} + Q^{-2} + \delta. \quad (19)$$

*Proof sketch.* Without loss of generality, assume  $\|\mathbf{x}\|_2 = 1 + \varepsilon$ , where  $0 < \varepsilon < \delta$ . Then,

$$\mu_p = \mathbb{E}[f_{Q(1)}((1 + \varepsilon)g)g] = \sum_{i=1}^Q \int_{\tau_i}^{s_i} m_i \tau p_g d\tau + \int_{s_i}^{\tau_{i+1}} m_{i+1} \tau p_g d\tau$$

where  $s_i = \tau_{i+1}/(1 + \varepsilon) = \tau_{i+1}(1 - \varepsilon) + \mathcal{O}(\varepsilon^2)$ . We can bound the difference  $|\mu - \mu_p|$  as

$$\begin{aligned} |\mu - \mu_p| &\leq \sum_{i=1}^Q \int_{s_i}^{\tau_{i+1}} |m_{i+1} - m_i| \tau p_g d\tau \\ &= \frac{\varepsilon}{\sqrt{2\pi}} \sum_{i=1}^Q |m_{i+1} - m_i| \tau_{i+1}^2 e^{-\tau_{i+1}^2/2} + \mathcal{O}(\varepsilon^2) \\ &\lesssim \frac{\varepsilon}{\sqrt{2\pi}} \sum_{i=1}^Q |\tau_{i+1}|^3 e^{-\tau_{i+1}^2/2} + \mathcal{O}(\varepsilon^2). \end{aligned}$$

Since Gaussian functions lie in the Schwartz space, this sum converges absolutely. Thus,  $|\mu - \mu_p| = \mathcal{O}(\varepsilon)$ , implying that  $|\mu_p - (1 + \varepsilon)| \leq |\mu - 1| + |\mu - \mu_p| + \varepsilon = e_Q + \mathcal{O}(\varepsilon)$ .

Next, since  $\sigma^2 = \mathbb{E}[f_{Q(1)}(g)^2] - \mu^2$ ,  $\sigma_p^2 = \mathbb{E}[f_{Q(1)}((1 + \varepsilon)g)^2] - \mu_p^2$ , it suffices to find an upper bound for  $U_\sigma = |\mathbb{E}[f_{Q(1)}(g)^2] - \mathbb{E}[f_{Q(1)}((1 + \varepsilon)g)^2]|$ . By a similar argument,

$$U_\sigma \leq \frac{4\varepsilon}{\sqrt{\pi}} \sum_{i=1}^Q |\tau_{i+1}|^3 e^{-\tau_{i+1}^2} + \mathcal{O}(\varepsilon^2) = \mathcal{O}(\varepsilon), \quad (20)$$

and,

$$\begin{aligned} |\sigma^2 - \sigma_p^2| &\leq |\mathbb{E}[f_{Q(1)}(g)^2] - \mathbb{E}[f_{Q(1)}((1 + \varepsilon)g)^2]| + |\mu^2 - \mu_p^2| \\ &= \mathcal{O}(\varepsilon) + |\mu + \mu_p| |\mu - \mu_p| \\ &= \mathcal{O}(\varepsilon). \end{aligned}$$

Finally,  $\eta_p \leq \sqrt{3}(2 - e_Q) + \mathcal{O}(\varepsilon)$ , since  $U_\eta := |\mathbb{E}[f_Q((1 + \varepsilon)g)^4] - \mathbb{E}[f_Q(g)^4]|$  is bounded as

$$\begin{aligned} U_\eta &= \frac{\varepsilon}{\sqrt{\pi}} \sum_{i=1}^Q |m_{i+1}^4 - m_i^4| \tau_{i+1} e^{-\tau_{i+1}^2} + \mathcal{O}(\varepsilon^2) \\ &\leq \frac{16\varepsilon}{\sqrt{\pi}} \sum_{i=1}^Q |\tau_{i+1}|^5 e^{-\tau_{i+1}^2} + \mathcal{O}(\varepsilon^2) \\ &= \mathcal{O}(\varepsilon). \end{aligned}$$

Corollary 3.4 and taking  $e_Q \simeq Q^{-2}$  yields the desired result.  $\square$

## 4 Numerical Experiments

Figure 1a plots reconstruction errors under the assumption that  $\|\mathbf{x}\|_2 = 1$ . All quantizers are computed using the Lloyd-Max algorithm [12]. We display the error  $\|\mathbf{x}' - \mathbf{x}\|_2$  and normalized error  $\|\frac{\mathbf{x}'}{\|\mathbf{x}'\|_2} - \mathbf{x}\|_2$  for each reconstruction method. The dimension  $N$  of the signal  $\mathbf{x}$  is 200 and the number of measurements is  $M = 50000$ . Observe that the  $K$ -Lasso with restriction to  $\mu = 1$  gives much better reconstruction than that with no restriction.

Figure 1b compares the reconstruction error of signals with unit norm and perturbed norms (1.05 here) under the same quantization levels. We only consider the recovery error of the type  $\|\mathbf{x}' - \mathbf{x}\|_2$  for simplicity.



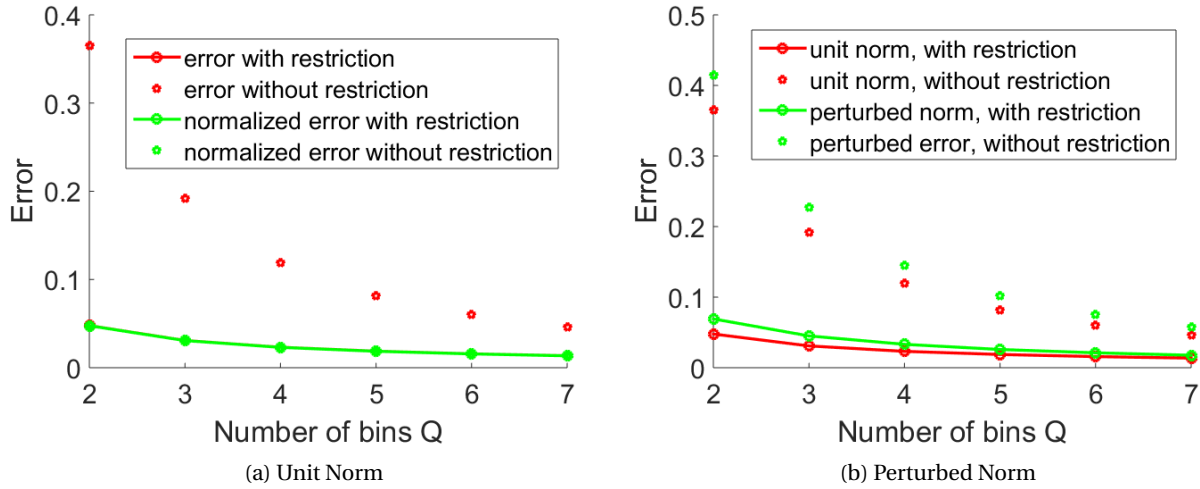


Figure 1: Left: Plot of recovery error vs number of bins for  $k = 7$ ,  $M = 50000$  and  $N = 200$ . The first, third and last lines overlap. Right: The recovery of unit norm signal and perturbed norm signal.

## 5 Conclusion

This letter extends existing work on the non-linear Lasso problem [1] to quantized Compressed Sensing with two different assumptions on the signal norm. When the signal norm is known, we show that the recovered signal converges to the actual signal with a quadratic rate in the quantization level. We also show that the quantizer obtained from restricting  $\mu = 1$  gives a better recovery rate than the conventional Lloyd-Max quantizer. When the norm is slightly perturbed, we show that the recovery rate of the conventional Lloyd-Max quantizer is inversely proportional to the level of quantization with quadratic rate, and also linearly proportional to the degree of perturbation.

## 6 Acknowledgement

This work was supported by NSF CAREER #1348721, NSF DMS #1045536, NSERC 22R23068, and the Alfred P. Sloan Foundation.

## References

- [1] Y. Plan, R. Vershynin, The generalized lasso with non-linear observations, submitted (2015).
- [2] V. Goyal, A. Fletcher, S. Rangan, Compressive sampling and lossy compression, *IEEE Signal Proc. Mag.* 25 (2008) 48–56.
- [3] W. Dai, H. V. Pham, O. Milenkovic, Quantized compressive sensing, preprint (2009).
- [4] E. J. Candès, J. K. Romberg, T. Tao, Stable signal recovery from incomplete and inaccurate measurements, *Comm. Pure Appl. Math.* 59 (8) (2006) 1207–1223.

- [5] W. Dai, O. Milenkovic, Subspace pursuit for compressive sensing signal reconstruction, *IEEE T. Inform. Theory* 55 (2009) 2230–2249.
- [6] D. Needell, J. A. Tropp, CoSaMP: Iterative signal recovery from incomplete and inaccurate samples, *Appl. Comput. Harmon. A.* 26 (3) (2009) 301–321.
- [7] U. S. Kamilov, V. K. Goyal, S. Rangan, Optimal quantization for compressive sensing under message passing reconstruction, *IEEE Int. Symp. Inform. Theory* (2011) 390–394.
- [8] E. J. Candès, J. K. Romberg, Encoding the  $\ell_p$  ball from limited measurements, *Proc. Data Compression Conf. (DDC)* (2006) 28–30.
- [9] P. T. Boufounos, R. G. Baraniuk, Quantization of sparse representations, in: Rice Univ. ECE Dept. Tech. Report 0701., 2007.
- [10] P. T. Boufounos, R. G. Baraniuk, 1-bit compressive sensing, in: 42nd Ann. Conf. Inform. Sciences and Systems (CISS), IEEE, 2008, pp. 16–21.
- [11] C. Thrampoulidis, E. Abbasi, B. Hassibi, Lasso with non-linear measurements is equivalent to one with linear measurements, in: *Advances in Neural Inform. Proc. Systems*, 2015, pp. 3402–3410.
- [12] S. P. Lloyd, Least squares quantization in PCM, *IEEE T. Inform. Theory* IT-28 (1982) 129–137.
- [13] R. M. Gray, D. L. Neuhoff, Quantization, *IEEE T. Inform. Theory* 44 (6) (1998) 2325–2383.