

Optimization for Compressed Sensing

Robert J. Vanderbei

2014 March 21

Dept. of Industrial & Systems Engineering
University of Florida

<http://www.princeton.edu/~rvdb>

Lasso Regression

The problem is to solve a sparsity-encouraging “regularized” regression problem:

$$\text{minimize } \|A\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

My reaction:

Why not replace *least squares* (LS) with *least absolute deviations* (LAD)?

LAD is to LS as median is to mean. Median is a more robust statistic (i.e., insensitive to outliers).

The LAD version can be recast as a *linear programming* (LP) problem.

If the solution is expected to be sparse, then the *simplex method* can be expected to solve the problem very quickly.

The parameter λ in the LAD version is *unitless*. This is good.

In the LS version it has units. This is bad.

No one knows the “correct” value of the parameter λ . The *parametric simplex method* can solve the problem for *all values of λ* from $\lambda = \infty$ to a small value of λ in the same (fast) time it takes the standard simplex method to solve the problem for one choice of λ .

Lasso Regression

The problem is to solve a sparsity-encouraging “regularized” regression problem:

$$\text{minimize } \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

My reaction:

Why not replace *least squares* (LS) with *least absolute deviations* (LAD)?

LAD is to LS as median is to mean. Median is a more robust statistic (i.e., insensitive to outliers).

The LAD version can be recast as a *linear programming* (LP) problem.

If the solution is expected to be sparse, then the *simplex method* can be expected to solve the problem very quickly.

The parameter λ in the LAD version is *unitless*. This is good.

In the LS version it has units. This is bad.

No one knows the “correct” value of the parameter λ . The *parametric simplex method* can solve the problem for *all values of λ* from $\lambda = \infty$ to a small value of λ in the same (fast) time it takes the standard simplex method to solve the problem for one choice of λ .

This talk is not about Lasso.

Compressed Sensing

Compressed sensing aims to recover a sparse signal from a small number of measurements.

Let $\mathbf{x}^0 = (x_1^0, \dots, x_n^0)^T \in \mathbb{R}^n$ denote a signal to be recovered.

We assume n is large and that \mathbf{x}^0 is sparse.

Let \mathbf{A} be a given (or chosen) $m \times n$ matrix with $m \ll n$.

The *compressed sensing problem* is to recover \mathbf{x}^0 assuming only that we know $\mathbf{y} = \mathbf{A}\mathbf{x}^0$ and that \mathbf{x}^0 is sparse.

Since \mathbf{x}^0 is a sparse vector, one can hope that it is the sparsest solution to the underdetermined linear system and therefore can be recovered from \mathbf{y} by solving

$$(P_0) \quad \min_{\mathbf{x}} \|\mathbf{x}\|_0 \text{ subject to } \mathbf{A}\mathbf{x} = \mathbf{y},$$

where

$$\|\mathbf{x}^0\|_0 = \#\{i : x_i \neq 0\}.$$

This problem is NP-hard due to the nonconvexity of the 0-pseudo-norm.

Basis Pursuit

Basis pursuit is one way to circumvent NP-hardness.

Replace $\|\mathbf{x}\|_0$ with $\|\mathbf{x}\|_1 = \sum_j |x_j|$:

$$(P_1) \quad \min_{\mathbf{x}} \|\mathbf{x}\|_1 \text{ subject to } \mathbf{A}\mathbf{x} = \mathbf{y}.$$

Various conditions exist under which the solutions to (P_0) and (P_1) are unique.

One key question is: under what conditions are the solutions to (P_0) and (P_1) the same?

Various sufficient conditions have been discovered.

Let \mathbf{A}_{*S} denote the submatrix of \mathbf{A} with columns from a subset $S \subset \{1, \dots, n\}$.

We say that \mathbf{A} has the *k -restricted isometry property (k -RIP)* with constant δ if for any S with cardinality k ,

$$(1 - \delta)\|\mathbf{v}\|_2^2 \leq \|\mathbf{A}_{*S}\mathbf{v}\|_2^2 \leq (1 + \delta)\|\mathbf{v}\|_2^2 \text{ for any } \mathbf{v} \in \mathbb{R}^k,$$

where $\|\mathbf{v}\|_2 = \sqrt{\sum_j v_j^2}$.

Let $\delta_k(\mathbf{A})$ denote the smallest value of δ for which the matrix \mathbf{A} has the k -RIP property.

Theorem. Let $k = \|\mathbf{x}^0\|_0$. Suppose that $k \ll n$, \mathbf{A} satisfies the k -RIP condition, and $\delta_k(\mathbf{A}) < 1/3$. Then the solutions to (P_0) and (P_1) are the same.

Previous Approaches

Existing algorithms for solving the convex program (P_1) include

- *interior-point methods*,
- *projected gradient methods*, and
- *Bregman iterations*.

Besides solving the convex program (P_1) “exactly”, several greedy algorithms have been proposed:

- *matching pursuit* and its many variants,
- *HHS pursuit*, and
- *sub-linear Fourier transform*.

Two “New” Ideas

The *first idea* is motivated by the fact that the desired solution is *sparse* and therefore should require only a relatively small number of simplex pivots to find, starting from an appropriately chosen starting point—the zero vector.

Using the *parametric simplex method* it is easy to take the zero vector as the starting basic solution.

The *second idea* requires a modified sensing scheme:

stack the signal vector \mathbf{x} into a matrix \mathbf{X} and then multiply the matrix signal on both the left and the right sides to get a compressed matrix signal.

This idea allows one to formulate the linear programming problem in such a way that the constraint matrix is very sparse and therefore the problem can be solved very efficiently.

An Equivalent Parametric Formulation

Consider the following parametric perturbation to (P_1) :

$$\begin{aligned} \hat{\mathbf{x}} &= \underset{\mathbf{x}}{\operatorname{argmin}} \quad \mu \|\mathbf{x}\|_1 + \|\boldsymbol{\epsilon}\|_1 \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} + \boldsymbol{\epsilon} = \mathbf{y}. \end{aligned} \tag{1}$$

For large values of μ , the optimal solution has $\hat{\mathbf{x}} = \mathbf{0}$ and $\hat{\boldsymbol{\epsilon}} = \mathbf{y}$.
For values of μ close to zero, the situation reverses: $\hat{\boldsymbol{\epsilon}} = \mathbf{0}$.

We refer to this problem as the *vector compressed sensing problem*.

Looks like Lasso Regression!

Dealing with Absolute Values

One way to reformulate the optimization problem (1) as a linear programming problem is to split each variable into a difference between two nonnegative variables,

$$\mathbf{x} = \mathbf{x}^+ - \mathbf{x}^- \quad \text{and} \quad \boldsymbol{\epsilon} = \boldsymbol{\epsilon}^+ - \boldsymbol{\epsilon}^-,$$

where the entries of $\mathbf{x}^+, \mathbf{x}^-, \boldsymbol{\epsilon}^+, \boldsymbol{\epsilon}^-$ are all nonnegative.

The next step is to replace $\|\mathbf{x}\|_1$ with $\mathbf{1}^T(\mathbf{x}^+ + \mathbf{x}^-)$ and to make a similar substitution for $\|\boldsymbol{\epsilon}\|_1$.

In general, the sum does not equal the absolute value but it is easy to see that it does at optimality.

Linear Programming Formulation

Here is the reformulated linear programming problem:

$$\begin{aligned} \min_{\mathbf{x}^+, \mathbf{x}^-, \boldsymbol{\epsilon}^+, \boldsymbol{\epsilon}^-} \quad & \mu \mathbf{1}^T (\mathbf{x}^+ + \mathbf{x}^-) + \mathbf{1}^T (\boldsymbol{\epsilon}^+ + \boldsymbol{\epsilon}^-) \\ \text{subject to} \quad & \mathbf{A}(\mathbf{x}^+ - \mathbf{x}^-) + (\boldsymbol{\epsilon}^+ - \boldsymbol{\epsilon}^-) = \mathbf{y} \\ & \mathbf{x}^+, \mathbf{x}^-, \boldsymbol{\epsilon}^+, \boldsymbol{\epsilon}^- \geq 0. \end{aligned}$$

For μ large, the *optimal solution* has $\mathbf{x}^+ = \mathbf{x}^- = 0$, and $\boldsymbol{\epsilon}^+ - \boldsymbol{\epsilon}^- = \mathbf{y}$.

And, given that these latter variables are required to be nonnegative, it follows that

$$\begin{aligned} y_i > 0 &\implies \epsilon_i^+ > 0 \text{ and } \epsilon_i^- = 0 \\ y_i < 0 &\implies \epsilon_i^+ = 0 \text{ and } \epsilon_i^- > 0 \\ y_i = 0 &\implies \epsilon_i^+ = 0 \text{ and } \epsilon_i^- = 0 \end{aligned}$$

With these choices, the solution is feasible for all μ and is optimal for large μ .

Furthermore, declaring the nonzero variables to be *basic* variables and the zero variables to be *nonbasic*, we see that this optimal solution is also a basic solution and can therefore serve as a starting point for the *parametric simplex method*.

Kronecker Compressed Sensing

Unlike the vector compressed sensing problem, Kronecker compressed sensing is used for sensing multidimensional signals (e.g., matrices or tensors).

For example, given a sparse matrix signal $\mathbf{X}^0 \in \mathbb{R}^{n_1 \times n_2}$, we can use two sensing matrices $\mathbf{A} \in \mathbb{R}^{m_1 \times n_1}$ and $\mathbf{B} \in \mathbb{R}^{m_2 \times n_2}$ and try to recover \mathbf{X}^0 from knowledge of $\mathbf{Y} = \mathbf{A}\mathbf{X}^0\mathbf{B}^T$ by solving the *Kronecker compressed sensing* problem:

$$(P_2) \quad \hat{\mathbf{X}} = \operatorname{argmin} \|\mathbf{X}\|_1 \quad \text{subject to} \quad \mathbf{A}\mathbf{X}\mathbf{B}^T = \mathbf{Y}.$$

When the signal is multidimensional, Kronecker compressed sensing is more natural than classical vector compressed sensing.

Kroneckerizing Vector Problems

Sometimes, even when facing vector signals, it is beneficial to use Kronecker compressed sensing due to its added computational efficiency.

More specifically, even though the target signal is a vector $\mathbf{x}^0 \in \mathbb{R}^n$, if we assume that n can be factored into $n_1 \times n_2$, then we can first stack \mathbf{x}^0 into a matrix $\mathbf{X}^0 \in \mathbb{R}^{n_1 \times n_2}$ by putting each length n_1 sub-vector of \mathbf{x}^0 into a column of \mathbf{X}^0 .

We then multiply the matrix signal \mathbf{X}^0 on both the left and the right by sensing matrices \mathbf{A} and \mathbf{B} to get a compressed matrix signal \mathbf{Y}^0 .

We will show that we are able to solve this Kronecker compressed sensing problem much more efficiently than the corresponding vector compressed sensing problem.

Vectorizing the Kroneckerization

Standard LP solvers expect to see a vector of variables, not a matrix.

The naive vectorization goes like this...

Let $\mathbf{x} = \text{vec}(\mathbf{X})$ and $\mathbf{y} = \text{vec}(\mathbf{Y})$, where, as usual, the $\text{vec}(\cdot)$ operator takes a matrix and concatenates its elements column-by-column to build one large column-vector containing all the elements of the matrix.

In terms of \mathbf{x} and \mathbf{y} , problem (P_2) can be rewritten as an equivalent *vector compressed sensing* problem:

$$\text{vec}(\hat{\mathbf{X}}) = \underset{\mathbf{x}}{\text{argmin}} \|\mathbf{x}\|_1 \quad \text{subject to} \quad \mathbf{U}\mathbf{x} = \mathbf{y},$$

where \mathbf{U} is given by the $(m_1 m_2) \times (n_1 n_2)$ Kronecker product of \mathbf{A} and \mathbf{B} :

$$\mathbf{U} = \mathbf{B} \otimes \mathbf{A} = \begin{bmatrix} \mathbf{A}b_{11} & \cdots & \mathbf{A}b_{1n_2} \\ \vdots & \ddots & \vdots \\ \mathbf{A}b_{m_2 1} & \cdots & \mathbf{A}b_{m_2 n_2} \end{bmatrix}.$$

The matrix \mathbf{U} is fully dense. *This is bad.*

Sparsifying the Constraint Matrix

The key to an *efficient* algorithm for solving the linear programming problem associated with the Kronecker sensing problem lies in noting that the dense matrix \mathbf{U} can be factored into a product of two very sparse matrices:

$$\mathbf{U} = \begin{bmatrix} \mathbf{A}b_{11} & \cdots & \mathbf{A}b_{1n_2} \\ \vdots & \ddots & \vdots \\ \mathbf{A}b_{m_21} & \cdots & \mathbf{A}b_{m_2n_2} \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{A} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{A} \end{bmatrix} \begin{bmatrix} b_{11}\mathbf{I} & b_{12}\mathbf{I} & \cdots & b_{1n_2}\mathbf{I} \\ b_{21}\mathbf{I} & b_{22}\mathbf{I} & \cdots & b_{2n_2}\mathbf{I} \\ \vdots & \vdots & \ddots & \vdots \\ b_{m_21}\mathbf{I} & b_{m_21}\mathbf{I} & \cdots & b_{m_2n_2}\mathbf{I} \end{bmatrix} =: \mathbf{V}\mathbf{W},$$

where \mathbf{I} denotes an $n_1 \times n_1$ identity matrix and $\mathbf{0}$ denotes an $m_1 \times m_1$ zero matrix.

The constraints on the problem are

$$\mathbf{U}\mathbf{x} + \boldsymbol{\epsilon} = \mathbf{y}.$$

Exploiting the Sparsification

The matrix \mathbf{U} is usually completely dense.

But, it is a product of two very sparse matrices: \mathbf{V} and \mathbf{W} .

Hence, introducing some new variables, call them \mathbf{z} , we can rewrite the constraints like this:

$$\begin{aligned}\mathbf{z} - \mathbf{W}\mathbf{x} &= 0 \\ \mathbf{V}\mathbf{z} + \boldsymbol{\epsilon} &= \mathbf{y}.\end{aligned}$$

And, as before, we can split \mathbf{x} and $\boldsymbol{\epsilon}$ into a difference between their positive and negative parts to convert the problem to a linear program:

$$\begin{aligned}& \min_{\mathbf{x}^+, \mathbf{x}^-, \boldsymbol{\epsilon}^+, \boldsymbol{\epsilon}^-} \mu \mathbf{1}^T (\mathbf{x}^+ + \mathbf{x}^-) + \mathbf{1}^T (\boldsymbol{\epsilon}^+ + \boldsymbol{\epsilon}^-) \\ \text{subject to } & \mathbf{z} - \mathbf{W}(\mathbf{x}^+ - \mathbf{x}^-) = 0 \\ & \mathbf{V}\mathbf{z} + (\boldsymbol{\epsilon}^+ - \boldsymbol{\epsilon}^-) = \mathbf{y} \\ & \mathbf{x}^+, \mathbf{x}^-, \boldsymbol{\epsilon}^+, \boldsymbol{\epsilon}^- \geq 0.\end{aligned}$$

This formulation has more variables and more constraints.

But, the constraint matrix is *very sparse*.

For linear programming, sparsity of the constraint matrix is the key to algorithm efficiency.

Numerical Results

For the *vector* sensor, we generated random problems using $m = 1,122 = 33 \times 34$ and $n = 20,022 = 141 \times 142$.

We varied the number of nonzeros k in signal \mathbf{x}^0 from 2 to 150.

We solved the straightforward linear programming formulations of these instances using an interior-point solver called *LOQO*.

We also solved a large number of instances of the parametrically formulated problem using the parametric simplex method as outlined above.

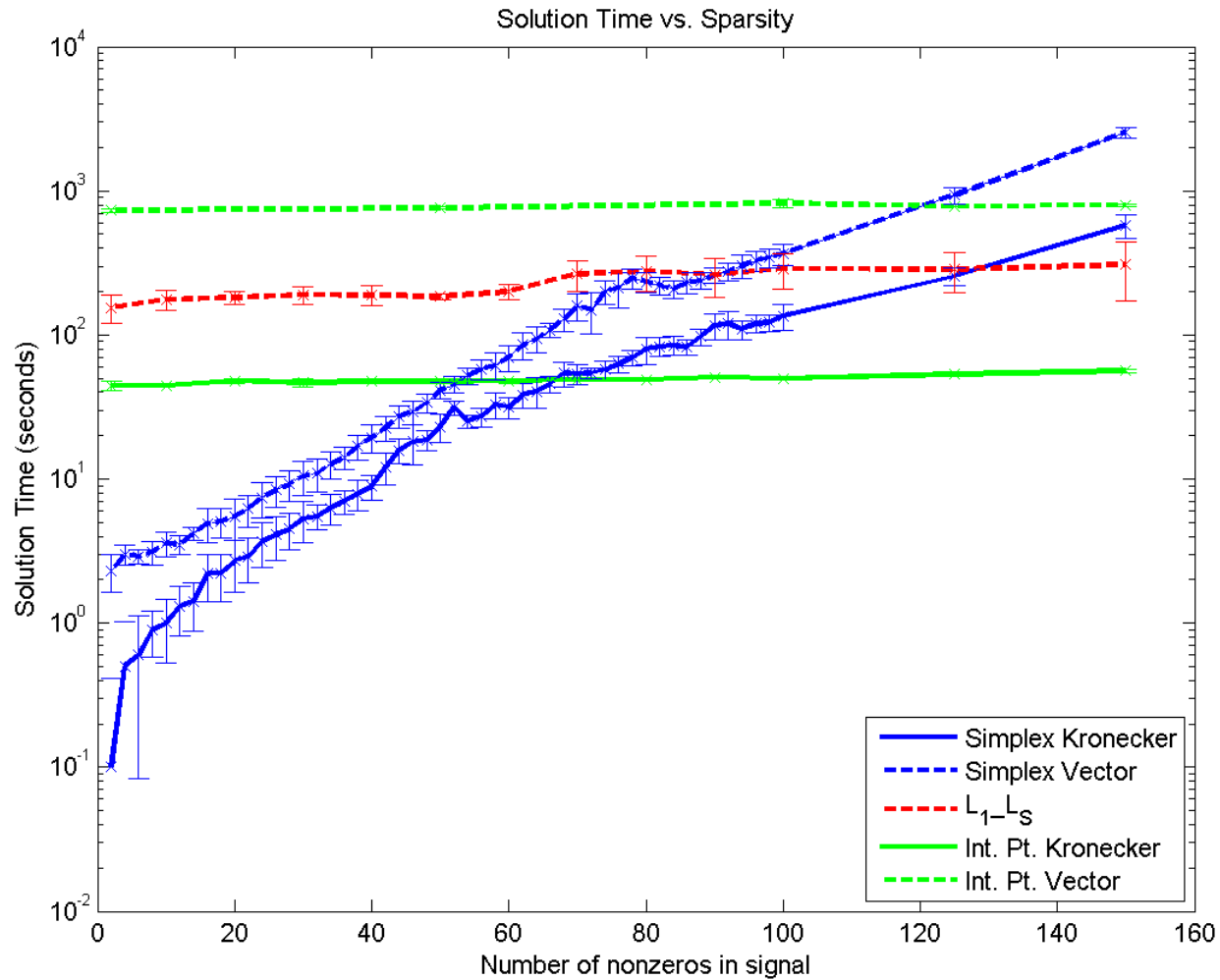
We followed a similar plan for the *Kronecker* sensor.

For these problems, we used $m_1 = 33$, $m_2 = 34$, $n_1 = 141$, $n_2 = 142$, and various values of k . Again, the straightforward linear programming problems were solved by *LOQO* and the parametrically formulated versions were solved by a custom developed parametric simplex method.

For the Kronecker sensing problems, the matrices \mathbf{A} and \mathbf{B} were generated so that their elements are independent standard Gaussian random variables.

For the vector sensing problems, the corresponding matrix \mathbf{U} was used.

We also ran the publicly-available, state-of-the-art l_1 - l_s code by Boyd et al.



The error bars have lengths equal to one standard deviation based on the multiple trials.

Conclusions

The interior-point solver (LOQO) applied to the Kronecker sensing problem is uniformly faster than both l_1 - l_s and the interior-point solver applied to the vector problem (the three horizontal lines in the plot).

For very sparse problems, the parametric simplex method is best.

In particular, for $k \leq 70$, the parametric simplex method applied to the Kronecker sensing problem is the fastest method.

It can be two or three orders of magnitude faster than l_1 - l_s .

But, as explained earlier, the Kronecker sensing problem involves changing the underlying problem being solved.

If one is required to stick with the vector problem, then it too is the best method for $k \leq 80$ after which the l_1 - l_s method wins.

Instructions for downloading and running the various codes/algorithms described herein can be found at:

http://www.orfe.princeton.edu/~rvdb/tex/CTS/kronecker_sim.html

Reference

The work described here is joint work with Han Liu and Kevin Lin.

The paper describing the work can be found at

http://orfe.princeton.edu/~rvdb/tex/CTS/cts_MPC5.pdf

Moral

“The simplex method is 200 times faster than the simplex method.” – John Forrest

Thank You!

Supplementary Info

Theoretically, KCS involves a tradeoff between computational complexity and informational complexity: it gains computational advantages at the price of requiring more measurements (i.e., larger m). More specifically, in later sections, we show that, using sub-Gaussian random sensing matrices, whenever

$$m \geq 225k^2(\log(n/k^2))^2, \quad (2)$$

we recover the true signal with probability at least $1 - 4\exp(-0.1\sqrt{m})$. It is easy to see that this scaling of (m, n, k) is tight by considering the special case when all the nonzero entries of \mathbf{x} form a continuous block.

To analyze the properties of this Kronecker sensing approach, we recall the definition of the restricted isometry constant for a matrix. For any $m \times n$ matrix \mathbf{U} , the k -restricted isometry constant $\delta_k(\mathbf{U})$ is defined as the smallest nonnegative number such that for any k -sparse vector $\mathbf{h} \in \mathbb{R}^n$,

$$(1 - \delta_k(\mathbf{U}))\|\mathbf{h}\|_2^2 \leq \|\mathbf{U}\mathbf{h}\|_2^2 \leq (1 + \delta_k(\mathbf{U}))\|\mathbf{h}\|_2^2. \quad (3)$$

Based on the results in Cai (2012), we have

Lemma. Suppose $k = \|\mathbf{X}^0\|_0$ is the sparsity of matrix \mathbf{X}^0 . Then if $\delta_k(\mathbf{U}) < 1/3$, we have $\text{vec}(\hat{\mathbf{X}}) = \mathbf{x}^0$ or equivalently $\hat{\mathbf{X}} = \mathbf{X}^0$.

For the value of $\delta_k(\mathbf{U})$, by lemma 2 of Duarte et al. (2012), we know that

$$1 + \delta_k(\mathbf{U}) \leq (1 + \delta_k(\mathbf{A}))(1 + \delta_k(\mathbf{B})). \quad (4)$$

In addition, we define strictly a sub-Gaussian distribution as follows:

Definition. We say a mean-zero random variable X follows a **strictly sub-Gaussian distribution** with variance $1/m$ if it satisfies

- $\mathbb{E}X^2 = \frac{1}{m},$
- $\mathbb{E} \exp(tX) \leq \exp\left(\frac{t^2}{2m}\right)$ for all $t \in \mathbb{R}.$

It is obvious that the Gaussian distribution with mean 0 and variance $1/m^2$ satisfies the above definition. The next theorem provides sufficient conditions that guarantees perfect recovery of the KCS problem with a desired probability.

Theorem. *Suppose matrices \mathbf{A} and \mathbf{B} are both generated by independent strictly sub-Gaussian entries with variance $1/m$. Let $C > 28.1$ be a constant. Whenever*

$$m_1 \geq C \cdot k \log(n_1/k) \quad \text{and} \quad m_2 \geq C \cdot k \log(n_2/k), \quad (5)$$

the convex program (P_2) attains perfect recovery with probability

$$\mathbb{P}\left(\widehat{\mathbf{X}} = \mathbf{X}^0\right) \geq 1 - \underbrace{2 \exp\left(-\left(0.239 - \frac{6.7}{C}\right)m_1\right) - 2 \exp\left(-\left(0.239 - \frac{6.7}{C}\right)m_2\right)}_{\rho(m_1, m_2)}. \quad (6)$$

Proof. From Equation (4) and Lemma , it suffices to show that

$$\mathbb{P}\left(\delta_k(\mathbf{A}) < \frac{2}{\sqrt{3}} - 1 \text{ and } \delta_k(\mathbf{B}) < \frac{2}{\sqrt{3}} - 1\right) \geq 1 - \rho(m_1, m_2). \quad (7)$$

This result directly follows from Theorem 3.6 of Baraniuk (2010) with a careful calculation of constants. \square

From the above theorem, we see that for $m_1 = m_2 = \sqrt{m}$ and $n_1 = n_2 = \sqrt{n}$, whenever the number of measurements satisfies

$$m \geq 225k^2(\log(n/k^2))^2, \quad (8)$$

we have $\hat{\mathbf{X}} = \mathbf{X}^0$ with probability at least $1 - 4 \exp(-0.1\sqrt{m})$.

Here we compare the above result to that of vector compressed sensing, i.e., instead of stacking the original signal $\mathbf{x}^0 \in \mathbb{R}^n$ into a matrix, we directly use a strictly sub-Gaussian sensing matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ to multiply on \mathbf{x}^0 to get $\mathbf{y} = \mathbf{A}\mathbf{x}^0$. We then plug \mathbf{y} and \mathbf{A} into the convex program (P_1) in Equation (??) to recover \mathbf{x}^0 . Following the same argument as in Theorem , whenever

$$m \geq 30k \log(n/k), \quad (9)$$

we have $\hat{\mathbf{x}} = \mathbf{x}^0$ with probability at least $1 - 2 \exp(-0.1m)$. Comparing (9) with (8), we see that KCS requires more stringent conditions for perfect recovery.