


Sparsity Matters



Robert J. Vanderbei

2017 September 20

IDA: Center for Communications Research
Princeton NJ

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

“The simplex method is 200 times faster...”

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

More Generally...

Any optimization algorithm is 200 times faster...

More Generally...

**Any optimization algorithm is 200 times faster...
than that same algorithm.**

What Can Go Wrong?

- Sometimes the straightforward/default numerical *linear algebra* isn't the best way to solve a particular type of problem.
- Sometimes the obvious natural formulation of an optimization problem isn't very easy to solve but there is a *mathematically equivalent variant* that one can solve quickly.
- Sometimes the real-world problem to be solved isn't precisely defined and an *alternate formulation* might be vastly easier to solve.

I will discuss some examples of each of these scenarios.

Numerical Linear Algebra

Linear Programming (LP)

Primal:

$$\begin{aligned} & \text{maximize} && c^T x \\ & \text{subject to} && Ax = b \\ & && x \geq 0 \end{aligned}$$

Dual:

$$\begin{aligned} & \text{minimize} && b^T y \\ & \text{subject to} && A^T y - z = c \\ & && z \geq 0 \end{aligned}$$

Interior-Point Algorithm

initialize x, y, z such that x and z are positive

while (not optimal) {

$$\rho = b - Ax$$

$$\sigma = c - A^T y + z$$

$$\mu = \delta z^T x / n$$

solve:

$$A\Delta x = \rho$$

$$A^T \Delta y + X^{-1} Z \Delta x = \sigma + \mu X^{-1} e - z$$

$$\Delta z = X^{-1} Z \Delta x - \mu X^{-1} e + z$$

$$\theta = r \left(\max_j \left\{ -\Delta x_j / x_j, -\Delta z_j / z_j \right\} \right)^{-1}$$

$$x \leftarrow x + \theta \Delta x,$$

$$y \leftarrow y + \theta \Delta y, \quad z \leftarrow z + \theta \Delta z$$

}

Linear Programming (LP)

Primal:

$$\begin{aligned} & \text{maximize} && c^T x \\ & \text{subject to} && Ax = b \\ & && x \geq 0 \end{aligned}$$

Dual:

$$\begin{aligned} & \text{minimize} && b^T y \\ & \text{subject to} && A^T y - z = c \\ & && z \geq 0 \end{aligned}$$

Interior-Point Algorithm

primal infeasibility \Rightarrow

dual infeasibility \Rightarrow

step length \Rightarrow

initialize x, y, z such that x and z are positive

while (not optimal) {

$$\rho = b - Ax$$

$$\sigma = c - A^T y + z$$

$$\mu = \delta z^T x / n$$

solve:

$$A\Delta x = \rho$$

$$A^T \Delta y + X^{-1} Z \Delta x = \sigma + \mu X^{-1} e - z$$

$$\Delta z = X^{-1} Z \Delta x - \mu X^{-1} e + z$$

$$\theta = r \left(\max_j \{ -\Delta x_j / x_j, -\Delta z_j / z_j \} \right)^{-1}$$

$$x \leftarrow x + \theta \Delta x,$$

$$y \leftarrow y + \theta \Delta y, \quad z \leftarrow z + \theta \Delta z$$

}

Note: $X = \text{Diag}(x)$ and $Z = \text{Diag}(z)$.

Reduced KKT System

$$A\Delta x = \rho$$

$$A^T \Delta y + X^{-1}Z\Delta x = \sigma + \mu X^{-1}e - z$$

Normal Equations

Use second set of equations above to solve for Δx ,

$$\Delta x = -DA^T \Delta y + D\sigma + \mu Z^{-1}e - x$$

where

$$D = XZ^{-1}$$

and eliminate Δx from the first set of equations to arrive at the *normal equations*:

$$-ADA^T \Delta y = \rho - AD\sigma - \mu AZ^{-1}e + Ax$$

Problems with Dense Columns

Consider this LP:

$$\begin{array}{ll} \text{maximize} & c^T x + c_0 x_0 \\ \text{subject to} & Ax + a x_0 = b \\ & x, x_0 \geq 0 \end{array}$$

There are practical real-world problems with this structure. Also, the “Big-M” method for algorithms that require a feasible initial starting point involves such a structure.

The computationally intensive part of the original approach to finding the step directions in interior-point methods was to solve the so-called “normal equations” for the dual step direction Δy :

$$\begin{bmatrix} A & a \end{bmatrix} \begin{bmatrix} D^2 & 0 \\ 0 & d^2 \end{bmatrix} \begin{bmatrix} A^T \\ a^T \end{bmatrix} \Delta y = \dots$$

where D is a diagonal matrix and d is a scalar. Multiplying out, we get:

$$(AD^2A^T + d^2 aa^T) \Delta y = \dots$$

The matrix AD^2A^T is *sparse* but aa^T is *dense*. Hence, the sum is dense and solving the system of equations as formulated is highly inefficient.

There are three ways to address this problem:

- Use the *Sherman-Morrison-Woodbury-AddYourNameHere* formula to express solutions to the dense system in terms of solutions to the system without the dense (rank one) term.
- Solve the *reduced KKT system* instead of the normal equations.
(*Symmetric Indefinite Systems for Interior Point Methods*, Math. Prog., 58:1-32, 1993)
- Re-express the problem using several variables in place of the single variable x_0 .
(*Splitting Dense Columns in Sparse Linear Systems*, Linear Algebra and Applications, 152:107-117, 1991)

The *dense columns* example is a bit dated but it illustrates an important point:

The number of constraints, m , and the number of variables, n , often have little to do with how long it takes to solve a problem.

The *sparsity* of the constraint matrix plays a huge role.

Example:

I recently solved a generalized network flow problem using AMPL/LOQO.

It has 195,503 variables and 123,571 constraints.

The problem solved in 3.9 seconds on my laptop computer.

For comparison, I solved a dense problem with 1/100-th as many variables and 1/100-th as many constraints.

It took 59 seconds to solve.

The *dense columns* example is a bit dated but it illustrates an important point:

The number of constraints, m , and the number of variables, n , often have little to do with how long it takes to solve a problem.

The *sparsity* of the constraint matrix plays a huge role.

Example:

I recently solved a generalized network flow problem using AMPL/LOQO.

It has 195,503 variables and 123,571 constraints.

The problem solves in 3.9 seconds on my laptop computer.

For comparison, I solved a dense problem with 1/100-th as many variables and 1/100-th as many constraints.

It took 59 seconds to solve.

SPARSITY MATTERS

Mathematically Equivalent Problems

Markowitz Model for Portfolio Selection

$$\begin{array}{ll} \text{minimize} & \mu x^T \Sigma x - r^T x \\ \text{subject to} & e^T x = 1 \\ & x \geq 0 \end{array}$$

Here, r denotes the vector of expected returns, Σ denotes the matrix of covariances of the returns, e is the vector of all ones and μ is the *risk aversion* parameter.

Equivalent Reformulation

The matrix Σ is positive semidefinite and therefore can be expressed (and probably was defined) as a product $\Sigma = U^T U$.

$$\begin{array}{ll} \text{minimize} & \mu y^T y - r^T x \\ \text{subject to} & y - Ux = 0 \\ & e^T x = 1 \\ & x \geq 0 \end{array}$$

Example. A problem with 10,000 market assets and a covariance matrix based on data from 100 time periods in the past, the first problem has 10,000 variables and just one constraint. It solves in 1480 seconds. The second formulation has 10,100 variables and 101 constraints. It solves in 10.3 seconds—a speedup by a factor of 144.

Alternate Formulations

Compressed sensing

Aims to recover a sparse signal from a small number of measurements.

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

We assume n is large and that x^0 is sparse.

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

The *compressed sensing problem* is to recover x^0 assuming only that we know $y = Ax^0$ and that x^0 is sparse.

Since 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 y by solving

$$(P_0) \quad \min_x \|x\|_0 \text{ subject to } Ax = y,$$

where

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

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 $\|x\|_0$ with $\|x\|_1 = \sum_j |x_j|$:

$$(P_1) \quad \min_x \|x\|_1 \text{ subject to } Ax = y.$$

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:

reshape the signal vector x into a matrix 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.

Details to follow...

An Equivalent Parametric Formulation

Consider the following parametric perturbation to (P_1) :

$$\begin{aligned} (\hat{x}, \hat{\varepsilon}) = \operatorname{argmin}_{x, \varepsilon} \quad & \mu \|x\|_1 + \|\varepsilon\|_1 \\ \text{subject to} \quad & Ax + \varepsilon = y. \end{aligned} \tag{1}$$

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

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

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,

$$x = x^+ - x^- \quad \text{and} \quad \varepsilon = \varepsilon^+ - \varepsilon^-,$$

where the entries of $x^+, x^-, \varepsilon^+, \varepsilon^-$ are all nonnegative.

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

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:

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

For μ large, the *optimal solution* has $x^+ = x^- = 0$, and $\varepsilon^+ - \varepsilon^- = y$.

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

$$\begin{array}{ll} y_i > 0 & \implies \varepsilon_i^+ > 0 \text{ and } \varepsilon_i^- = 0 \\ y_i < 0 & \implies \varepsilon_i^+ = 0 \text{ and } \varepsilon_i^- > 0 \\ y_i = 0 & \implies \varepsilon_i^+ = 0 \text{ and } \varepsilon_i^- = 0 \end{array}$$

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

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

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 $x^0 \in \mathbb{R}^n$, if we assume that n can be factored into $n_1 \times n_2$, then we can first reshape x^0 into a matrix $X^0 \in \mathbb{R}^{n_1 \times n_2}$ by putting successive length n_1 sub-vectors of x^0 into columns of X^0 .

We then multiply the matrix signal X^0 on both the left and the right by sensing matrices A and B to get a compressed matrix signal 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 $x = \text{vec}(X)$ and $y = \text{vec}(Y)$, 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 x and y , problem (P_2) can be rewritten as an equivalent *vector compressed sensing* problem:

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

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

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

The matrix 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 U can be factored into a product of two very sparse matrices:

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

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

The constraints on the problem are

$$Ux = y.$$

Exploiting the Sparsification

The matrix U is usually completely dense.

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

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

$$\begin{aligned} z - Wx &= 0 \\ Vz &= y. \end{aligned}$$

And, as before, we can split x into a difference between their positive and negative parts to convert the problem to a linear program:

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

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 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.

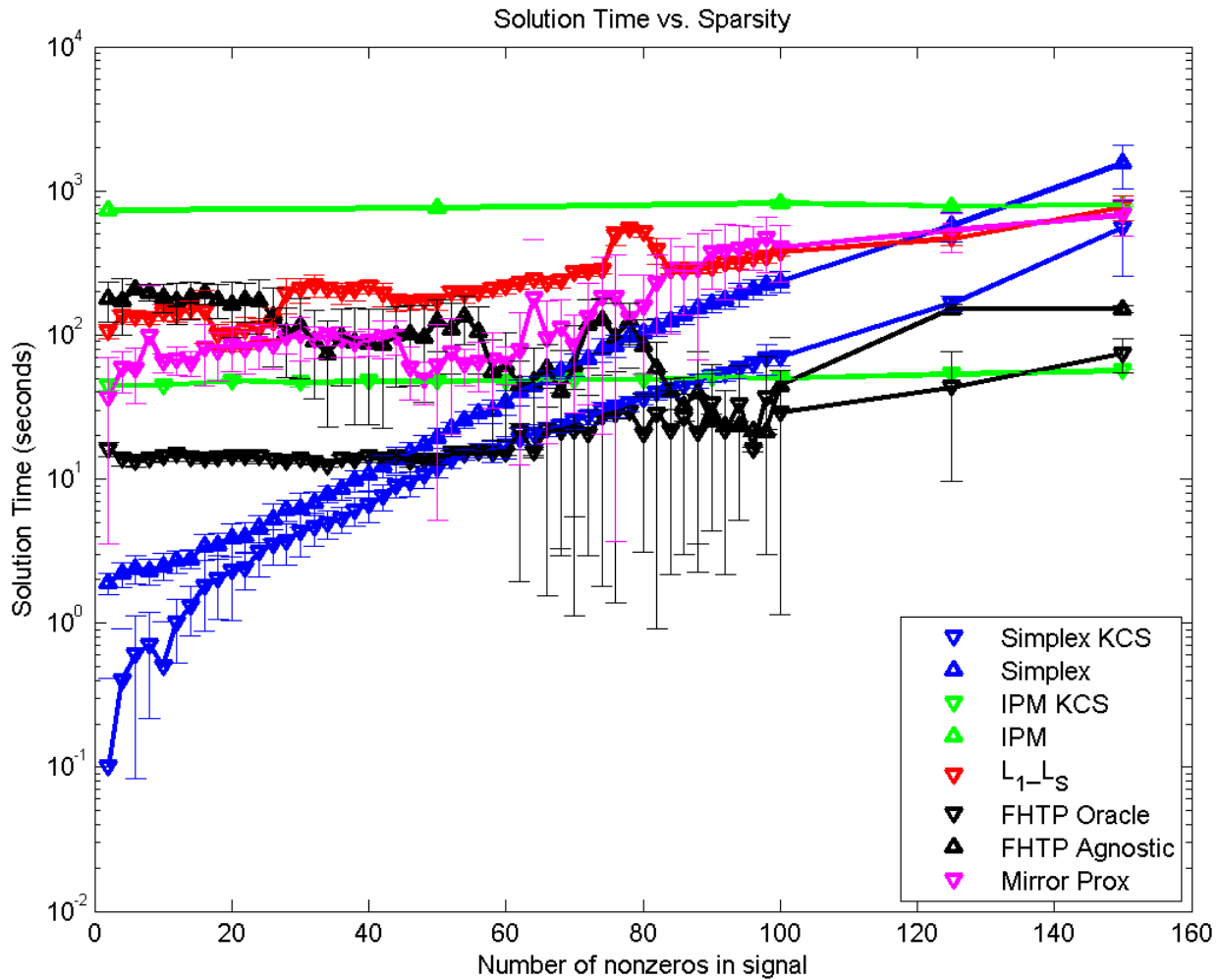
We followed a similar plan for the *Kronecker* sensor.

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.

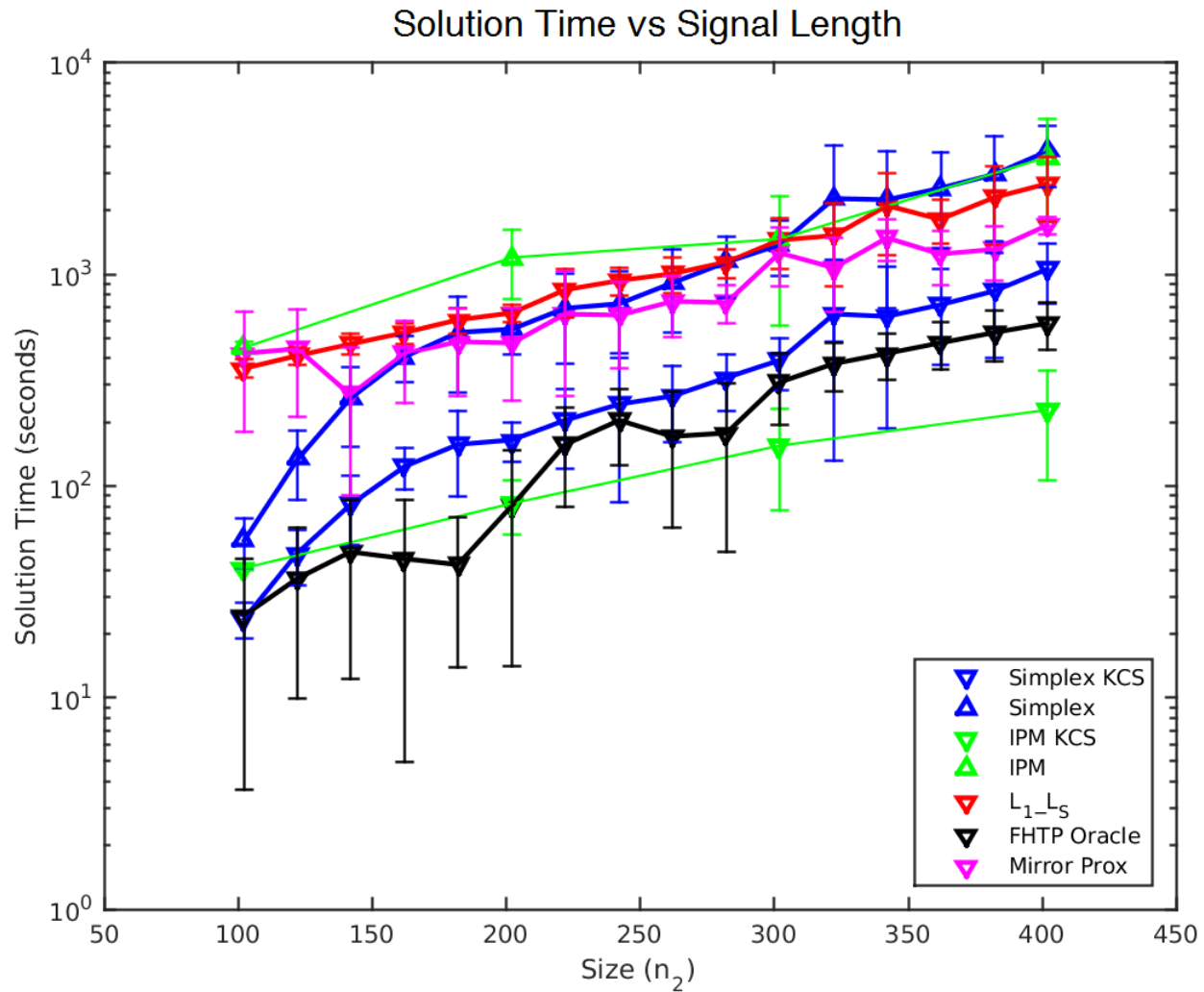
For the Kronecker sensing problems, the matrices A and B were generated so that their elements are independent standard Gaussian random variables.

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

We also ran some publicly-available, state-of-the-art codes: *L₁-L_s*, *FHTP*, and *Mirror Prox*.



$m = 1,222, n = 20,022$. Error bars represent one standard deviation.



$$m = 1,122, k = 100, n = 141 \times n_2.$$

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

Thank You!

And Let The Fun Begin...











