

A Primer on Compressive Sensing

Phil Schniter



THE OHIO STATE UNIVERSITY

Duke
UNIVERSITY



Univ. of North Carolina, Chapel Hill
12/1/2016

Traditional sensing

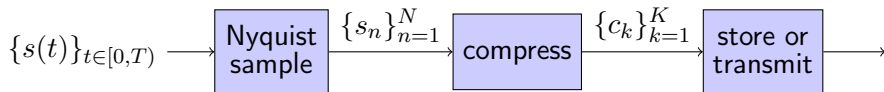
- We'd like to capture **analog** signals from the physical world and store them **digitally** on computers for subsequent processing, transmission, or reconstruction.
- Examples of “**signals**” include
 - speech or audio waveforms,
 - images (i.e., 2D waveforms),
 - video (i.e., 3D waveforms).
- The **Nyquist theorem** says that any bandlimited (i.e., smooth) signal can be sampled (giving a sequence) and then **perfectly** reconstructed.
- The **Nyquist rate** is the minimum sampling rate (i.e., # samples per unit time) needed for perfect reconstruction.

Traditional compression

Some signals are **intrinsically simple** and thus can be **compressed** without much loss of quality.

- Audio: MP3 gives roughly 10:1 compression relative to CD (=Nyquist)
- Images: JPEG gives roughly 25:1 compression relative to Nyquist
- Videos: MPEG gives roughly 100:1 compression relative to Nyquist

Compression facilitates efficient storage or transmission:



Compressive sensing

- Sometimes Nyquist sampling is **too expensive**.
- For compressible signals, Nyquist sampling is **overkill**.
- Can we do “**compressive**” sampling? Yes!
- Typical ingredients are:
 - randomly designed linear measurements
 - sparse signal representation
 - sophisticated signal reconstruction

Motivation

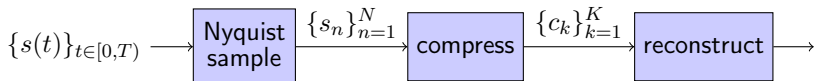
In some applications, measurements are costly:

- Magnetic resonance imaging:
 - scan time \approx 30 minutes
 - scan time proportional to # samples taken
- Imaging outside visible spectrum:
 - CMOS doesn't work
 - high cost per pixel
- Wireless communication:
 - pilots inserted to measure channel
 - more pilots means less payload

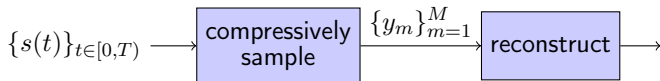


System architecture

■ Classical approach:

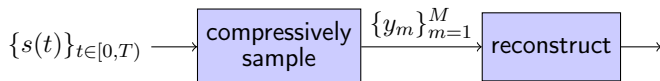


■ New approach:



$$\text{Nyquist rate } \frac{N}{T} \gg \text{compressive sampling rate } \frac{M}{T} \gtrsim \text{information rate } \frac{K}{T}$$

Principal challenges in compressive sensing



- 1 Design of the compressive-measurement scheme
- 2 Reconstruction from the compressed measurements
 - We focus on recovering the Nyquist-rate signal samples $\{s_n\}_{n=1}^N$
 - Could easily reconstruct analog $\{s(t)\}_{t \in [0, T]}$ from Nyquist samples.

Simplifying assumptions

- 1 For now, assume noiseless **linear** measurements, e.g.,

$$y_m = \int_0^T \phi_m(t) s(t) dt, \quad m = 1, \dots, M$$

- 2 Also assume signal $s(t)$ is **bandlimited**, in which case Nyquist says

$$s(t) = \sum_{n=1}^N s_n \operatorname{sinc} \left(\frac{t}{T_s} - n + 1 \right), \quad t \in [0, T).$$

Putting these together, we get the convenient **discrete** representation

$$y_m = \sum_{n=1}^N s_n \underbrace{\int_0^T \phi_m(t) \operatorname{sinc} \left(\frac{t}{T_s} - n + 1 \right) dt}_{\triangleq \Phi_{m,n}}$$

or, in **matrix/vector form**, $\boxed{\mathbf{y} = \Phi \mathbf{s}}$ for $\mathbf{s} \in \mathbb{R}^N$ and $\mathbf{y} \in \mathbb{R}^M$.

Design of linear measurements

Goal: design the matrix $\Phi \in \mathbb{R}^{M \times N}$ so that

- 1 any signal s in class \mathcal{S} can be reconstructed from $\mathbf{y} = \Phi s$,
- 2 the number of measurements M is minimal.

Key challenge:

There are fewer measurements M than unknowns N .

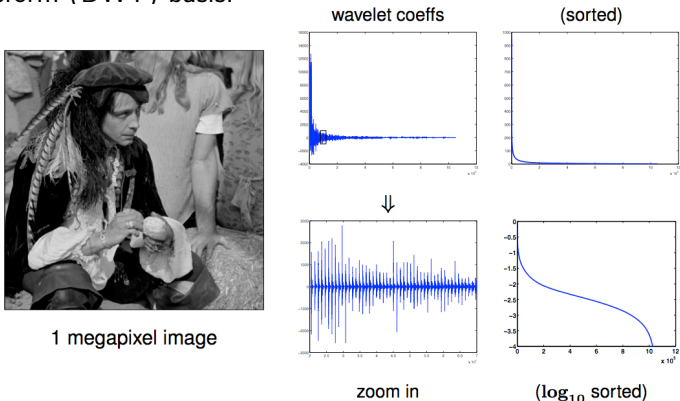
⇒ Many s satisfy the equation $\mathbf{y} = \Phi s$. How to find the correct s ?

Solution:

- If the signals in class \mathcal{S} are sufficiently **structured**, only one of the s satisfying “ $\mathbf{y} = \Phi s$ ” will be valid!
- Examples of structured signals include **sparse** signals, signals on **manifolds**, signals that can be expressed as **low-rank matrices**, etc.

Sparsity

- Many real-world signals are **approximately sparse in a known basis**.
- For example, natural images are sparse in the discrete wavelet transform (DWT) basis:



Typically: 99% signal energy captured by only 2.5% of DWT coefficients!

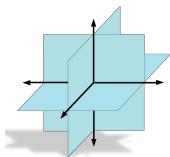
K -sparse in the dictionary Ψ

- We say that a signal class \mathcal{S} is K -sparse in the dictionary Ψ if each $s \in \mathcal{S}$ can be written as

$$s = \Psi x$$

for some K -sparse vector x (i.e., x has at most K nonzero elements).

- Usually orthonormal dictionaries Ψ are used (e.g., DWT, DCT, DFT), but overcomplete dictionaries may also be considered.
- Geometrically, a K -sparse vector $x \in \mathbb{R}^N$ lives in a union of $\binom{N}{K}$ subspaces, each of dimension K :



Merging sparsity with linear compression

Recall...

- Linear measurement model: $\mathbf{y} = \Phi \mathbf{s}$ for $\Phi \in \mathbb{R}^{M \times N}$
- Sparse signal model: $\mathbf{s} = \Psi \mathbf{x}$ for K -sparse $\mathbf{x} \in \mathbb{R}^N$

Together...

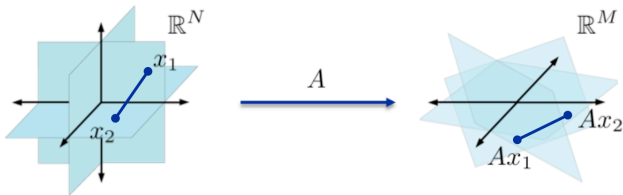
- Compressive sensing model: $\mathbf{y} = \underbrace{\Phi \Psi}_{\triangleq \mathbf{A}} \mathbf{x}$ for $\mathbf{A} \in \mathbb{R}^{M \times N}$

Questions:

- 1 What properties of \mathbf{A} ensure the recovery of \mathbf{x} ?
- 2 Given dictionary Ψ , how can we design Φ to ensure a good \mathbf{A} ?

Restricted isometry property

- Recall model: $\mathbf{y} = \mathbf{A}\mathbf{x}$ for $\mathbf{A} \in \mathbb{R}^{M \times N}$ and K -sparse $\mathbf{x} \in \mathbb{R}^N$.
- Note: if signals $\mathbf{x}_1 \neq \mathbf{x}_2$ map to the same \mathbf{y} , they can't be recovered!



- In general, for our measurement system to be **information preserving**, we want that $\|\mathbf{x}_1 - \mathbf{x}_2\|_2 \approx \|\mathbf{A}\mathbf{x}_1 - \mathbf{A}\mathbf{x}_2\|_2$ for all K -sparse $\mathbf{x}_1, \mathbf{x}_2$, or

$$1 - \delta \leq \frac{\|\mathbf{A}\mathbf{d}\|_2^2}{\|\mathbf{d}\|_2^2} \leq 1 + \delta \quad \text{for all } 2K\text{-sparse } \mathbf{d}. \quad \text{“RIP”}$$

Ensuring RIP with randomness

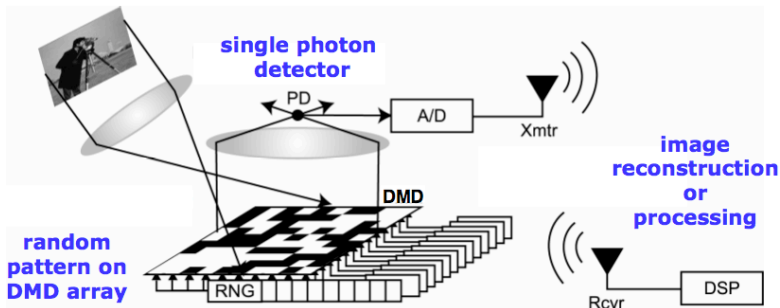
- Testing a given matrix for RIP is an NP-hard (combinatorial) problem.
- Fortunately, if \mathbf{A} is **randomly** drawn with **independent zero-mean sub-Gaussian** entries (e.g., normal, ± 1), then *with high probability* it will satisfy RIP if

$$M \geq O\left(K \log \frac{N}{K}\right).$$

- Similarly, if Φ is **constructed randomly** in the same way, then $\mathbf{A} = \Phi\Psi$ will satisfy RIP for *any* orthonormal Ψ .
- In practice, **semi-random** Φ are preferable, e.g.,

Create $\Phi = \mathbf{J}\mathbf{F}\mathbf{D}$, where \mathbf{D} is a diagonal matrix with random ± 1 s, \mathbf{F} is the N -FFT matrix, and \mathbf{J} randomly selects M outputs.

Example: Single-pixel camera (Rice Univ.)



target
65536 pixels



11000 measurements
(16%)

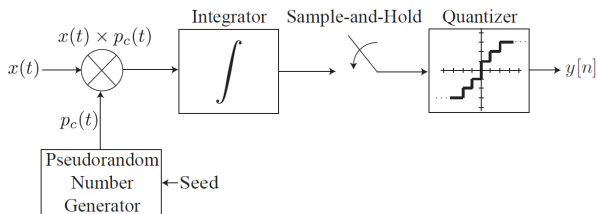


1300 measurements
(2%)

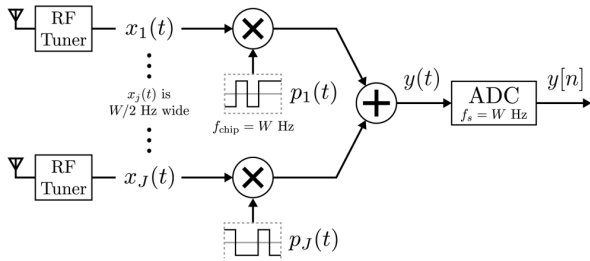


Other examples

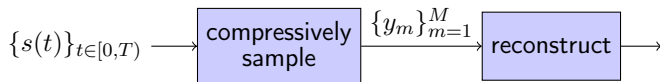
Random demodulator:



Compressive multiplexor:



Signal recovery from compressive measurements



- So far we've talked about the design of the compressive sampler. Now we'll shift focus to **signal reconstruction from compressed y** .
- In particular, we'll talk about how to **reconstruct the Nyquist-rate signal samples s** from

$$\begin{aligned}
 \mathbf{y} &= \Phi \mathbf{s} + \mathbf{w} && \text{with additive measurement noise } \mathbf{w}! \\
 &= \Phi \Psi \mathbf{x} + \mathbf{w} \\
 &= \mathbf{A} \mathbf{x} + \mathbf{w} && \text{where } \mathbf{x} \text{ is approximately } K\text{-sparse}
 \end{aligned}$$

In fact, **recovering \mathbf{x}** is enough, since we can then construct $\mathbf{s} = \Psi \mathbf{x}$.

Sparse reconstruction

Goal: estimate $x \in \mathbb{R}^N$ from $y = Ax + w \in \mathbb{R}^M$ where

- x is approximately K -sparse (although K is unknown)
- $M \ll N$ but $M \geq K$
- A is RIP-like (all subsets of K columns are nearly orthonormal)

Popular methods:

- Convex methods based on ℓ_1 -regularization
- Greedy search
- Bayesian inference

Best sparse fit — the ℓ_0 technique

Find the sparsest \mathbf{x} that explains \mathbf{y} up to a specified tolerance of ϵ :

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \underbrace{\|\mathbf{x}\|_0}_{\# \text{ nonzero coefs}} \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \epsilon.$$

Unfortunately, this is **NP-hard**; we'd need to check all $\binom{N}{K} \approx N^K$ possible supports!

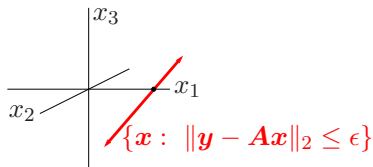
Let's think about this problem geometrically...

A toy example

Consider $\mathbf{y} = \mathbf{A} \mathbf{x} + \mathbf{w}$ with 1-sparse \mathbf{x} .

$$\begin{bmatrix} \bullet \\ \bullet \end{bmatrix} = \begin{bmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{bmatrix} \begin{bmatrix} \bullet \\ \bullet \\ \bullet \end{bmatrix} + \begin{bmatrix} \bullet \\ \bullet \end{bmatrix} \quad \begin{cases} M = 2 \\ N = 3 \\ K = 1 \end{cases}$$

- The set of \mathbf{x} such that $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \epsilon$ is described by an ϵ -thin rod.



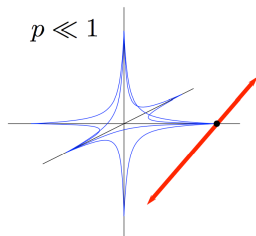
- The ℓ_0 technique would check increasingly large support hypotheses until it finds one whose signal subspace intersects the ϵ -rod. In this example, it would recover the true \mathbf{x} if $\epsilon = 0$.

The geometry of constrained ℓ_p -minimization

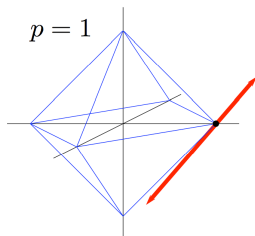
Now consider, for some fixed $p > 0$, the optimization problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \underbrace{\|\mathbf{x}\|_p}_{\sqrt[p]{\sum_n |x_n|^p}} \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \epsilon.$$

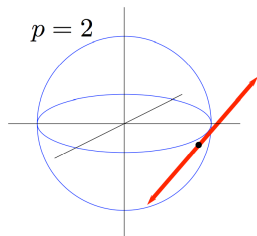
The solution can be found by **growing the ℓ_p -ball until it touches the ϵ -rod**:



Solution definitely sparse
but problem is **NP hard**.



Solution usually sparse
and problem is **convex**!



Solution is **not sparse**;
 \Leftrightarrow LS when $\epsilon = 0$.

This suggests to use the ℓ_1 norm as a surrogate for the ℓ_0 norm!

LASSO

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \epsilon$$

- Convex! Can be solved very efficiently.
- For \mathbf{A} satisfying $2K$ -RIP, LASSO guarantees that

$$\|\hat{\mathbf{x}} - \mathbf{x}\|_2 \leq \frac{C_1}{\sqrt{K}} \|\mathbf{x} - \mathbf{x}_K\|_1 + C_2 \|\mathbf{w}\|_2$$

where \mathbf{x}_K is the best K -sparse approximation of \mathbf{x} and C_1, C_2 are constants that depend on the RIP δ . Wow!

- In the special case when \mathbf{x} is K -sparse, this simplifies to

$$\|\hat{\mathbf{x}} - \mathbf{x}\|_2 \leq C_2 \|\mathbf{w}\|_2.$$

Greedy search

Main ideas:

- If we can correctly recover the **support** Λ of \mathbf{x} (i.e., the locations of nonzeros), then determining the non-zero amplitudes is easy, e.g.,

$$\mathbf{x}_\Lambda = (\mathbf{A}_\Lambda^H \mathbf{A}_\Lambda)^{-1} \mathbf{A}_\Lambda^H \mathbf{y}$$

(least squares)

The diagram shows a vertical vector \mathbf{y} on the left, a matrix \mathbf{A}_Λ in the middle, and a vertical vector \mathbf{x} on the right. The matrix \mathbf{A}_Λ is a grid of colored blocks representing columns. The equation $\mathbf{y} = \mathbf{A}_\Lambda \mathbf{x}$ is indicated by an equals sign between \mathbf{y} and the product of \mathbf{A}_Λ and \mathbf{x} . The matrix \mathbf{A}_Λ is labeled $M \times K$ below it.

- Estimate the support **sequentially**:
 - Find the column of \mathbf{A} most “similar” to \mathbf{y} and store its index.
 - Subtract the effect of this column from \mathbf{y} .
 - Repeat (until residual is sufficiently small)!

Famous algorithms include MP, OMP, IHT, CoSaMP, Subspace Pursuit

Bayesian Methods

In the Bayesian approach, one . . .

- models the signal using a **prior** pdf $p(\mathbf{x})$,
- models the measurement process using a **likelihood** function $p(\mathbf{y}|\mathbf{x})$,
- performs inference via **Bayes rule**, yielding the **posterior** pdf

$$p(\mathbf{x}|\mathbf{y}) = Z^{-1}p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) \quad \text{where } Z \text{ is a scaling constant,}$$

- often summarizes the posterior pdf by a **point estimate** like

$$\hat{\mathbf{x}} = \int \mathbf{x} p(\mathbf{x}|\mathbf{y}) d\mathbf{x} \quad \text{MMSE estimate}$$

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) \quad \text{MAP estimate}$$

and possible other statistics that quantify **estimate uncertainty**.

Bayesian interpretation of LASSO

If we assume ...

- additive white Gaussian noise of variance σ^2
- i.i.d Laplacian signal with rate λ/σ^2

then

- likelihood: $p(\mathbf{y}|\mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{M/2}} \exp(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{Ax}\|_2^2)$
- prior: $p(\mathbf{x}) = \frac{1}{(2\sigma^2/\lambda)^M} \exp(-\frac{\lambda}{\sigma^2} \|\mathbf{x}\|_1)$

for which the maximum a posteriori (MAP) estimate is

$$\begin{aligned} \hat{\mathbf{x}} &= \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) = \arg \max_{\mathbf{x}} \log (Z^{-1} p(\mathbf{y}|\mathbf{x}) p(\mathbf{x})) \\ &= \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2 + \lambda \|\mathbf{x}\|_1 \end{aligned}$$

which is an unconstrained version of the LASSO problem.

The relevance vector machine (RVM)

- The RVM is based on the *conditionally Gaussian* priors

$$p(\mathbf{x}|\boldsymbol{\alpha}) = \prod_{n=1}^N \mathcal{N}(x_n; 0, \alpha_n^{-1}) \quad \text{and} \quad p(\boldsymbol{\alpha}) = \prod_{n=1}^N \Gamma(\alpha_n; 0, 0)$$

$$p(\mathbf{w}|\boldsymbol{\beta}) \sim \prod_{m=1}^M \mathcal{N}(w_m; 0, \beta^{-1}) \quad \text{and} \quad \beta \sim \Gamma(0, 0)$$

Note that, as “precision” $\alpha_n \rightarrow \infty$, the coefficient x_n is zeroed.

- The *conditional* posterior is (due to Gaussianity) simply

$$p(\mathbf{x}|\mathbf{y}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \text{for} \quad \begin{cases} \boldsymbol{\mu} = \boldsymbol{\beta} \boldsymbol{\Sigma} \mathbf{A}^T \mathbf{y} \\ \boldsymbol{\Sigma} = (\boldsymbol{\beta} \mathbf{A}^T \mathbf{A} + \mathcal{D}(\boldsymbol{\alpha}))^{-1}. \end{cases}$$

- In practice, $\{\boldsymbol{\alpha}, \boldsymbol{\beta}\}$ are estimated using the EM algorithm and then **plugged into** $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ to *approximate* the posterior $p(\mathbf{x}|\mathbf{y})$.
- The RVM (also known as “SBL” and “BCS”) is relatively slow.

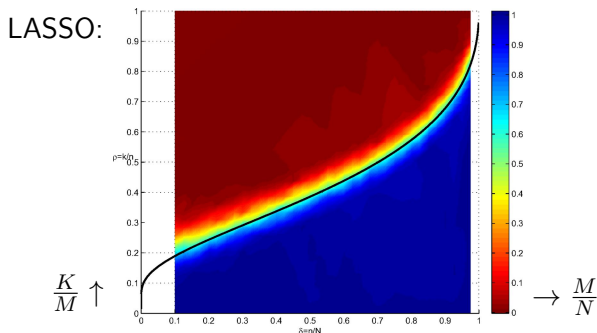
Other Bayesian methods

- Bayesian matching pursuits:
 - Greedy methods that use probabilistic support selection.

- Approximate message passing (AMP):
 - Inspired by methods from statistical physics and information theory.
 - Near-optimal in terms of speed and accuracy if \mathbf{A} is large & random.

Phase transition curve (PTC) under large random \mathbf{A}

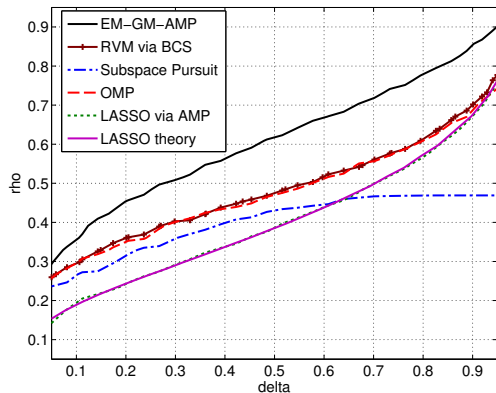
When examining a given algorithm's performance as a function of **sampling ratio** $\frac{M}{N}$ and **sparsity ratio** $\frac{K}{M}$, one finds a very sharp transition between perfect success and complete failure as $N, M, K \rightarrow \infty$.



In some cases (e.g., LASSO), the PTC can be determined *analytically*.

Algorithm comparison 1

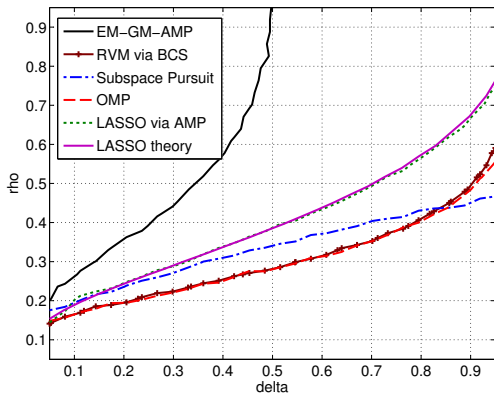
Recall: higher PTC = better algorithm.



Here, the non-zero elements of x were drawn independent [zero-mean Gaussian](#).

Algorithm comparison 1

Recall: higher PTC = better algorithm.



Here, the non-zero elements of \mathbf{x} were $= 1$.

More structure \Rightarrow *possibility* for better performance.

Conclusions

Compressive sensing . . .

- merges sampling and signal compression into a single operation
- is motivated by applications where cost-per-sample is high
- uses random linear measurements
- exploits the inherent sparsity of natural signals
- requires sophisticated algorithms for signal reconstruction.