

# Stable sparse expansions via non-convex optimization

Özgür Yılmaz

The University of British Columbia

February 19, 2008

## **Joint work with:**

- ▶ Rayan Saab (UBC)
- ▶ Rick Chartrand (Los Alamos)

To be presented at ICASSP 2008.

# Motivation: Digital Signal Processing

**Inherently analog signals:** Audio, images, seismic, etc.

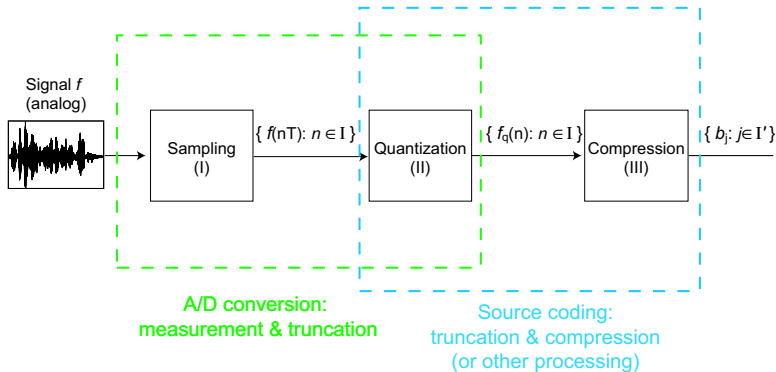
**Objective:** Use digital technology to store and process analog signals – find efficient digital representation of analog signals.

# Motivation: Digital Signal Processing

**Inherently analog signals:** Audio, images, seismic, etc.

**Objective:** Use digital technology to store and process analog signals – find efficient digital representation of analog signals.

How is this done - classical approach



# Classical Approach (ctd)

## Stage I (Sampling)

- ▶ samples obtained on a dense temporal/spatial grid,
- ▶ an appropriate sampling theorem ties resolution of “reconstruction” with the grid density.

### Example.

- ▶ Audio signals... bandlimited, thus perfect reconst. from samples taken at Nyquist rate or higher via Shannon-Nyquist Sampling Theorem. Phones: 8kHz, CDs: 44.1 kHz.
- ▶ Images... Sampling imposes a bandlimit although images are **not bandlimited**. So, sampling on denser grids in principle improves quality. (Some low-pass filtering happens in the human visual system...)

# Classical Approach (ctd)

## Stage II (Quantization)

- ▶ round-off (in a clever way) after sampling,
- ▶ can be combined with Stage III,
- ▶ not our emphasis today...theory is rich...

See the AIM Workshop (August 18-22) on:

*“Frames for the finite world: Sampling, coding, and quantization”*  
co-organized by Gunturk, Pfander, Rauhut, and OY.

# Classical Approach (ctd)

## Stage III (Compression or “Transform Coding”)

- ▶ Sampled (and quantized) signal lives in  $\mathbb{R}^N$  ( $N$  very large).
- ▶ Find a “nice” basis for  $\mathbb{R}^N$ .

# Classical Approach (ctd)

## Stage III (Compression or “Transform Coding”)

- ▶ Sampled (and quantized) signal lives in  $\mathbb{R}^N$  ( $N$  very large).
- ▶ Find a “nice” basis for  $\mathbb{R}^N$ .  
“nice” := only a few basis coef. relatively large in magnitude, i.e., basis coef. of signals of interest are approx. **sparse**.
- ▶ Exploit this sparsity and discard small coefficients.
- ▶ This is called **transform coding**.
- ▶ Example: “jpeg” format for images...



# Classical Approach (ctd)

## Stage III (Compression or “Transform Coding”)

**Formally.**  $F \subset \mathbb{R}^N$ : signals of interest, e.g., seismic signals, images, audio.

$B = [b_1 | \dots | b_N]$  an orthonormal basis for  $\mathbb{R}^N$ .

# Classical Approach (ctd)

## Stage III (Compression or “Transform Coding”)

**Formally.**  $F \subset \mathbb{R}^N$ : signals of interest, e.g., seismic signals, images, audio.

$B = [b_1 | \dots | b_N]$  an orthonormal basis for  $\mathbb{R}^N$ .

- ▶  $f \in F \Rightarrow f = Bx$  where  $x = B^T f$  is the coef. vector.

# Classical Approach (ctd)

## Stage III (Compression or “Transform Coding”)

**Formally.**  $F \subset \mathbb{R}^N$ : signals of interest, e.g., seismic signals, images, audio.

$B = [b_1 | \dots | b_N]$  an orthonormal basis for  $\mathbb{R}^N$ .

- ▶  $f \in F \Rightarrow f = Bx$  where  $x = B^T f$  is the coef. vector.
- ▶ **Choose**  $B$  such that  $x = B^T f$  is sparse for  $f \in F$ , e.g., seismic: curvelets, images: wavelets, audio: Gabor...

# Classical Approach (ctd)

## Stage III (Compression or “Transform Coding”)

**Formally.**  $F \subset \mathbb{R}^N$ : signals of interest, e.g., seismic signals, images, audio.

$B = [b_1 | \dots | b_N]$  an orthonormal basis for  $\mathbb{R}^N$ .

- ▶  $f \in F \Rightarrow f = Bx$  where  $x = B^T f$  is the coef. vector.
- ▶ **Choose**  $B$  such that  $x = B^T f$  is sparse for  $f \in F$ , e.g., seismic: curvelets, images: wavelets, audio: Gabor...

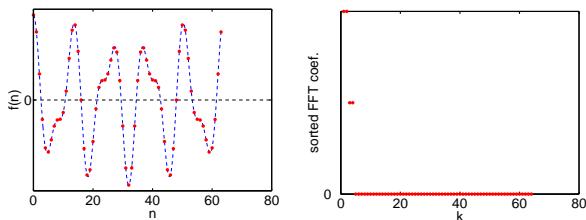
**Definitions.** Let  $x \in \mathbb{R}^N$ .

1. “0-norm”  $\|x\|_0 := \#\text{non-zero entries of } x$ .
2.  $x$  is **S-sparse** if  $\|x\|_0 \leq S$ .
3.  $x$  is **compressible** if sorted entries of  $x$  decay fast.

# Classical Approach (ctd)

## Example

Suppose  $f$  = linear combination of discrete sinusoids with period  $T$  in  $\mathbb{R}^N$ . Then use  $B_{\text{DFT}} := N$ -point DFT matrix.

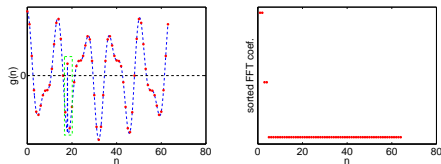


**Note.**  $f$  has 62 non-zero coef. wrt. standard basis for  $\mathbb{R}^{64}$ . On the other hand, it has only 4 non-zero coef. wrt. DFT basis.

# Classical Approach: Shortcomings (Transform Coding)

Sparsity is the key! One can obtain **much sparser** representations by using an appropriate **redundant dictionary** rather than a basis.

**Example.** Let  $g(n) = f(n) + \delta(n - 19)$ .



- ▶  $g = B_{\text{DFT}}x$ , unique solution  $x = \hat{g}$  has 64 non-zero entries.
- ▶  $A := [B_{\text{DFT}} \mid \text{Id}_{64}]$ .  $g = A\tilde{x}^*$  where  $\tilde{x}^*$  has 5 non-zero terms (4 terms for the sinusoids, one term for the dirac)!
- ▶ Difficulty:  $A$  is **not a basis**, so there are **infinitely many**  $\tilde{x}$  that solve the equation

$$g = [B_{\text{DFT}} \mid \text{Id}_{64}] \tilde{x}$$

Among the solutions, we want the **sparsest**.

# Classical Approach: Shortcomings (Sampling)

How many samples (measurements) for required resolution?

**Example.** Consider a 1 Megapixel image ( $1024 \times 1024$ ).

- ▶ Need  $2^{20}$  pixel values (sensors)  $\rightsquigarrow$  file size in the order of megabytes. (Situation is more extreme in seismology!)
- ▶ Use transform coding (images are sparse in DCT, so transform and discard the small coef.)  $\rightsquigarrow$  reduced file size: order of kilobytes (up to 90% smaller).
- ▶ In essence, collect huge amount of data in the sampling stage, throw most of it away in the compression/processing stage!
- ▶ Can we possibly reconstruct the original signal from using fewer measurements? Alternatively: can we combine sampling and compression stages to one compressive sampling stage?

## Combine sampling and compression: compressive sampling

Signal  $f \in \mathbb{R}^n$ , want to collect information on  $f$ . Take “generalized samples” or “measurements”  $b_j = \langle \mu_j, f \rangle$  where  $\mu_j \in \mathbb{R}^n$ , i.e.,

$$b = Mf, \quad \mu_j^T: \text{rows of the “measurement matrix” } M.$$



# Combine sampling and compression: compressive sampling

Signal  $f \in \mathbb{R}^n$ , want to collect information on  $f$ . Take “generalized samples” or “measurements”  $b_j = \langle \mu_j, f \rangle$  where  $\mu_j \in \mathbb{R}^n$ , i.e.,

$$b = Mf, \quad \mu_j^T: \text{rows of the “measurement matrix” } M.$$

## Remarks.

- ▶ If  $M$  is an invertible square matrix,  $n$  measurements ( $b$ ) determine  $f$  via  $f = M^{-1}b$ .
- ▶ Can we get  $f$  back from fewer than  $n$  measurements? Need additional information on  $f$ ...
- ▶ If  $f$  admits a sparse representation wrt. some known basis  $B$ ,

$$Mf = M \underbrace{Bx}_f = \underbrace{MB}_A x = b.$$

We again have an undetermined system with infinitely many solutions. But we know  $x$  should be sparse.

# Sparse recovery problem

In both cases above (transform coding with redundant dictionaries and compressive sampling), we need to solve:

## Sparse Recovery Problem

Find a sparse / the sparsest  $x$  that satisfies

$$b = Ax + r.$$

- ▶  $A \in \mathbb{R}^{m \times n}$ , with  $m < n$ ,
- ▶  $b \in \mathbb{R}^m$ : signal in transform coding, measurements in compressive sampling,
- ▶  $r \in \mathbb{R}^m$ : additive noise
- ▶  $x$ : sparse coef. vector for the signal (wrt a redundant dictionary in transform coding, wrt a basis in compressive sampling).

# Sparse recovery problem – applications to inverse problems

Let  $f \in \mathbb{R}^n$  be a signal of interest.

- ▶ Choose a basis (or dictionary)  $\Phi$ , so that  $f = \Phi x$  where  $x$  is “sparse”.
- ▶ Next, take  $n$  measurements ( $M$  below is square, invertible, possibly identity):

$$Mf = M\Phi x.$$

- ▶ Now, restrict the measurement matrix – drop some rows of  $M$ :

$$R_m M \underbrace{f}_{\text{want to find}} = \underbrace{R_m M \Phi}_A x = \underbrace{b}_{\text{accessible}}$$

- ▶ Again,  $A$  is  $m \times n$ ,  $m < n$ . To find  $f$ , we need to solve the sparse recovery problem: find sparse  $x$  that solves the underdetermined system  $Ax = b$ .

**Note.** If we can find  $x$ , we can reconstruct  $f$ .

# Sparse recovery problem – fundamental questions

Consider the undetermined system  $Ax = b + r$ . Want to find sparse(st) solution  $x$ .

## Fundamental problems.

1. Is there a **unique sparsest solution**, in particular when  $r = 0$ ?
2. Can one find a sparse / the sparsest solution in a **computationally tractable** way?
3. **Robustly** in the noisy setting (when  $r \neq 0$ )?
4. **Fast algorithm** that gives solutions guaranteed to be sparse (in some sense)?
5. How should we **choose  $A$**  so that we have a **favorable scenario**?

# Sparse recovery problem – optimization problems

1. Ideally: Choose the solution that has **smallest 0-norm**.

$$P_0^\sigma : \min_x \|x\|_0 \quad \text{subject to} \quad \|Ax - b\|_2 \leq \sigma$$

This problem is combinatorial and **NP hard**. Need alternatives!

2. Choose the solution that has **smallest 2-norm**.

$$P_2^\sigma : \min_x \|x\|_2 \quad \text{subject to} \quad \|Ax - b\|_2 \leq \sigma$$

This is the classical LS problem. The solution is **not sparse**.

3. Choose the solution that has **smallest 1-norm**.

$$P_1^\sigma : \min_x \|x\|_1 \quad \text{subject to} \quad \|Ax - b\|_2 \leq \sigma$$

This can be formulated as a convex program.

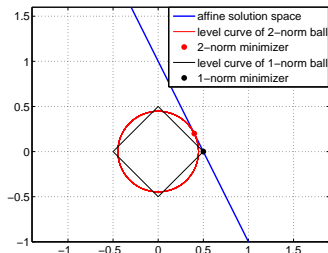
Moreover, unlike 2-norm, **1-norm promotes sparsity**.

# Sparse recovery problem: 1-norm vs. 2-norm

**Toy example.** Solve

$$P_q : \min_{x,y} \|[x \ y]^T\|_0 \quad \text{subject to} \quad [2 \ 1][x \ y]^T = 1$$

- ▶ **q=0** Sparsest solutions (not unique) of  $2x + y = 1$ :  
 $(x, y) = (1/2, 0)$  and  $(x, y) = (0, 1)$ .
- ▶ **q=2** The LS solution is  $(x, y) = (2/5, 1/5)$ , clearly not sparse.
- ▶ **q=1** The solution is  $(x, y) = (1/2, 0)$ , one of the two sparsest solutions.



# Sparse recovery by $P_1$

Recent exciting developments show that

$$P_0^\sigma : \min_x \|x\|_0 \quad \text{subject to} \quad \|Ax - b\|_2 \leq \sigma$$

can be solved in a computationally tractable way in certain cases.

**Theorem.**[Candès et al., Donoho et al.]  $P_0^\sigma$  is “equivalent” to  $P_1^\sigma$  provided:

- (i)  $\exists$  a “sufficiently sparse” solution,
- (ii)  $A$  is “sufficiently similar” to an orthonormal matrix.

# Candès-Tao-Romberg Theory – Conditions on $A$

Next, we want to specify precise conditions on  $A$  that ensure successful sparse recovery via  $P_1$ .

## Restricted isometry constants

Let  $A = [a_1 | a_2 | \dots | a_n]$  where  $a_j \in \mathbb{R}^m$ , thus  $A \in \mathbb{R}^{m \times n}$ . Suppose  $\delta_S > 0$  such that  $\forall c \in \mathbb{R}^n$ ,  $\|c\|_0 \leq S$ ,

$$(1 - \delta_S)\|c\|_2^2 \leq \|Ac\|_2^2 \leq (1 + \delta_S)\|c\|_2^2.$$

Intuitively,  $m \times S$  submatrices of  $A$  are like isometries.

### Note.

- ▶ The closer  $\delta_S$  to 0, the better the analogy.
- ▶  $A$  is orthogonal  $\Rightarrow \delta_S = 0$ .



# Candès-Romberg-Tao Theory – Exact Recovery ( $\sigma = 0$ )

**Theorem 1.**[Candès et al.] Assume  $\|x\|_0 \leq S$ , and  $b = Ax$ . Solving  $P_1$  recovers  $x$  exactly if for some  $k$

$$(A1) \quad \delta_{kS} + k\delta_{(k+1)S} < k - 1.$$

**Theorem 2.**[Candès et al.] Let  $A$  be an  $m \times n$  Gaussian matrix: each entry of  $A$  is i.i.d.  $N(0, 1/m)$ . Then  $A$  satisfies the above condition for  $S$  w.o.p. if

$$S \sim Cm / \log(n/m).$$

## Remark and Question.

- ▶ **Checking (A1)** numerically is **intractable** as  $n$  and  $m$  grow.
- ▶ Can one **construct deterministic measurement matrices** which obey (A1) for  $S \sim m / \log(n/m)$  (optimal)?
- ▶ Known constructions ( DeVore ) have  $S \sim \sqrt{m}$ .

## C-R-T Theory – Compressible and Noisy Cases.

**Theorem.**[Candès et al.] Assume  $x \in \mathbb{R}^n$  is arbitrary, and  $b = Ax + r$ . Suppose  $\delta_{kS} + k\delta_{(k+1)S} < k - 1$ . for some  $k$ . Then the solution  $x^*$  to  $P_1^\sigma$  obeys

$$\|x^* - x\|_2 \leq C_{1,S}\sigma + C_{2,S} \frac{\|x - x_S\|_1}{\sqrt{S}}.$$

Here  $x_S$  is the truncated vector, obtained by keeping  $S$  largest-in-magnitude entries of  $x$ .

### Remarks.

- ▶  $x$  is  $S$ -sparse and  $\sigma = 0 \Rightarrow$  we get the previous theorem.
- ▶  $x$  is  $S$ -sparse  $\Rightarrow$  solution is accurate within the noise level.
- ▶  $x$  is “compressible”  $\Rightarrow$  solution comparable to best  $S$ -term approx.

# C-R-T Theory – Restrictions

**Objective.** Recover  $x$  from  $b = Ax$  if  $\exists$  unique sparsest solution.

**Theorem.**  $\exists$  a unique sparsest solution if  $\|x\|_0 < S$  and  $\delta_{2S} < 1$ .

**Compare.** Sufficient condition for  $P_1^\sigma$  to recover  $x$ :  $\|x\|_0 < S$   
where

$$\delta_{kS} + k\delta_{(k+1)S} < k - 1.$$

Set  $k = 2$ , then we need  $\delta_{2S} + 2\delta_{3S} < 1$ . This cannot hold unless  $\delta_{2S} < 1/3$ . There is a **huge gap!**

**Question.** Can we **shrink this gap** by considering optimization problems other than  $P_1$ , possibly non-convex?

**Rest of the talk.** The answer is YES. We will consider

$$P_q^\sigma : \min_x \|x\|_q \quad \text{subject to} \quad \|Ax - b\|_2 \leq \sigma, \quad 0 < q < 1.$$

## Sparse recovery with $P_q^\sigma$

**Theorem.**[Saab, Chartrand, OY] Assume  $x \in \mathbb{R}^n$  is arbitrary, and  $b = Ax + r$ . Suppose for some  $k$

$$(A2) \quad \delta_{kS} + k^{2/p-1} \delta_{(k+1)S} < k^{2/p-1} - 1.$$

Then the solution  $x_q^*$  to  $P_q^\sigma$  obeys

$$\|x_q^* - x\|_2 \leq C_{q,k,S}^1 \sigma^p + C_{q,k,S}^2 \frac{\|x - x_S\|_p^p}{S^{1-p/2}}.$$

Here  $x_S$  is the truncated vector, obtained by keeping  $S$  largest entries of  $x$ .

### Remarks.

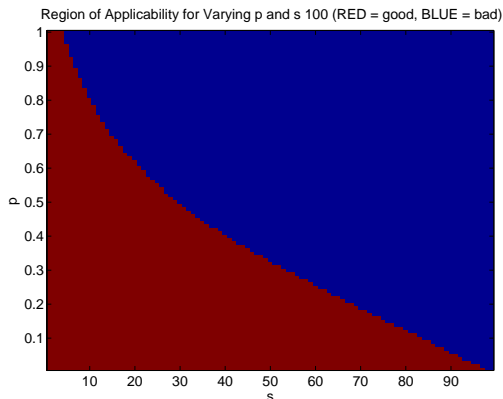
- ▶  $x$  is  $S$ -sparse and  $\sigma = 0 \Rightarrow (A2)$  implies exact reconstruction.
- ▶  $x$  is  $S$ -sparse  $\Rightarrow$  solution is accurate within the noise level.
- ▶  $x$  is “compressible”  $\Rightarrow$  solution comparable to best  $S$ -term approx.

# Significance

Main reasons why going to  $P_q$ ,  $q < 1$  pays off.

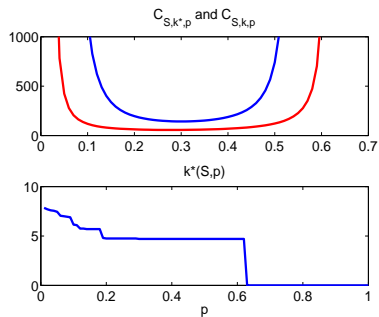
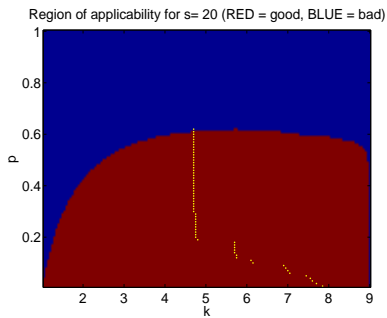
Reason 1. (A2) is less restrictive than (A1).

Take a  $256 \times 1024$  Gaussian matrix...

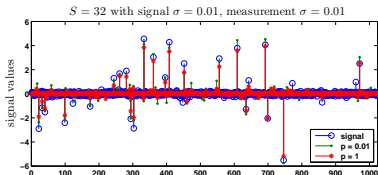
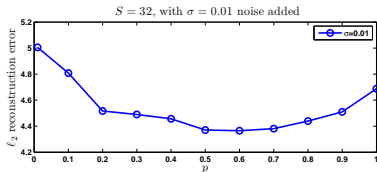
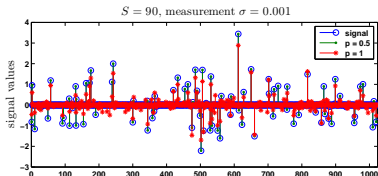
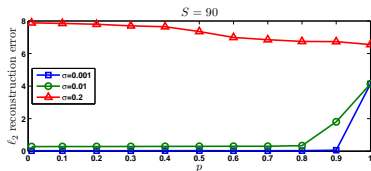
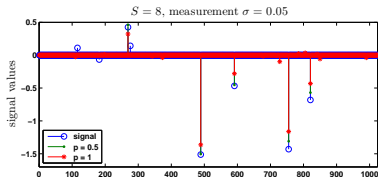
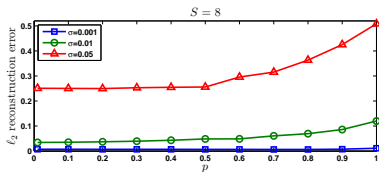


# Significance

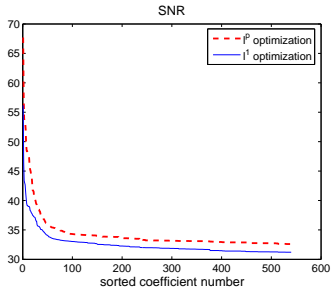
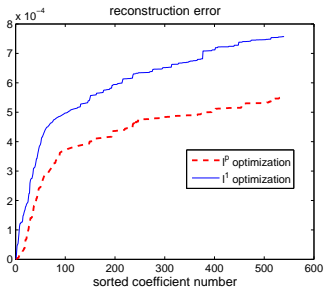
## Reason 2. Better constants, smaller error



## Reason 2. Better constants, smaller error (ctd.)



### Reason 3. Significant improvement with compressible signals





# Acknowledgements

This presentation was carried out as part of the SINBAD project with financial support, secured through ITF, from the following organizations: BG, BP, Chevron, ExxonMobil, and Shell. SINBAD is part of the collaborative research & development (CRD) grant number 334810-05 funded by the Natural Science and Engineering Research Council (NSERC).