

Algorithms for sparse optimization

Michael P. Friedlander

SLIM 
University of British Columbia

Outline

- motivation
- example applications and formulations
- main algorithmic approaches

projected gradient, iterative soft-thresholding, proximal point, augmented Lagrangian, Bregman, alternating direction of method of multipliers, split Bregman

- sample research topics
- software

MOTIVATION

Sparse solutions

$$\boxed{A} \quad x \approx b \quad \text{or} \quad \boxed{A} \quad x \approx b$$

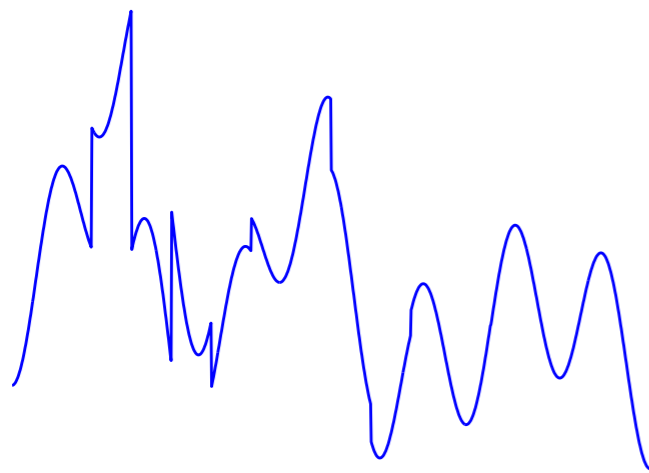
Problem: find sparse solution x

$$\text{minimize } \text{nnz}(x) \quad \text{subj to } Ax \approx b$$

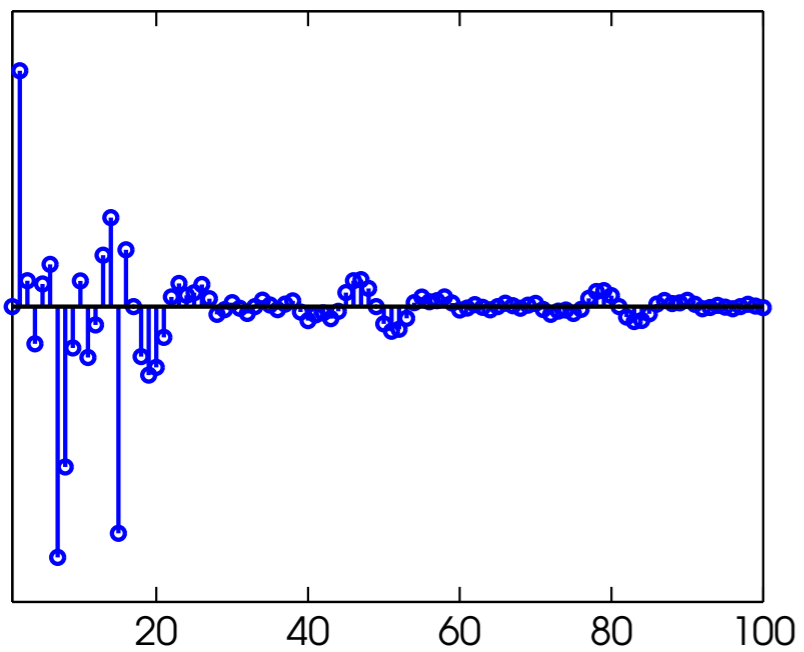
Drawbacks

- combinatorial problem
- tractable for only trivial problems

Sparse representations

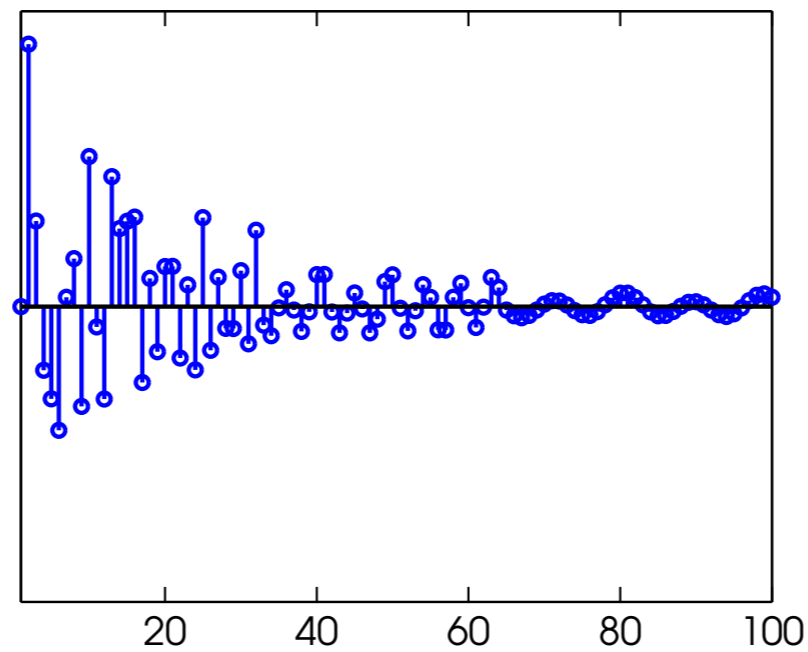


$$y = Hx$$



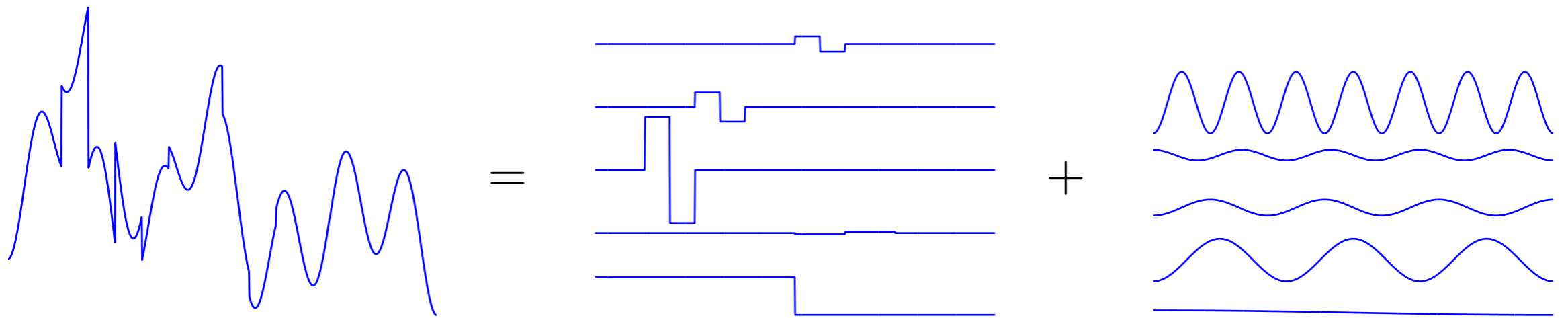
Haar

$$y = Dx$$

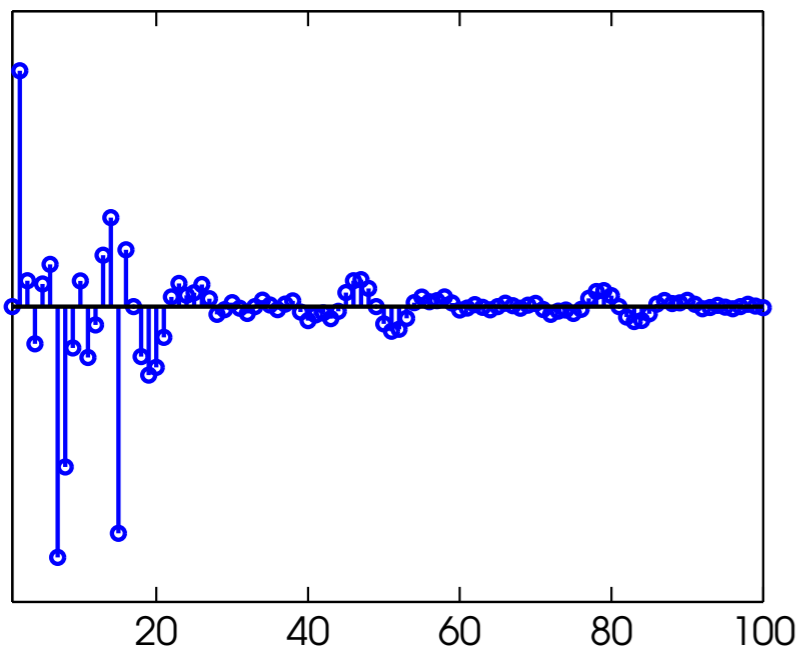


DCT

Sparse representations

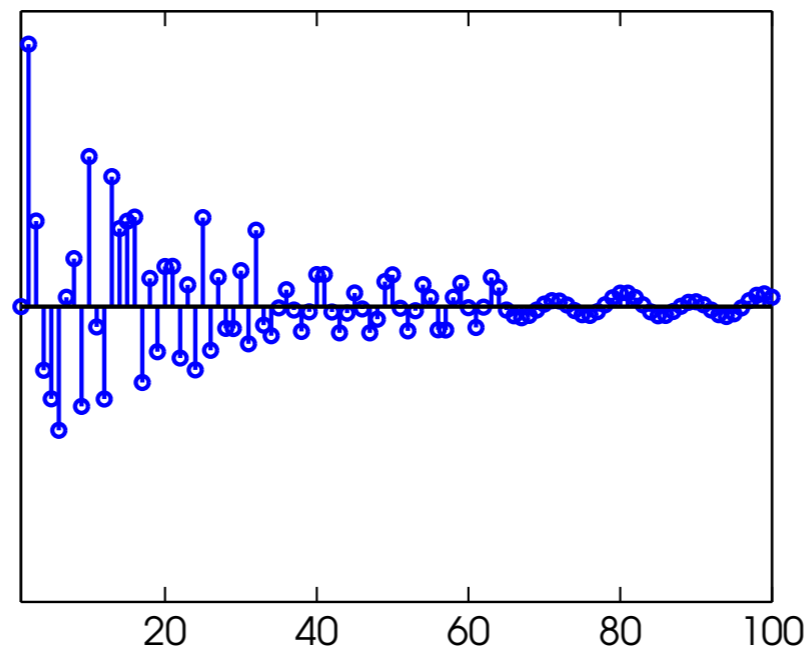


$$y = Hx$$



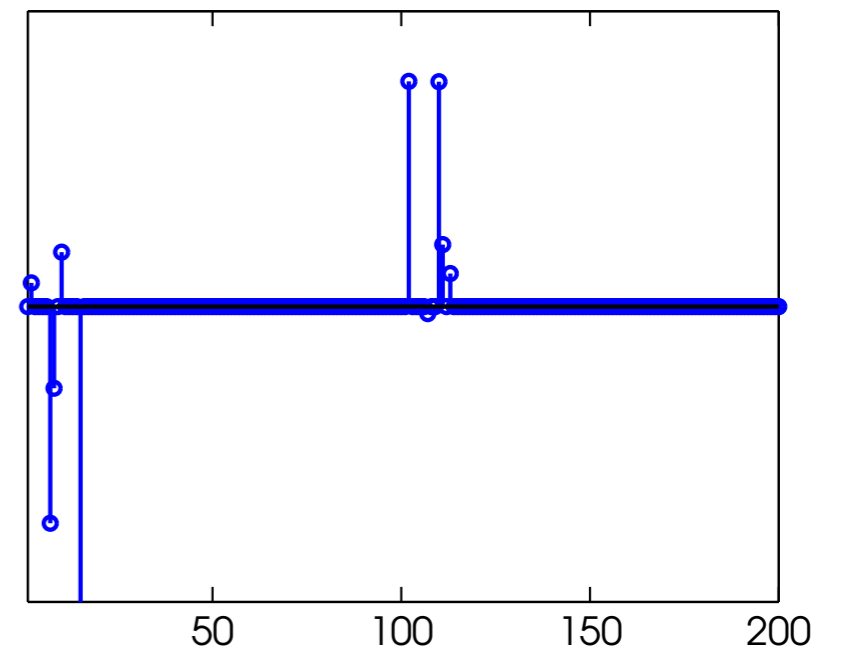
Haar

$$y = Dx$$



DCT

$$y = [H \quad D] x$$



Haar/DCT dictionary

Sparsity via a transform

- represent a signal y as a superposition of elementary signal atoms:

$$y = \sum_j \phi_j x_j \quad \text{ie,} \quad y = \Phi x$$

Sparsity via a transform

- represent a signal y as a superposition of elementary signal atoms:

$$y = \sum_j \phi_j x_j \quad \text{ie,} \quad y = \Phi x$$

- orthogonal bases Φ have unique representation $x = \Phi^T y$

Sparsity via a transform

- represent a signal y as a superposition of elementary signal atoms:

$$y = \sum_j \phi_j x_j \quad \text{ie,} \quad y = \Phi x$$

- orthogonal bases Φ have unique representation $x = \Phi^T y$
- overcomplete dictionaries, eg, $A = [\Phi_1 \ \Phi_2]$ have more flexibility, but representation is not unique. One approach:

$$\text{minimize } \text{nnz}(x) \quad \text{subj to } Ax \approx y$$

Sparsity via a transform

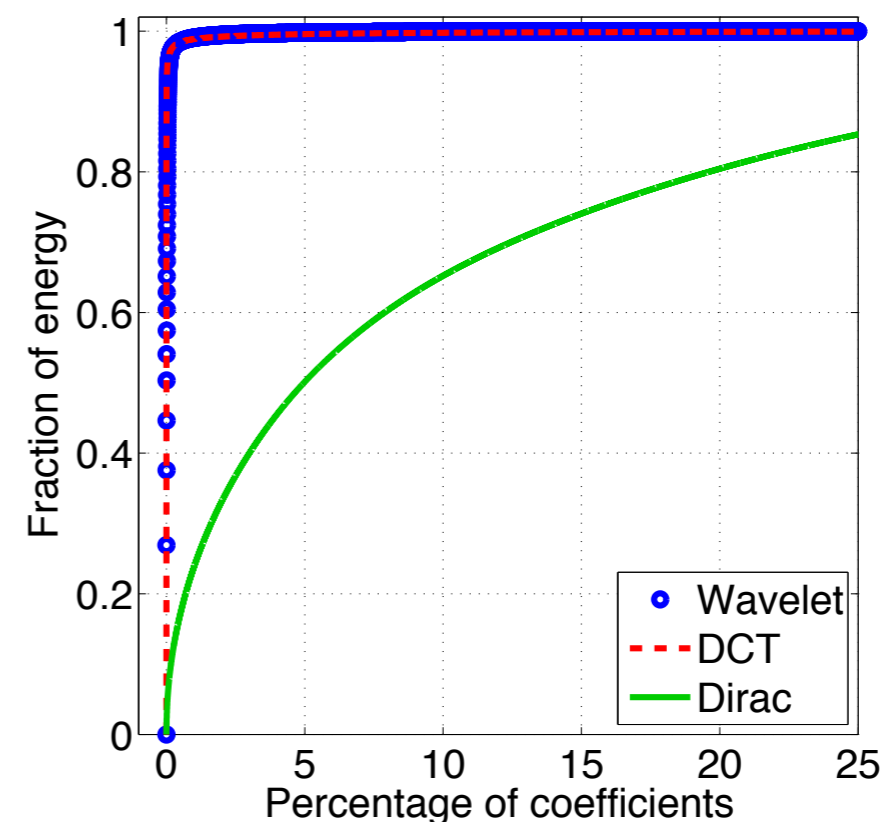
- represent a signal y as a superposition of elementary signal atoms:

$$y = \sum_j \phi_j x_j \quad \text{ie,} \quad y = \Phi x$$

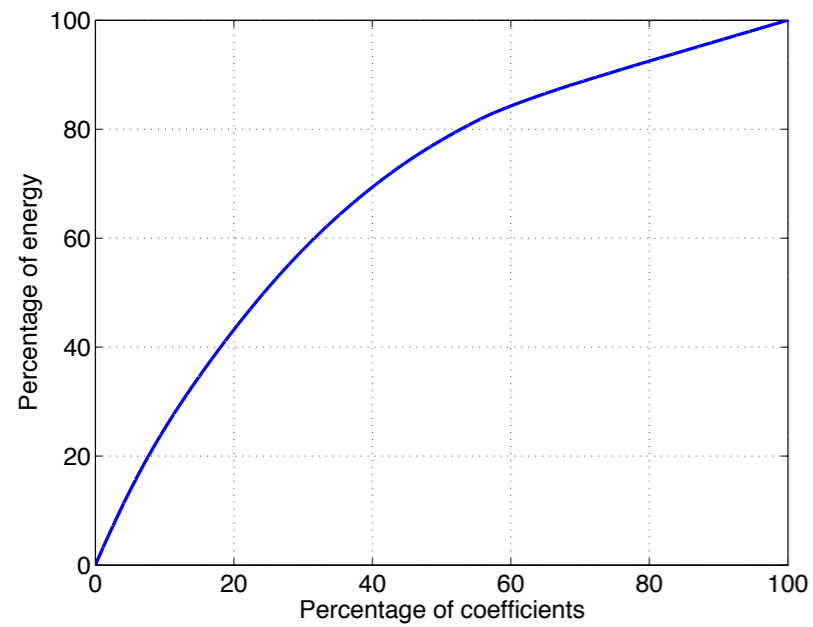
- orthogonal bases Φ have unique representation $x = \Phi^T y$
- overcomplete dictionaries, eg, $A = [\Phi_1 \ \Phi_2]$ have more flexibility, but representation is not unique. One approach:

minimize $\text{nnz}(x)$ subj to $Ax \approx y$

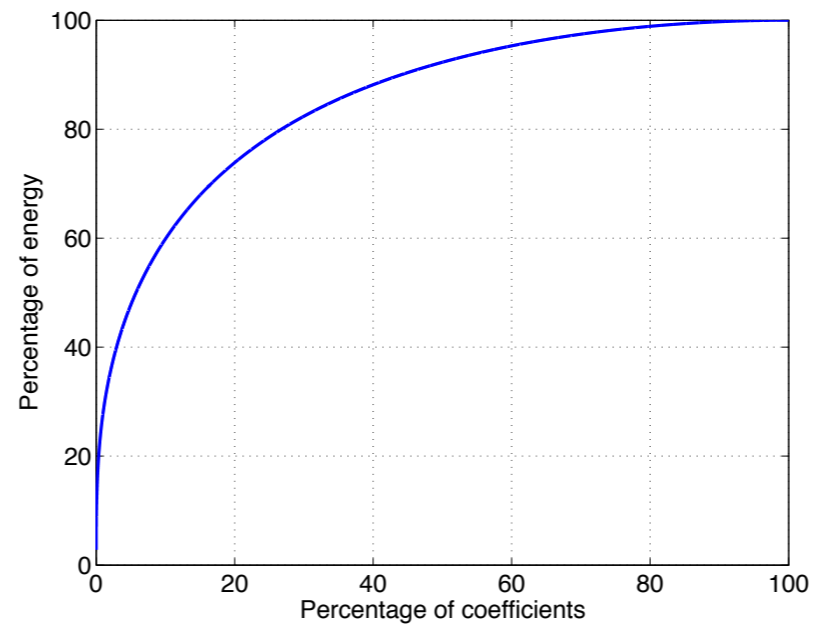
- “efficient” transforms lead to fast decay in coefficients x_j



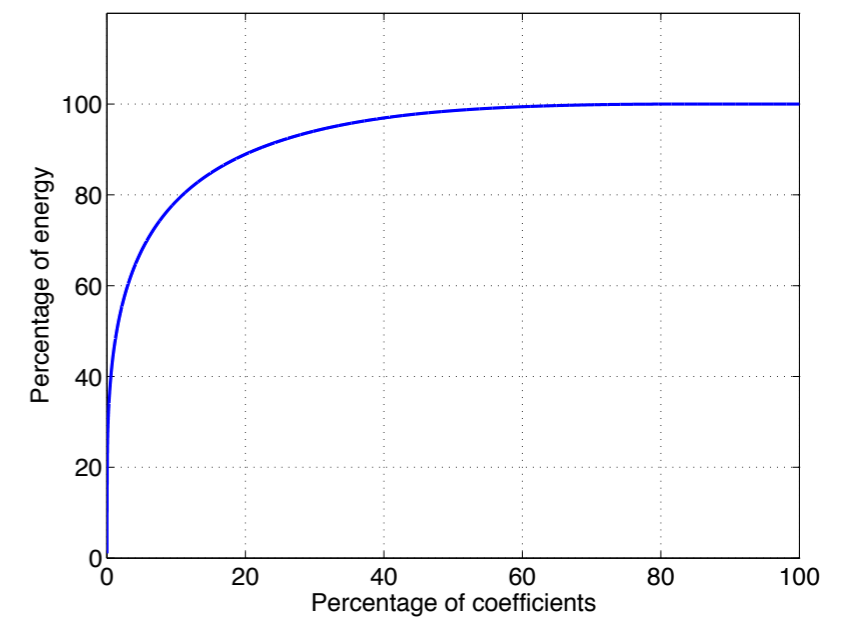
Approximate sparsity



Identity

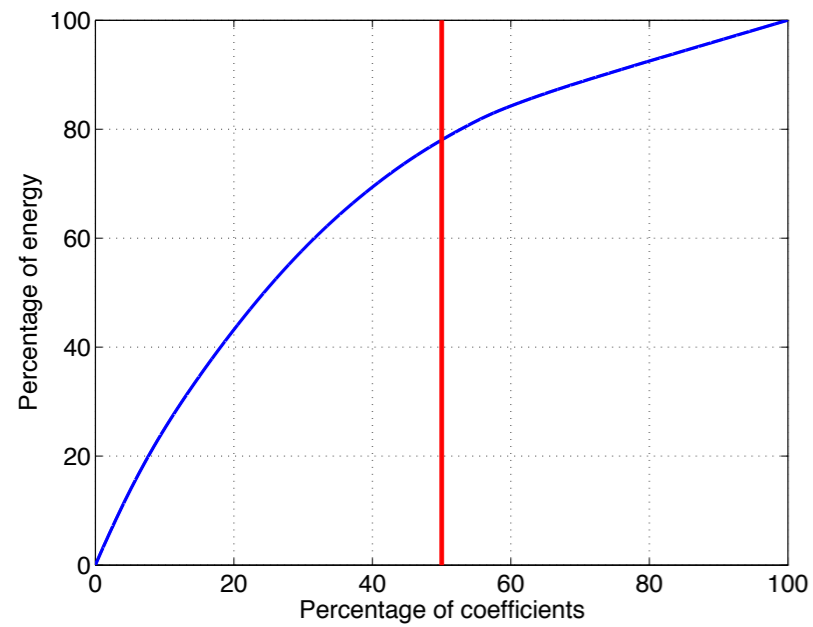


Fourier

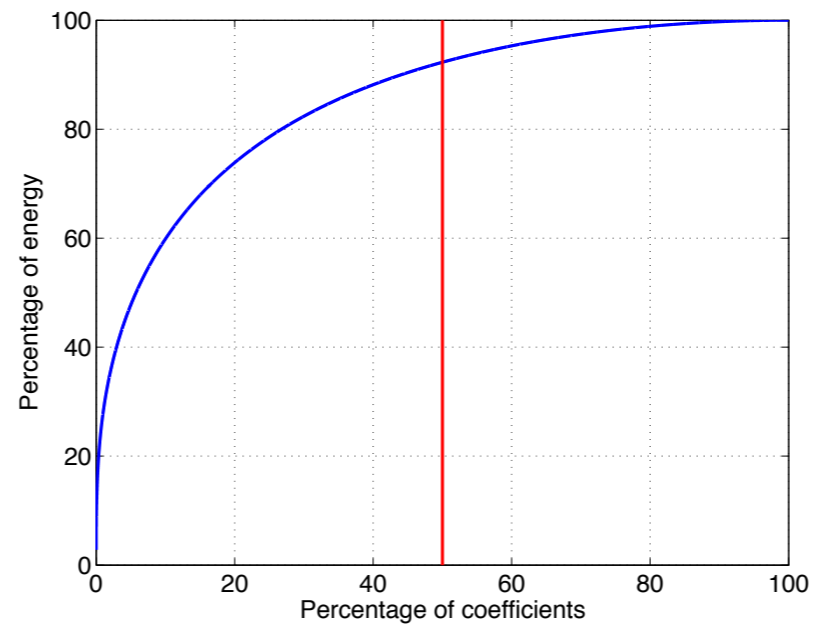


Wavelet

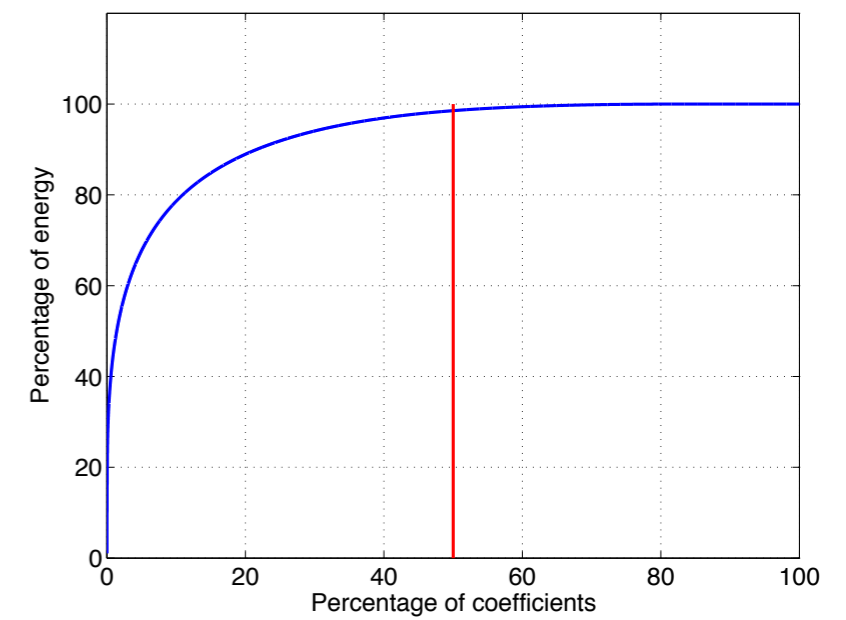
Approximate sparsity



Identity

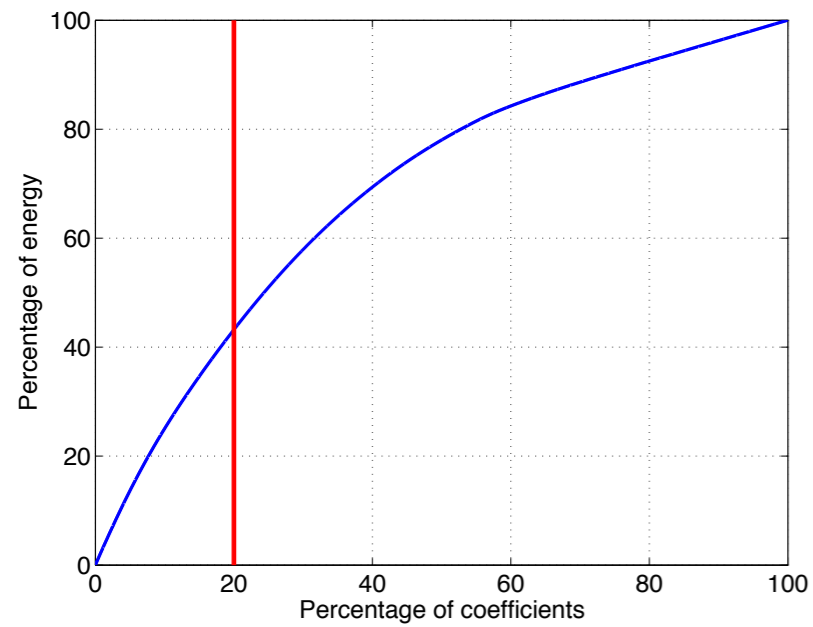


Fourier

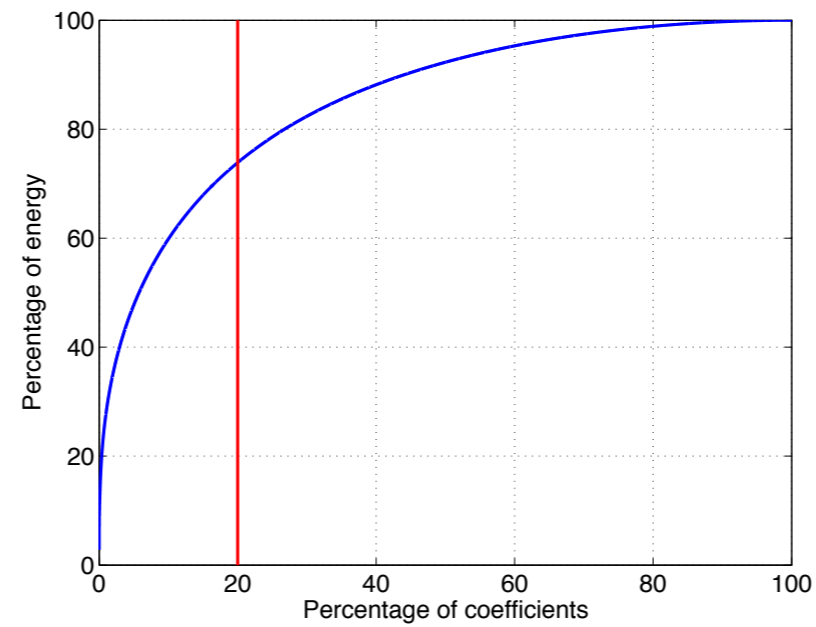


Wavelet

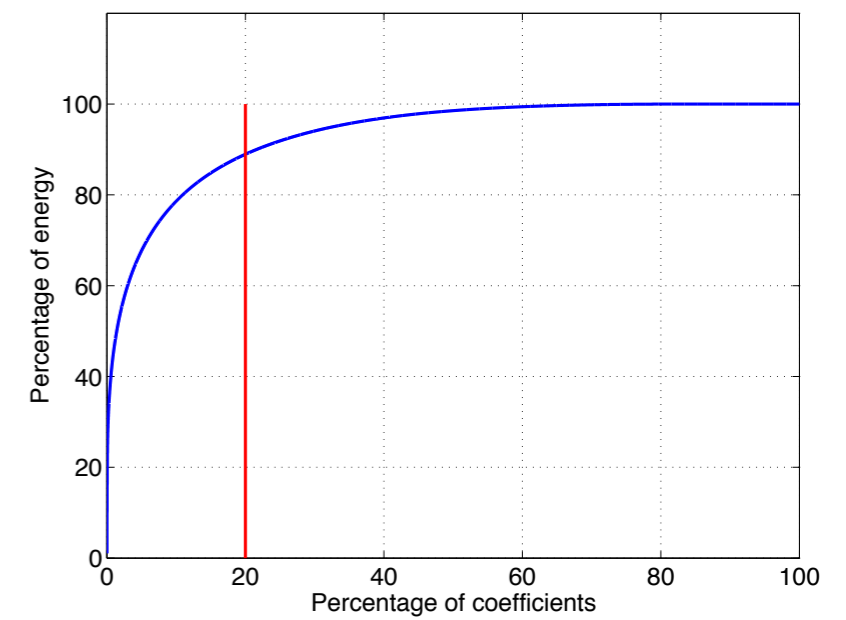
Approximate sparsity



Identity

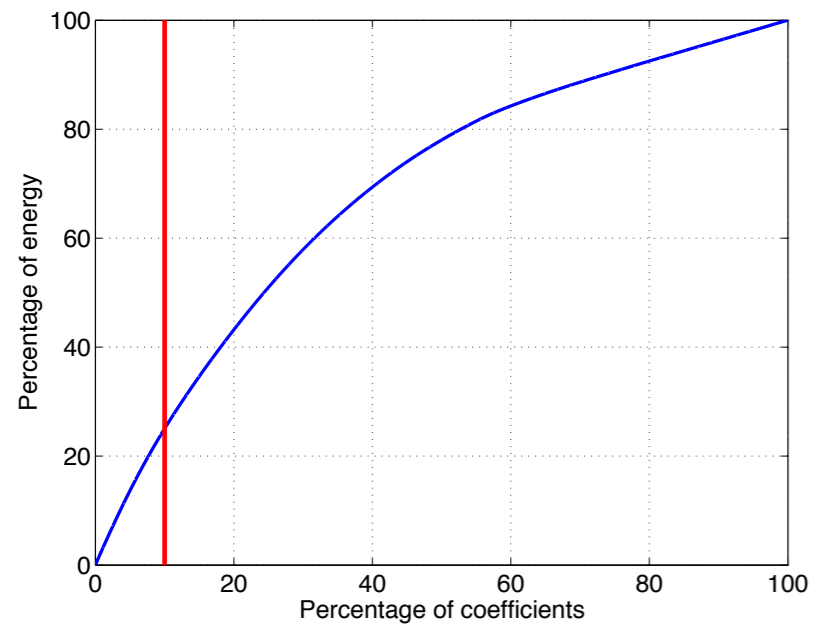


Fourier

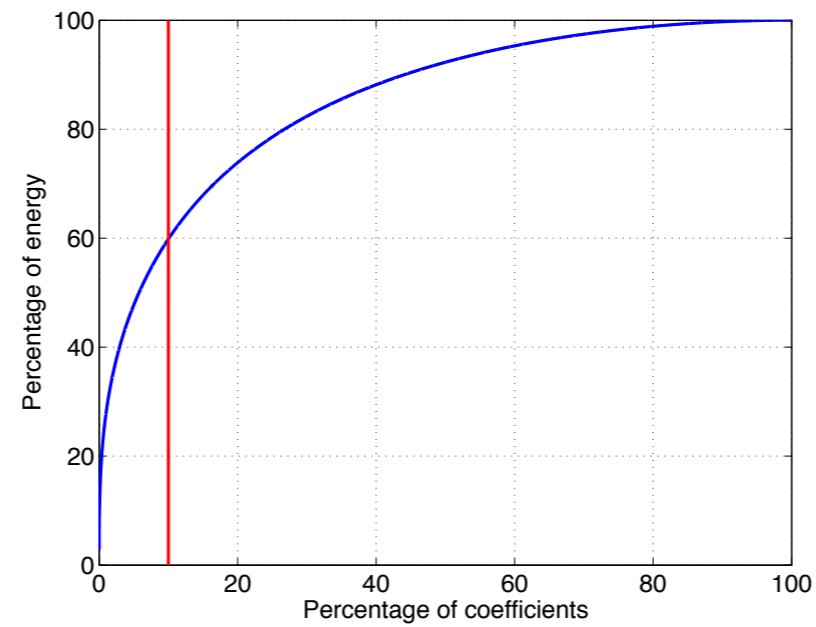


Wavelet

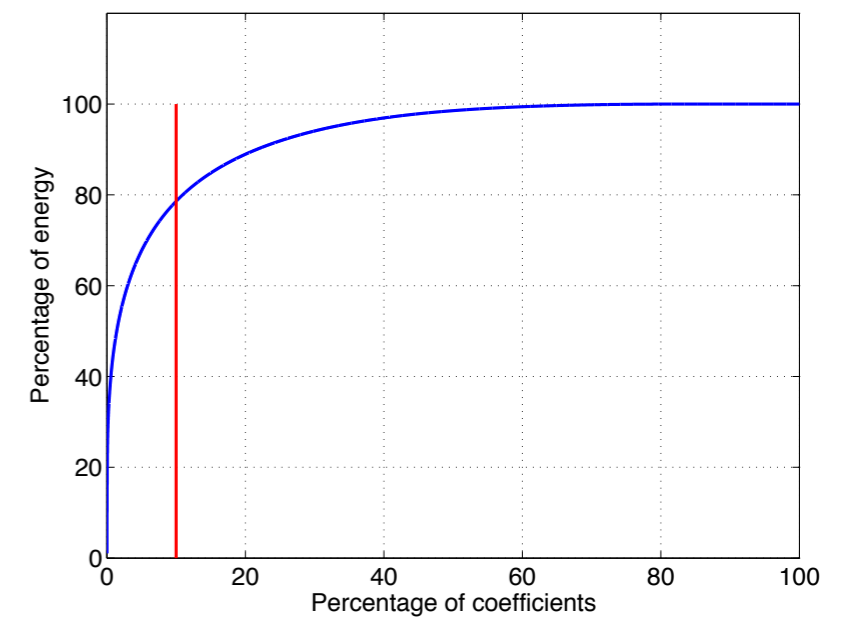
Approximate sparsity



Identity

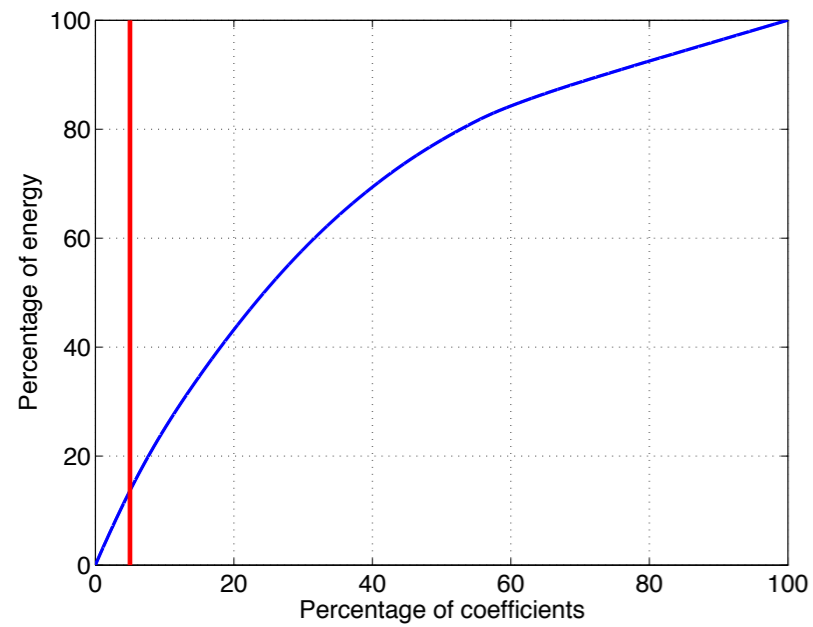


Fourier

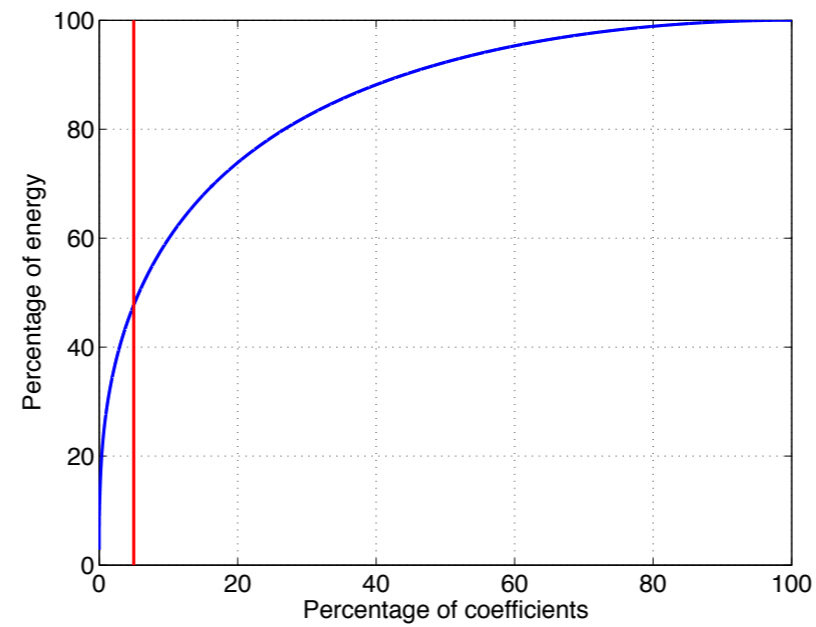


Wavelet

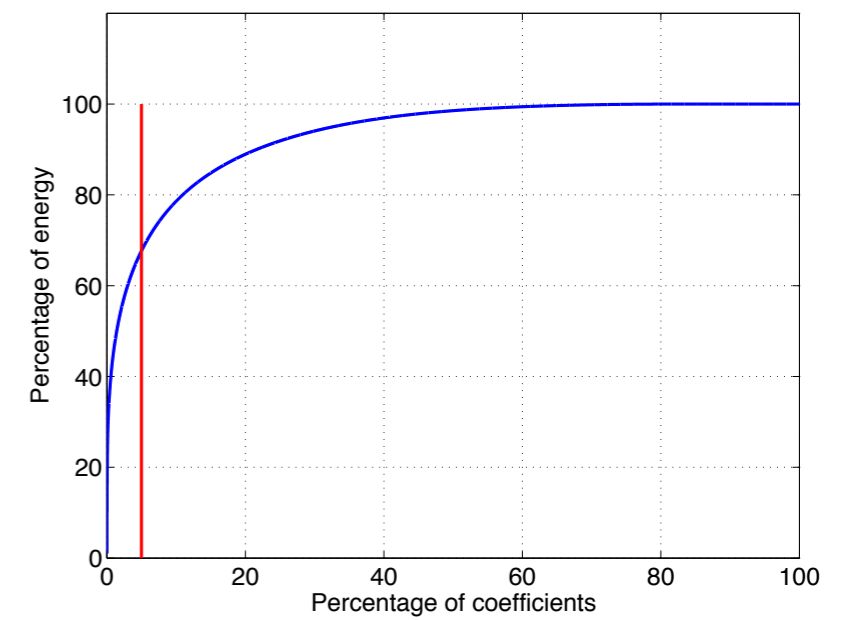
Approximate sparsity



Identity

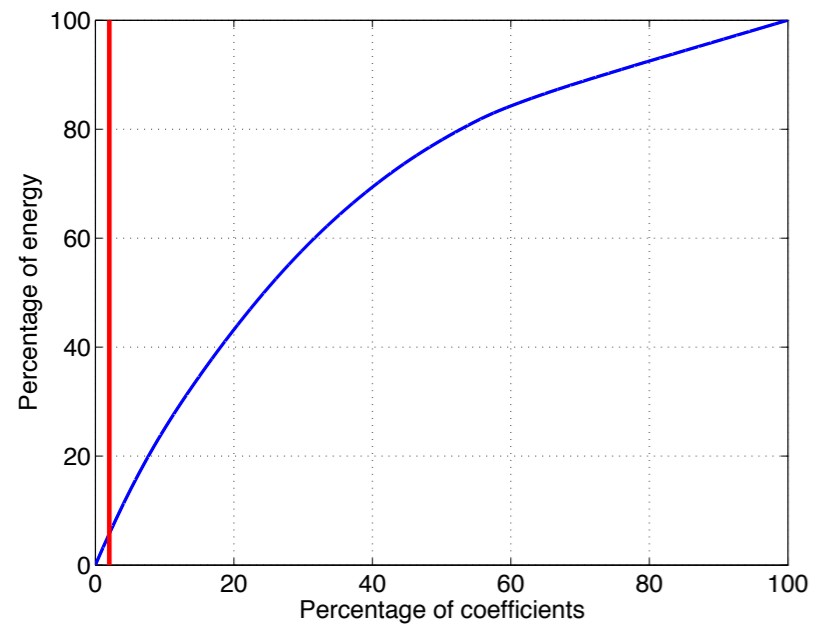


Fourier

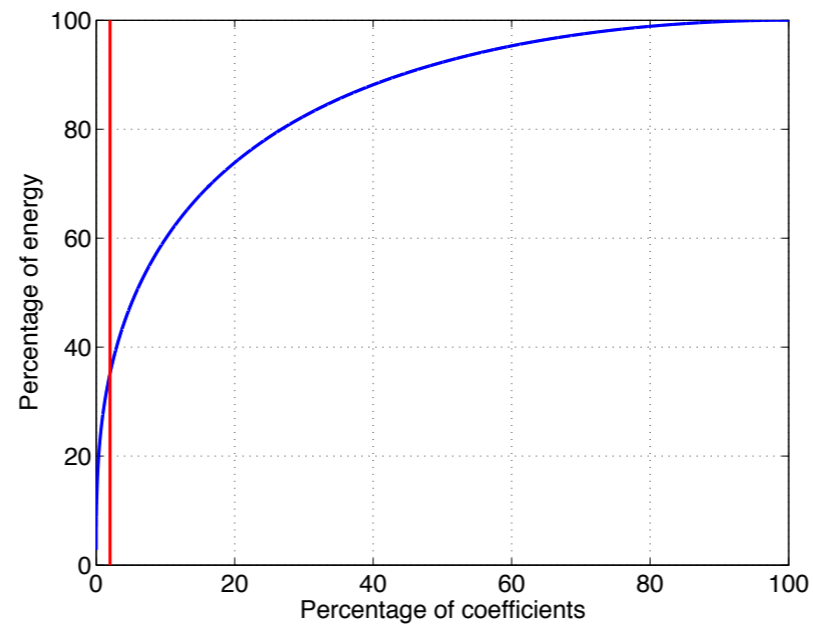


Wavelet

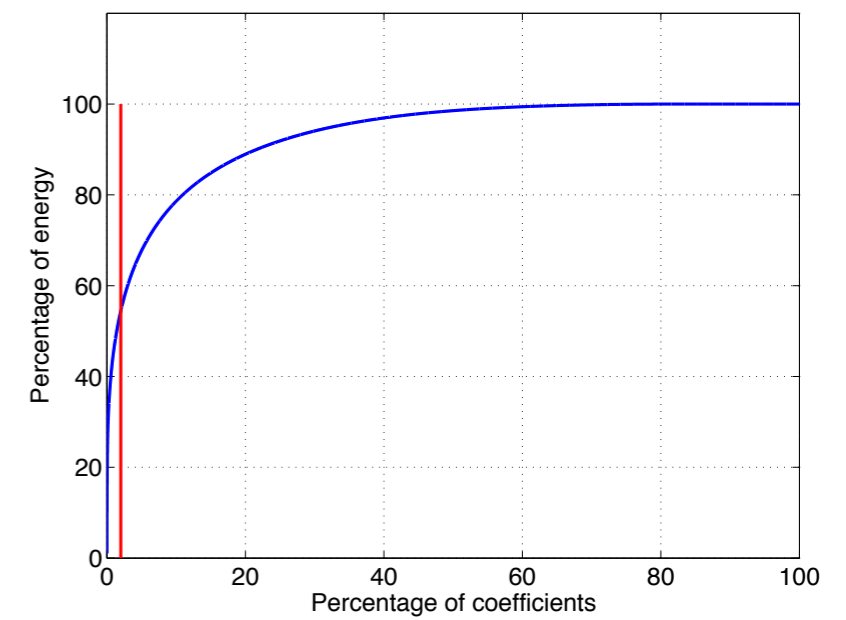
Approximate sparsity



Identity



Fourier



Wavelet

Approaches

Greedy algorithms choose nonzero x_j

- *one by one*: OMP $\min \|S_k x - b\|_2, \quad S_{k+1} \leftarrow [S_k \ a_j]$
- *batches*: CoSaMP $\min \|S_k x - b\|_2, \quad S_{k+1} \leftarrow [A_{\{s \text{ elements}\}}]$

[Tropp-Gilbert 07; Tropp-Needell 08]

Convex relaxations

- minimize $\|x\|_1$ subj to $Ax \approx b$

[Donoho 05; Candès-Tao 06; Candès 06]

Nonconvex heuristics

- minimize $\|x\|_p$ subj to $Ax \approx b, \quad p < 1$

[Saab-Yilmaz 08; Chartrand 07]

A sample of solvers

$$\text{minimize } \|x\|_1 \quad \text{subj to } \|Ax - b\|_2 \leq \sigma$$

ℓ_1 -Magic Candès & Romberg 05
SPGL1 van den Berg & Friedlander 07
Homotopy Osborne, Presnell & Turlach 00

interior, CG
root-finding
active-set, all λ

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

PDCO Saunders 97, 02
L1_LS Kim, Koh, et al 07
BCR Sardy, Bruce & Tseng 98
GPSR Figueiredo, Nowak & Wright 07
FPC Hale, Yin & Zhang 07
SpaRSA Wright, Figueiredo & Nowak 07

interior, LSQR
primal barrier, CG
block Gauss-Seidel
gradient projection
operator splitting
shrinkage

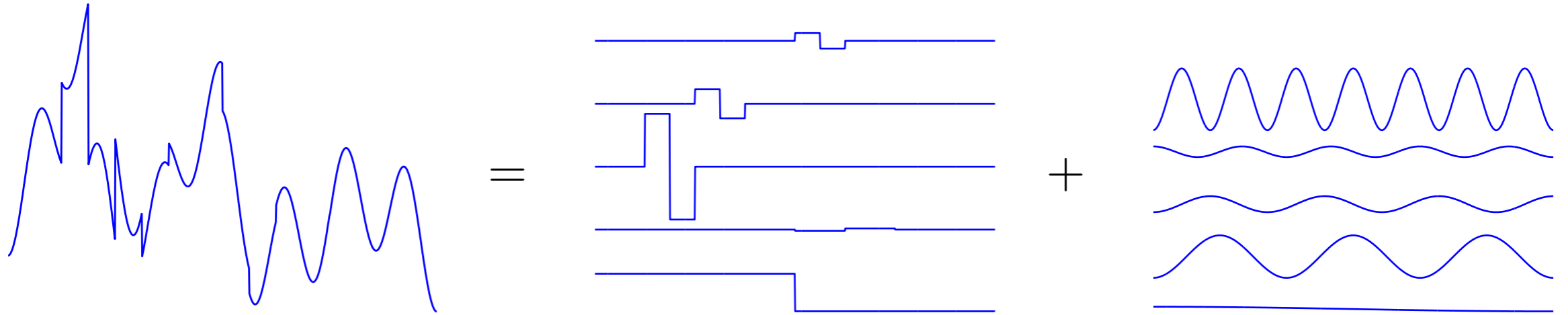
$$\text{Sparse } Ax \approx b$$

OMP Pati...93, Davis...97; Tropp & Gilbert 07
LARS Efron, Hastie & Tibshirani 04
StOMP Donoho, Tsaig, et al 06
ROMP Needell & Vershynin 07
CoSaMP Needell & Tropp 08

greedy
greedy active-set
greedy++
greedy+/-
greedy+/-

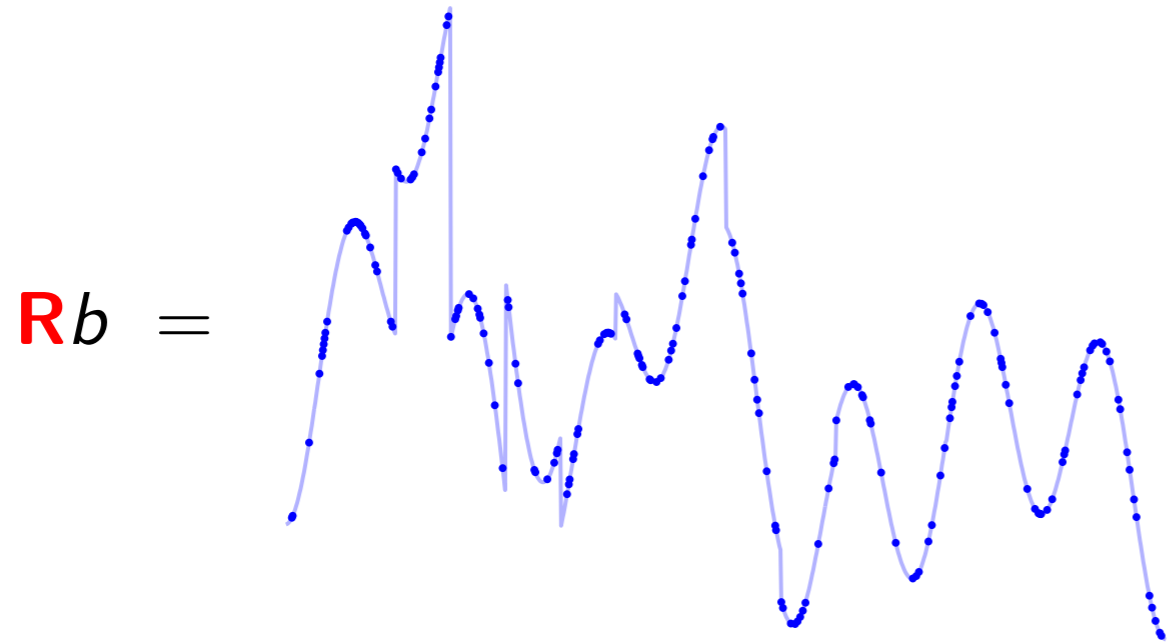
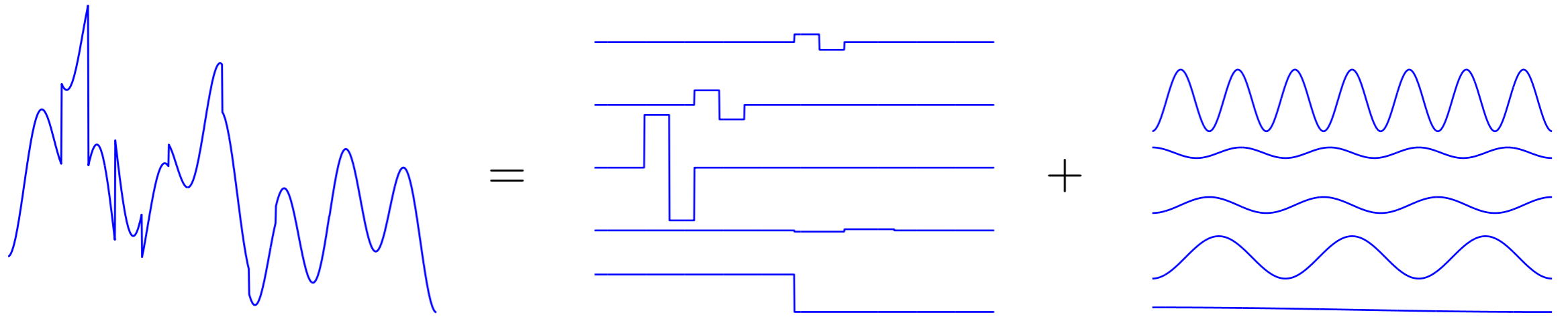
APPLICATIONS and FORMULATIONS

Missing observations



$$y \approx \begin{array}{|c|c|} \hline \text{Haar} & \text{DCT} \\ \hline \end{array} x$$

Missing observations



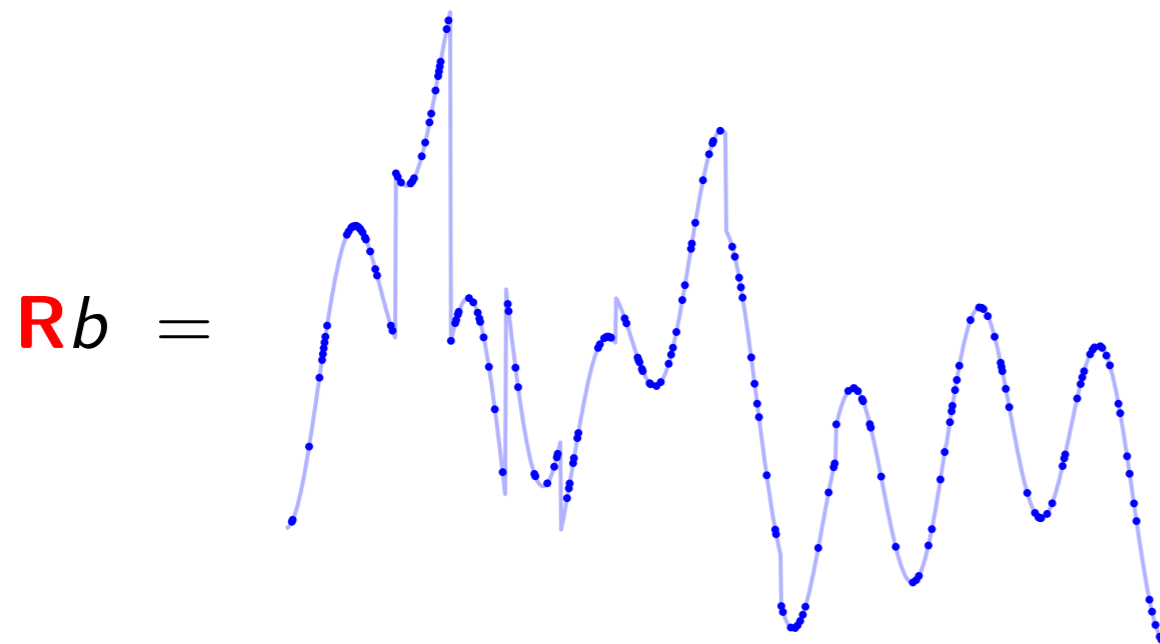
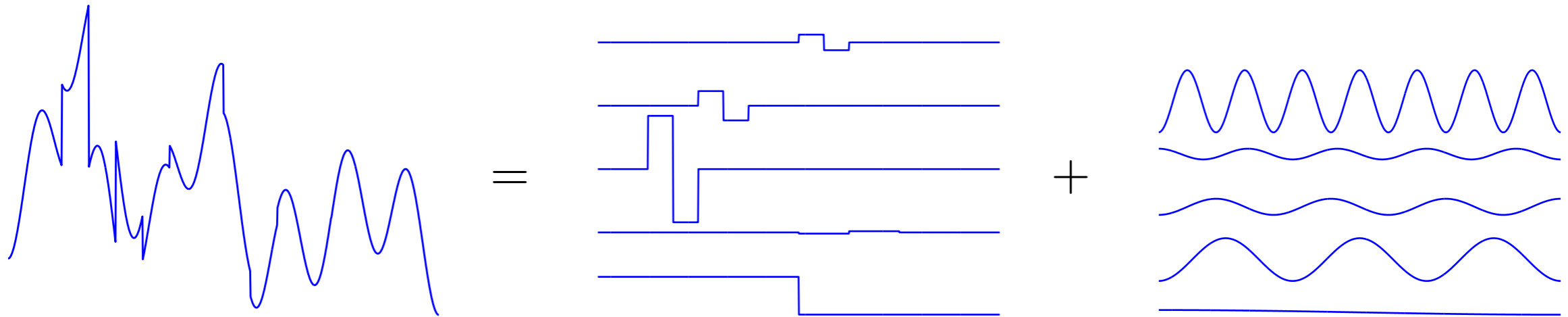
$R =$ Restriction

$Ry \approx R$

Haar	DCT
------	-----

 x

Missing observations

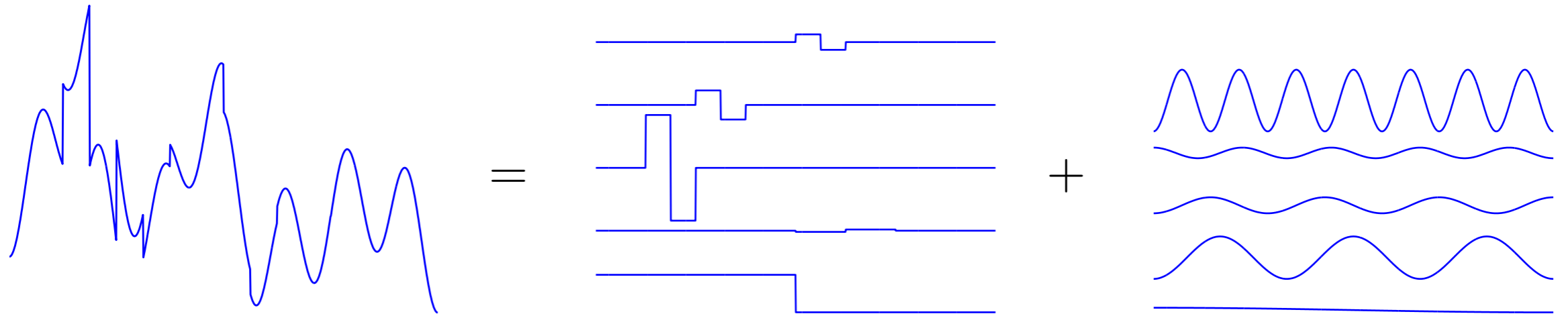


$R =$ Restriction

$$Ry \approx R \begin{bmatrix} \text{Haar} & \text{DCT} \end{bmatrix} x \Rightarrow b \approx \begin{bmatrix} A \end{bmatrix} x$$

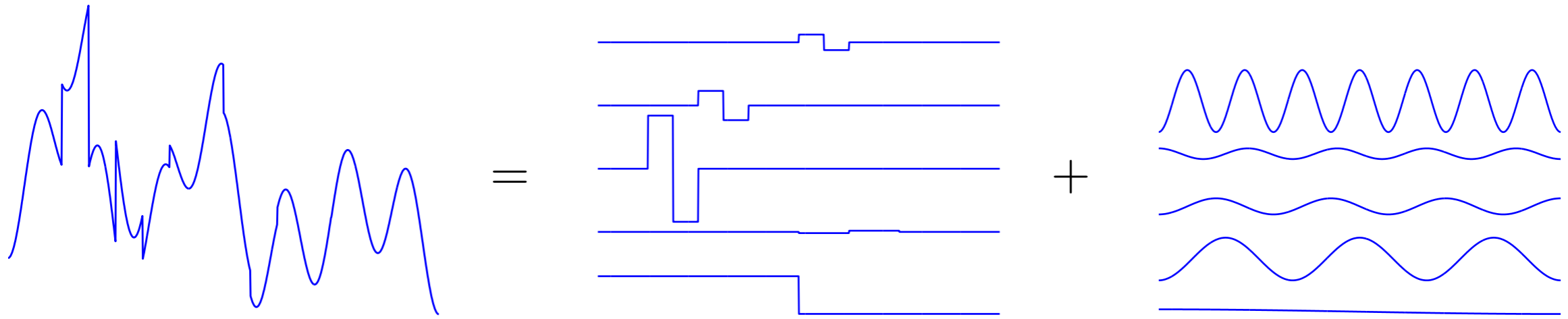
minimize $\frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1$

Compressed sensing



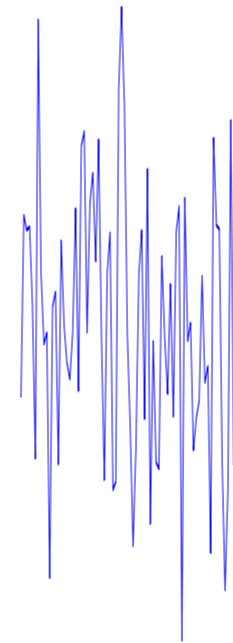
$$y \approx \begin{array}{|c|c|} \hline \text{Haar} & \text{DCT} \\ \hline \end{array} x$$

Compressed sensing



$$\langle \phi_1, y \rangle, \dots, \langle \phi_m, y \rangle$$

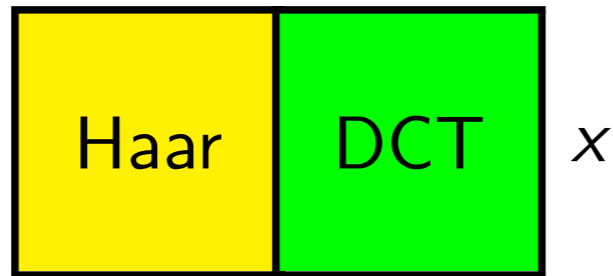
$$\Phi y =$$



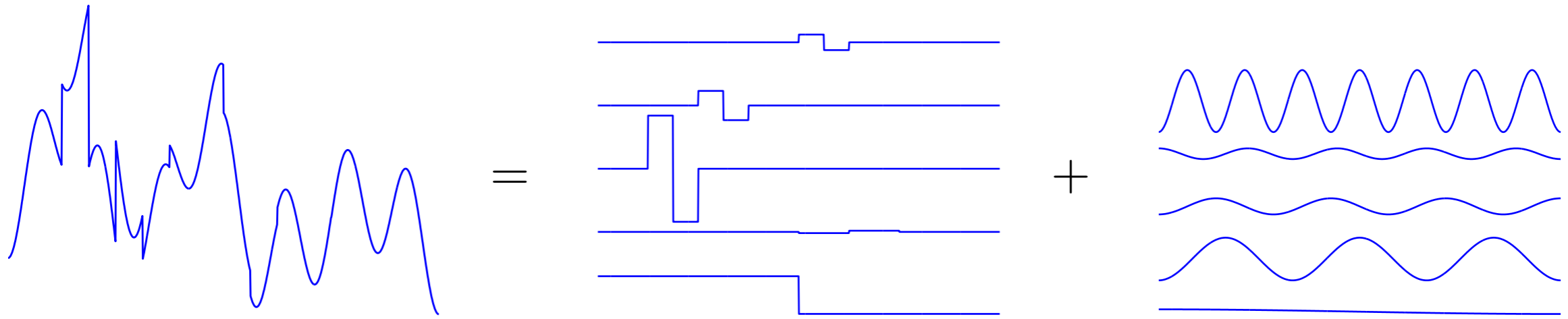
$$\Phi =$$

Gaussian

$$\Phi y \approx \Phi$$

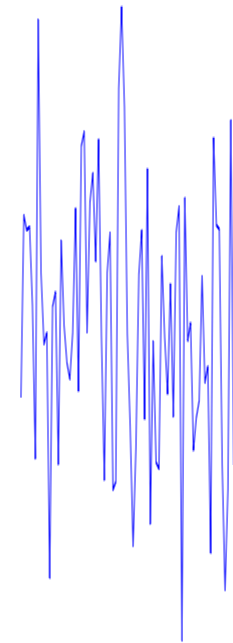


Compressed sensing



$$\langle \phi_1, y \rangle, \dots, \langle \phi_m, y \rangle$$

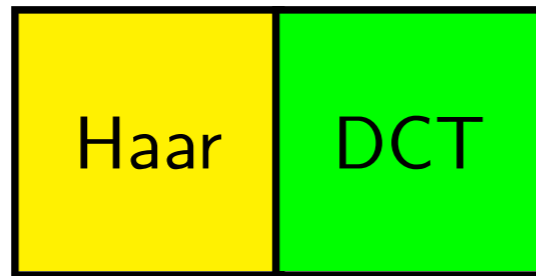
$$\Phi y =$$



$$\Phi =$$

Gaussian

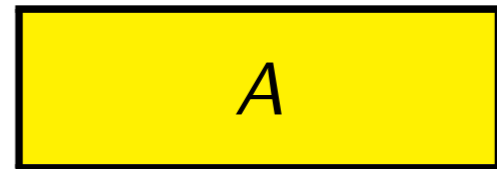
$$\Phi y \approx \Phi$$



x

\Rightarrow

$$b \approx$$



x

$$\text{minimize } \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1$$

Image deblurring

Observe $b = My$, $y = \text{image}$, $M = \text{blurring operator}$



Image deblurring

Observe $b = My$, $y = \text{image}$, $M = \text{blurring operator}$

Recover significant coeff's via

$$\text{minimize } \frac{1}{2} \|Mx - b\|^2 + \lambda \|x\|_1$$



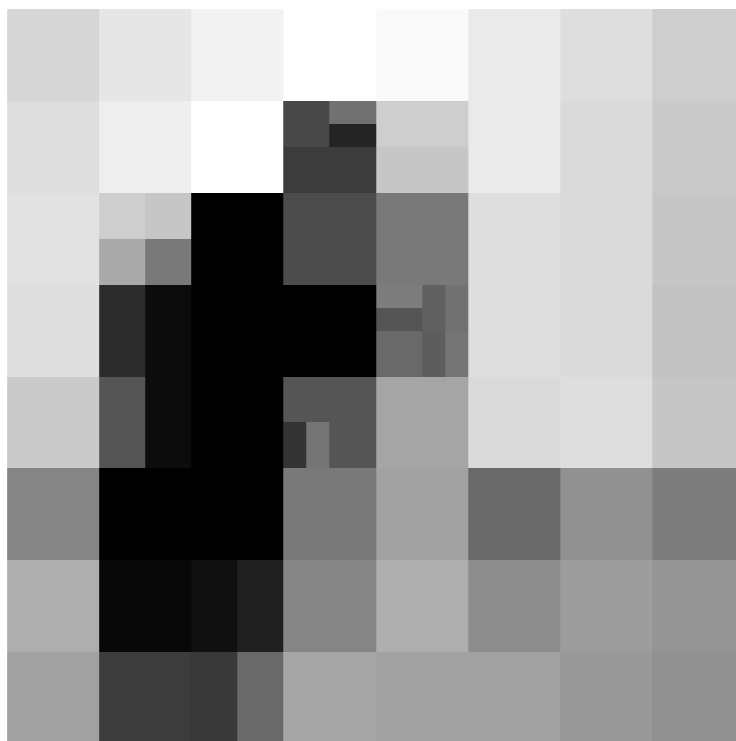
Image deblurring



Observe $b = My$, $y = \text{image}$, $M = \text{blurring operator}$

Recover significant coeff's via

$$\text{minimize } \frac{1}{2} \|Mx - b\|^2 + \lambda \|x\|_1$$



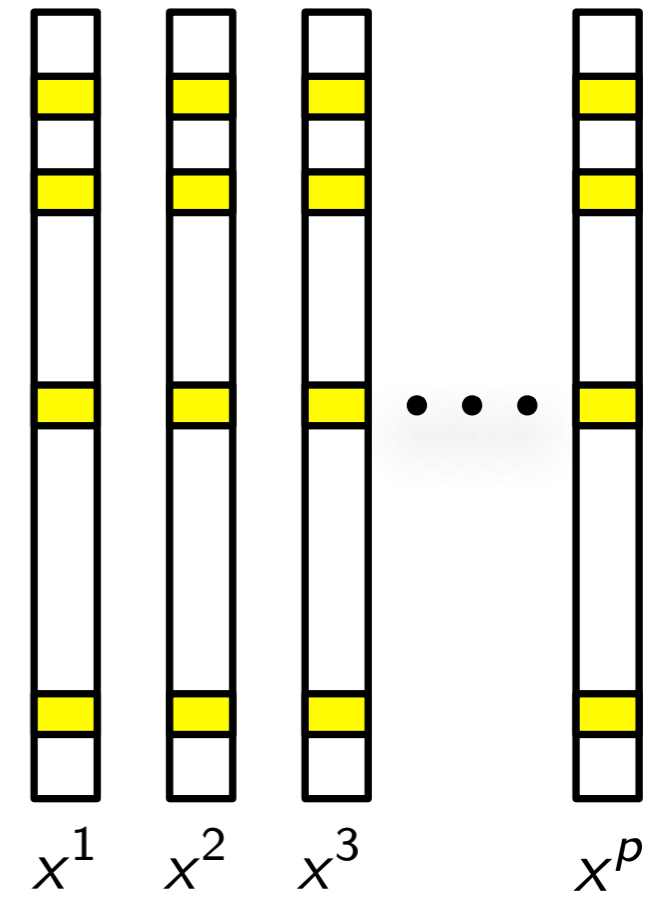
λ	2.3e+0	2.3e-2	2.3e-4
$\ Ax - b\ $	2.6e+2	3.0e+0	3.1e-2
$\text{nnz}(x)$	72	1,741	28,484

Sparco problem [blurrycam](#) [Figueiredo-Nowak-Wright 07]

Sum-of-norms and joint sparsity

Multiple measurements

$$Ax^1 \approx b^1, \quad Ax^2 \approx b^2, \quad \dots, \quad Ax^p \approx b^p$$



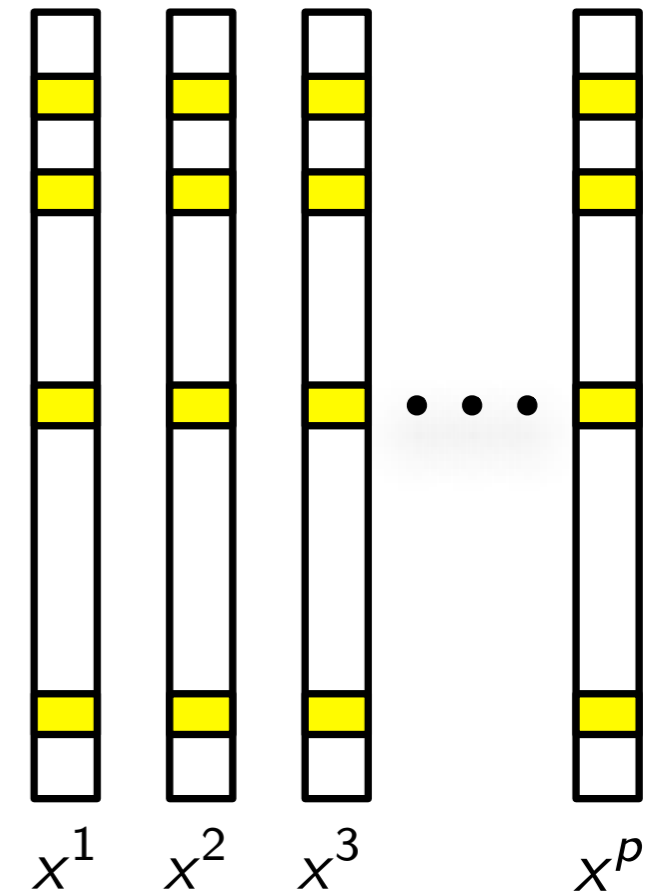
Sum-of-norms and joint sparsity

Multiple measurements

$$Ax^1 \approx b^1, \quad Ax^2 \approx b^2, \quad \dots, \quad Ax^p \approx b^p$$

Minimize sum of row norms

$$\begin{aligned} & \underset{X \in \mathbb{R}^{n \times k}}{\text{minimize}} && \|X\|_{1,2} := \sum \|X^{j \rightarrow}\|_2 \\ & \text{subj to} && \|AX - B\|_F \leq \sigma \end{aligned}$$



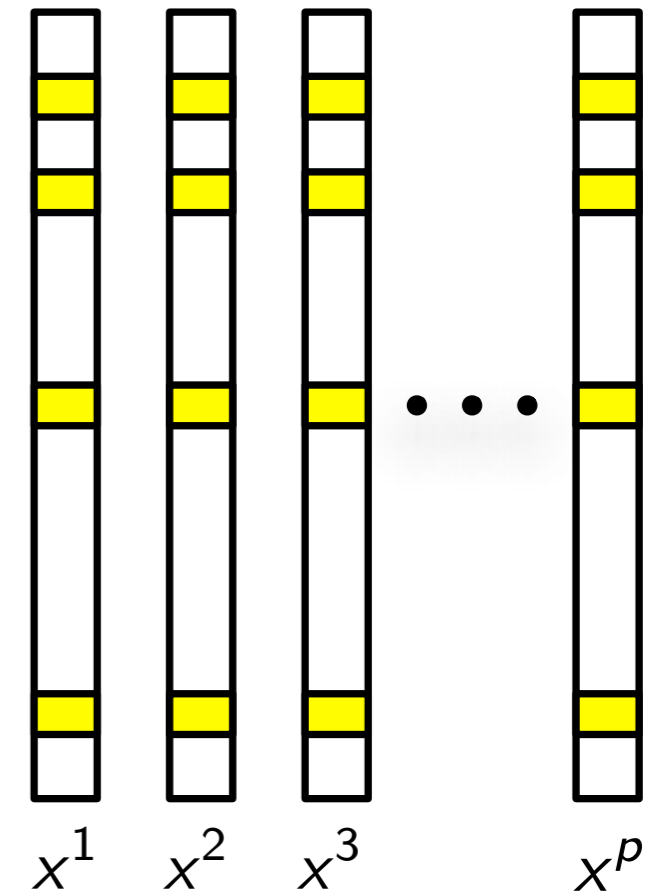
Sum-of-norms and joint sparsity

Multiple measurements

$$Ax^1 \approx b^1, \quad Ax^2 \approx b^2, \quad \dots, \quad Ax^p \approx b^p$$

Minimize sum of row norms

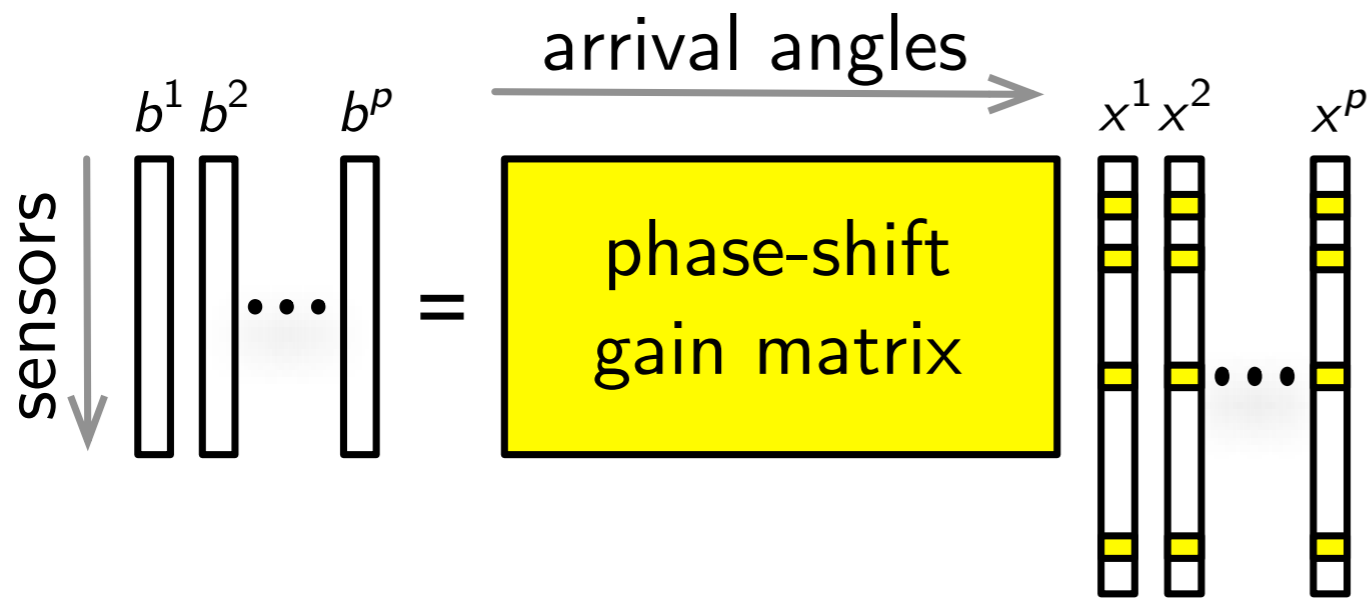
$$\begin{aligned} & \underset{X \in \mathbb{R}^{n \times k}}{\text{minimize}} && \|X\|_{1,2} := \sum \|X^{j \rightarrow}\|_2 \\ & \text{subj to} && \|AX - B\|_F \leq \sigma \end{aligned}$$



Projection onto $\{X \mid \|X\|_{1,2} \leq \tau\}$:

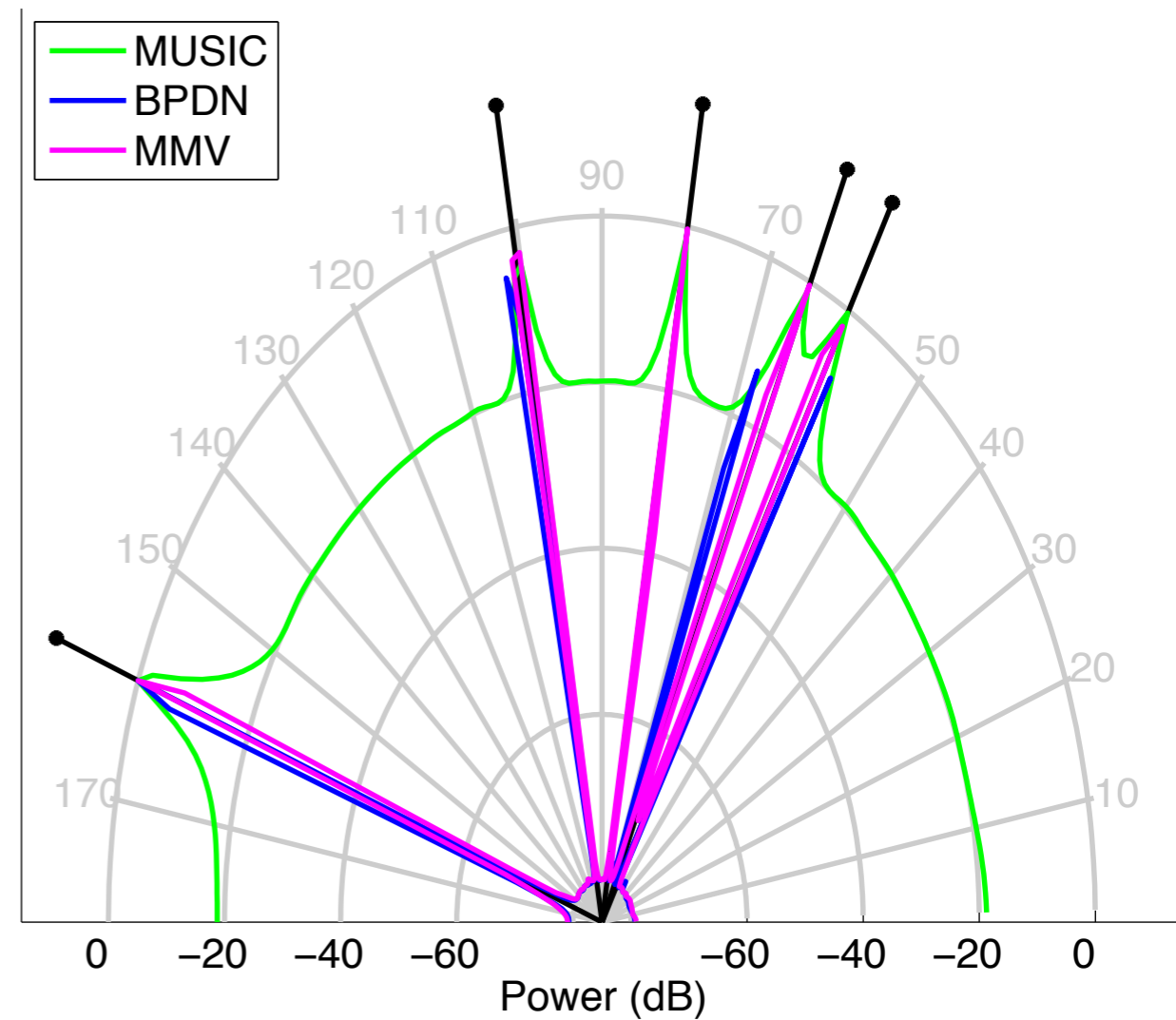
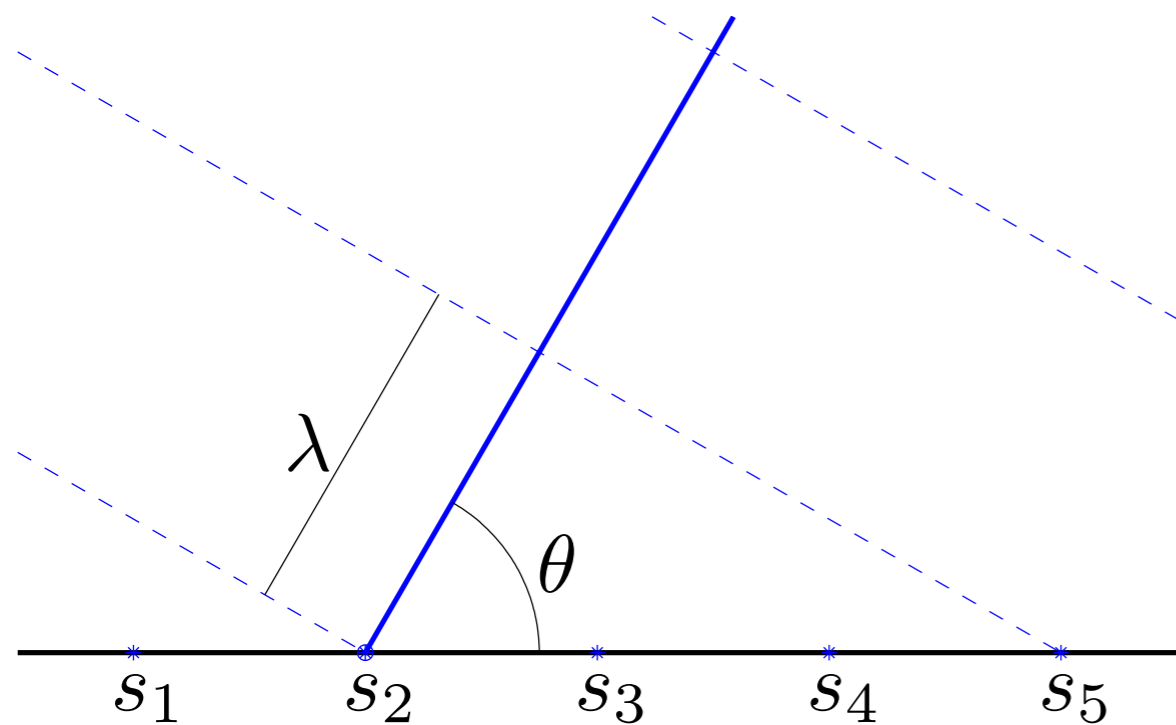
1. compute vector of row norms: $r \leftarrow (\|X^{1 \rightarrow}\|_2, \dots, \|X^{n \rightarrow}\|_2)$
2. **one-norm projection:** $\bar{r} \leftarrow \mathcal{P}_\tau[r]$
3. scale rows by \bar{r} : $\mathcal{P}_\tau[X] \leftarrow \text{diag}(\bar{r}./r) X$

Joint sparsity: source localization

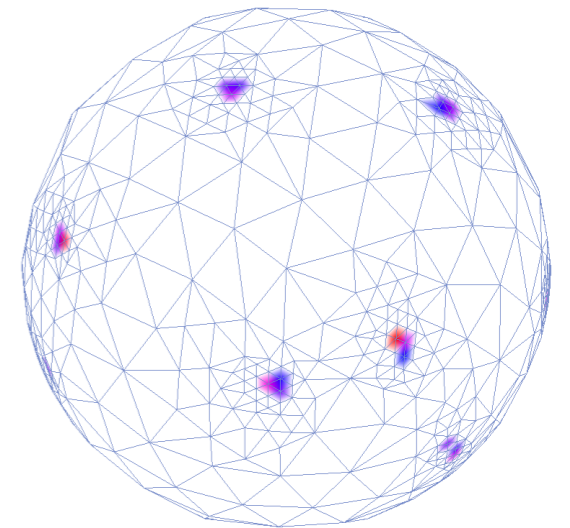
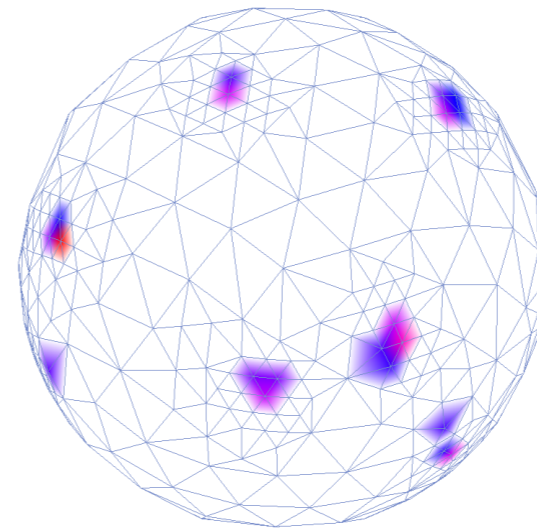
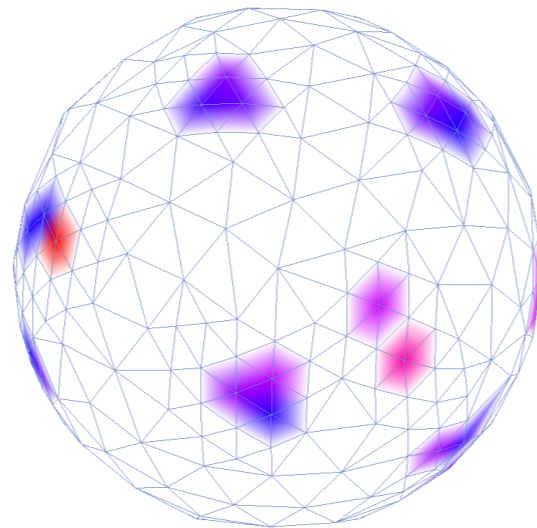
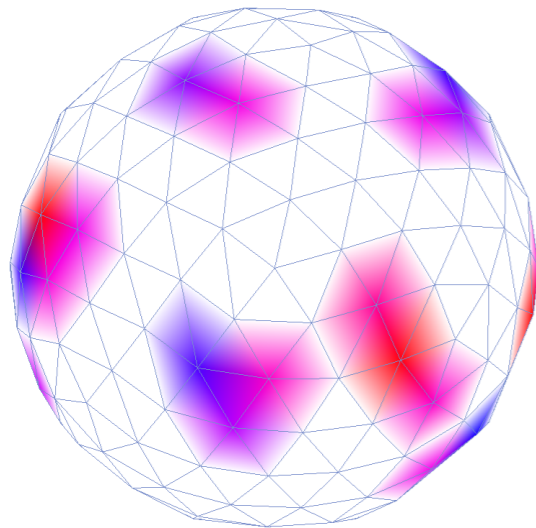
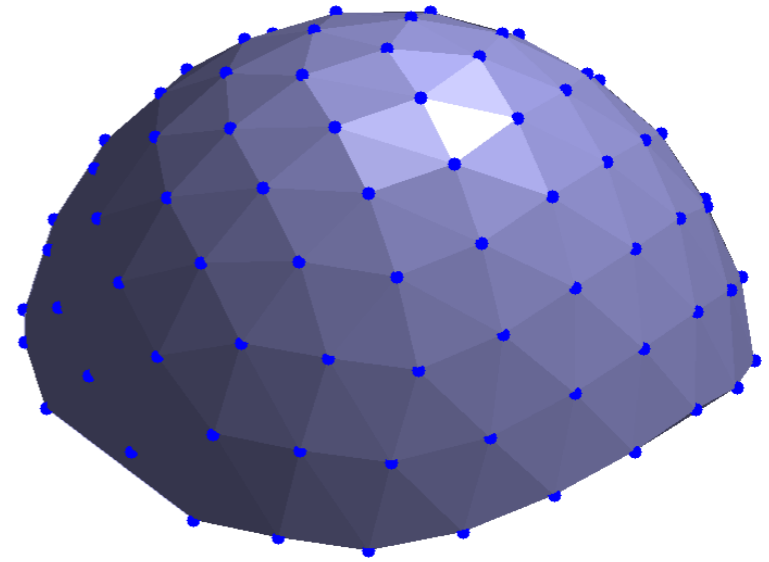
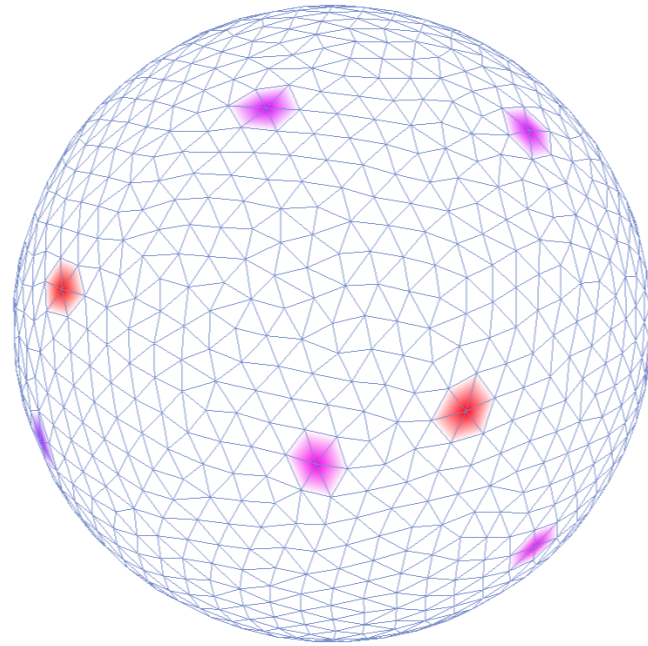
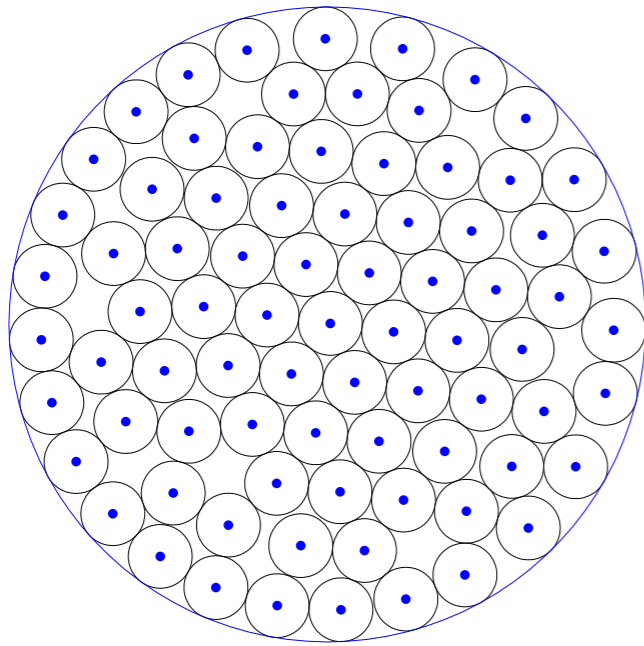


[Malioutov–Çetin–Willsky 03]

minimize $\|X\|_{1,2}$
 subj to $\|B - AX\|_F \leq \sigma$

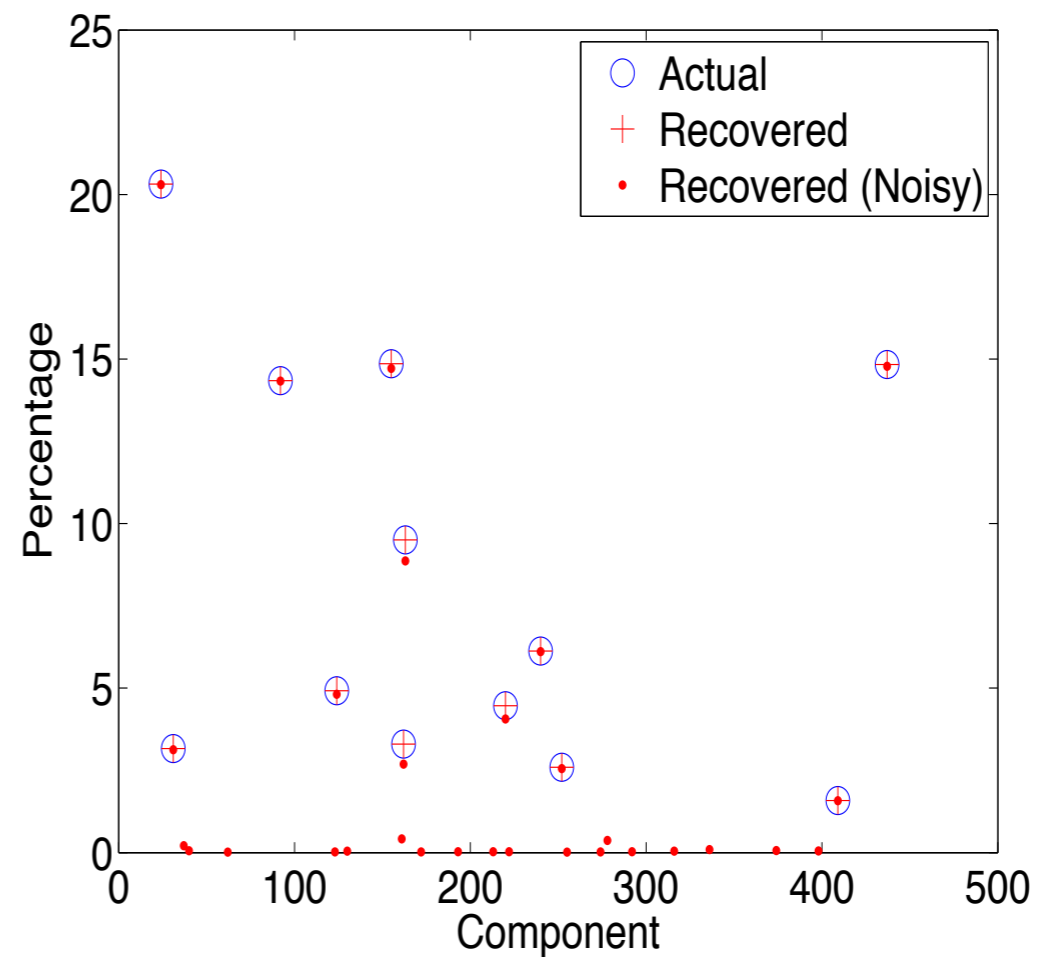
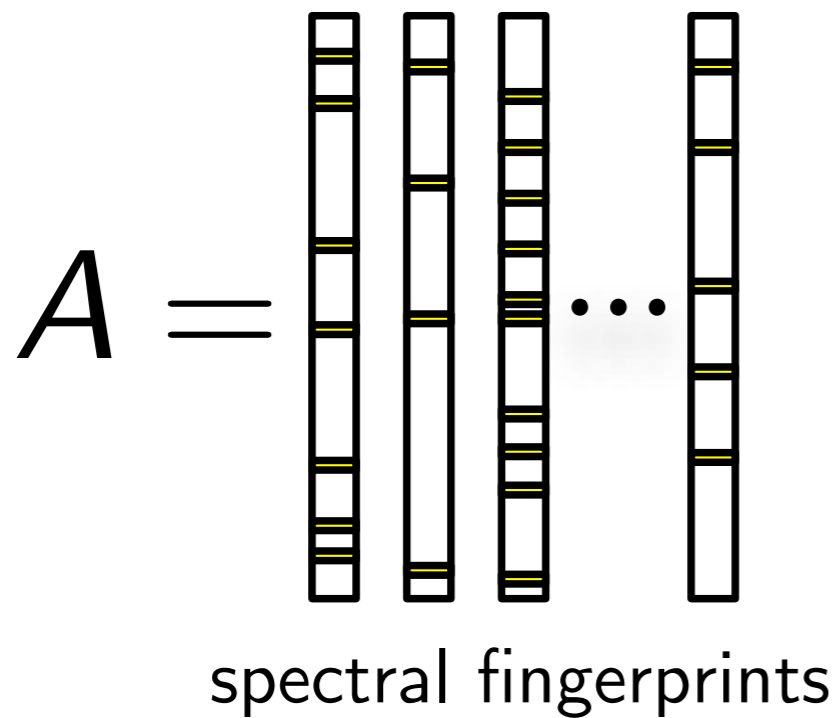
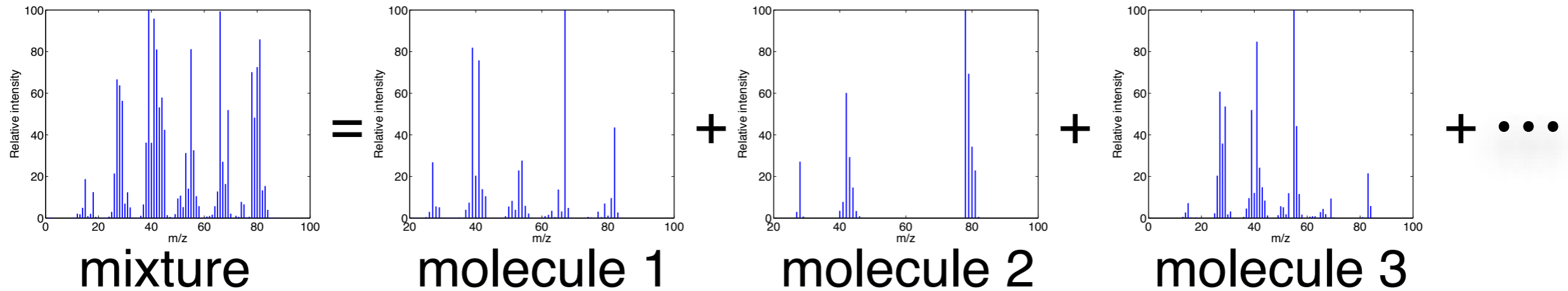


...and in 3D



[van den Berg's thesis 09]

Mass spectrometry – separating mixtures



$$\underset{x \geq 0}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\|^2 + \lambda \underbrace{e^T x}_{\equiv \|x\|_1}$$

[Du-Angeletti 06, vdBerg-Friedlander 10]

SPARSE OPTIMIZATION ALGORITHMS

**greedy algorithms, projected gradient, iterative soft-thresholding,
fixed point, proximal point, augmented Lagrangian, Bregman, split
Bregman, alternating direction of method of multipliers, . . .**

Canonical greedy approach

Orthogonal matching pursuit (OMP) for sparse $Ax = b$

Canonical greedy approach

Orthogonal matching pursuit (OMP) for sparse $Ax = b$

0. Initialize $r \leftarrow b, S \leftarrow []$

Canonical greedy approach

Orthogonal matching pursuit (OMP) for sparse $Ax = b$

0. Initialize $r \leftarrow b, S \leftarrow []$
1. Largest correlation find j st $|a_j^T r| = \|A^T r\|_\infty$

Canonical greedy approach

Orthogonal matching pursuit (OMP) for sparse $Ax = b$

0. Initialize $r \leftarrow b, \quad S \leftarrow []$
1. Largest correlation find j st $|a_j^T r| = \|A^T r\|_\infty$
2. Get new column $S \leftarrow [S \quad a_j]$

Canonical greedy approach

Orthogonal matching pursuit (OMP) for sparse $Ax = b$

0. Initialize $r \leftarrow b, \quad S \leftarrow []$
1. Largest correlation find j st $|a_j^T r| = \|A^T r\|_\infty$
2. Get new column $S \leftarrow [S \quad a_j]$
3. Least squares $\min \|b - Sx\|_2$

Canonical greedy approach

Orthogonal matching pursuit (OMP) for sparse $Ax = b$

0. Initialize $r \leftarrow b, \quad S \leftarrow []$
1. Largest correlation find j st $|a_j^T r| = \|A^T r\|_\infty$
2. Get new column $S \leftarrow [S \quad a_j]$
3. Least squares $\min \|b - Sx\|_2$
4. Update residual $r \leftarrow b - Sx$

Canonical greedy approach

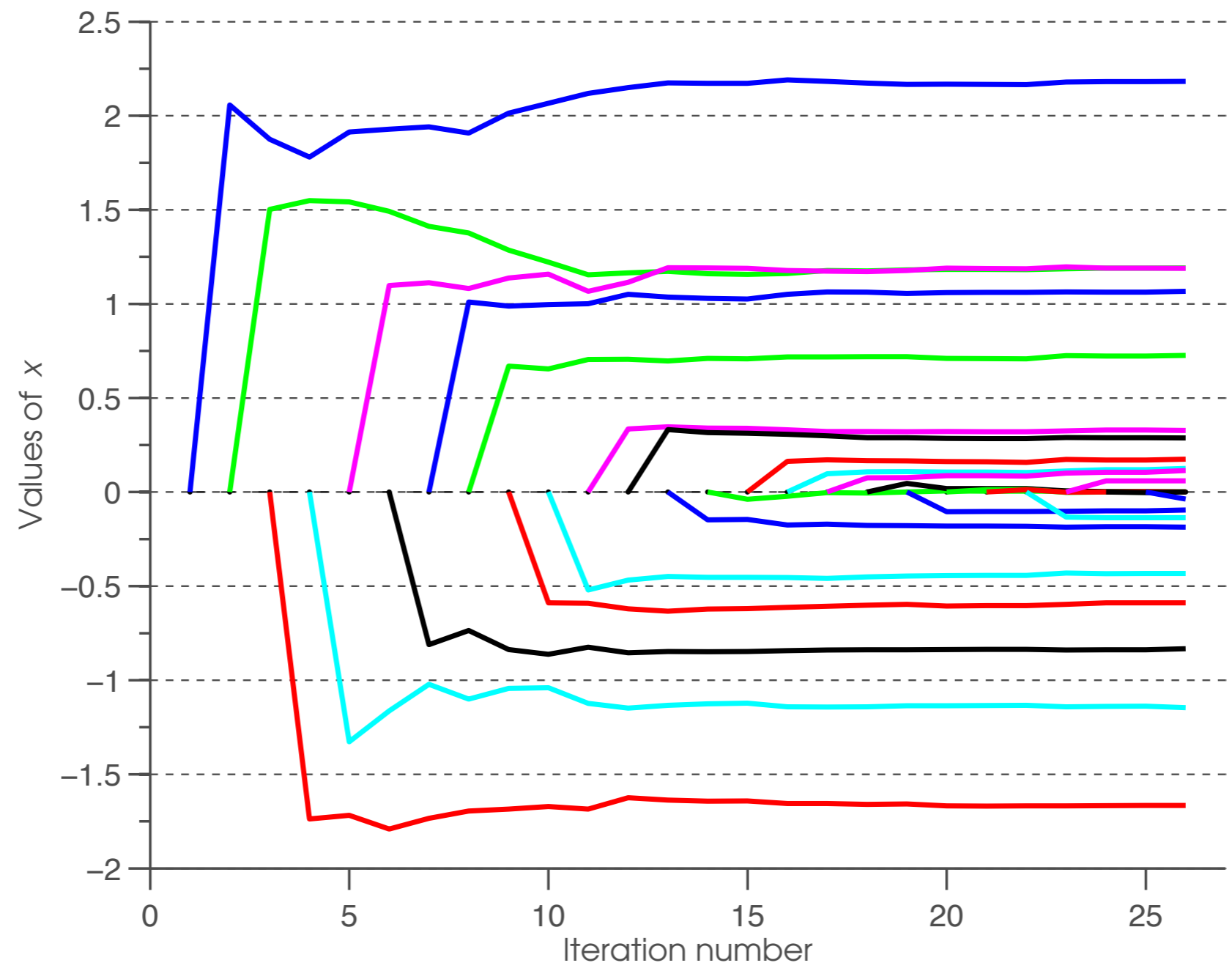
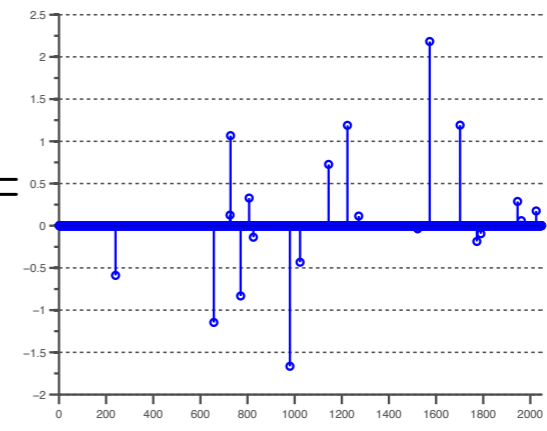
Orthogonal matching pursuit (OMP) for sparse $Ax = b$

0. Initialize $r \leftarrow b, \quad S \leftarrow []$
1. Largest correlation find j st $|a_j^T r| = \|A^T r\|_\infty$
2. Get new column $S \leftarrow [S \quad a_j]$
3. Least squares $\min \|b - Sx\|_2$
4. Update residual $r \leftarrow b - Sx$

- **Often** recovers sparsest solution
- Not optimizing a particular problem

[Tropp & Gilbert '07]

$$Ax = b, \quad A = \boxed{200 \times 2048}, \quad x =$$



Forward-backward splitting

General problem:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) + P(x)$$

where

f is smooth, P is sparsifying/nonsmooth

Two common sparsifying P :

- $P(x) = \lambda \|x\|_1$ minimize $f(x) + \lambda \|x\|_1$
- $P(x) = \delta_{\mathcal{C}}(x)$ with $\mathcal{C} = \{x \mid \|x\|_1 \leq \tau\}$ minimize $f(x)$ st $\|x\|_1 \leq \tau$

Many related methods:

projected gradient, iterative soft-thresholding, fixed point, proximal point, augmented Lagrangian, Bregman, split Bregman, alternating direction method of multipliers, etc

Proximity operators

Proximity operator of P :

$$\mathbf{prox}_P(x) \text{ solves } \underset{z \in \mathbb{R}^n}{\text{minimize}} \frac{1}{2} \|x - z\|^2 + P(z)$$

Generalizes the **projection** operator of $\mathcal{C} \subseteq \mathbb{R}^n$:

$$\mathbf{proj}_{\mathcal{C}}(x) \text{ solves } \underset{z \in \mathbb{R}^n}{\text{minimize}} \frac{1}{2} \|x - z\|^2 + \delta_{\mathcal{C}}(z)$$

where $\delta_{\mathcal{C}}(x)$ is the indicator function for \mathcal{C} :

$$\delta_{\mathcal{C}}(x) = \begin{cases} 0 & \text{if } x \in \mathcal{C} \\ \infty & \text{otherwise} \end{cases}$$

Soft thresholding is proximity operator of the one-norm:

$$P(z) = \|z\|_1 \quad \rightarrow \quad \mathbf{prox}_P(x) := \mathcal{S}_{\lambda}(x) := \text{sgn}(x) \cdot \max\{|x| - \lambda, 0\}$$

Forward-backward splitting algorithm

$$\text{minimize } f(x) + P(x)$$

Optimality: x^* is a solution if and only if

$$x^* = \text{prox}_{\alpha P}(x^* - \alpha \nabla f(x^*)) \quad \text{for all } \alpha > 0$$

Algorithm:

$$x^{k+1} = \underbrace{\text{prox}_{\alpha^k P}}_{\text{backward step}} \left[\underbrace{x^k - \alpha^k \nabla f(x^k)}_{\text{forward step}} \right]$$

where

$$\alpha^k \in (\mu, 2/L), \quad \mu > 0, \quad \|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$$

Convergence:

Every sequence x^k generated by this algorithm converges to a solution

Iterative soft-thresholding (IST)

$$\text{minimize } f(x) + \lambda \|x\|_1$$

For $f(x) = \frac{1}{2} \|r\|^2$, $r \equiv Ax - b$,

$$\begin{aligned} x^{k+1} &= \text{prox}_{\lambda \|\cdot\|_1}(x^k - \alpha^k A^T r_k) \\ &= \mathcal{S}_\lambda(x^k - \alpha^k A^T r_k) \end{aligned}$$

where

$$\alpha^k \in (0, 2/\|A\|^2) \quad \text{and} \quad \mathcal{S}_\lambda(x) = \text{sgn}(x) \cdot \max\{|x| - \lambda, 0\}$$

[aka, iterative shrinkage, iterative Landweber]

Derived from various viewpoints:

- expectation-maximization [Figueiredo & Nowak 03]
- surrogate functionals [Daubechies et al 04]
- forward-backward splitting: $\min f(x) + P(x)$ [Combettes & Wajs 05]
- fixed-point w/continuation: $\min f(x) + \lambda \|x\|_1$ [Hale et al 07]

Research & Development

Research and development

- Hybrid stochastic-deterministic methods
Mark Schmidt, PDF
- Weighted one-norm minimization
Hassan Mansour, PDF
- Newton methods for sparse optimization
Ives Macedo, incoming PhD
- Spot: A linear-operator toolbox for Matlab
Ewout van den Berg (2009 PhD), et al.

Email

- `mpf@cs.ubc.ca`

Web

- `http://www.cs.ubc.ca/~mpf`