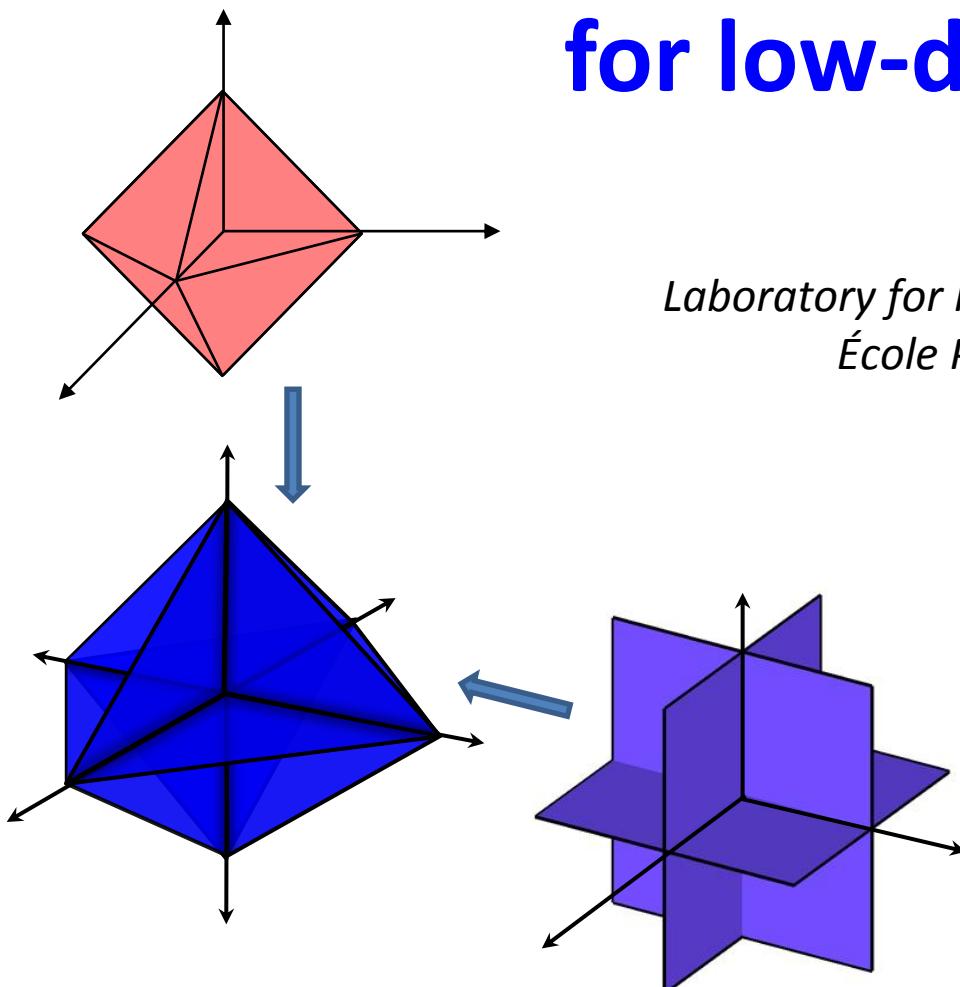


# Convex and non-convex approaches for low-dimensional models



*Volkan Cevher*

*Laboratory for Information and Inference Systems (LIONS)*  
*École Polytechnique Fédérale de Lausanne (EPFL)*  
*Switzerland*  
<http://lions.epfl.ch>

*Mário Figueiredo*

*Instituto de Telecomunicações (IT)*  
*Instituto Superior Técnico (IST)*  
*Portugal*  
<http://wwwlx.it.pt/~mtf/>

# Linear Inverse Problems

$$u = \Phi x$$

Diagram illustrating the linear inverse problem equation  $u = \Phi x$ . The variables are represented as:

- $u$ : A vertical vector of size  $M \times 1$ , composed of  $M$  colored blocks.
- $\Phi$ : A square matrix of size  $M \times N$  ( $M < N$ ), represented by a grid of colored squares.
- $x$ : A vertical vector of size  $N \times 1$ , composed of  $N$  colored blocks.

The diagram shows the multiplication of the matrix  $\Phi$  by the vector  $x$  to produce the vector  $u$ .

Myriad applications involve linear dimensionality reduction  
**deconvolution to data mining**  
**compression to compressive sensing**  
**geophysics to medical imaging**

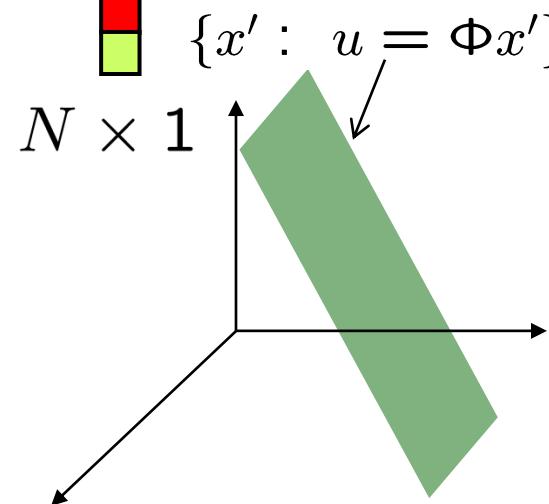
# Linear Inverse Problems

$$u = \Phi x$$

Diagram illustrating the linear inverse problem  $u = \Phi x$ . On the left, a vertical vector  $u$  of size  $M \times 1$  is shown as a stack of colored blocks. In the middle, an equals sign ( $=$ ) is followed by a large square matrix  $\Phi$  of size  $M \times N$  ( $M < N$ ). The matrix  $\Phi$  is filled with a variety of colored blocks. On the right, a vertical vector  $x$  of size  $N \times 1$  is shown as a stack of colored blocks, matching the colors of the blocks in  $u$ .

- **Challenge:** Null space of  $\Phi$ :  $\mathcal{N}(\Phi)$

$$\Phi x' = \Phi(x + v) = u, \quad \forall v \in \mathcal{N}(\Phi)$$



# Linear Inverse Problems



**Deterministic**

**Probabilistic**

**Prior**

 parsity

distribution

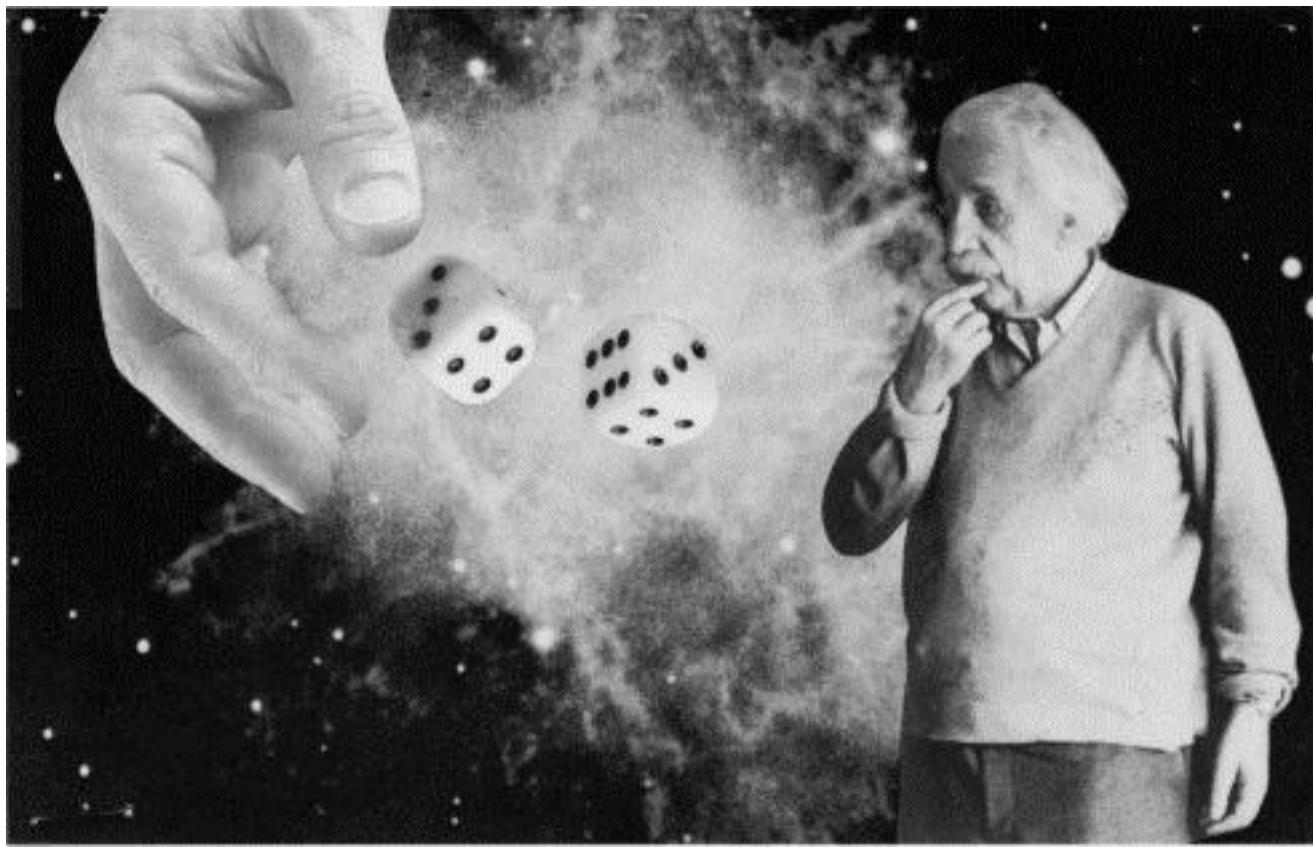
**Metric**

$\ell_p$ -norm\*

likelihood/  
posterior

$$*: \|x\|_p = (\sum_i |x_i|^p)^{1/p}$$

# Deterministic Low-Dimensional Models



# Sparse representations

- **Sparse** signal  $\alpha$

only  $K$  out of  $N$   
coordinates nonzero

$$K \ll N$$

**support:**

$$\mathcal{S} = \{i : x_i \neq 0\}$$

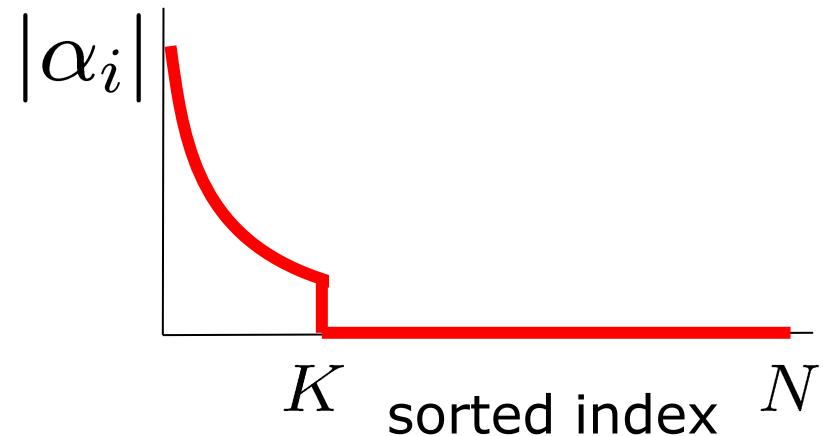
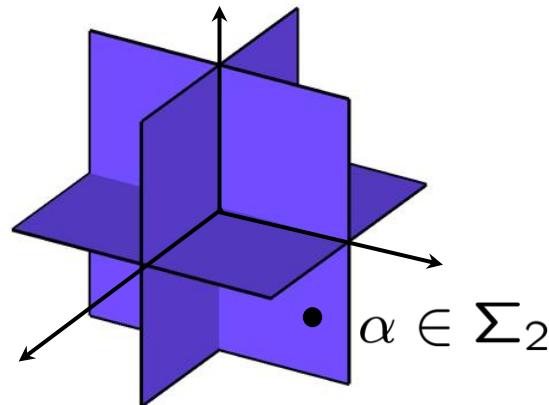
$$\|\alpha\|_0 = |\mathcal{S}| = K$$



$\alpha$

$$K = 2$$

$\mathbb{R}^3$



# Sparse representations

- **Sparse** signal  $x$   
only K out of N  
coordinates nonzero  
in an **appropriate  
representation**
- Sparse representations  
**sparse** transform  
coefficients  $\alpha$ 
  - Basis representations  
 $\Psi \in \mathbb{R}^{N \times N}$ 
    - **Wavelets**, DCT...
  - Frame representations  
 $\Psi \in \mathbb{R}^{N \times L}, L > N$ 
    - Gabor, curvelets, shearlets...
  - Other **dictionary** representations...

$$x = \Psi \times \alpha$$

The diagram shows the mathematical equation for sparse representation. On the left, a vertical vector  $x$  is shown as a column of colored blocks (red, green, blue, magenta, cyan, yellow, purple, etc.). In the center, there is an equals sign ( $=$ ). To the right of the equals sign is a large square matrix  $\Psi$ , which is also composed of colored blocks. To the right of  $\Psi$  is a times symbol ( $\times$ ). To the right of the times symbol is a vertical vector  $\alpha$ , represented by a column of colored blocks.



# Sparse representations

- Sparse signal:

only  $K$  out of  $N$  coordinates nonzero

$$K \ll N$$

- Sparse representations:

*sparse* transform coefficients

$$\begin{matrix} & \\ & \\ & = \\ \begin{matrix} & \\ & \\ & \end{matrix} & \begin{matrix} & \\ & \\ & \end{matrix} & \begin{matrix} & \\ & \\ & \end{matrix} \\ x & = & \Psi & \times & \alpha \end{matrix}$$

- A fundamental impact:

$$\begin{matrix} u & = & \Phi & x \\ \begin{matrix} & \\ & \\ & \end{matrix} & = & \begin{matrix} & \\ & \\ & \end{matrix} & \begin{matrix} & \\ & \\ & \end{matrix} \\ & & & \end{matrix}$$

# Sparse representations

- Sparse signal:

only K out of N  
coordinates nonzero

$$K \ll N$$

- Sparse representations:

*sparse* transform  
coefficients

- A fundamental impact:

$$\begin{matrix} & & \\ \textcolor{red}{\boxed{\cdot}} & = & \textcolor{red}{\boxed{\cdot}} & \textcolor{blue}{\boxed{\cdot}} & \times & \textcolor{green}{\boxed{\cdot}} \\ & & & & & \\ \textcolor{red}{x} & = & \Psi & \times & \alpha & \\ & & & & & \end{matrix}$$

$$\begin{matrix} u & = & \Phi & \Psi & \alpha \\ \textcolor{red}{\boxed{\cdot}} & = & \textcolor{red}{\boxed{\cdot}} & \textcolor{blue}{\boxed{\cdot}} & \textcolor{green}{\boxed{\cdot}} \\ & & & & \\ x & & & & \end{matrix}$$

# Sparse representations

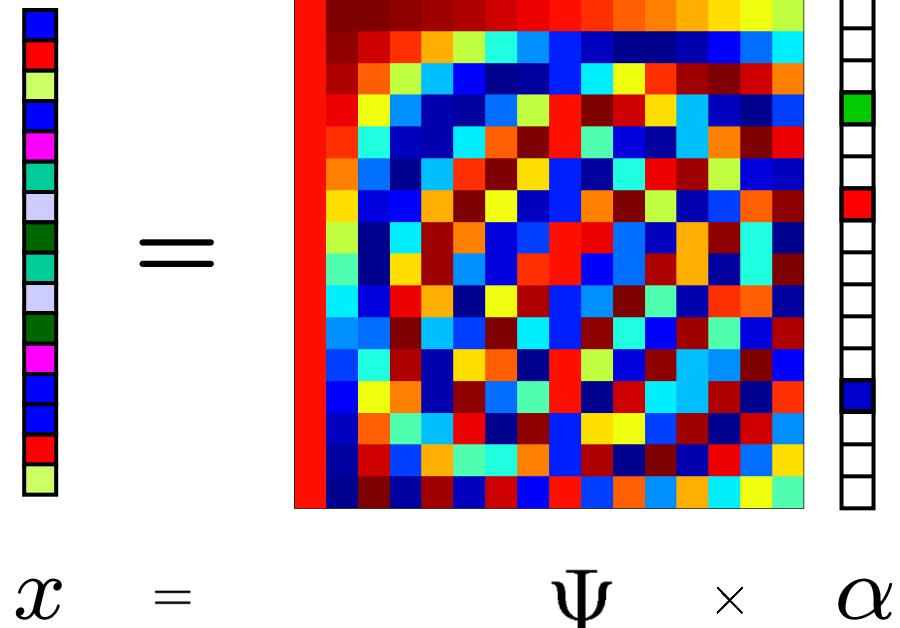
- Sparse signal:

only  $K$  out of  $N$  coordinates nonzero

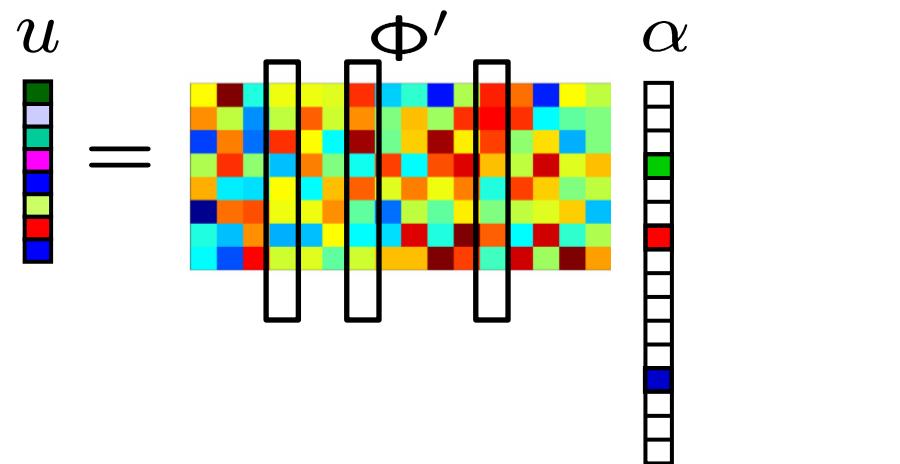
$$K \ll N$$

- Sparse representations:

*sparse* transform coefficients

$$\begin{matrix} & \\ & \\ \textcolor{blue}{\alpha} & = & \textcolor{red}{\Psi} & \times & \textcolor{green}{x} \\ & \\ & & & & \end{matrix}$$


- A fundamental impact:

$$\begin{matrix} u & = & \Phi' & \times & \alpha \\ \textcolor{blue}{\alpha} & = & \textcolor{red}{\Phi}' & \times & \textcolor{green}{u} \\ & & & & \end{matrix}$$


# Sparse representations

- Sparse signal:

only  $K$  out of  $N$  coordinates nonzero

$$K \ll N$$

- Sparse representations:

*sparse* transform coefficients

- A fundamental impact:

$$\Phi$$

becomes effectively low dimensional\*

$$M \times K$$

The diagram illustrates the decomposition of a sparse signal  $x$  into a product of a sparse matrix  $\Psi$  and a sparse vector  $\alpha$ . On the left, a vertical vector  $x$  is shown with only  $K$  non-zero entries highlighted in various colors (red, green, blue, yellow). This is followed by an equals sign. To the right of the equals sign is a large, square, multi-colored matrix  $\Psi$ , which is highly sparse. Another equals sign follows, and then a vertical vector  $\alpha$  is shown, which is also very sparse, with only  $K$  non-zero entries.

$$x = \Psi \times \alpha$$

The diagram illustrates the decomposition of a sparse vector  $u$  into a product of a sparse matrix  $\Phi'$  and a sparse vector  $\alpha$ . On the left, a vertical vector  $u$  is shown with only  $K$  non-zero entries highlighted in various colors. This is followed by an equals sign. To the right of the equals sign is a small, square, multi-colored matrix  $\Phi'$ , which is highly sparse. To the right of  $\Phi'$  is another equals sign, and then a vertical vector  $\alpha$  is shown, which is also very sparse, with only  $K$  non-zero entries.

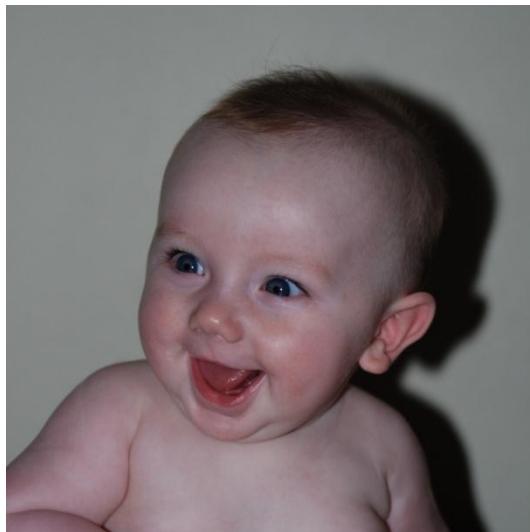
$$u = \Phi' \times \alpha$$

$$M > K$$

\*: If we knew the locations of the coefficients. **More on this later.**

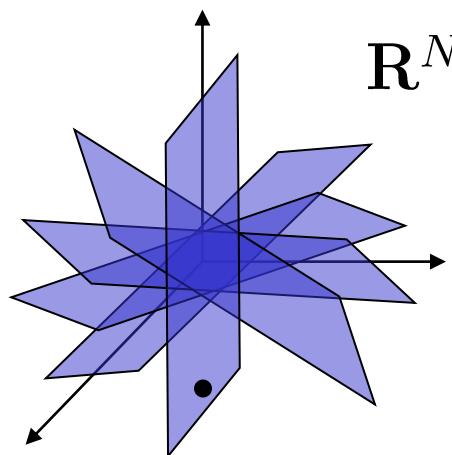
# Low-dimensional signal models

$N$   
pixels

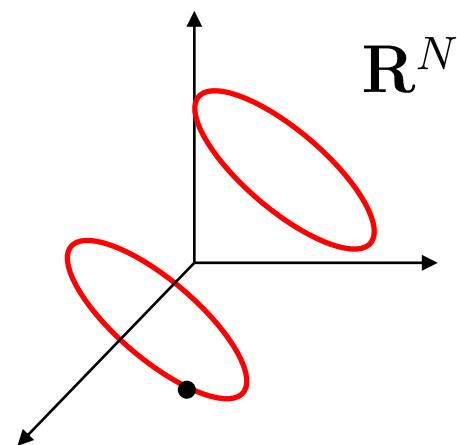


**Information level:**

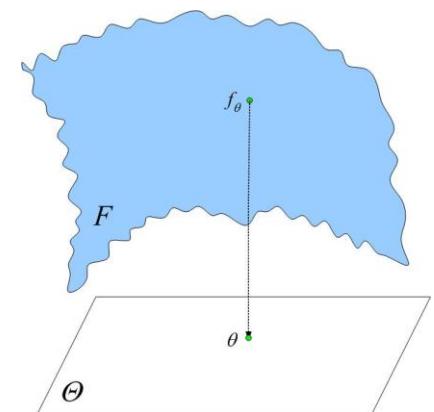
$K \ll N$   
large  
wavelet  
coefficients  
(blue = 0)



sparse  
signals



low-rank  
matrices

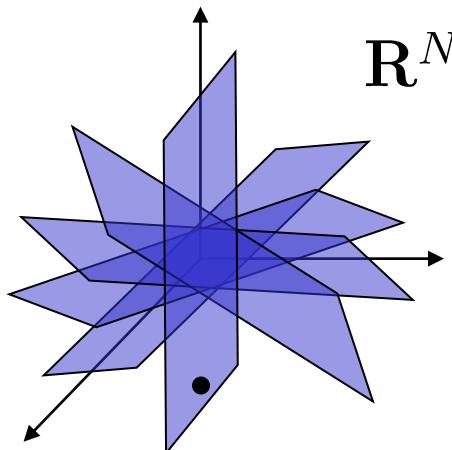


nonlinear  
models

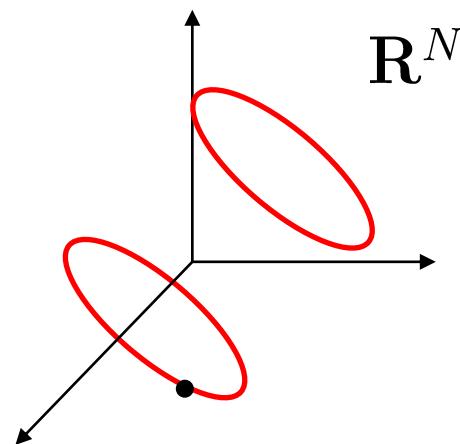
# Low-dimensional signal models

- This tutorial:

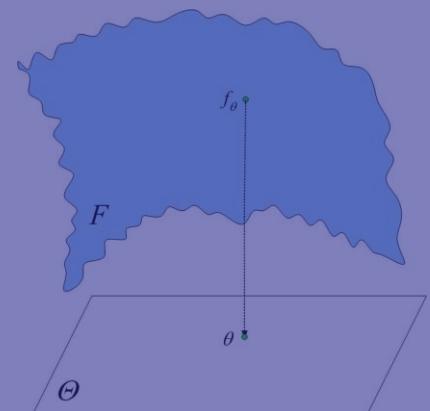
**Low-dimensional models  
based on  
linear representations**



sparse  
signals



low-rank  
matrices



nonlinear  
models

# Linear representation of low-dimensional models

- A key notion in sparse representation
  - synthesis of the signal using a few vectors

$$\begin{matrix} | & | \\ \textcolor{black}{x} & = & \textcolor{red}{\Psi} & \times & \textcolor{black}{\alpha} \end{matrix}$$

- A slightly different mathematical formalism for generalization

**Synthesis model:**  $x = \sum_{i=1}^{|\mathcal{A}|} a_i c_i \quad a_i \in \mathcal{A}, c_i \geq 0$

$a_i$ : atoms  
 $\mathcal{A}$ : atomic set

i.e., linear (positive) combination of elements from an atomic set

[Chandrasekaran et al. 2010]

# Linear representation of low-dimensional models

- A key notion in sparse representation
  - synthesis of the signal using a few vectors

$$\begin{matrix} | & | \\ \textcolor{black}{\bullet} & \textcolor{black}{\bullet} \\ | & | \end{matrix} = \begin{matrix} | & | \\ \textcolor{red}{\bullet} & \textcolor{blue}{\bullet} \\ | & | \end{matrix} \times \begin{matrix} | & | \\ \textcolor{black}{\bullet} & \textcolor{black}{\bullet} \\ | & | \end{matrix}$$

$x = \Psi \times \alpha$

- Sparse representations via the atomic formulation

$$x = \sum_{i=1}^{|A|} a_i c_i \quad \begin{aligned} a_i &\in \mathcal{A}, c_i \geq 0 \\ a_i &\text{: atoms} \\ \mathcal{A} &\text{: atomic set} \end{aligned}$$

- Example:

$$\Psi = [\psi_1, \dots, \psi_L] \quad \mathcal{A} = \{\psi_1, \dots, \psi_L, -\psi_1, \dots, -\psi_L\}$$

$$\text{rank}(\Psi) = N$$

$$c_i = \begin{cases} \alpha_i, & \alpha_i > 0; \\ 0, & \text{otherwise.} \end{cases} \quad i = 1, \dots, L$$
$$c_{i+L} = \begin{cases} -\alpha_i, & \alpha_i < 0; \\ 0, & \text{otherwise.} \end{cases}$$

# Linear representation of low-dimensional models

- Basic definitions on **low-dimensional** atomic representations

$$x = \sum_{i=1}^{|A|} a_i c_i \quad a_i \in A, c_i \geq 0 \quad K \ll N$$

$\|c_i\|_0 \leq K$

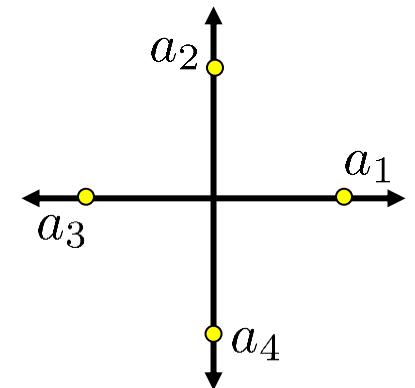
# Linear representation of low-dimensional models

- Basic definitions on low-dimensional *atomic representations*

$$x = \sum_{i=1}^{|A|} a_i c_i \quad a_i \in A, c_i \geq 0 \quad \|c_i\|_0 \leq K \quad K \ll N$$

- $\text{conv}(A)$ : convex hull of atoms in  $A$

$$A = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right\}$$



$$\text{conv}(A) = \left\{ \sum_i a_i \beta_i : a_i \in A, \beta_i \in \mathbb{R}_+, \sum_{i=1}^n \beta_i = 1, n = 1, 2, \dots, |A| \right\}$$

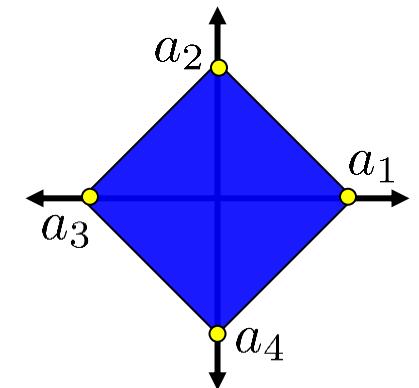
# Linear representation of low-dimensional models

- Basic definitions on low-dimensional *atomic representations*

$$x = \sum_{i=1}^{|A|} a_i c_i \quad a_i \in A, c_i \geq 0 \quad \|c_i\|_0 \leq K \quad K \ll N$$

- $\text{conv}(A)$ : convex hull of atoms in  $A$

$$A = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right\}$$



**atomic ball**

$$\text{conv}(A) = \left\{ \sum_i a_i \beta_i : a_i \in A, \beta_i \in \mathbb{R}_+, \sum_{i=1}^n \beta_i = 1, n = 1, 2, \dots, |A| \right\}$$

# Linear representation of low-dimensional models

- Basic definitions on low-dimensional *atomic representations*

$$x = \sum_{i=1}^{|A|} a_i c_i$$

$$a_i \in A, c_i \geq 0$$

$$\|c_i\|_0 \leq K$$

$$K \ll N$$

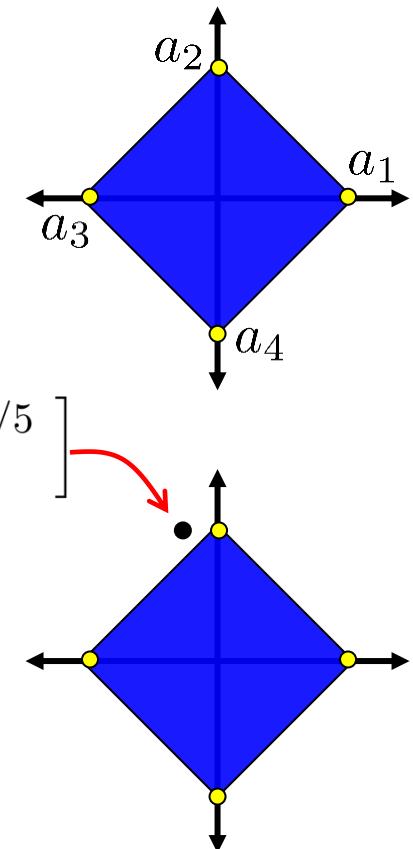
- $\text{conv}(A)$ : convex hull of atoms in  $A$

$$A = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right\}$$

- $\|x\|_A$ : atomic norm\*

$$\|x\|_A = \inf\{t > 0 : x \in t \times \text{conv}(A)\}$$

$$x = \begin{bmatrix} -1/5 \\ 1 \end{bmatrix}$$



\*: requires  $A$  to be centrally symmetric

# Linear representation of low-dimensional models

- Basic definitions on low-dimensional *atomic representations*

$$x = \sum_{i=1}^{|A|} a_i c_i$$

$$a_i \in A, c_i \geq 0$$

$$\|c_i\|_0 \leq K$$

$$K \ll N$$

- $\text{conv}(A)$ : convex hull of atoms in  $A$

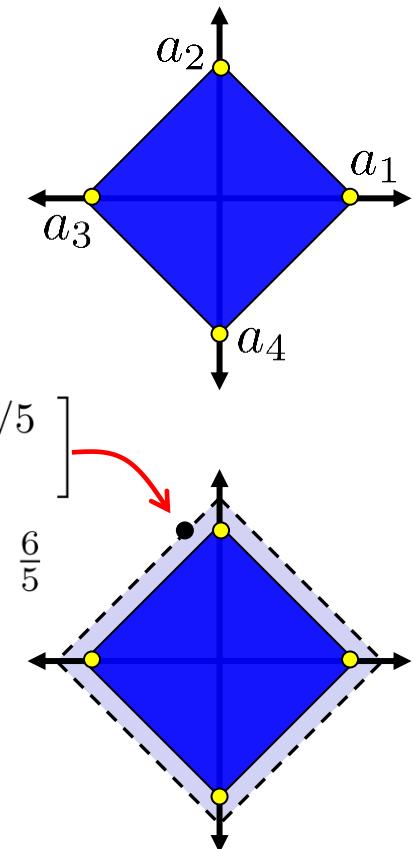
$$A = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right\}$$

- $\|x\|_A$ : atomic norm\*

$$\|x\|_A = \inf\{t > 0 : x \in t \times \text{conv}(A)\}$$

$$x = \begin{bmatrix} -1/5 \\ 1 \end{bmatrix}$$

$$\|x\|_A = \frac{6}{5}$$



\*: requires  $A$  to be centrally symmetric

# Linear representation of low-dimensional models

- Basic definitions on low-dimensional *atomic representations*

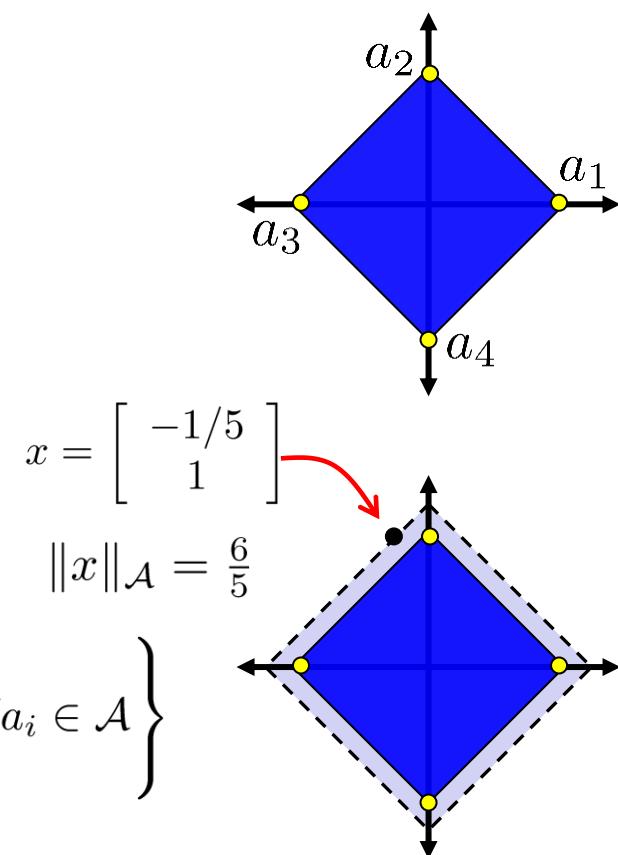
$$x = \sum_{i=1}^{|A|} a_i c_i \quad a_i \in A, c_i \geq 0 \quad \|c_i\|_0 \leq K \quad K \ll N$$

- $\text{conv}(A)$ : convex hull of atoms in  $A$

$$A = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right\}$$

- $\|x\|_A$ : atomic norm\*

$$\|x\|_A = \inf\{t > 0 : x \in t \times \text{conv}(A)\}$$



**Alternative:**  $\|x\|_A = \inf \left\{ \sum_{i=1}^{|A|} c_i : x = \sum_{i=1}^{|A|} a_i c_i, c_i \geq 0, \forall a_i \in A \right\}$

\*: requires  $A$  to be centrally symmetric

# Linear representation of low-dimensional models

Examples with easy forms:

- *sparse vectors*

$$\mathcal{A} = \{\pm e_i\}_{i=1}^N$$

$\text{conv}(\mathcal{A}) = \text{cross-polytope}$

$$\|x\|_{\mathcal{A}} = \|x\|_1$$

- *low-rank matrices*

$$\mathcal{A} = \{A : \text{rank}(A) = 1, \|A\|_F = 1\}$$

$\text{conv}(\mathcal{A}) = \text{nuclear norm ball}$

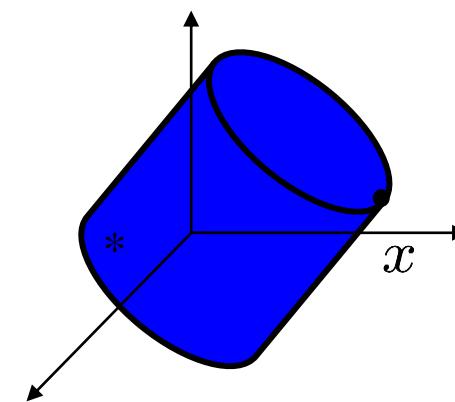
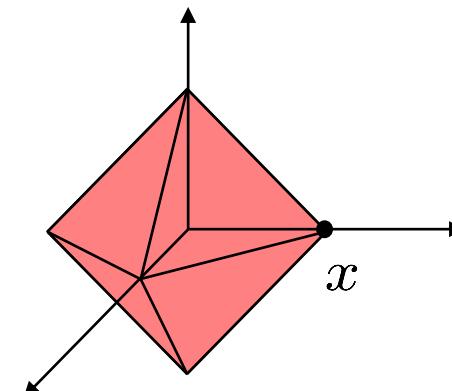
$$\|x\|_{\mathcal{A}} = \|x\|_{\star}$$

- *binary vectors*

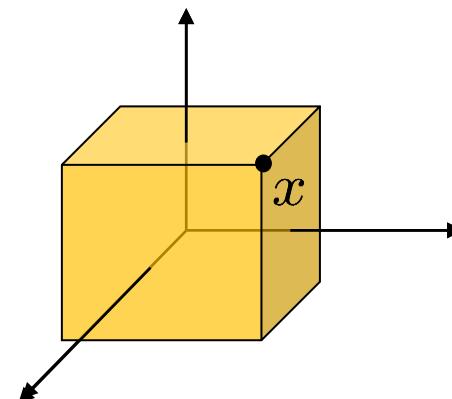
$$\mathcal{A} = \{\pm 1\}^N$$

$\text{conv}(\mathcal{A}) = \text{hypercube}$

$$\|x\|_{\mathcal{A}} = \|x\|_{\infty}$$



\*symmetric  
matrices



# Linear representation of low-dimensional models

Examples with easy forms:

- *sparse vectors*

$$\mathcal{A} = \{\pm e_i\}_{i=1}^N$$

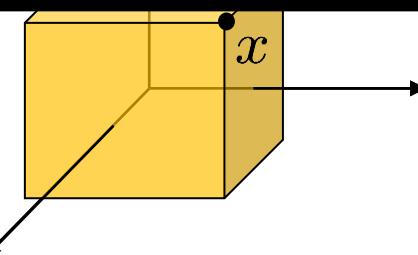
## **Examples with no-so-easy forms:**

- *lo* ✓ $\mathbf{A}$  : infinite set of unit-norm rank-one tensors
- *b* ✓ $\mathbf{A}$  : finite (but large) set of permutation matrices
  - ✓ $\mathbf{A}$  : infinite set of orthogonal matrices
  - ✓ $\mathbf{A}$  : infinite set of matrices constrained by eigenvalues
  - ✓ $\mathbf{A}$  : infinite set of measures
- *b* ✓ $\mathbf{A}$  : finite (but large) set of cut matrices

[Chandrasekaran et al. 2010]

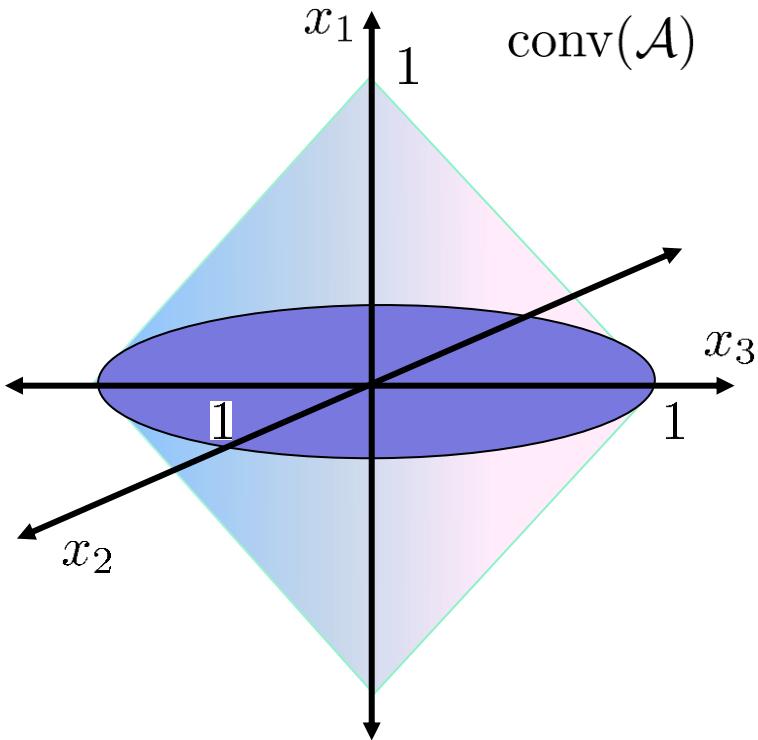
$$\text{conv}(\mathcal{A}) = \text{hypercube}$$

$$\|x\|_{\mathcal{A}} = \|x\|_{\infty}$$



# Linear representation of low-dimensional models

Pop-quiz:

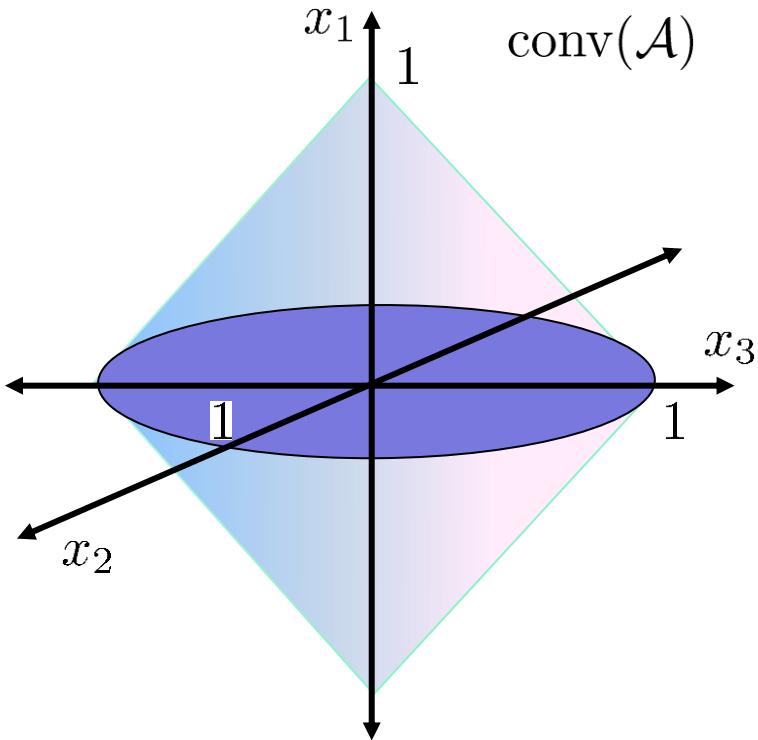


What is  $\|x\|_{\mathcal{A}}$ ?

$$\|x\|_{\mathcal{A}} = \inf\{t > 0 : x \in t \times \text{conv}(\mathcal{A})\}$$

# Linear representation of low-dimensional models

Pop-quiz:



**HINT:**

$$\mathcal{A} = \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \\ 0 \end{bmatrix}, \|x_G\|_2 = 1 \right\}$$

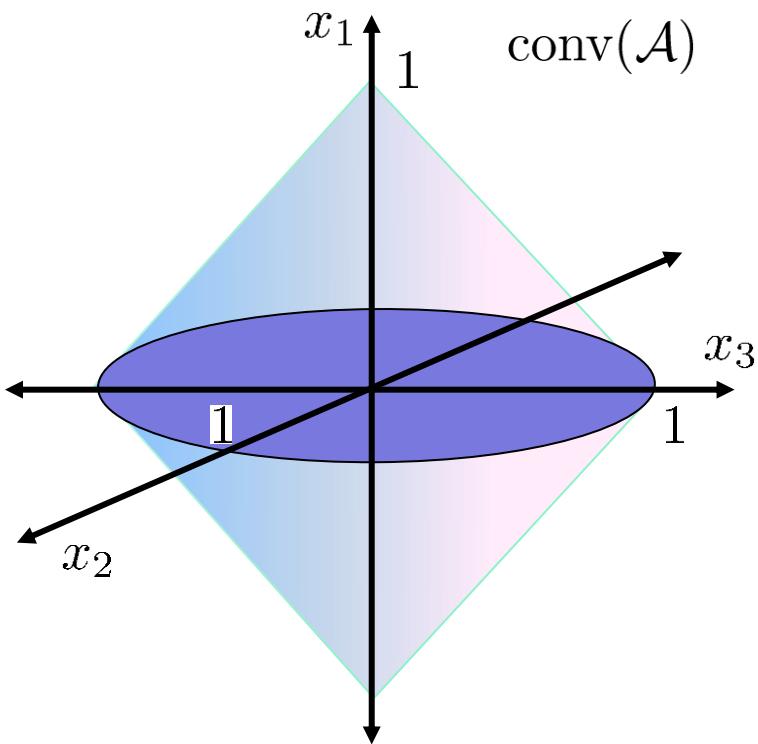
$$G = \{2, 3\}$$

**What is**  $\|x\|_{\mathcal{A}}$ ?

$$\|x\|_{\mathcal{A}} = \inf\{t > 0 : x \in t \times \text{conv}(\mathcal{A})\}$$

# Linear representation of low-dimensional models

Pop-answer:



$$\mathcal{A} = \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \\ 0 \end{bmatrix}, \|x_G\|_2 = 1 \right\}$$

**What is**  $\|x\|_{\mathcal{A}}$  ?

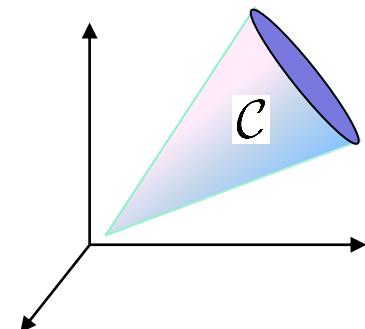
$$\|x\|_{\mathcal{A}} = |x_1| + \|x_G\|_2$$

$$G = \{2, 3\}$$

# Towards algorithms: a geometric perspective

## Other key concepts:

- Cone  $\mathcal{C}$ :  $x, y \in \mathcal{C} \Rightarrow tx + \omega y \in \mathcal{C}, \forall t, \omega \in \mathbb{R}_+$

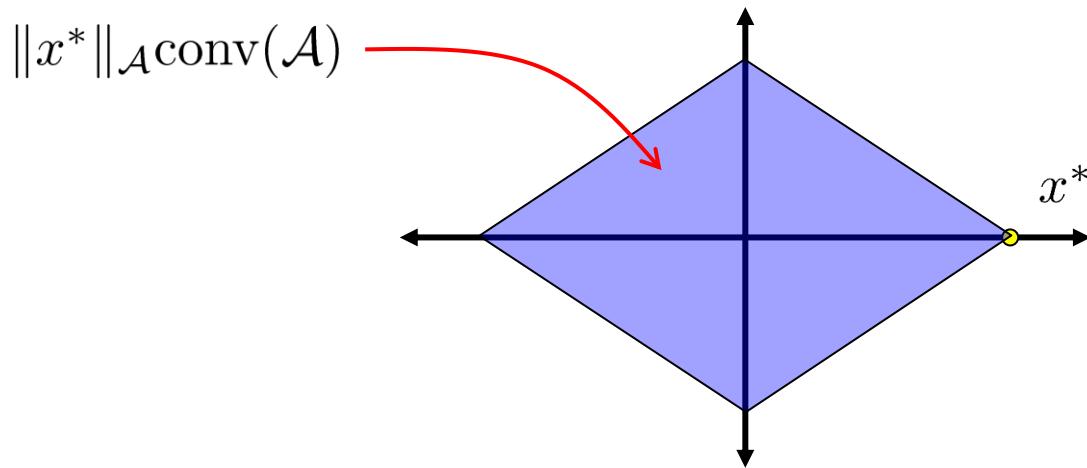
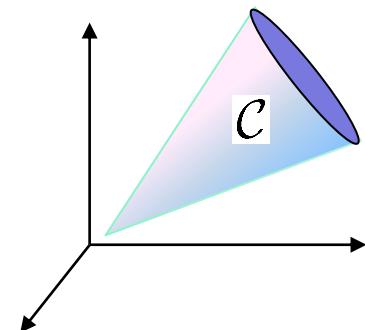


# Towards algorithms: a geometric perspective

## Other key concepts:

- Cone  $\mathcal{C}$ :  $x, y \in \mathcal{C} \Rightarrow tx + \omega y \in \mathcal{C}, \forall t, \omega \in \mathbb{R}_+$
- Tangent cone of  $x^*$  with respect to  $\|x^*\|_{\mathcal{A}\text{conv}}(\mathcal{A})$ :

$$T_{\mathcal{A}}(x^*) = \text{cone}\{z - x^* : \|z\|_{\mathcal{A}} \leq \|x^*\|_{\mathcal{A}}\}$$

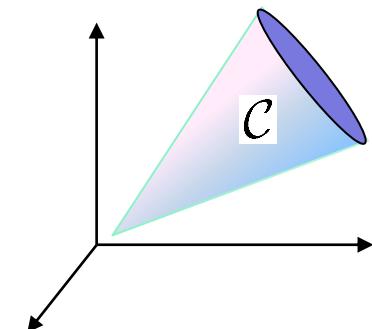
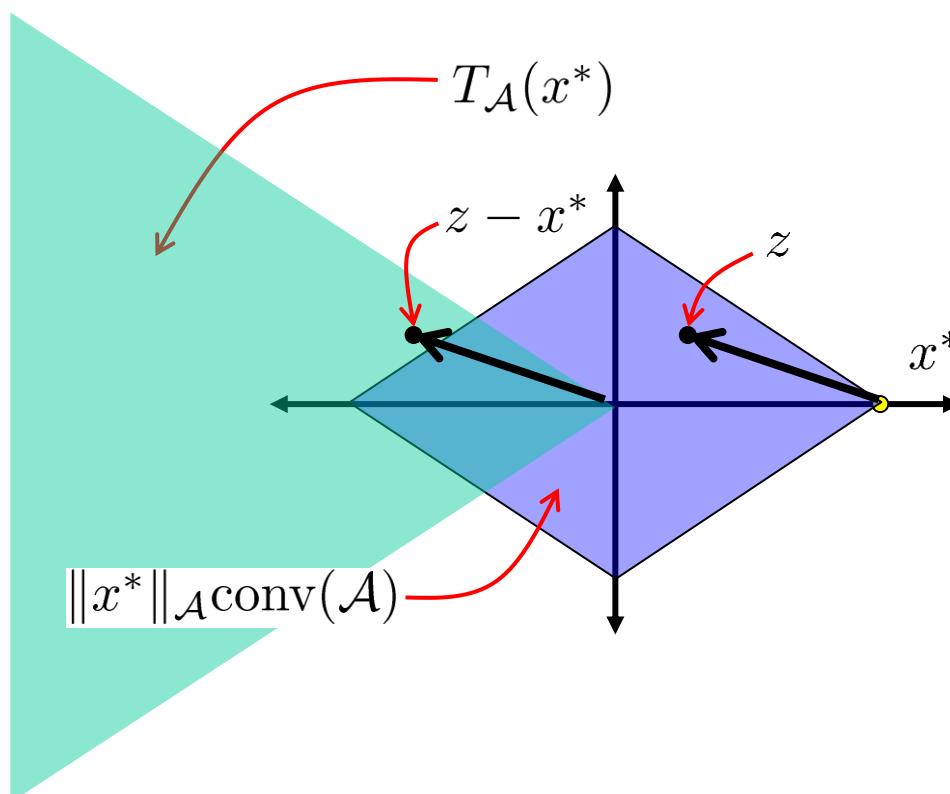


# Towards algorithms: a geometric perspective

## Other key concepts:

- Cone  $\mathcal{C}$ :  $x, y \in \mathcal{C} \Rightarrow tx + \omega y \in \mathcal{C}, \forall t, \omega \in \mathbb{R}_+$
- Tangent cone of  $x^*$  with respect to  $\|x^*\|_{\mathcal{A}\text{conv}}(\mathcal{A})$ :

$$T_{\mathcal{A}}(x^*) = \text{cone}\{z - x^* : \|z\|_{\mathcal{A}} \leq \|x^*\|_{\mathcal{A}}\}$$



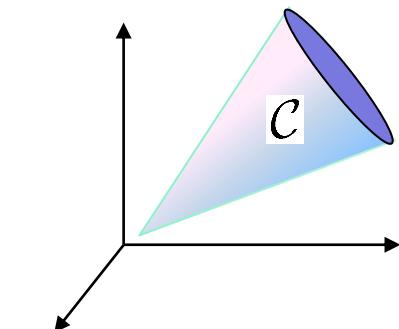
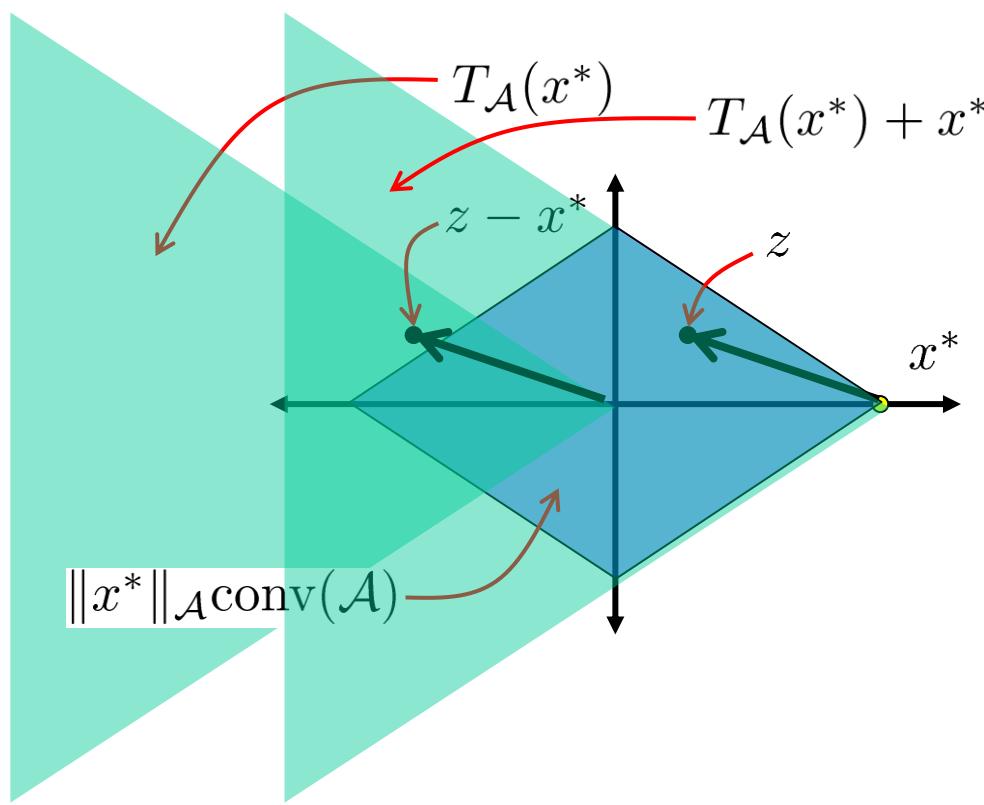
**Tangent cone**  
is the set of descent directions where you do not increase the atomic norm.

# Towards algorithms: a geometric perspective

## Other key concepts:

- Cone  $\mathcal{C}$ :  $x, y \in \mathcal{C} \Rightarrow tx + \omega y \in \mathcal{C}, \forall t, \omega \in \mathbb{R}_+$
- Tangent cone of  $x^*$  with respect to  $\|x^*\|_{\mathcal{A}\text{conv}}(\mathcal{A})$ :

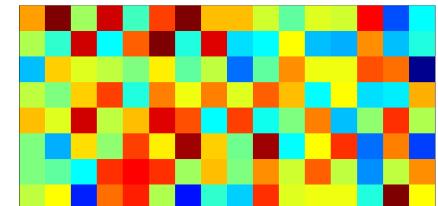
$$T_{\mathcal{A}}(x^*) = \text{cone}\{z - x^* : \|z\|_{\mathcal{A}} \leq \|x^*\|_{\mathcal{A}}\}$$



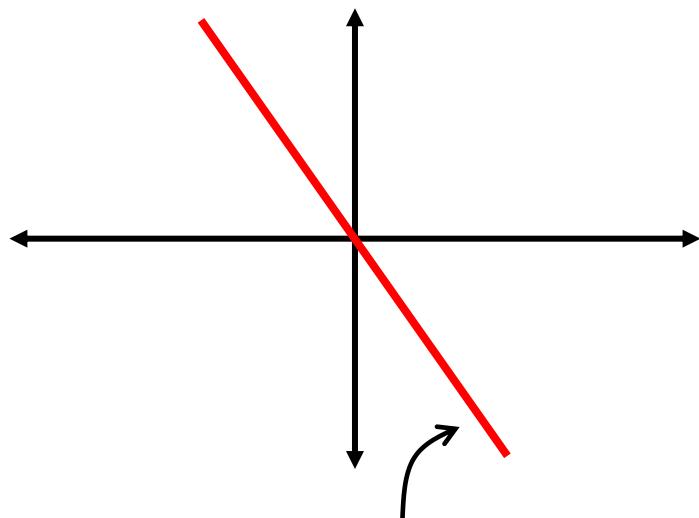
**Tangent cone**  
is the set of descent directions where you do not increase the atomic norm.

# Towards algorithms: a geometric perspective

$\Phi$



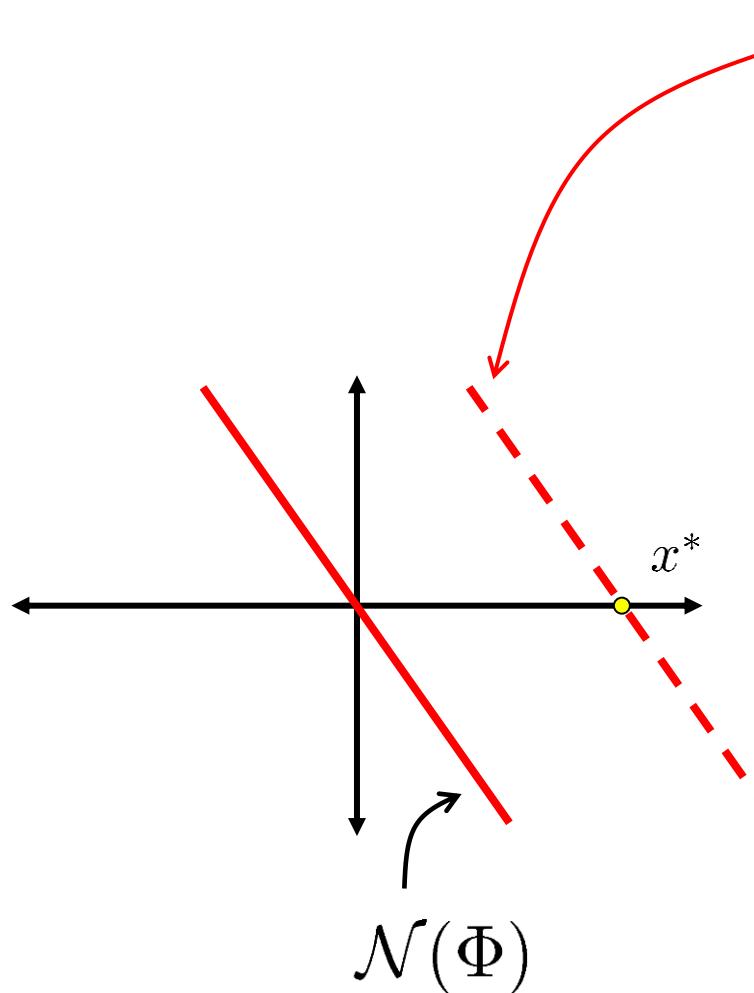
$M \times N$  ( $M < N$ )



Null space of  $\Phi$ :  $\mathcal{N}(\Phi)$

$$\Phi v = 0, \quad \forall v \in \mathcal{N}(\Phi)$$

# Towards algorithms: a geometric perspective



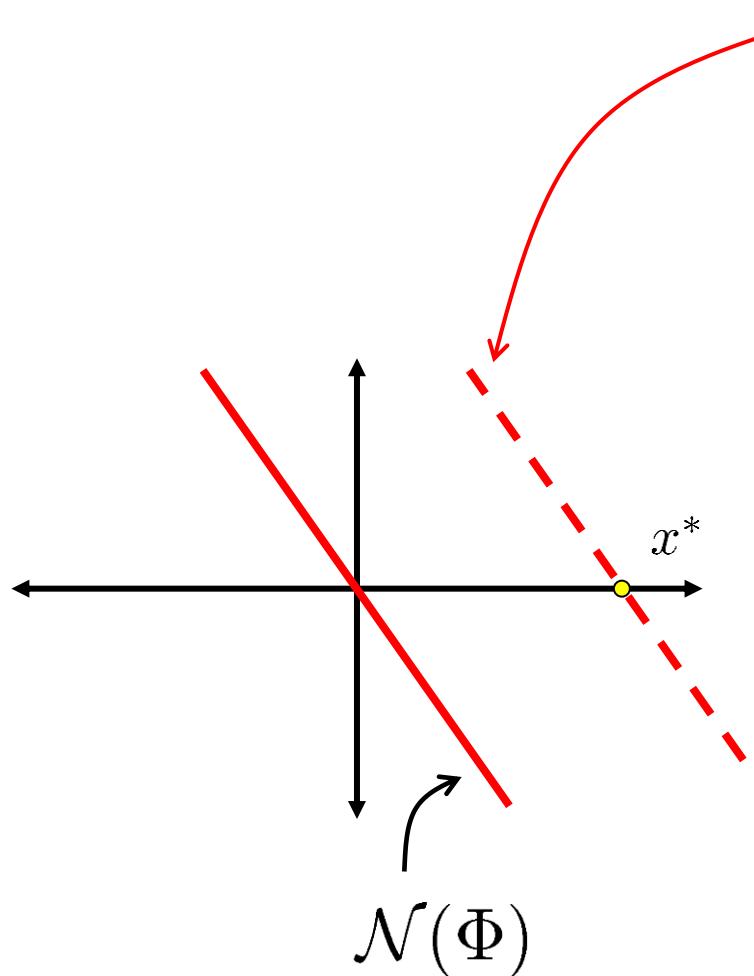
$$u = \Phi x^*$$

Diagram illustrating the matrix equation  $u = \Phi x^*$  using colored vectors and a matrix:

- The vector  $u$  is a vertical column of colors, labeled  $M \times 1$ .
- The matrix  $\Phi$  is a grid of colored squares, labeled  $M \times N$  ( $M < N$ ).
- The vector  $x^*$  is a vertical column of colors, labeled  $N \times 1$ .

The diagram shows the matrix  $\Phi$  as a grid of colored squares, where each square's color represents its value in the matrix. The vector  $u$  is a vertical column of colors, and the vector  $x^*$  is a vertical column of colors. The equation  $u = \Phi x^*$  is represented by an equals sign between the vector  $u$  and the product of the matrix  $\Phi$  and the vector  $x^*$ .

# Towards algorithms: a geometric perspective



$$u = \Phi x^* \\ M \times 1 \quad M \times N \quad (M < N) \\ N \times 1$$

The diagram illustrates the relationship between vectors and matrices. On the left, a vertical vector  $u$  of size  $M \times 1$  is shown. In the center, a matrix  $\Phi$  of size  $M \times N$  ( $M < N$ ) is represented as a grid of colored squares. To the right, a vertical vector  $x^*$  of size  $N \times 1$  is shown.

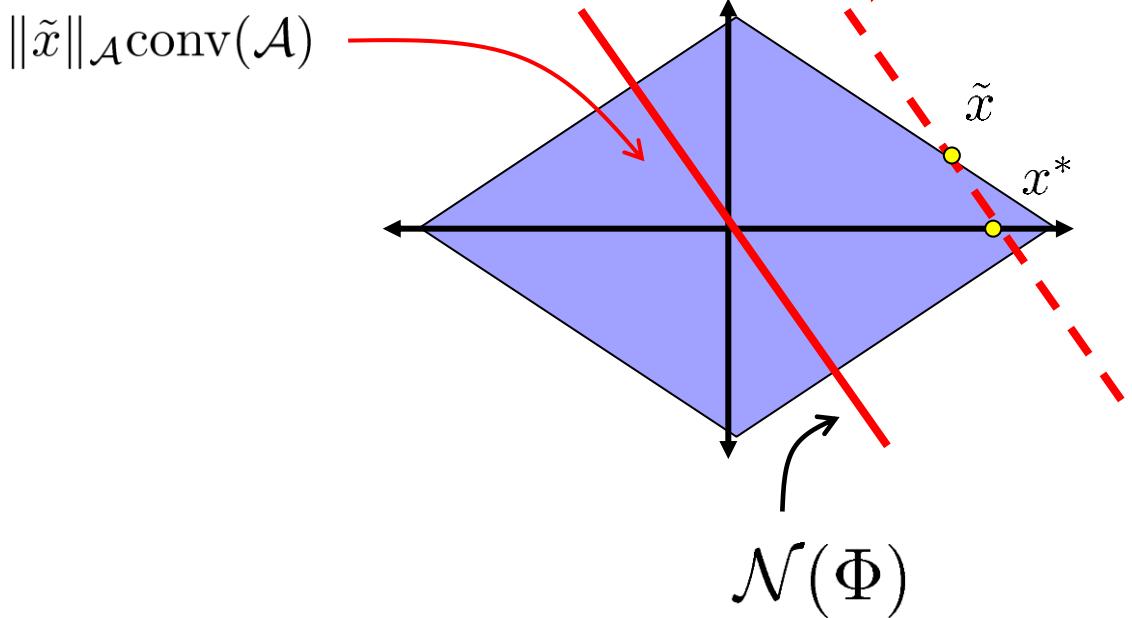
**Consider the criteria:**

$$\hat{x} = \arg \min_{x: u = \Phi x} \|x\|_{\mathcal{A}}$$

# Towards algorithms: a geometric perspective

$$u \quad \Phi \quad x^*$$
$$M \times 1 \quad M \times N \quad (M < N) \quad N \times 1$$

A diagram illustrating matrix multiplication. On the left, there is a vertical vector labeled  $u$  with dimensions  $M \times 1$ . In the center, there is an equals sign followed by a matrix labeled  $\Phi$  with dimensions  $M \times N$  where  $M < N$ . To the right of the equals sign is a vertical vector labeled  $x^*$  with dimensions  $N \times 1$ . The matrix  $\Phi$  is shown as a grid of colored squares (red, green, blue, yellow) with  $M$  rows and  $N$  columns.



**Consider the criteria:**

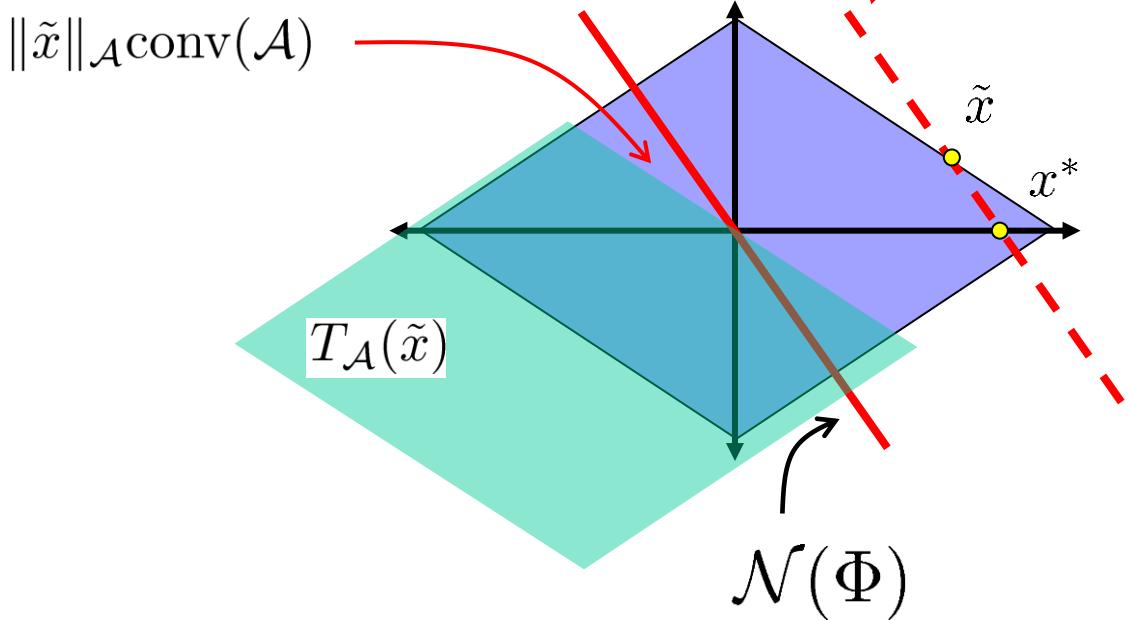
$$\hat{x} = \arg \min_{x: u = \Phi x} \|x\|_{\mathcal{A}}$$

# Towards algorithms: a geometric perspective

$$u \quad \Phi \quad x^*$$

$$M \times 1 \quad M \times N \quad (M < N) \quad N \times 1$$

=



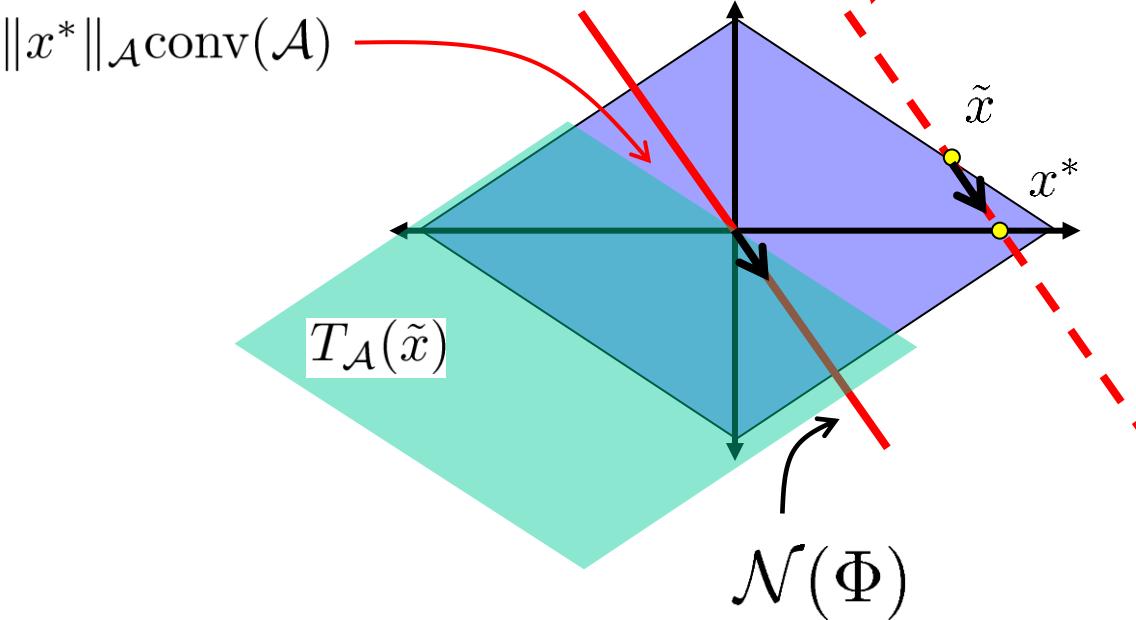
**Consider the criteria:**

$$\hat{x} = \arg \min_{x: u = \Phi x} \|x\|_{\mathcal{A}}$$

# Towards algorithms: a geometric perspective

$$u \quad \Phi \quad x^*$$
$$M \times 1 \quad M \times N \quad (M < N) \quad N \times 1$$

A diagram illustrating the relationship between vectors  $u$ ,  $\Phi$ , and  $x^*$ .  $u$  is a vertical vector of size  $M \times 1$ .  $\Phi$  is a matrix of size  $M \times N$  ( $M < N$ ).  $x^*$  is a vertical vector of size  $N \times 1$ . The equation  $u = \Phi x^*$  is shown, where  $\Phi$  is represented as a grid of colored squares.



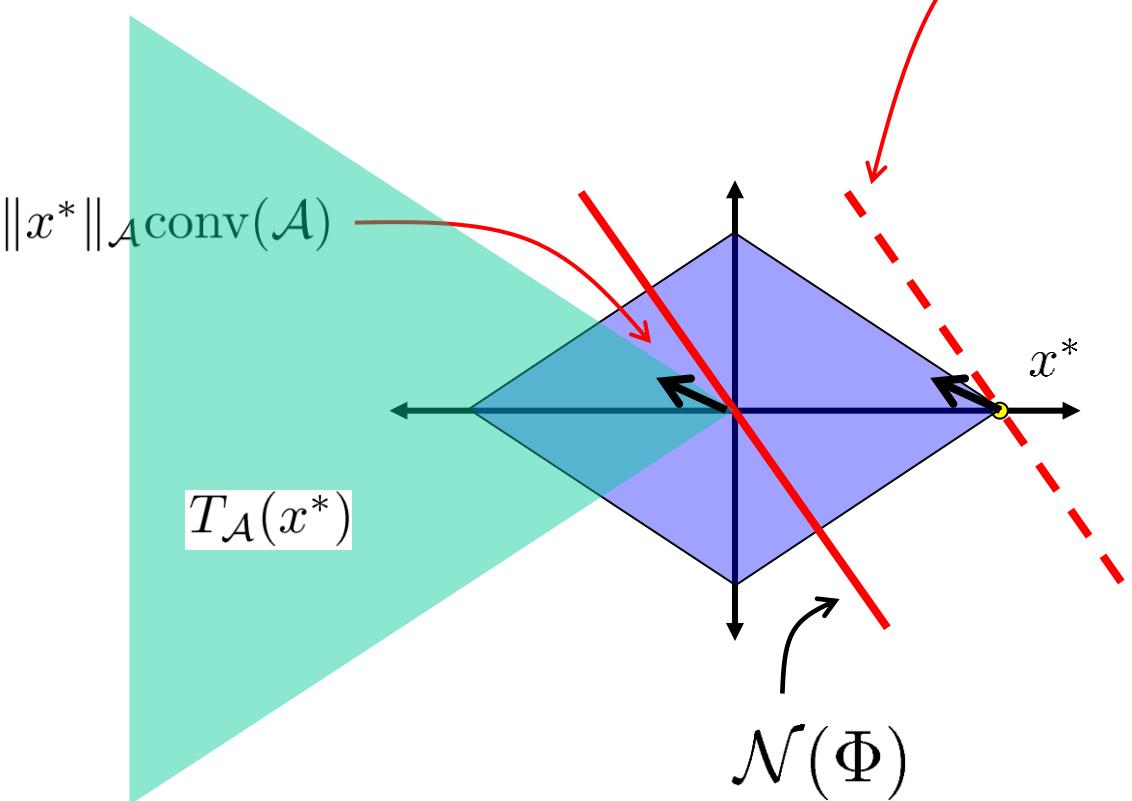
**Consider the criteria:**

$$\hat{x} = \arg \min_{x: u = \Phi x} \|x\|_A$$

# Towards algorithms: a geometric perspective

$$u \quad \Phi \quad x^*$$
$$M \times 1 \quad M \times N \quad (M < N) \quad N \times 1$$

The diagram illustrates the dimensions of the vectors involved in the equation  $u = \Phi x^*$ . On the left, vector  $u$  is shown as a vertical column of height  $M$ , labeled  $M \times 1$ . In the center, matrix  $\Phi$  is represented as a grid of size  $M \times N$  where  $M < N$ . To the right, vector  $x^*$  is shown as a vertical column of height  $N$ , labeled  $N \times 1$ .



**Key observation:**

$$\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) = \{0\} \Rightarrow x^* = \arg \min_{x: u = \Phi x} \|x\|_{\mathcal{A}}$$

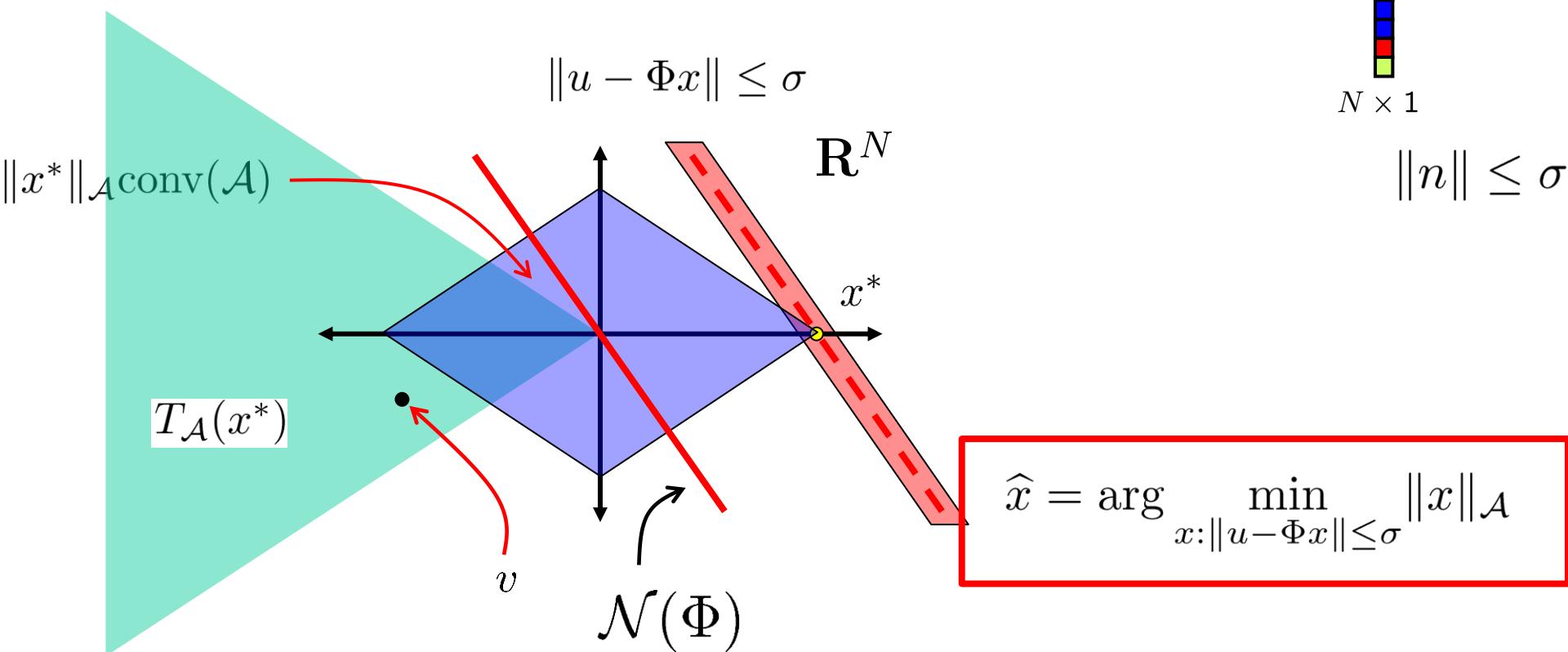
# Towards algorithms: a geometric perspective

How about noise?

$$u = \Phi x^* + n$$

$M \times 1$        $M \times N \ (M < N)$        $N \times 1$

$+ \quad M \times 1$



# Towards algorithms: a geometric perspective

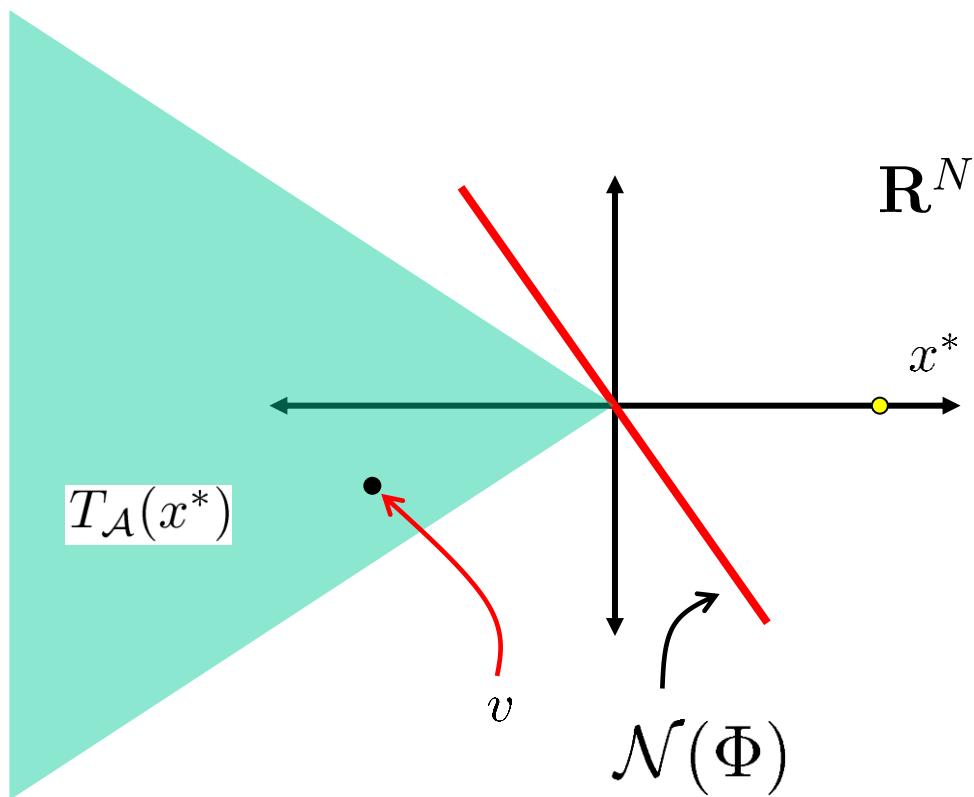
How about noise?

$$u = \Phi x^* + n$$

where

$$\begin{aligned} u &= M \times 1 \\ \Phi &= M \times N \quad (M < N) \\ x^* &= N \times 1 \\ n &= M \times 1 \end{aligned}$$

+



**Stability assumption:**

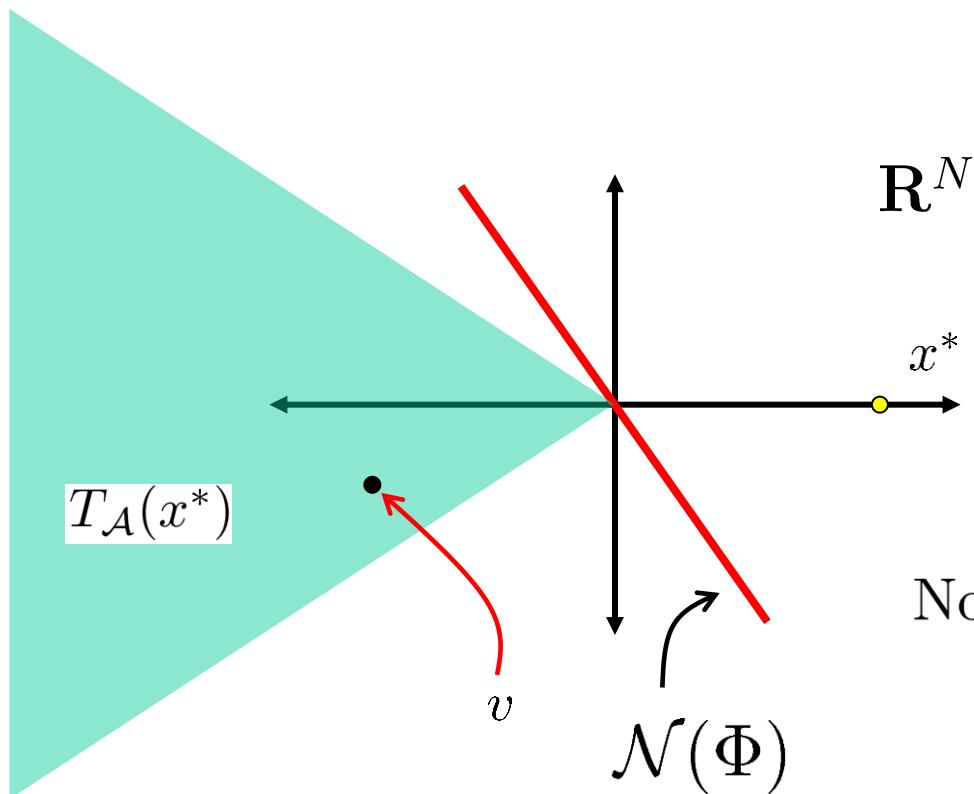
$$\|\Phi v\| \geq \epsilon \|v\|, \forall v \in T_{\mathcal{A}}(x^*)$$

# Towards algorithms: a geometric perspective

How about noise?

$$u = \Phi x^* + n$$

where  $\Phi$  is an  $M \times N$  matrix ( $M < N$ ) with  $M \times 1$  columns  $x^*$  and  $M \times 1$  rows  $n$ .



Note that if  $\mathcal{N}(\Phi) \cap T_A(x^*) = \{0\}$   
 $\Rightarrow \|\Phi v\| > 0, \forall v \neq 0$

**Stability assumption:**

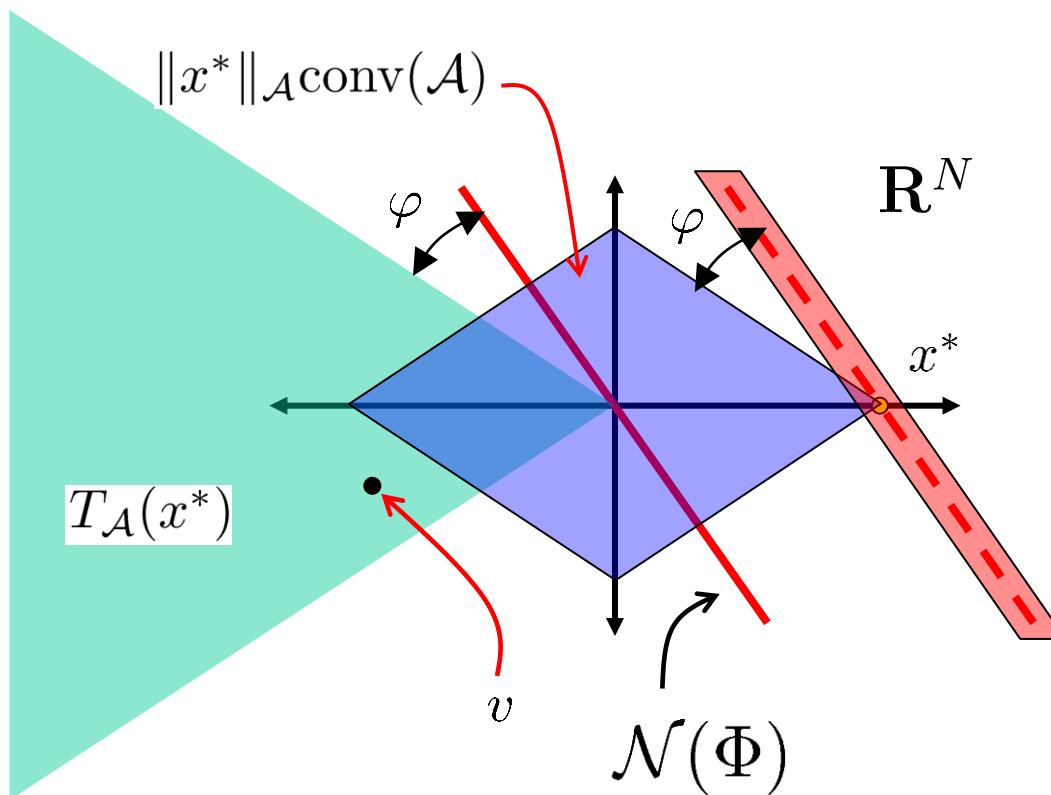
$$\|\Phi v\| \geq \epsilon \|v\|, \forall v \in T_A(x^*)$$

# Towards algorithms: a geometric perspective

How about noise?

$$u = \Phi x^* + n$$

$M \times 1$        $M \times N \ (M < N)$        $N \times 1$   
 $+ \quad \quad \quad +$   
 $M \times 1$        $N \times 1$



**Stability assumption:**  
 $\|\Phi v\| \geq \epsilon \|v\|, \forall v \in T_A(x^*)$

want epsilon large  
to minimize overlap  
between  $\|x^*\|_{A\text{-conv}(A)}$   
and  $\|u - \Phi x\| \leq \sigma$

For this 2D example:  $\|\Phi v\| \geq \|v\| \sin(\varphi) \min_i \|\Phi(i, :)\|$

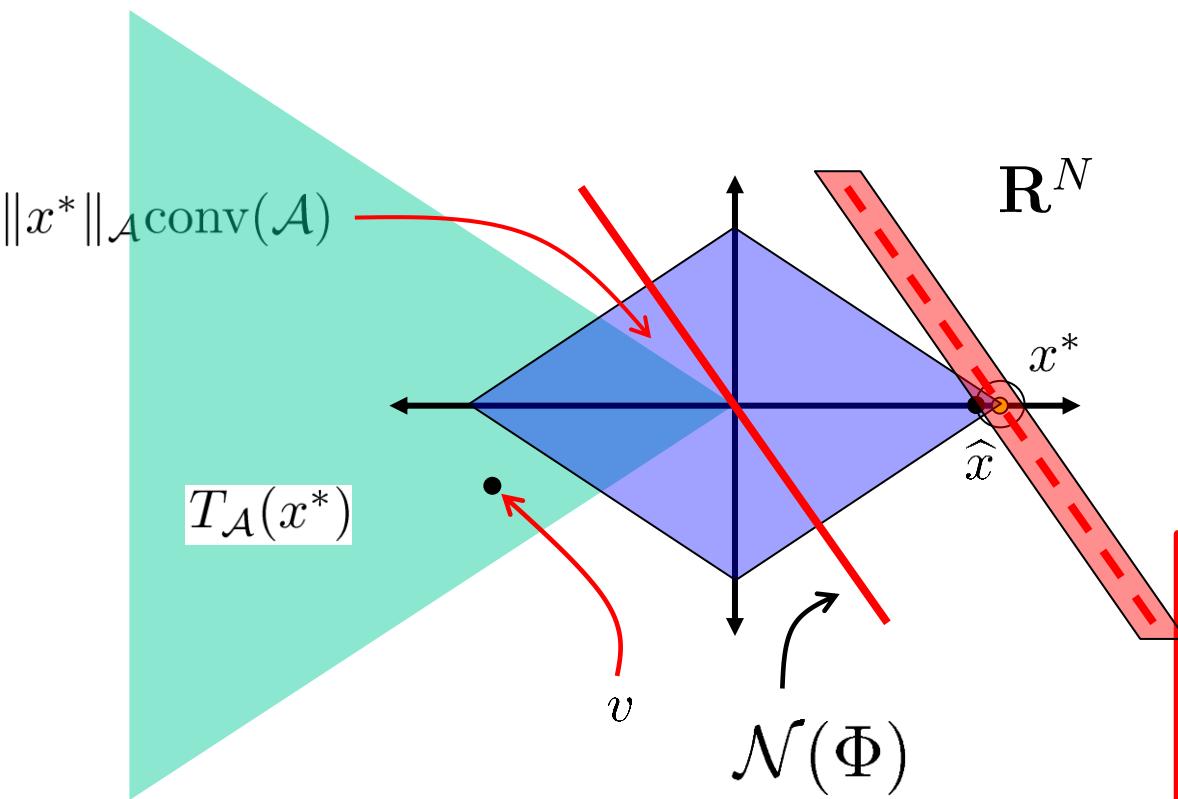
Matlab notation

# Towards algorithms: a geometric perspective

**How about noise?**

$$u = \Phi x^* + n$$

$M \times 1$        $M \times N \ (M < N)$        $N \times 1$



**Stability assumption:**  
 $\|\Phi v\| \geq \epsilon \|v\|, \forall v \in T_{\mathcal{A}}(x^*)$

$$\hat{x} = \arg \min_{x: \|u - \Phi x\| \leq \sigma} \|x\|_{\mathcal{A}}$$

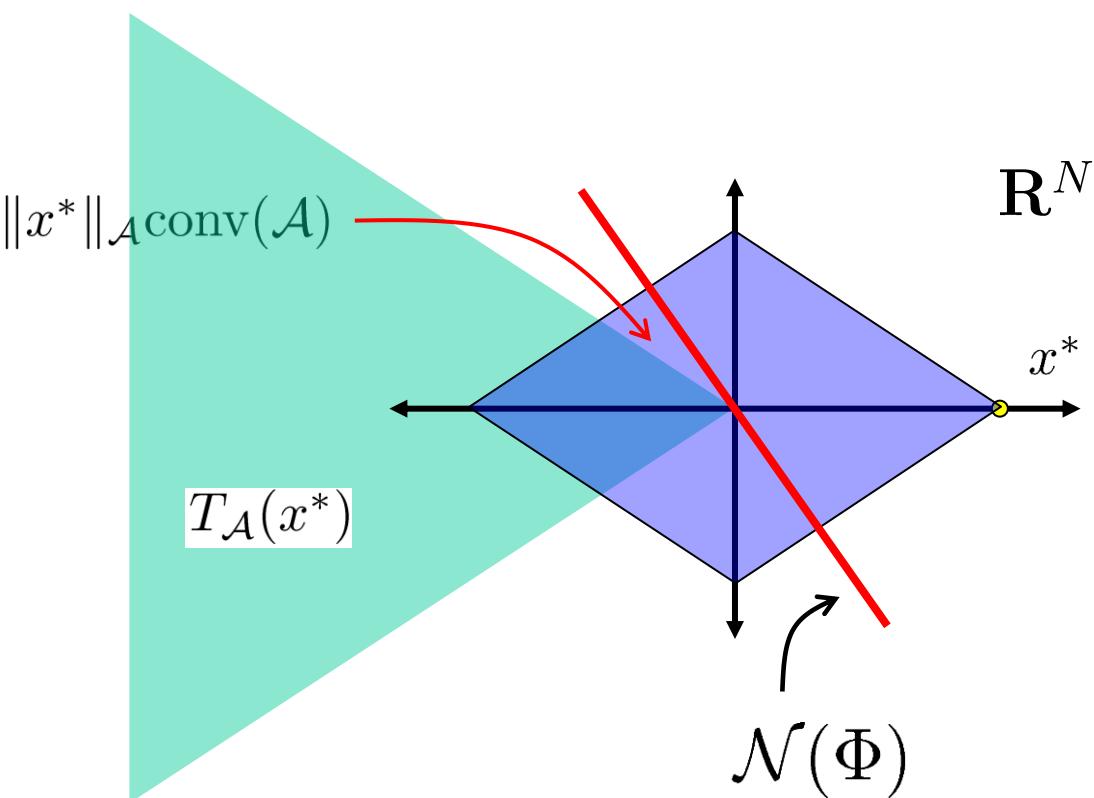
$$\Rightarrow \|x^* - \hat{x}\| \leq \frac{2\sigma}{\epsilon}$$

# Towards algorithms: a geometric perspective

Can we guarantee the following?\*

$$\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) = \{0\}$$

$$u = \begin{matrix} \text{color bar} \\ M \times 1 \end{matrix} = \begin{matrix} \text{color grid} \\ M \times N \quad (M < N) \end{matrix} \Phi \begin{matrix} \text{color bar} \\ N \times 1 \end{matrix} = x^*$$



\*without knowing  $x^*$

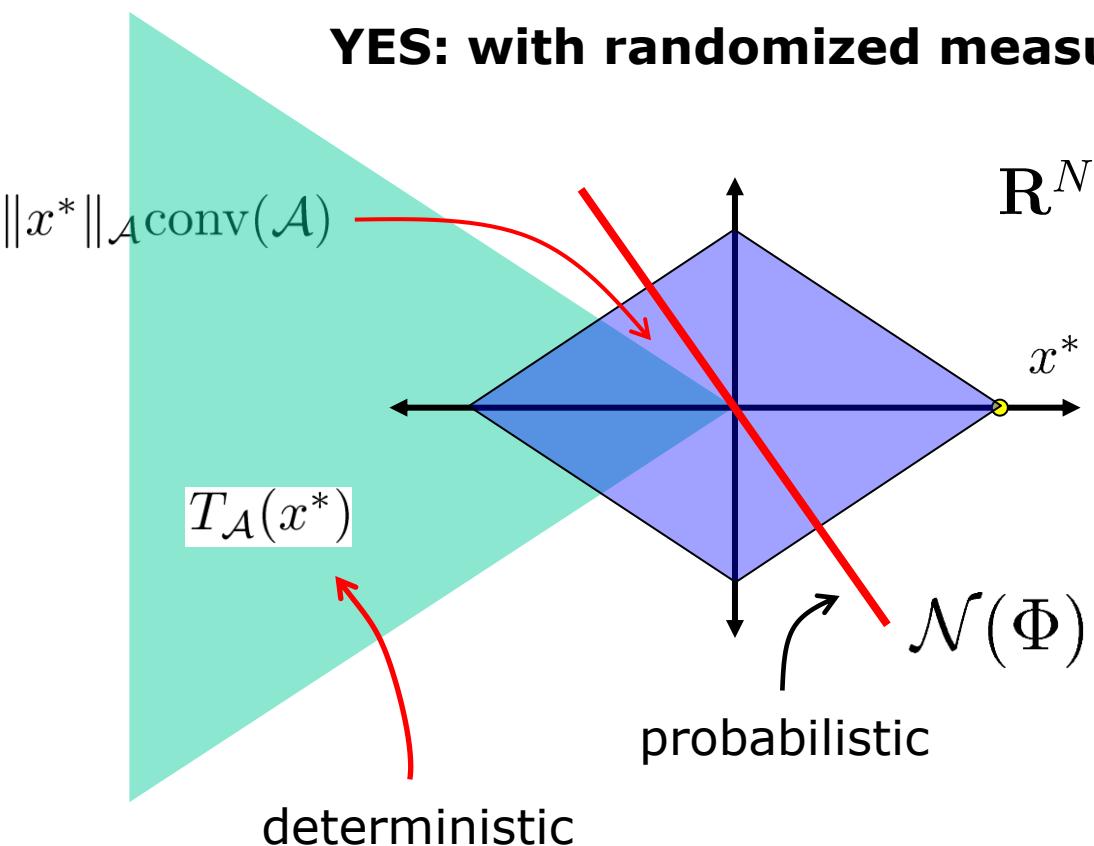
# Towards algorithms: a geometric perspective

Can we guarantee the following?\*

$$\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) = \{0\}$$

$$u = \Phi x^* \\ M \times 1 \quad M \times N \quad (M < N) \\ N \times 1$$

**YES: with randomized measurements!**



Gordon's Minimum Restricted Singular Values Theorem has a probabilistic characterization.

$$\text{Prob}(\min_v \|\Phi v\| \geq \epsilon)$$

$$\forall v \in T_{\mathcal{A}}(x^*), \|v\| = 1$$

[Gordon 1988]

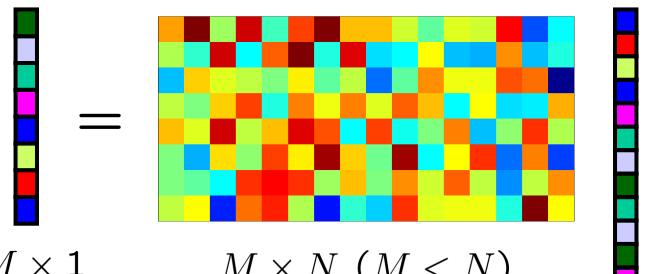
We have a lot more to say on this later.

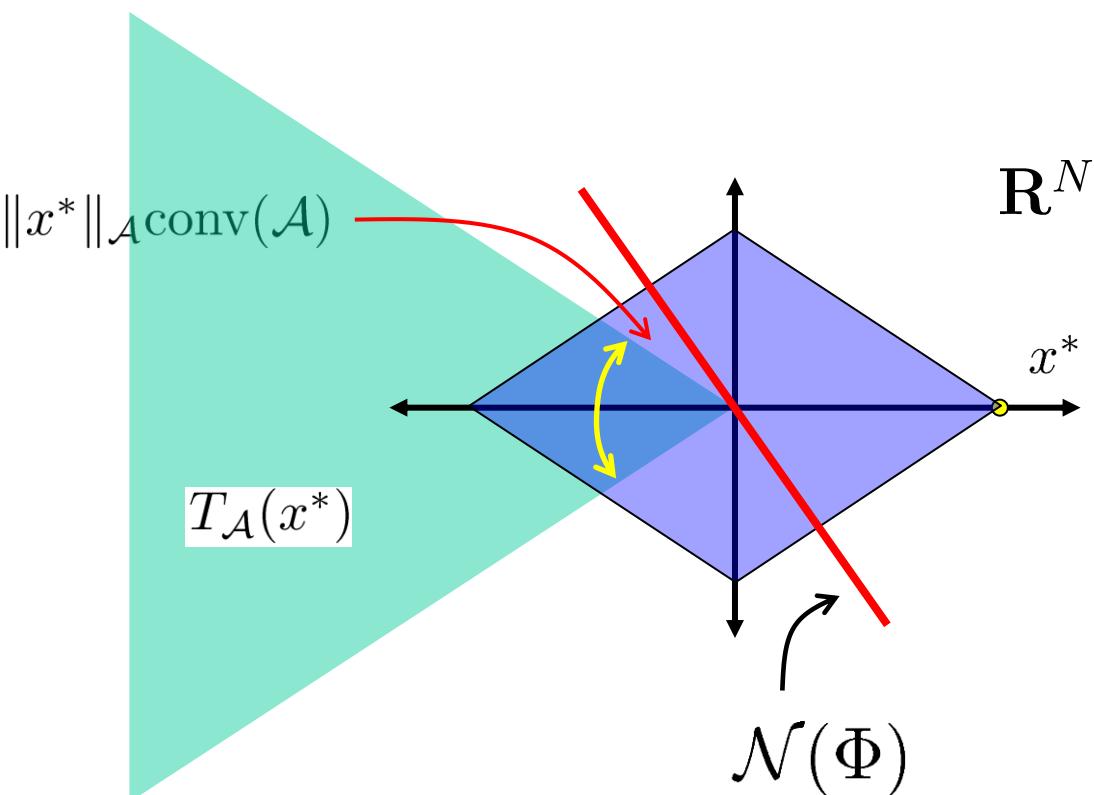
\*without knowing  $x^*$

# Towards algorithms: a geometric perspective

Can we guarantee the following?\*

$$\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) = \{0\}$$

$$u = \Phi x^* \\ M \times 1 \quad M \times N \quad (M < N) \\ N \times 1$$




Gordon's Minimum Restricted Singular Values Theorem has a probabilistic characterization.

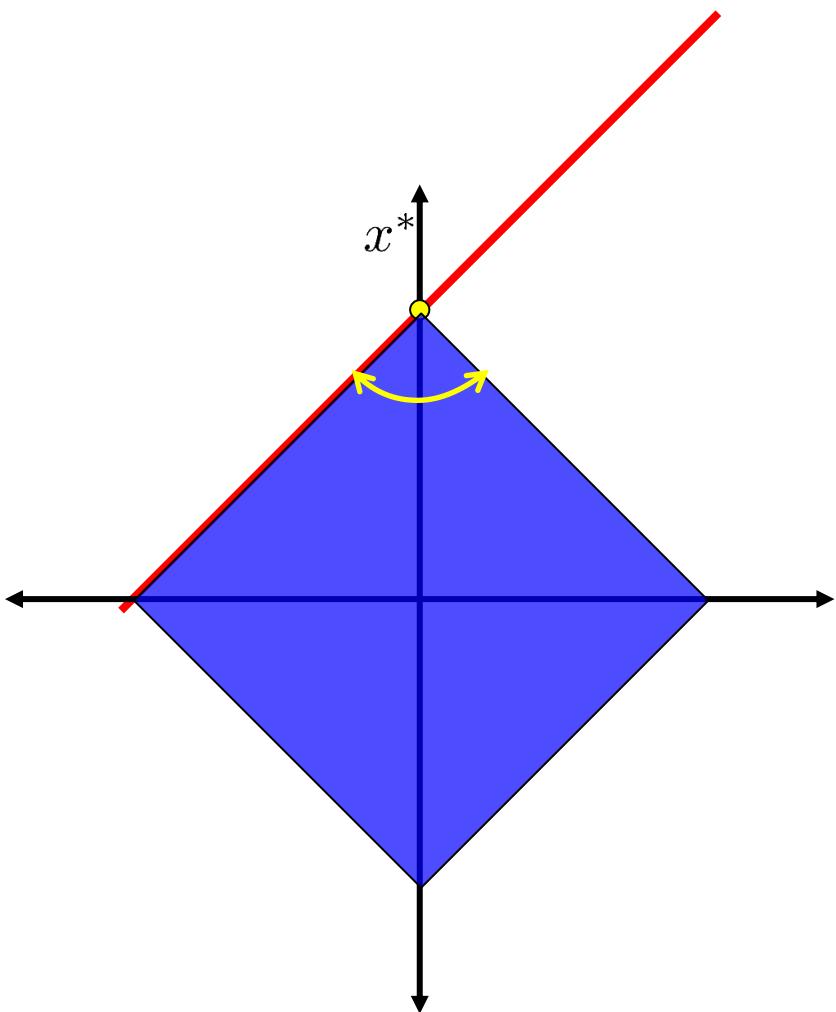
Key concept:  
**width of the tangent cone!**

\*without knowing  $x^*$

# Towards algorithms: a geometric perspective

Can we guarantee the following?\*

$$\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) = \{0\}$$



$$\mathcal{A} = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right\}$$

$$\begin{aligned}\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) &= \{0\} \text{ w.p. } 1/2 \\ \Rightarrow x^* &= \arg \min_{x: u=\Phi x} \|x\|_1\end{aligned}$$

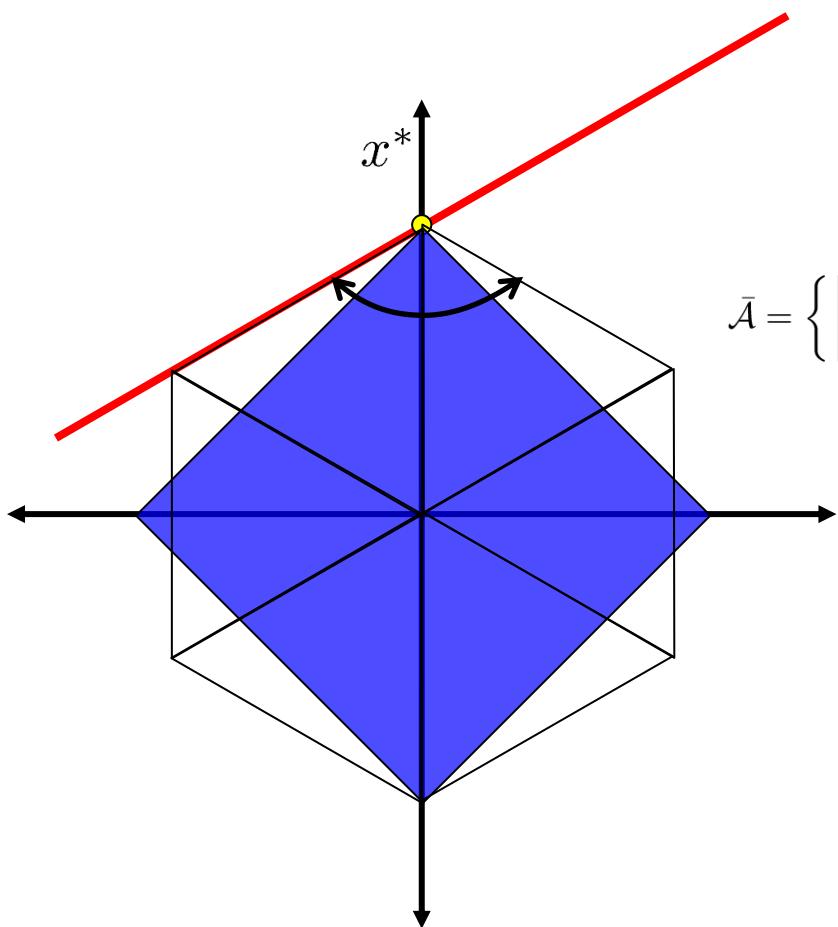
\*without knowing 1-sparse  $x^*$  and 1-random measurement

# Towards algorithms: a geometric perspective

Can we guarantee the following?\*

$$\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) = \{0\}$$

$$\mathcal{A} = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right\}$$



$$\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) = \{0\} \text{ w.p. } 1/2$$

$$\Rightarrow x^* = \arg \min_{x: u = \Phi x} \|x\|_1$$

$$\bar{\mathcal{A}} = \left\{ \begin{bmatrix} \sqrt{3}/2 \\ 1/2 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -\sqrt{3}/2 \\ 1/2 \end{bmatrix}, \begin{bmatrix} -\sqrt{3}/2 \\ -1/2 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \begin{bmatrix} \sqrt{3}/2 \\ -1/2 \end{bmatrix} \right\}$$

$$\mathcal{N}(\Phi) \cap T_{\bar{\mathcal{A}}}(x^*) = \{0\} \text{ w.p. } 1/3$$

$$\Rightarrow x^* = \arg \min_{x: u = \Phi x} \|x\|_{\bar{\mathcal{A}}}$$

\*without knowing 1-sparse  $x^*$  and 1-random measurement

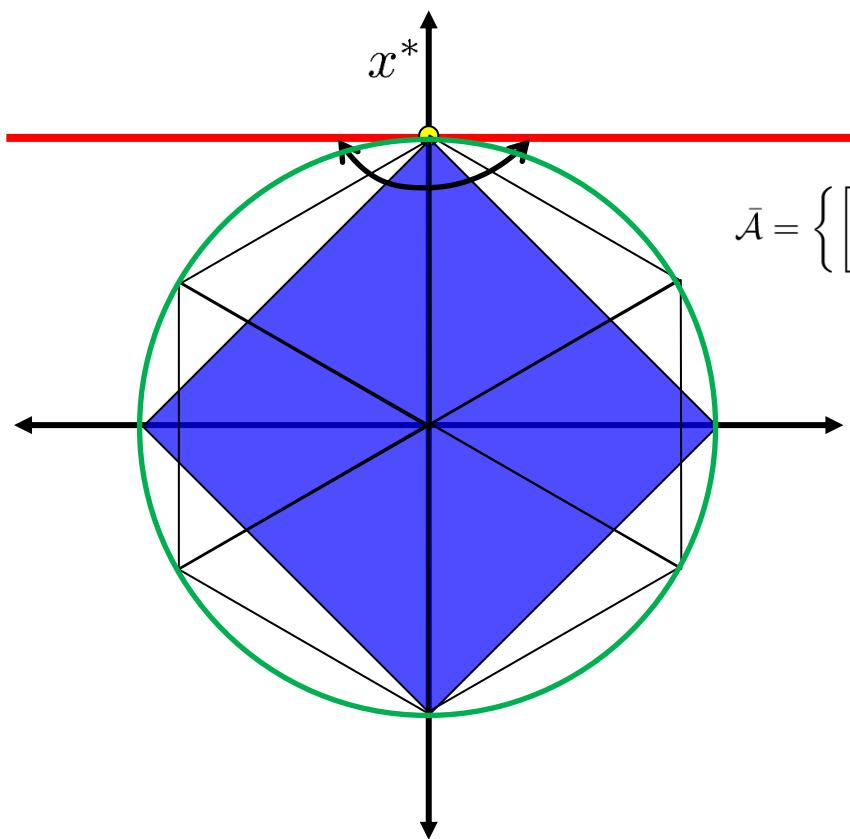
# Towards algorithms: a geometric perspective

Can we guarantee the following?\*

$$\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) = \{0\}$$

$$\mathcal{A} = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right\}$$

$$\begin{aligned} \mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) &= \{0\} \text{ w.p. } 1/2 \\ \Rightarrow x^* &= \arg \min_{x: u=\Phi x} \|x\|_1 \end{aligned}$$



$$\bar{\mathcal{A}} = \left\{ \begin{bmatrix} \sqrt{3}/2 \\ 1/2 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -\sqrt{3}/2 \\ 1/2 \end{bmatrix}, \begin{bmatrix} -\sqrt{3}/2 \\ -1/2 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \begin{bmatrix} \sqrt{3}/2 \\ -1/2 \end{bmatrix} \right\}$$

$$\begin{aligned} \mathcal{N}(\Phi) \cap T_{\bar{\mathcal{A}}}(x^*) &= \{0\} \text{ w.p. } 1/3 \\ \Rightarrow x^* &= \arg \min_{x: u=\Phi x} \|x\|_{\bar{\mathcal{A}}} \end{aligned}$$

$$\tilde{\mathcal{A}} = \{\|x\|_2 = 1\}$$

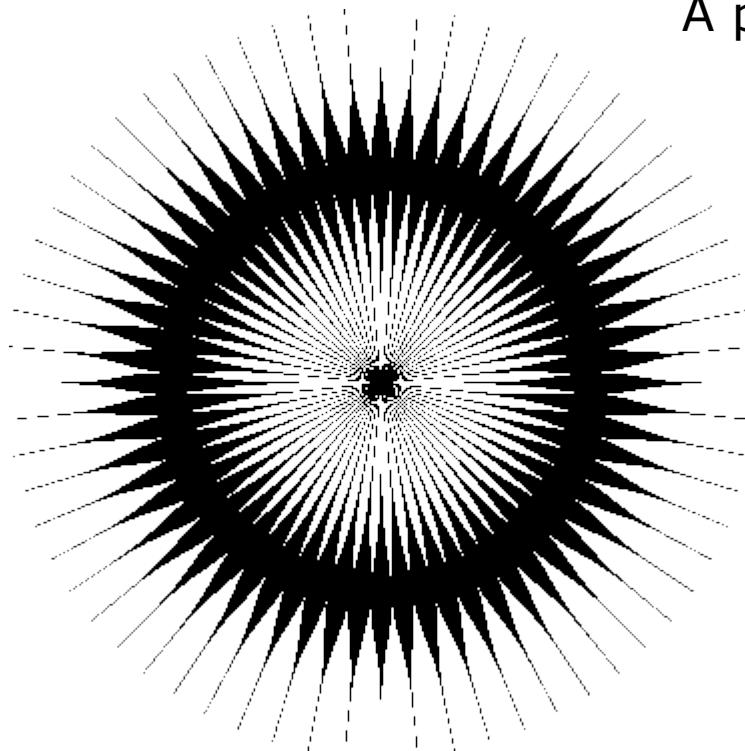
$$\begin{aligned} \mathcal{N}(\Phi) \cap T_{\tilde{\mathcal{A}}}(x^*) &= \{0\} \text{ w.p. } 0 \\ \Rightarrow x^* &= \arg \min_{x: u=\Phi x} \|x\|_2 \end{aligned}$$

\*without knowing 1-sparse  $x^*$  and 1-random measurement

# Towards algorithms: a geometric perspective

Can we guarantee the following?\*

$$\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) = \{0\}$$



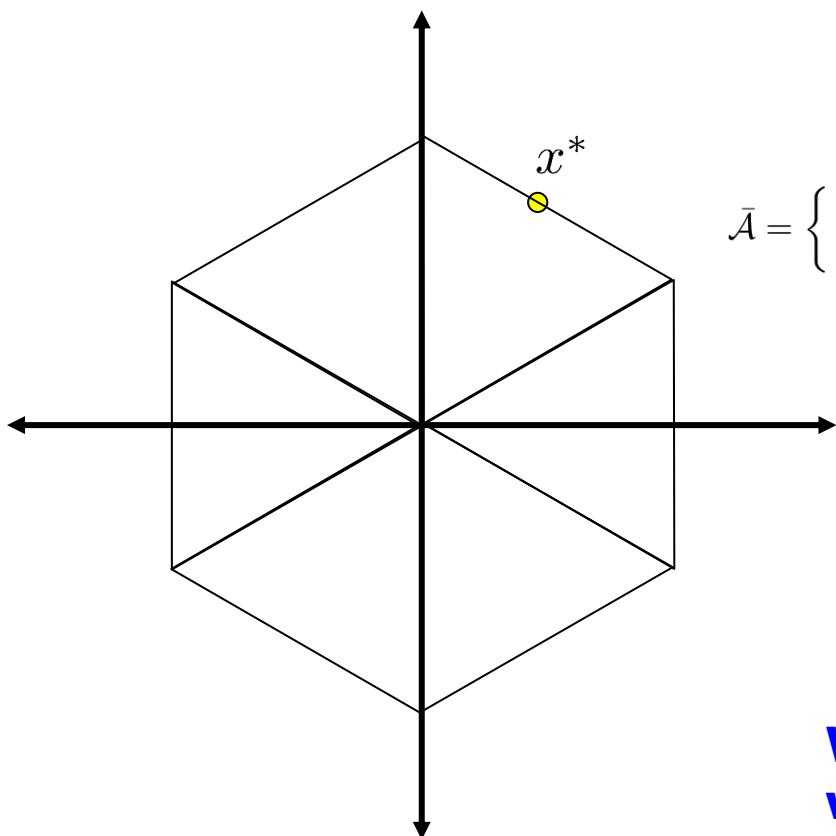
A projected 6D hypercube with 64 vertices

**Blessing-of-dimensionality!**

# Towards algorithms: a geometric perspective

Pop-quiz:

$$\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) = \{0\}$$



$$\bar{\mathcal{A}} = \left\{ \begin{bmatrix} \sqrt{3}/2 \\ 1/2 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -\sqrt{3}/2 \\ 1/2 \end{bmatrix}, \begin{bmatrix} -\sqrt{3}/2 \\ -1/2 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \begin{bmatrix} \sqrt{3}/2 \\ -1/2 \end{bmatrix} \right\}$$

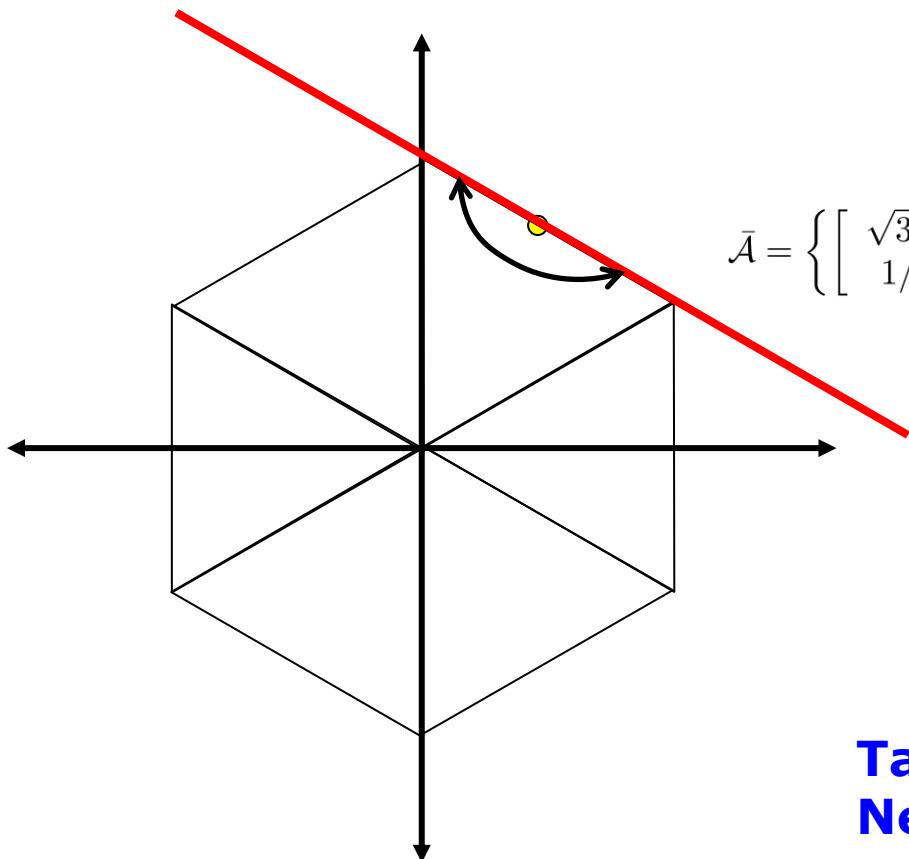
$$\begin{aligned} \mathcal{N}(\Phi) \cap T_{\bar{\mathcal{A}}}(x^*) &= \{0\} \text{ w.p. ???} \\ \Rightarrow x^* &= \arg \min_{x: u = \Phi x} \|x\|_{\bar{\mathcal{A}}} \end{aligned}$$

What is the probability that  
we can determine a 2-sparse  $x^*$   
with 1-random measurement?

# Towards algorithms: a geometric perspective

Pop-answer:

$$\mathcal{N}(\Phi) \cap T_{\mathcal{A}}(x^*) = \{0\}$$



$$\bar{\mathcal{A}} = \left\{ \begin{bmatrix} \sqrt{3}/2 \\ 1/2 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -\sqrt{3}/2 \\ 1/2 \end{bmatrix}, \begin{bmatrix} -\sqrt{3}/2 \\ -1/2 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \begin{bmatrix} \sqrt{3}/2 \\ -1/2 \end{bmatrix} \right\}$$

$$\begin{aligned} \mathcal{N}(\Phi) \cap T_{\bar{\mathcal{A}}}(x^*) &= \{0\} \text{ w.p. 0} \\ \Rightarrow x^* &= \arg \min_{x: u = \Phi x} \|x\|_{\bar{\mathcal{A}}} \end{aligned}$$

**Tangent cone is too wide!  
Need at least 2 measurements!**

# Take home messages

Underlying Model	Atomic Norm	Gaussian Measurements
$K$ -sparse vector in $\mathbb{R}^N$	$\ell_1$ -norm	$(2K + 1) \log(N - K)$
$N \times N$ rank- $R$ matrix	nuclear norm	$3R(2N - R) + 2(N - R - R^2)$
sign vector $\{\pm 1\}^N$	$\ell_\infty$ -norm	$N/2$
$N \times N$ -perm. matrix	Birkoff polytope norm	$9N \log(N)$
$N \times N$ orth. matrix	spectral norm	$(3N^2 - N)/4$

[Chandrasekaran et al. 2010]

convex polytope                     $<>$                     atomic norm

- geometry (and algebra) of representations in **high dimensions**

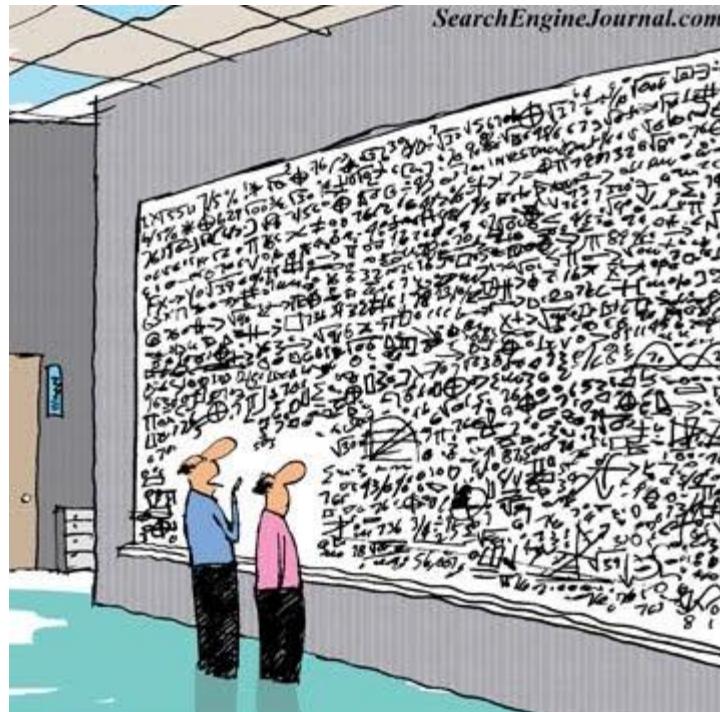
geometric perspective                     $<>$                     convex criteria

- convex optimization algorithms in **high dimensions**

tangent cone width                     $<>$                     # of randomized samples

- probabilistic concentration-of-measures in **high dimensions**

# Convex Algorithms for Low-Dimensional Models



*...And, this is how you solve  
huge-dimensional problems*

# The classical problem templates

$$\|x\|_1 = \sum_{i=1}^N |[x]_i|$$

Criteria seen above have the form

Basis pursuit (BP)

[Chen, Donoho, Saunders, 1998]

$$\min_x \|x\|_1 \quad \text{s.t.} \quad \Phi x = u$$

BP denoising (BPDN):

[Chen, Donoho, Saunders, 1998]

$$\min_x \|x\|_1 \quad \text{s.t.} \quad \|\Phi x - u\|_2^2 \leq \varepsilon$$

Also well known: LASSO (least absolute shrinkage/selection operator):  
[Tibshirani, 1996]

$$\min_x \|\Phi x - u\|_2^2 \quad \text{s.t.} \quad \|x\|_1 \leq \tau$$

All can be written as  $\hat{x} \in \arg \min_{x \in \mathbb{R}^N} f_1(x) + f_2(x)$

# Convex optimization and proximal algorithms

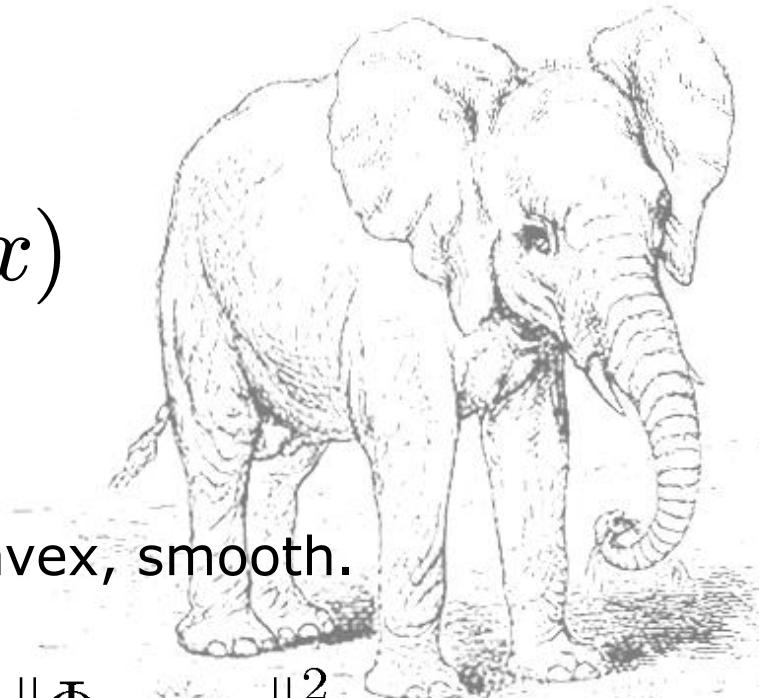
$$\hat{x} \in \arg \min_{x \in \mathbb{R}^N} f_1(x) + f_2(x)$$

$f_1 : \mathbb{R}^N \rightarrow \mathbb{R}$  data fidelity term; convex, smooth.

typically:  $f_1(x) = \frac{1}{2} \|\Phi x - u\|_2^2$

$f_2 : \mathbb{R}^N \rightarrow \bar{\mathbb{R}} = \mathbb{R} \cup \{+\infty\}$  Convex regularizer  
(maybe non-smooth; e.g.  $\ell_1$ )  
(non-convex, later...).

Difficulties: **non-smoothness** and **large dimension** ( $N \gg 1$ )



# Constrained vs unconstrained formulations

Constrained optimization formulations

$$\begin{array}{ll} \widehat{x} \in \arg \min_{x \in \mathbb{R}^N} f_1(x) & \widehat{x} \in \arg \min_{x \in \mathbb{R}^N} f_2(x) \\ \text{s.t. } h(x) \leq \nu & \text{s.t. } g(x) \leq \tau \end{array} \quad (*)$$

can be written as  $\widehat{x} \in \arg \min_{x \in \mathbb{R}^N} f_1(x) + f_2(x)$

...using indicator functions:  $\iota_S(x) = \begin{cases} 0 & \Leftarrow x \in S \\ +\infty & \Leftarrow x \notin S \end{cases}$

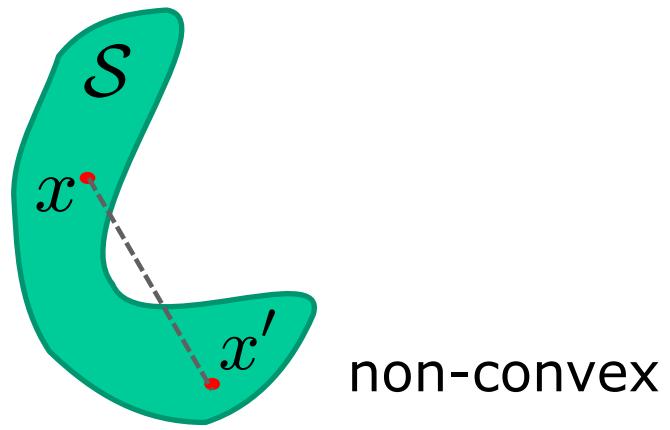
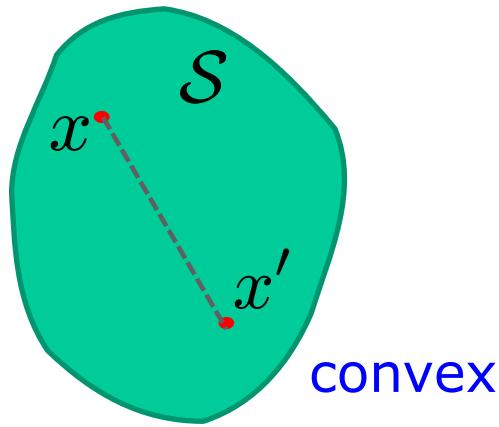
Example: (\*) same as

$$\widehat{x} \in \arg \min_{x \in \mathbb{R}^N} f_1(x) + \iota_{\{x:g(x) \leq \nu\}}(x)$$

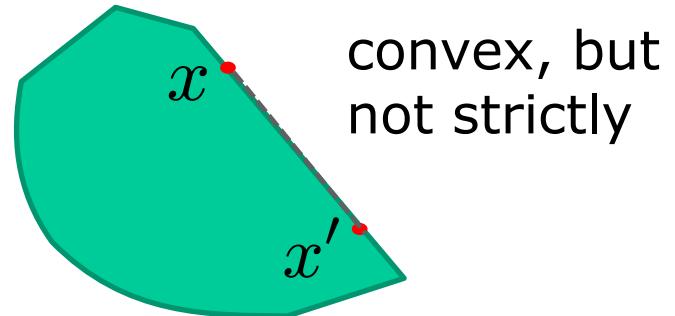
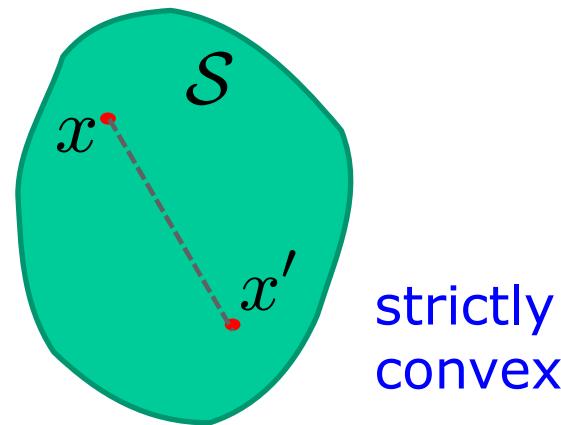
Classical example: the LASSO:  $\min \|\Phi x - u\|_2^2$  s.t.  $\|x\|_1 \leq \tau$

# Convex and strictly convex sets

$\mathcal{S}$  is **convex** if  $x, x' \in \mathcal{S} \Rightarrow \forall \lambda \in [0, 1] \ \lambda x + (1 - \lambda)x' \in \mathcal{S}$



$\mathcal{S}$  is **strictly convex** if  $x, x' \in \mathcal{S} \Rightarrow \forall \lambda \in (0, 1) \ \lambda x + (1 - \lambda)x' \in \text{int}(\mathcal{S})$



# Convex and strictly convex functions

Extended real valued function:  $f : \mathbb{R}^N \rightarrow \bar{\mathbb{R}} = \mathbb{R} \cup \{+\infty\}$

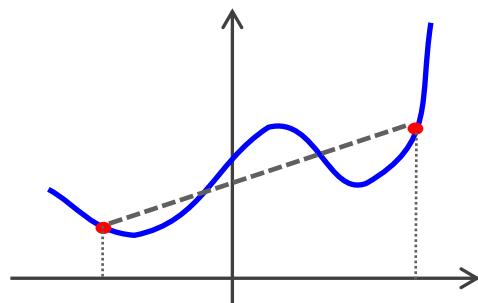
Domain of a function:  $\text{dom}(f) = \{x : f(x) \neq +\infty\}$

$f$  is a **convex function** if

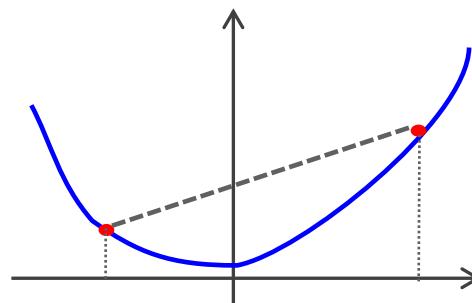
$$\forall \lambda \in [0, 1], x, x' \in \text{dom}(f) \quad f(\lambda x + (1 - \lambda)x') \leq \lambda f(x) + (1 - \lambda)f(x')$$

$f$  is a **strictly convex function** if

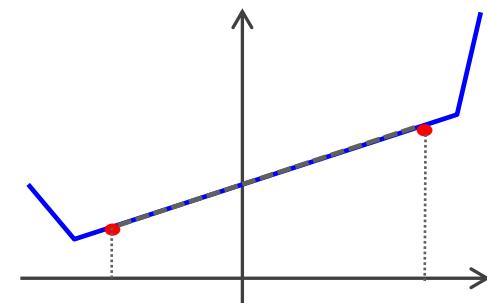
$$\forall \lambda \in (0, 1), x, x' \in \text{dom}(f) \quad f(\lambda x + (1 - \lambda)x') < \lambda f(x) + (1 - \lambda)f(x')$$



non-convex



convex  
strictly convex



convex, not strictly

# Convexity, coercivity, and minima

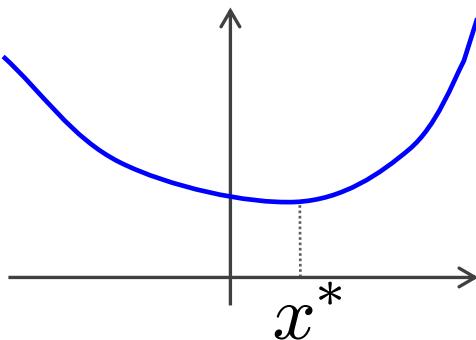
$$f : \mathbb{R}^N \rightarrow \bar{\mathbb{R}} = \mathbb{R} \cup \{+\infty\}$$

$f$  is coercive if  $\lim_{\|x\| \rightarrow +\infty} f(x) = +\infty$

if  $f$  is coercive, then  $G \equiv \arg \min_x f(x)$  is a non-empty set

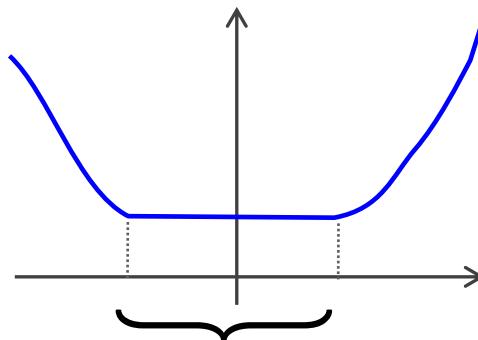
if  $f$  is strictly convex, then  $G$  has at most one element

coercive and  
strictly convex



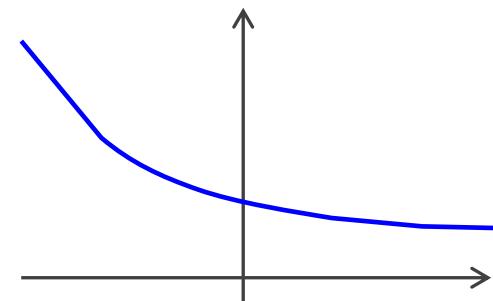
$$G = \{x^*\}$$

coercive, not  
strictly convex



$$G$$

convex, not  
coercive



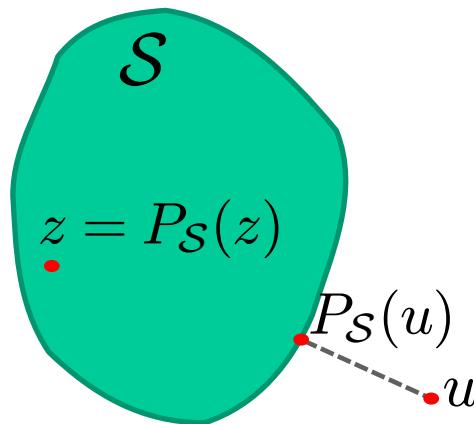
$$G = \emptyset$$

# Euclidean projections on convex sets

Our problem:  $\hat{x} \in \arg \min_{x \in \mathbb{R}^n} f_1(x) + f_2(x)$

consider  $f_2(x) = \iota_S(x) = \begin{cases} 0 & \Leftarrow x \in S \\ +\infty & \Leftarrow x \notin S \end{cases}$   
(convex if  $S$  is convex)

and  $f_1(x) = \frac{1}{2} \|u - x\|_2^2$  (strictly convex)



$$\begin{aligned}\hat{x} &= \arg \min_{x \in \mathbb{R}^n} f_1(x) + f_2(x) \\ &= \arg \min_{x \in S} \|u - x\|_2^2 \\ &\equiv P_S(u) \quad (\text{Euclidean projection})\end{aligned}$$

# Projected gradient algorithm

Our problem:  $\hat{x} \in \arg \min_{x \in \mathbb{R}^n} f_1(x) + f_2(x)$

with  $f_2(x) = \iota_{\mathcal{S}}(x)$  ( $\mathcal{S}$  is a convex set)

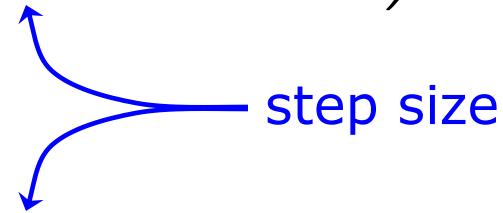
and  $f_1$  some function, e.g.,  $f_1(x) = \frac{1}{2} \|\Phi x - u\|_2^2$

Projected gradient algorithm:

$$x_{k+1} = P_{\mathcal{S}} \left( x_k - \beta_k \nabla f_1(x_k) \right)$$

if  $f_1(x) = \frac{1}{2} \|\Phi x - u\|_2^2$

$$x_{k+1} = P_{\mathcal{S}} \left( x_k - \beta_k \Phi^T (\Phi x_k - u) \right)$$

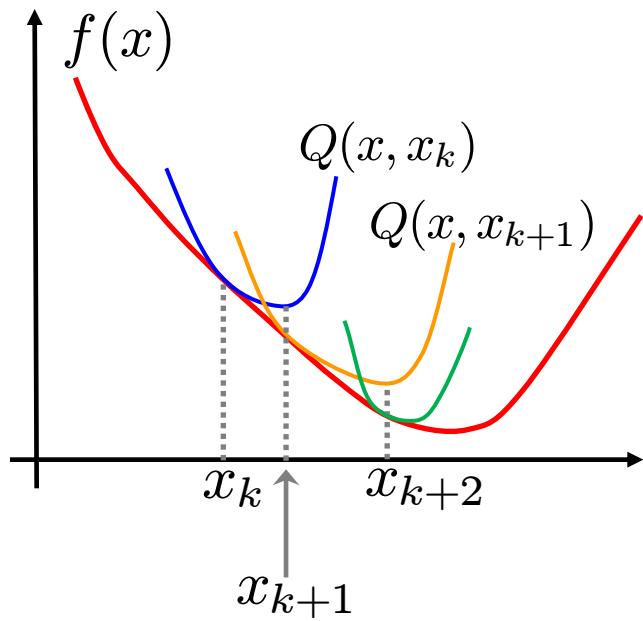


# Detour: majorization-minimization (MM)

Problem:  $\widehat{x} \in \arg \min_{x \in \mathbb{R}^n} f(x)$

$Q(x, x_k)$  is a **majorizer** of  $f$  at  $x_k$

$$Q(x, x_k) \geq f(x), \quad Q(x_k, x_k) = f(x_k)$$



MM algorithm:

$$x_{k+1} = \arg \min_x Q(x, x_k)$$

monotonicity:

$$\begin{aligned} f(x_{k+1}) &\leq Q(x_{k+1}, x_k) \\ &\leq Q(x_k, x_k) \\ &= f(x_k) \end{aligned}$$

# Projected gradient from majorization-minimization

Our problem:  $\hat{x} \in \arg \min_{x \in \mathbb{R}^n} f_1(x) + f_2(x)$

with  $f_2(x) = \iota_{\mathcal{S}}(x)$  ( $\mathcal{S}$  is a convex set)

and  $f_1$  has  $L$ -Lipschitz gradient

$$\|\nabla f_1(x) - \nabla f_1(x')\| \leq L\|x - x'\|$$

e.g.  $f_1(x) = \frac{1}{2}\|\Phi x - u\|_2^2 \Rightarrow L = \lambda_{\max}(\Phi^T \Phi) = \|\Phi\|_2^2$

Hessian of  $f_1$  

...a separable approximation of  $f_1$

$$Q(x, x_k) = f_1(x_k) + (x - x_k)^T \nabla f_1(x_k) + \frac{1}{2\beta_k} \|x - x_k\|_2^2$$

# Projected gradient from majorization-minimization

Our problem:  $\hat{x} \in \arg \min_{x \in \mathbb{R}^n} f_1(x) + \iota_S(x)$

Separable approximation of  $f_1$

$$Q(x, x_k) = f_1(x_k) + (x - x_k)^T \nabla f_1(x_k) + \frac{1}{2\beta_k} \|x - x_k\|_2^2$$

$Q(x, x_k)$  is a majorizer of  $f_1$ , if  $\boxed{\beta_k < \frac{1}{L}}$

$Q(x, x_k) + \iota_S(x)$  is a majorizer  $f_1(x) + \iota_S(x)$

MM algorithm:

$$\begin{aligned} x_{k+1} &= \arg \min_x Q(x, x_k) + \iota_S(x) \\ &= \arg \min_x \frac{1}{2\beta_k} \left\| x - (x_k - \beta_k \nabla f_1(x_k)) \right\|_2^2 + \iota_S(x) \\ &= P_S \left( x_k - \beta_k \nabla f_1(x_k) \right) \quad \text{...projected gradient.} \end{aligned}$$

# Proximity operators

Our problem:  $\hat{x} \in \arg \min_{x \in \mathbb{R}^n} f_1(x) + f_2(x)$

with  $f_2$  a convex function

and  $f_1(x) = \frac{1}{2} \|u - x\|_2^2$  (strictly convex)

$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|u - x\|_2^2 + f_2(x) \equiv \text{prox}_{f_2}(u)$$

Proximity operator [Moreau 62], [Combettes 01].

Generalizes the notion of Euclidean projection.

# Proximity operators (linear)

$$\text{prox}_f(u) = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|u - x\|_2^2 + f(x) \quad (\mathbb{R}^N \rightarrow \mathbb{R}^N)$$

Classical cases: squared  $\ell_2$  regularizer  $f(x) = \frac{\tau}{2} \|x\|_2^2$

$$\text{prox}_f(u) = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|u - x\|_2^2 + \frac{\tau}{2} \|x\|_2^2 = \frac{u}{1 + \tau}$$

squared  $\ell_2$  regularizer with “analysis” operator  $f(x) = \frac{\tau}{2} \|Dx\|_2^2$

$$\begin{aligned} \text{prox}_f(u) &= \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|u - x\|_2^2 + \frac{\tau}{2} \|Dx\|_2^2 \\ &= (I + \tau D^T D)^{-1} u \end{aligned}$$

if  $D$  is a circulant matrix,  $O(N \log N)$  cost using the FFT

# Proximity operator of the $\ell_1$ norm



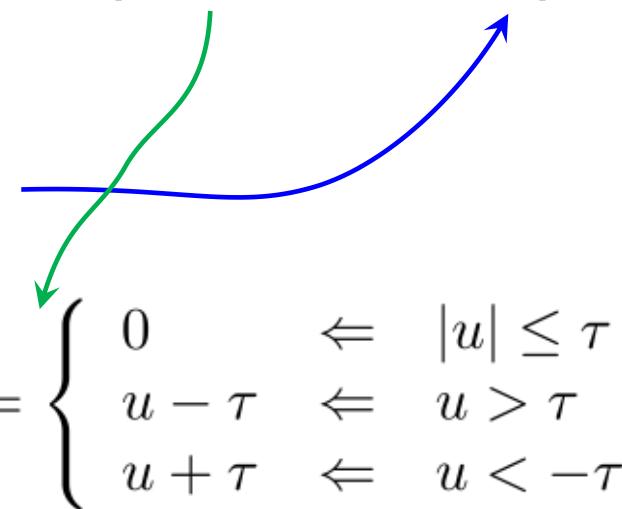
$$\text{prox}_{\tau \|\cdot\|_1}(u) = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|u - x\|_2^2 + \tau \|x\|_1$$

**Separable:** solve w.r.t. each component:  $\min_x \tau|x| + 0.5(x - u)^2$

Possible approach: write  $|x| = \max_{|z| \leq 1} zx$

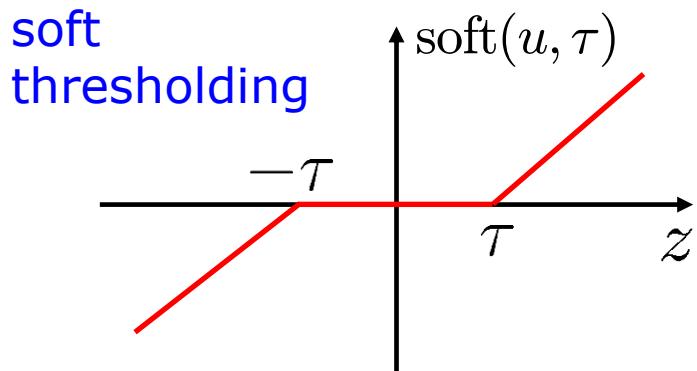
$$\begin{aligned} \min_x \max_{|z| \leq 1} \tau zx + 0.5(x - u)^2 &= \max_{|z| \leq 1} \min_x \tau zx + 0.5(x - u)^2 \\ &= \max_{|z| \leq 1} -0.5\tau^2 z^2 + \tau zu \quad (\text{for } x = u - \tau z) \end{aligned}$$

$$\arg \max_{|z| \leq 1} -0.5\tau^2 z^2 + \tau zu = \begin{cases} u/\tau & \Leftarrow |u| \cdot \tau \\ 1 & \Leftarrow u > \tau \\ -1 & \Leftarrow u < -\tau \end{cases}$$



# Proximity operator of the $\ell_1$ norm: the “soft”

$$\arg \min_x \tau|x| + 0.5(x - u)^2 = \begin{cases} 0 & \Leftarrow |u| \leq \tau \\ u - \tau & \Leftarrow u > \tau \\ u + \tau & \Leftarrow u < -\tau \end{cases}$$



$$= \text{sign}(u) \max\{0, |u| - \tau\}$$
$$\equiv \text{soft}(u, \tau) = \text{prox}_{\tau|\cdot|}$$

(for vectors,  $\text{soft}(u, \tau)$  is applied component-wise)

$p$ -th power of  $\ell_p$  norms  $\|x\|_p^p = \sum_i |[x]_i|^p$

closed form prox for  $p \in \left\{1, 2, \frac{4}{3}, \frac{3}{2}, 3, 4\right\}$

[Combettes, Wajs, 2005]

# Dual norms, proximity operators, and projections

Dual norm: some norm,  $\|\cdot\| : \mathbb{R}^N \rightarrow \mathbb{R}_+$

its dual norm:  $\|x\|^* = \max_{\|z\| \leq 1} \langle x, z \rangle$

Dual norm of  $\|\cdot\|_p$  is  $\|\cdot\|_q$ , where  $\frac{1}{p} + \frac{1}{q} = 1$  Hölder conjugates

... simple corollary of Hölder's inequality:  $x^T z \leq \|x\|_p \|z\|_q$

Examples of Hölder conjugates:  $(2, 2)$ ,  $(1, +\infty)$ ,  $(3/2, 3)$ , ...

These concepts are related through:

$$\text{prox}_{\|\cdot\|}(u) = u - P_{\{x: \|x\|^* \leq 1\}}(u)$$

[Combettes, Wajs, 2005]

# Dual norms, proximity operators, and projections

$$\text{prox}_{\tau \|\cdot\|}(u) = u - P_{\{x: \|x\|^* \leq \tau\}}(u)$$

This relation underlies our earlier derivation of  $\text{prox}_{\|\cdot\|_1}$

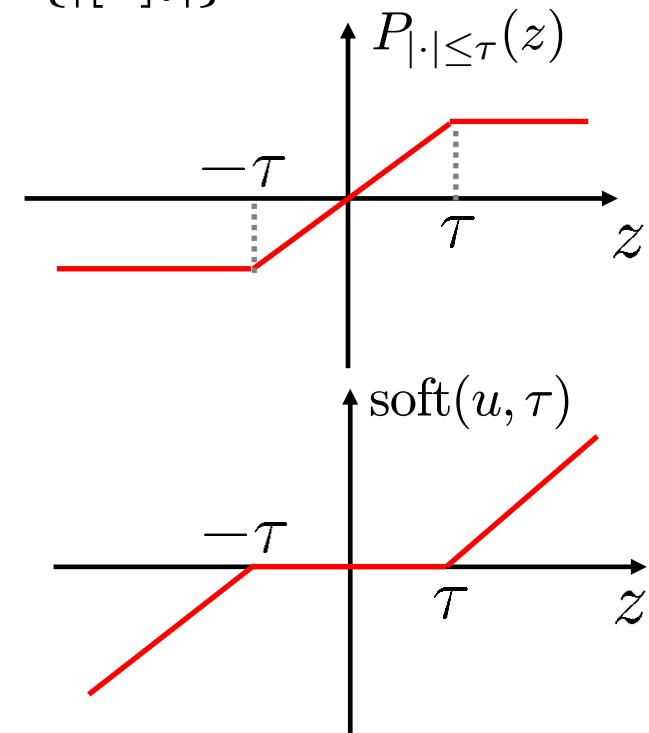
$$\text{prox}_{\tau \|\cdot\|_1}(u) = u - P_{\{x: \|x\|_\infty \leq \tau\}}(u)$$

$$\|x\|_\infty = \max\{|[x]_i|\}$$

It's all separable,

$$\text{prox}_{\tau | \cdot |}(u) = u - P_{\{x: |x| \leq \tau\}}(u)$$

$$= u - \begin{cases} u & \Leftarrow |u| \leq \tau \\ -\tau & \Leftarrow u < -\tau \\ \tau & \Leftarrow u > \tau \end{cases}$$
$$= \text{soft}(u, \tau)$$



# Dual norms, proximity operators, and projections

$$\text{prox}_{\tau \|\cdot\|}(u) = u - P_{\{x: \|x\|^* \leq \tau\}}(u)$$

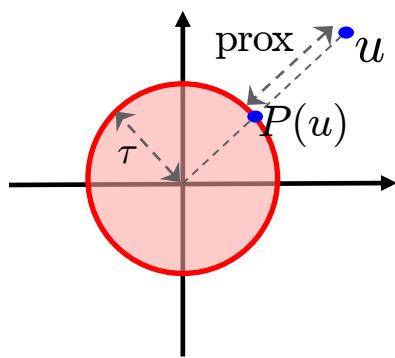
This relation allows deriving  $\text{prox}_{\|\cdot\|_\infty}$  and  $\text{prox}_{\|\cdot\|_2}$

$$\text{prox}_{\|\cdot\|_\infty}(u) = u - P_{\{x: \|x\|_1 \leq \tau\}}(u)$$

projection on the  $\ell_1$  ball of radius  $\tau$

$O(n \log n)$

$$\text{prox}_{\|\cdot\|_2}(u) = u - P_{\{x: \|x\|_2 \leq \tau\}}(u)$$



$$= u - \begin{cases} u & \Leftarrow \|u\|_2 \leq \tau \\ \tau u / \|u\|_2 & \Leftarrow \|u\|_2 > \tau \end{cases}$$

$$= \frac{u}{\|u\|_2} \max\{0, \|u\|_2 - \tau\}$$

vector soft thresholding

# Proximity operators of atomic norms

$$\text{prox}_{\tau \|\cdot\|}(u) = u - P_{\{x: \|x\|^* \leq \tau\}}(u)$$

These relation allows deriving prox operators of **atomic norms**:

$$\|x\|_{\mathcal{A}} = \inf\{t > 0 : x \in t \text{ conv}(\mathcal{A})\}$$

The dual of an atomic norm ball:

$$\begin{aligned}\|x\|_{\mathcal{A}}^* &= \max_{\|z\|_{\mathcal{A}} \leq 1} \langle z, x \rangle = \max_{z \in \text{conv}(\mathcal{A})} \langle z, x \rangle \\ &= \max\{\langle a, x \rangle, a \in \mathcal{A}\}\end{aligned}$$

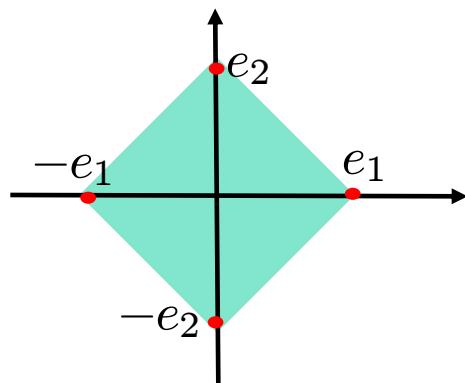
$$P_{\{x: \|x\|_{\mathcal{A}}^* \leq \tau\}}(u) = \arg \min_{\langle a, x \rangle \leq \tau, \forall a \in \mathcal{A}} \|u - x\|_2^2$$

$$\text{prox}_{\tau \|\cdot\|_{\mathcal{A}}}(u) = u - \arg \min_{\langle a, x \rangle \leq \tau, \forall a \in \mathcal{A}} \|u - x\|_2^2$$

# Proximity operators of atomic norms: $\ell_1$

Deriving  $\text{prox}_{\tau \|\cdot\|_1}$  from the atomic norm view

$$\|x\|_1 = \|x\|_{\mathcal{A}}$$



$$\begin{aligned}\mathcal{A} &= \left\{ \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ \vdots \\ -1 \end{bmatrix} \right\} \\ &= \{e_1, e_2, \dots, e_N, -e_1, \dots, -e_N\}\end{aligned}$$

$$|\mathcal{A}| = 2N$$

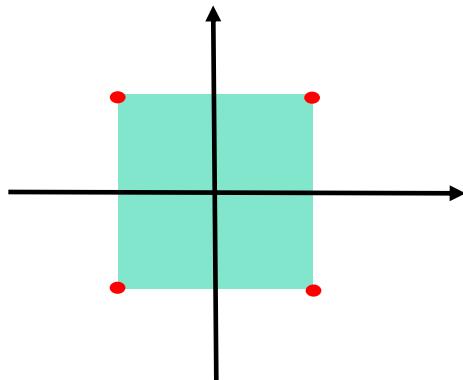
$$\|x\|_{\mathcal{A}}^* = \max\{\langle a, x \rangle, a \in \mathcal{A}\} = \max\{|[x]_i|\} = \|x\|_{\infty}$$

$$\begin{aligned}\text{prox}_{\tau \|\cdot\|_1}(u) &= u - P_{\{x: \|x\|_{\infty} \leq \tau\}}(u) \\ &= \text{soft}(x, \tau)\end{aligned}$$

# Proximity operators of atomic norms: $\ell_\infty$

Deriving  $\text{prox}_{\tau \|\cdot\|_\infty}$  from the atomic norm view

$$\|x\|_\infty = \|x\|_{\mathcal{A}}$$



$$\mathcal{A} = \left\{ \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ -1 \\ \vdots \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ -1 \\ \vdots \\ 1 \end{bmatrix}, \dots, \begin{bmatrix} -1 \\ -1 \\ \vdots \\ -1 \end{bmatrix} \right\}$$

$$= \{-1, +1\}^N$$

$$|\mathcal{A}| = 2^N$$

$$\|x\|_{\mathcal{A}}^* = \max\{\langle a, x \rangle, a \in \mathcal{A}\} = \sum_{i=1}^N |[x]_i| = \|x\|_1$$

$$\text{prox}_{\tau \|\cdot\|_\infty}(u) = u - P_{\{x: \|x\|_1 \leq \tau\}}(u)$$

# Proximity of atomic norms: matrix nuclear norm

Matrix nuclear norm:  $\|X\|_* = \sum_i \sigma_i(X) = \sum_i \sqrt{\lambda_i(X^T X)}$

$$\|X\|_* = \|X\|_{\mathcal{A}} \quad \mathcal{A} = \{Z : \text{rank}(Z) = 1, \|Z\|_F = 1\}$$

$$\text{rank}(Z) = |\{\sigma_i(Z) \neq 0\}|$$

Frobenius norm  $\|Z\|_F^2 = \sum_{ij} [Z]_{ij}^2 = \sum_i \sigma_i^2(Z)$

$$\|X\|_{\mathcal{A}}^* = \max\{\langle Z, X \rangle, Z \in \mathcal{A}\}$$

$$= \max \left\{ \sum_i \sigma_i(Z) \sigma_i(X), \text{rank}(Z) = 1, \sum_i \sigma_i^2(Z) = 1 \right\}$$

$$= \sigma_{\max}(X) = \|X\|_2 \quad \text{spectral norm}$$

# Proximity of atomic norms: matrix nuclear norm

Euclidean matrix projection:  $P_{\mathcal{S}}(X) = \arg \min_{Z \in \mathcal{S}} \|Z - X\|_F^2$

Note: for any unitary matrix  $U$  ( $U^T U = I, UU^T = I$ )

$$\|UM\|_F^2 = \text{trace}(M^T U^T UM) = \text{trace}(M^T A) = \|M\|_F^2$$

$$\text{prox}_{\tau \|\cdot\|_*}(X) = X - P_{\{Z: \|Z\|_2 \leq \tau\}}(X)$$

singular value  
diagonal matrix

$$= U\Lambda V^T - P_{\{Z: \sigma_{\max}(Z) \leq \tau\}}(U\Lambda V^T)$$

[Lewis, Malick, 2009]

$$= U \text{diag}(\text{diag}(\Lambda) - P_{\{x: \|x\|_\infty \leq \tau\}}(\text{diag}(\Lambda))) V^T$$

$$= U \text{soft}(\Lambda, \tau) V^T \quad \text{singular value thresholding (svt)}$$

# Proximity of atomic norms: matrix spectral norm

Matrix spectral norm:  $\|X\|_2 = \sigma_{\max}(X)$

$$\|X\|_2 = \|X\|_{\mathcal{A}} \quad \mathcal{A} = \{Z : Z^T Z = I\} = \{Z : \sigma_i(Z) = 1, \forall_i\}$$

orthogonal matrices

$$\|X\|_{\mathcal{A}}^* = \max\{\langle Z, X \rangle, Z \in \mathcal{A}\}$$

$$= \max \left\{ \sum_i \sigma_i(Z) \sigma_i(X), \sigma_i(Z) = 1, \forall_i \right\}$$

$$= \sum_i \sigma_i(X) = \|X\|_* \quad \text{nuclear norm}$$

# Proximity of atomic norms: matrix spectral norm



$$\text{prox}_{\tau \|\cdot\|_2}(X) = X - P_{\{Z: \|Z\|_* \leq \tau\}}(X)$$

$$= U \Lambda V^T - P_{\{Z: \|Z\|_* \leq \tau\}}(U \Lambda V^T)$$

singular value  
diagonal matrix

$$= U \left( \Lambda - P_{\{Z: \sum_i \sigma_i(Z) \leq \tau\}}(\Lambda) \right) V^T$$

$$= U \text{diag}(\text{diag}(\Lambda) - P_{\{x: \|x\|_1 \leq \tau\}}(\text{diag}(\Lambda))) V^T$$

residual of projection of the singular  
values on an  $\ell_1$  ball of radius  $\tau$

# Proximity and atomic sets: vectors vs matrices

vectors			matrices		
norm	prox	atomic set	norm	prox	atomic set
$\ell_1$ $\ x\ _1$	component soft thresholding	$\mathcal{A} = \{\pm e_i\}$ $ \mathcal{A}  = 2N$	nuclear $\ X\ _*$	singular value thresholding	$\mathcal{A} = \text{set of all rank 1, norm 1 matrices}$
$\ell_\infty$ $\ x\ _\infty$	residual of projection on $\ell_1$ ball	$\mathcal{A} = \{\pm 1\}^N$ $ \mathcal{A}  = 2^N$	spectral $\ X\ _2$	residual of s.v. proj. on $\ell_1$ ball	$\mathcal{A} = \text{set of all orthogonal matrices}$
$\ell_2$ $\ x\ _2$	vector soft thresholding	$\mathcal{A} = \text{set of all vectors with norm 1}$ $ \mathcal{A}  = \infty$	Frobenius $\ X\ _F$	matrix soft threshold.	$\mathcal{A} = \text{all matrices of unit Frobenius norm.}$

# Proximal algorithms

Back to the problem:  $\widehat{x} \in \arg \min_{x \in \mathbb{R}^n} f_1(x) + f_2(x)$

with  $f_2$  a proper convex function

and  $f_1$  has a  $L$ -Lipschitz gradient; e.g.  $f_1(x) = \frac{1}{2}\|\Phi x - u\|_2^2$   
with  $L = \lambda_{\max}(\Phi^* \Phi)$

separable majorizer ( $\beta_k < 1/L$ )

$$Q(x, x_k) = f_1(x_k) + (x - x_k)^T \nabla f_1(x_k) + \frac{1}{2\beta_k} \|x - x_k\|_2^2$$

majorization-minimization algorithm

$$\begin{aligned} x_{k+1} &= \arg \min_x Q(x, x_k) + f_2(x) \\ &= \arg \min_x \frac{1}{2\beta_k} \left\| x - (x_k - \beta_k \nabla f_1(x_k)) \right\|_2^2 + f_2(x) \end{aligned}$$

$$x_{k+1} = \text{prox}_{\beta_k f_2} \left( x_k - \beta_k \nabla f_1(x_k) \right)$$

# Proximal algorithms: convergence

Problem:  $\hat{x} \in \arg \min_{x \in \mathbb{R}^n} f_1(x) + f_2(x)$

$$f(x)$$

$f_1$  has a  $L$ -Lipschitz gradient; e.g.  $f_1(x) = \frac{1}{2}\|\Phi x - u\|_2^2$

Iterative shrinkage/thresholding (IST)  
(or forward-backward)

$$L = \lambda_{\max}(\Phi^* \Phi)$$

$$x_{k+1} = \text{prox}_{\beta_k f_2} \left( x_k - \beta_k \nabla f_1(x_k) \right)$$

if  $\beta_k < \frac{1}{L}$ , IST is a majorization-minimization algorithm, thus

$$f(x_{k+1}) \leq f(x_k)$$

$f(x) \geq 0$ , thus  $(f(x_1), f(x_2), \dots, f(x_k), \dots)$  converges.

Attention: this does **not** imply convergence of  $(x_1, \dots, x_k, \dots)$

# Proximal algorithms: convergence

$$\hat{x} \in G = \arg \min_{x \in \mathbb{R}^n} f_1(x) + f_2(x)$$

IST algorithm:  $x_{k+1} = \text{prox}_{\beta_k f_2} \left( x_k - \beta_k \nabla f_1(x_k) \right)$

if  $0 < \beta_k < \frac{2}{L}$ , then  $(x_1, x_2, \dots, x_k, \dots)$   
converges to a point in  $G$

Inexact version:

$$x_{k+1} = \text{prox}_{\beta_k f_2} \left( x_k - (\beta_k \nabla f_1(x_k) + b_k) \right) + a_k$$

convergence still guaranteed if

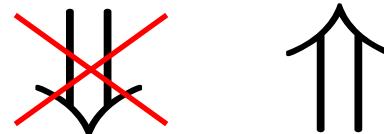
$$\sum_{k=1}^{\infty} \|a_k\| < \infty$$

$$\sum_{k=1}^{\infty} \|b_k\| < \infty$$

Results and proofs in [Combettes and Wajs, 2005]

# Proximal algorithms: convergence

Convergence of function values  $(f(x_1), \dots, f(x_k), \dots) \rightarrow f(\hat{x})$



Convergence of iterates  $(x_1, x_2, \dots, x_k, \dots) \rightarrow \hat{x}$

Convergence rates (for function values) [Beck, Teboulle, 2009]:

$$f(x_k) - f(\hat{x}) \leq \frac{L\|x_0 - \hat{x}\|_2^2}{2k}$$

Convergence rate for the iterates require further assumptions on  $f$

# Proximal algorithms: convergence of iterates

$$\widehat{x} = \arg \min_x \frac{1}{2} \|\Phi x - u\|_2^2 + f_2(x)$$

With  $L = \lambda_{\max}(\Phi^* \Phi)$   $l = \lambda_{\min}(\Phi^* \Phi) > 0 \Rightarrow G = \{\widehat{x}\}$   
 $\kappa = l/L$  (condition number) (unique minimizer)

Under- ( $\gamma < 1$ ) or **over-relaxed** ( $\gamma > 1$ ) IST

$$x_{k+1} = (1 - \gamma)x_k + \gamma \text{prox}_{f_2} \left( x_k - \beta \Phi^T (\Phi x_k - u) \right)$$

Optimal choice  $\gamma = \frac{2}{L + l}$   $\rho = \frac{1 - \kappa}{1 + \kappa}$

Q-linear convergence  $\|x_{k+1} - \widehat{x}\| \leq \rho \|x_k - \widehat{x}\|$

Small  $l \Rightarrow \rho \lesssim 1 \Rightarrow$  slow convergence!

[F, Bioucas-Dias, 2007]

# Proximal algorithms: convergence of iterates

$$\hat{x} \in G = \arg \min_x \frac{1}{2} \|\Phi x - u\|_2^2 + \tau \|x\|_1$$

With  $L = \lambda_{\max}(\Phi^* \Phi)$ ; using a step-size  $\beta < 2/L$ ,

$$x_{k+1} = \text{soft}\left(x_k - \beta \Phi^T (\Phi x_k - u), \beta \tau\right)$$

$\mathcal{Z} \subseteq \{1, 2, \dots, n\}$  such that  $\hat{x} \in G \Rightarrow [\hat{x}]_{\mathcal{Z}} = 0$

Then, after a finite number of iterations:  $[x_k]_{\mathcal{Z}} = [\hat{x}]_{\mathcal{Z}} = 0$

After this, Q-linear convergence:  $l = \lambda_{\min}(\Phi_{\bar{\mathcal{Z}}}^* \Phi_{\bar{\mathcal{Z}}}) > 0$

Optimal choice  $\beta = \frac{2}{L+l}$ ,

$$\rho = \frac{1-\kappa}{1+\kappa}$$

$$\|x_{k+1} - \hat{x}\| \leq \rho \|x_k - \hat{x}\|$$

[Hale, Yin, Zhang, 2008]

# Slowness and acceleration of IST

Problem:  $\hat{x} \in G = \arg \min_x \frac{1}{2} \|\Phi x - u\|_2^2 + \tau \|x\|_1$

IST algorithm:  $x_{k+1} = \text{soft}\left(x_k - \beta \Phi^T(\Phi x_k - u), \beta \tau\right)$

IST is **slow**, if  $\Phi$  is very ill-conditioned and/or  $\tau$  is very small!

Several proposals for accelerated variants of IST

Methods with memory (TwIST, FISTA)

Quasi-Newton methods (SpaRSA)

Continuation, i.e., use a varying  $\tau$  (FPC, SpaRSA)

# Memory-based variants of IST: FISTA

Fast IST algortihm (FISTA); based on Nesterov's work (1980's)  
[Beck, Teboulle, 2009]

FISTA

$$t_{k+1} = \frac{1 + \sqrt{1 + 4 t_k^2}}{2}$$
$$z_{k+1} = x_k + \frac{t_k - 1}{t_{k+1}} (x_k - x_{k-1})$$
$$x_{k+1} = \text{soft}\left(z_k - \beta \Phi^T(\Phi z_k - u), \beta \tau\right)$$

IST:

$$f(x_k) - f(\hat{x}) = O\left(\frac{1}{k}\right) \quad \left( \leq \frac{L \|x_0 - \hat{x}\|_2^2}{2 k} \right)$$

FISTA:

$$f(x_k) - f(\hat{x}) = O\left(\frac{1}{k^2}\right)$$

# Memory-based variants of IST: twist

Inspired by 2-step methods for linear systems

[Frankel, 1950], [Axelsson, 1996]

TwIST (two-step IST):

[Bioucas-Dias, F, 2007]

$$x_{k+1} = (\alpha - \beta)x_k + (1 - \alpha)x_{k-1} + \beta \text{prox}_{f_2}(x_k - \Phi^T(\Phi x_k - u))$$

$$\kappa = \frac{\lambda_{\min}(\Phi^T \Phi)}{\lambda_{\max}(\Phi^T \Phi)}$$

$$\text{Q-linear convergence } \|x_{k+1} - \hat{x}\| \leq \rho \|x_k - \hat{x}\|$$

$$\rho = \frac{1 - \sqrt{\kappa}}{1 + \sqrt{\kappa}} \quad \text{TwIST}$$

$$\rho = \frac{1 - \kappa}{1 + \kappa} \quad \text{IST}$$

# Memory-based variants of IST: twist

original



Blurred ( $B$ ), 9x9, 40db noise

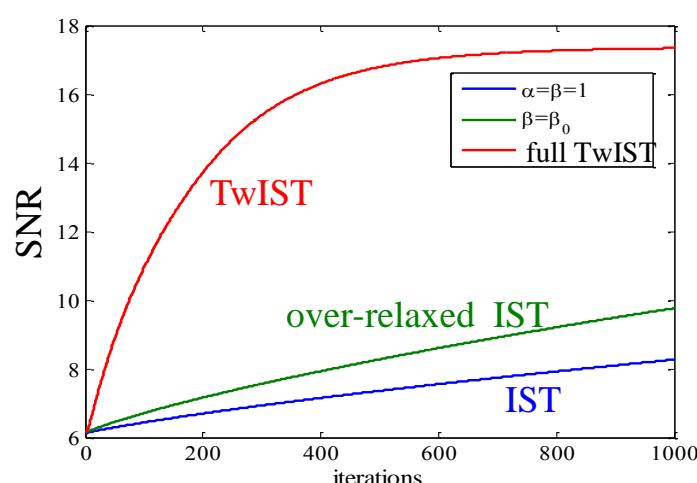
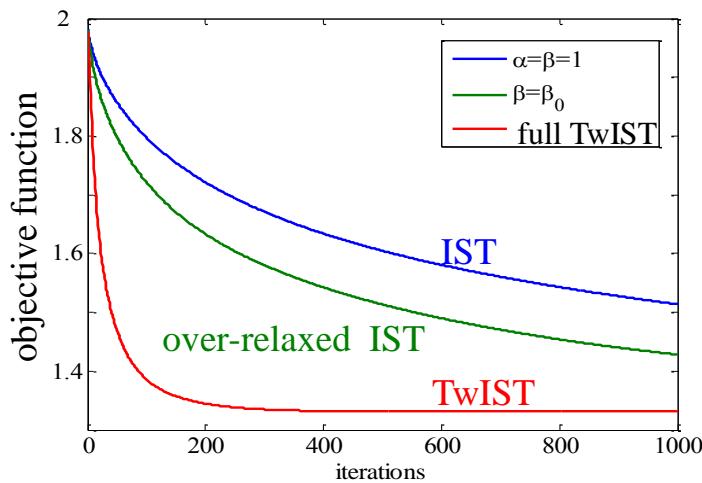


restored



$$\hat{x} \in \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|B\Psi x - u\|_2^2 + \tau \|x\|_1$$

representation coefficients  
dictionary (e.g, wavelet basis, frame, ...)



# Quasi-newton acceleration of IST: SpaRSA

$$\text{IST: } x_{k+1} = \text{prox}_{\beta_k f_2} \left( x_k - \beta_k \nabla f_1(x_k) \right)$$

A Newton step (instead of gradient descent) would be:

$$x_{k+1} = \text{prox}_{\beta_k f_2} \left( x_k - [H(x_k)]^{-1} \nabla f_1(x_k) \right)$$

...computationally **too expensive!**

  
Hessian  
(matrix of second derivatives)

Barzilai-Borwein approach:

[Barzilai-Borwein, 1988], [Wright, Nowak, F, 2009]

$$\boxed{\frac{1}{\beta_k} I \simeq H(x_k)}$$

$$\frac{1}{\beta_k} = \arg \min_{\alpha} \|\alpha(x_k - x_{k-1}) - (\nabla f(x_k) - \nabla f(x_{k-1}))\|_2^2$$

$$\text{If } f_1(x) = \frac{1}{2} \|\Phi x - u\|_2^2, \text{ then } \beta_k = \frac{\|x_k - x_{k-1}\|_2^2}{\|\Phi(x_k - x_{k-1})\|_2^2}$$

# Acceleration via continuation

$$\text{IST: } x_{k+1} = \text{soft}\left(x_k - \beta \Phi^T(\Phi x_k - u), \beta\tau\right)$$

**Slow**, if  $\tau$  is small.

Observation: IST (as SpaRSA) benefits from  
“warm-starting” (being initialized *close* to the minimizer)

Continuation: start with large  $\tau$   
slowly decrease  $\tau$  while tracking the solution.  
[F, Nowak, Wright, 2007], [Hale, Yin, Zhang, 2007]

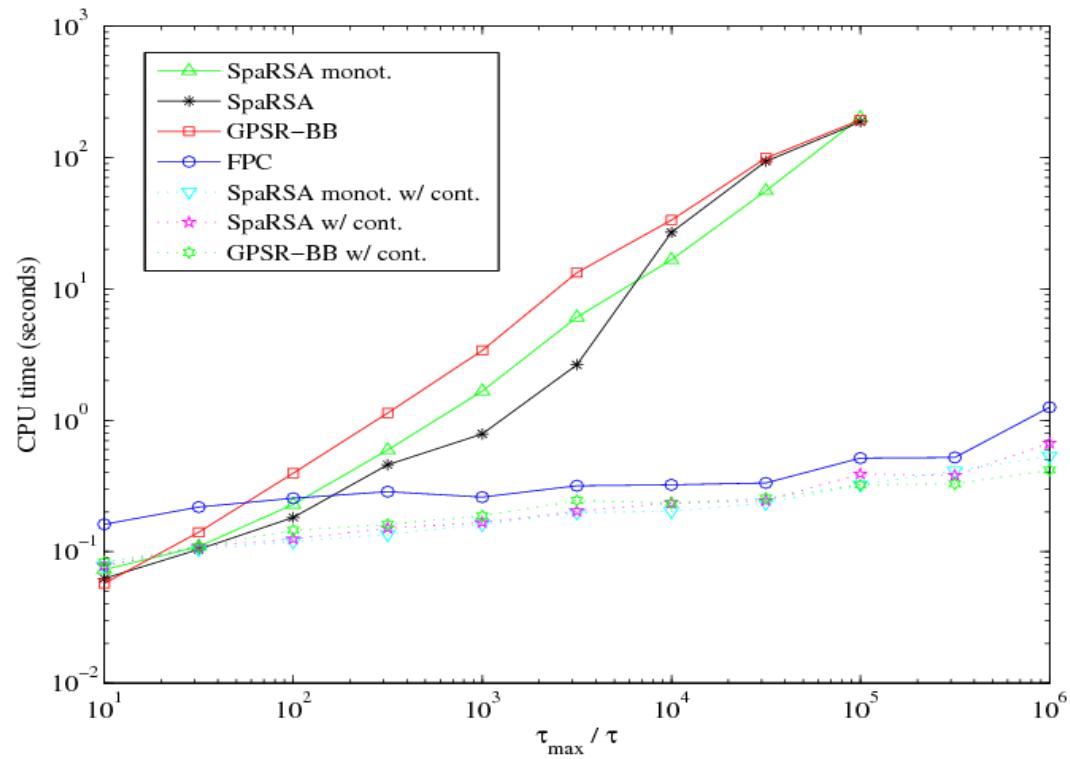
IST + continuation = fixed point continuation (FPC)  
[Hale, Yin, Zhang, 2007]

# Acceleration via continuation

$$\hat{x} \in G = \arg \min_x \frac{1}{2} \|\Phi x - u\|_2^2 + \tau \|x\|_1$$

$1024 \times 4096$

$$u = \Phi x^* + n$$



$$\tau_{\max} = \|\Phi^T \mathbf{y}\|_\infty \quad (\tau \geq \tau_{\max} \Rightarrow \hat{x} = 0)$$

# Some speed comparisons

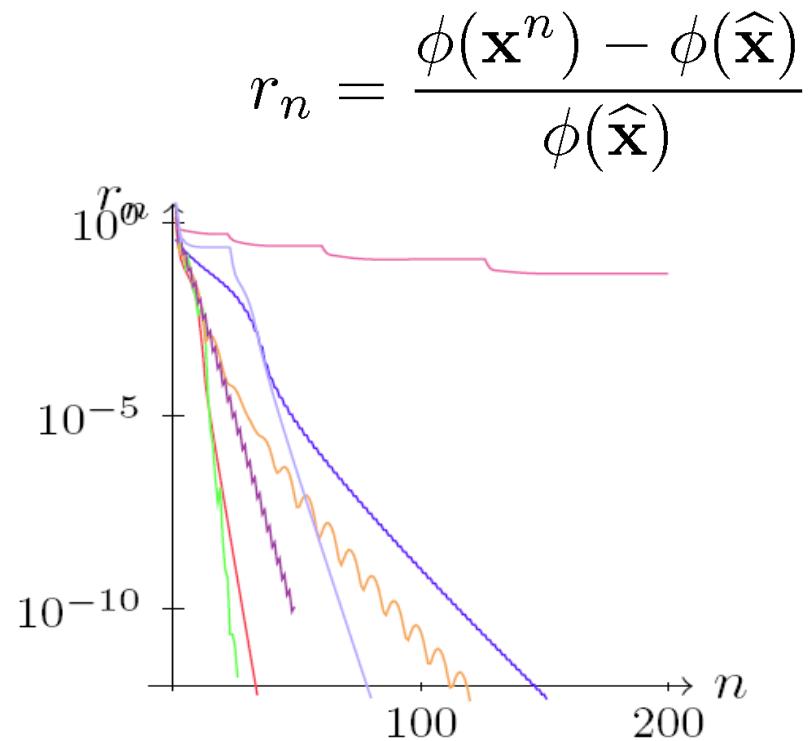
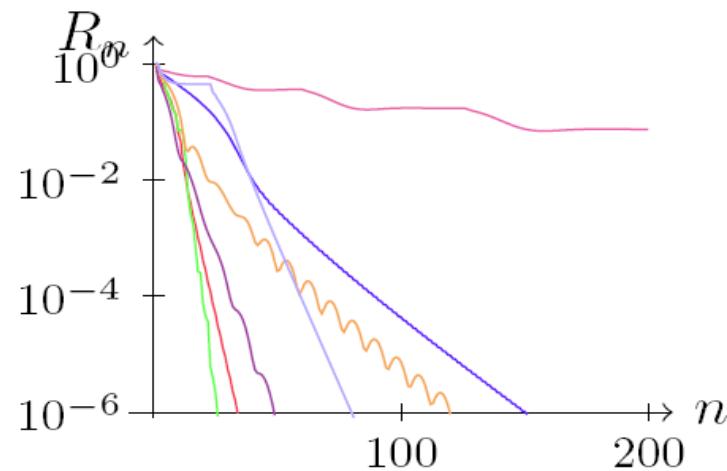
from [Lorenz, 2011]

$$\hat{x} = \arg \min_x \frac{1}{2} \|\Phi x - u\|_2^2 + \tau \|x\|_1$$

$$\Phi = [I \ U \ R] \quad \tau = 0.1$$

$$(512 \times 1536) \quad \hat{x} \text{ with 120 non-zeros}$$

$$R_n = \frac{\|\mathbf{x}^n - \hat{\mathbf{x}}\|_2}{\|\hat{\mathbf{x}}\|_2}$$



IST, GPSR, SpaRSA, FISTA, YALL1, NESTA, fpc

# Proximal algorithms for matrices

$$\widehat{M} \in \arg \min_{M \in \mathbb{R}^{n \times n}} \frac{1}{2} \|\Phi(M) - U\|_F^2 + \mu \|M\|_*$$

↑ linear operator  
↗ ...its adjoint

The proximal algorithm (ISTA) is as before:

$$X_{k+1} = \text{svt}_{\mu \beta_k} \left( X_k - \beta_k \Phi^*(\Phi(X_k) - U) \right)$$

Matrix completion:  $\Phi(X) = X$  (subset of entries)  $|\Omega| = p$

Unknown M				IST			APG (FISTA)		
$n/r$	$p$	$p/d_r$	$\mu$	iter	#sv	error	iter	#sv	error
100/10	5666	3	8.21e-03	7723	61	1.88e-01	655	13	1.06e-03
200/10	15665	4	1.05e-02	12180	96	2.45e-01	812	12	1.02e-03
500/10	49471	5	1.21e-02	10900	203	5.91e-01	1132	16	7.63e-04

Unknown M				FPC (continuation)			APG + continuation		
$n/r$	$p$	$p/d_r$	$\mu$	iter	#sv	error	iter	#sv	error
100/10	5666	3	8.21e-03	429	32	1.06e-03	74	10	1.46e-04
200/10	15665	4	1.05e-02	278	49	4.38e-04	73	10	1.02e-04
500/10	49471	5	1.21e-02	484	125	5.50e-04	72	10	8.06e-05

from [Toh, Yun, 2009]

...the importance of acceleration!

# Another class of methods: augmented Lagrangian

The problem:  $\min_x f(x)$

s.t.  $\Phi x = u$

The augmented Lagrangian (AL)

$$L_\mu(x, \lambda) = f(x) + \lambda^T(\Phi x - u) + \frac{\mu}{2} \|\Phi x - u\|_2^2$$

Penalty parameter  
↓

The “AL method” (ALM)  
(a.k.a. method of multipliers)  
[Hestenes, Powell, 1969]

$$\begin{aligned} x_{k+1} &= \arg \min_x L_\mu(x, \lambda_k) \\ \lambda_{k+1} &= \lambda_k + \mu(\Phi x_{k+1} - u) \end{aligned}$$

Can be  
written as:

$$\begin{aligned} x_{k+1} &= \arg \min_x f(x) + \frac{\mu}{2} \|\Phi x - u - d_k\|_2^2 \\ d_{k+1} &= d_k - (\Phi x_{k+1} - u) \end{aligned}$$

Similar to Bregman method [Osher, Burger, Goldfarb, Xu, Yin, 2005]  
[Yin, Osher, Goldfarb, Darbon, 2008]

# Augmented Lagrangian for variable splitting

The problem:  $\min_x f_1(\Phi x) + f_2(x)$

Equivalent constrained formulation  $\begin{aligned} \min_x & f_1(z) + f_2(x) \\ \text{s.t. } & \Phi x - z = 0 \end{aligned}$

Can be written as

$$\begin{aligned} \min_y & f(y) \\ \text{s.t. } & \Psi y = 0 \end{aligned} \quad \text{with} \quad \begin{aligned} y &= \begin{bmatrix} x \\ z \end{bmatrix} \\ \Psi &= [\Phi \quad -I] \end{aligned}$$

ALM:

$$(x_{k+1}, z_{k+1}) = \arg \min_{x, z} f_1(z) + f_2(x) + \frac{\mu}{2} \|\Phi x - z - d_k\|_2^2$$

$$d_{k+1} = d_k - (\Phi x_{k+1} - z_{k+1})$$

# Augmented Lagrangian for variable splitting

It may be hard to solve

$$(x_{k+1}, z_{k+1}) = \arg \min_{x, z} f_1(z) + f_2(x) + \frac{\mu}{2} \|\Phi x - z - d_k\|_2^2$$

Alternative:

$$\begin{aligned} x_{k+1} &= \arg \min_x f_2(x) + \frac{\mu}{2} \|\Phi x - z_k - d_k\|_2^2 \\ z_{k+1} &= \arg \min_z f_1(z) + \frac{\mu}{2} \|\Phi x_{k+1} - z - d_k\|_2^2 \\ d_{k+1} &= d_k - (\Phi x_{k+1} - z_{k+1}) \end{aligned}$$

Alternating directions method of multipliers (ADMM)

[Glowinsky, Marrocco, 1975], [Gabay, Mercier, 1976], [Eckstein, Bertsekas, 1992]

When applied to  $\hat{x} = \arg \min_x \frac{1}{2} \|\Phi x - u\|_2^2 + \tau \|x\|_1$

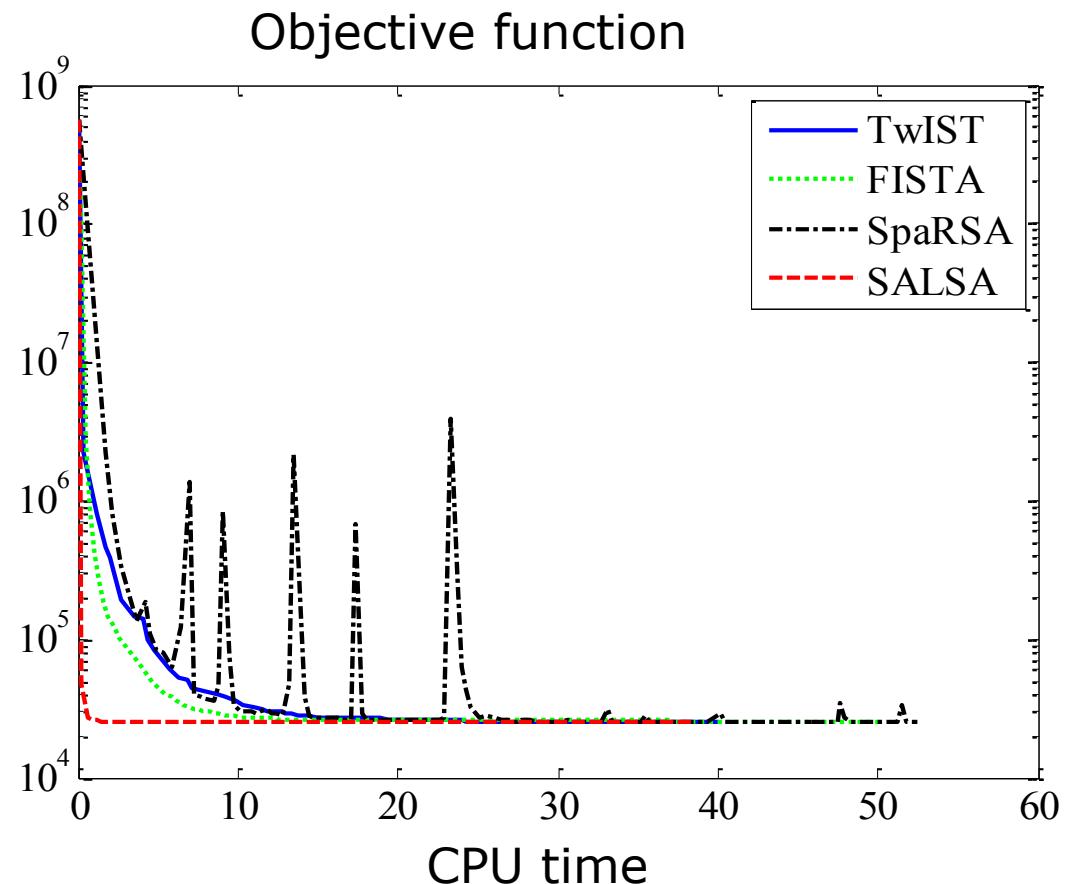
split augmented Lagrangian shrinkage algorithm (SALSA)  
[F. Bioucas-Dias, Afonso, 2009]

# Augmented Lagrangian for variable splitting

Testing ADMM/SALSA on a typical image deblurring problem



$$\hat{x} \in \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|B\Psi x - u\|_2^2 + \tau \|x\|_1$$



# Handling more than two functions

$$\hat{x} \in \arg \min_{x \in \mathbb{R}^n} f_0(x) + f_1(x) + \cdots + f_n(x)$$

$f_0$  has a  $L$ -Lipschitz gradient       $f_1, \dots, f_n$  are convex

Possible uses: multiple regularizers, positivity constraints, ...

Generalized forward-backward algorithm [Raguet, Fadili, Peyré, 2011]

Parameters:  $\omega_1, \dots, \omega_n \in (0, 1)$ , s.t.  $\sum_j \omega_j = 1$

Initialization:  $k = 0$ ,  $z_0^1, \dots, z_0^n$ ,  $x_0 = \sum_{j=1}^n \omega_j z_0^j$

repeat until convergence

```
for i = 1 : n
    zk+1i = zki + proxβkfi/ωi(2xk - zki - βk ∇fi(xk)) - xk
    xk+1 = ∑i=1n ωi zk+1i
    k ← k + 1
```

# Handling more than two functions

$$\hat{x} \in \arg \min_{x \in \mathbb{R}^n} f_1(x) + \cdots + f_n(x)$$

$f_1, \dots, f_n$  arbitrary convex functions

ADMM-based method [F and Bioucas-Dias, 2009], [Setzer, Steidl, Teuber, 2009]

Parameter:  $\gamma$

Initialization:  $k = 0, z_0^1, \dots, z_0^n, y_0^1, \dots, y_0^n$

repeat until convergence

$$x_{k+1} = (1/n) \sum_{i=1}^n (y_k^i - z_k^i)$$

for  $i = 1 : n$

$$y_{k+1}^i = \text{prox}_{\gamma f_i}(x_k - z_k^i)$$

$$z_{k+1}^i = z_k^i + x_k - y_{k+1}^i$$

$$k \leftarrow k + 1$$

# **Non-Convex Algorithms for Low-Dimensional Models**



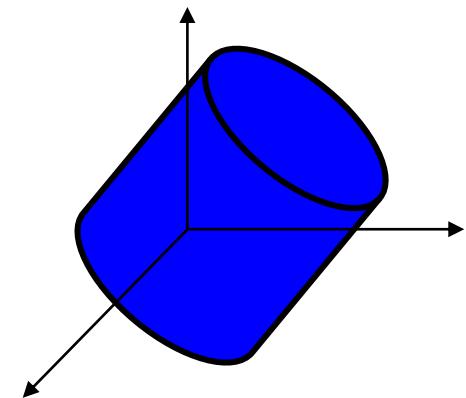
# Motivation

Discrete descriptions of low-dimensional models

$$x = \sum_{i=1}^{|A|} a_i c_i \quad a_i \in A, \|c_i\|_0 \leq K$$

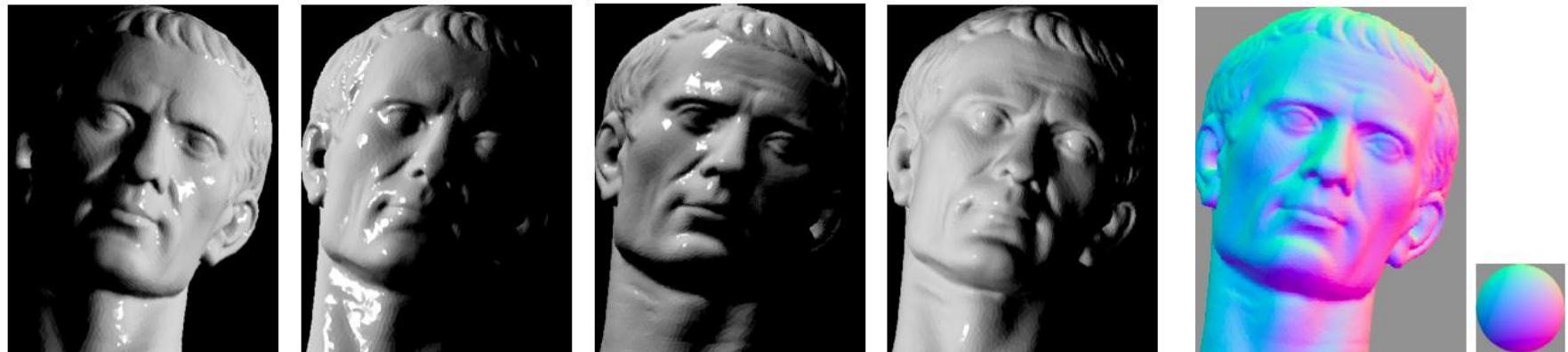
$a_i$ : atoms  
 $A$ : atomic set

$$A = \{A : \text{rank}(A) = 1, \|A\|_F = 1\}$$



Example: reflectivity of Lambertian surfaces

[Basri and Jacobs 2001]



$$K \leq 9$$

$$\text{Intensity} = \rho \max\{\langle n, l \rangle, 0\}$$

# Motivation

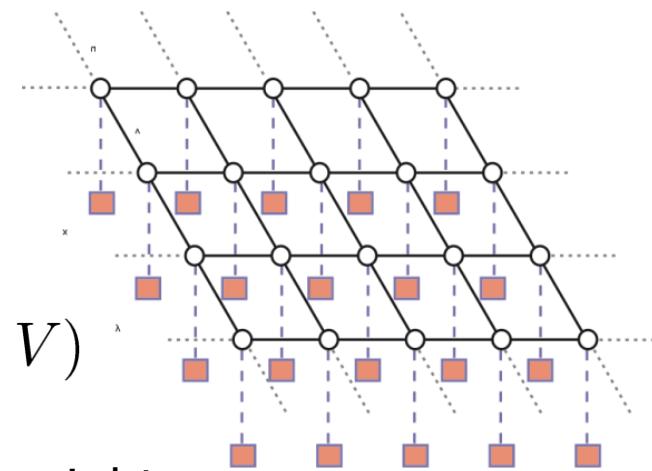
Discrete descriptions of low-dimensional models

$$x = \sum_{i=1}^{|A|} a_i c_i \quad a_i \in A, \|c_i\|_0 \leq K$$

$a_i$ : atoms  
 $A$ : atomic set

$$A = \{\pm e_i\}_{i=1}^N$$

$$\mathcal{G} = (E, V)$$



Example: graphical model selection

$$u = E_i \quad \Phi = E_{\setminus i} \quad \alpha = V_i$$

vertex weights  
of the  $i$ -th edge

$$= \begin{matrix} \text{color bar} \\ \text{matrix} \end{matrix}$$

$M \times 1$

$M \times N \ (M < N)$

$$= \begin{matrix} \text{color bar} \\ \text{matrix} \end{matrix}$$

$N \times 1$

Gauss-Markov graph  
<>  
linear regression

$K \leq$  node degree of  $\mathcal{G}$

[Meinshausen and Bühlman 2006]

# Motivation

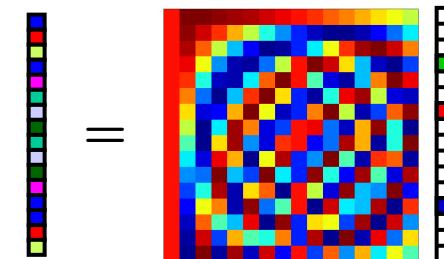
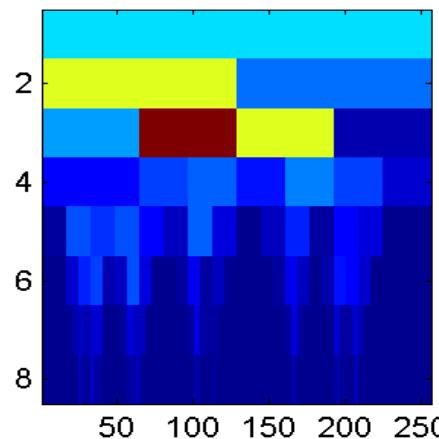
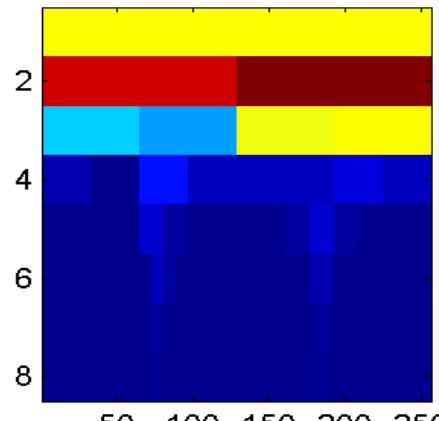
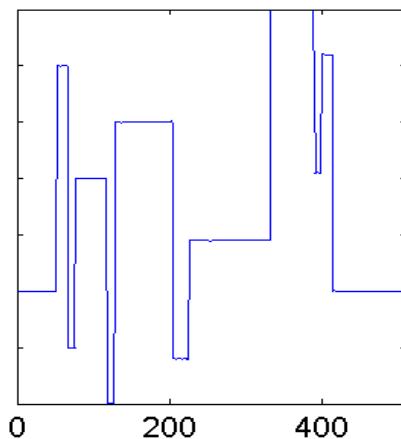
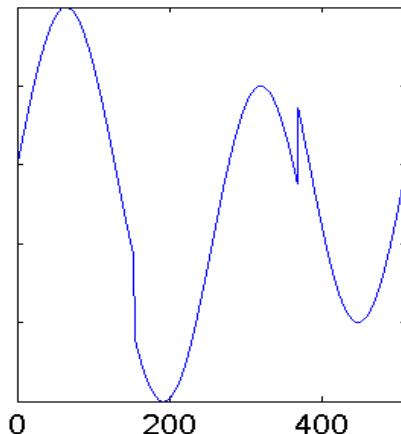
Discrete descriptions of **structure** in low-dimensional models

$$x = \sum_{i=1}^{|A|} a_i c_i$$

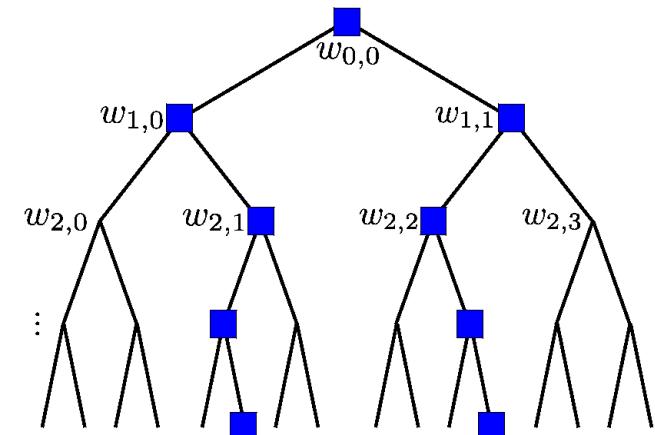
$$a_i \in A, \|c_i\|_0 \leq K$$

$a_i$ : atoms

$A$ : atomic set



$$x = \Psi \times \alpha$$



Typical of wavelet transforms of natural signals and images (piecewise smooth)

# Motivation

Non-convex criteria beyond atomic norms

- 1-bit compressive sensing

$$u = \text{sign}(\Phi x)$$

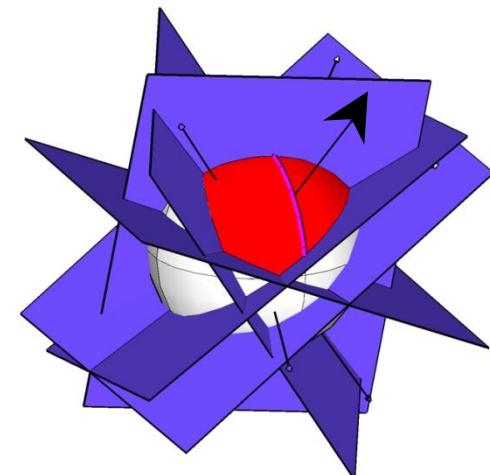
- optimization criteria  $\arg \min_{x: \|x\|_0 \leq K} f(x)$

$$f(x) = -\langle u, \text{sign}(\Phi x) \rangle$$

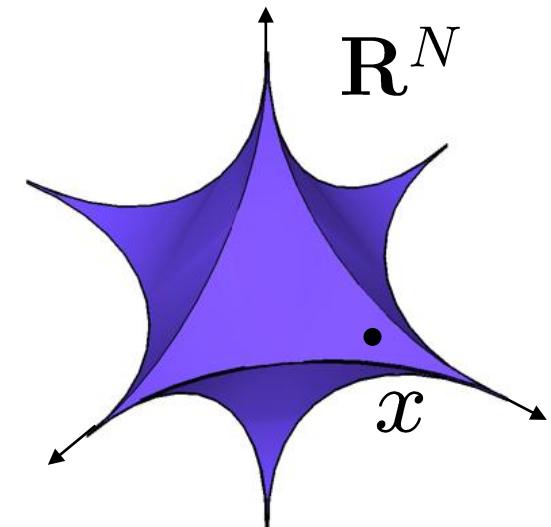
- Compressible signals in weak  $\ell_q, q < 1$

$$|x|_{(i)} \leq R i^{-1/q}$$

- optimization criteria  $\arg \min_{x: u = \Phi x} \|x\|_q$



[Boufounos and Baraniuk 2008]



[Chartrand and Yin 2008]

# Non-convexity in this tutorial

- Anything **not** convex      <>      too big to cover

***convexity is in general a rare condition***

- Active research topic with great depth

[Attouch et al. 2010]

***Key lesson:***

**convergence of the projected gradient-descent algorithm**

- This tutorial                          <>                          *a special subset*

$$\hat{x} \in \arg \min_{x \in \mathbb{R}^N} f_1(x) + f_2(x) \quad (\mathcal{S} \text{ is non-convex})$$

with       $f_2(x) = \begin{cases} g(x) & \Leftarrow x \in \mathcal{S} \\ +\infty & \Leftarrow x \notin \mathcal{S} \end{cases}$

***Assumptions:***

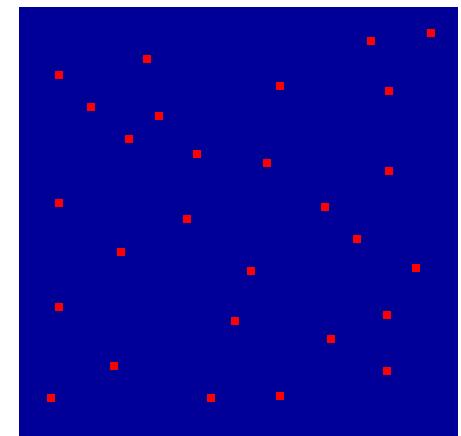
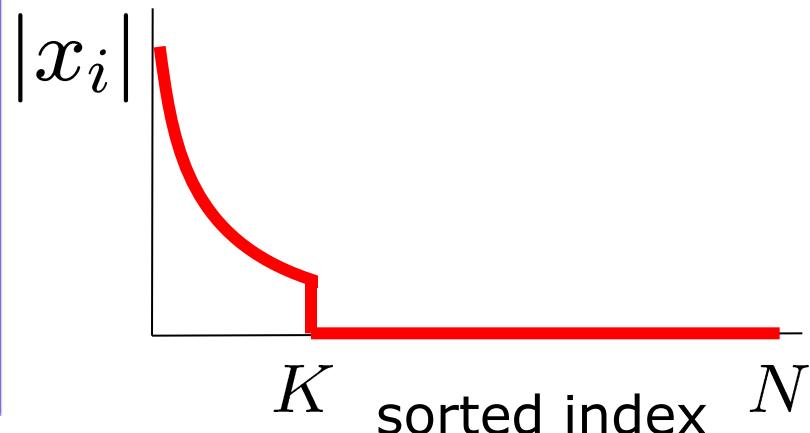
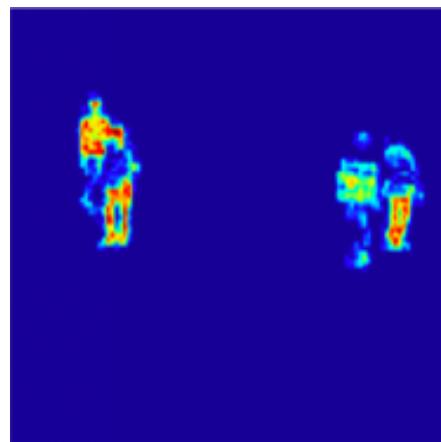
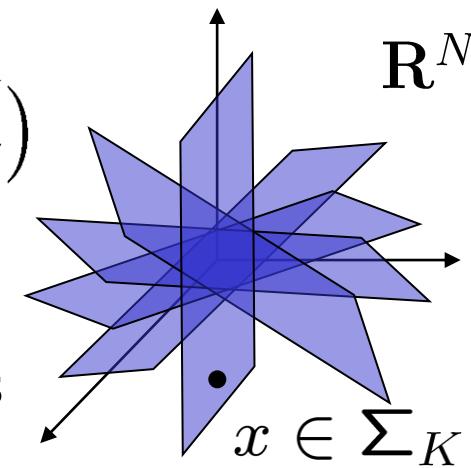
1. ***access to the gradient of convex***  $f_1$
2. ***tractable/approximate prox of non-convex***  $f_2$



# Can we project onto non-convex sets?

## Running examples

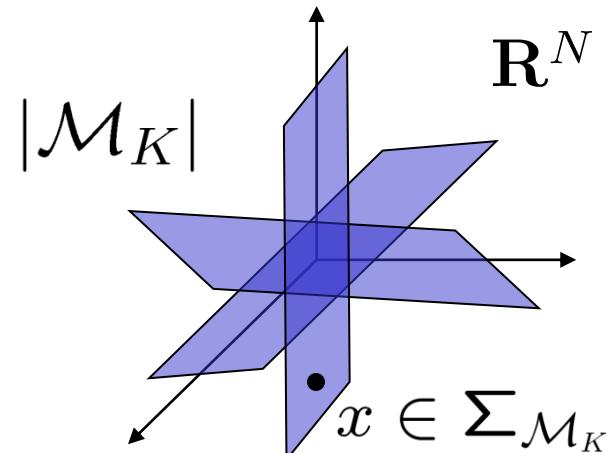
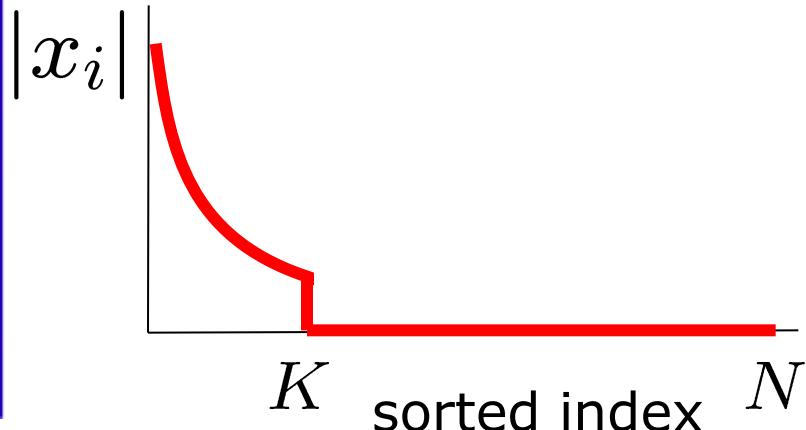
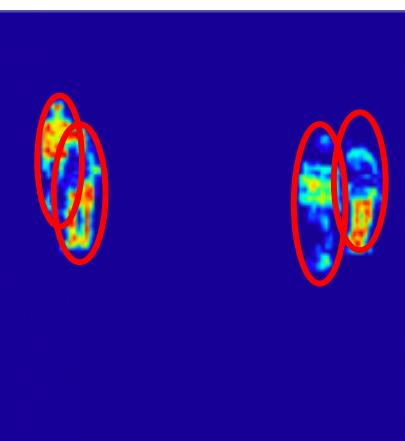
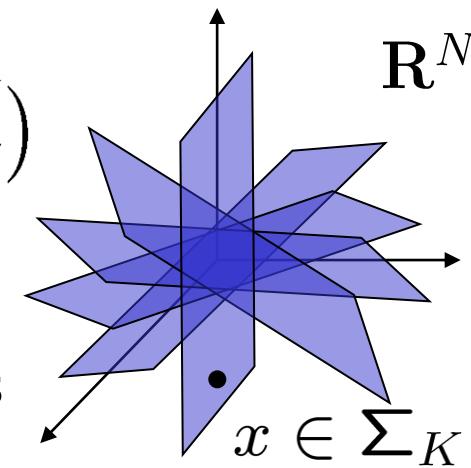
- Sparse signal: only  $K$  out of  $N$  coordinates nonzero
  - model: union of all  $K$ -dimensional subspaces aligned w/ coordinate axes
- **Structured** sparse signal: reduced set of subspaces (or model-sparse)
  - model: a particular union of subspaces  
ex: clustered or dispersed sparse patterns



# Can we project onto non-convex sets?

## Running examples

- Sparse signal: only  $K$  out of  $N$  coordinates nonzero
  - model: union of all  $K$ -dimensional subspaces aligned w/ coordinate axes
- **Structured** sparse signal: reduced set of subspaces (or model-sparse)
  - model: a particular union of subspaces  
ex: clustered or dispersed sparse patterns



# Can we project onto non-convex sets?

$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|y - x\|_2^2 + f_2(x) \equiv \text{prox}_{f_2}(y)$$

- Analysis of the prox for *structured* sparse sets  $g(x) = 0$

$$\text{prox}_{f_2}(y) = \arg \min_{x: x \in \Sigma_{\mathcal{M}_K}} \|x - y\|$$

***support of the solution*** <> ***modular approximation problem***

$$\text{supp}(\arg \min_{x: \text{supp}(x) \in \mathcal{M}_K} \|x - y\|_2^2) = \arg \min_{\mathcal{S}: \mathcal{S} \in \bar{\mathcal{M}}_K} \|(y)_{\mathcal{S}} - y\|_2^2$$

indexing set

# Can we project onto non-convex sets?

$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|y - x\|_2^2 + f_2(x) \equiv \text{prox}_{f_2}(y)$$

- Analysis of the prox for *structured* sparse sets  $g(x) = 0$

$$\text{prox}_{f_2}(y) = \arg \min_{x: x \in \Sigma_{\mathcal{M}_K}} \|x - y\|$$

support of the solution <> modular approximation problem

$$\text{supp}(\arg \min_{x: \text{supp}(x) \in \mathcal{M}_K} \|x - y\|_2^2) = \arg \min_{\mathcal{S}: \mathcal{S} \in \bar{\mathcal{M}}_K} \|(y)_{\mathcal{S}} - y\|_2^2$$

$$= \arg \max_{\mathcal{S}: \mathcal{S} \in \bar{\mathcal{M}}_K} \{ \|y\|^2 - \|(y)_{\mathcal{S}} - y\|_2^2 \}$$

# Can we project onto non-convex sets?

$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|y - x\|_2^2 + f_2(x) \equiv \text{prox}_{f_2}(y)$$

- Analysis of the prox for *structured* sparse sets  $g(x) = 0$

$$\text{prox}_{f_2}(y) = \arg \min_{x: x \in \Sigma_{\mathcal{M}_K}} \|x - y\|$$

support of the solution <> modular approximation problem

$$\text{supp}(\arg \min_{x: \text{supp}(x) \in \mathcal{M}_K} \|x - y\|_2^2) = \arg \min_{\mathcal{S}: \mathcal{S} \in \bar{\mathcal{M}}_K} \|(y)_{\mathcal{S}} - y\|_2^2$$

$$= \arg \max_{\mathcal{S}: \mathcal{S} \in \bar{\mathcal{M}}_K} \{ \|y\|^2 - \|(y)_{\mathcal{S}} - y\|_2^2 \}$$

$$= \arg \max_{\mathcal{S}: \mathcal{S} \in \bar{\mathcal{M}}_K} \|(y)_{\mathcal{S}}\|^2$$

# Can we project onto non-convex sets?

$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|y - x\|_2^2 + f_2(x) \equiv \text{prox}_{f_2}(y)$$

- Analysis of the prox for *structured* sparse sets  $g(x) = 0$

$$\text{prox}_{f_2}(y) = \arg \min_{x: x \in \Sigma_{\mathcal{M}_K}} \|x - y\|$$

support of the solution <> modular approximation problem

$$\text{supp}(\arg \min_{x: \text{supp}(x) \in \mathcal{M}_K} \|x - y\|_2^2) = \arg \max_{S: S \in \bar{\mathcal{M}}_K} F(S; y)$$

where  $F(S; y) = \sum_{i \in S} |y_i|^2$ .

# Can we project onto non-convex sets?

$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|y - x\|_2^2 + f_2(x) \equiv \text{prox}_{f_2}(y)$$

- Analysis of the prox for *structured* sparse sets  $g(x) = 0$

$$\text{prox}_{f_2}(y) = \arg \min_{x: x \in \Sigma_{\mathcal{M}_K}} \|x - y\|$$

support of the solution  $\leftrightarrow$  modular approximation problem

$$\text{supp} \left( \arg \min_{x: \text{supp}(x) \in \mathcal{M}_K} \|x - y\|_2^2 \right) = \arg \max_{S: S \in \bar{\mathcal{M}}_K} F(S; y)$$

***underlying optimization problem*  $\leftrightarrow$  *integer linear program***

$$\text{supp} \left( \arg \min_z \left\{ \rho^T z : z \in \Sigma_{\mathcal{M}_K} \right\} \right)$$

$$z_i \in \{0, 1\}: \text{support indicator variables} \quad \rho_i = -|y_i|^2$$



# Can we project onto non-convex sets?

$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|y - x\|_2^2 + f_2(x) \equiv \text{prox}_{f_2}(y)$$

- Analysis of the prox for *structured* sparse sets  $g(x) = 0$

$$\text{prox}_{f_2}(y) = \arg \min_{x: x \in \Sigma_{\mathcal{M}_K}} \|x - y\|$$

support of the solution  $\leftrightarrow$  modular approximation problem

$$\text{supp}(\arg \min_{x: \text{supp}(x) \in \mathcal{M}_K} \|x - y\|_2^2) = \arg \max_{S: S \in \bar{\mathcal{M}}_K} F(S; y)$$

underlying optimization problem  $\leftrightarrow$  integer linear program

**Class of problems we can tractably solve:**

**PMAP**

- **Polynomial time modular epsilon-approximation property**

$$F(\hat{\mathcal{S}}_\epsilon; y) \geq (1 - \epsilon) \max_{S: S \in \bar{\mathcal{M}}_K} F(S; y)$$

[Kyrillidis and C, 2011]

# Can we project onto non-convex sets?

PMAP-0:

- Matroid structured sparse models:

$$\mathcal{M} = (\mathcal{N}, \mathcal{I} \subseteq 2^{\mathcal{N}}), \mathcal{N} = \{1, \dots, N\}$$

$\mathcal{N}$ : ground set

$\mathcal{I}$ : base set

Definition:

**non-emptiness**

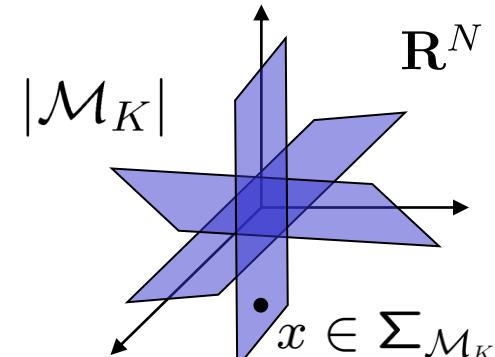
$$1. \emptyset \in \mathcal{I}$$

**heredity**

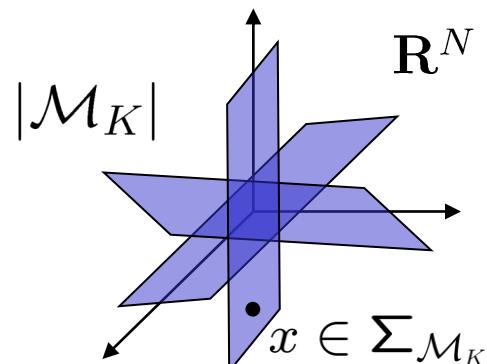
$$2. A \in \mathcal{I} \text{ and } B \subseteq A \Rightarrow B \in \mathcal{I}$$

**exchange**

$$3. A, B \in \mathcal{I} \text{ and } |A| > |B| \Rightarrow \exists e \in A \setminus B \text{ such that } B \cup \{e\} \in \mathcal{I}$$



# Can we project onto non-convex sets?



PMAP-0:

- Matroid structured sparse models:

$$\mathcal{M} = (\mathcal{N}, \mathcal{I} \subseteq 2^{\mathcal{N}}), \mathcal{N} = \{1, \dots, N\}$$

$\mathcal{N}$ : ground set

$\mathcal{I}$ : base set

Definition:

**non-emptiness**

$$1. \emptyset \in \mathcal{I}$$

**heredity**

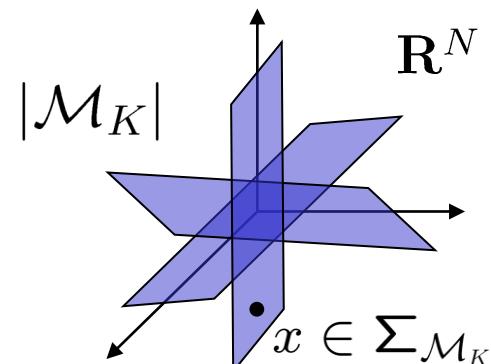
$$2. A \in \mathcal{I} \text{ and } B \subseteq A \Rightarrow B \in \mathcal{I}$$

**exchange**

$$3. A, B \in \mathcal{I} \text{ and } |A| > |B| \Rightarrow \exists e \in A \setminus B \text{ such that } B \cup \{e\} \in \mathcal{I}$$

Let  $\mathcal{N} = \{1, 2, 3, 4\}$ . The smallest matroid that contains  $\{1, 2\}$  and  $\{3, 4\}$  is ???

# Can we project onto non-convex sets?



PMAP-0:

- Matroid structured sparse models:

$$\mathcal{M} = (\mathcal{N}, \mathcal{I} \subseteq 2^{\mathcal{N}}), \mathcal{N} = \{1, \dots, N\}$$

$\mathcal{N}$ : ground set

$\mathcal{I}$ : base set

Definition:

- |                      |   |
|----------------------|---|
| <b>non-emptiness</b> | 1. $\emptyset \in \mathcal{I}$  |
| <b>heredity</b>      | 2. $A \in \mathcal{I}$ and $B \subseteq A \Rightarrow B \in \mathcal{I}$  |
| <b>exchange</b>      | 3. $A, B \in \mathcal{I}$ and $ A  >  B  \Rightarrow \exists e \in A \setminus B \text{ such that } B \cup \{e\} \in \mathcal{I}$ |

Let  $\mathcal{N} = \{1, 2, 3, 4\}$ . The smallest matroid that contains  $\{1, 2\}$  and  $\{3, 4\}$

$\mathcal{I} = \{ \emptyset, \dots \}$  *by the non-emptiness property*  
 $\{1\}, \{2\}, \{3\}, \{4\}, \{1,2\}, \{3,4\}, \dots$  *by the heredity property*  
 $\{1,3\}, \{1,4\}, \{2,3\}, \{2,4\}$  *by the exchange property*  
}

# Can we project onto non-convex sets?

PMAP-0:

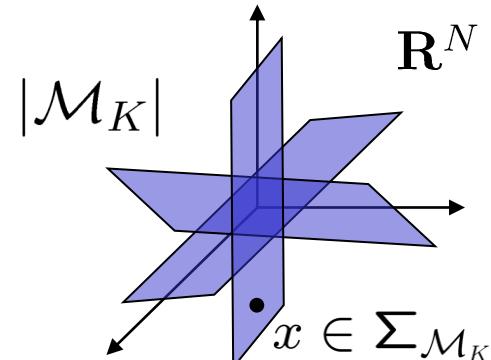
- Matroid structured sparse models:

$$\mathcal{M} = (\mathcal{N}, \mathcal{I} \subseteq 2^{\mathcal{N}}), \mathcal{N} = \{1, \dots, N\}$$

Definition: 1.  $\emptyset \in \mathcal{I}$

2.  $A \in \mathcal{I}$  and  $B \subseteq A \Rightarrow B \in \mathcal{I}$

3.  $A, B \in \mathcal{I}$  and  $|A| > |B| \Rightarrow \exists e \in A \setminus B$  such that  $B \cup \{e\} \in \mathcal{I}$



**Greedy basis algorithm efficiently solves**

$$\arg \max_{\mathcal{S}: \mathcal{S} \in \mathcal{M}} \sum_{i \in \mathcal{S}} w_i^2$$

sort  $\mathcal{N}$  in decreasing order by weight  $w_i^2$

start with empty set:  $\mathcal{S}_0 = \emptyset$

1.  $\mathcal{R}_i = \{r_i \in \mathcal{N} \setminus \mathcal{S}_i\}$  while keeping the order

2.  $r = \arg \max_j \{w_j^2 : (j \in \mathcal{R}_i) \wedge (\mathcal{S}_i \cup \{j\} \in \mathcal{I})\}$

3.  $\mathcal{S}_{i+1} = \mathcal{S}_i \cup \{r\}$

# Can we project onto non-convex sets?

PMAP-0:

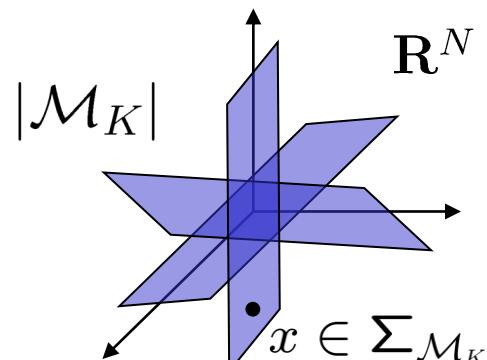
- Matroid structured sparse models:

$$\mathcal{M} = (\mathcal{N}, \mathcal{I} \subseteq 2^{\mathcal{N}}), \mathcal{N} = \{1, \dots, N\}$$

Definition: 1.  $\emptyset \in \mathcal{I}$

2.  $A \in \mathcal{I}$  and  $B \subseteq A \Rightarrow B \in \mathcal{I}$

3.  $A, B \in \mathcal{I}$  and  $|A| > |B| \Rightarrow \exists e \in A \setminus B$  such that  $B \cup \{e\} \in \mathcal{I}$



**Greedy basis algorithm efficiently solves matroid constrained problems**

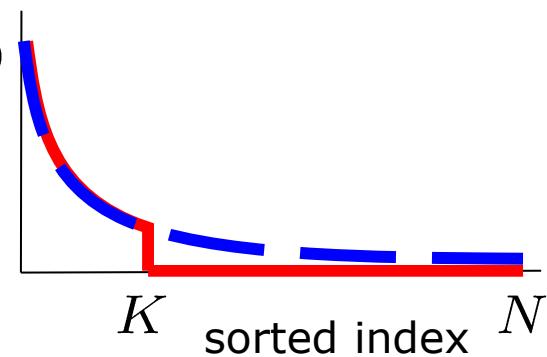
Examples:

1. uniform matroid:  $\mathcal{I} = \{\mathcal{S} : \mathcal{S} \subseteq \mathcal{N}, |\mathcal{S}| \leq K\}$

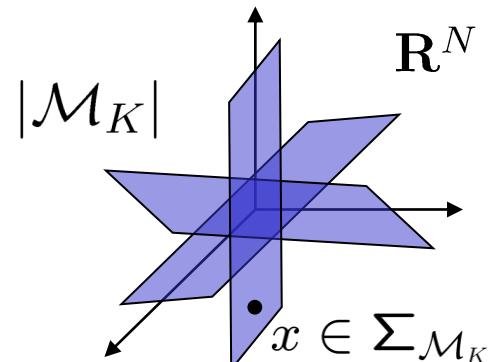
$$\text{prox}_{f_2}(y) = \arg \min_{x: x \in \Sigma_K} \|x - y\|$$

**hard thresholding!**

$$H_K(y)$$



# Can we project onto non-convex sets?



PMAP-0:

- Matroid structured sparse models:

$$\mathcal{M} = (\mathcal{N}, \mathcal{I} \subseteq 2^{\mathcal{N}}), \mathcal{N} = \{1, \dots, N\}$$

Definition: 1.  $\emptyset \in \mathcal{I}$

2.  $A \in \mathcal{I}$  and  $B \subseteq A \Rightarrow B \in \mathcal{I}$

3.  $A, B \in \mathcal{I}$  and  $|A| > |B| \Rightarrow \exists e \in A \setminus B$  such that  $B \cup \{e\} \in \mathcal{I}$

**Greedy basis algorithm efficiently solves matroid constrained problems**

Examples:

[Kyrillidis and C, 2011]

1. uniform matroid                       $\leftrightarrow$                       simple sparsity

***intersection with the following matroids (result is still a matroid!\*)***

2. partition matroid                       $\leftrightarrow$                       distributed sparsity

3. graphic matroid                       $\leftrightarrow$                       spanning tree sparsity

4. matching matroid                       $\leftrightarrow$                       graph matching sparsity

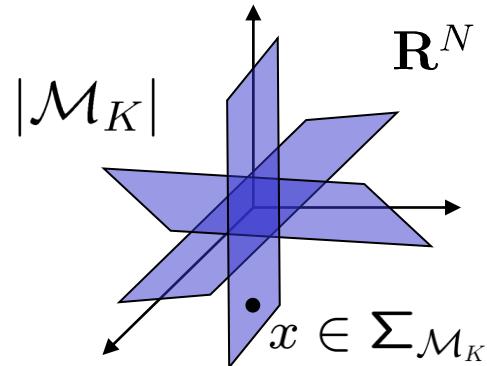
\*: in general, the intersection of two matroids is not a matroid.

# Can we project onto non-convex sets?

PMAP-0:

- Linear support constraints:

Definition:  $\Sigma_{\mathcal{M}_K} = \bigcup_{\forall z \in \mathfrak{Z}} \text{supp}(z)$ , where  $\mathfrak{Z} := \{z \in \{0, 1\}^N : Az \leq b\}$

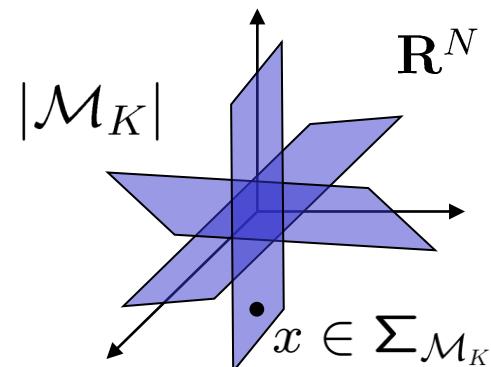


A and b	<>	integral
first row of A	<>	all 1's
first entry of b	<>	K

# Can we project onto non-convex sets?

PMAP-0:

- Linear support constraints:



Definition:  $\Sigma_{\mathcal{M}_K} = \bigcup_{\forall z \in \mathfrak{Z}} \text{supp}(z)$ , where  $\mathfrak{Z} := \{z \in \{0, 1\}^N : Az \leq b\}$

Example: neuronal spike model

$z \in \{0, 1\}^N$ : binary support variables

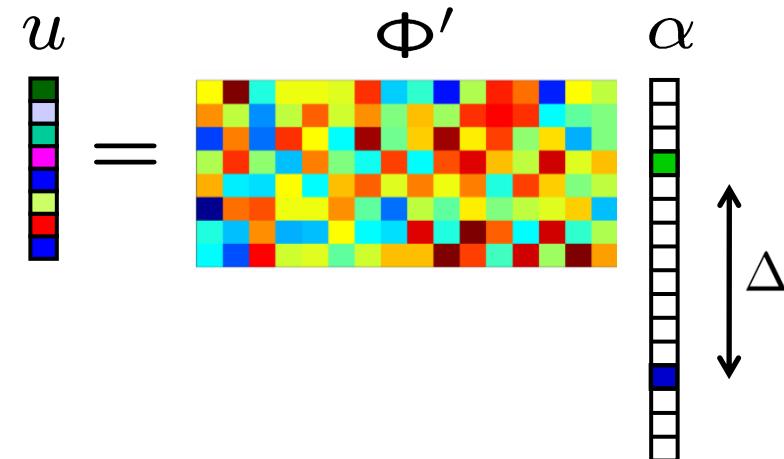
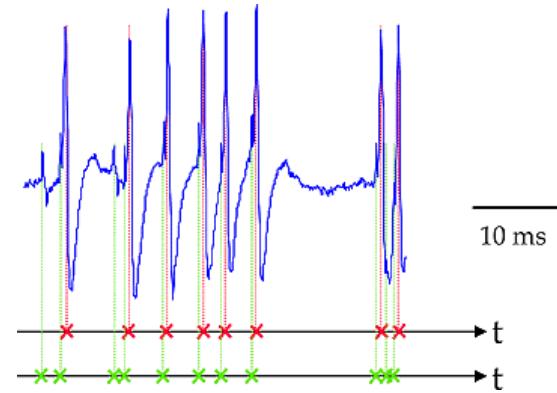
$$z_1 + z_2 + \dots + z_N \leq K$$

$$z_1 + z_2 + \dots + z_\Delta \leq 1$$

$$z_2 + z_3 + \dots + z_{\Delta+1} \leq 1$$

⋮

$$z_{N-\Delta+1} + z_{N-\Delta+2} + \dots + z_N \leq 1$$

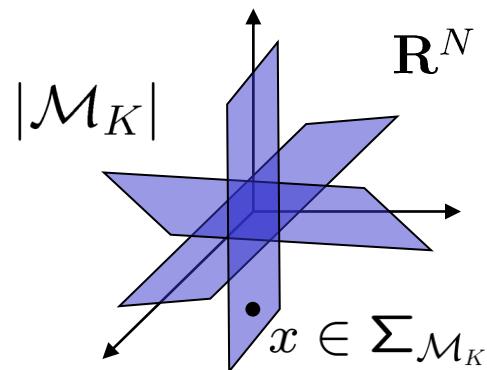


# Can we project onto non-convex sets?

PMAP-0:

- Linear support constraints:

Definition:  $\Sigma_{\mathcal{M}_K} = \bigcup_{\forall z \in \mathfrak{Z}} \text{supp}(z)$ , where  $\mathfrak{Z} := \{z \in \{0, 1\}^N : Az \leq b\}$

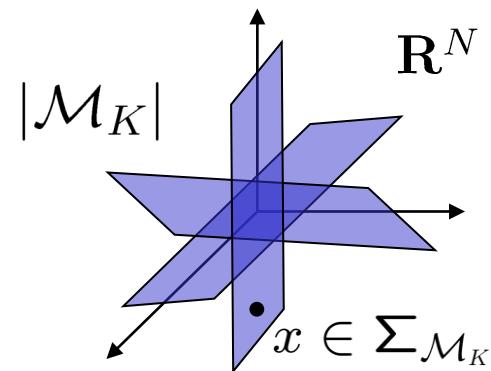


We can use LP can relax the LS constrained ILPs:

$$\arg \min_z \left\{ \rho^T z : z \in [0, 1]^N, Az \leq b \right\} \quad \rho_i = -|y_i|^2$$

...but, when is the result binary?

# Can we project onto non-convex sets?



PMAP-0:

- Linear support constraints:

Definition:  $\Sigma_{\mathcal{M}_K} = \bigcup_{\forall z \in \mathfrak{Z}} \text{supp}(z)$ , where  $\mathfrak{Z} := \{z \in \{0, 1\}^N : Az \leq b\}$

**LP can exactly solve the LS constrained ILPs:**

$$\arg \min_z \left\{ \rho^T z : z \in [0, 1]^N, Az \leq b \right\} \quad \rho_i = -|y_i|^2$$

**...when A is totally unimodular (TU)\*!**

[Nemhauser and Wolsey, 1999]

- the determinant of each square submatrix is {-1,0,1}

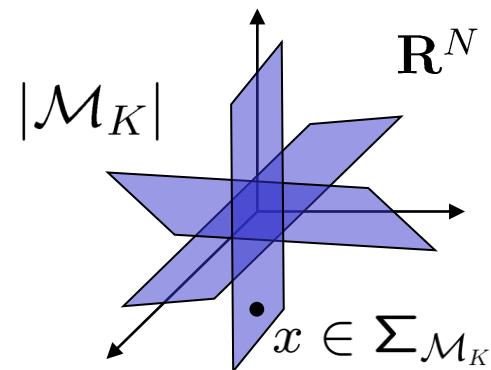
Examples: interval matrices, perfect matrices, network matrices

\*: if we want LP relaxation to work for all b, TU is a necessary condition.

# Can we project onto non-convex sets?

PMAP-0:

- Linear support constraints:

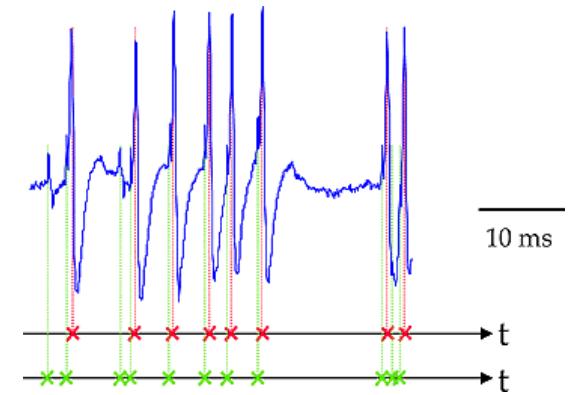


Definition:  $\Sigma_{\mathcal{M}_K} = \bigcup_{\forall z \in \mathfrak{Z}} \text{supp}(z)$ , where  $\mathfrak{Z} := \{z \in \{0, 1\}^N : Az \leq b\}$

Example: neuronal spike model

$z \in \{0, 1\}^N$ : binary support variables

$$\left. \begin{array}{l} z_1 + z_2 + \dots + z_N \leq K \\ z_1 + z_2 + \dots + z_\Delta \leq 1 \\ z_2 + z_3 + \dots + z_{\Delta+1} \leq 1 \\ \vdots \\ z_{N-\Delta+1} + z_{N-\Delta+2} + \dots + z_N \leq 1 \end{array} \right\} \text{TU}$$



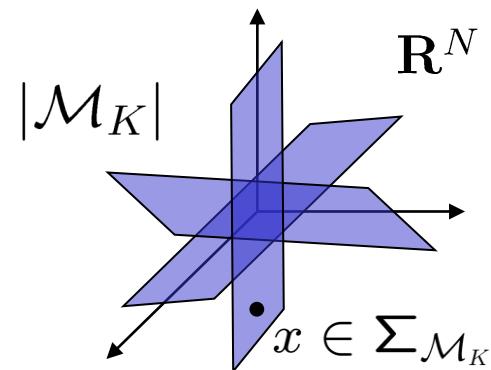
$$u = \Phi' \alpha$$

Diagram illustrating the relationship between input vector  $u$ , transformation matrix  $\Phi'$ , and output vector  $\alpha$ . The input  $u$  is a vertical vector with colored segments. The transformation  $\Phi'$  is represented by a grid of colored squares. The output  $\alpha$  is a vertical vector with black segments. A double-headed arrow labeled  $\Delta$  indicates the width of the input vector  $u$ .

# Can we project onto non-convex sets?

PMAP-0:

- prox-sparse models



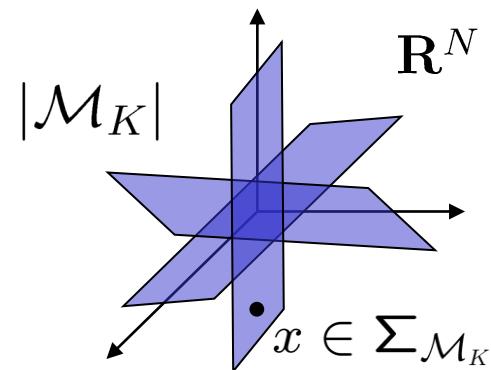
Definition: define algorithmically!

$$\text{prox}_{f_2}(y) = \arg \min_{x:x \in \Sigma_{\mathcal{M}_K}} \|x - y\| \quad g(x) = 0$$

# Can we project onto non-convex sets?

PMAP-0:

- prox-sparse models

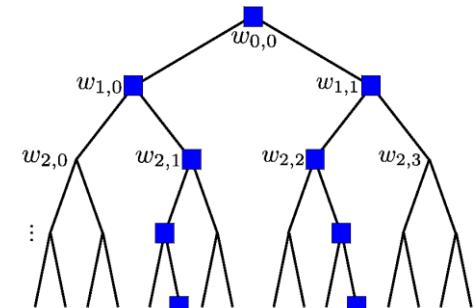


Definition: define algorithmically!

$$\text{prox}_{f_2}(y) = \arg \min_{x: x \in \Sigma_{\mathcal{M}_K}} \|x - y\| \quad g(x) = 0$$

Example: clustered sparsity models

- tree-sparse      <>      dynamic program
- clustered sparse    <>      dynamic program



# Can we project onto non-convex sets?

**Pop-quiz: A prox with convex and non-convex terms**

Let us consider  $f_2(x) = \|x\|_1 + \iota_{\{x: \|x\|_0 \leq K\}}(x)$   $g(x) = \|x\|_1$

$$\text{prox}_{f_2}(y) = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|y - x\|_2^2 + f_2(x)$$

Is it PMAP-0?

# Can we project onto non-convex sets?

**Pop-answer: A prox with convex and non-convex terms**

Let us consider  $f_2(x) = \|x\|_1 + \iota_{\{x: \|x\|_0 \leq K\}}(x)$   $g(x) = \|x\|_1$

$$\text{prox}_{f_2}(y) = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|y - x\|_2^2 + f_2(x)$$

$$\text{supp}(\text{prox}_{f_2}(y)) = \arg \max_{\mathcal{S}: |\mathcal{S}| \leq K} F(\mathcal{S}; y)$$

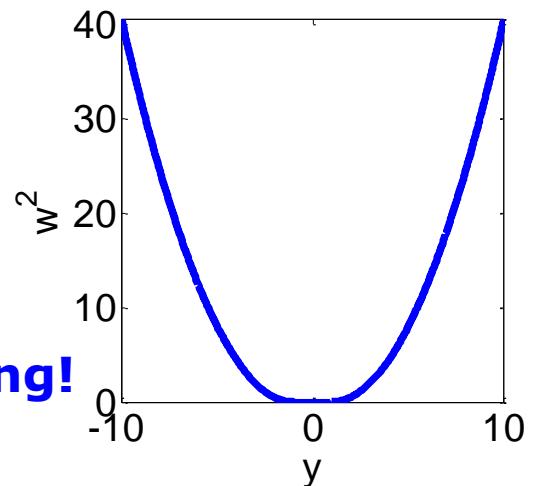


$$F(\mathcal{S}; y) = \frac{1}{2} \|y\|^2 - \min_{x: \text{supp}(x) = \mathcal{S}} \frac{1}{2} \|y - x\|_2^2 + \|x\|_1$$

$$\Rightarrow F(\mathcal{S}; y) = \sum_{i \in \mathcal{S}} w_i^2$$

$$w_i^2 = y_i \times \text{soft}(y_i, 1) - \frac{1}{2} |\text{soft}(y_i, 1)|^2 - |\text{soft}(y_i, 1)|$$

**Hard thresholding followed by soft thresholding!**



**YES: certified PMAP-0**

# Can we project onto non-convex sets?

PMAP-epsilon:  $F(\hat{\mathcal{S}}_\epsilon; y) \geq (1 - \epsilon) \max_{\mathcal{S} \in \bar{\mathcal{M}}_K} F(\mathcal{S}; y)$

- **Knapsack**

multi-knapsack constraints

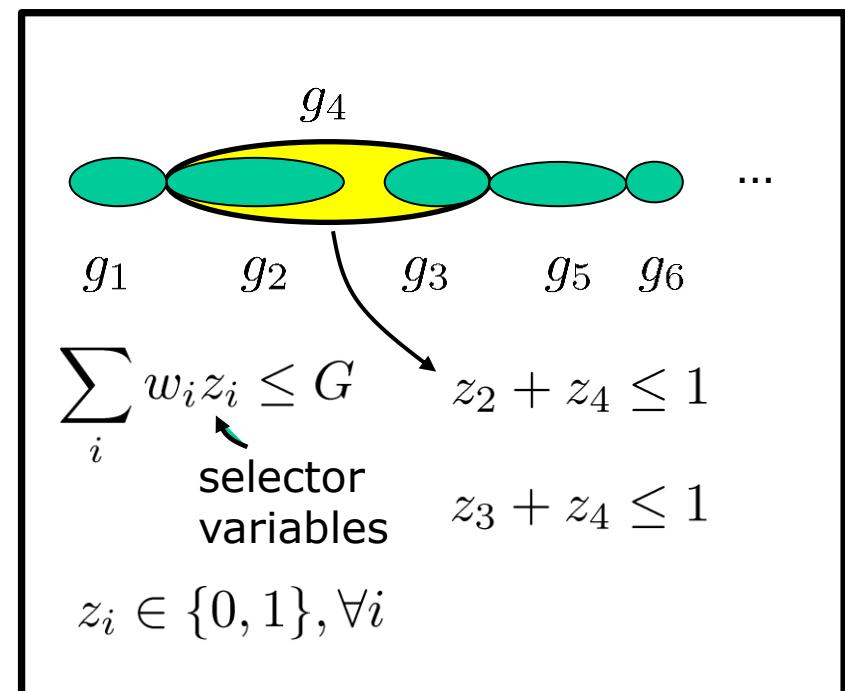
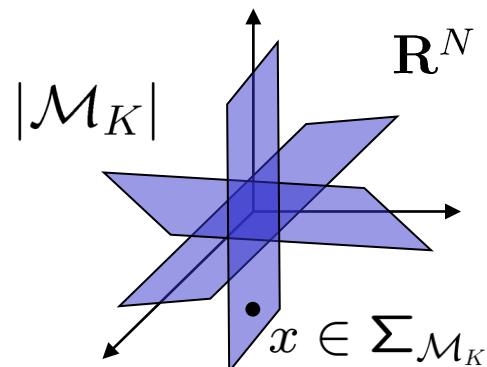
weighted multi-knapsack

Ex: Nested group sparse problems

quadratically-constrained

- **Define algorithmically!**

approximate solutions for computational reasons



# Can we project onto non-convex sets?

PMAP-epsilon:  $F(\hat{S}_\epsilon; y) \geq (1 - \epsilon) \max_{S \in \bar{\mathcal{M}}_K} F(S; y)$

- **Knapsack**

multi-knapsack constraints

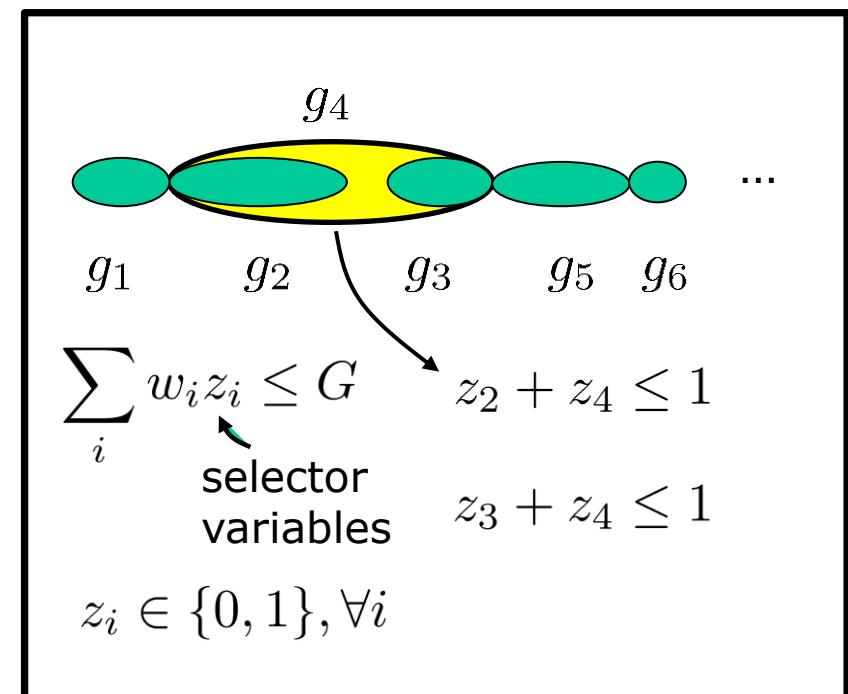
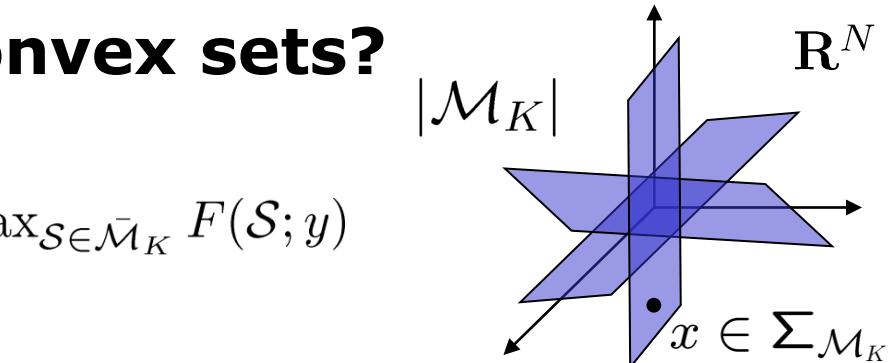
weighted multi-knapsack

Ex: Nested group sparse problems

quadratically-constrained

- **Define algorithmically!**

approximate solutions for computational reasons



- **Pairwise overlapping groups** <> quadratic binary w/ cardinality cons.

$$\max_{S: S \in \bar{\mathcal{M}}_K} F(S; y) = - \min \left\{ \sum_{i>j} \|(y)_{g_i \cap g_j}\|_2^2 z_i z_j - \sum_i \|(y)_{g_i}\|_2^2 z_i : \sum_i z_i \leq G \right\}.$$

**we can only approximate... and epsilon is large!**

# Can we project onto non-convex sets?

PMAP-epsilon:  $F(\hat{\mathcal{S}}_\epsilon; y) \geq (1 - \epsilon) \max_{\mathcal{S} \in \bar{\mathcal{M}}_K} F(\mathcal{S}; y)$

- **Knapsack**

- multi-knapsack constraints

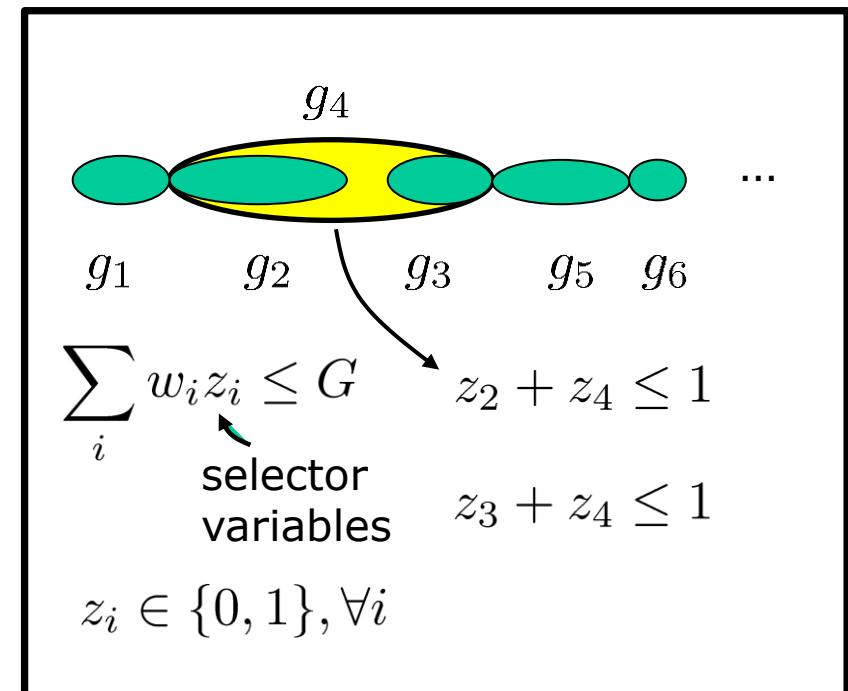
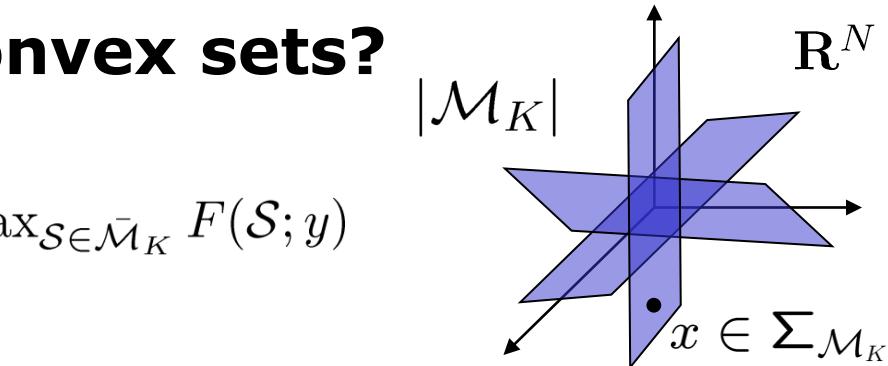
- weighted multi-knapsack

- Ex: Nested group sparse problems

- quadratically-constrained

- **Define algorithmically!**

- approximate solutions for computational reasons



- **Pairwise overlapping groups** <> quadratic binary w/ cardinality cons.

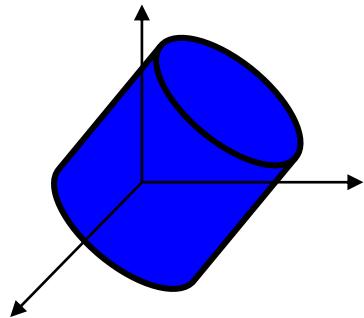
- **Multi-knapsack + multi-matroids**

[Lee et al., 2009]

**we can only approximate... and epsilon is large!**

# Can we project onto non-convex sets?

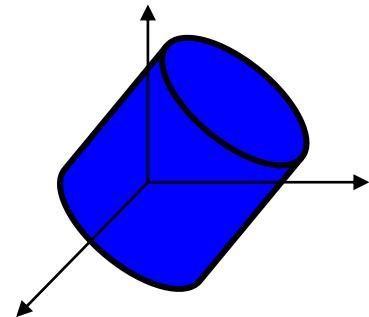
Matrix examples!



- Rank constrained projections  $\text{prox}_{f_2}(Y) = \arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F$

# Can we project onto non-convex sets?

Matrix examples!



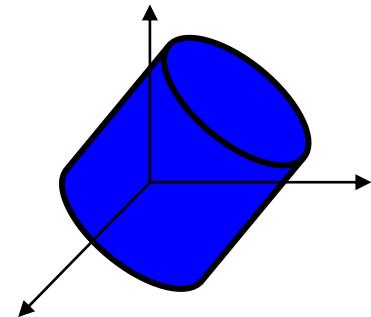
- Rank constrained projections  $\text{prox}_{f_2}(Y) = \arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F$

$$\arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F = \arg \min_{X:\text{rank}(X) \leq R} \|X - U\Lambda_Y V^T\|_F \begin{matrix} \text{singular value} \\ \text{decomposition} \end{matrix}$$



# Can we project onto non-convex sets?

Matrix examples!



- Rank constrained projections  $\text{prox}_{f_2}(Y) = \arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F$

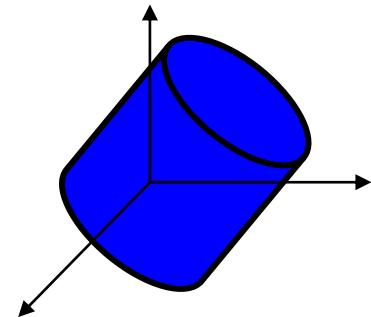
$$\arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F = \arg \min_{X:\text{rank}(X) \leq R} \|X - U\Lambda_Y V^T\|_F$$

$$= \arg \min_{X:\text{rank}(X) \leq R} \|U^T X V - \Lambda_Y\|_F \begin{matrix} \text{invariance to} \\ \text{unitary transform} \end{matrix}$$



# Can we project onto non-convex sets?

Matrix examples!



- Rank constrained projections  $\text{prox}_{f_2}(Y) = \arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F$

$$\arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F = \arg \min_{X:\text{rank}(X) \leq R} \|X - U\Lambda_Y V^T\|_F$$

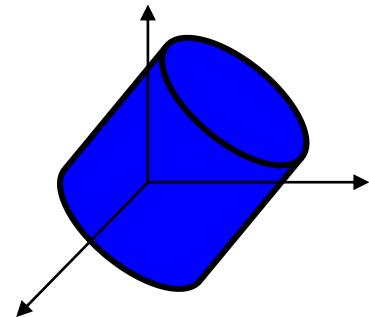
$$= U \left( \arg \min_{\tilde{X}:\text{rank}(\tilde{X}) \leq R} \|\tilde{X} - \Lambda_Y\|_F \right) V^T$$

sparse approximation problem!



# Can we project onto non-convex sets?

Matrix examples!



- Rank constrained projections  $\text{prox}_{f_2}(Y) = \arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F$

$$\arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F = \arg \min_{X:\text{rank}(X) \leq R} \|X - U\Lambda_Y V^T\|_F$$

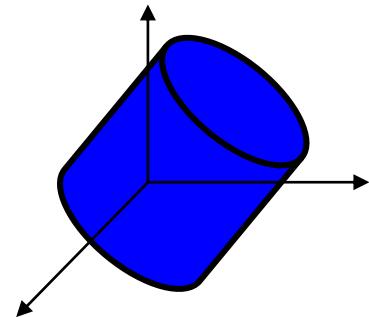
$$= U H_R(\Lambda_Y) V^T$$

singular value (hard) thresholding



# Can we project onto non-convex sets?

Matrix examples!



- Rank constrained projections  $\text{prox}_{f_2}(Y) = \arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F$

$$\arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F = \arg \min_{X:\text{rank}(X) \leq R} \|X - U\Lambda_Y V^T\|_F$$

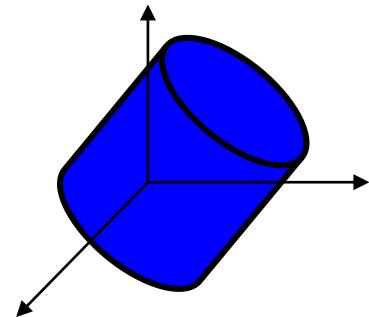
$$= U H_R(\Lambda_Y) V^T$$

singular value (hard) thresholding

- Non-convex spectral projections    <>    sets described by their eigenvalue properties
  - exact projections                  >>           basic operations on eigenvalues

# Can we project onto non-convex sets?

Matrix examples!



- Rank constrained projections  $\text{prox}_{f_2}(Y) = \arg \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F$
- Non-convex spectral projections <> sets described by their eigenvalue properties
- epsilon-approximate projections (note the difference with PMAP)  
$$\|\text{prox}_{f_2}^\epsilon(Y) - Y\|_F \leq (1 + \epsilon) \min_{X:\text{rank}(X) \leq R} \|X - Y\|_F$$

## Two highlights:

- structure from randomness/power methods [Halko, Martinsson, Tropp, 2010]
- column subset selection approaches [Boutsidis, Mahoney, Drineas, 2010]

# Recovery algorithms for low-dimensional models

**Now that we have projections...**

	Non-convex $\binom{N}{K}$	Convex	Probabilistic
Encoding	combinatorial / manifolds	atomic norm / convex relaxation	compressible / sparse priors

*A common criteria covering a broad set of applications:*

$$\min_X \|u - \Phi(X)\|^2 \text{ s.t. } X = S + L, \|S\|_0 \leq K, \text{rank}(L) \leq R$$

- *affine rank minimization, matrix completion, robust PCA...*

[Candes and Recht 2009; Waters, Sankaranayanan, Baraniuk, 2011]

*A common algorithm:*

**projected gradient**

$$\|S\|_0 = \#\{S_i \neq 0\}$$

# Recovery algorithms for low-dimensional models

**To highlight the salient differences, we will consider**

	Non-convex $\binom{N}{K}$	Convex	Probabilistic
Encoding	combinatorial / manifolds	atomic norm / convex relaxation	compressible / sparse priors

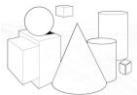
*compressive sensing recovery*

$$\min_{x: \|x\|_0 \leq K} \|u - \Phi x\|^2$$

*A common algorithm:*

***projected gradient***

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



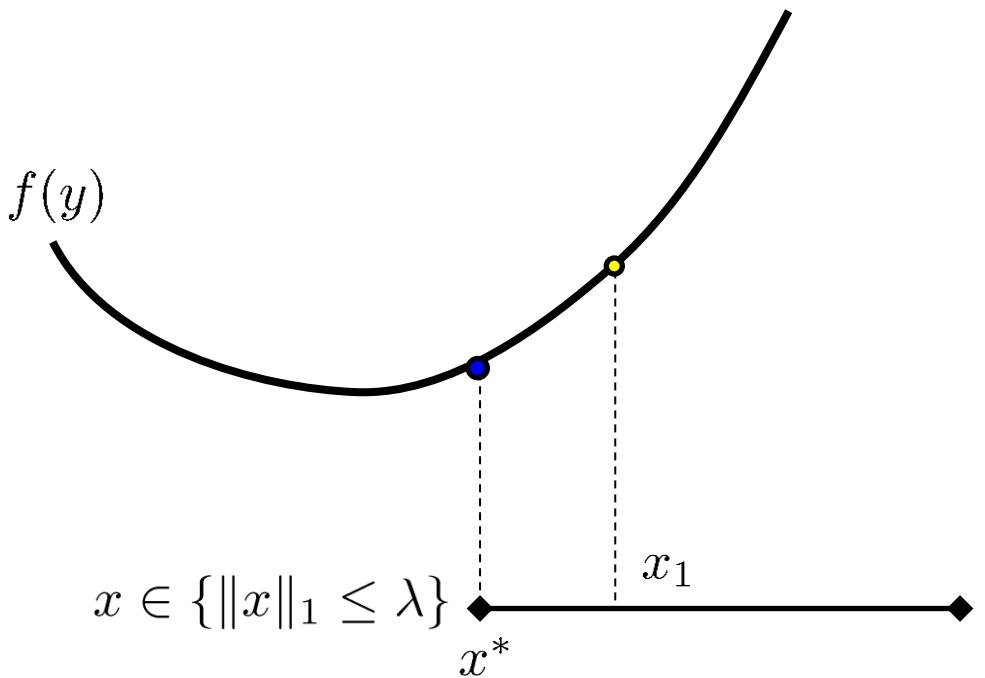
# A tale of two algorithms

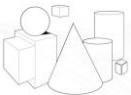
- Soft thresholding

$$f(x) = \|u - \Phi x\|^2$$

$$\min_{x: \|x\|_1 \leq \lambda} f(x)$$

---



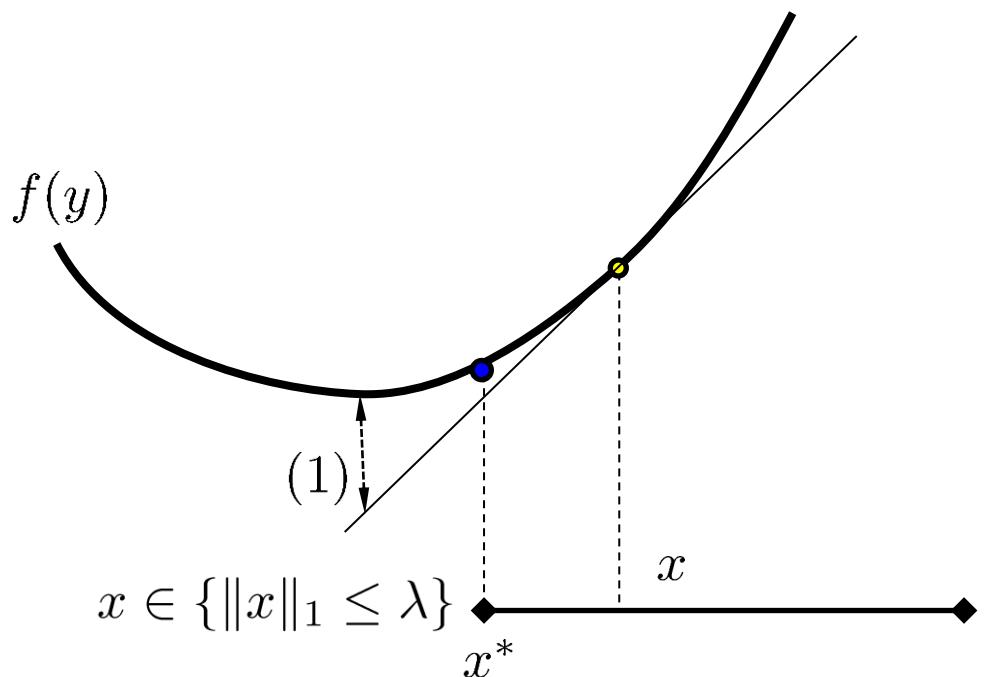


# A tale of two algorithms

- Soft thresholding

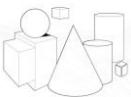
$$f(x) = \|u - \Phi x\|^2$$

$$\min_{x: \|x\|_1 \leq \lambda} f(x)$$



## Structure in optimization:

$$(1) \quad f(y) - f(x) - \langle \nabla f(x), y - x \rangle = \| \Phi(y - x) \|^2 \quad \forall x, y \in \mathcal{R}^N,$$



# A tale of two algorithms

- Soft thresholding

$$f(x) = \|u - \Phi x\|^2$$

$$\min_{x: \|x\|_1 \leq \lambda} f(x)$$

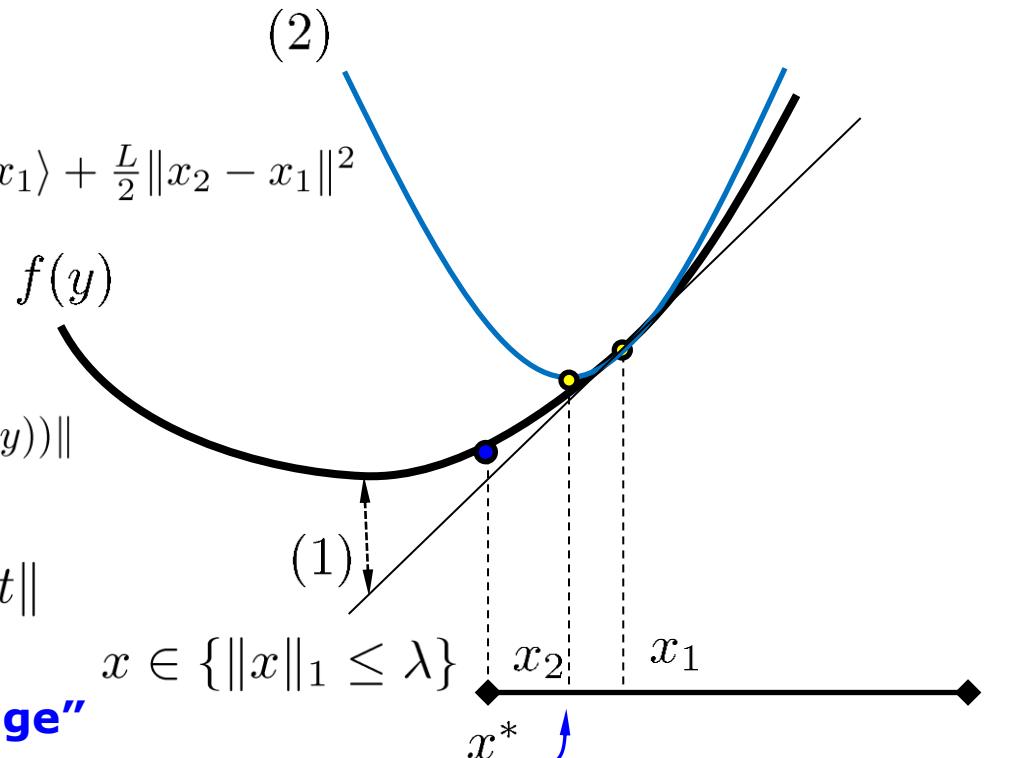
$$U(x_2, x_1) = f(x_1) + \langle \nabla f(x_1), x_2 - x_1 \rangle + \frac{L}{2} \|x_2 - x_1\|^2$$

majorization-minimization

$$\arg \min_{\|x\|_1 \leq \lambda} U(x, y) = \arg \min_{\|x\|_1 \leq \lambda} \|x - (y - \frac{1}{L} \nabla f(y))\|$$

$$\text{St}_{\{\|x\|_1 \leq \lambda\}}(t) = \arg \min_{\|x\|_1 \leq \lambda} \|x - t\|$$

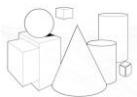
**Key actor: “least absolute shrinkage”**



$$x_{i+1} = \text{St}_{\{\|x\|_1 \leq \lambda\}} \left( x_i - \frac{1}{L} \nabla f(x_i) \right)$$

$$(1) \quad f(y) - f(x) - \langle \nabla f(x), y - x \rangle = \|\Phi(y - x)\|^2 \quad \forall x, y \in \mathcal{R}^N,$$

$$(2) \quad f(y) - f(x) - \langle \nabla f(x), y - x \rangle \leq \frac{L}{2} \|y - x\|^2 \quad L = 2\|\Phi\|^2, \forall x, y \in \mathcal{R}^N,$$



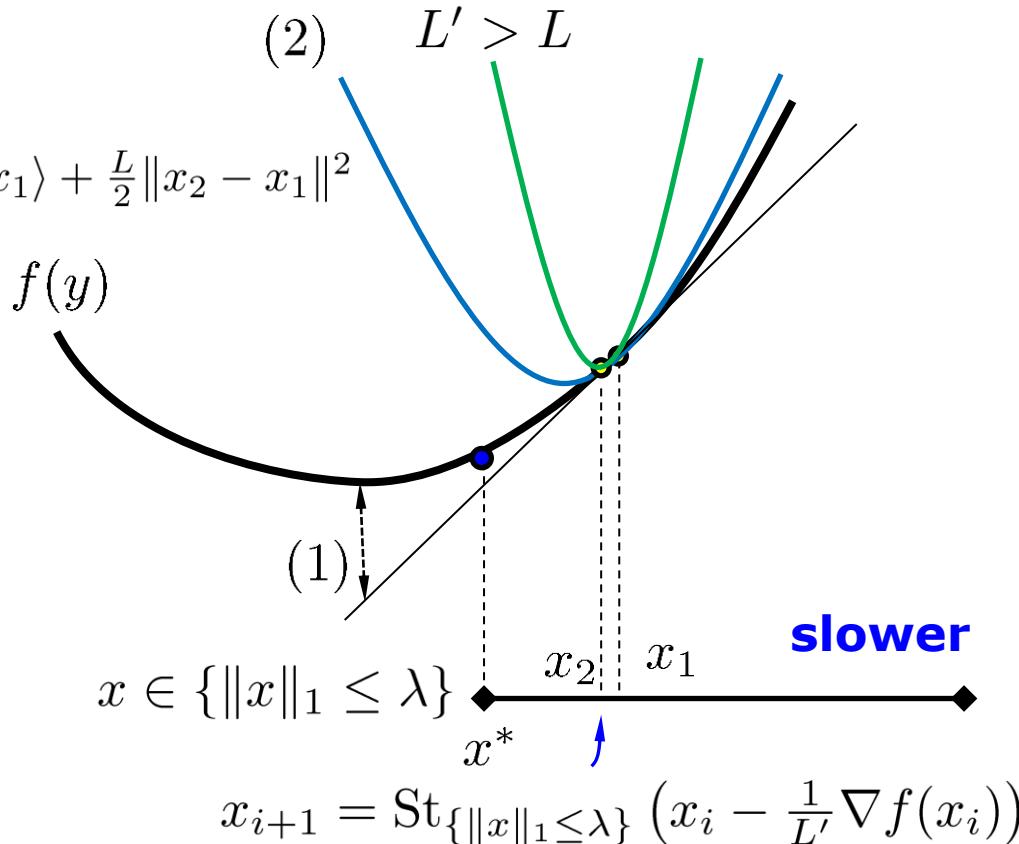
# A tale of two algorithms

- Soft thresholding

$$f(x) = \|u - \Phi x\|^2$$

$$\min_{x: \|x\|_1 \leq \lambda} f(x)$$

$$U(x_2, x_1) = f(x_1) + \langle \nabla f(x_1), x_2 - x_1 \rangle + \frac{L}{2} \|x_2 - x_1\|^2$$



$$(1) \quad f(y) - f(x) - \langle \nabla f(x), y - x \rangle = \|\Phi(y - x)\|^2 \quad \forall x, y \in \mathcal{R}^N,$$

$$(2) \quad f(y) - f(x) - \langle \nabla f(x), y - x \rangle \leq \frac{L}{2} \|y - x\|^2 \quad L = 2\|\Phi\|, \forall x, y \in \mathcal{R}^N,$$



# A tale of two algorithms

- Soft thresholding

$$\min_{x: \|x\|_1 \leq \lambda} f(x)$$

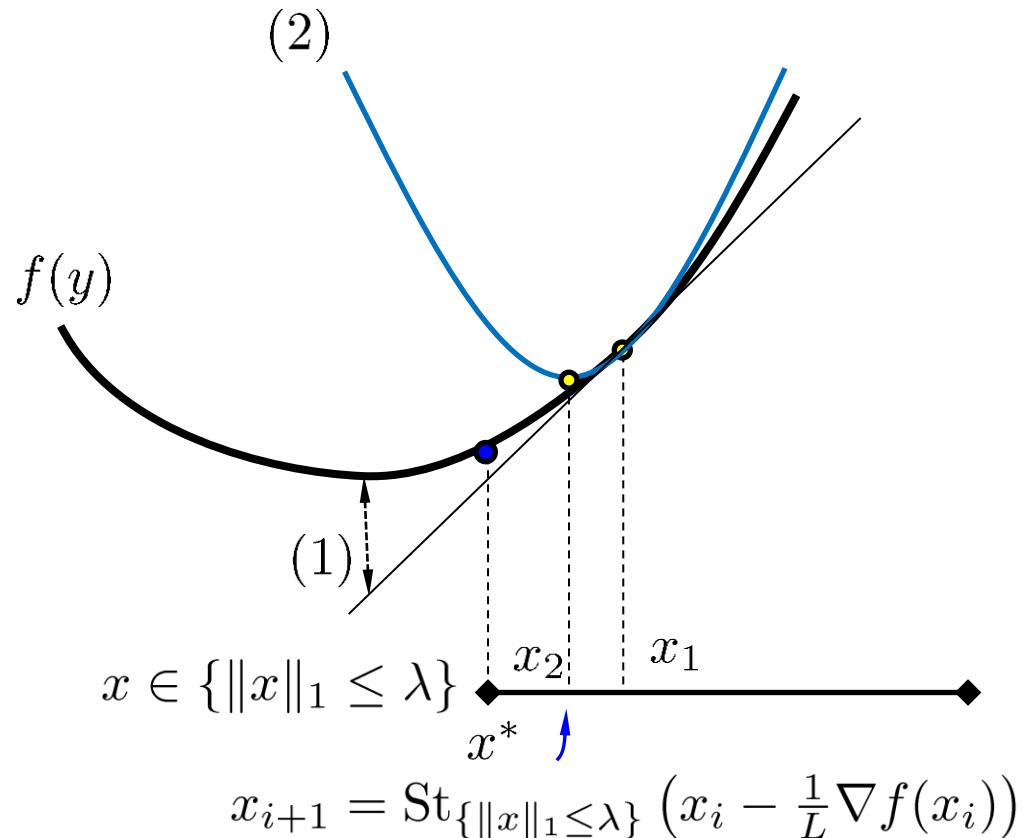
$$f(x) = \|u - \Phi x\|^2$$

- Is  $x^*$  what we are looking for?

local “unverifiable” assumptions:

- ERC/URC/RSC condition
- coherence based conditions ...

(local  $\rightarrow$  global / dual certification)



# A tale of two algorithms

$$\binom{N}{K}$$

- Hard thresholding

$$\min_{x: \|x\|_0 \leq K} f(x)$$

$$f(x) = \|u - \Phi x\|^2$$

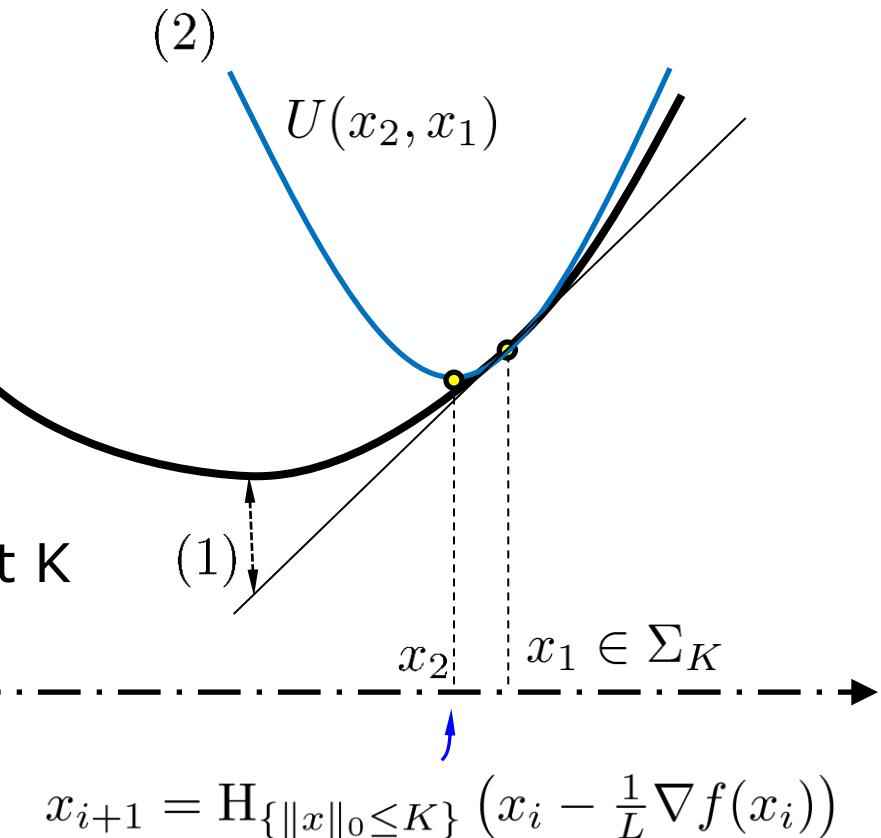
$$\arg \min_{\|x\|_0 \leq K} U(x, y) = \arg \min_{\|x\|_0 \leq K} \|x - (y - \frac{1}{L} \nabla f(y))\|$$

$$H_{\{\|x\|_0 \leq K\}}(t) = \arg \min_{\|x\|_0 \leq K} \|x - t\|$$

**Key actor: “hard thresholding”**

ALGO: sort and pick the largest K

$$y \in \Sigma_K$$



# A tale of two algorithms

$$\binom{N}{K}$$

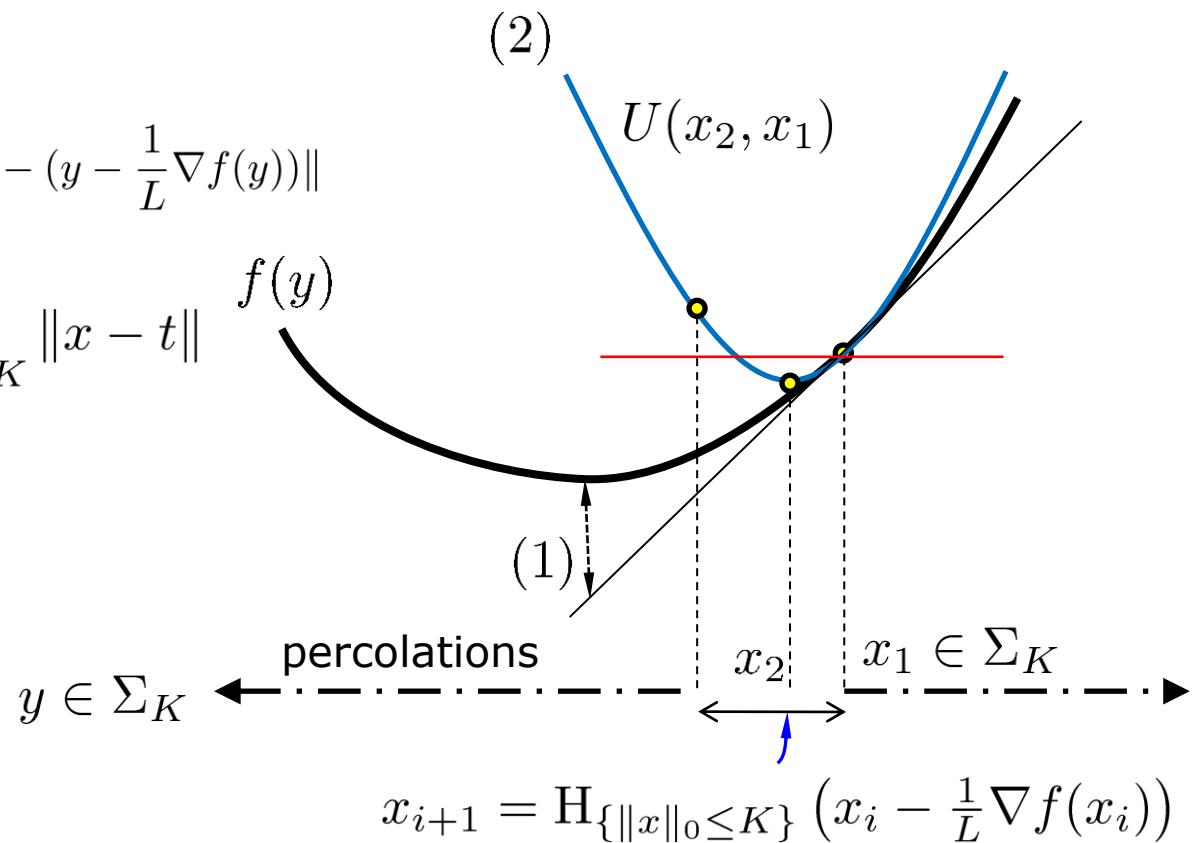
- Hard thresholding

$$f(x) = \|u - \Phi x\|^2$$

$$\min_{x: \|x\|_0 \leq K} f(x)$$

$$\arg \min_{\|x\|_0 \leq K} U(x, y) = \arg \min_{\|x\|_0 \leq K} \|x - (y - \frac{1}{L} \nabla f(y))\|$$

$$H_{\{\|x\|_0 \leq K\}}(t) = \arg \min_{\|x\|_0 \leq K} \|x - t\|$$



What could possibly go wrong with this naïve approach?

# A tale of two algorithms

 $\binom{N}{K}$ 

- Hard thresholding

$$f(x) = \|u - \Phi x\|^2$$

$$\min_{x: \|x\|_0 \leq K} f(x)$$

Global “unverifiable” assumption:

$$(1 - \delta_K) \leq \frac{\|\Phi x\|_2^2}{\|x\|_2^2} \leq (1 + \delta_K), \quad \forall x \in \Sigma_K$$

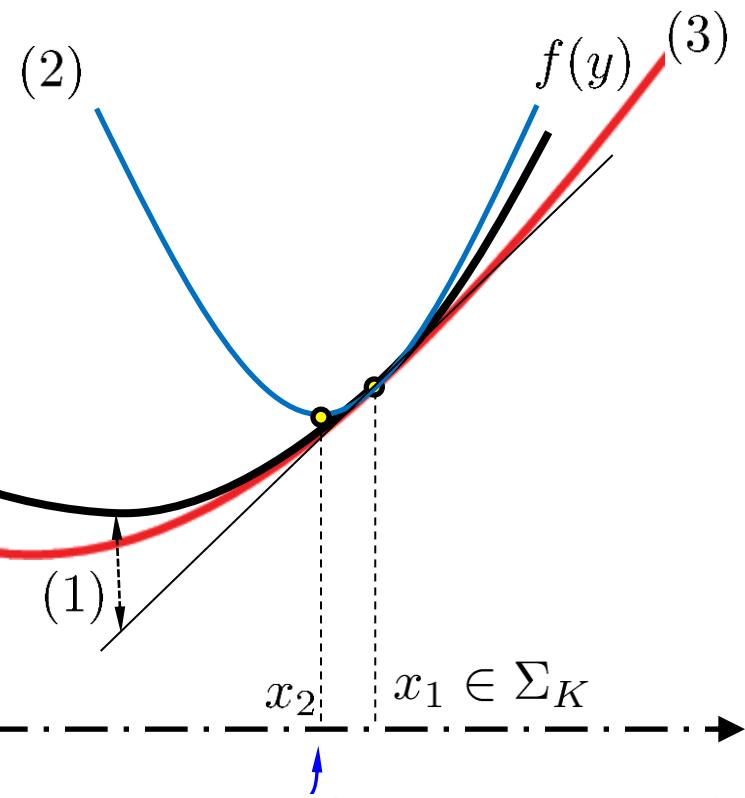
**RIP condition**  $M = O(K \log(N/K))$

⇒ we can tiptoe among percolations!

$$\underline{\delta_{2K} < 1/3}$$

$$y \in \Sigma_K$$

another variant has  $\delta_{3K} < 1/2$



$$x_{i+1} = H_{\{\|x\|_0 \leq K\}} \left( x_i - \frac{1}{L_{2K}} \nabla f(x_i) \right)$$

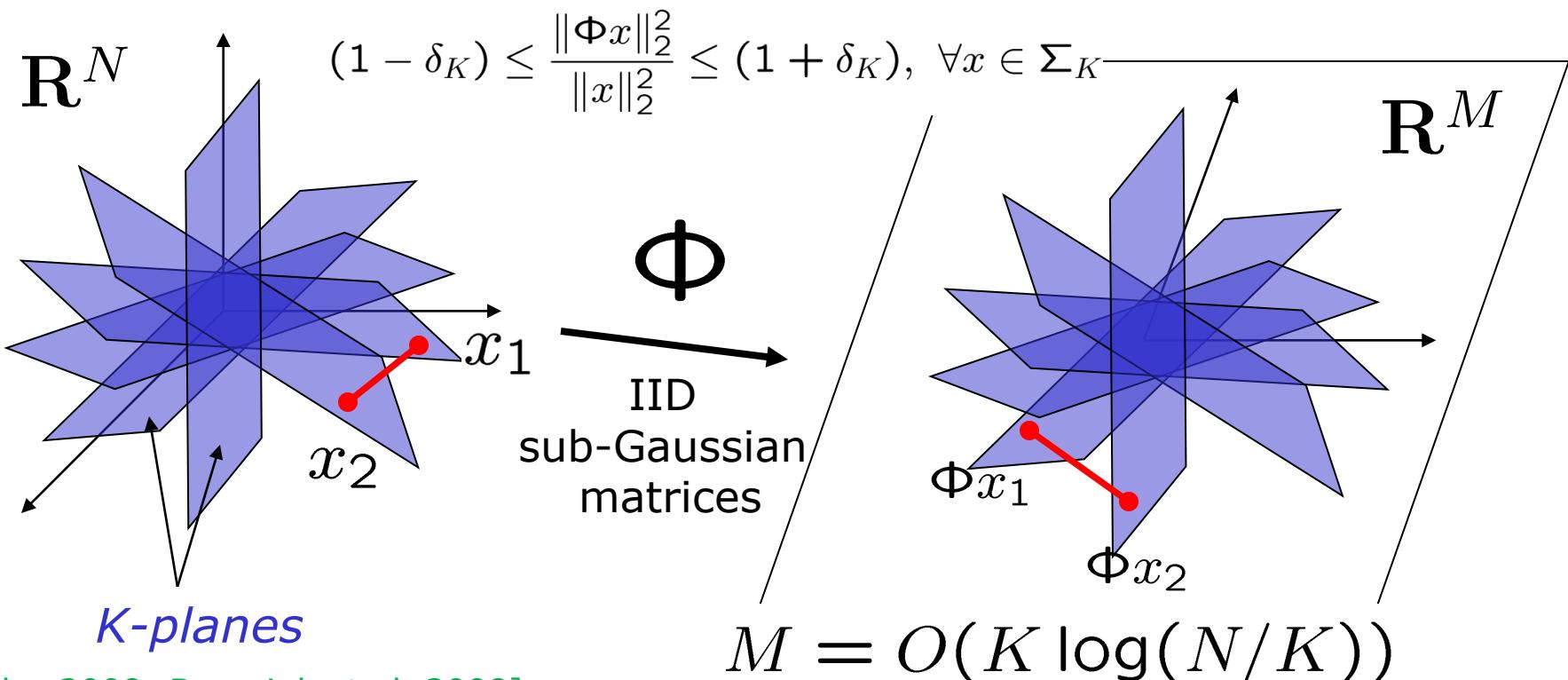
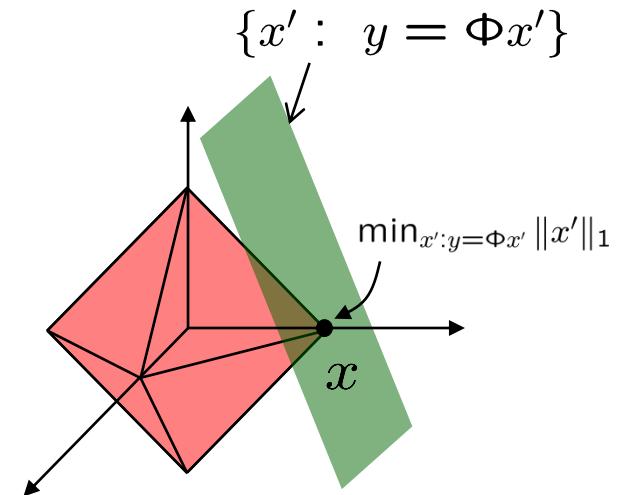
- (1)  $f(y) - f(x) - \langle \nabla f(x), y - x \rangle = \|\Phi(y - x)\|^2 \quad \forall x, y \in \mathcal{R}^N,$
- (2)  $f(y) - f(x) - \langle \nabla f(x), y - x \rangle \leq \frac{L_{2K}}{2} \|y - x\|^2 \quad L_{2K} = 2(1 + \delta_{2K}), \forall x, y \in \Sigma_K,$
- (3)  $f(y) - f(x) - \langle \nabla f(x), y - x \rangle \geq \frac{\mu_{2K}}{2} \|y - x\|^2 \quad \mu_{2K} = 2(1 - \delta_{2K}), \forall x, y \in \Sigma_K,$

# Restricted Isometry Property

- **Model:**  $K$ -sparse coefficients

**Remark:** implies convergence of convex relaxations also  
e.g.,  $\delta_{2K} < .465$  is sufficient for BP

- **RIP:** stable embedding

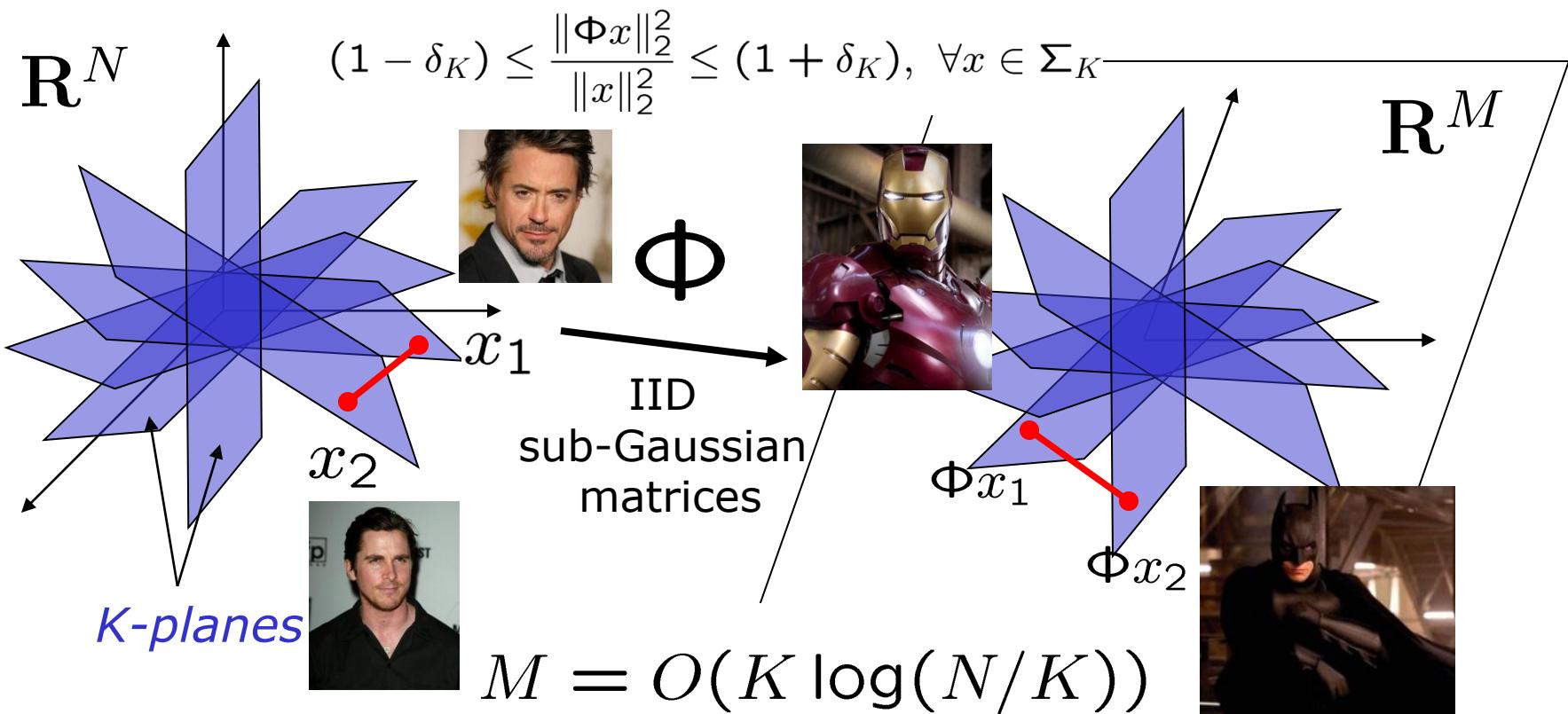
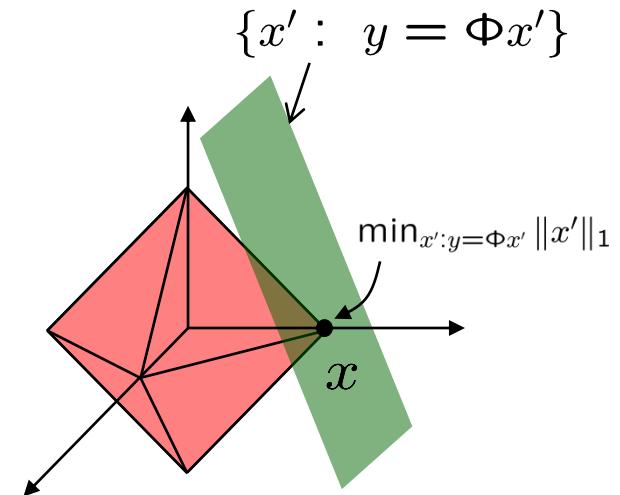


# Restricted Isometry Property

- **Model:**  $K$ -sparse coefficients

**Remark:** implies convergence of convex relaxations also  
e.g.,  $\delta_{2K} < .465$  is sufficient for BP

- **RIP:** stable embedding

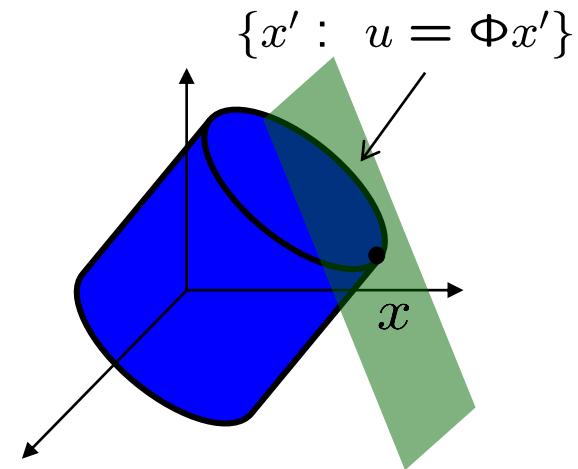


# Restricted Isometry Property for Matrices!

- **Model:** rank- $R$  matrices

**Remark:** bi-Lipschitz embedding  
of low-rank matrices

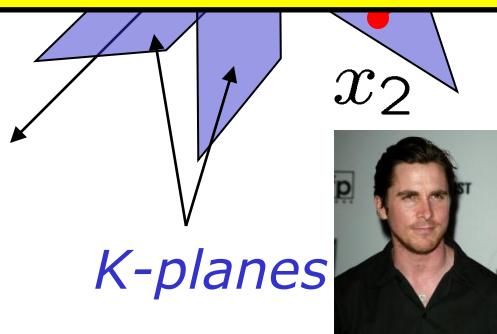
- **RIP:** stable embedding



$$\|\Phi\|_F^2$$

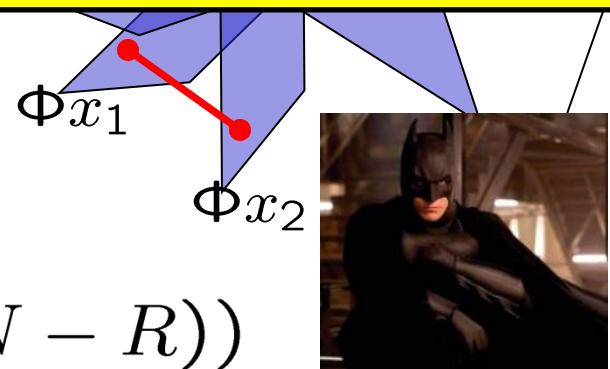
$$(1 - \delta_R) \leq \frac{\|\Phi X\|_F^2}{\|X\|_F^2} \leq (1 + \delta_R), \quad \forall X : \text{rank}(X) \leq R$$

[Plan 2011]



sub-Gaussian  
matrices

$$M = O(R(2N - R))$$



# Projected gradient method for non-convex sets

- Model-based hard thresholding  $f(x) = \|u - \Phi x\|^2$

$$\min_{x: x \in \Sigma_{\mathcal{M}_K}} f(x)$$

Global “unverifiable” assumption:

$$(1 - \delta_{\mathcal{M}_K}) \leq \frac{\|\Phi x\|_2^2}{\|x\|_2^2} \leq (1 + \delta_{\mathcal{M}_K}), \quad \forall x \in \Sigma_{\mathcal{M}_K}$$

$$H_{\Sigma_{\mathcal{M}_K}}(t) = \arg \min_{x: x \in \Sigma_{\mathcal{M}_K}} \|x - t\|$$

[Baraniuk, C, Duarte, Hegde 2010]

**Key actor: non-convex projector**

$$\underline{\delta_{\mathcal{M}_{2K}} < 1/3}$$

$$x_{i+1} = H_{\Sigma_{\mathcal{M}_K}} \left( x_i - \frac{1}{L_{\mathcal{M}_{2K}}} \nabla f(x_i) \right)$$

$$(1) \quad f(y) - f(x) - \langle \nabla f(x), y - x \rangle = \|\Phi(y - x)\|^2 \quad \forall x, y \in \mathbb{R}^N,$$

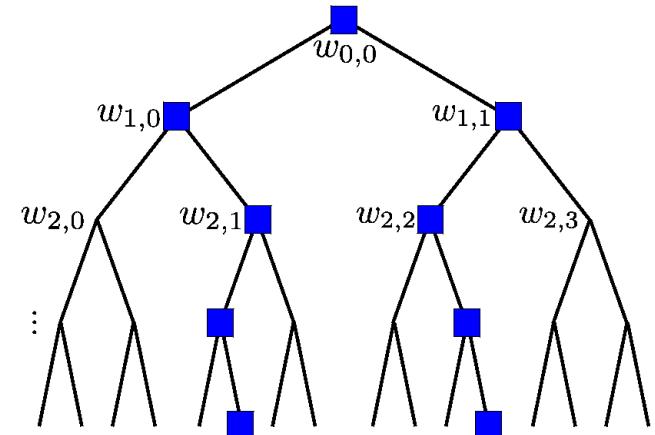
$$(2) \quad f(y) - f(x) - \langle \nabla f(x), y - x \rangle \leq \frac{L_{\mathcal{M}_{2K}}}{2} \|y - x\|^2 \quad L_{\mathcal{M}_{2K}} = 2(1 + \delta_{\mathcal{M}_{2K}}), \forall x, y \in \Sigma_{\mathcal{M}_{2K}},$$

$$(3) \quad f(y) - f(x) - \langle \nabla f(x), y - x \rangle \geq \frac{\mu_{\mathcal{M}_{2K}}}{2} \|y - x\|^2 \quad \mu_{\mathcal{M}_{2K}} = 2(1 - \delta_{\mathcal{M}_{2K}}), \forall x, y \in \Sigma_{\mathcal{M}_{2K}},$$

# Example: tree-sparse recovery

- **Model:**  $K$ -sparse coefficients

- + significant coefficients lie on a rooted subtree



- **Sparse approx:** find **best set** of coefficients

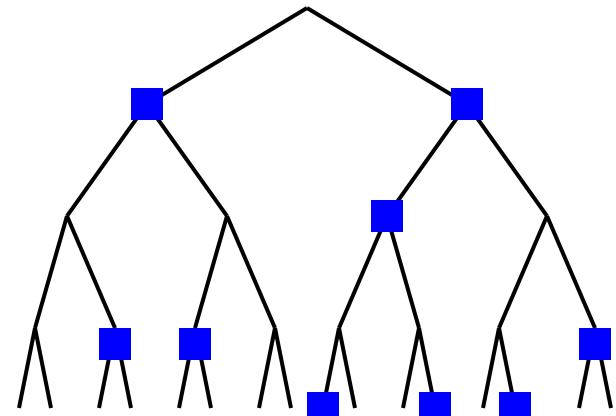
- sorting
  - hard thresholding

- **Tree-sparse approx:** find **best rooted subtree** of coefficients

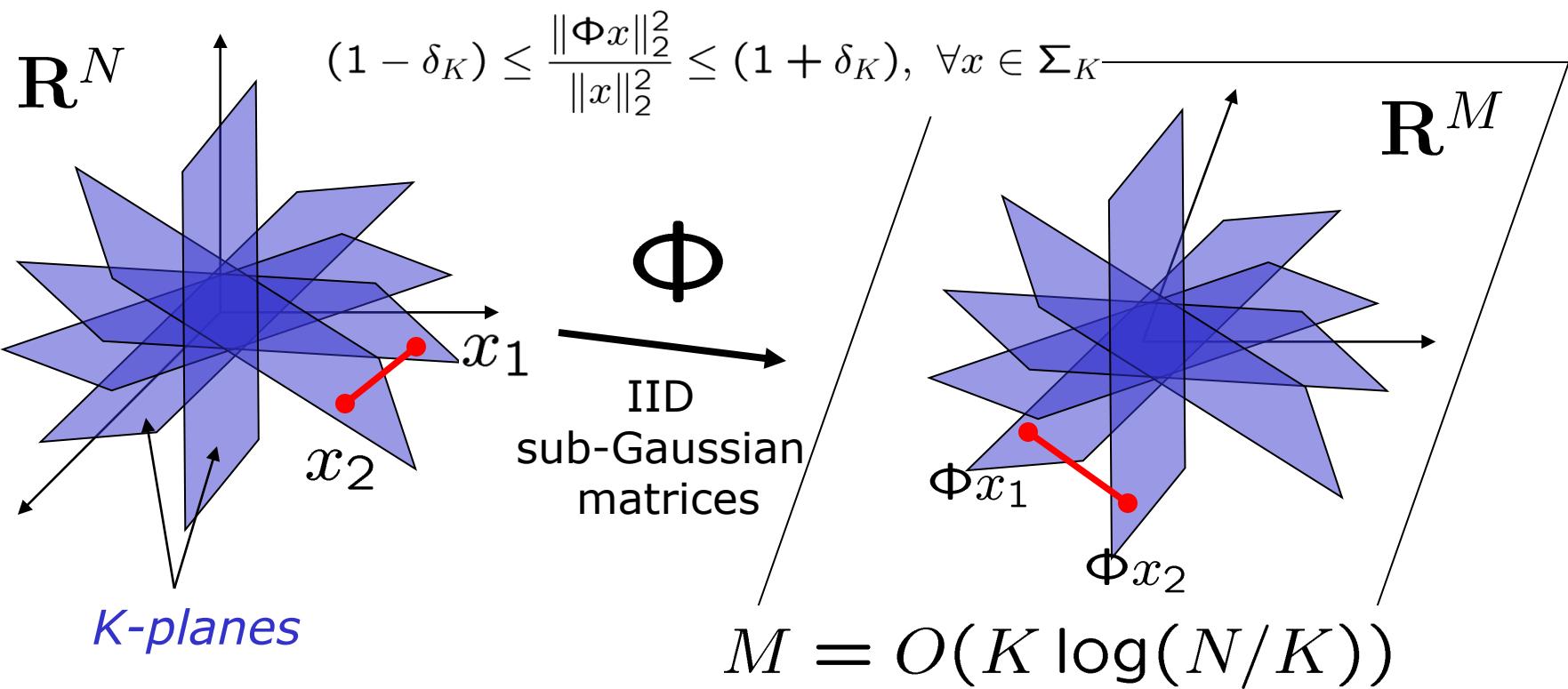
- condensing sort and select [Baraniuk]
  - dynamic programming [Donoho]

# Example: tree-sparse recovery

- **Model:**  $K$ -sparse coefficients

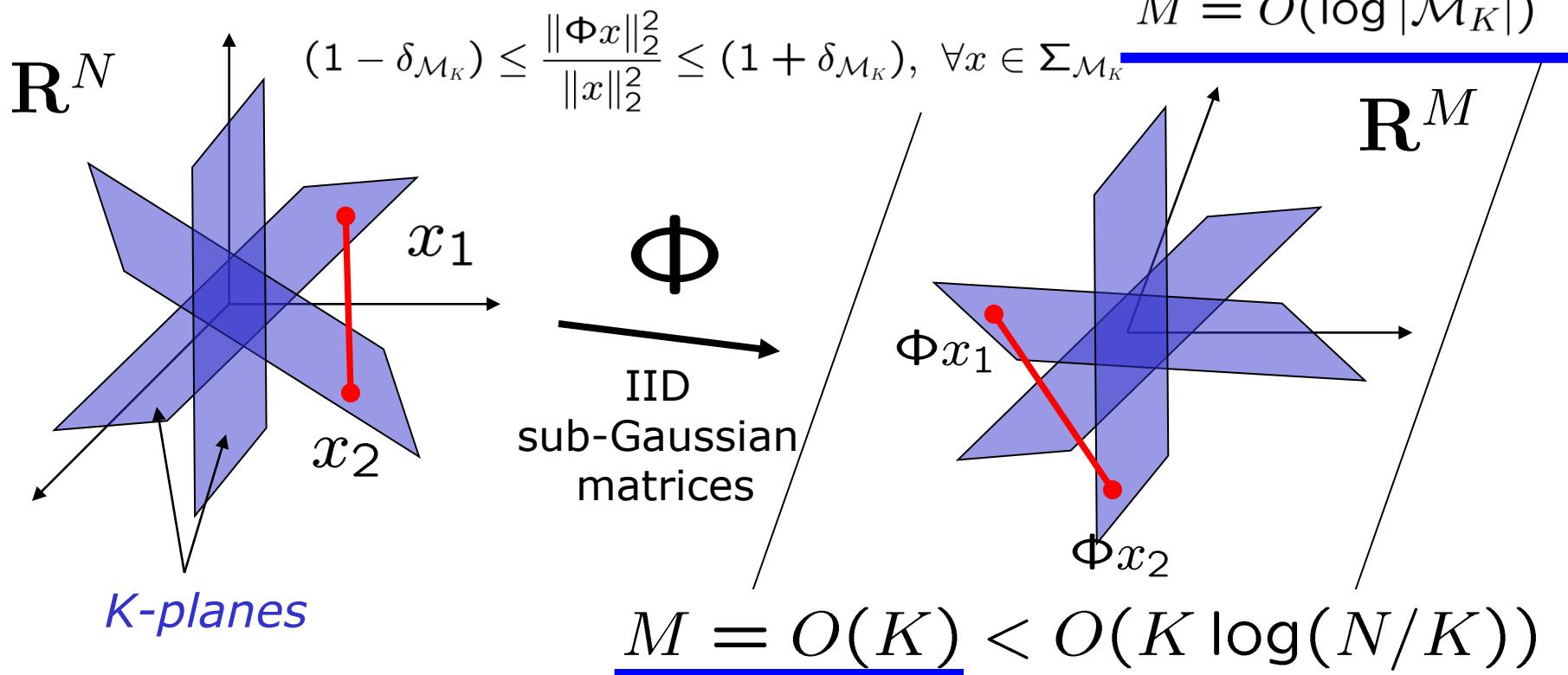
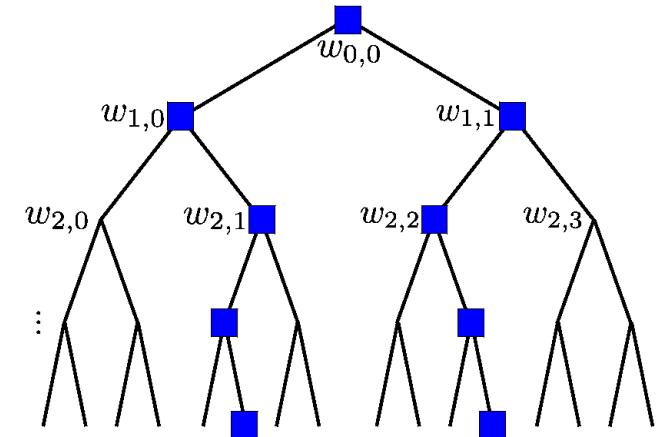


- **RIP:** stable embedding



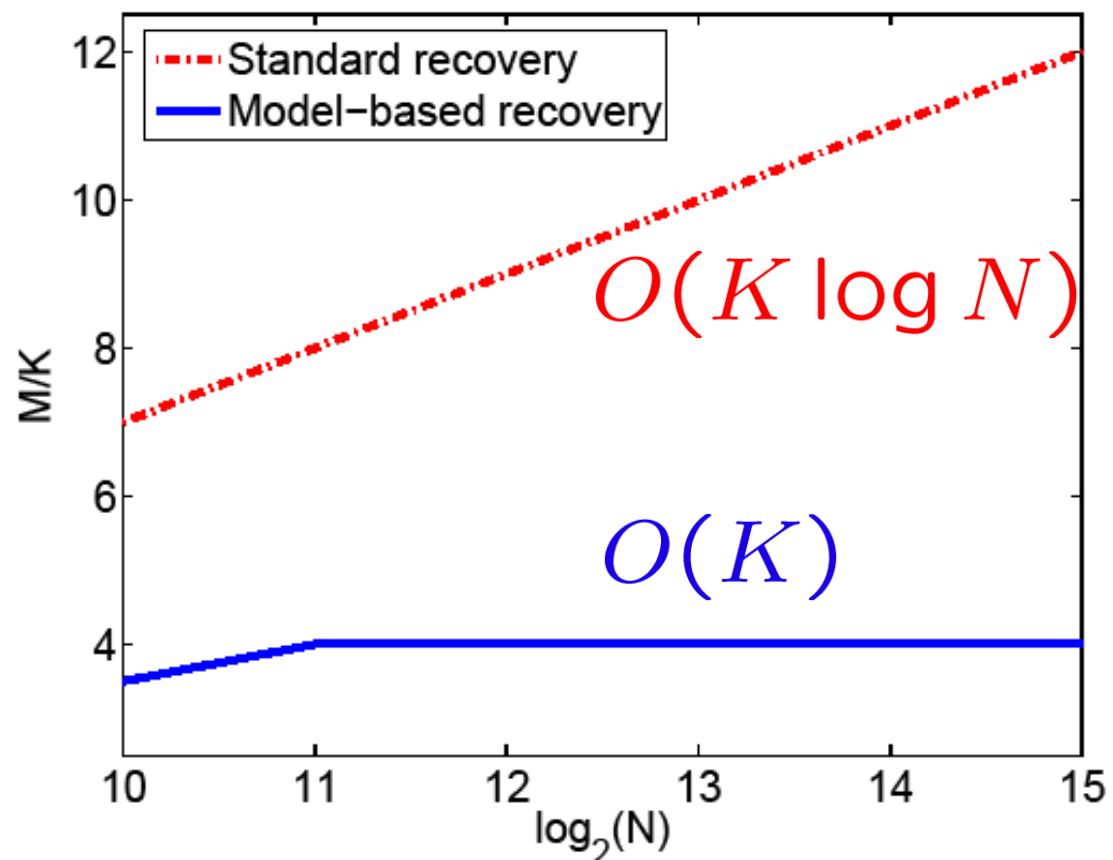
# Example: tree-sparse recovery

- **Model:**  $K$ -sparse coefficients
  - + significant coefficients lie on a rooted subtree
- **Tree-RIP:** stable embedding



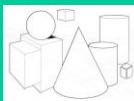
# Example: tree-sparse recovery

- Number samples for correct recovery
- Piecewise cubic signals + wavelets
- Models/algorithms:
  - compressible (CoSaMP)
  - tree-compressible (tree-CoSaMP)



# Recovery algorithms for low-dimensional models

## The Clash Operator

	<b>Non-convex</b> $\binom{N}{K}$	<b>Convex</b> 	<b>Probabilistic</b> 
Encoding	combinatorial / manifolds	atomic norm / convex relaxation	compressible / sparse priors
Example	$\min_{x: \ x\ _0 \leq K} \ u - \Phi x\ ^2$	$\min_{x: \ x\ _1 \leq \lambda} \ u - \Phi x\ ^2$	$E\{x u\}$
Algorithm	IHT, CoSaMP, SP, ALPS, OMP...	Basis pursuit, Lasso, basis pursuit denoising...	Variational Bayes, EP, Approximate message passing (AMP)...

$$\hat{x}_{\text{Clash}} = \arg \min_{x: \|x\|_0 \leq K, \|x\|_1 \leq \lambda} \|u - \Phi x\|^2$$

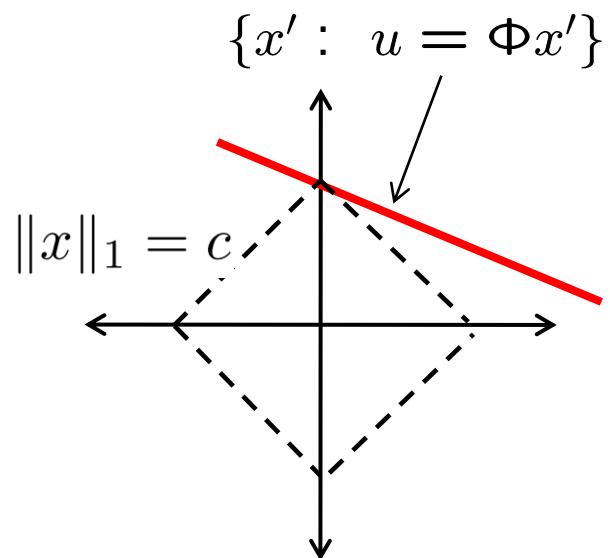
$$\|x\|_0 = \#\{x_i \neq 0\}$$

# Recovery algorithms for low-dimensional models

$$\hat{x} = \arg \min \|x\|_0 \text{ s.t. } u = \Phi x$$



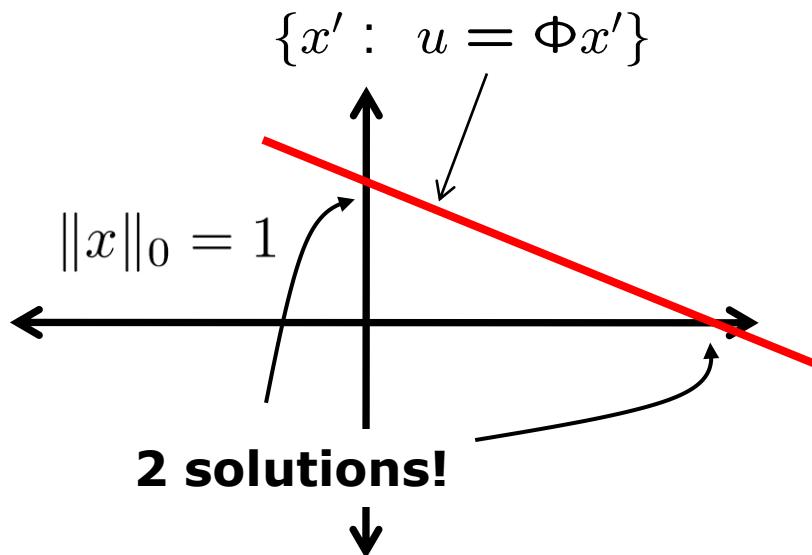
$$\hat{x} = \arg \min \|x\|_1 \text{ s.t. } u = \Phi x$$



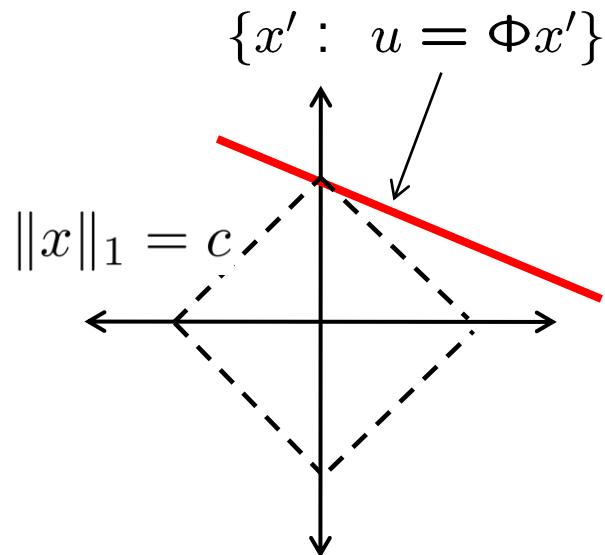
# Recovery algorithms for low-dimensional models

- A subtle issue

$$\hat{x} = \arg \min \|x\|_0 \text{ s.t. } u = \Phi x$$



$$\hat{x} = \arg \min \|x\|_1 \text{ s.t. } u = \Phi x$$

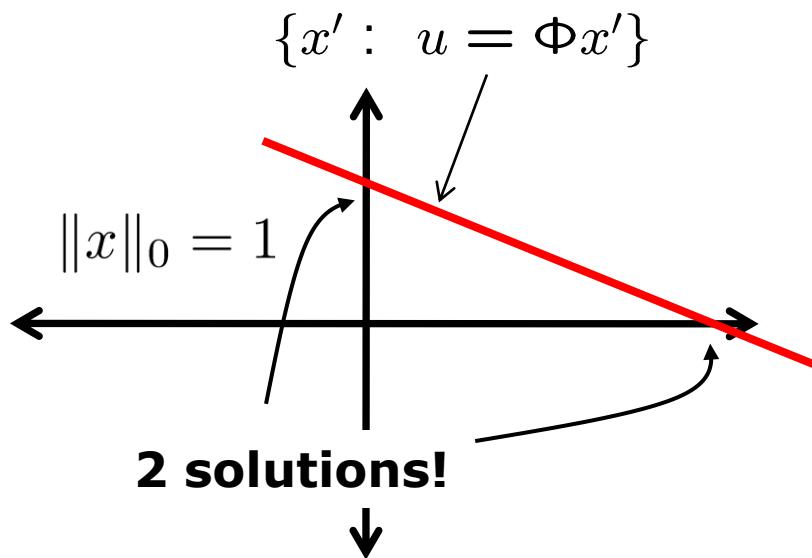


Which one is correct?

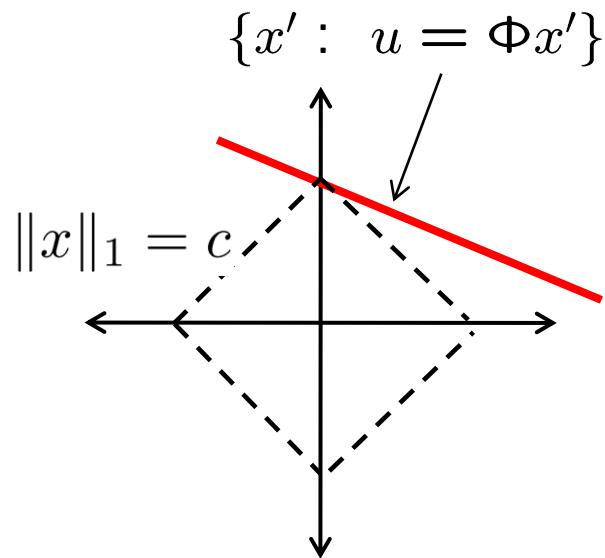
# Recovery algorithms for low-dimensional models

- A subtle issue

$$\hat{x} = \arg \min \|x\|_0 \text{ s.t. } u = \Phi x$$



$$\hat{x} = \arg \min \|x\|_1 \text{ s.t. } u = \Phi x$$

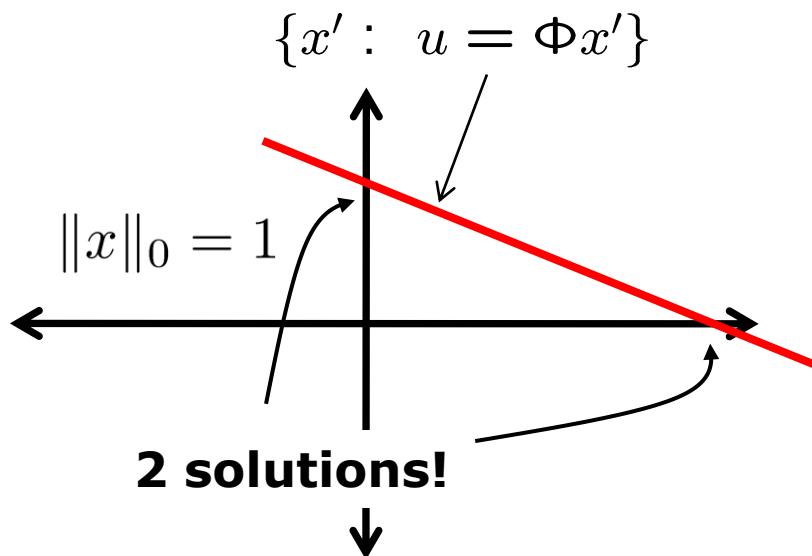


“Greed is good.” – Joel Tropp 2004

# Recovery algorithms for low-dimensional models

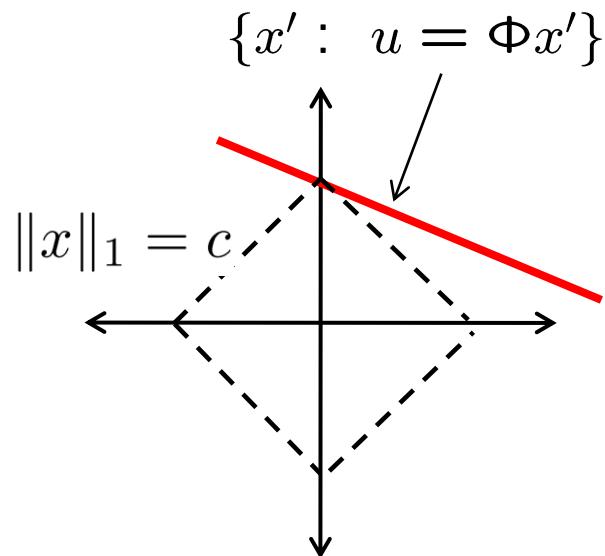
- A subtle issue

$$\hat{x} = \arg \min \|x\|_0 \text{ s.t. } u = \Phi x$$



Which one is correct?

$$\hat{x} = \arg \min \|x\|_1 \text{ s.t. } u = \Phi x$$

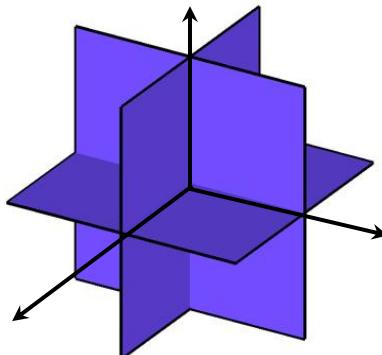


# The CLASH algorithm

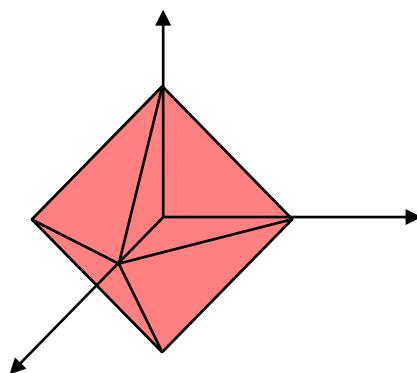
combinatorial selection  
+  
least absolute shrinkage

$$H_{\{\|x\|_0 \leq K\}}(t) = \arg \min_{\|x\|_0 \leq K} \|x - t\|$$

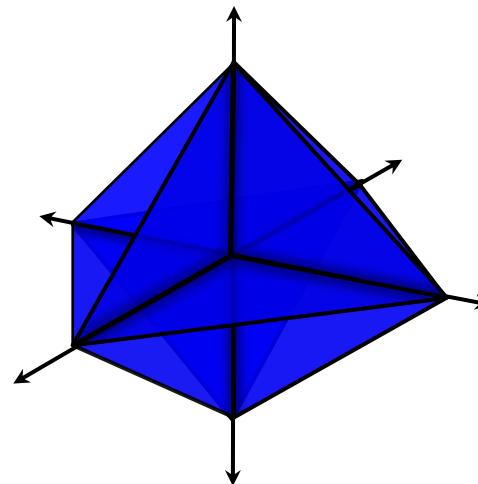
$$St_{\{\|x\|_1 \leq \lambda\}}(t) = \arg \min_{\|x\|_1 \leq \lambda} \|x - t\|$$



+

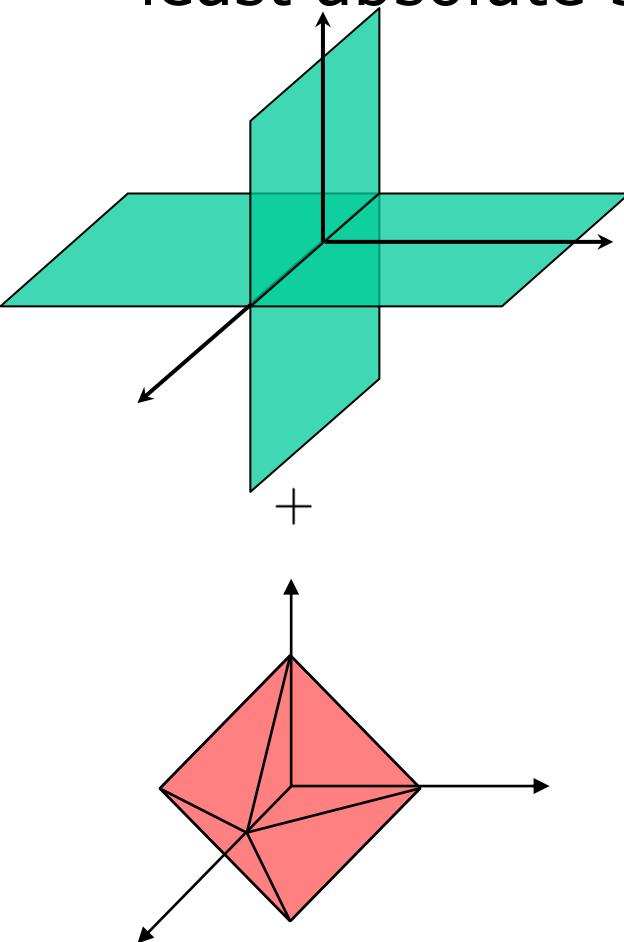


$\approx$



# The CLASH algorithm

combinatorial selection  
+  
least absolute shrinkage



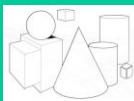
$$H_{\Sigma_{\mathcal{M}_K}}(y) = \arg \min_{x: x \in \Sigma_{\mathcal{M}_K}} \|x - y\|$$

$$St_{\{\|x\|_1 \leq \lambda\}}(t) = \arg \min_{\|x\|_1 \leq \lambda} \|x - t\|$$

**combinatorial origami**

# Recovery algorithms for low-dimensional models

## The Clash Operator

	<b>Non-convex</b> $\binom{N}{K}$	<b>Convex</b> 	<b>Probabilistic</b> 
Encoding	combinatorial / manifolds	atomic norm / convex relaxation	compressible / sparse priors
Example	$\min_{x: \ x\ _0 \leq K} \ u - \Phi x\ ^2$	$\min_{x: \ x\ _1 \leq \lambda} \ u - \Phi x\ ^2$	$E\{x u\}$
Algorithm	IHT, CoSaMP, SP, ALPS, OMP...	Basis pursuit, Lasso, basis pursuit denoising...	Variational Bayes, EP, Approximate message passing (AMP)...

$$\hat{x}_{\text{Clash}} = \arg \min_{x: \|x\|_0 \leq K, \|x\|_1 \leq \lambda} \|u - \Phi x\|^2$$

The idea is much more general

$$\hat{x}_{\text{Normed Pursuit}} = \arg \min_{x: \|x\|_0 \leq K, \|x\|_* \leq \lambda} \|u - \Phi x\|^2$$

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

# Recovery algorithms for low-dimensional models

- Using projected gradient with exact non-convex projections  
***with RIP/ERC/URC/RSC...***

- **Exact low-dimensional model**

- noise-free measurements: exact recovery
  - noisy measurements: stable recovery

- **Approximately low-dimensional model**

- recovery as good as  $K$ -model-sparse approximation

$$\frac{\|x - \hat{x}\|_{\ell_2}}{\text{recovery error}} \leq C_1 \log \left( \frac{N}{K} \right) \frac{\|x - x_{\mathcal{M}_K}\|_{\ell_1}}{K^{1/2}} + C_2 \epsilon$$

signal  $K$ -term  
model approx error

noise

# Recovery algorithms for low-dimensional models

- Using projected gradient with exact non-convex projections  
***with RIP/ERC/URC/RSC...***

- **Exact low-dimensional model**

- noise-free measurements: exact recovery
  - noisy measurements: stable recovery

- **Approximately low-dimensional model**

- recovery as good as  $K$ -model-sparse approximation

$$\frac{\|x - \hat{x}\|_{\ell_2}}{\text{recovery error}} \leq C_1 \log \left( \frac{N}{K} \right) \frac{\|x - x_{\mathcal{M}_K}\|_{\ell_1}}{K^{1/2}} + \frac{C_2 \epsilon}{\text{noise}}$$

signal  $K$ -term  
model approx error

- the bound remains qualitatively the same for other models!!!

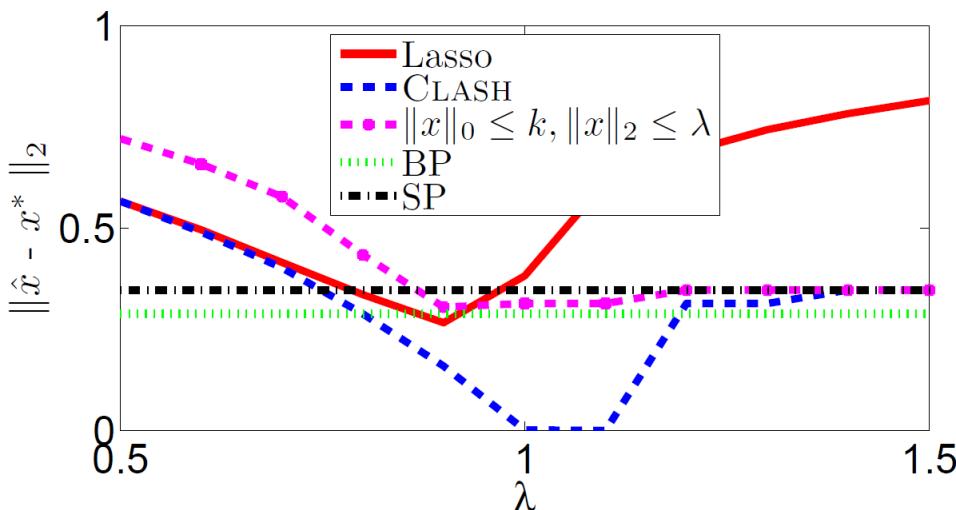
# Recovery algorithms for low-dimensional models

- Projected gradient with (non)exact non-convex projections  
*without RIP/ERC/URC/RSC...*
- **Not much!**
  - convergence to stationary point with *Kurdyka-Łojasiewicz*  
[Attouch et al., 2010]

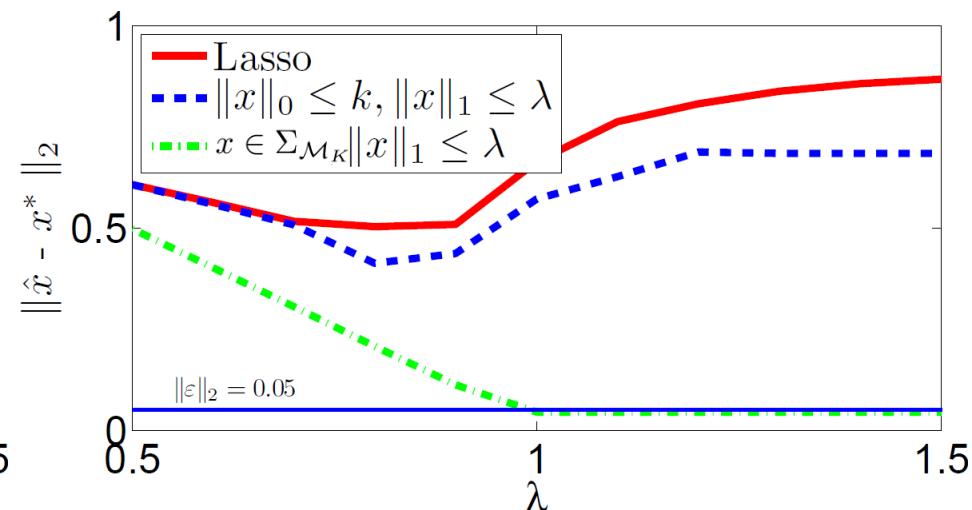
# Examples

$$\hat{x} = \arg \min_{x: \text{supp}(x) \in \Sigma_{\mathcal{M}_K}, \|x\|_* \leq \lambda} \|u - \Phi x\|^2$$

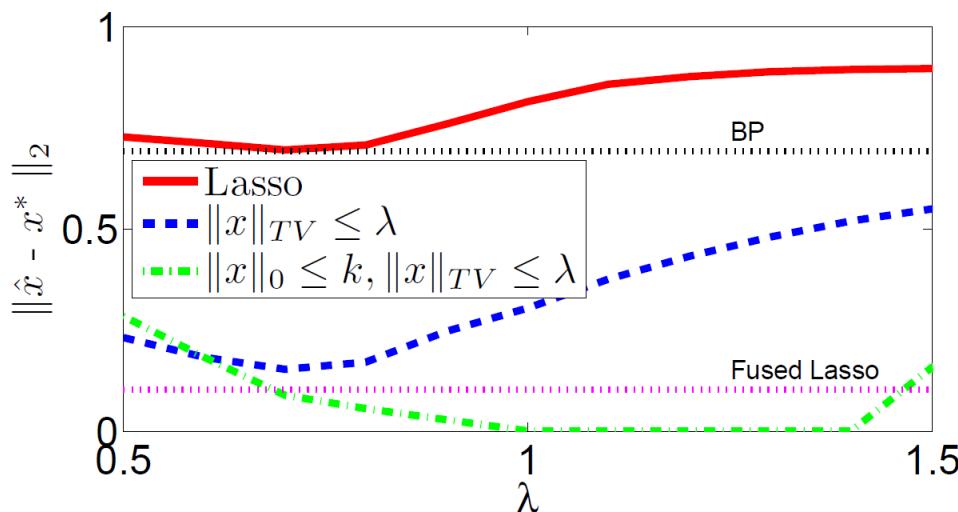
CLASH



Structured Sparsity

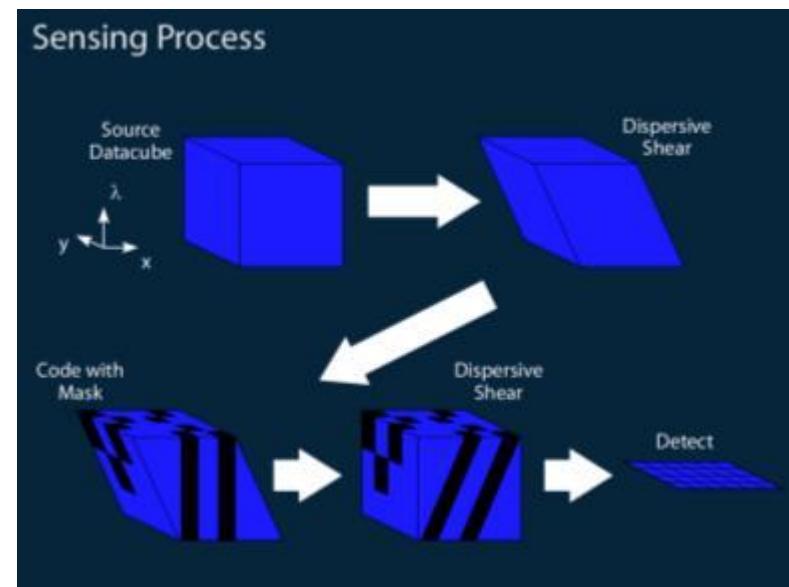
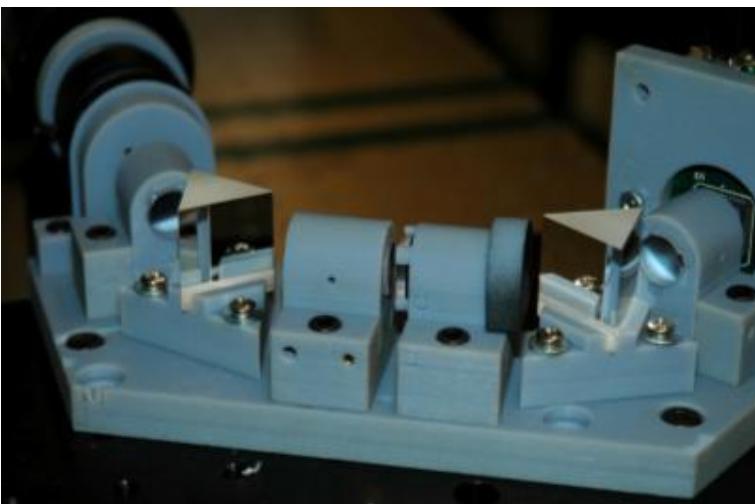
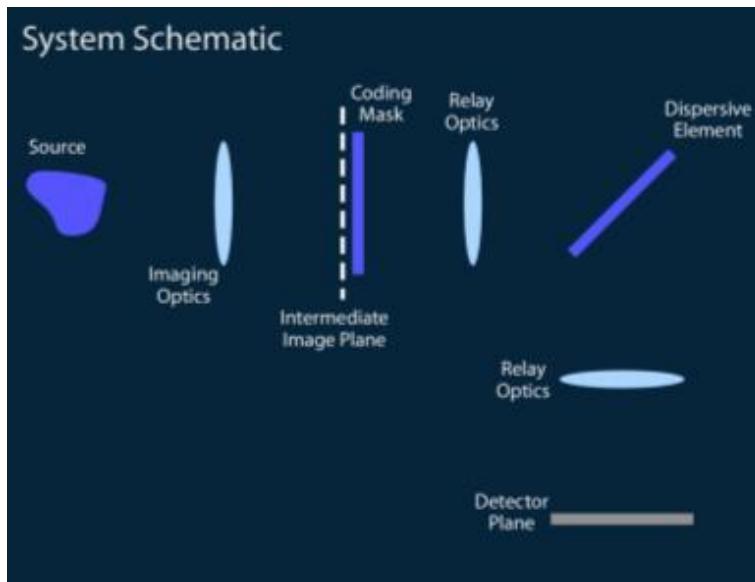


Norm Constraints



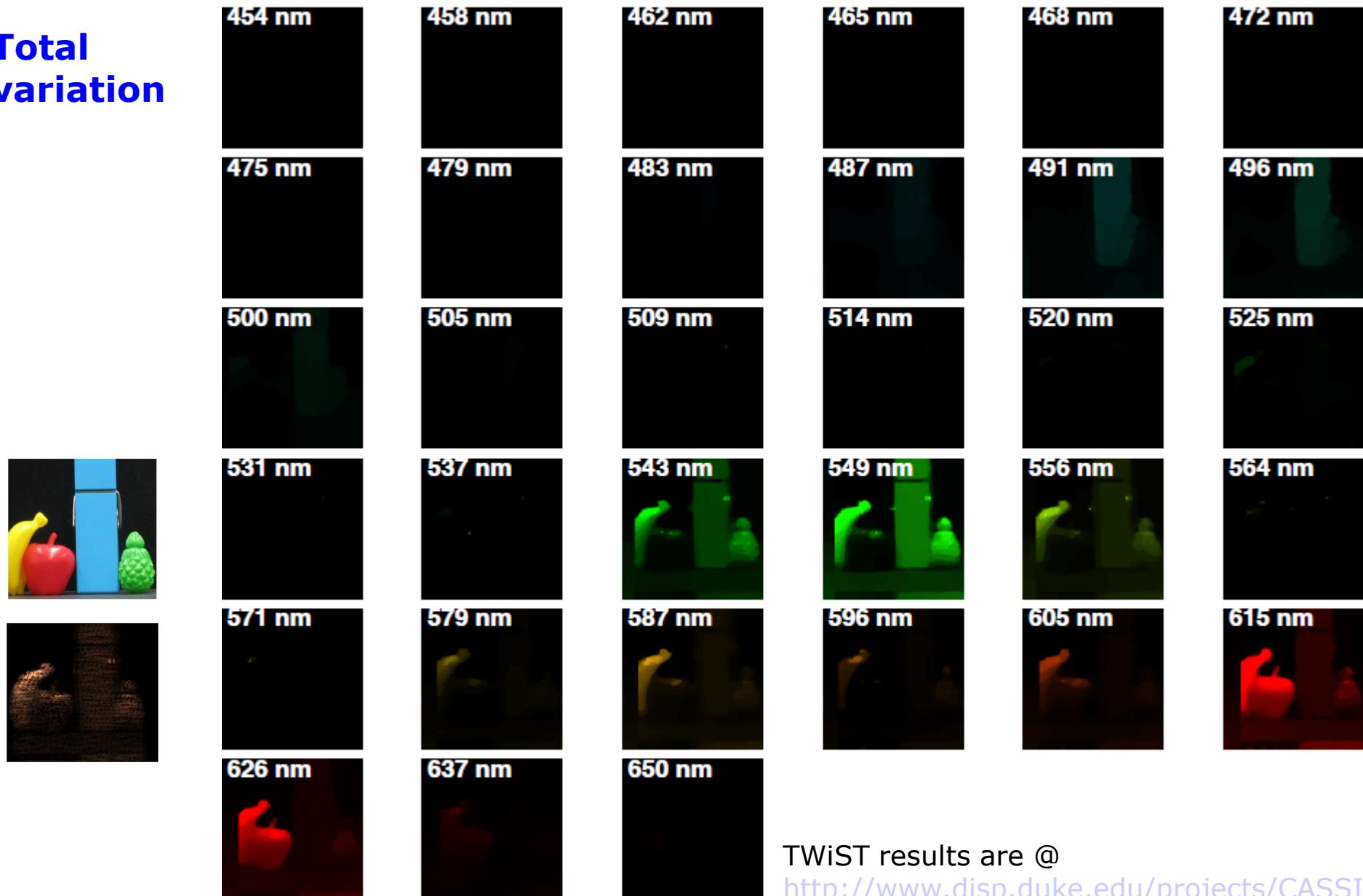
# Examples

## Coded Aperture Snapshot Spectral Imager



# Examples

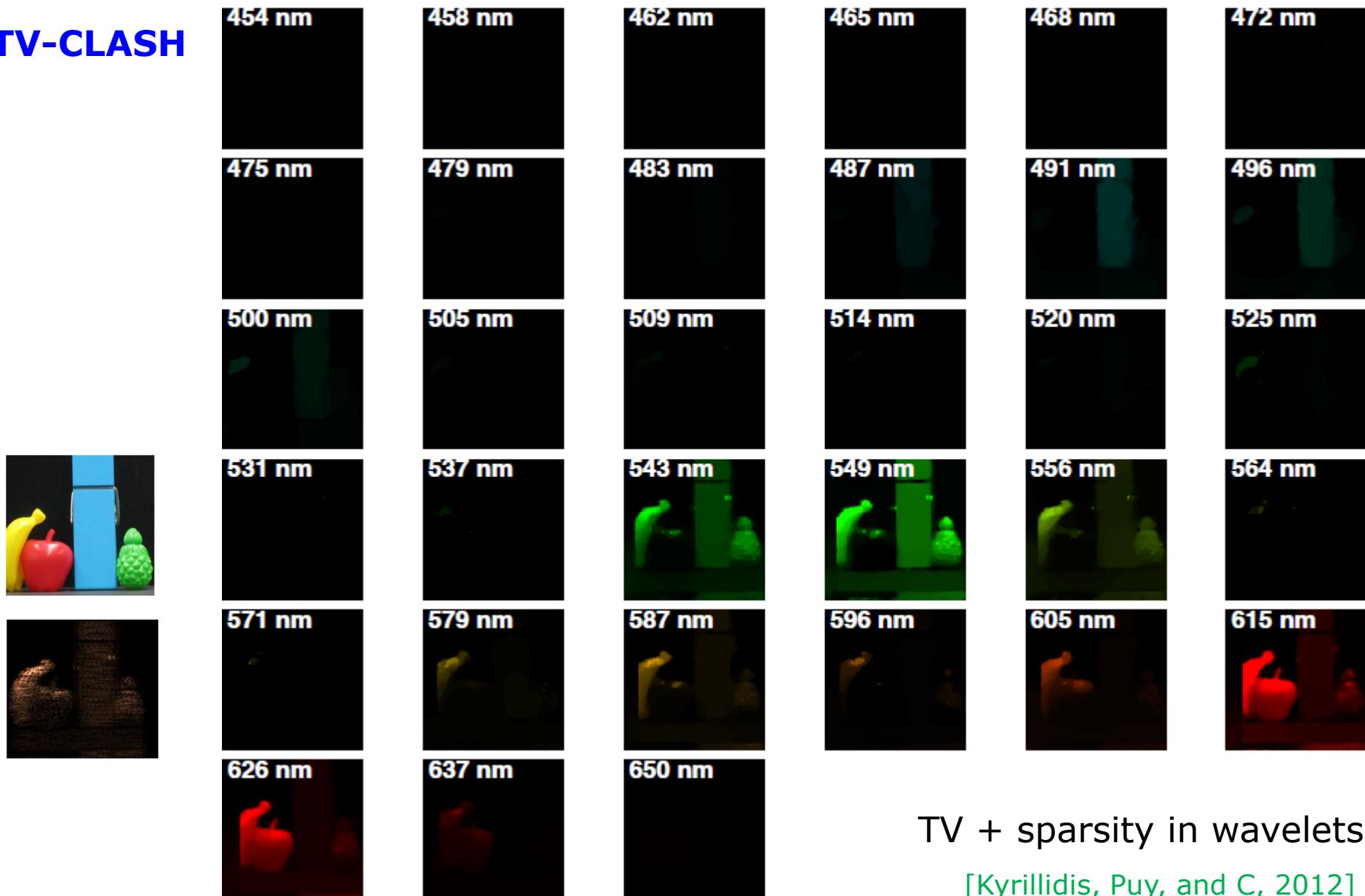
Total variation



TWiST results are @  
<http://www.disp.duke.edu/projects/CASSI/>

# Examples

TV-CLASH



TV + sparsity in wavelets  
[Kyrillidis, Puy, and C, 2012]

# Acceleration of non-convex algorithms

- Several approaches

**step-size selection**

**memory based** methods

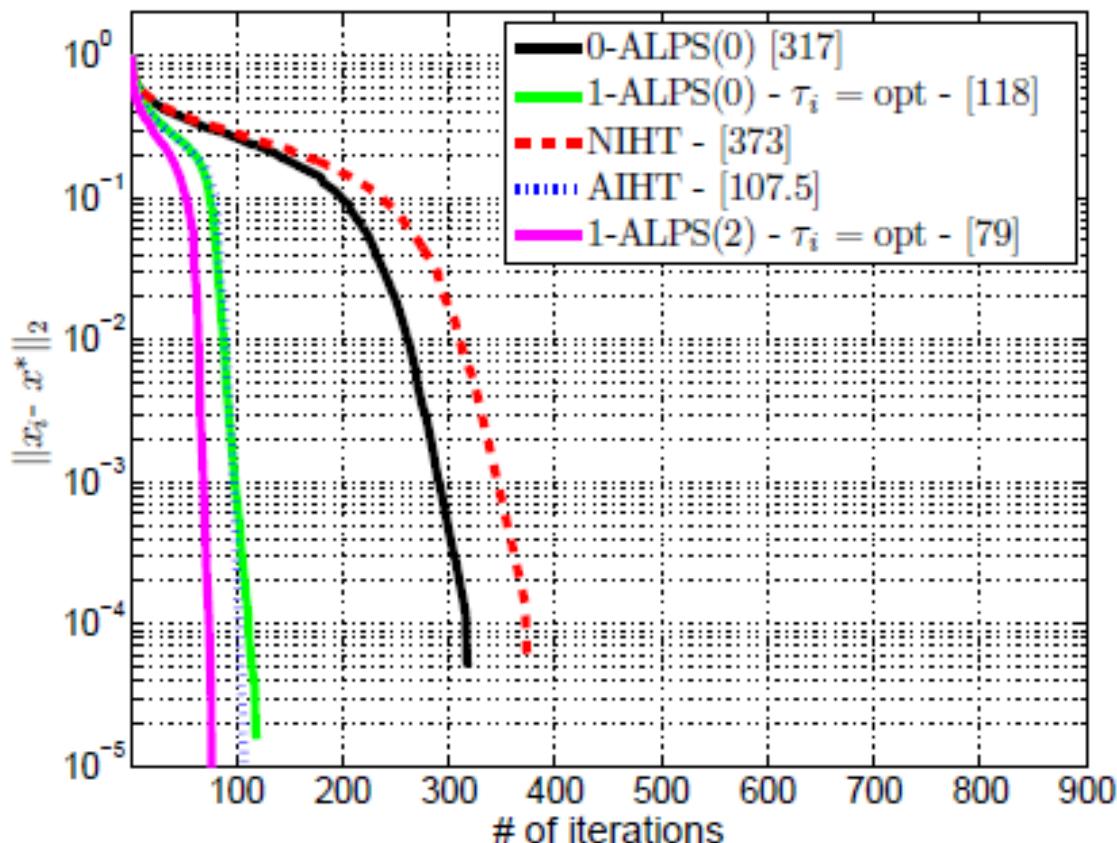
similar to Nesterov  
acceleration / double  
overrelaxation

non-convex splitting

**(adaptive) block  
coordinate descent**

epsilon-approximate  
projections

$$x_{i+1} = H_{\Sigma M_K}(y_i - \mu_i \nabla f(y_i))$$
$$y_{i+1} = x_{i+1} + \tau_i(x_{i+1} - x_i)$$



# Acceleration of non-convex algorithms

- Several approaches

step-size selection

**memory based** methods  
similar to Nesterov  
acceleration / double  
overrelaxation

34.8s

**non-convex splitting**

(adaptive) block  
coordinate descent

**epsilon-approximate  
projections**

15.8s

144 x 176 x 200

Original



Low rank



Sparse



GoDec



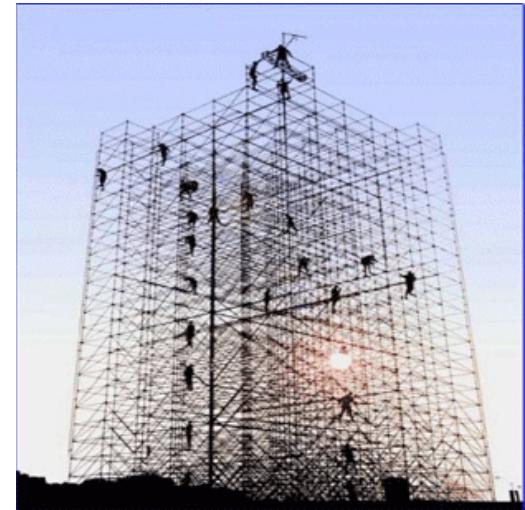
MATRIX ALPS

# Final remarks

- non-convex algorithms    <>    low-dimensional scaffold

- possible performance gains
- non-convexifiable priors
- matching prox operator with optimal space/time bounds

*complexity of structured approximation*



- non-convex algorithms                      vs.              convex algorithms

- no clear winner / scenario dependent
- decades of research in both



# References

- M. Afonso, J. Bioucas-Dias, M. Figueiredo, "Fast image recovery using variable splitting and constrained optimization", *IEEE Transactions on Image Processing*, vol. 19, 2010.
- H. Attouch, J. Bolte, P. Redont, A. Soubeyran, "Proximal alternating minimization and projection methods for nonconvex problems: an approach based on the Kurdyka-Łojasiewicz inequality", *Math. Oper. Research*, 2010
- O. Axelsson, *Iterative Solution Methods*, Cambridge University Press, 1996.
- R. Baraniuk, V. Cevher, M. Duarte, C. Hegde, "Model-based compressive sensing", *IEEE Transactions on Information Theory*, vol. 56, 2010.
- R. Baraniuk, M. Davenport, R. de Vore, M. Wakin, "A Simple Proof of the Restricted Isometry Property for Random Matrices", *Constructive Approximation*, 2008.
- R. Baraniuk, V. Cevher, M. Wakin, "Low-dimensional models for dimensionality reduction and signal recovery: A geometric perspective", *Proceedings of the IEEE*, vol. 98, 2010.
- J. Barzilai and J. Borwein, "Two point step size gradient methods," *IMA Journal of Numer. Anal.*, vol. 8, 1988.
- R. Basri, D. Jacobs, "Lambertian Reflectance and Linear Subspaces", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, 2003.
- A. Beck, M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM Journal on Imaging Science*, vol. 2, 2009.
- J. Bioucas-Dias, M. Figueiredo, "A new TwIST: Two-step iterative shrinkage/thresholding algorithms for image restoration," *IEEE Transactions on Image Processing*, vol. 16, 2007.
- P. Boufounos, R. Baraniuk, "1-bit compressed sensing", *Proceedings of the Conference on Information Science and Systems*, Princeton, 2008.
- C. Boutsidis, M. Mahoney, P. Drineas, "An improved approximation algorithm for the column subset selection problem", *Proc. 20th Annual ACM/SIAM Symposium on Discrete Algorithms*, New York, NY, 2008.

# References

- P. Bühlmann, S. van der Geer, *Statistics for High-Dimensional Data*, Springer, 2011.
- E. Candès, "The restricted isometry property and its implications for compressed sensing", *Comptes Rendus Mathematique*, vol. 346, 2008.
- E. Candès and B. Recht, "Exact matrix completion via convex optimization", *Foundations of Computational Mathematics*, vol. 9, 2009.
- L. Carin, R. Baraniuk, V. Cevher, D. Dunson, M. Jordan, G. Sapiro, M. Wakin, "Learning low-dimensional signal models", *IEEE Signal Processing Magazine*, vol. 28, 2010.
- V. Cevher, "An ALPS view of sparse recovery", *Proc. ICASSP*, 2011.
- V. Chandrasekaran, B. Recht, P. Parrilo, A. Willsky, "The convex geometry of linear inverse problems", submitted, 2010.
- R. Chartrand, W. Yin, "Iteratively reweighted algorithms for compressive sensing", *Proc. ICASSP*, 2008
- S. Chen, D. Donoho, M. Saunders, "Atomic decomposition by Basis Pursuit", *SIAM Review*, vol. 43, 2001.
- P. Combettes, V. Wajs, "Signal recovery by proximal forward-backward splitting", *SIAM Journal Multiscale Modeling and Simulation*, vol. 4, 2005.
- J. Eckstein, D. Bertsekas, "On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators", *Mathematical Programming*, vol. 5, 1992.
- M. Figueiredo and J. Bioucas-Dias, "Restoration of Poissonian images using alternating direction optimization", *IEEE Transactions on Image Processing*, vol. 19, 2010.
- M. Figueiredo, R. Nowak, S. Wright, "Gradient projection for sparse reconstruction: application to compressed sensing and other inverse problems", *IEEE Journal of Selected Topics in Signal Processing*, vol. 1, 2007.

# References

- D. Gabay, B. Mercier, "A dual algorithm for the solution of nonlinear variational problems via finite-element approximations", *Computers and Mathematics with Application*, vol. 2, 1976.
- R. Glowinski, A. Marroco, "Sur l'approximation, par éléments finis d'ordre un, et la résolution, par penalisation-dualité d'une classe de problèmes de Dirichlet non linéaires," *Rev. Française d'Automatique*, 1975.
- Y. Gordon, "On Milman's inequality and random subspaces which escape through a mesh in  $\mathbb{R}^n$ ", in *Geometric Aspects of Functional Analysis*, Springer, 1988.
- E. Hale, W. Yin, Y. Zhang, "Fixed-point continuation for l1-minimization: Methodology and convergence", *SIAM Journal on Optimization*, vol. 19, 2008.
- N. Halko, P.-G. Martinsson, J. Tropp, "Finding structure with randomness: stochastic algorithms for constructing approximate matrix decompositions", *SIAM Review*, vol. 53, 2011.
- C. Hegde, M. Duarte, V. Cevher, "Compressive sensing recovery of spike trains using a structured sparsity model", *Proceedings of SPARS'09*, Saint-Malo, France, 2009.
- M. Hestenes, "Multiplier and gradient methods", *Journal of Optimization Theory and Applications*, vol. 4, 1969.
- A. Kyrillidis, V. Cevher, "Recipes for hard thresholding methods", Tech. Rep., EPFL, 2011.
- J. Lee, V. Mirrokni, V. Nagarajan, M. Sviridenko, "Non-monotone submodular maximization under matroid and knapsack constraints", *Proc. 41st Annual ACM Symposium on Theory of Computing*, Bethesda, MD, 2009.
- A. Lewis, J. Malick, "Alternating projections on manifolds", *Math. of Operations Research*, vol. 33, 2008.
- D. Lorenz, "Constructing test instances for basis pursuit denoising", submitted, 2011.
- N. Meinshausen, P. Bühlmann, "High-dimensional graphs and variable selection with the lasso", *The Annals of Statistics*, vol. 34, pp. 1436-1462, 2006.
- J.-J. Moreau, "Proximité et dualité dans un espace hilbertien," *Bull. Soc. Mathématiques de France*, vol. 93, 1965.

# References

- G. Nemhauser, L. Wolsey, *Integer and combinatorial optimization*, Wiley, 1988.
- S. Osher, M. Burger, D. Goldfarb, J. Xu, W. Yin, "An iterative regularization method for total variation-based image restoration", *SIAM Journal on Multiscale Modeling and Simulation*, vol. 4, 2005.
- Y. Plan, "Compressed sensing, sparse approximation, low-rank matrix estimation", PhD Thesis, Caltech, 2011
- M. Powell, "A method for nonlinear constraints in minimization problems", in *Optimization*, Academic Press, 1969.
- H. Raguet, J. Fadili, G. Peyré, "Generalized Forward-Backward splitting", Tech. report, Hal-00613637, 2011.
- S. Setzer, G. Steidl, T. Teuber, "Deblurring Poissonian images by split Bregman techniques," *Journal of Visual Communication and Image Representation*, 2010.
- K.-C. Toh , S. Yun, "An accelerated proximal gradient algorithm for nuclear norm regularized least squares Problems", *Pacific Journal of Optimization*, vol. 6, 2010.
- A. Waters, A. Sankaranarayanan, R. Baraniuk, "SpaRCS: recovering low-rank and sparse matrices from compressive measurements", *Neural Information Processing Systems*, 2011.
- S. Wright, R. Nowak, M. Figueiredo, "Sparse reconstruction by separable approximation," *IEEE Transactions on Signal Processing*, vol. 57, 2009.
- W. Yin, S. Osher, D. Goldfarb, J. Darbon, "Bregman iterative algorithms for l1-minimization with applications to compressed sensing", *SIAM Journal on Imaging Science*, vol. 1, 2008.
- T. Zhou, D. Tao, "Godec: randomized low-rank & sparse matrix decomposition in noisy case," *Proc. International Conference on Machine Learning*, Bellevue, WA, 2011.