

Theory of Sparse and Low-Rank Recovery

John Wright

Electrical Engineering
Columbia University

UNDERDETERMINED LINEAR SYSTEMS

Observation $y \in \mathbb{R}^m$

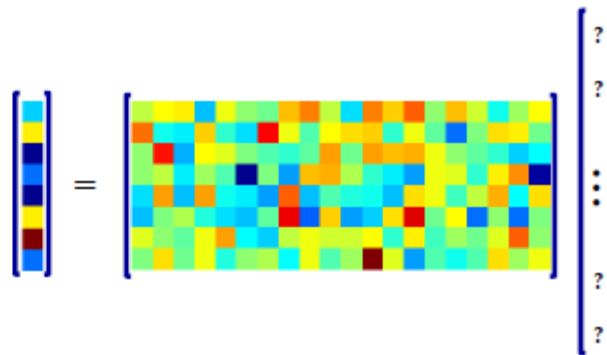
$A \in \mathbb{R}^{m \times n}$

Unknown $x \in \mathbb{R}^n$

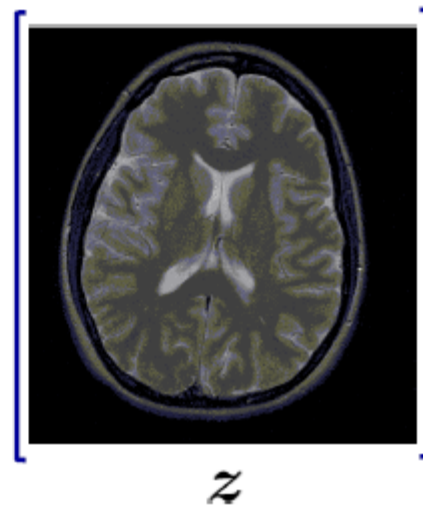
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$



Signal acquisition



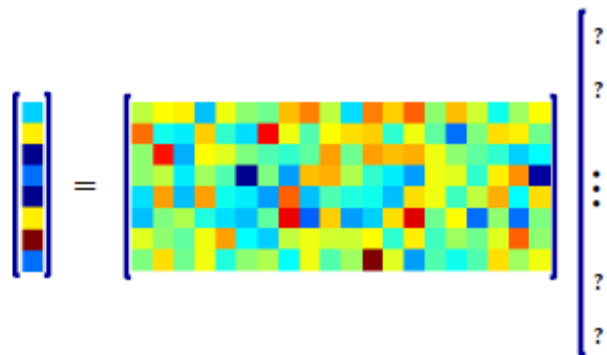
z

Image to be sensed

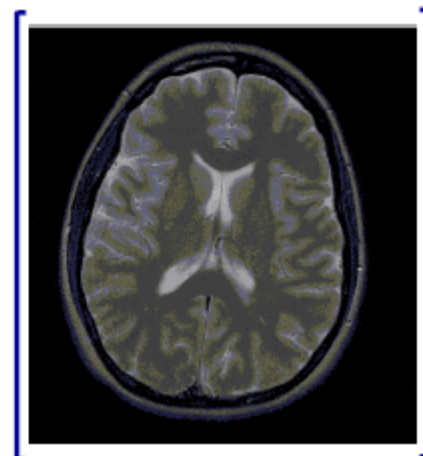
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$



Signal acquisition



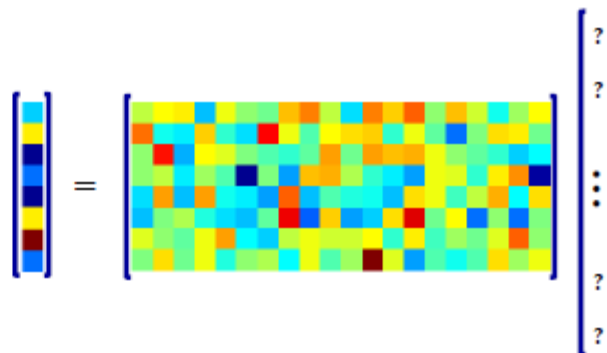
z

Image to be sensed

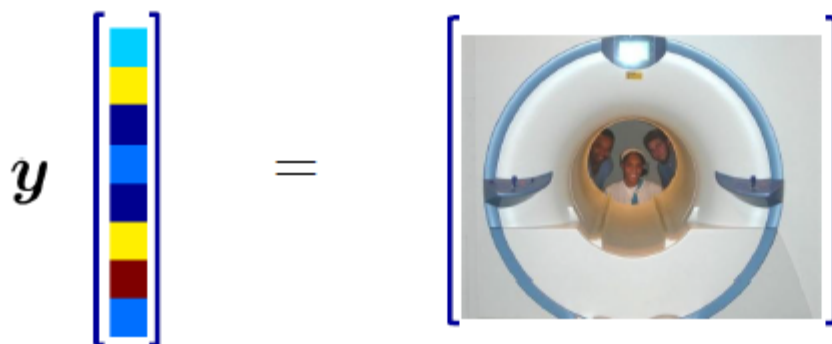
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$

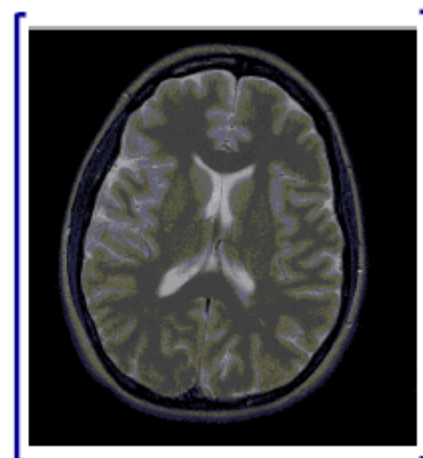


Signal acquisition



$$y_i = \int_{\mathbf{u}} z(\mathbf{u}) \exp(-2\pi j \mathbf{k}(t_i)^* \mathbf{u}) d\mathbf{u}$$

Observations are Fourier coefficients!



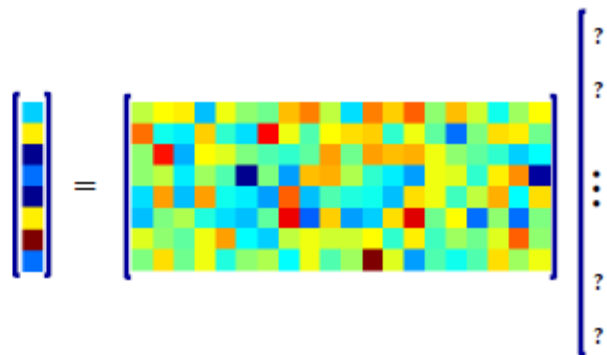
z

Image to be sensed

UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$



Signal acquisition

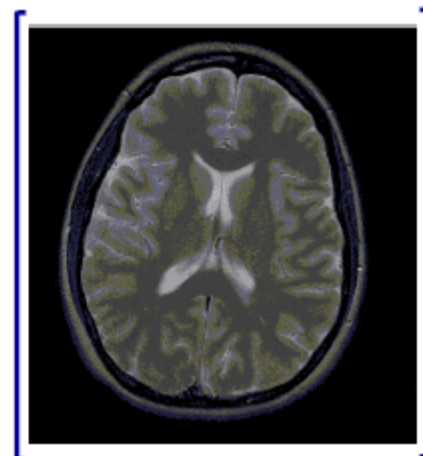


A few Fourier
coefficients

=



F_{Ω}



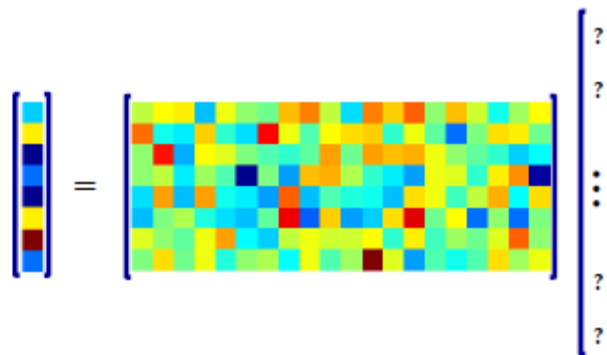
z

Image to be sensed

UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$



Signal acquisition

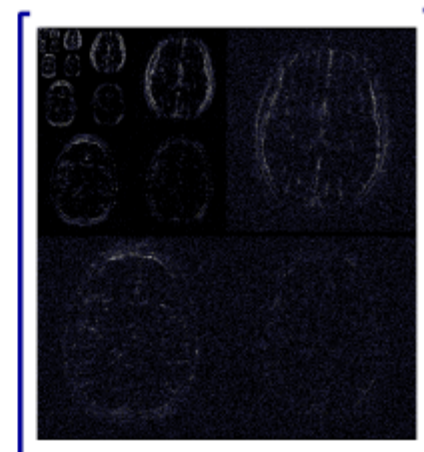


=



F_{Ω}

Ψ



x

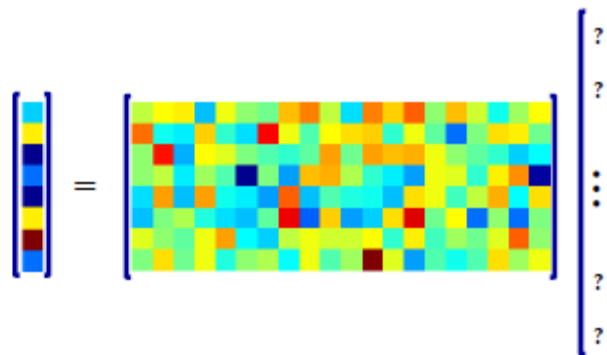
A few Fourier coefficients

Wavelet coefficients: $z = \Psi x$

UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$

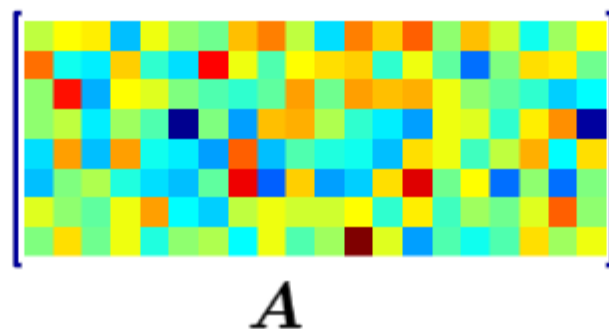


Signal acquisition

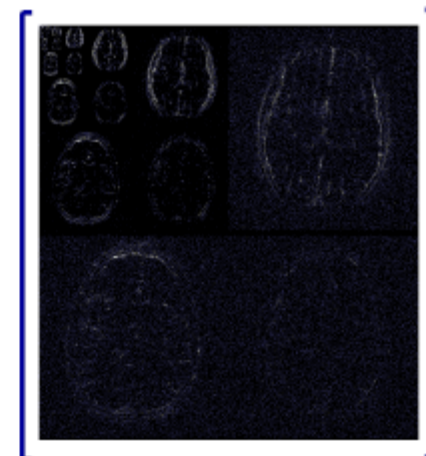


A few Fourier
coefficients

=



A



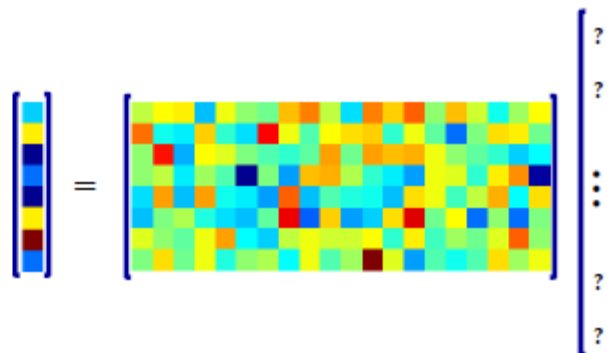
x

Wavelet coefficients

UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$

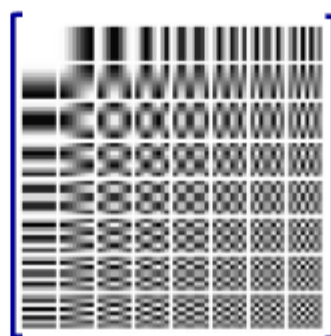


Compression – JPEG



(Patches of) ...
input image

\approx



A DCT basis

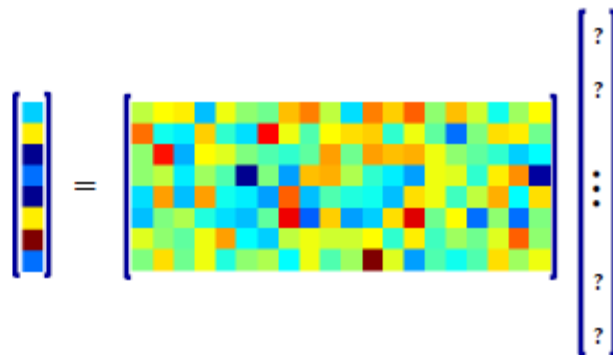


x coefficients

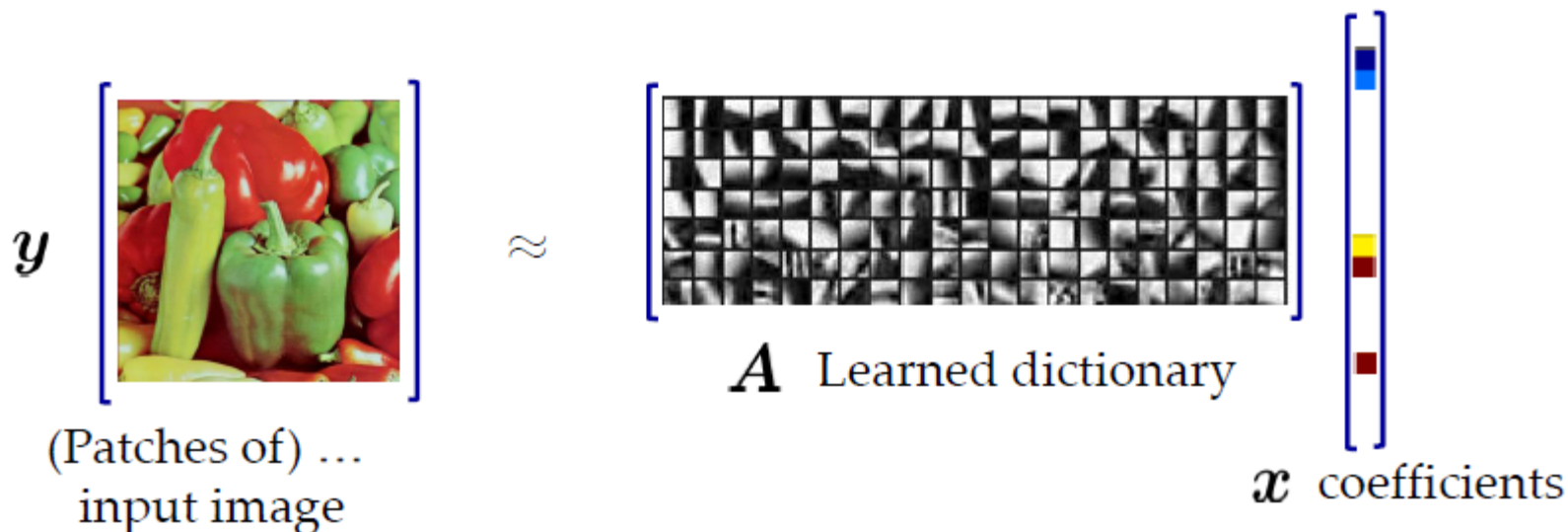
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$



Compression – Learned dictionary

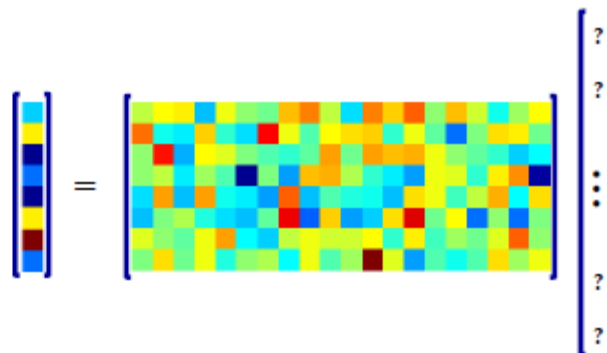


See [Elad+Bryt '08], [Horev et. Al., '12] ... Image: [Aharon+Elad '05]

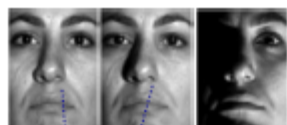
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

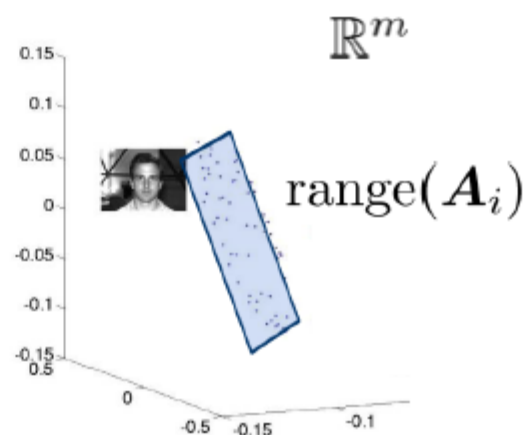
$$y = Ax$$



Recognition



$$A_i = \begin{bmatrix} | & | & \dots \\ | & | & \dots \\ | & | & \dots \\ | & | & \dots \end{bmatrix} \in \mathbb{R}^{m \times n_i}$$

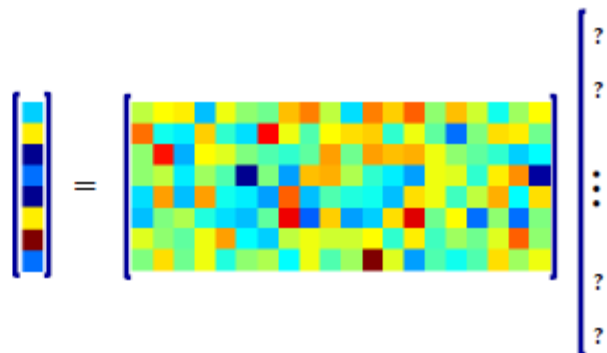


Linear subspace model for images of same face under varying lighting.

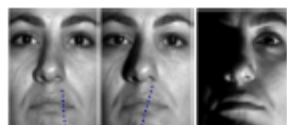
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

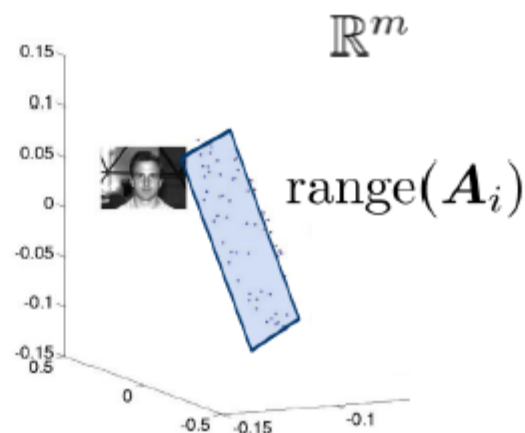
$$y = Ax$$



Recognition



$$A_i = [\begin{array}{|c|} \hline \bullet \\ \hline \bullet \\ \hline \bullet \\ \hline \bullet \\ \hline \end{array} \quad \begin{array}{|c|} \hline \bullet \\ \hline \bullet \\ \hline \bullet \\ \hline \bullet \\ \hline \end{array} \quad \dots] \in \mathbb{R}^{m \times n_i}$$

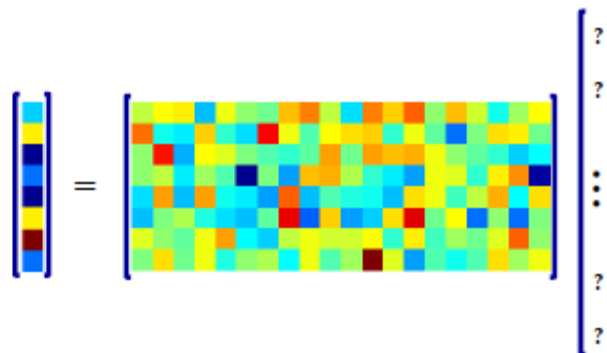


$$y \text{ (face image)} \approx x_{i,1} \text{ (face image)} + x_{i,2} \text{ (face image)} + \dots + x_{i,n} \text{ (face image)} = A_i x_i$$

UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$

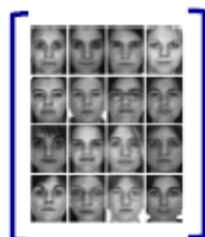


Recognition



$y \in \mathbb{R}^m$
Test image

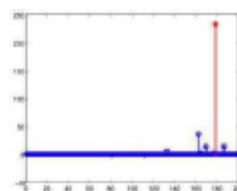
=



$$A = [A_1 \mid A_2 \mid \cdots \mid A_k]$$

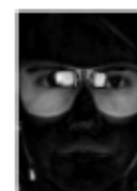
Combined
training
dictionary

×



$x \in \mathbb{R}^n$
coefficients

+



$e \in \mathbb{R}^m$
corruption,
occlusion

UNDERDETERMINED LINEAR SYSTEMS

Observation $y \in \mathbb{R}^m$

$A \in \mathbb{R}^{m \times n}$

Unknown $x \in \mathbb{R}^n$

In all of these examples,

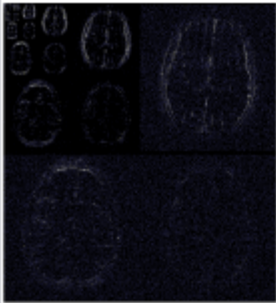



$$\underbrace{m}_{\text{\#observations}} \ll \underbrace{n}_{\text{\#unknowns}}$$

Solution is **not unique** ... is there any hope?

WHAT DO WE KNOW ABOUT x ?

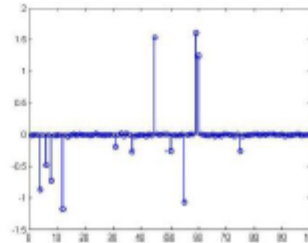
Underdetermined system

$$y = Ax$$

| Signal acquisition | Image compression | Face Recognition |
|---|---|--|
|  <p>x^* contains just a few significant wavelet coefficients.</p> |  <p>x^* uses just a few dictionary elements.</p> |  <p>x^* uses just a few training faces.</p>  <p>e^* corrects a few gross errors.</p> |

SPARSITY – More formally

A vector $\mathbf{x} \in \mathbb{R}^n$ is **sparse** if
only a few entries are nonzero:

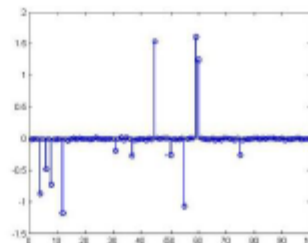


The **number of nonzeros** is called the ℓ^0 -“norm” of \mathbf{x} :

$$\|\mathbf{x}\|_0 \doteq \#\{i \mid x_i \neq 0\}.$$

SPARSITY – More formally

A vector $\mathbf{x} \in \mathbb{R}^n$ is **sparse** if only a few entries are nonzero:



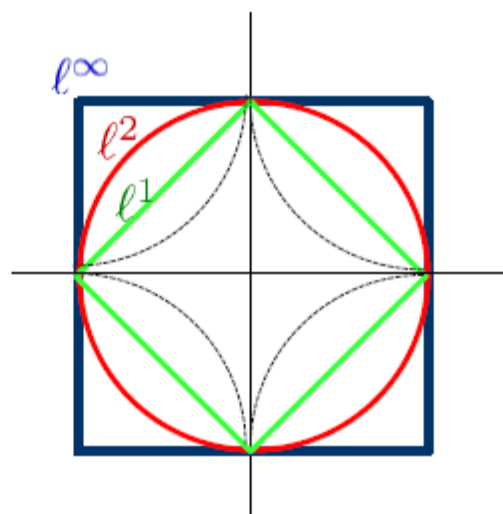
The **number of nonzeros** is called the ℓ^0 -“norm” of \mathbf{x} :

$$\|\mathbf{x}\|_0 \doteq \#\{i \mid x_i \neq 0\}.$$

Geometrically

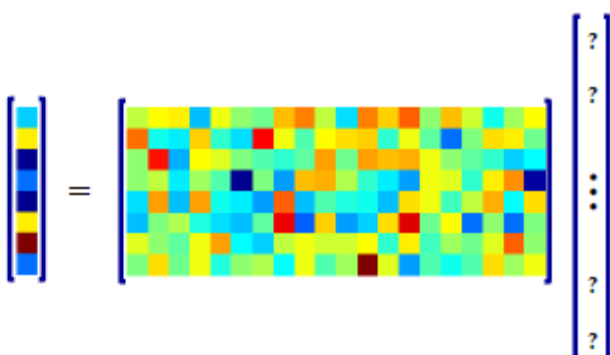
$$\|\mathbf{x}\|_p = (\sum_i |x_i|^p)^{1/p}$$

$$\|\mathbf{x}\|_0 = \lim_{p \searrow 0} \|\mathbf{x}\|_p^p.$$



THE SPARSEST SOLUTION

Underdetermined system

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$


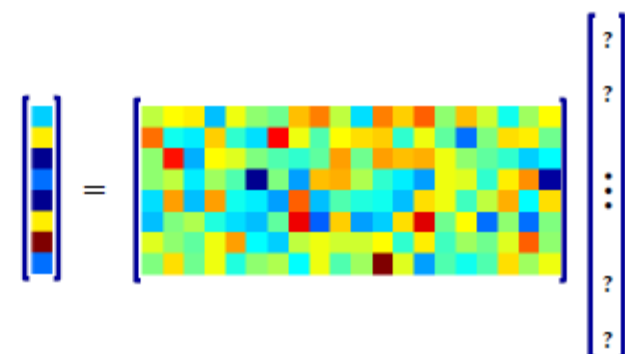
Look for the sparsest \mathbf{x} that agrees with our observation:

$$\text{minimize } \|\mathbf{x}\|_0 \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{y}.$$

[Demo]

THE SPARSEST SOLUTION

Underdetermined system

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$


Look for the sparsest \mathbf{x} that agrees with our observation:

$$\text{minimize } \|\mathbf{x}\|_0 \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{y}.$$

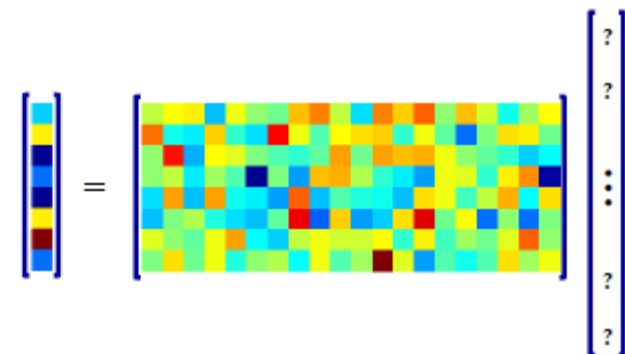
Theorem 1 (Gorodnitsky+Rao '97) .

Suppose $\mathbf{y} = \mathbf{A}\mathbf{x}_0$, and let $k = \|\mathbf{x}_0\|_0$. If $\text{null}(\mathbf{A})$ contains no $2k$ -sparse vectors, \mathbf{x}_0 is the unique optimal solution to

$$\text{minimize } \|\mathbf{x}\|_0 \quad \text{subject to } \mathbf{y} = \mathbf{A}\mathbf{x}.$$

THE SPARSEST SOLUTION

Underdetermined system

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$


Look for the sparsest \mathbf{x} that agrees with our observation:

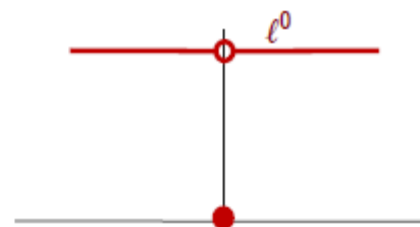
~~minimize $\|\mathbf{x}\|_0$ subject to $\mathbf{A}\mathbf{x} = \mathbf{y}$.~~

INTRACTABLE

RELAX!

~~minimize $\|\mathbf{x}\|_0$ subject to $\mathbf{Ax} = \mathbf{y}$.~~

The cardinality $\|\mathbf{x}\|_0$ is **nonconvex**:



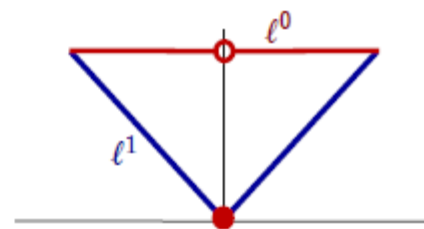
RELAX!

~~minimize $\|\mathbf{x}\|_0$ subject to $\mathbf{Ax} = \mathbf{y}$.~~

The cardinality $\|\mathbf{x}\|_0$ is **nonconvex**:

Its **convex envelope*** is

the ℓ^1 norm: $\|\mathbf{x}\|_1 = \sum_i |x_i|$



* Over the set $\{\mathbf{x} \mid |x_i| \leq 1 \forall i\}$

RELAX!

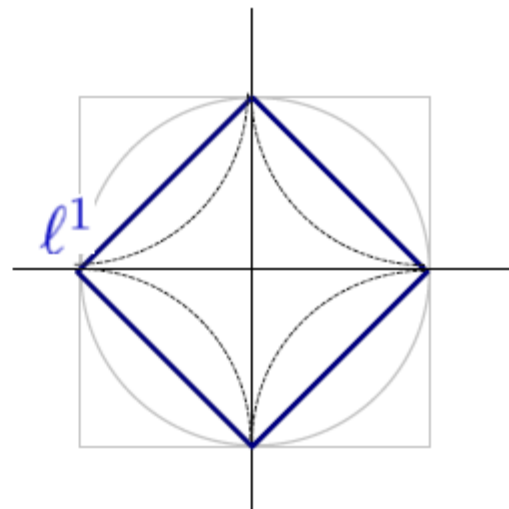
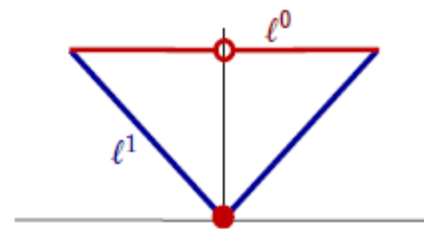
~~minimize $\|\mathbf{x}\|_0$ subject to $\mathbf{Ax} = \mathbf{y}$.~~

The cardinality $\|\mathbf{x}\|_0$ is **nonconvex**:

Its **convex envelope*** is

the ℓ^1 norm: $\|\mathbf{x}\|_1 = \sum_i |x_i|$

* Over the set $\{\mathbf{x} \mid |x_i| \leq 1 \forall i\}$



RELAX!

minimize $\|\mathbf{x}\|_0$ subject to $\mathbf{Ax} = \mathbf{y}$.

NP-hard, hard to appx.
[Natarjan '95],
[Amaldi+Kann '97]



minimize $\|\mathbf{x}\|_1$ subject to $\mathbf{Ax} = \mathbf{y}$.

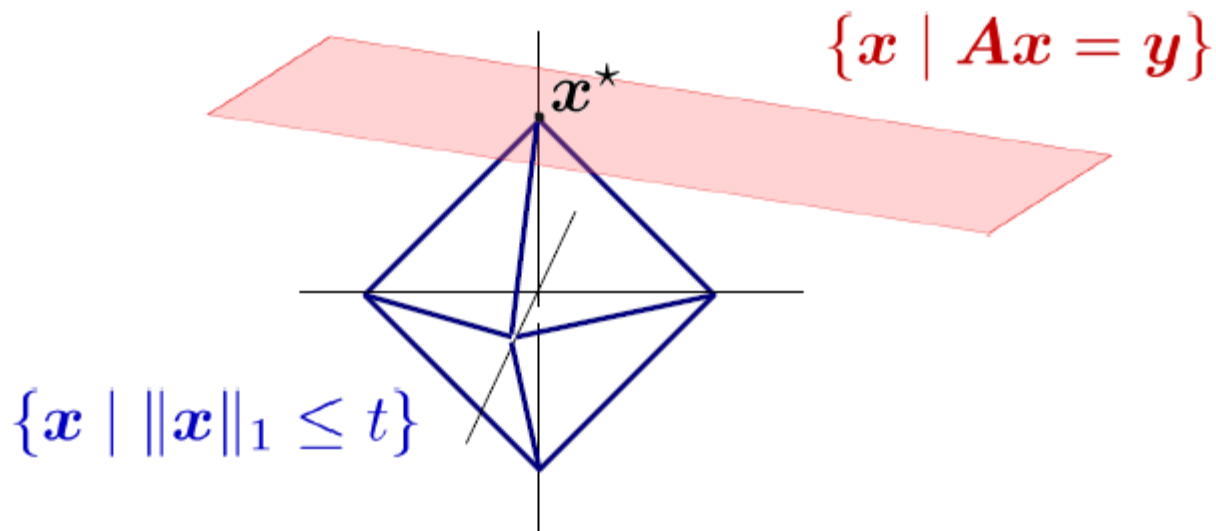
Efficiently solvable

Have we lost anything? [demo]

WHY DOES THIS WORK? Geometric intuition

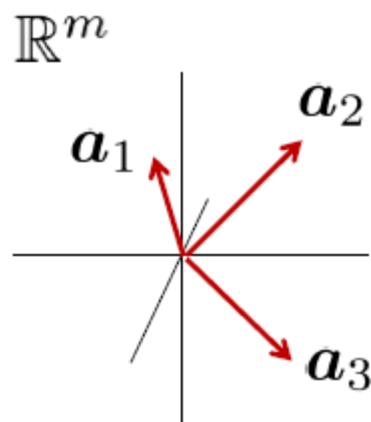
minimize $\|\mathbf{x}\|_1$ subject to $\mathbf{Ax} = \mathbf{y}$.

\mathbb{R}^n



WHY DOES THIS WORK? More formally...

We see: $\mathbf{y} = \mathbf{A}\mathbf{x} = \sum_{i \in \text{supp}(\mathbf{x})} \mathbf{a}_i x_i$



Intuition: Recovering \mathbf{x} is “easier” if the \mathbf{a}_i are not too similar...

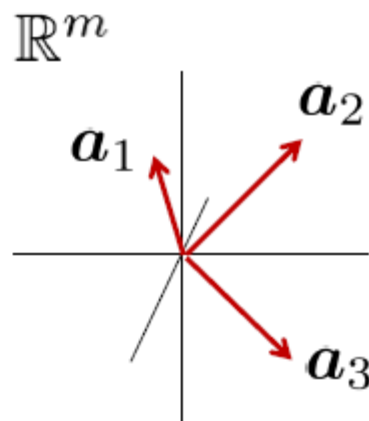
Mutual coherence $\mu(\mathbf{A}) \doteq \max_{i \neq j} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|$

Smaller is better!

WHY DOES THIS WORK? More formally...

Mutual coherence

$$\mu(\mathbf{A}) \doteq \max_{i \neq j} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|$$



Theorem 2 (Gribonval+Nielsen '03, Donoho+Elad '03) .

Suppose $\mathbf{y} = \mathbf{A}\mathbf{x}_0$ with

$$\|\mathbf{x}_0\|_0 < \frac{1}{2}(1 + 1/\mu(\mathbf{A})).$$

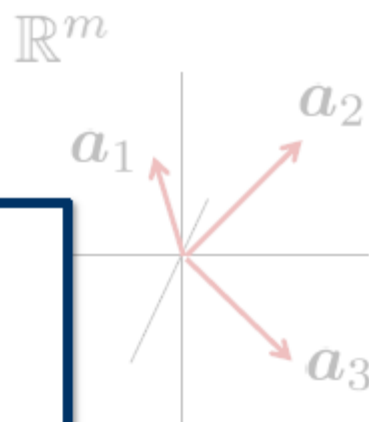
Then \mathbf{x}_0 is the unique optimal solution to

$$\text{minimize } \|\mathbf{x}\|_1 \quad \text{subject to } \mathbf{y} = \mathbf{A}\mathbf{x}.$$

WHY DOES THIS WORK? More formally...

Mutual coherence

The target solution \mathbf{x}_0 is
sufficiently structured (sparse!).



Theorem 2 (Gribonval+Nielsen '03, Donoho+Elad '03) .
Suppose $\mathbf{y} = \mathbf{A}\mathbf{x}_0$ with

$$\|\mathbf{x}_0\|_0 < \frac{1}{2}(1 + 1/\mu(\mathbf{A})).$$

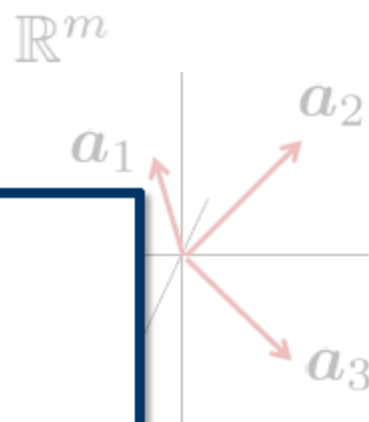
Then \mathbf{x}_0 is the unique optimal solution to

$$\text{minimize } \|\mathbf{x}\|_1 \quad \text{subject to } \mathbf{y} = \mathbf{A}\mathbf{x}.$$

WHY DOES THIS WORK? More formally...

Mutual coherence

The matrix \mathbf{A} is **incoherent** – and so, preserves sparse \mathbf{x} .



Theorem 2 (Gribonval+Nielsen '03, Donoho+Elad '03) .

Suppose $\mathbf{y} = \mathbf{A}\mathbf{x}_0$ with

$$\|\mathbf{x}_0\|_0 < \frac{1}{2}(1 + 1/\mu(\mathbf{A}))$$

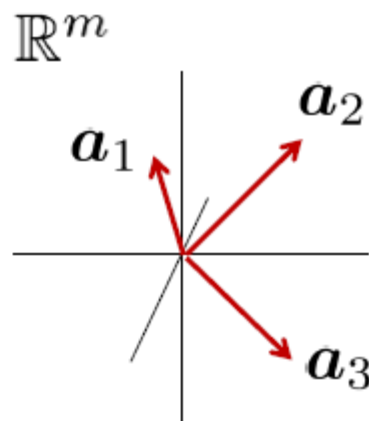
Then \mathbf{x}_0 is the unique optimal solution to

$$\text{minimize } \|\mathbf{x}\|_1 \quad \text{subject to } \mathbf{y} = \mathbf{A}\mathbf{x}.$$

WHY DOES THIS WORK? More formally...

Mutual coherence

$$\mu(\mathbf{A}) \doteq \max_{i \neq j} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|$$



Theorem 2 (Gribonval+Nielsen '03, Donoho+Elad '03) .

Suppose $\mathbf{y} = \mathbf{A}\mathbf{x}_0$ with

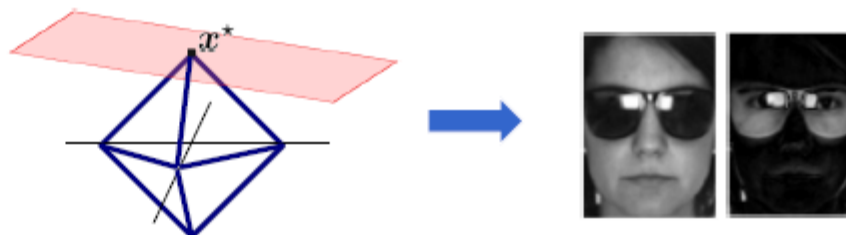
$$\|\mathbf{x}_0\|_0 < \frac{1}{2}(1 + 1/\mu(\mathbf{A})).$$

Then \mathbf{x}_0 is the unique optimal solution to

$$\text{minimize } \|\mathbf{x}\|_1 \quad \text{subject to } \mathbf{y} = \mathbf{A}\mathbf{x}.$$

WHY CARE ABOUT THE THEORY?

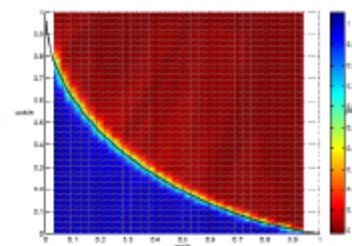
Motivates applications



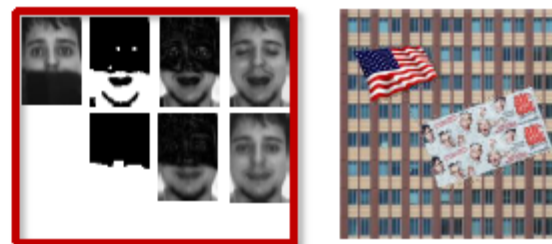
... but be careful: need to justify (and modify) the basic models

Template for stronger results

... predictions can be very sharp in high dimensions.

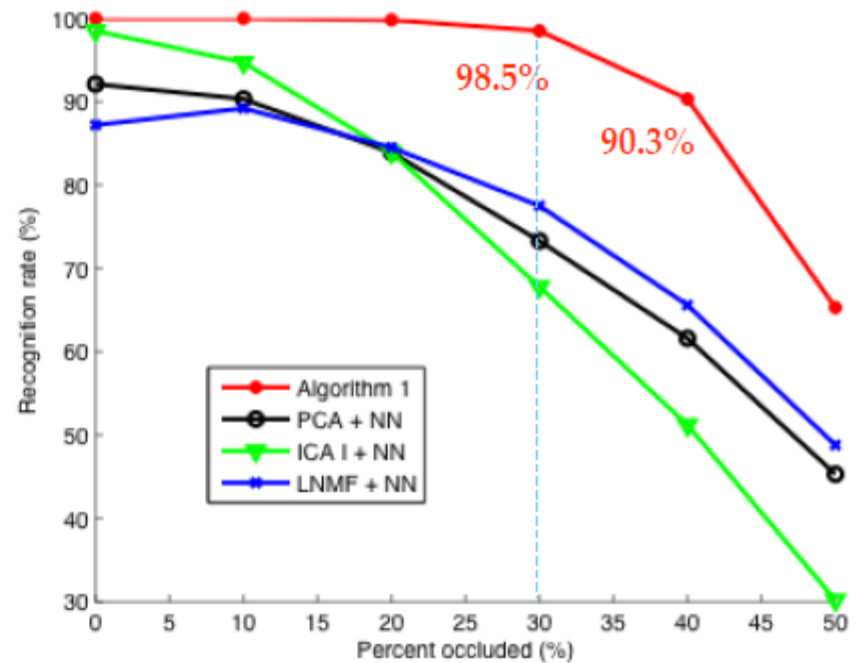
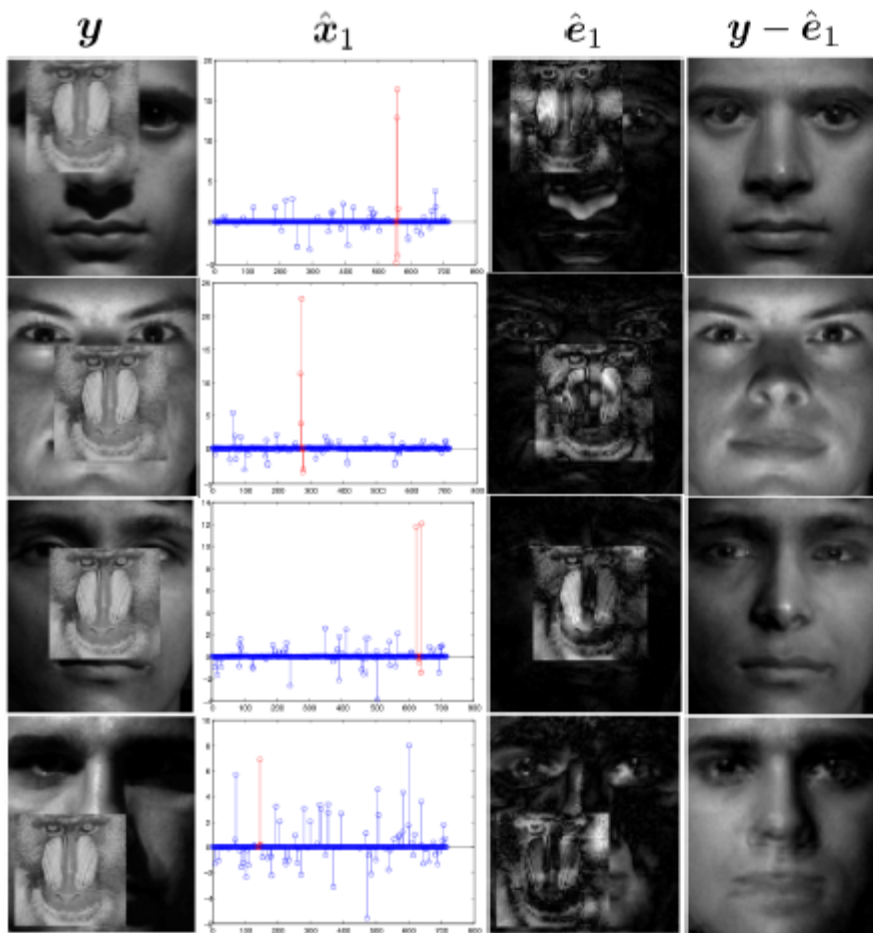


Generalizes to many other types
of low-dimensional structure

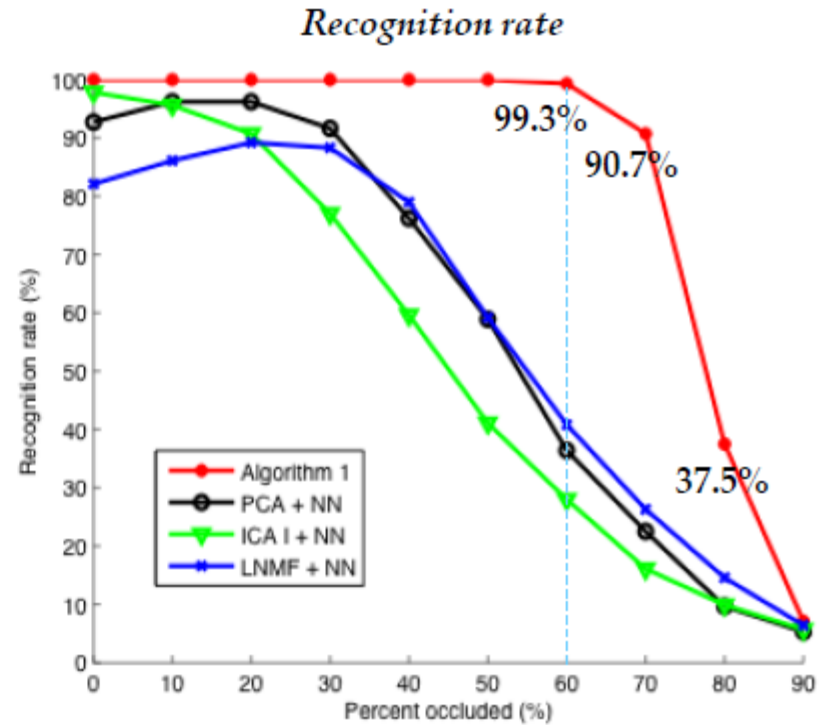
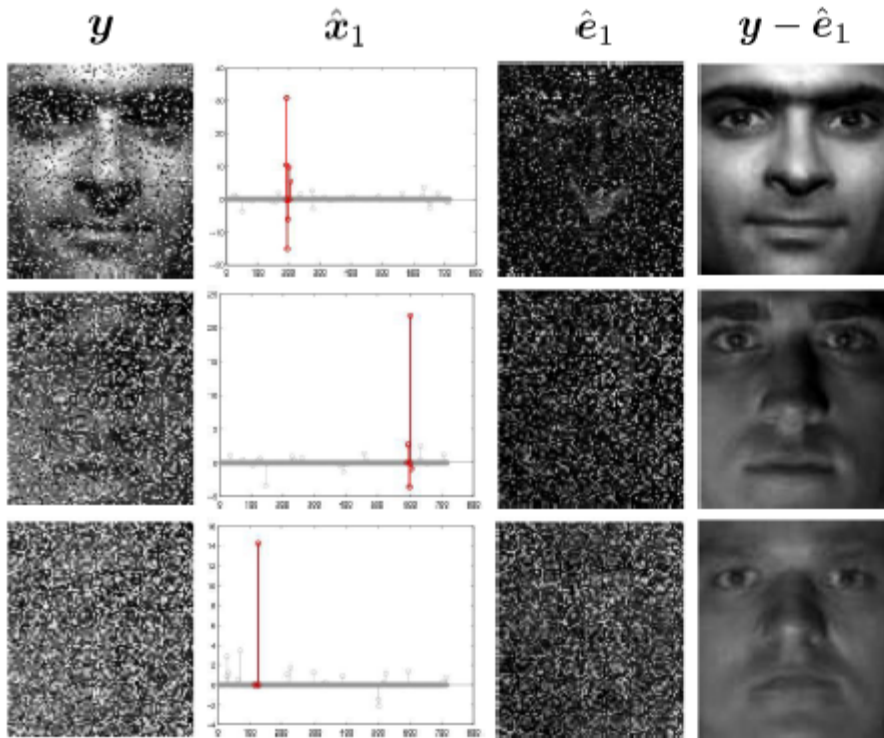


... structured sparsity, low-rank recovery

MOTIVATING APPLICATIONS – Face Recognition

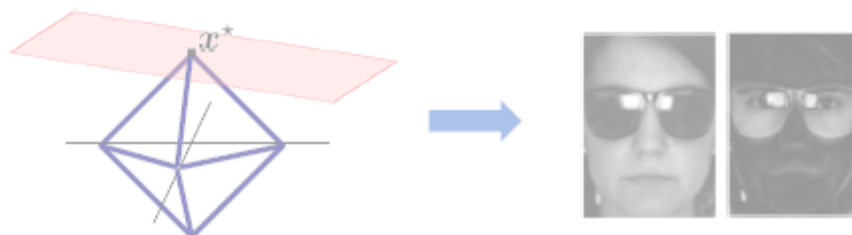


MOTIVATING APPLICATIONS – Face Recognition



WHY CARE ABOUT THE THEORY?

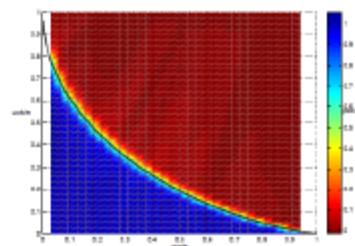
Motivates applications



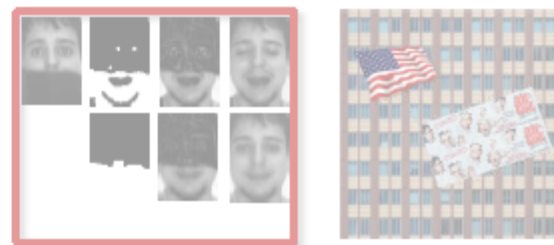
... but be careful: need to justify (and modify) the basic models [Lecture 2].

Template for stronger results

... predictions can be very sharp in high dimensions.



Generalizes to many other types
of low-dimensional structure



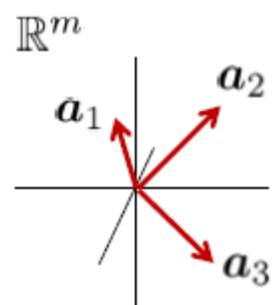
... structured sparsity, low-rank recovery

LIMITATIONS OF COHERENCE?

For any $m \times n$ \mathbf{A} $\mu(\mathbf{A}) \geq \sqrt{\frac{n-m}{m(n-1)}}$

Prev. result therefore requires

$$\|\mathbf{x}_0\|_0 < \frac{1}{2}(1 + \mu(\mathbf{A})^{-1}) = O(\sqrt{m})$$

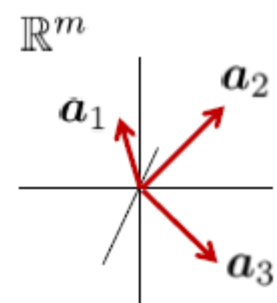


LIMITATIONS OF COHERENCE?

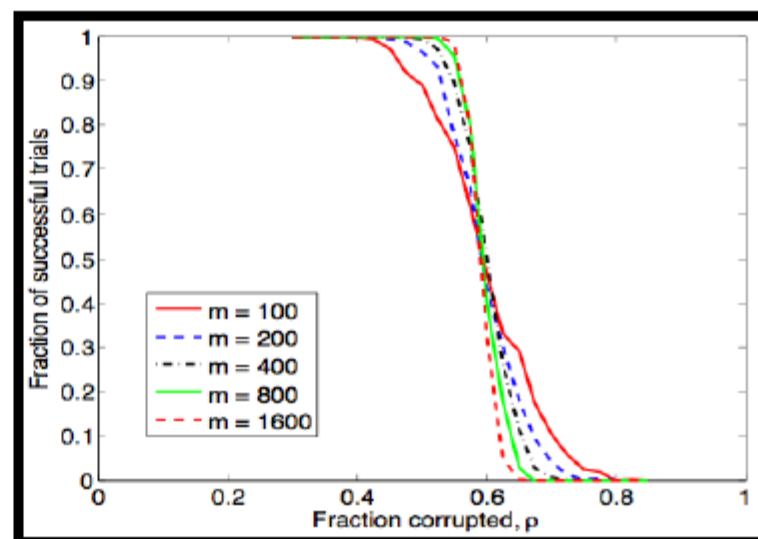
For any $m \times n$ \mathbf{A} $\mu(\mathbf{A}) \geq \sqrt{\frac{n-m}{m(n-1)}}$

Prev. result therefore requires

$$\|\mathbf{x}_0\|_0 < \frac{1}{2}(1 + \mu(\mathbf{A})^{-1}) = O(\sqrt{m})$$



Truth is often **much better**:



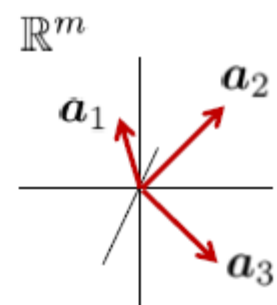
Plot: Fraction of correct recovery
vs. fraction of nonzeros $\|\mathbf{x}_0\|_0/m$

LIMITATIONS OF COHERENCE?

For any $m \times n$ \mathbf{A} $\mu(\mathbf{A}) \geq \sqrt{\frac{n-m}{m(n-1)}}$

Prev. result therefore requires

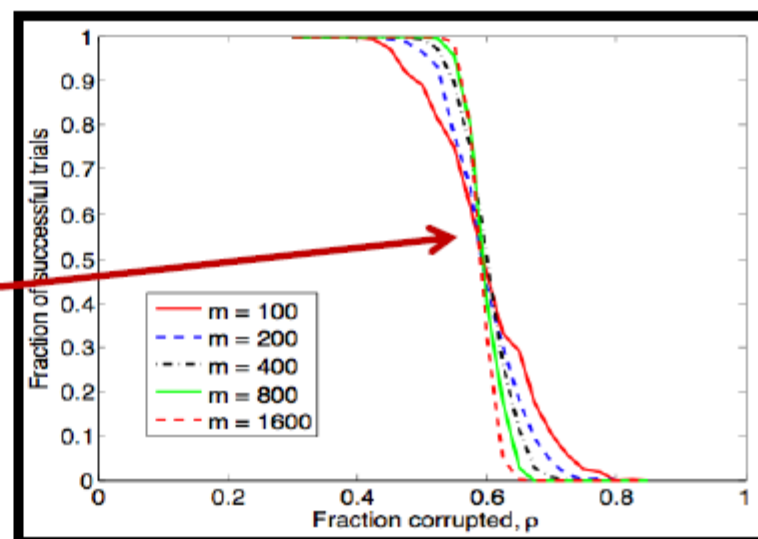
$$\|\mathbf{x}_0\|_0 < \frac{1}{2}(1 + \mu(\mathbf{A})^{-1}) = O(\sqrt{m})$$



Truth is often much better:

Phase transition at

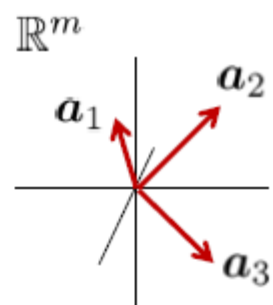
$$\|\mathbf{x}_0\|_0 = \alpha^* m$$



Plot: Fraction of correct recovery vs. fraction of nonzeros $\|\mathbf{x}_0\|_0/m$

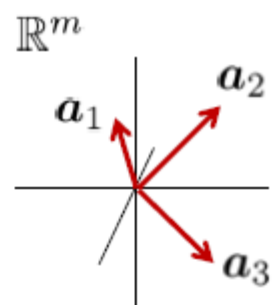
STRENGTHENING THE BOUND – the RIP

Incoherence: Each pair $\mathbf{A}_{i,j} = [\mathbf{a}_i \mid \mathbf{a}_j]$ spread.



STRENGTHENING THE BOUND – the RIP

Incoherence: Each pair $\mathbf{A}_{i,j} = [\mathbf{a}_i \mid \mathbf{a}_j]$ spread.



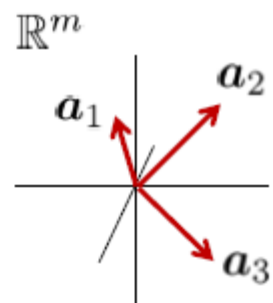
Generalize to subsets of size k :

\mathbf{A}_I well-spread (almost orthonormal) for all I of size k

$$\implies \text{all } k\text{-sparse } \mathbf{x}, \quad \|\mathbf{A}\mathbf{x}\|_2 \approx \|\mathbf{x}\|_2$$

STRENGTHENING THE BOUND – the RIP

Incoherence: Each pair $\mathbf{A}_{i,j} = [\mathbf{a}_i \mid \mathbf{a}_j]$ spread.



Generalize to subsets of size k :

\mathbf{A}_I well-spread (almost orthonormal) for all I of size k

$$\implies \text{all } k\text{-sparse } \mathbf{x}, \quad \|\mathbf{A}\mathbf{x}\|_2 \approx \|\mathbf{x}\|_2$$

\mathbf{A} satisfies the **Restricted Isometry Property** of order k with constant δ if for all k -sparse \mathbf{x} ,

$$(1 - \delta)\|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta)\|\mathbf{x}\|_2^2.$$

IMPLICATIONS OF RIP

Good sparse recovery

Theorem 2 (Candès+Tao '05, Candès '08) .

Suppose $\mathbf{y} = \mathbf{A}\mathbf{x}_0$ with

$$\delta_{2\|\mathbf{x}_0\|_0} < \sqrt{2} - 1.$$

Then \mathbf{x}_0 is the unique optimal solution to

$$\text{minimize } \|\mathbf{x}\|_1 \quad \text{subject to } \mathbf{y} = \mathbf{A}\mathbf{x}.$$

IMPLICATIONS OF RIP

Good sparse recovery

Theorem 2 (Candès+Tao '05, Candès '08) .
Suppose $\mathbf{y} = \mathbf{A}\mathbf{x}_0$ with

$$\delta_{2\|\mathbf{x}_0\|_0} < \sqrt{2} - 1.$$

Then \mathbf{x}_0 is the unique optimal solution to

$$\text{minimize } \|\mathbf{x}\|_1 \quad \text{subject to } \mathbf{y} = \mathbf{A}\mathbf{x}.$$

Again, if ... \mathbf{x}_0 is “structured” and \mathbf{A} is “nice”

we exactly recover \mathbf{x}_0 .

Compare condition to condition $\|\mathbf{x}_0\|_0 < \frac{1}{2}(1 + \mu(\mathbf{A})^{-1})$

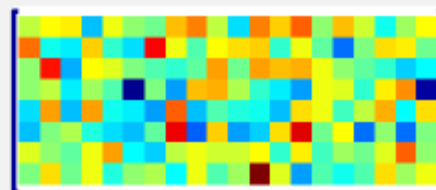
IMPLICATIONS OF RIP

Random \mathbf{A} are great:

If $\mathbf{A} \sim_{iid} \mathcal{N}(0, m^{-1/2})$ then

\mathbf{A} has RIP of order k with

high probability, when $m \geq Ck \log(n/k)$.

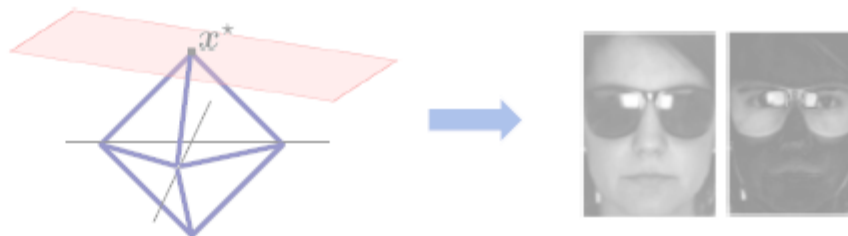


For random \mathbf{A} , ℓ^1 works even when $\|\mathbf{x}_0\|_0 \sim m$.

Useful property for designing sampling operators
(*Compressed sensing*).

WHY CARE ABOUT THE THEORY?

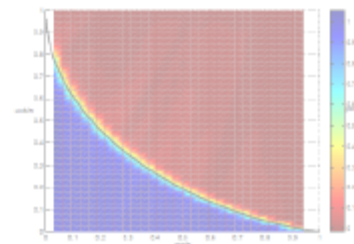
Motivates applications



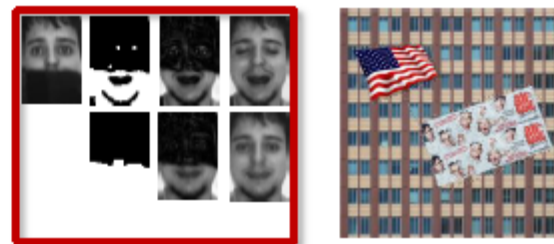
... but be careful: need to justify (and modify) the basic models

Template for stronger results

... predictions can be very sharp in high dimensions.



Generalizes to many other types
of low-dimensional structure



... structured sparsity, low-rank recovery

GENERALIZATIONS – From Sparse to Low-Rank

So far: Recovering *a single sparse vector*:

$$\begin{array}{c} \text{[Image of person with sunglasses]} \\ \mathbf{y} \end{array} = \begin{array}{c} \left[\begin{array}{c} \text{[Image of person with glasses]} \\ \dots \\ \text{[Image of person with glasses]} \end{array} \right] \mathbf{x} + \begin{array}{c} \text{[Image of person with glasses and sunglasses]} \\ \mathbf{e} \end{array}$$

Next: Recovering *low-rank matrix (many correlated vectors)*:

$$\begin{array}{c} \left[\begin{array}{c} \text{[Image of person with sunglasses]} \\ \dots \\ \text{[Image of noise]} \end{array} \right] \\ \mathbf{Y} \end{array} = \begin{array}{c} \left[\begin{array}{c} \text{[Image of person with glasses]} \\ \dots \\ \text{[Image of person with glasses]} \end{array} \right] \\ \mathbf{X} \end{array} + \begin{array}{c} \left[\begin{array}{c} \text{[Image of person with glasses and sunglasses]} \\ \dots \\ \text{[Image of noise]} \end{array} \right] \\ \mathbf{E} \end{array}$$

FORMULATION – Robust PCA?

$$\begin{bmatrix} \text{Face with sunglasses} & \dots & \text{Noisy image} \end{bmatrix} = \begin{bmatrix} \text{Face with glasses} & \dots & \text{Face with glasses} \end{bmatrix} + \begin{bmatrix} \text{Face with dark spots} & \dots & \text{Noisy image} \end{bmatrix}$$

\mathbf{Y} \mathbf{X} \mathbf{E}

Given $\mathbf{Y} = \mathbf{X} + \mathbf{E}$, with \mathbf{X} low-rank, \mathbf{E} sparse, recover \mathbf{X} .

Numerous approaches to **robust PCA** in the literature:

- Multivariate trimming [*Gnanadeskian + Kettering '72*]
- Random sampling [*Fischler + Bolles '81*]
- Alternating minimization [*Ke + Kanade '03*]
- Influence functions [*de la Torre + Black '03*]

Can we give an efficient, provably correct algorithm?

RELATED SOLUTIONS – Matrix recovery

Classical PCA/SVD – low rank + noise [Hotelling '35, Karhunen+Loeve '72,...]

Given $\mathbf{Y} = \mathbf{X} + \mathbf{Z}$, recover \mathbf{X} .

Stable, efficient algorithm, theoretically optimal \rightarrow huge impact

Matrix Completion – low rank, missing data

[Candès + Recht '08,

Candès + Tao '09,

Keshevan, Oh, Montanari '09,

Gross '09,

Ravikumar and Wainwright '10]

From $\mathbf{Y} = \mathcal{P}_\Omega[\mathbf{X}]$, recover \mathbf{X} .

Increasingly well-understood; solvable if \mathbf{X} is low rank and Ω large enough.

WHY IS THE PROBLEM HARD?

Some very sparse matrices are also low-rank:

$$\begin{bmatrix} \blacksquare & & \\ & \blacksquare & \\ & & \blacksquare \end{bmatrix} \rightarrow \begin{bmatrix} \blacksquare & & \\ & \blacksquare & \\ & & \blacksquare \end{bmatrix} + \begin{bmatrix} \blacksquare & & \\ & \blacksquare & \\ & & \blacksquare \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} \blacksquare & & \\ & \blacksquare & \\ & & \blacksquare \end{bmatrix} + \begin{bmatrix} \blacksquare & & \\ & \blacksquare & \\ & & \blacksquare \end{bmatrix}$$

$\mathbf{Y} = \mathbf{1}_{ij}$ $\mathbf{X} = \mathbf{1}_{ij}$ $\mathbf{E} = \mathbf{0}$ $\mathbf{X} = \mathbf{0}$ $\mathbf{E} = \mathbf{1}_{ij}$

Can we recover \mathbf{X} that are incoherent with the standard basis?

Certain sparse error patterns \mathbf{E} make recovering \mathbf{X} impossible:

$$\begin{bmatrix} \text{blue} & \text{red} & \text{green} & \text{yellow} \\ \text{yellow} & \text{red} & \text{green} & \text{yellow} \\ \text{blue} & \text{red} & \text{green} & \text{yellow} \\ \text{yellow} & \text{red} & \text{green} & \text{yellow} \end{bmatrix} + \begin{bmatrix} \text{black} & \text{black} & \text{black} & \text{black} \\ \text{white} & \text{white} & \text{white} & \text{white} \\ \text{black} & \text{black} & \text{black} & \text{black} \\ \text{black} & \text{black} & \text{black} & \text{black} \end{bmatrix} = \begin{bmatrix} \text{blue} & \text{red} & \text{green} & \text{yellow} \\ \text{yellow} & \text{red} & \text{green} & \text{yellow} \\ \text{blue} & \text{red} & \text{green} & \text{yellow} \\ \text{yellow} & \text{red} & \text{green} & \text{yellow} \end{bmatrix}$$

\mathbf{X} $\mathbf{E} = \mathbf{e}_i \mathbf{v}^*$ $\mathbf{Y} = \mathbf{X} + \mathbf{E}$

Can we correct \mathbf{E} whose support is not adversarial?

WHEN IS THERE HOPE? Again, (in)coherence

Can we recover \mathbf{X} that are *incoherent* with the standard basis from *almost all* errors \mathbf{E} ?

Incoherence condition on singular vectors, **singular values arbitrary**:

Singular vectors of \mathbf{X} not too spiky:
$$\begin{cases} \max_i \|\mathbf{U}_i\|^2 \leq \mu r / m. \\ \max_i \|\mathbf{V}_i\|^2 \leq \mu r / n. \end{cases}$$

not too cross-correlated:
$$\|\mathbf{UV}^*\|_\infty \leq \sqrt{\mu r / mn}$$

Uniform model on error support, **signs and magnitudes arbitrary**:


$$\text{support}(\mathbf{E}) \sim \text{uni} \binom{[m] \times [n]}{\rho mn}$$

... AND HOW SHOULD WE SOLVE IT?

Naïve optimization approach

Look for a low-rank \mathbf{X} that agrees with the data up to some sparse error \mathbf{E} :

$$\min \text{rank}(\mathbf{X}) + \gamma \|\mathbf{E}\|_0 \quad \text{subj } \mathbf{X} + \mathbf{E} = \mathbf{Y}.$$


$$\text{rank}(\mathbf{X}) = \#\{\sigma_i(\mathbf{X}) \neq 0\}. \quad \|\mathbf{E}\|_0 = \#\{\mathbf{E}_{ij} \neq 0\}.$$

... AND HOW SHOULD WE SOLVE IT?

Naïve optimization approach

Look for a low-rank \mathbf{X} that agrees with the data up to some sparse error \mathbf{E} :

$$\min \text{rank}(\mathbf{X}) + \gamma \|\mathbf{E}\|_0 \quad \text{subj } \mathbf{X} + \mathbf{E} = \mathbf{Y}.$$

$$\text{rank}(\mathbf{X}) = \#\{\sigma_i(\mathbf{X}) \neq 0\}. \quad \|\mathbf{E}\|_0 = \#\{\mathbf{E}_{ij} \neq 0\}.$$

INTRACTABLE

... AND HOW SHOULD WE SOLVE IT?

Naïve optimization approach

Look for a low-rank \mathbf{X} that agrees with the data up to some sparse error \mathbf{E} :

$$\min \text{rank}(\mathbf{X}) + \gamma \|\mathbf{E}\|_0 \quad \text{subj } \mathbf{X} + \mathbf{E} = \mathbf{Y}.$$

Convex relaxation

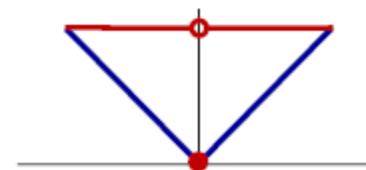
$$\text{rank}(\mathbf{X}) = \#\{\sigma_i(\mathbf{X}) \neq 0\}. \quad \|\mathbf{E}\|_0 = \#\{\mathbf{E}_{ij} \neq 0\}.$$



$$\|\mathbf{X}\|_* = \sum_i \sigma_i(\mathbf{X}).$$



$$\|\mathbf{E}\|_1 = \sum_{ij} |\mathbf{E}_{ij}|.$$



Nuclear norm heuristic: [Fazel, Hindi, Boyd '01], see also [Recht, Fazel, Parillo '08]

[Chandrasekharan et. al. '11]

MAIN RESULT – Correct recovery

Theorem 1 (Principal Component Pursuit). *If $\mathbf{X}_0 \in \mathbb{R}^{m \times n}$, $m \geq n$ has rank*

$$r \leq \rho_r \frac{n}{\mu \log^2(m)}$$

and \mathbf{E}_0 has Bernoulli support with error probability $\rho \leq \rho_s^$, then with very high probability*

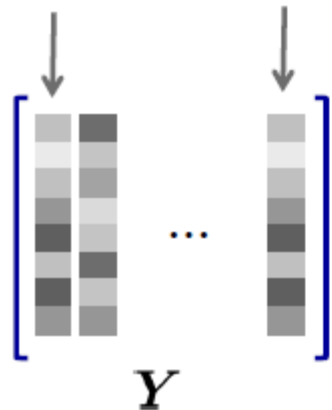
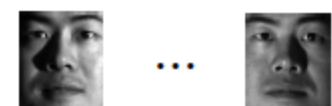
$$(\mathbf{X}_0, \mathbf{E}_0) = \arg \min \|\mathbf{X}\|_* + \frac{1}{\sqrt{m}} \|\mathbf{E}\|_1 \quad \text{subj} \quad \mathbf{X} + \mathbf{E} = \mathbf{X}_0 + \mathbf{E}_0,$$

and the minimizer is unique.

“Convex optimization recovers matrices of rank $O\left(\frac{n}{\log^2 m}\right)$ from errors corrupting $O(mn)$ entries”

EXAMPLE – Faces under varying illumination

58 images of one person under varying lighting:



RPCA

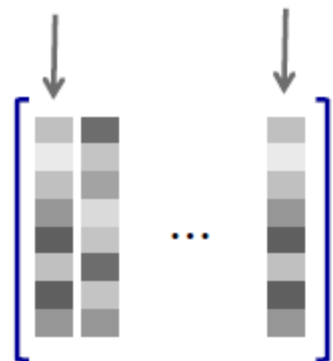


APPLICATIONS – Background modeling from video

Static camera
surveillance video

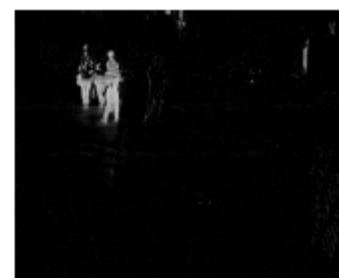
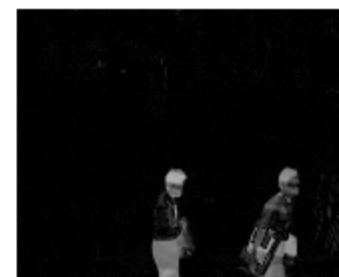
200 frames,
144 x 172 pixels,

Significant foreground
motion



Y

Video Y = Low-rank appx. X + Sparse error E

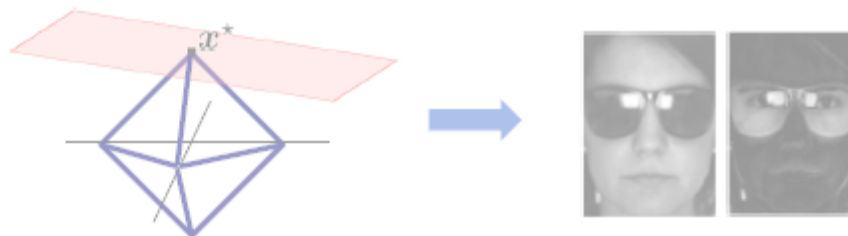


BIG PICTURE – Parallelism of Sparsity and Low-Rank

| | <i>Sparse Vector</i> | <i>Low-Rank Matrix</i> |
|--------------------|-------------------------|---------------------------|
| Degeneracy of | individual signal | correlated signals |
| Measure | L_0 norm $\ x\ _0$ | $\text{rank}(X)$ |
| Convex Surrogate | L_1 norm $\ x\ _1$ | Nuclear norm $\ X\ _*$ |
| Compressed Sensing | $y = Ax$ | $Y = A(X)$ |
| Error Correction | $y = Ax + e$ | $Y = A(X) + E$ |
| Domain Transform | $y \circ \tau = Ax + e$ | $Y \circ \tau = A(X) + E$ |
| Mixed Structures | $Y = A(X) + B(E) + Z$ | |

WHY CARE ABOUT THE THEORY?

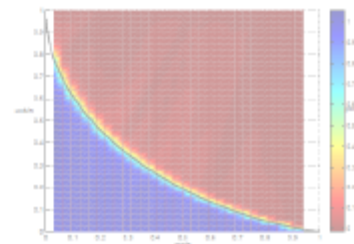
Motivates applications



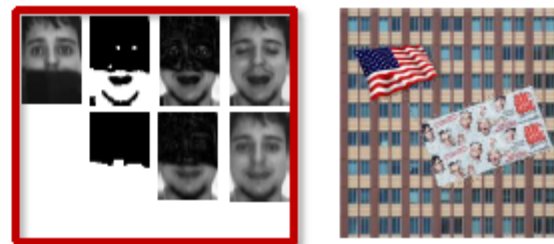
... but be careful: need to justify (and modify) the basic models

Template for stronger results

... predictions can be very sharp in high dimensions.



Generalizes to many other types
of low-dimensional structure



... structured sparsity, low-rank recovery

General theory: constructing norms

Atomic norm: choose a set of atoms \mathcal{A} . Write

$$\|\mathbf{x}\|_{\diamond} = \inf \left\{ \sum_i c_i \mid \sum_i c_i \mathbf{a}_i = \mathbf{x}, c_i > 0, \mathbf{a}_i \in \mathcal{A} \right\}$$

[Chandrasekharan et. al. '12]

General theory: constructing norms

Atomic norm: choose a set of atoms \mathcal{A} . Write

$$\|\mathbf{x}\|_{\diamond} = \inf \left\{ \sum_i c_i \mid \sum_i c_i \mathbf{a}_i = \mathbf{x}, c_i > 0, \mathbf{a}_i \in \mathcal{A} \right\}$$

E.g., **sparsity** $\mathcal{A} = \{\mathbf{e}_i \mid i = 1 \dots n\}$, $\|\mathbf{x}\|_{\diamond} = \|\mathbf{x}\|_{\ell^1}$

low-rank $\mathcal{A} = \{\mathbf{u}\mathbf{v}^* \mid \|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1\}$, $\|\mathbf{x}\|_{\diamond} = \|\mathbf{x}\|_*$

[Chandrasekharan et. al. '12]


General theory: constructing norms

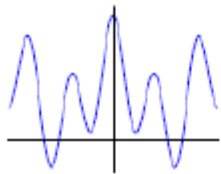
Atomic norm: choose a set of atoms \mathcal{A} . Write

$$\|\mathbf{x}\|_{\diamond} = \inf \left\{ \sum_i c_i \mid \sum_i c_i \mathbf{a}_i = \mathbf{x}, c_i > 0, \mathbf{a}_i \in \mathcal{A} \right\}$$

E.g., sparsity $\mathcal{A} = \{\mathbf{e}_i \mid i = 1 \dots n\}$, $\|\mathbf{x}\|_{\diamond} = \|\mathbf{x}\|_{\ell^1}$

low-rank $\mathcal{A} = \{\mathbf{u}\mathbf{v}^* \mid \|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1\}$, $\|\mathbf{x}\|_{\diamond} = \|\mathbf{x}\|_*$

column sparsity  $\mathcal{A} = \{\mathbf{u}\mathbf{e}_i^* \mid \|\mathbf{u}\|_2 = 1, i = 1 \dots n\}$
e.g., [Xu+Caramanis+Sanghavi'12]

sinusoids  $\mathcal{A} = \{e^{2\pi f t + \xi} \mid f \in [0, 1], \xi \in [0, 2\pi)\}$
[Tang + Recht '12]
[Candes + Fernandez-Garza '12]

General theory: constructing norms

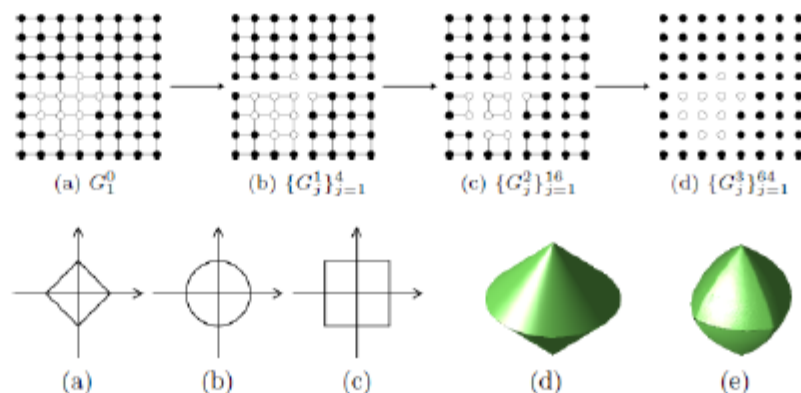
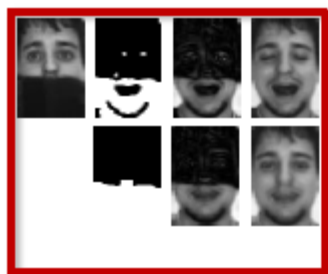
Atomic norm: choose a set of atoms \mathcal{A} . Write

$$\|\mathbf{x}\|_{\diamond} = \inf \left\{ \sum_i c_i \mid \sum_i c_i \mathbf{a}_i = \mathbf{x}, c_i > 0, \mathbf{a}_i \in \mathcal{A} \right\}$$

E.g., sparsity $\mathcal{A} = \{\mathbf{e}_i \mid i = 1 \dots n\}$, $\|\mathbf{x}\|_{\diamond} = \|\mathbf{x}\|_{\ell^1}$

low-rank $\mathcal{A} = \{\mathbf{u}\mathbf{v}^* \mid \|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1\}$, $\|\mathbf{x}\|_{\diamond} = \|\mathbf{x}\|_*$

spatial sparsity



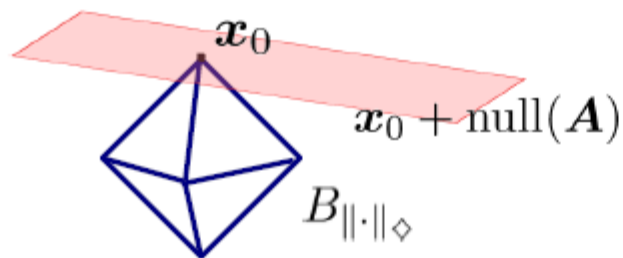
[Bach '11]
[Jia et. '12]

General theory: recovering a single structure

Observe: $y = Ax_0$ with $A \sim \mathcal{N}(0, 1)$ random. When does

$$\min \|x\|_{\diamond} \quad \text{s.t.} \quad Ax = y$$

uniquely recover x_0 ?

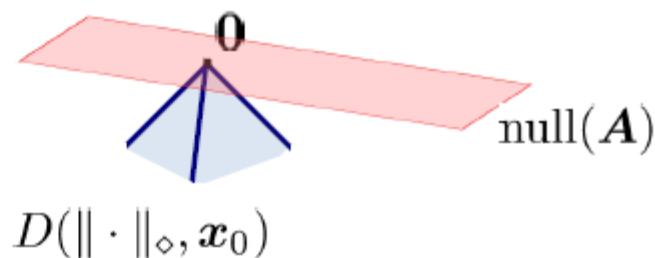
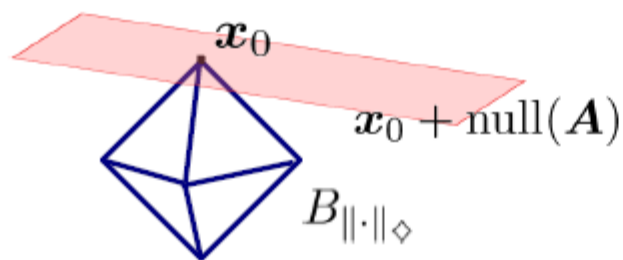


General theory: recovering a single structure

Observe: $y = Ax_0$ with $A \sim \mathcal{N}(0, 1)$ random. When does

$$\min \|x\|_{\diamond} \quad \text{s.t.} \quad Ax = y$$

uniquely recover x_0 ?

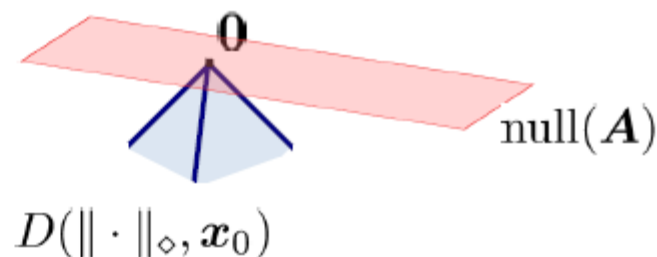
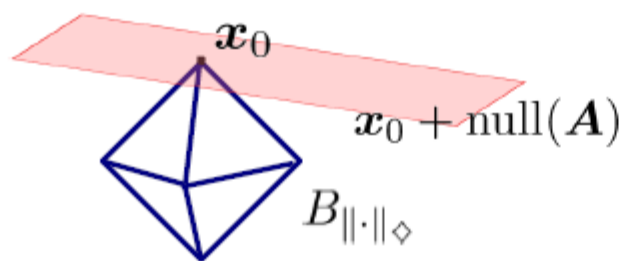


General theory: recovering a single structure

Observe: $\mathbf{y} = \mathbf{A}\mathbf{x}_0$ with $\mathbf{A} \sim \mathcal{N}(0, 1)$ random. When does

$$\min \|\mathbf{x}\|_{\diamond} \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{y}$$

uniquely recover \mathbf{x}_0 ?



Recovery iff the **descent cone**

$$D(\|\cdot\|_{\diamond}, \mathbf{x}_0) = \{\mathbf{v} \mid \|\mathbf{x}_0 + t\mathbf{v}\|_{\diamond} \leq \|\mathbf{x}_0\|_{\diamond} \text{ for some } t > 0\}$$

has $D(\|\cdot\|_{\diamond}, \mathbf{x}_0) \cap \text{null}(\mathbf{A}) = \{\mathbf{0}\}$.

More likely if descent cone is "small". Can we make this precise?

General theory: recovering a single structure

Observe: $\mathbf{y} = \mathbf{A}\mathbf{x}_0$ with $\mathbf{A} \sim \mathcal{N}(0, 1)$ random. When does

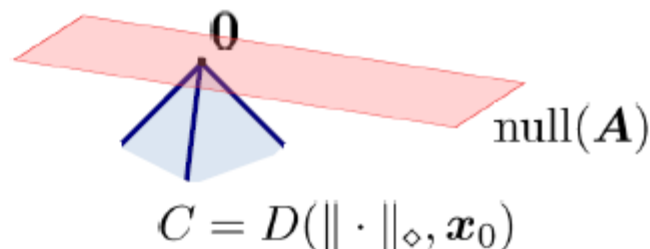
$$\min \|\mathbf{x}\|_{\diamond} \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{y}$$

uniquely recover \mathbf{x}_0 ?

The **statistical dimension** of a cone C is

$$\delta(C) = \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(0,1)} \left[\|P_C \mathbf{g}\|^2 \right].$$

Many nice properties. E.g., if C a subspace, $\delta(C) = \dim(C)$.



General theory: recovering a single structure

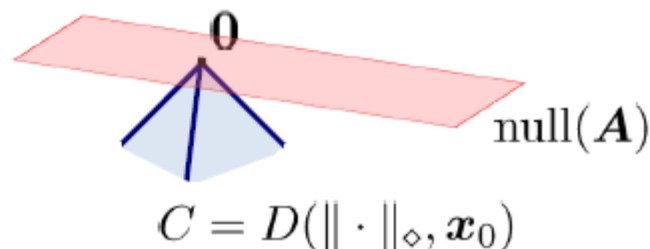
Observe: $\mathbf{y} = \mathbf{A}\mathbf{x}_0$ with $\mathbf{A} \sim \mathcal{N}(0, 1)$ random. When does

$$\min \|\mathbf{x}\|_{\diamond} \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{y}$$

uniquely recover \mathbf{x}_0 ?

The **statistical dimension** of a cone C is

$$\delta(C) = \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(0,1)} \left[\|\mathbf{P}_C \mathbf{g}\|^2 \right].$$



Many nice properties. E.g., if C a subspace, $\delta(C) = \dim(C)$.

Sharp **phase transition** at $m = \delta(C)$:

$$\begin{aligned} m > \delta(C) &\implies \mathbb{P}[\text{recovery}] > 1 - \exp(-c(m - \delta(C))^2/n) \\ m < \delta(C) &\implies \mathbb{P}[\text{recovery}] < \exp(-c(m - \delta(C))^2/n) \end{aligned}$$

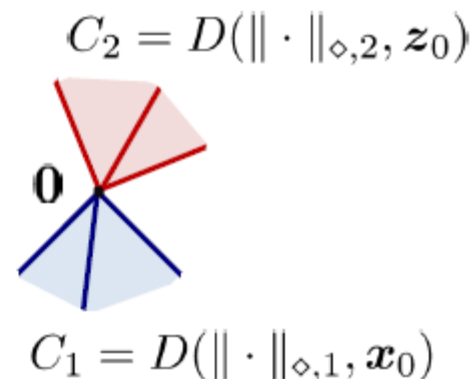
General theory: decomposing two structures

Observe: $\mathbf{y} = \mathbf{x}_0 + \mathbf{z}_0$ with regularizers $\|\mathbf{x}\|_{\diamond,1}$, $\|\mathbf{z}\|_{\diamond,2}$. Does

$$\min \|\mathbf{x}\|_{\diamond,1} + \|\mathbf{z}\|_{\diamond,2} \quad \text{s.t.} \quad \mathbf{x} + \mathbf{z} = \mathbf{y}$$

uniquely recover $\mathbf{x}_0, \mathbf{z}_0$?

Variant: $\min \|\mathbf{x}\|_{\diamond,1} \quad \text{s.t.} \quad \|\mathbf{z}\|_{\diamond,2} \leq 1, \mathbf{x} + \mathbf{z} = \mathbf{y}$



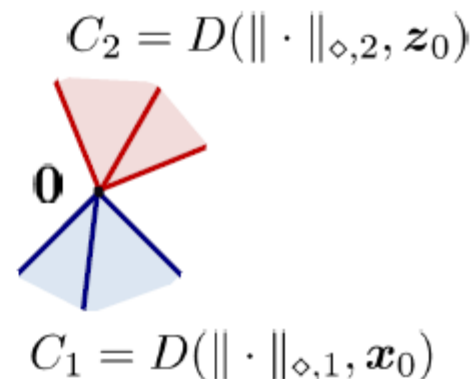
General theory: decomposing two structures

Observe: $\mathbf{y} = \mathbf{x}_0 + \mathbf{z}_0$ with regularizers $\|\mathbf{x}\|_{\diamond,1}$, $\|\mathbf{z}\|_{\diamond,2}$. Does

$$\min \|\mathbf{x}\|_{\diamond,1} + \|\mathbf{z}\|_{\diamond,2} \quad \text{s.t.} \quad \mathbf{x} + \mathbf{z} = \mathbf{y}$$

uniquely recover $\mathbf{x}_0, \mathbf{z}_0$?

Variant: $\min \|\mathbf{x}\|_{\diamond,1} \quad \text{s.t.} \quad \|\mathbf{z}\|_{\diamond,2} \leq 1, \mathbf{x} + \mathbf{z} = \mathbf{y}$



In a random incoherence model (C_2 randomly rotated), **phase transition** at

$$\delta(C_1) + \delta(C_2) = n$$

$$n > \delta(C_1) + \delta(C_2) \implies \mathbb{P}[\text{recovery}] > 1 - \exp(-c(n - \delta(C_1) - \delta(C_2))^2/n)$$

$$n < \delta(C_1) + \delta(C_2) \implies \mathbb{P}[\text{recovery}] < \exp(-c(n - \delta(C_1) - \delta(C_2))^2/n)$$

General theory: statistical estimation

Observe: noisy measurements $\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{z}$. Noise-aware optimization:

$$\min \|\mathbf{x}\|_{\diamond} + \frac{\gamma}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$$

E.g., **Basis pursuit denoising:** $\min \|\mathbf{x}\|_1 + \frac{\lambda}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$

Noise-aware RPCA: $\min \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \frac{\gamma}{2} \|\mathbf{L} + \mathbf{S} - \mathbf{D}\|_F^2$

When does $\min \|\mathbf{x}\|_{\diamond} + \frac{\gamma}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$ produce $\hat{\mathbf{x}} \approx \mathbf{x}_0$?

General theory for **decomposable regularizers** $\|\cdot\|_{\diamond}$

[Negahbhan, Agarwal, Yu, Wainwright '12]

A suite of models and theoretical guarantees

For robust recovery of a family of low-dimensional structures:

- [Zhou et. al. '09] **Spatially contiguous** sparse errors via MRF
- [Bach '10] – structured relaxations from **submodular functions**
- [Negahban+Yu+Wainwright '10] – **geometric analysis** of recovery
- [Becker+Candès+Grant '10] – **algorithmic templates**
- [Xu+Caramanis+Sanghavi '11] **column sparse errors** $L_{2,1}$ norm
- [Recht+Parillo+Chandrasekaran+Wilsky '11] – **compressive sensing** of various structures
- [Candes+Recht '11] – **compressive sensing of decomposable structures**

$$X^0 = \arg \min \|X\|_{\diamond} \quad \text{s.t.} \quad \mathcal{P}_Q(X) = \mathcal{P}_Q(X^0)$$

- [McCoy+Tropp'11] – **decomposition of sparse and low-rank structures**

$$(X_1^0, X_2^0) = \arg \min \|X_1\|_{(1)} + \lambda \|X_2\|_{(2)} \quad \text{s.t.} \quad X_1 + X_2 = X_1^0 + X_2^0$$

- [W.+Ganesh+Min+Ma, I&I'13] – **superposition of decomposable structures**

$$(X_1^0, \dots, X_k^0) = \arg \min \sum \lambda_i \|X_i\|_{(i)} \quad \text{s.t.} \quad \mathcal{P}_Q(\sum_i X_i) = \mathcal{P}_Q(\sum_i X_i^0)$$

Take home message: Let the data and application tell you the structure...