

Sparse and Low-Rank Representation

Lecture I: Motivation and Theory

Yi Ma

MSRA and UIUC

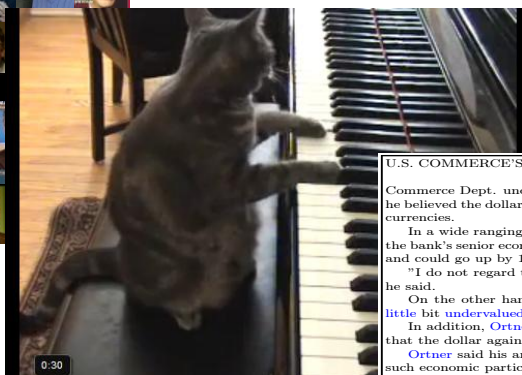
Allen Yang

UC Berkeley

John Wright

Columbia University

CONTEXT – Data increasingly massive, high-dimensional...



U.S. COMMERCE'S ORTNER SAYS YEN UNDERVALUED

Commerce Dept. undersecretary of economic affairs Robert Ortner said that he believed the dollar at current levels was fairly priced against most European currencies.

In a wide ranging address sponsored by the Export-Import Bank, Ortner, the bank's senior economist also said he believed that the yen was undervalued and could go up by 10 or 15 pct.

"I do not regard the dollar as undervalued at this point against the yen," he said.

On the other hand, Ortner said that he thought that "the yen is still a little bit undervalued," and "could go up another 10 or 15 pct."

In addition, Ortner, who said he was speaking personally, said he thought that the dollar against most European currencies was "fairly priced" on such economic particulars as wage rate differentiations.

Ortner said his analysis of the various exchange rate values was "fairly level" of the dollar because at the time of the Plaza Accord, the dollar was overvalued and that the first 15 pct decline had little impact.

He said there were indications now that the trade deficit was likely to level off.

Turning to Brazil and Mexico, Ortner made it clear that it would be almost impossible for those countries to earn enough foreign exchange to pay the service on their debts. He said the best way to deal with this was to use the policies outlined in Treasury Secretary James Baker's debt initiative.

Images

➤ 1M pixels

Videos

➤ 1B voxels

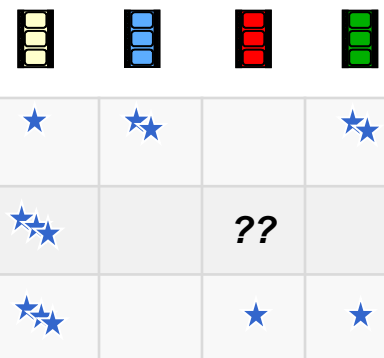
Web data

➤ 100B webpages



User data

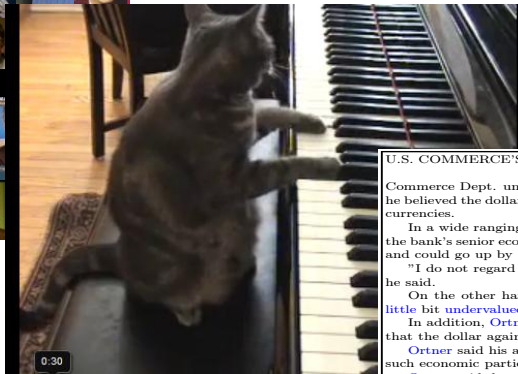
➤ 1B users



CONTEXT – Discovering knowledge from data



Images



Videos

U.S. COMMERCE'S ORTNER SAYS YEN UNDERVALUED

Commerce Dept. undersecretary of economic affairs Robert Ortner said that he believed the dollar at current levels was fairly priced against most European currencies.

In a wide ranging address sponsored by the Export-Import Bank, Ortner, the bank's senior economist also said he believed that the yen was undervalued and could go up by 10 or 15 pct.

"I do not regard the dollar as undervalued at this point against the yen," he said.

On the other hand, Ortner said that he thought that "the yen is still a little bit undervalued," and "could go up another 10 or 15 pct."

In addition, Ortner, who said he was speaking personally, said that the dollar against most European currencies was "fairly priced."

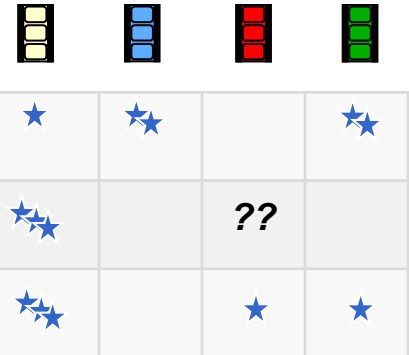
Ortner said his analysis of the various exchange rate values was based on such economic particulars as wage rate differentiations.

Ortner said there had been little impact on U.S. trade deficit by the overvaluation of the dollar because at the time of the Plaza Accord, the dollar was undervalued and that the first 15 pct decline had little impact.

He said there were indications now that the trade deficit was beginning to level off.

Turning to Brazil and Mexico, Ortner made it clear that it would be almost impossible for those countries to earn enough foreign exchange to pay the service on their debts. He said the best way to deal with this problem was the policies outlined in Treasury Secretary James Baker's debt initiative.

Web data



User data

How to extract **compact knowledge** from such **massive datasets**?

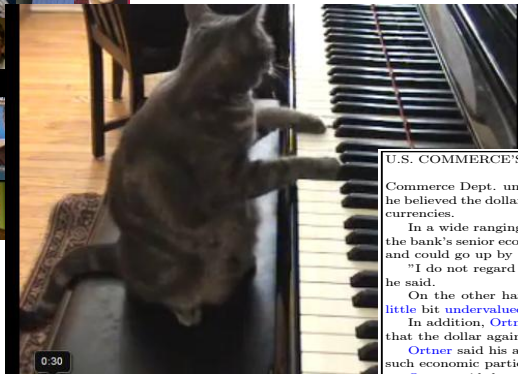
CONTEXT – Good solutions impact many applications



Images



Compression
Denoising
Superresolution
Recognition...



Videos



Streaming
Tracking
Stabilization...

U.S. COMMERCE'S ORTNER SAYS YEN UNDERVALUED

Commerce Dept. undersecretary of economic affairs Robert Ortner said that he believed the dollar at current levels was fairly priced against most European currencies.

In a wide ranging address sponsored by the Export-Import Bank, Ortner, the bank's senior economist also said he believed that the yen was undervalued, and could go up by 10 or 15 pct.

"I do not regard the dollar as undervalued at this point against the yen," he said.

On the other hand, Ortner said that he thought that "the yen is still a little bit undervalued," and "could go up another 10 or 15 pct."

In addition, Ortner, who said he was speaking personally, said that the dollar against most European currencies was "fairly priced."

Ortner said his analysis of the various exchange rate values was based on such economic particulars as wage rate differentiations.

Ortner said there had been little impact on U.S. trade deficit by the overvaluation and that the first 15 pct decline had little impact.

He said there were indications now that the trade deficit was beginning to level off.

Turning to Brazil and Mexico, Ortner made it clear that it would be almost impossible for those countries to earn enough foreign exchange to service on their debts. He said the best way to deal with this problem was the policies outlined in Treasury Secretary James Baker's debt initiative.

Web data



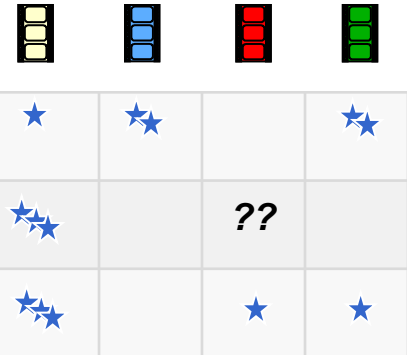
Indexing
Ranking
Search...



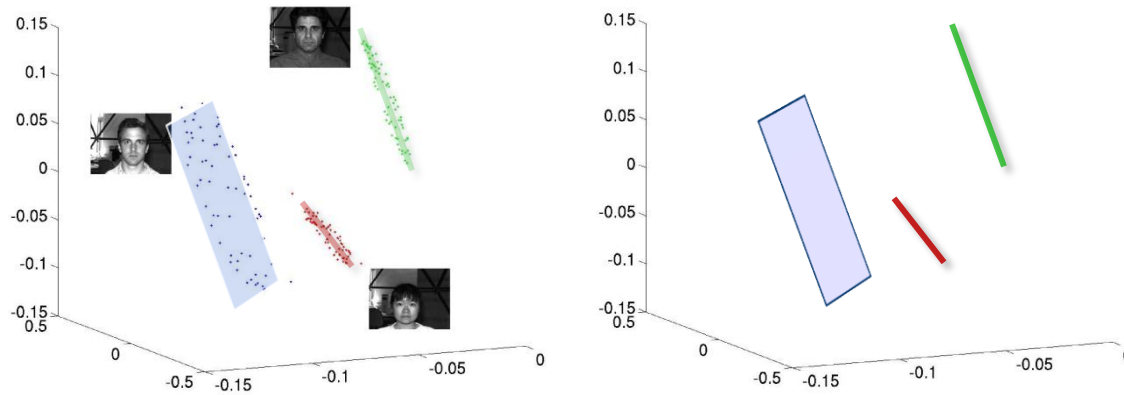
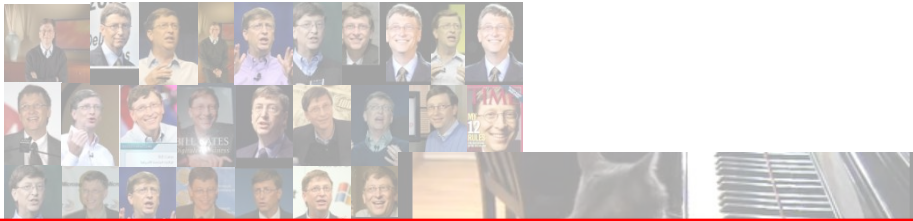
User data



Clustering
Classification
Collaborative filtering...



Low-dimensional structures in high-dimensional data

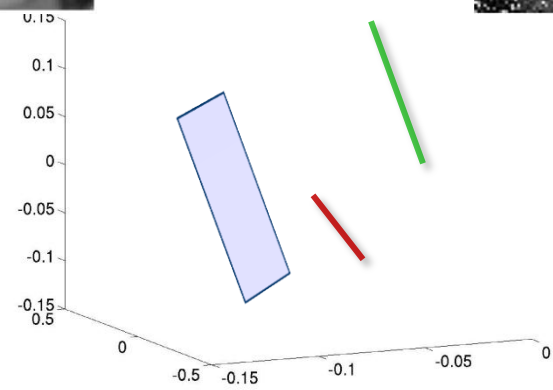
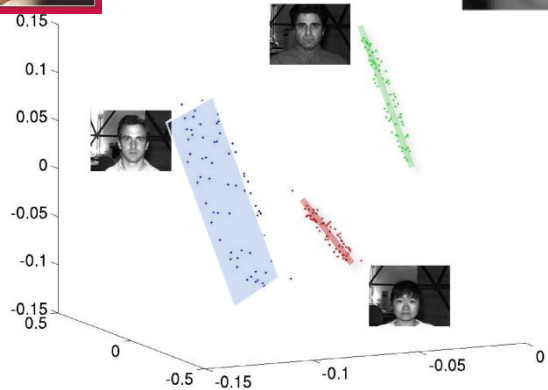


How can we learn and exploit low-dimensional structures in high-dimensional data?

But it is not so easy...



?



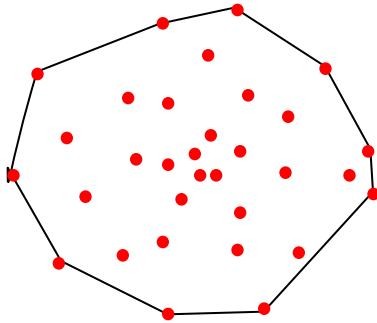
Real application data often contain missing observations, corruptions, or even malicious errors.

Classical methods (e.g., least squares, PCA) break down...

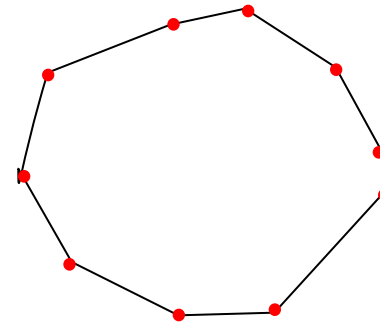
CONTEXT - *New Phenomena with High-Dimensional Data*

KEY CHALLENGE: efficiently and reliably recover low-dimensional structures from high-dimensional data, despite gross observation errors.

A sobering message: human intuition is **severely limited** in high-dimensional spaces:



Gaussian samples in 2D



As dimension grows **proportionally** with the number of samples...

A new regime of geometry, statistics, and computation...

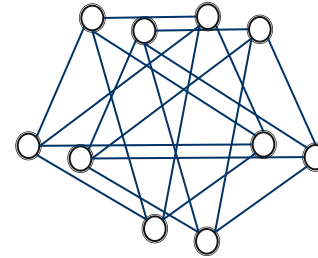
CONTEXT - *Massive High-Dimensional Data*



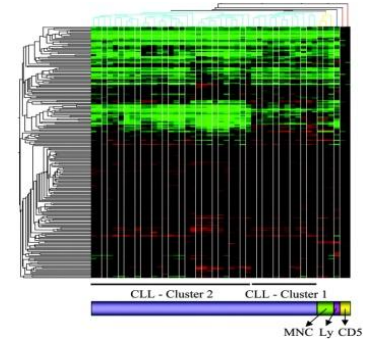
Recognition



Surveillance



Search and Ranking



Bioinformatics

The curse of dimensionality:

*...increasingly demand inference with **limited samples** for very **high-dimensional data**.*

The blessing of dimensionality:

*... real data highly concentrate on **low-dimensional, sparse, or degenerate structures** in the high-dimensional space.*

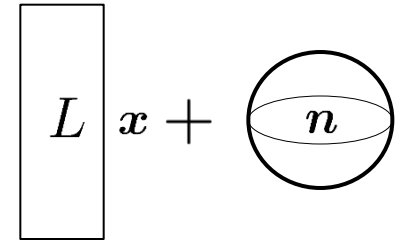
But nothing is free: ***Gross errors and irrelevant measurements** are now ubiquitous in massive cheap data.*

Everything old ...

A long and rich history of robust estimation with error correction and missing data imputation:



R. J. Boscovich. *De calculo probabilitatum que respondent diversis valoribus summe errorum post plures observationes ...* , before 1756



**over-determined
+ dense, Gaussian**

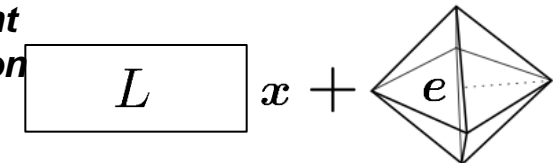


A. Legendre. *Nouvelles methodes pour la determination des orbites des cometes*, 1806

C. Gauss. *Theory of motion of heavenly bodies*, 1809



A. Beurling. *Sur les integrales de Fourier absolument convergentes et leur application a une transformation fonctionnelle*, 1938



**underdetermined
+ sparse, Laplacian**



B. Logan. *Properties of High-Pass Signals*, 1965

⋮

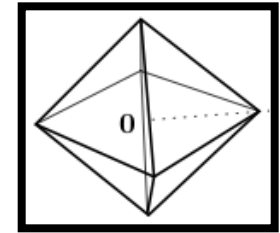
... IS NEW AGAIN

Today, robust estimation in high dimensions is more urgent, more tractable, and increasingly sharply understood.

Theory – **high-dimensional** geometry & statistics, measure concentration, combinatorics, coding theory...

Algorithms – **large scale** convex optimization, parallel and distributed computing....

Applications – **massive** data driven methods, sensing and hashing, denoising, superresolution, MRI, bioinformatics, image classification, recognition ...



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



Today's plan

Lecture I: Motivation and Theory

Lecture II: Efficient Optimization for Sparse Representation

Lecture III: Applications and Generalizations

Q&A and discussion

Lecture I

Theory of Sparse and Low-Rank Recovery

John Wright

Electrical Engineering

Columbia University

UNDERDETERMINED LINEAR SYSTEMS

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

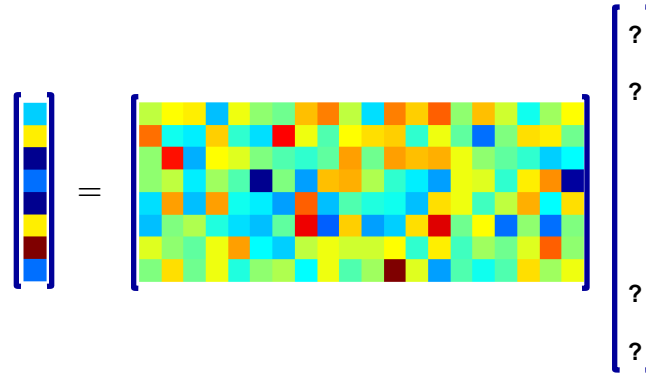
$A \in \mathbb{R}^{m \times n}$

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

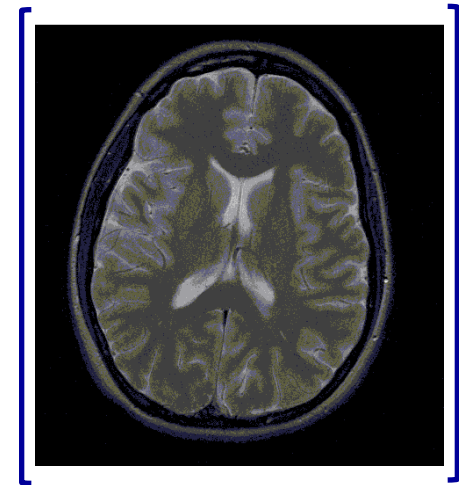
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$



Signal acquisition



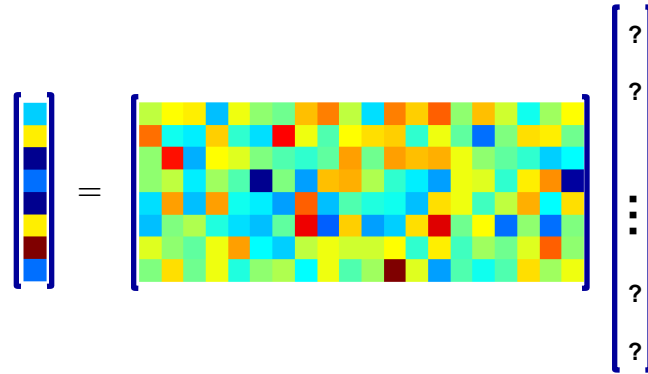
z

Image to be sensed

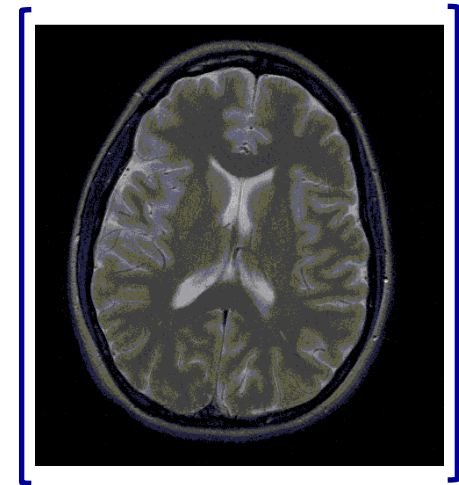
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$



Signal acquisition



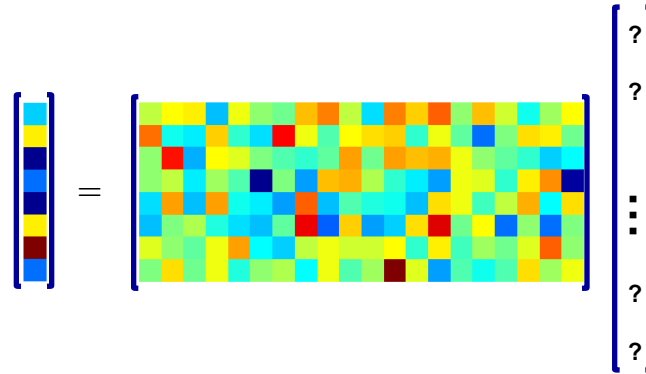
z

Image to be sensed

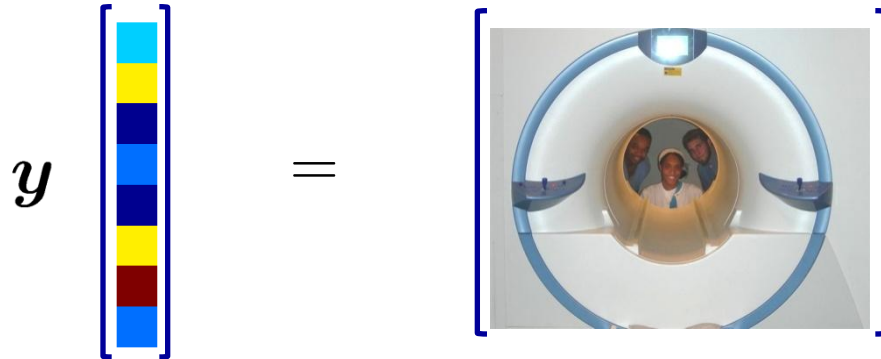
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$

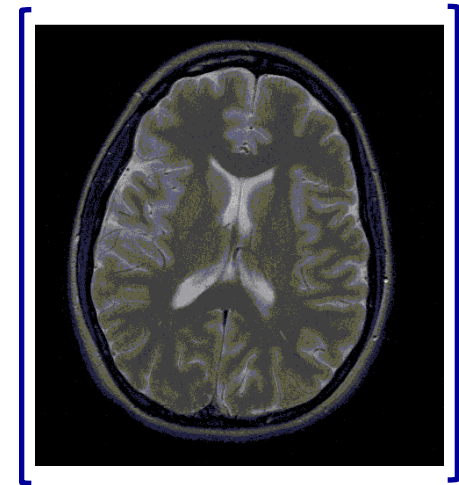


Signal acquisition



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

Observations are Fourier coefficients!



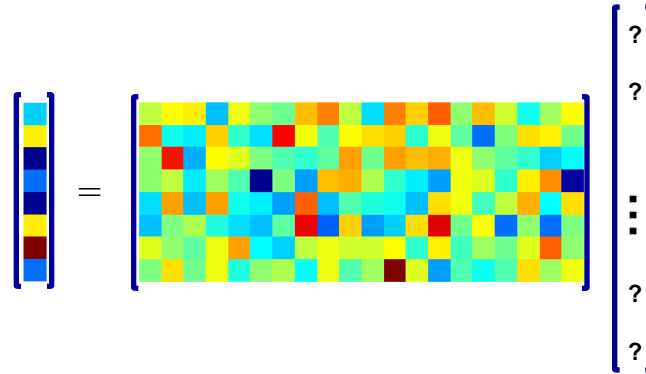
z

Image to be sensed

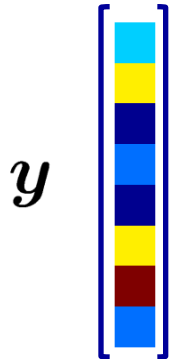
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

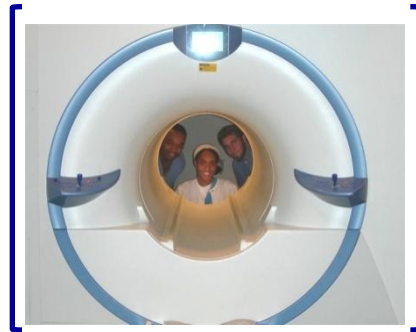
$$y = Ax$$



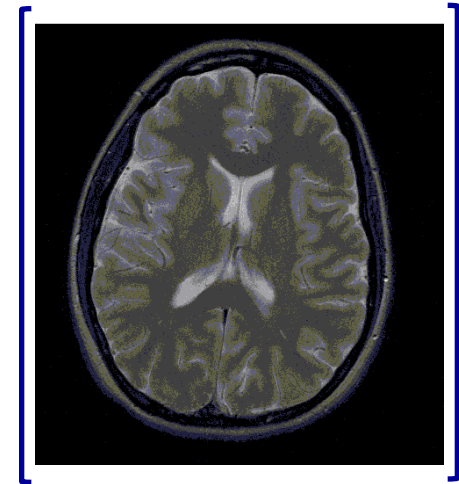
Signal acquisition



=



F_{Ω}



z

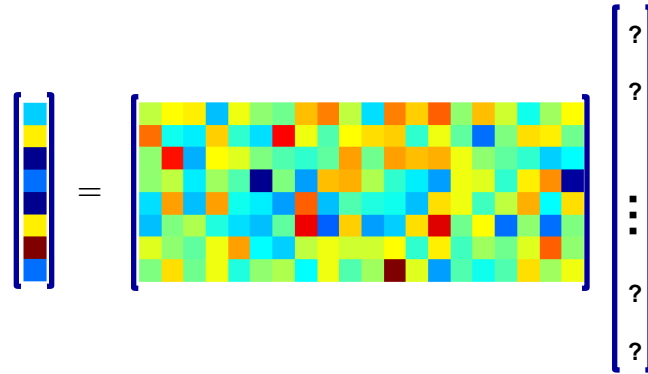
A few Fourier coefficients

Image to be sensed

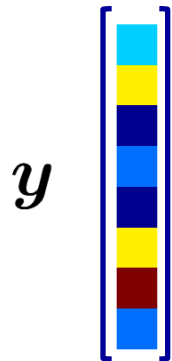
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

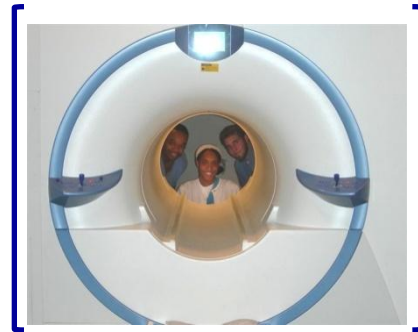
$$y = Ax$$



Signal acquisition

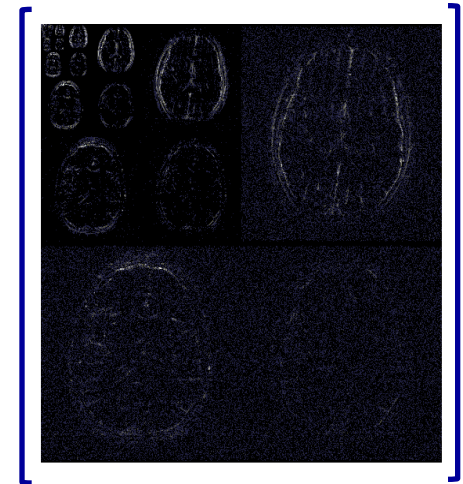


=



F_{Ω}

Ψ



x

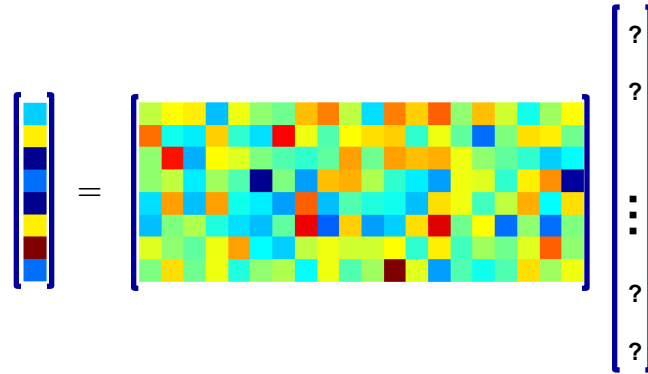
A few Fourier coefficients

Wavelet coefficients: $z = \Psi x$

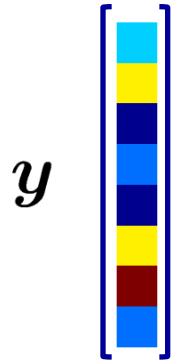
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$

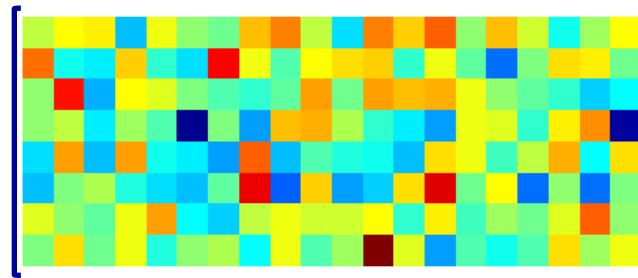


Signal acquisition

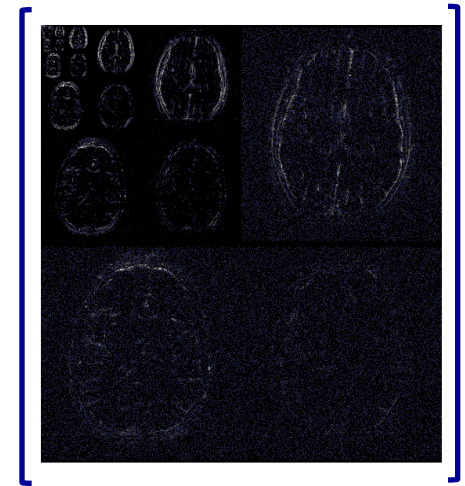


A few Fourier coefficients

=



A



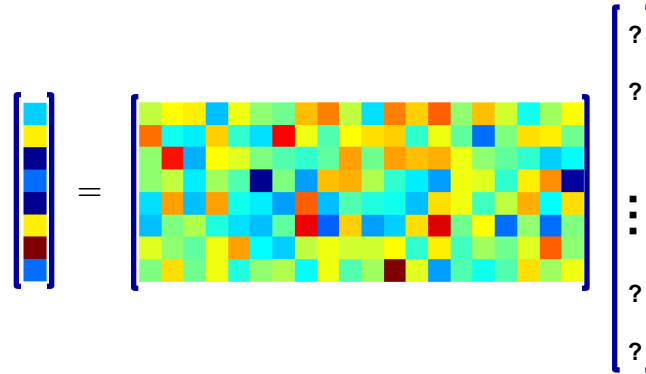
x

Wavelet coefficients

UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$



Compression

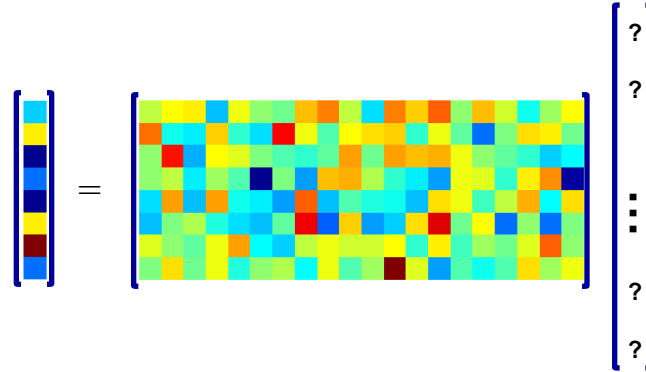


Image to be
compressed

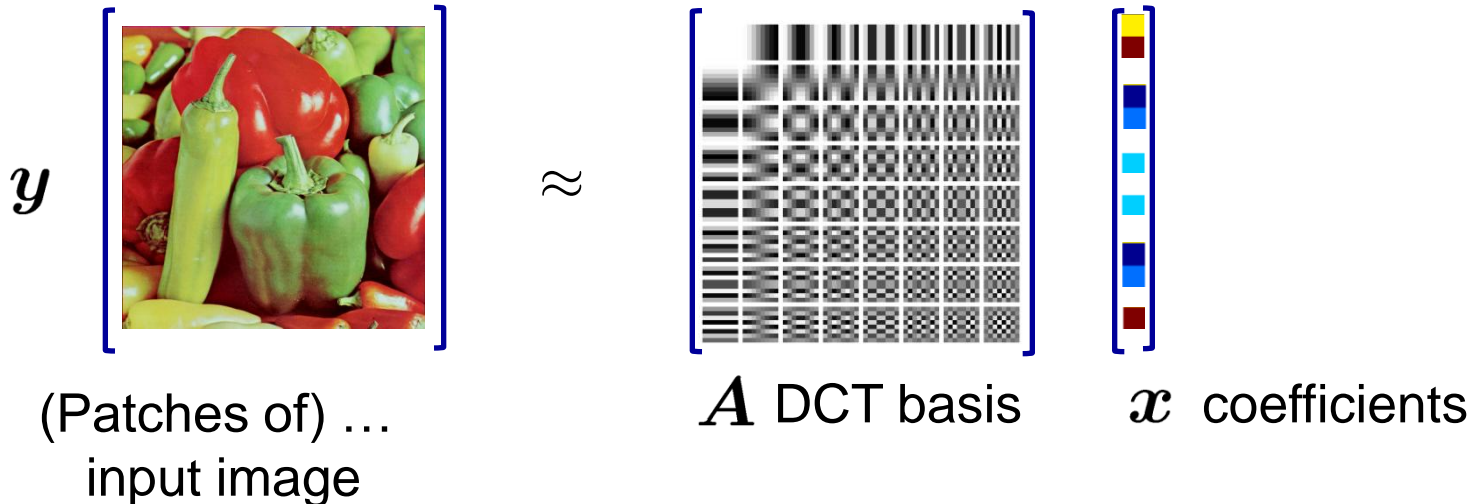
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$



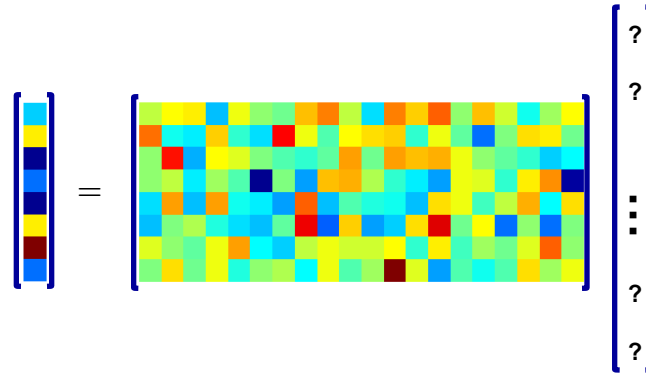
Compression – JPEG



UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$

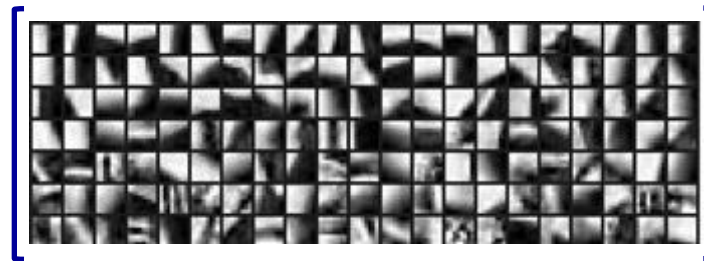


Compression – Learned dictionary



(Patches of) ...
input image

\approx



A Learned dictionary

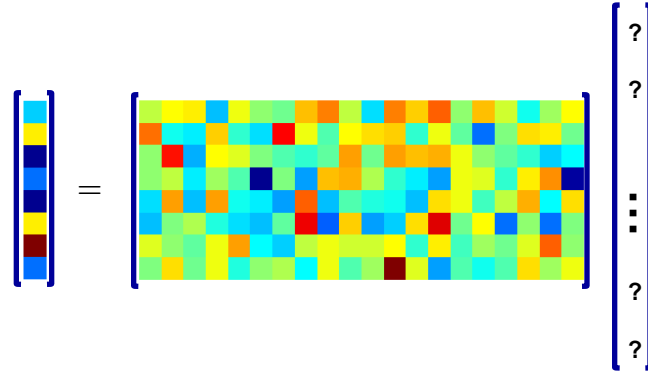


x coefficients

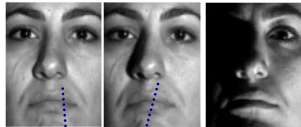
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

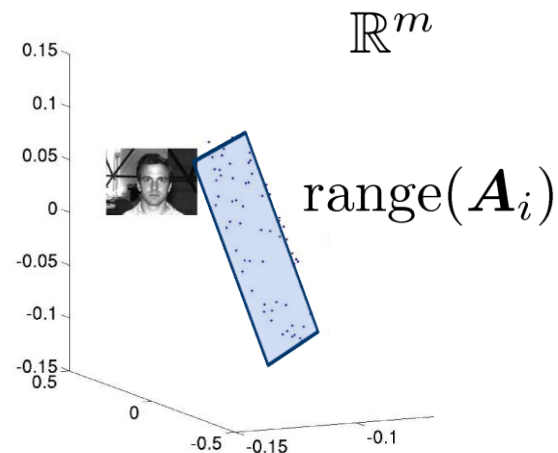
$$y = Ax$$



Recognition



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

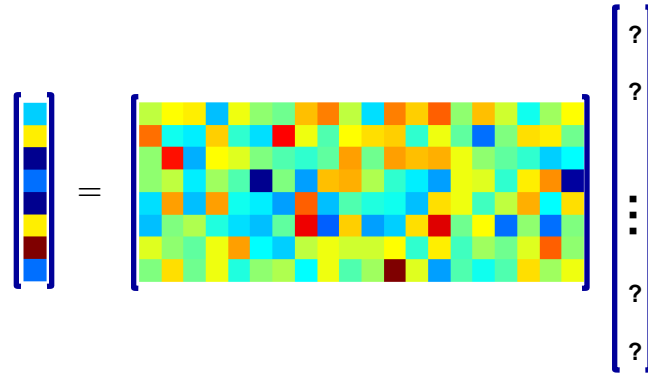


Linear subspace model for images of **same face** under **varying lighting**.

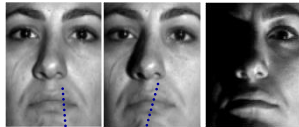
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

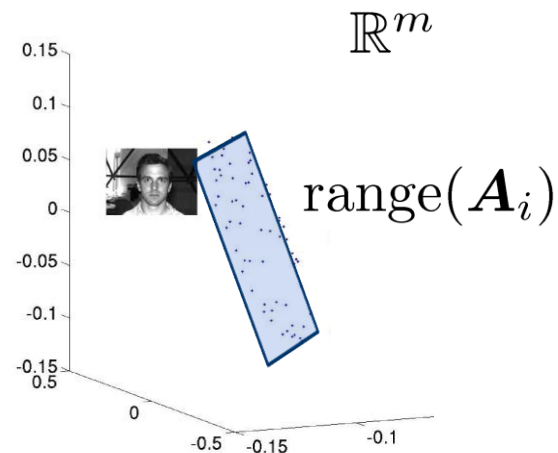
$$y = Ax$$



Recognition



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

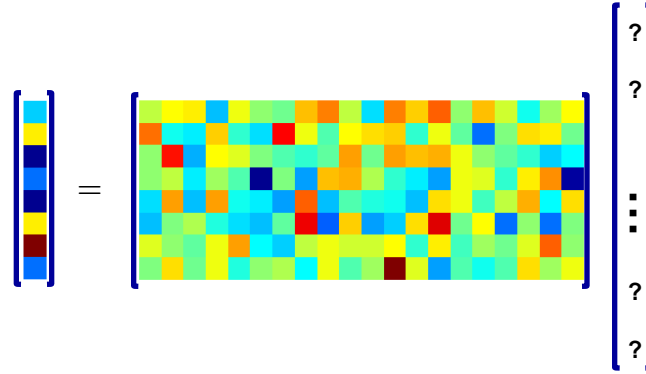


$$y \text{ (face image)} \approx x_{i,1} \text{ (face image)} + x_{i,2} \text{ (face image)} + \dots + x_{i,n} \text{ (face image)} = A_i x_i$$

UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$



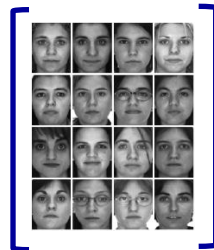
Recognition



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

Test image

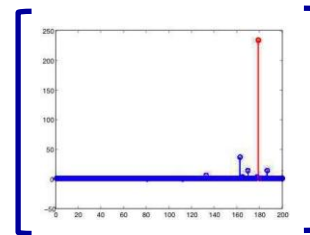
=



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

Combined
training
dictionary

×



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

coefficients

+



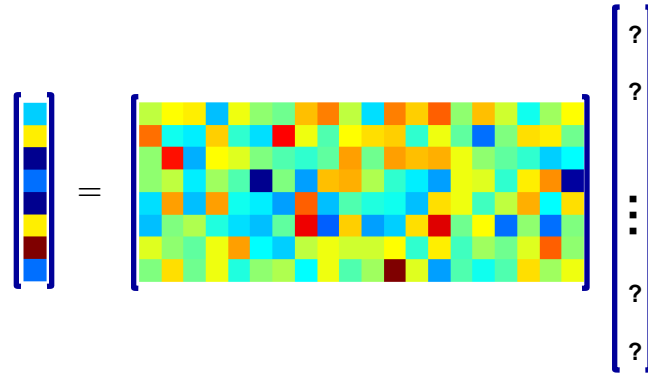
$$e \in \mathbb{R}^m$$

corruption,
occlusion

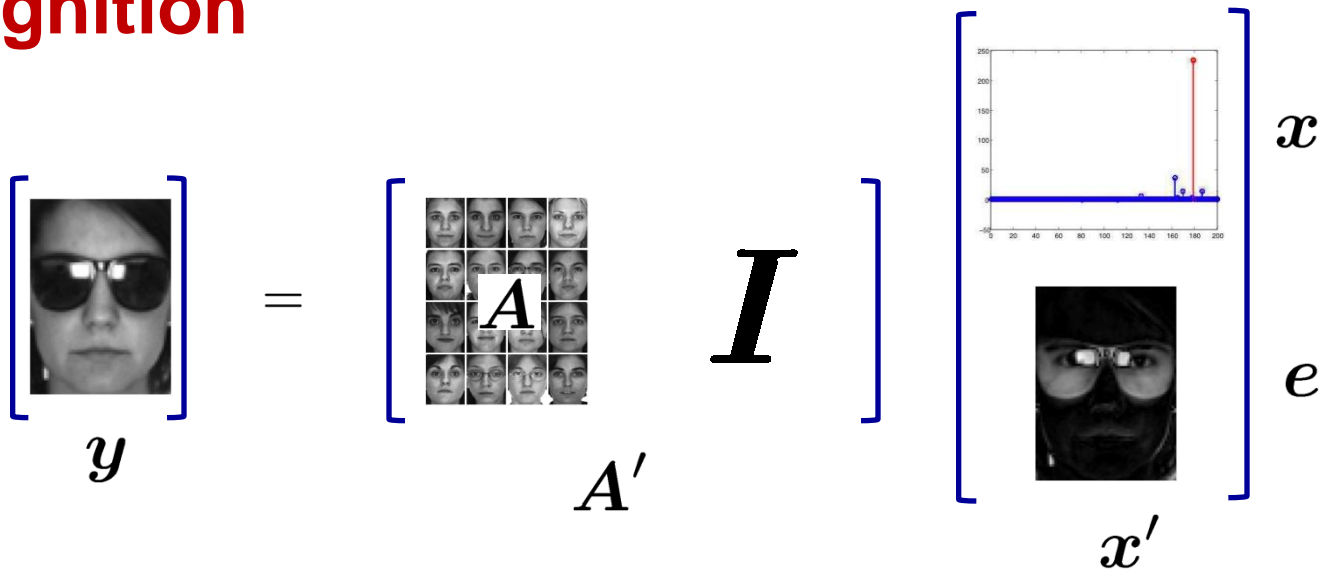
UNDERDETERMINED LINEAR SYSTEMS

Underdetermined system

$$y = Ax$$



Recognition



One large underdetermined system: $y = A'x'$.

UNDERDETERMINED LINEAR SYSTEMS

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

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

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

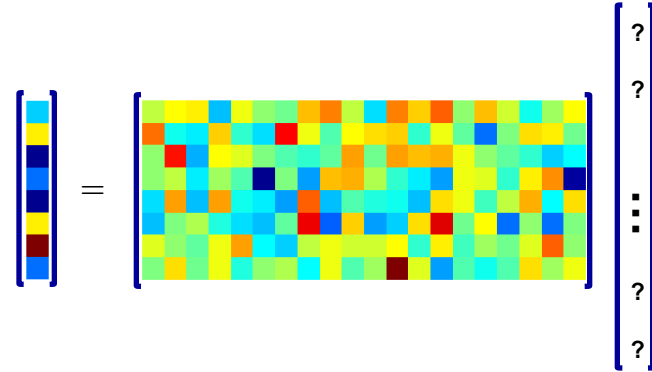
In all of these examples, $\underbrace{m}_{\text{\#observations}} \ll \underbrace{n}_{\text{\#unknowns}}$.

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

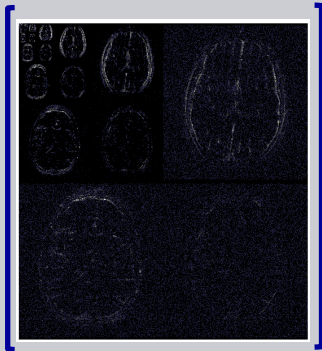
WHAT DO WE KNOW ABOUT x ?

Underdetermined system

$$y = Ax$$

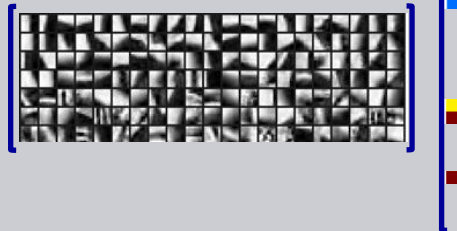


Signal acquisition



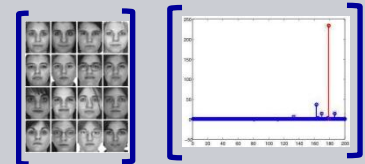
x^* contains **just a few** significant wavelet coefficients.

Image compression

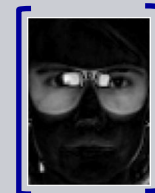


x^* uses **just a few** dictionary elements.

Face Recognition



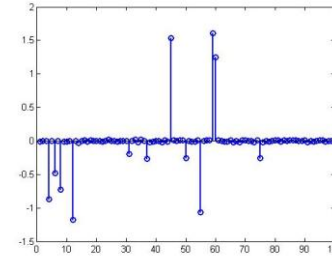
x^* uses **just a few** training faces.



e^* corrects **a few** gross errors.

SPARSITY – More formally

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

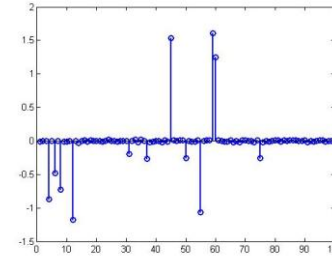


The **number of nonzeros** is called the ℓ^0 -“norm” of \mathbf{x} :

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

SPARSITY – More formally

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



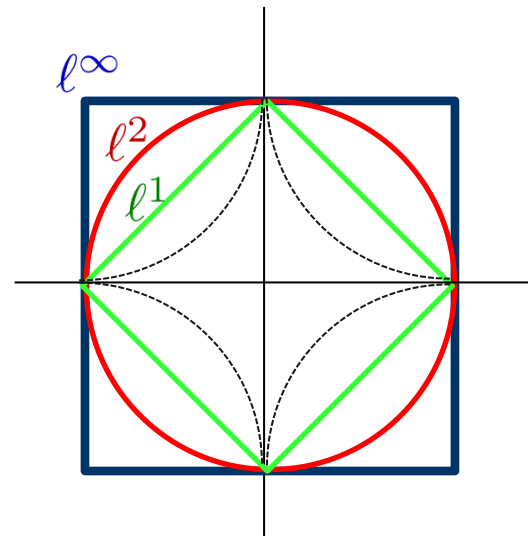
The **number of nonzeros** is called the ℓ^0 -“norm” of \mathbf{x} :

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

Geometrically

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

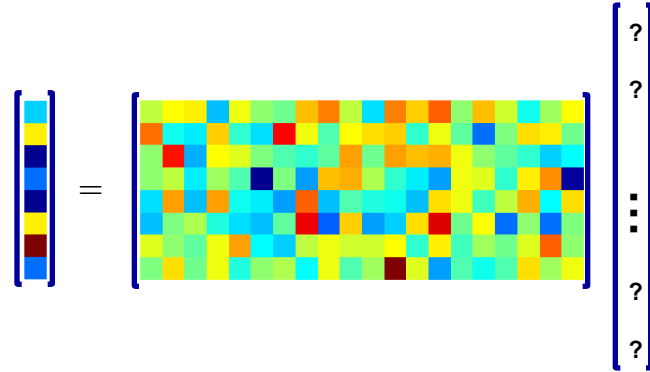
$$\|\mathbf{x}\|_0 = \lim_{p \searrow 0} \|\mathbf{x}\|_p^p.$$



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

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

A diagram illustrating the equation $\mathbf{y} = \mathbf{A}\mathbf{x}$. On the left is a small vertical vector \mathbf{y} with five colored elements (yellow, blue, blue, yellow, red). This is followed by an equals sign, a large heatmap matrix \mathbf{A} with 10 columns and 10 rows of colored squares. To the right of the matrix is a tall vertical vector \mathbf{x} with question marks in its top and bottom sections and a vertical ellipsis in the middle.

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) .

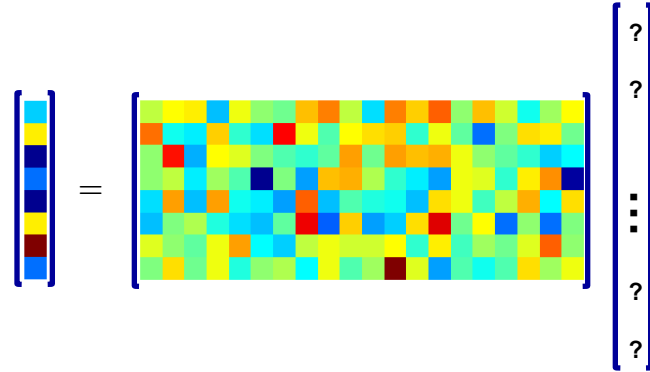
Suppose $\mathbf{y} = \mathbf{A}\mathbf{x}_0$, and let $k = \|\mathbf{x}_0\|_0$. If $\text{null}(\mathbf{A})$ contains no $2k$ -sparse vectors, \mathbf{x}_0 is the unique optimal solution to

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

THE SPARSEST SOLUTION

Underdetermined system

$$y = Ax$$



Look for the sparsest x that agrees with our observation:

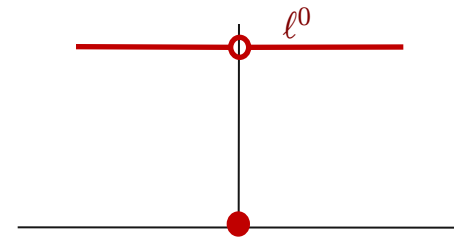
~~minimize $\|x\|_0$ subject to $Ax = y$.~~

INTRACTABLE

RELAX!

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

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



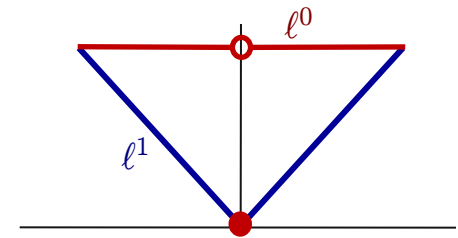
RELAX!

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

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!

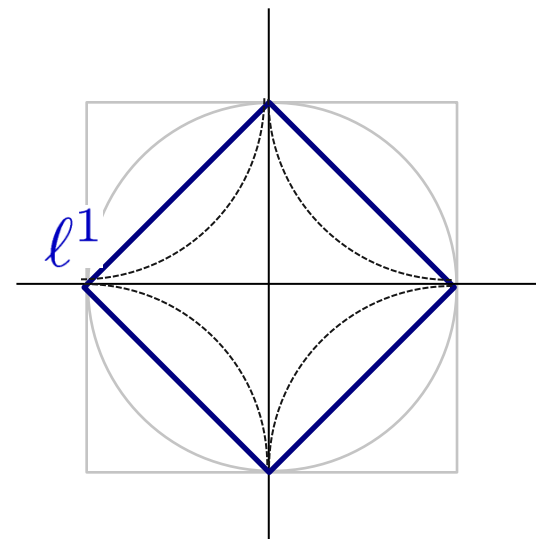
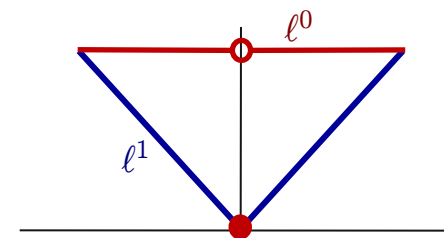
~~minimize $\|\mathbf{x}\|_0$ subject to $\mathbf{Ax} = \mathbf{y}$.~~

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

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



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

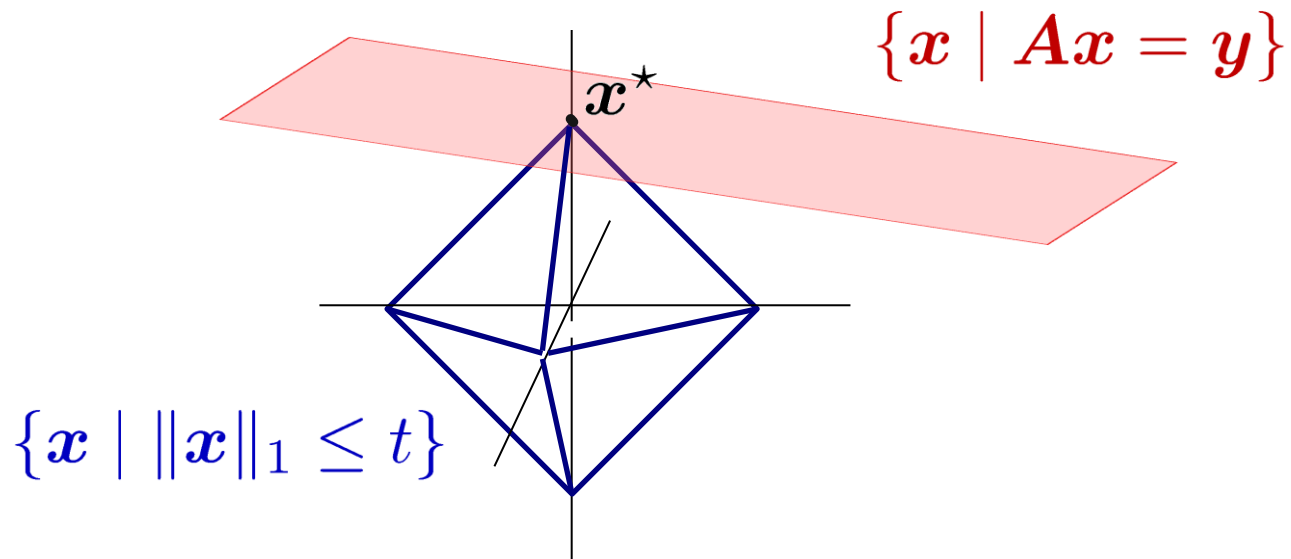
Efficiently solvable
– Lecture 2!

Have we lost anything? [demo]

WHY DOES THIS WORK? Geometric intuition

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

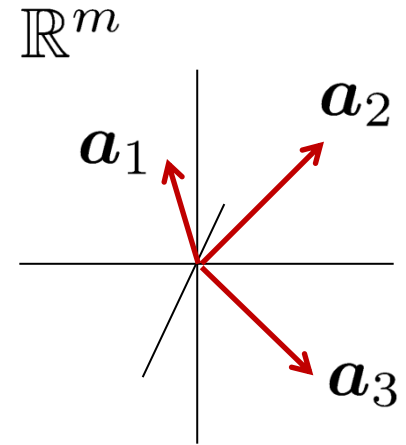
\mathbb{R}^n



WHY DOES THIS WORK? More formally...

We see: $y = Ax = \sum_{i \in \text{supp}(x)} a_i x_i$.

Intuition: Recovering x is “easier” if the a_i are not too similar...



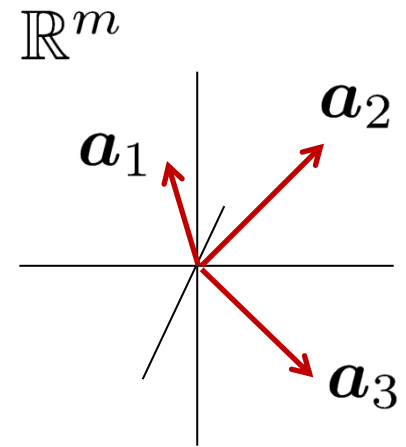
Mutual coherence $\mu(A) \doteq \max_{i \neq j} |\langle a_i, a_j \rangle|$.

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

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

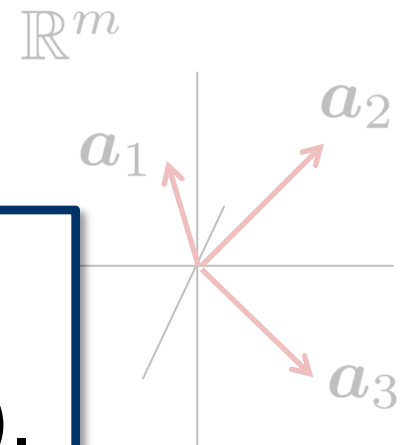
Then \mathbf{x}_0 is the unique optimal solution to

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

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

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

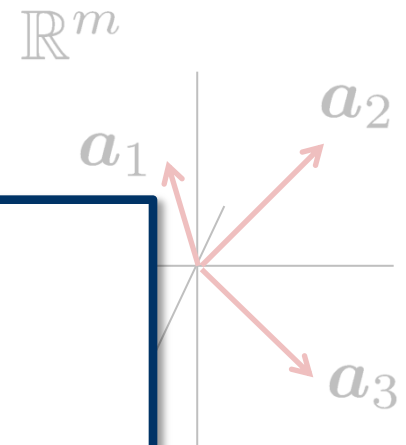
Then \mathbf{x}_0 is the unique optimal solution to

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

WHY DOES THIS WORK? More formally...

Mutual coherence

The matrix A is **incoherent** – and so, preserves sparse x .



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

Suppose $y = Ax_0$ with

$$\|x_0\|_0 < \frac{1}{2}(1 + 1/\mu(A)).$$

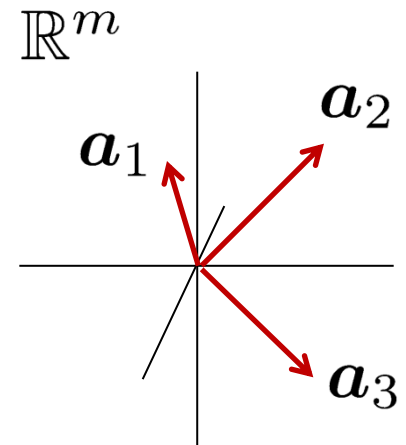
Then x_0 is the unique optimal solution to

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

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

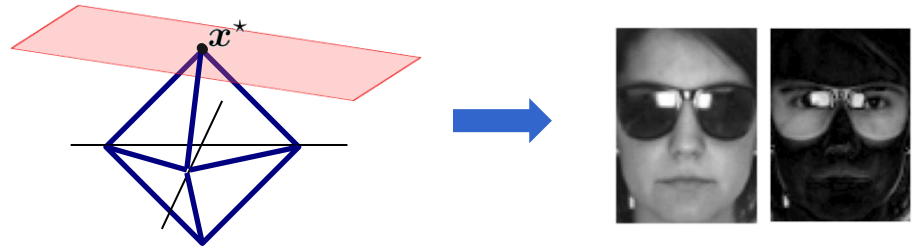
$$\|\mathbf{x}_0\|_0 < \frac{1}{2}(1 + 1/\mu(\mathbf{A})).$$

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?

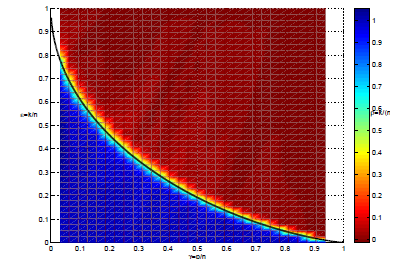
Motivates applications



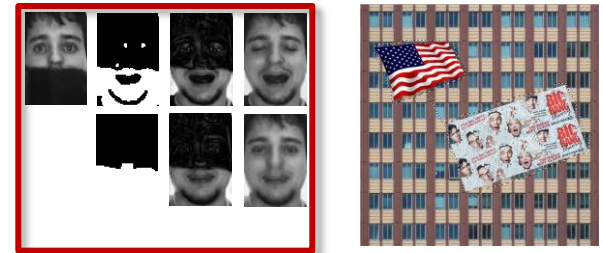
... but be careful: sometimes need to modify basic formulation [Lecture 3].

Template for stronger results

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

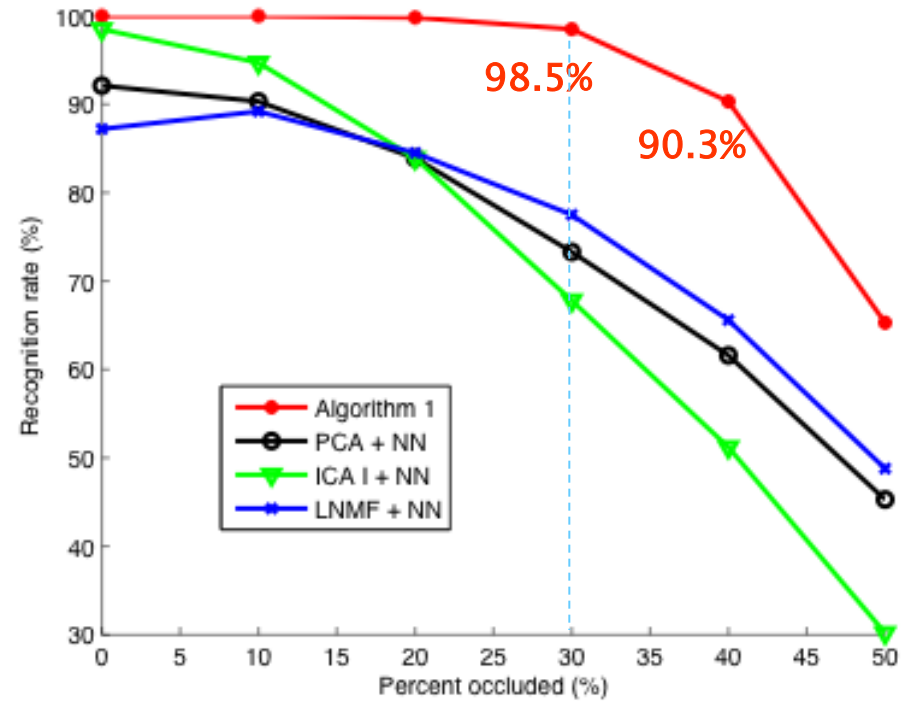
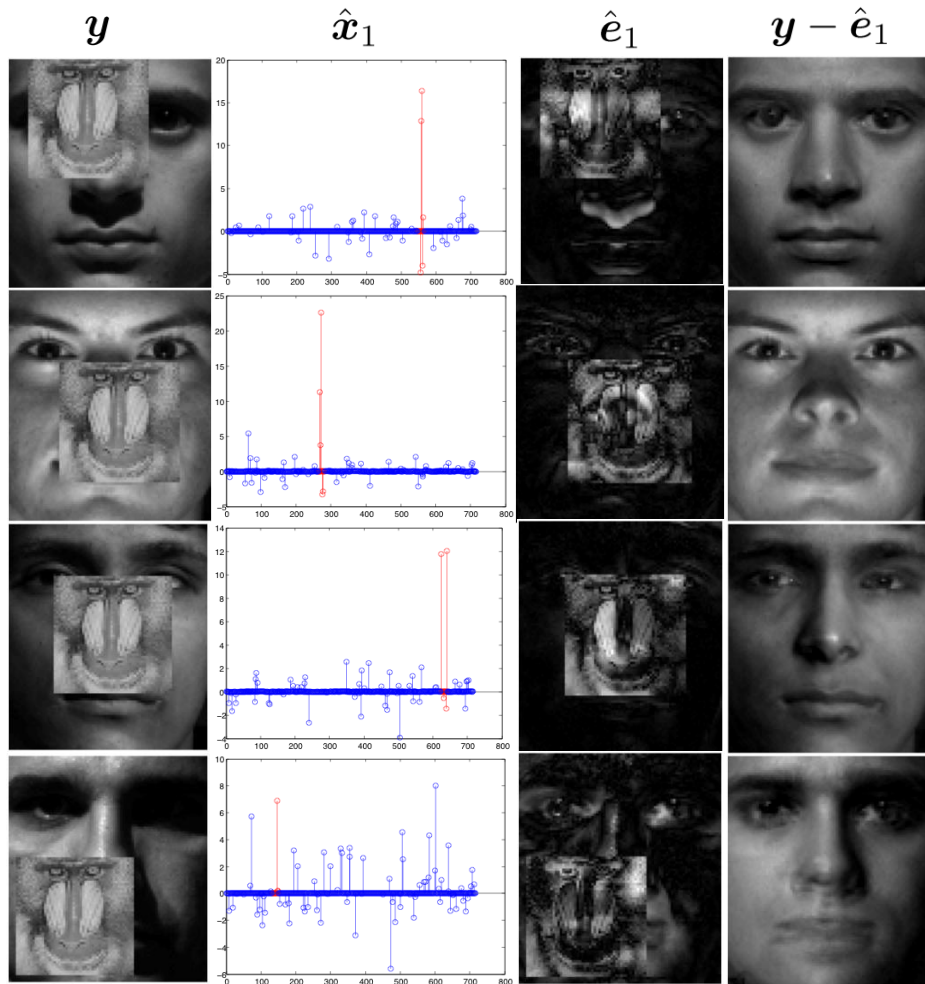


Generalizes to many other types
of low-dimensional structure

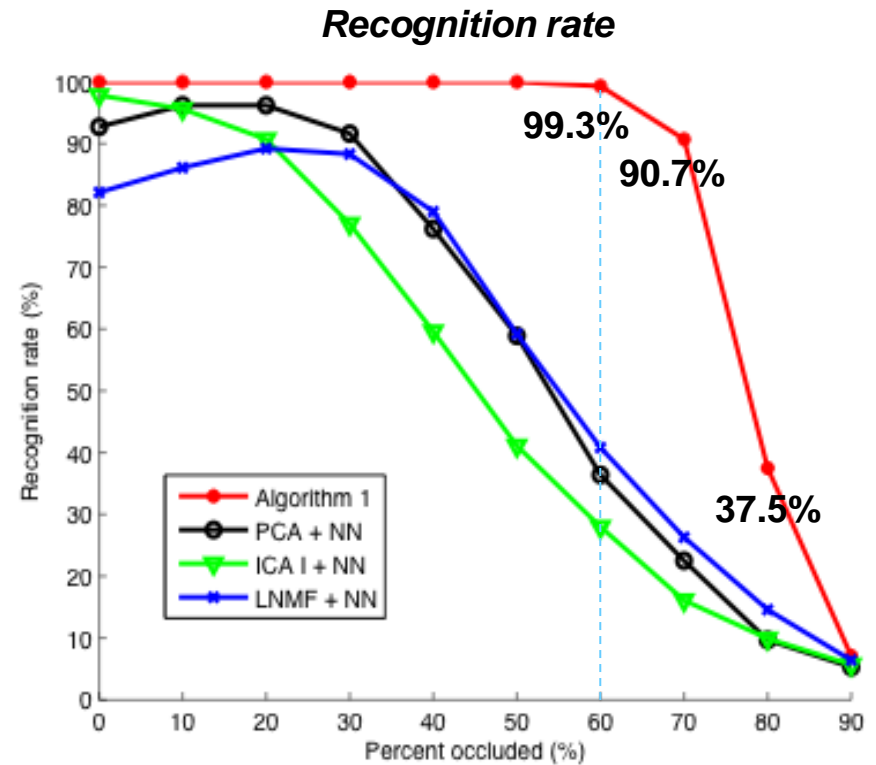
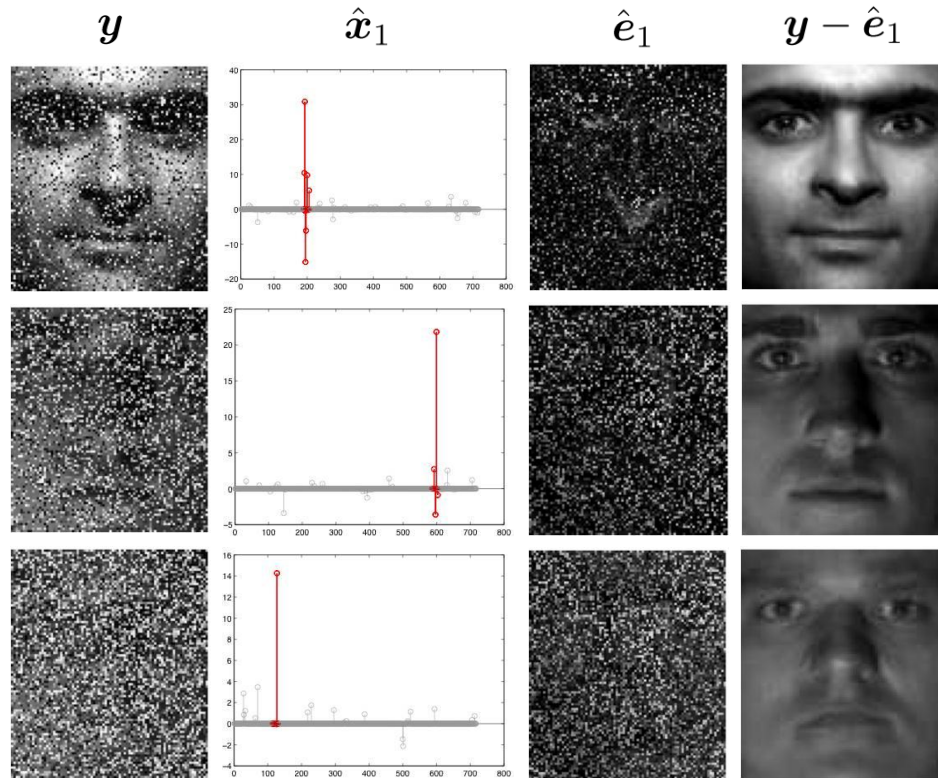


... structured sparsity [Lecture 2], low-rank recovery [later, Lecture 3].

THEORY TO APPLICATION – Face Recognition



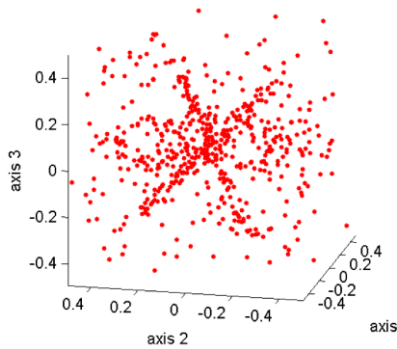
THEORY TO APPLICATION – Face Recognition



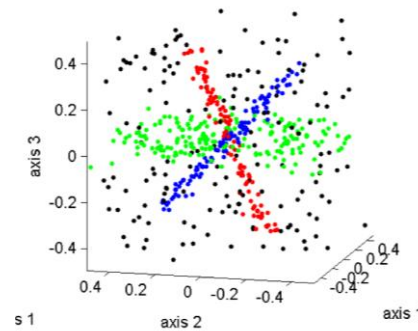
More practicalities in Lecture 3...

THEORY TO APPLICATION – Subspace Segmentation

Given $\mathbf{Y} = [\mathbf{y}_1 \mid \cdots \mid \mathbf{y}_N] \subset \mathcal{S}_1 \cup \cdots \cup \mathcal{S}_k$, determine (\mathcal{S}_i) .



Data \mathbf{Y}



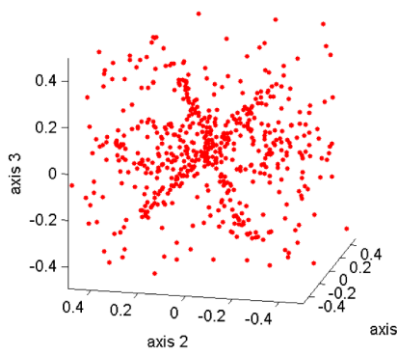
Segmentation

Applications include image segmentation, motion segmentation, hybrid system identification, and more.

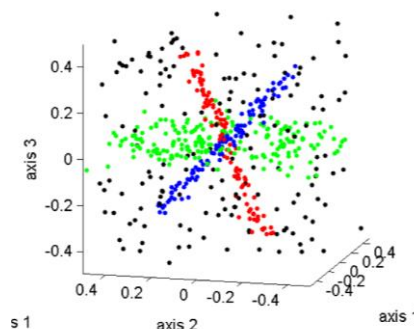


THEORY TO APPLICATION – Subspace Segmentation

Given $\mathbf{Y} = [\mathbf{y}_1 \mid \cdots \mid \mathbf{y}_N] \subset \mathcal{S}_1 \cup \cdots \cup \mathcal{S}_k$, determine (\mathcal{S}_i) .



Data \mathbf{Y}

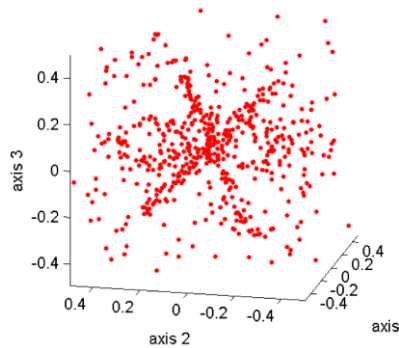


Segmentation

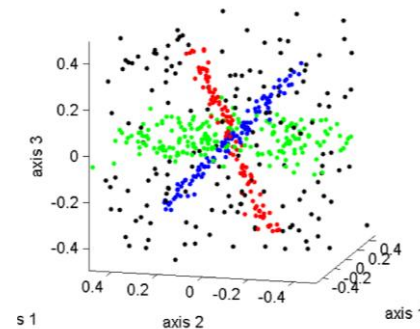
Each $\mathbf{y}_i \in \mathcal{S}_j$ can be expressed as a linear combination of just $\dim(\mathcal{S}_j)$ other points $\mathbf{y}'_i \in \mathcal{S}_j$.

THEORY TO APPLICATION – Subspace Segmentation

Given $\mathbf{Y} = [\mathbf{y}_1 \mid \cdots \mid \mathbf{y}_N] \subset \mathcal{S}_1 \cup \cdots \cup \mathcal{S}_k$, determine (\mathcal{S}_i) .



Data \mathbf{Y}



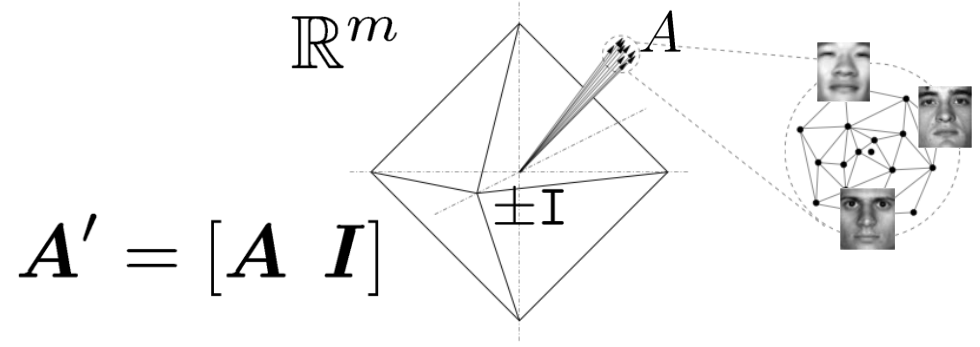
Segmentation

Each $\mathbf{y}_i \in \mathcal{S}_j$ can be expressed as a linear combination of just $\dim(\mathcal{S}_j)$ other points $\mathbf{y}'_i \in \mathcal{S}_j$.

minimize $\|\mathbf{X}\|_1$ subject to $\mathbf{Y}\mathbf{X} = \mathbf{Y}$, $\text{diag}(\mathbf{X}) = 0$.

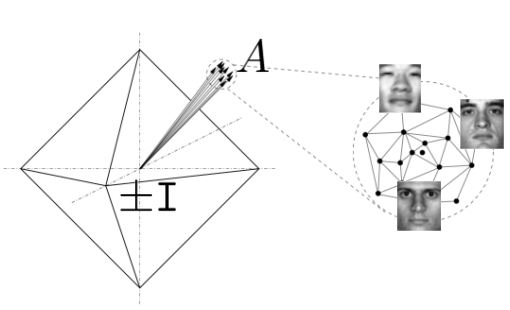
THEORY AND PRACTICE – Faces and Subspaces

In both applications, A can be coherent...

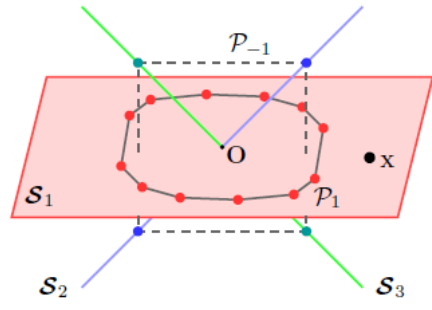


... ℓ^1 still exhibits **exact recovery** when x_0 is structured.

Theory extends by considering **problem-specific geometry**



[W.+Ma '10]

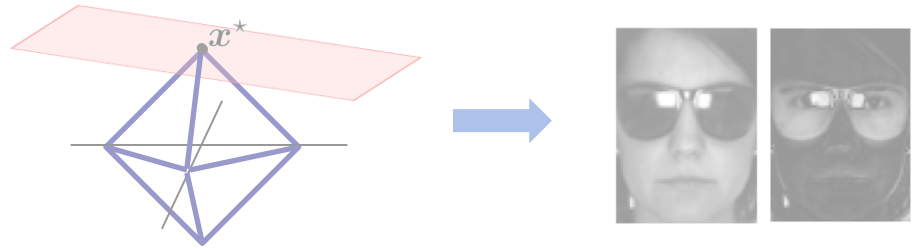


[Elhamifar+Vidal, '12]

[Soltanolkotabi and Candes, '11].

WHY CARE ABOUT THE THEORY?

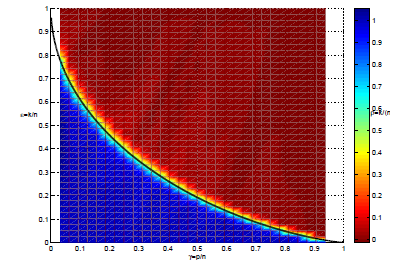
Motivates applications



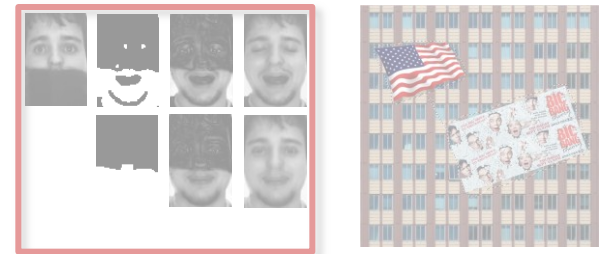
... but be careful: sometimes need to modify basic formulation [Lecture 3].

Template for stronger results

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



Generalizes to many other types
of low-dimensional structure



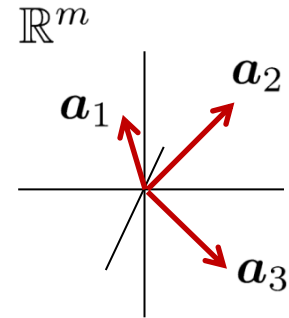
... structured sparsity [Lecture 2], low-rank recovery [later, Lecture 3].

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

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

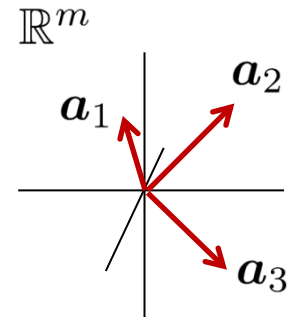


LIMITATIONS OF COHERENCE?

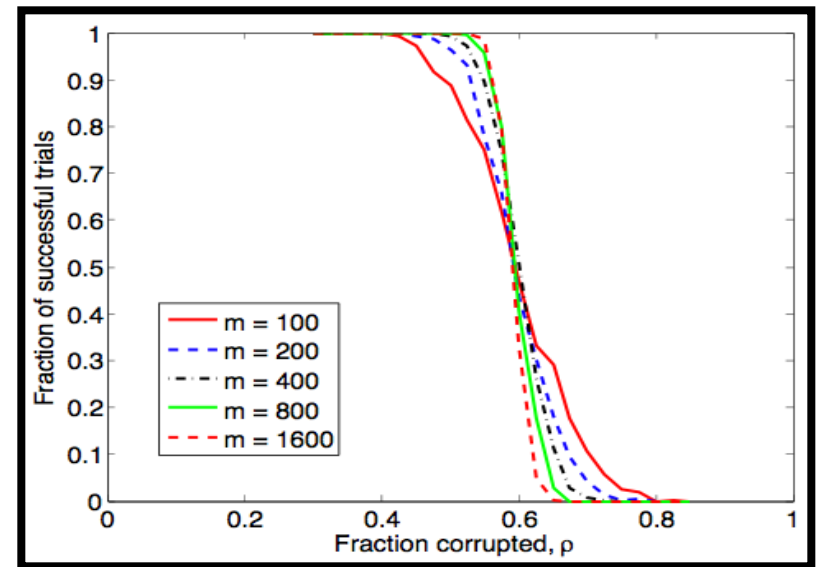
For any $m \times n$ \mathbf{A} , $\mu(\mathbf{A}) \geq \sqrt{\frac{n-m}{m(n-1)}}$.

Prev. result therefore requires

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



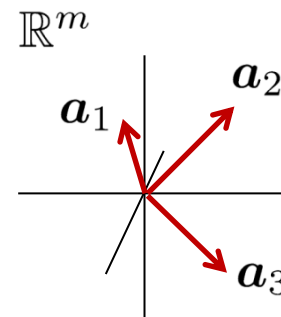
Truth is often **much better**:



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

LIMITATIONS OF COHERENCE?

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



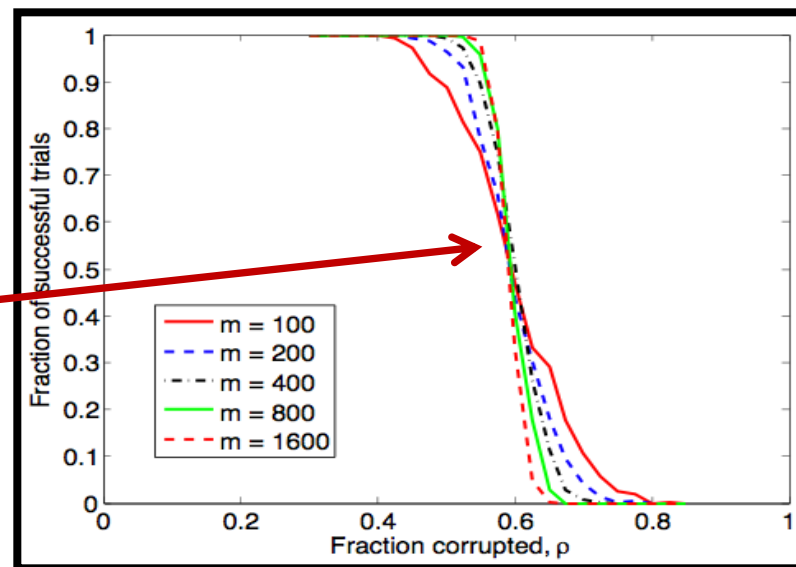
Prev. result therefore requires

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

Truth is often **much better**:

Phase transition at

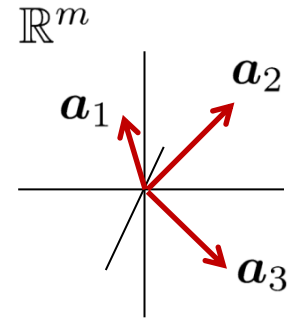
$$\|\mathbf{x}_0\|_0 = \alpha^* m$$



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

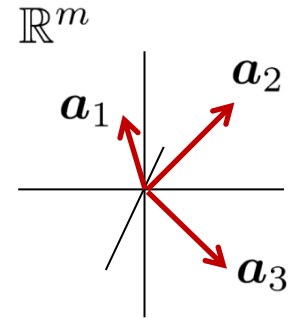
STRENGTHENING THE BOUND – the RIP

Incoherence: **Each pair** $A_{i,j} = [a_i \mid a_j]$ spread.



STRENGTHENING THE BOUND – the RIP

Incoherence: Each pair $A_{i,j} = [a_i \mid a_j]$ spread.



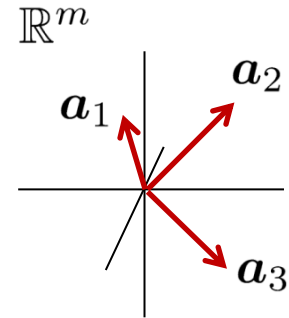
Generalize to **subsets of size k** :

A_I well-spread (almost orthonormal) for all I of size k

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

STRENGTHENING THE BOUND – the RIP

Incoherence: Each pair $A_{i,j} = [a_i \mid a_j]$ spread.



Generalize to **subsets of size k** :

A_I well-spread (almost orthonormal) for all I of size k

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

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

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

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:

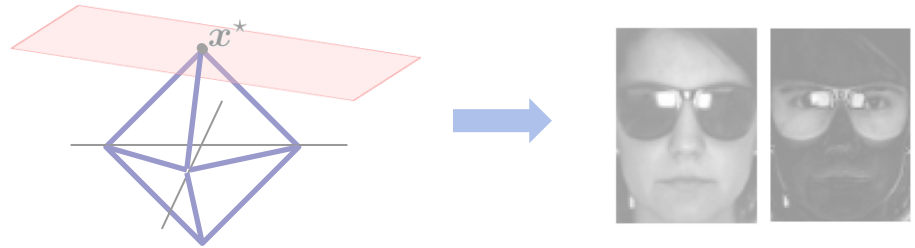
If $\mathbf{A} \sim_{iid} \mathcal{N}(0, m^{-1/2})$, then \mathbf{A} has RIP of order k with high probability, when $m \geq Ck \log(n/k)$.

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?

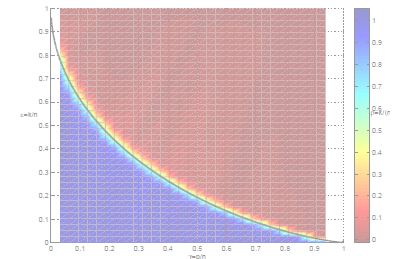
Motivates applications



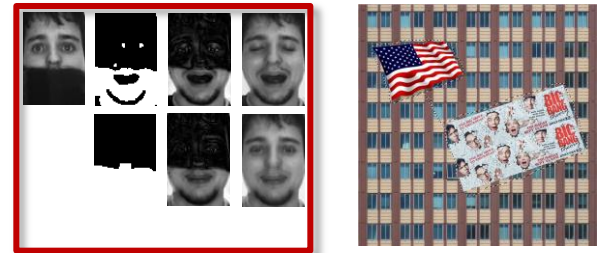
... but be careful: sometimes need to modify basic formulation [Lecture 3].

Template for stronger results

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



Generalizes to many other types
of low-dimensional structure



... structured sparsity [Lecture 2], low-rank recovery [later, Lecture 3].

GENERALIZATIONS – From Sparse to Low-Rank

So far: Recovering *a single sparse vector*:

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

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

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

FORMULATION – Robust PCA?

$$\begin{bmatrix} \text{Face with sunglasses} & \dots & \text{Noise} \end{bmatrix} = \begin{bmatrix} \text{Face with glasses} & \dots & \text{Face with glasses} \end{bmatrix} + \begin{bmatrix} \text{Face with mask} & \dots & \text{Noise} \end{bmatrix}$$

Y X E

Given $Y = X + E$, with X low-rank, E sparse, recover 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 $Y = X + Z$, recover X .

Stable, efficient algorithm, theoretically optimal → 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 $Y = \mathcal{P}_\Omega[X]$, recover X .

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

Our problem, with $Y = X + E$, looks more difficult...

WHY IS THE PROBLEM HARD?

Some very sparse matrices are also low-rank:

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

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

Can we recover X that are *incoherent* with the standard basis?

Certain sparse error patterns E make recovering X impossible:

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

X $E = e_i v^*$ $Y = X + E$

Can we correct E whose support is not *adversarial*?

WHEN IS THERE HOPE? Again, (in)coherence

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

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

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

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

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


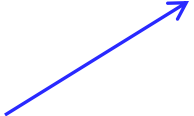
$$\text{support}(E) \sim \text{uni} \left(\begin{matrix} [m] \times [n] \\ \rho mn \end{matrix} \right)$$

... 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

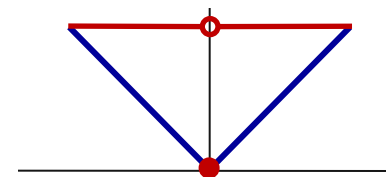
$$\text{rank}(\mathbf{X}) = \#\{\sigma_i(\mathbf{X}) \neq 0\}. \quad \|\mathbf{E}\|_0 = \#\{\mathbf{E}_{ij} \neq 0\}.$$



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



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



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

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*

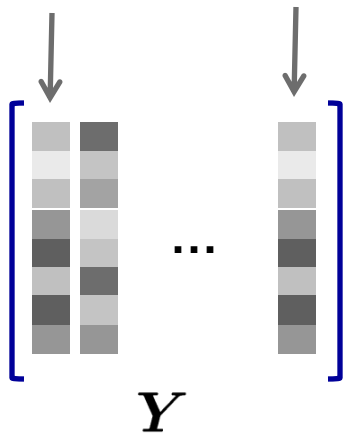
$$(\mathbf{X}_0, \mathbf{E}_0) = \arg \min \|\mathbf{X}\|_* + \frac{1}{\sqrt{m}} \|\mathbf{E}\|_1 \quad \text{subj} \quad \mathbf{X} + \mathbf{E} = \mathbf{X}_0 + \mathbf{E}_0,$$

and the minimizer is unique.

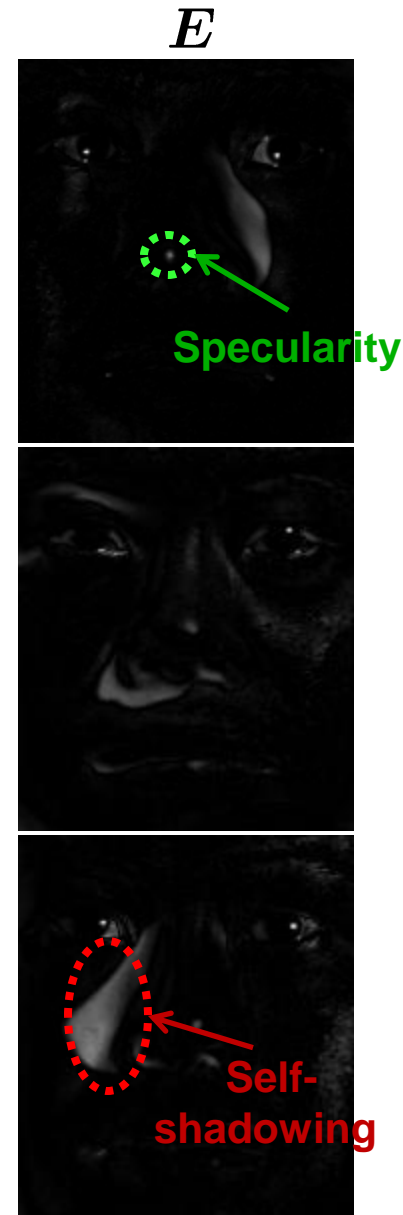
“Convex optimization recovers matrices of rank $O\left(\frac{n}{\log^2 m}\right)$ from errors corrupting $O(mn)$ entries”

EXAMPLE – Faces under varying illumination

58 images of one person under varying lighting:



RPCA →

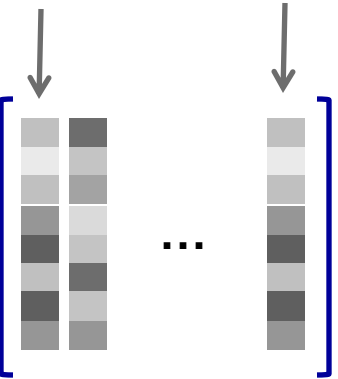


APPLICATIONS – Background modeling from video

Static camera
surveillance video

200 frames,
144 x 172 pixels,

Significant foreground
motion



Y

$RPCA$

$$\text{Video } Y = \text{Low-rank appx. } X + \text{Sparse error } E$$



BIG PICTURE – Parallelism of Sparsity and Low-Rank

	<i>Sparse Vector</i>	<i>Low-Rank Matrix</i>
Degeneracy of	one 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	$Y = A(X) + B(E) + Z$	

A SUITE OF POWERFUL REGULARIZERS ...

... for recovering various types of low-dimensional structure:

- [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**

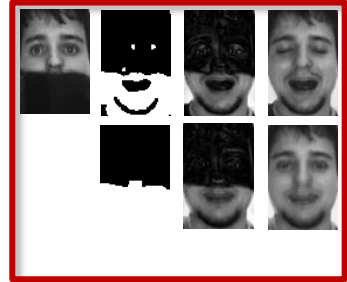
$$X^0 = \arg \min \|X\|_{\diamond} \quad \text{s.t.} \quad \mathcal{P}_Q(X) = \mathcal{P}_Q(X^0)$$

- [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, ISIT'12] – **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 / application tell you the structure...

THANK YOU!

Questions, please?

Next: Efficient, **scalable algorithms** for large ℓ^1 and nuclear norm problems

Later: More **applications** of these techniques

A FEW REFERENCES

General surveys:

Donoho, Elad and Bruckstein, *From Sparse Solutions of Systems of Equations to Sparse Modeling of Signals and Images*, SIAM Review '09 (see also Elad's book).

Davenport, Duarte, Eldar, Kityniok, *Introduction to Compressed Sensing*, Signal Processing Magazine '11

W., Ma, Sapiro, Mairal, Huang, Yan, *Sparse Representation for Computer Vision and Pattern Recognition*, Proc. IEEE '10

Hardness of L0 minimization:

Natarajan, *Sparse Approximate Solutions to Linear Systems*, SIAM Journal on Computing '95

Amaldi and Kann, *On the Approximability of Minimizing Nonzero variables or Unsatisfied Relations in Linear Systems*, Theoretical Computer Science '97

Uniqueness of sparse solutions:

Donoho and Elad, *Optimally Sparse Representations in General (nonorthogonal) Dictionaries via L1 Minimization*, PNAS '03

Gorodnitsky and Rao, *Sparse Signal Reconstruction from Limited Data using FOCUSS – A Reweighted Minimum Norm Algorithm*, IEEE Trans. Signal Processing '97

The L1 relaxation:

Tibshirani, *Regression shrinkage and selection via the LASSO*, Journal of the Royal Statistical Society Series B, '96

Chen, Donoho and Saunders, *Atomic Decomposition by Basis Pursuit*, SIAM Review '98

(In)-coherence and recovery guarantees:

Donoho and Elad, *Optimally Sparse Representations in General (nonorthogonal) Dictionaries via L1 Minimization*, PNAS '03

Gribonval and Nielsen, *Sparse Representations in Unions of Bases*, IEEE Info Thy '04

Fuchs, *On Sparse Representations in Arbitrary Redundant Bases*, IEEE Info Thy '04

A FEW MORE REFERENCES

RIP Based Guarantees for L1 Minimization

Candes and Tao, *Near Optimal Signal Recovery from Random Projections: Universal Encoding Strategies?* IEEE Trans Info Thy, '04

Candes and Tao, *Decoding by Linear Programming*, IEEE Trans Info Thy, '05

Candes, *The Restricted Isometry Property and Its Implications for Compressed Sensing*, '08

A Few Related Results for Sparse Recovery (not covered in lecture)

Tropp, *Just Relax: Convex Programming Methods for Recovering Sparse Signals in Noise*, IEEE Info Thy '06

Wainwright, *Sharp Thresholds for Noisy and High-Dimensional Recovery of Sparsity Using L1 Constrained Quadratic Programming (LASSO)*, IEEE Info Thy '09

Donoho and Tanner, *Counting Faces of Randomly Projected Polytopes when Projection Radically Lowers Dimension*, Journal of the AMS '09

Bayati, Lelarge and Montanari, *Universality of Polytope Phase Transitions and Message Passing Algorithms*, '12

A Few Applications in this Lecture (more in Lecture 3)

Lustig, Donoho, Santos and Pauly, *Compressive Sensing MRI*, Magnetic Resonance Medicine, '07

Horev, Bryt and Rubinstein, *Adaptive Image Compression using Sparse Dictionaries*, '12

W., Yang, Ganesh, Sastry, Ma, *Robust Face Recognition via Sparse Representation*, IEEE PAMI '09

Elhamifar and Vidal, *Sparse Subspace Clustering*, CVPR '09

A FEW *MORE* REFERENCES

Classical Matrix Approximations (PCA)

Pearson, *On lines and planes best fit to systems of points in space*, Philosophical Magazine 1901

Hotelling, *Analysis of a complex of statistical variables into principal components*, Journal of Educational Psychology, 1933

Nuclear Norm Minimization

Fazel, Hindi, Boyd, *A Rank Minimization Heuristic with Application to Minimum-Order System Identification*, ACC '01

Recht, Fazel, Parillo, *Guaranteed Minimum rank Solutions to Linear Matrix Equations via Nuclear Norm Minimization*, SIAM Review '10

Matrix Completion

Candes and Recht, *Exact Matrix Completion via Convex Optimization*, Foundations of Computational Mathematics '09

Gross, *Recovering Low-Rank Matrices from Few Coefficients in Any Basis*, IEEE Trans. Info. Thy., '10

Robust PCA and Matrix Decompositions

Candes, Li, Ma, W., *Robust Principal Component Analysis?* Journal of the ACM '11

Chandrasekaran, Sanghavi, Pararilo and Wilsky, *Rank-Sparsity Incoherence for Matrix Decomposition*, SIAM Journal on Optimization, '11

Agarwal, Negahban and Wainwright, *Noisy Matrix Decomposition via Convex Relaxation: Optimal Rates in High Dimensions* '12

This is a huge (and hugely active) area! Many important ideas and papers are missing from the above list. You can complement your reading by visiting ...

The UIUC Matrix Recovery site: <http://perception.csl.illinois.edu/matrix-rank/home.html>

The Rice Compressed Sensing Archive: <http://dsp.rice.edu/cs>

Nuit Blanche (a blog in this area): <http://nuit-blanche.blogspot.com/>