

# **Effective Representations for High-Dimensional Visual Data**

**John Wright**

**Electrical Engineering**

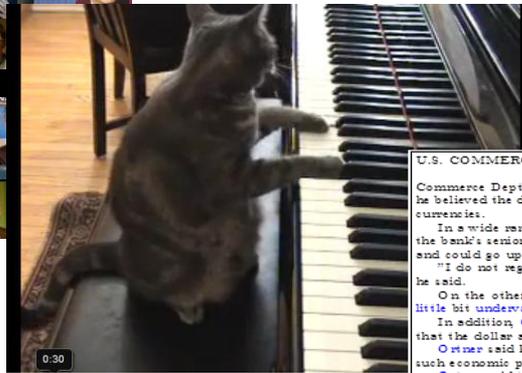
**Columbia University**

# Datasets are massive, high-dimensional...



Images

> 1M pixels



Videos

> 1B voxels

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 because of the decline 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 leveling off.

Turning to Brazil and Mexico, Ortner made it clear that it was 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 the policies outlined in Treasury Secretary James Baker's debt



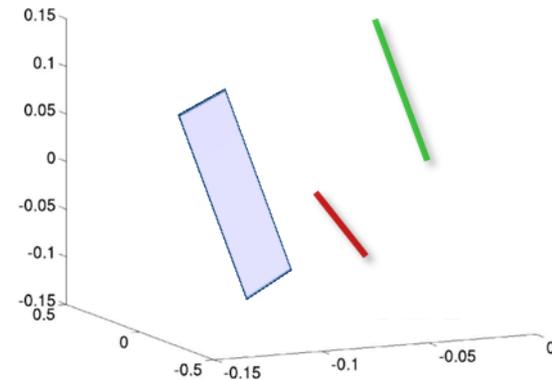
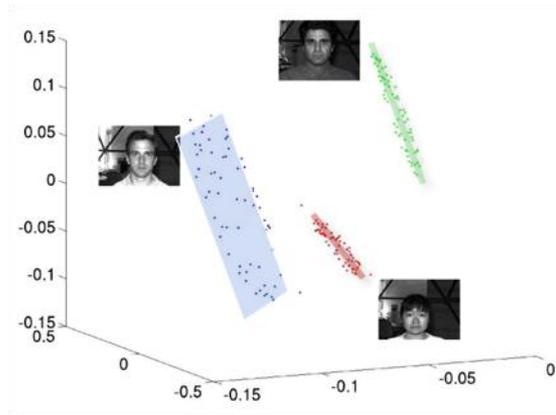
★	★		★
★★		??	
★★		★	★

Web data

> 10B+ websites

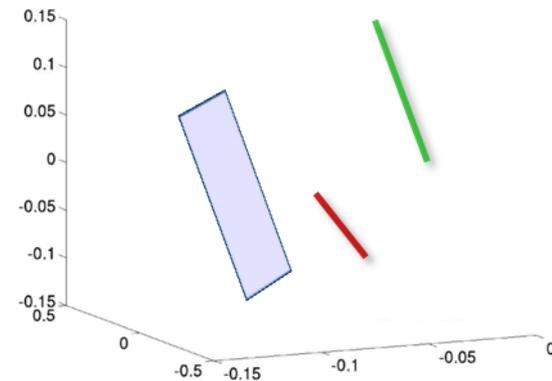
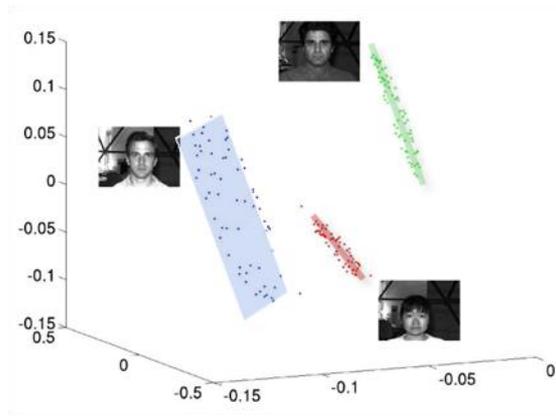


# ... intrinsic structures are low-dimensional



*How can we model low-dimensional structure  
in high-dimensional data?*

# ... intrinsic structures are low-dimensional



*How can we model low-dimensional structure in high-dimensional data?*

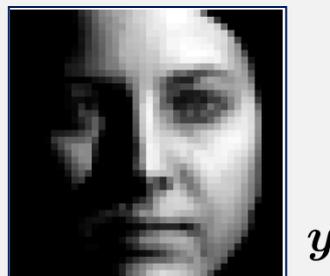
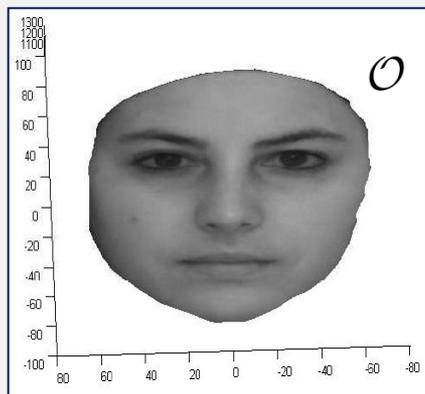
- 1. Analytically?*
- 2. Automatically?*

*... and what does this buy us for engineering applications?*

# The space/model of images of an object?

## Conceptual Problem:

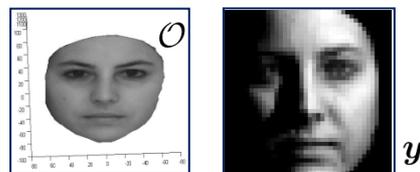
Given (information about) an object  $\mathcal{O}$ , can we *provably* detect or recognize the object from a new image  $y$  ?



## This talk:

Attempt this for *fixed pose*, variations in *illumination*, (and possibly *occlusion*).

# The goal – practical?



Many **engineering approaches** to mitigate illumination:

## *Quotient images*

[Riklin-Raviv, Shashua '99] [Wang, Li, Wang '04]

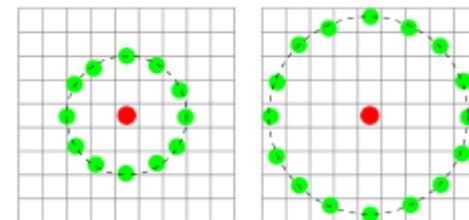


Wang, Li, Wang '04

## *Nonlinear features*

SIFT [Lowe '99], HOG [Dalal, Triggs '05]

LBP [Ojala, Pietikäinen, Harwood '94,  
Ahonen, Hadid, Pietikainen '06]



Ahonen, Hadid, Pietikainen '06

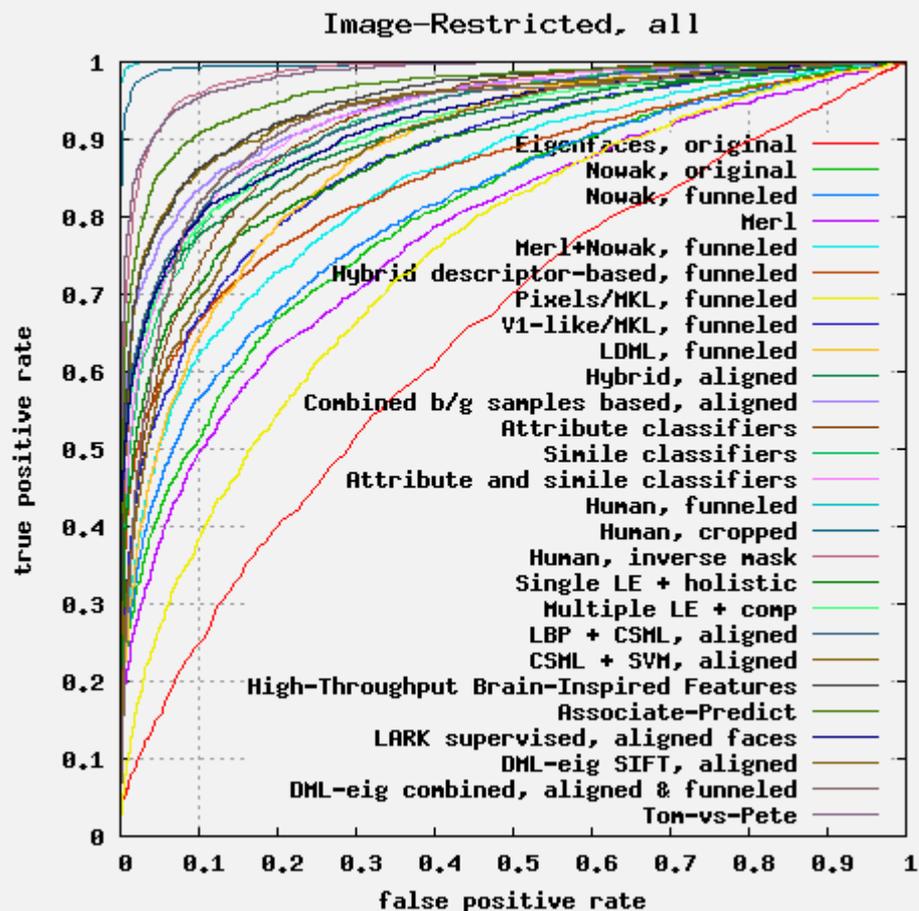
## *Total variation minimization*

[Yin et. Al. ]

## Or use **heuristic / intuitive physics** ...

[Murase, Nayar '96], [Belhumeur , Kriegman '98], [Basri, Jacobs '03],  
[Ramamoorthi '01], [Basri, Frolova '04], [ W. et. al. '09], many others...

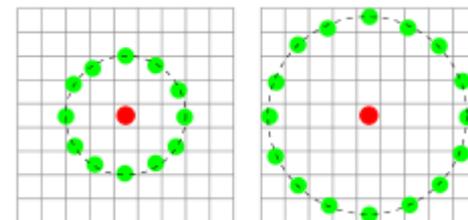
# The goal – practical?



ation:



Wang, Li, Wang '04



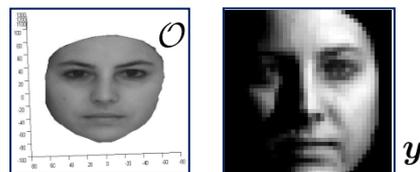
Ahonen, Hadid, Pietikainen '06

Labeled Faces in the Wild results

<http://vis-www.cs.umass.edu/lfw/results.html>

Jacobs '03],  
any others...

# The goal – practical?



Many **engineering approaches** to mitigate illumination:

## *Quotient images*

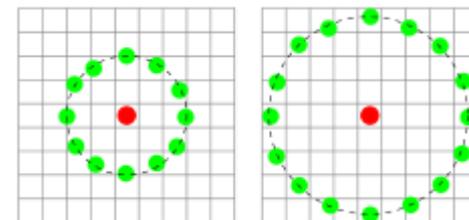
[Riklin-Raviv, Shashua '99] [Wang, Li, Wang '04]



Wang, Li, Wang '04

## *Nonlinear features*

SIFT [Lowe '99], HOG [Dalal, Triggs '05]  
LBP [Ojala, Pietikäinen, Harwood '94,  
Ahonen, Hadid, Pietikainen '06]



Ahonen, Hadid, Pietikainen '06

## *Total variation minimization*

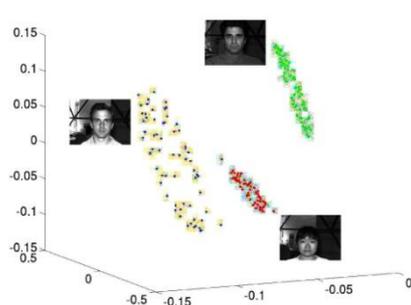
[Yin et. Al. ]

## Or use **heuristic / intuitive physics** ...

[Murase, Nayar '96], [Belhumeur , Kriegman '98], [Basri, Jacobs '03],  
[Ramamoorthi '01], [Basri, Frolova '04], [ W. et. al. '09], many others...

# Ex: Face recognition with $\ell^1$ -regression:

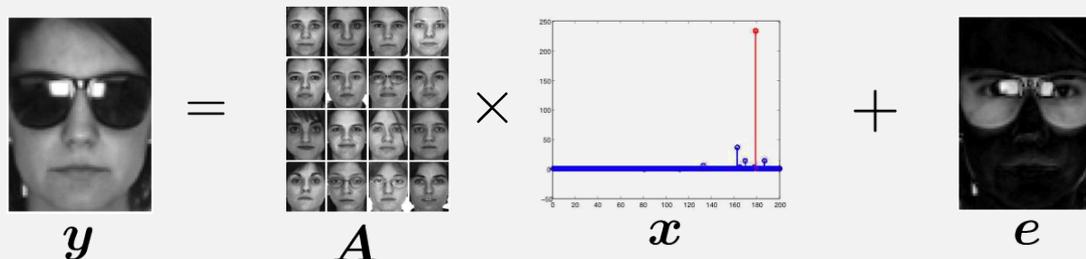
Linear subspace model for images of same face under varying illumination:



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

If test image  $\mathbf{y} \in \mathbb{R}^m$  is also of subject  $i$ , then  $\mathbf{y} = A_i \mathbf{x}_i$  for some  $\mathbf{x}_i \in \mathbb{R}^k$ .

Given  $\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{e}_0$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m \ll n$ , recover  $\mathbf{x}_0$  and  $\mathbf{e}_0$ .



# Ex: Face recognition with $\ell^1$ -regression:

*Underdetermined* system of linear equations in unknowns  $\mathbf{x}$ ,  $\mathbf{e}$ :

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e} = \begin{bmatrix} \mathbf{A} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{e} \end{bmatrix} \doteq \mathbf{B}\mathbf{w} \quad \mathbf{B} \in \mathbb{R}^{m \times m+n} \quad \mathbf{w} \in \mathbb{R}^{m+n}$$

Solution is not unique ... but

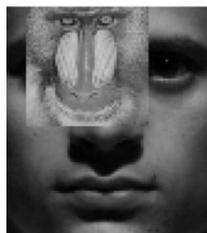
$\mathbf{x}$  should be *sparse*: ideally, only supported on images of the same subject  
 $\mathbf{e}$  expected to be *sparse*: occlusion only affects a subset of the pixels

Seek the *sparsest* solution:

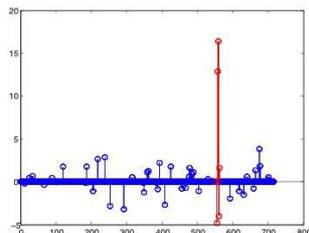
$$\min \|\mathbf{x}\|_0 + \|\mathbf{e}\|_0 \quad \text{subj } \mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}$$

$\Downarrow$  convex relaxation  $\Downarrow$

$$\min \|\mathbf{x}\|_1 + \|\mathbf{e}\|_1 \quad \text{subj } \mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}$$



$\mathbf{y}$



$\hat{\mathbf{x}}_1$



$\hat{\mathbf{e}}_1$

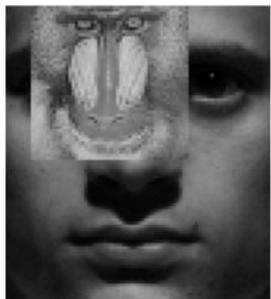


$\hat{\mathbf{y}} = \mathbf{y} - \hat{\mathbf{e}}_1$

# Ex: Face recognition with $\ell^1$ -regression:

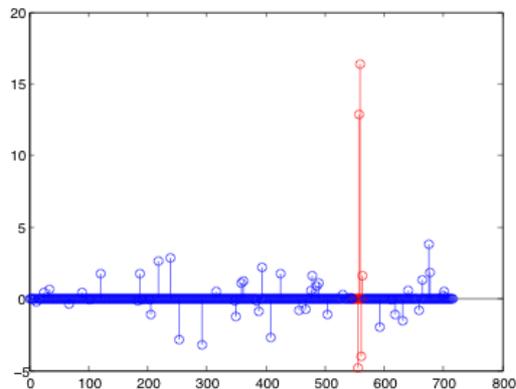
$$\hat{x}_1 = \arg \min \|x\|_1 + \|e\|_1.$$

Input:  $y \in \mathbb{R}^D$

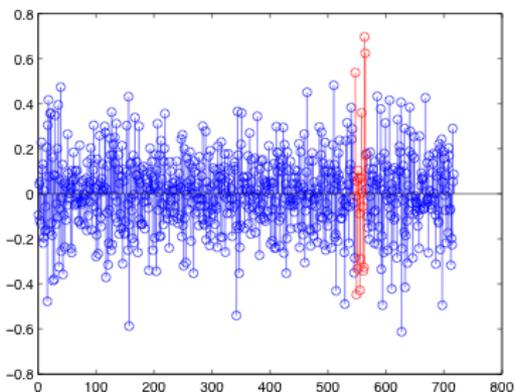


$$y = Ax + e$$

$$\hat{x}_2 = \arg \min_x \|y - Ax\|_2.$$



$$\hat{e}_1 \quad \hat{y}_0 = A\hat{x}_1$$



$$\hat{e}_2 \quad \hat{y}_0 = A\hat{x}_2$$



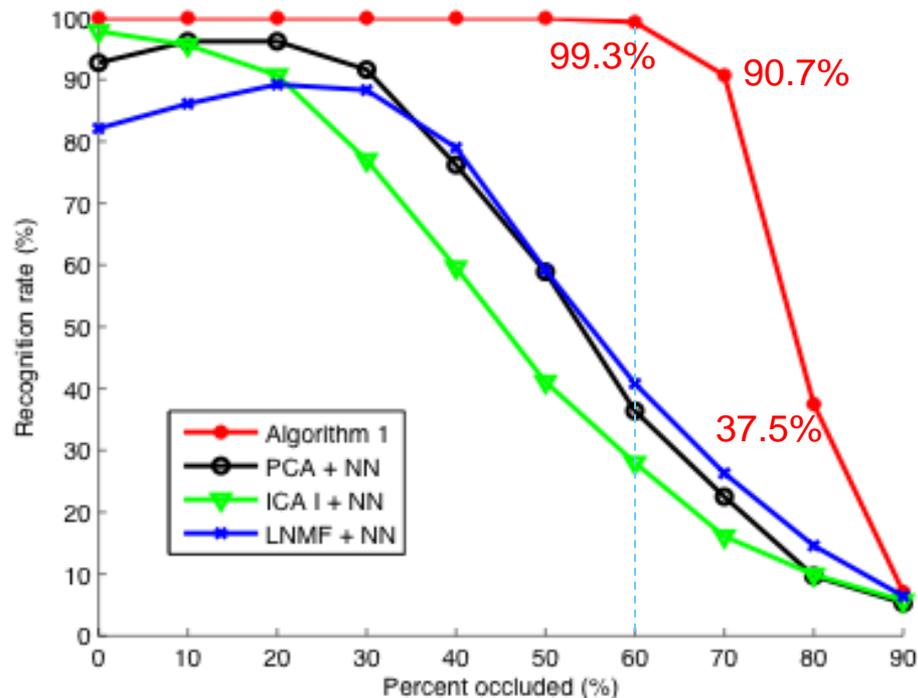
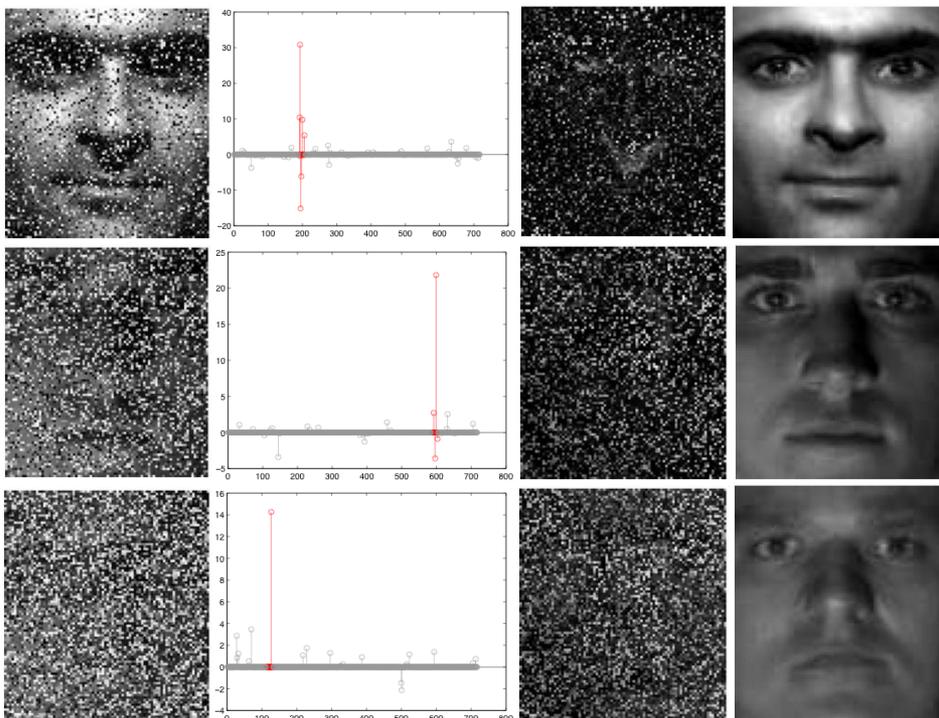
# Ex: Face recognition with $\ell^1$ -regression:

Extended Yale B Database  
(38 subjects)

**Training:** subsets 1 and 2 (717 images)

**Testing:** subset 3 (453 images)

$$y \quad \hat{x}_1 \quad \hat{e}_1 \quad \hat{y}_0 = A\hat{x}_1$$



# Ex: Face recognition with $\ell^1$ -regression:

**Theorem 1** (W. and Ma, '10). *For any  $\delta > 0, \rho < 1, \exists \nu_0(\delta, \rho) > 0$ , such that if  $\nu < \nu_0$  and  $\alpha < \alpha_0(\delta, \nu, \rho)$ , then with error support  $J$  and signs  $\sigma$  chosen uniformly at random,*

$$\lim_{m \rightarrow \infty} \mathbb{P}_{A, J, \sigma} [ \ell^1\text{-minimization recovers all } \alpha m\text{-sparse } \mathbf{x} ] = 1.$$

*If  $A$  is “nice” and the model  $\mathbf{y} = A\mathbf{x} + \mathbf{e}$  fits, can make strong statements about the performance of  $\ell^1$  regression.*

[W., Ma, Tr. IT '09]

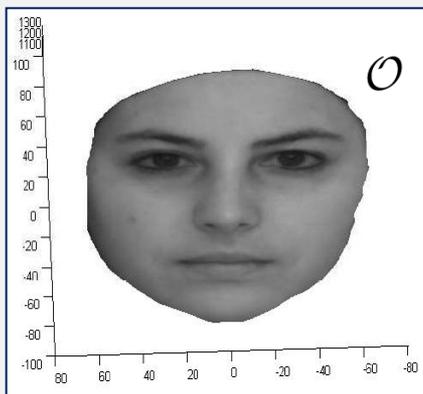
See also ... [Candes+Tao '05], [Li '11], [Ngyuen+Tran '12],

[McCoy+Tropp '13], [Foygel+Mackey '13]

# The space/model of face images of a person?

## Conceptual Problem:

Given (information about) an object  $\mathcal{O}$ , can we *provably* detect or recognize the object from a new image  $y$  ?



## This talk:

Attempt this for *fixed pose*, variations in *illumination*, (and possibly *occlusion*).

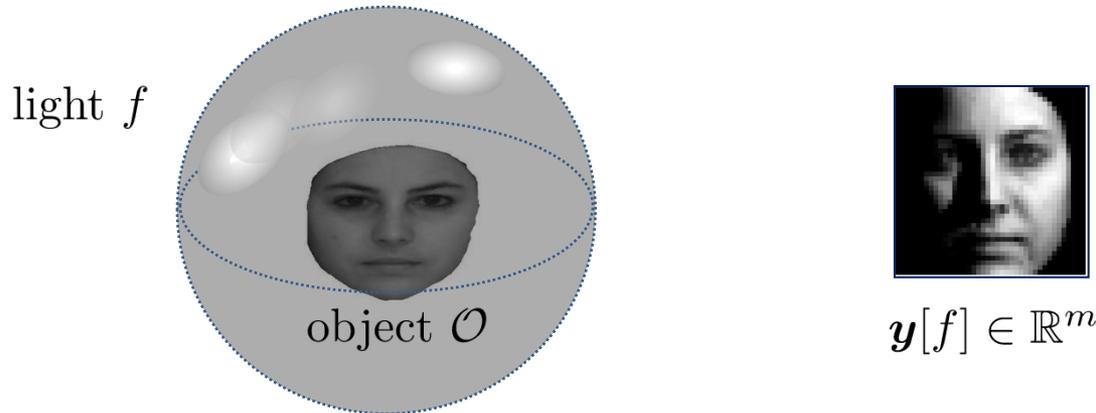
# Geometry of illumination variations

**Distant illumination** identified with a Riemann integrable, nonnegative function  $f : \mathbb{S}^2 \rightarrow \mathbb{R}_+$ .

Assume a **linear sensor response**. The **image**  $y[f]$  can often be written as

$$y[f] = \int_{\mathbf{u} \in \mathbb{S}^2} f(\mathbf{u}) \bar{\mathbf{y}}[\mathbf{u}] d\mathbf{u}$$

What is the set of possible images  $y[f]$  of  $\mathcal{O}$ ?



# Geometry of illumination variations

**Distant illumination** identified with a Riemann integrable, nonnegative function  $f : \mathbb{S}^2 \rightarrow \mathbb{R}_+$ .

Assume a **linear sensor response**. The **image**  $\mathbf{y}[f]$  can often be written as

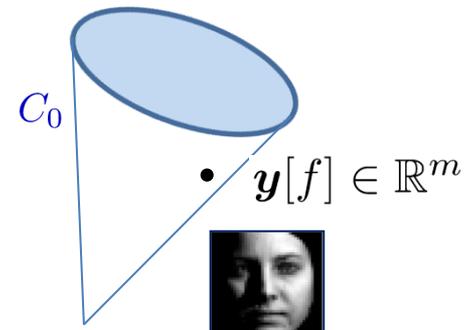
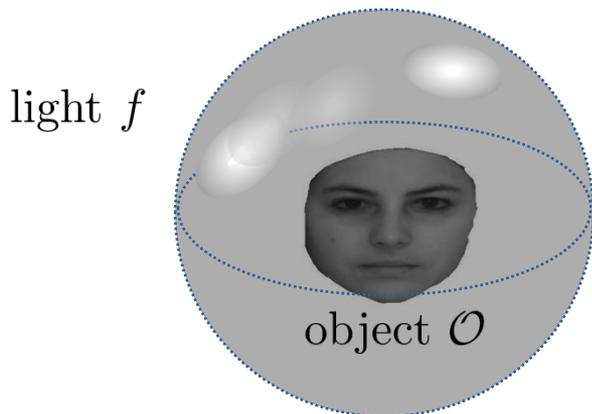
$$\mathbf{y}[f] = \int_{\mathbf{u} \in \mathbb{S}^2} f(\mathbf{u}) \bar{\mathbf{y}}[\mathbf{u}] d\mathbf{u}$$

What is the set of possible images  $\mathbf{y}[f]$  of  $\mathcal{O}$ ?

Let  $\mathcal{F}_0 \doteq \{f : \mathbb{S}^2 \rightarrow \mathbb{R}_+ \mid f \text{ Riemann integrable}\}$ , and

$$C_0 = \{\mathbf{y}[f] \mid f \in \mathcal{F}_0\} \subset \mathbb{R}^m.$$

*The set  $C_0$  is a convex cone.* [Belhumeur + Kriegman '98]

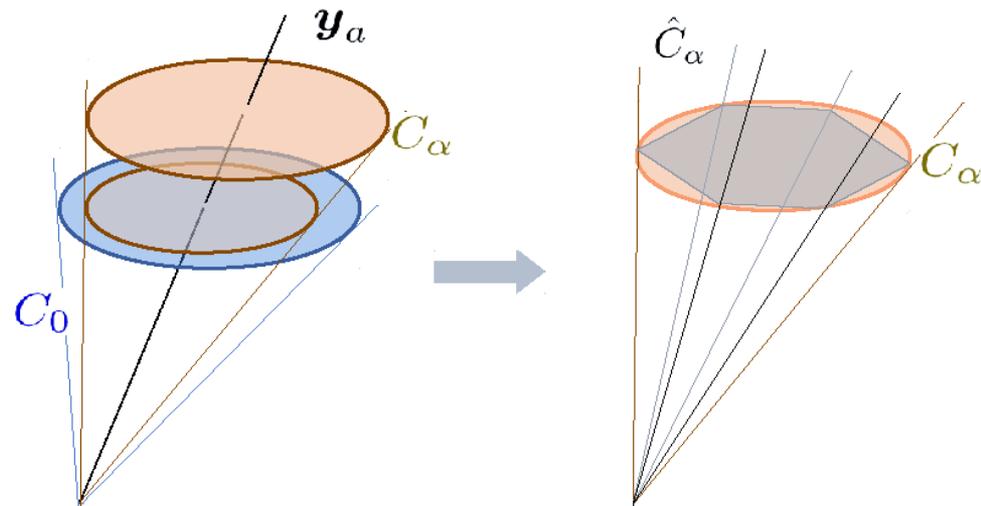
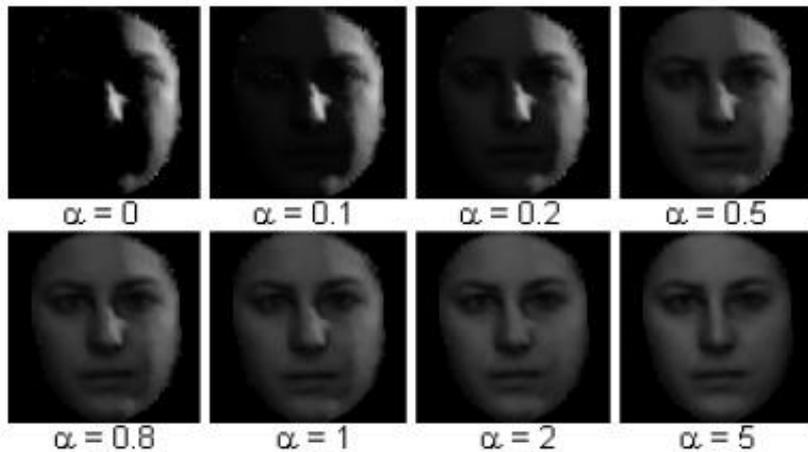


# Geometry of illumination variations

**Ambient cone model.** *Illumination is the sum of ambient and directional (arbitrary) components:*

$$\mathcal{F}_\alpha \doteq \{f_d + \alpha\omega \mid f_d : \mathbb{S}^2 \rightarrow \mathbb{R}_+, \text{Riemann integrable}, \|f_d\|_{L^1} \leq 1\}$$

*and corresponding cone of possible images:  $C_\alpha \doteq \text{cone}(\mathbf{y}[\mathcal{F}_\alpha])$*



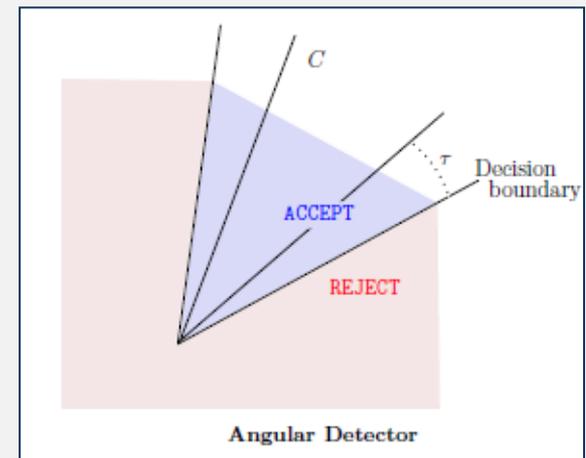
# Approximating a convex cone

**Angular detector:** is  $\mathbf{y}$  in or close to  $C_\alpha$ ?

*Accept any input  $\mathbf{y}$  that is an image of  $\mathcal{O}$  under some valid illumination.*

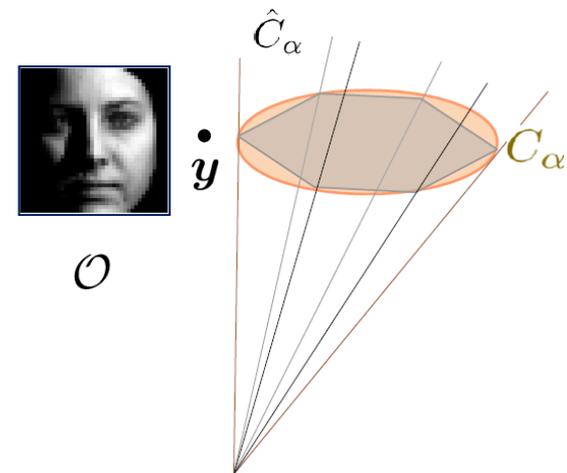
*Reject any input  $\mathbf{y}$  that is not.*

$$\mathcal{D}_\tau[\mathbf{y}] = \begin{cases} \text{ACCEPT} & \angle(\mathbf{y}, C) \leq \tau \\ \text{REJECT} & \angle(\mathbf{y}, C) > \tau \end{cases}$$



Work with a “Hausdorff distance” between cones:

$$\begin{aligned} \delta(C, C') &\doteq d_{\text{Hausdorff}}(C \cap B(0, 1), C' \cap B(0, 1)) \\ &= \max \left\{ \sup_{\mathbf{y} \in C \cap B(0, 1)} d(\mathbf{y}, C'), \sup_{\mathbf{y}' \in C' \cap B(0, 1)} d(\mathbf{y}', C) \right\} \end{aligned}$$



*If  $C'$  approximates  $C$  in Hausdorff sense, we lose little in working with  $C'$ .*

# Cone Approximation: Existing Results

Approximation of **general** convex bodies in high dimensions is a disaster:

**Theorem 2.** [Bronstein, Ivanov '76] *Let  $K \subset \mathbb{R}^m$  be a convex body. There exists an  $\varepsilon$ -approximation to  $K$  in Hausdorff distance, with  $O\left((\text{diam}(K)/\varepsilon)^{\frac{m-1}{2}}\right)$  vertices. For the unit sphere, this is optimal up to a constant.*

In vision literature, results for **average case**, for **convex objects**:

**Informal claim** [Basri and Jacobs '03, Basri and Frolova '04]:

*For an convex, Lambertian object  $\mathcal{O}$ , there exists a subspace  $\Gamma$ , with  $\dim(\Gamma) = 9$ , such that if  $\mathbf{u} \sim \text{uni}(\mathbb{S}^2)$  is a uniformly oriented random point source,*

$$\frac{\mathbb{E} \left[ \|\bar{\mathbf{y}}[\mathbf{u}] - \mathcal{P}_{\Gamma} \bar{\mathbf{y}}[\mathbf{u}]\|_2^2 \right]}{\mathbb{E} \left[ \|\bar{\mathbf{y}}[\mathbf{u}]\|_2^2 \right]} \leq .02$$

*Is there any special structure we can use here?*

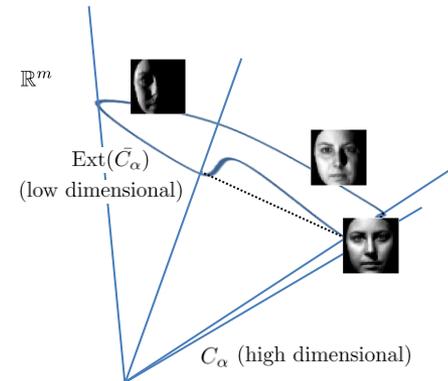
# Extreme rays

Image formation:  $y[f] = \int_{\mathbf{u} \in \mathbb{S}^2} f(\mathbf{u}) \bar{y}[\mathbf{u}] d\mathbf{u} .$

$\bar{y}[\mathbf{u}] \in \mathbb{R}^m$  are images under  
**directional illumination:**



*Extreme rays lie on a low dimensional submanifold of  $\mathbb{R}^m$  :*



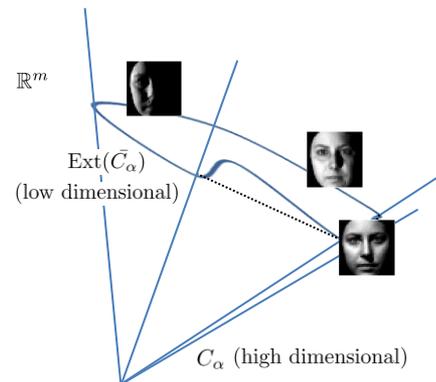
# Extreme rays

Image formation:  $\mathbf{y}[f] = \int_{\mathbf{u} \in \mathbb{S}^2} f(\mathbf{u}) \bar{\mathbf{y}}[\mathbf{u}] d\mathbf{u}$ .

$\bar{\mathbf{y}}[\mathbf{u}] \in \mathbb{R}^m$  are images under  
directional illumination:



*Extreme rays lie on a low dimensional submanifold of  $\mathbb{R}^m$ :*



**Lemma.** [Just approximate the extreme rays]. Let  $\mathbf{u}_1, \dots, \mathbf{u}_N \in \mathbb{S}^2$ . Then

$$\delta\left(C_\alpha, \text{cone}\{\bar{\mathbf{y}}[\mathbf{u}_i] + \alpha \mathbf{y}_a\}\right) \leq \frac{2}{\eta_\star \alpha \|\mathbf{y}_a\|} \times \sup_{\mathbf{u}} \min_i \|\bar{\mathbf{y}}[\mathbf{u}] - \bar{\mathbf{y}}[\mathbf{u}_i]\|_2,$$

where  $\eta_\star = \sup_{\|\mathbf{w}\|_2 \leq 1} \min_i \left\langle \mathbf{w}, \frac{\bar{\mathbf{y}}[\mathbf{u}_i]}{\|\bar{\mathbf{y}}[\mathbf{u}_i]\|_2} \right\rangle \geq 1/\sqrt{m}$ .

# Cone Approximation: New Results

**Theorem 3 (Zhang, Mu, Kuo, W. '12)** *(sketch)* Under our previously stated hypotheses, for all  $\mathbf{u}, \mathbf{u}' \in \mathbb{S}^2$ ,

$$\|\bar{\mathbf{y}}[\mathbf{u}] - \bar{\mathbf{y}}[\mathbf{u}']\|_2 \leq \frac{C_{\text{sensor}}}{1 - \nu_\star} \times \left( \text{area}(\partial\mathcal{O}) \|\mathbf{u} - \mathbf{u}'\|_2^2 + \chi_\star \text{diam}(\mathcal{O}) \|\mathbf{u} - \mathbf{u}'\|_2 \right)^{1/2}$$

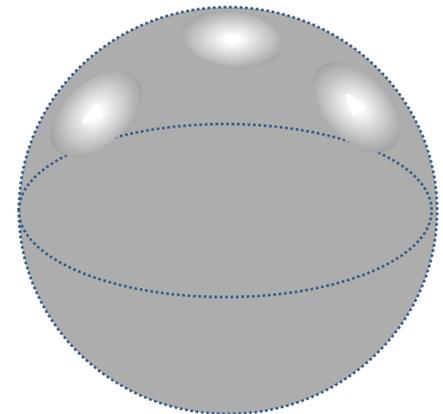
where  $C_{\text{sensor}}$  depends on the parameters of the imaging system.

If we define an **illumination covering number**

$$N(C, \gamma) = \min \{ \#V \mid \delta(C, \text{cone}(V)) \leq \gamma \},$$

This implies in general,  $N(C_\alpha, \gamma) \leq \frac{h(\mathcal{O}, \text{sensor})}{(\alpha\gamma)^4},$

and for convex objects  $N(C_\alpha, \gamma) \leq \frac{h(\mathcal{O}, \text{sensor})}{(\alpha\gamma)^2}.$

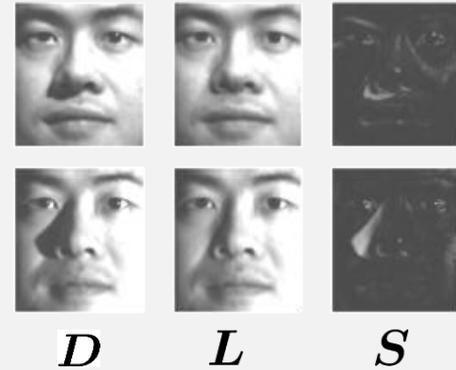


# Proof Sketch

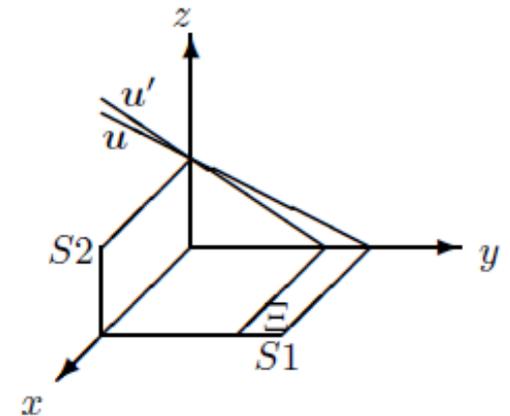
Control the complexity of cone approximation in terms of object convexity defect.

## The conceptual idea:

Separate the direct illumination operator into a **low-rank** part (due to smooth variations) and a **sparse** part (due to cast shadows):



The low-rank term can be bounded by direct calculation; the sparse term needs a more involved accounting.



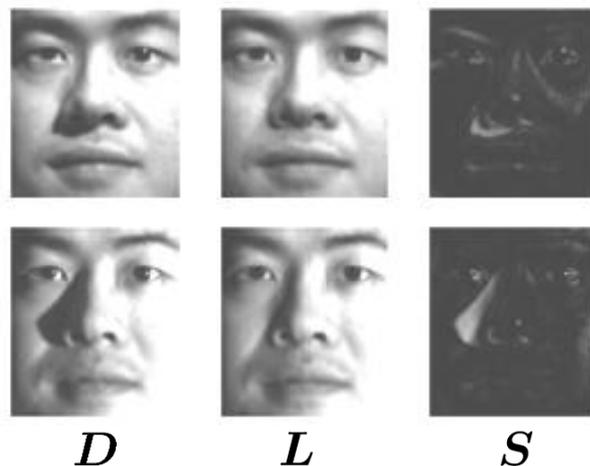
# LR+S captures physical intuition

Calculations with the Lambertian sphere suggest that sets of images of **smooth, near-convex** objects should be approximately **low-rank**:

[Basri+Jacobs '03,Ramamoorthi '04].

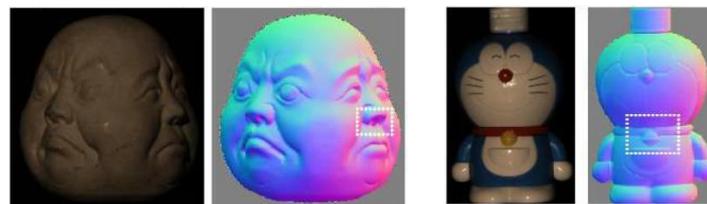
**Cast shadows** are often **sparse**

[W., Yang, ... Ma, '09] , [Candes, Li, Ma, W. '11]:



Observations used in  
photometric stereo:

[Wu et. Al. '11] :

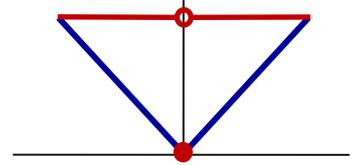


**How can we exploit these structures to obtain compact representation/approximation of the cone?**

# Low-rank and sparse recovery

**Convex relaxation:**  $\min \|L\|_* + \lambda \|S\|_1 \quad \text{s.t.} \quad L + S = D.$

[Chandrakaran et. al. '11]



Provably effective for recovery (and **statistical estimation**):

**Theorem 5** (Candès, Li, Ma, W. '11). If  $L_0 \in \mathbb{R}^{m \times n}$ ,  $m \geq n$  has rank

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

and  $S_0$  has Bernoulli support with error probability  $\rho \leq \rho_s^*$ , then w.h.p.,

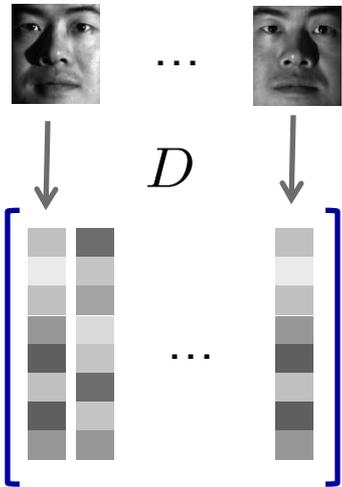
$$(L_0, S_0) = \arg \min \|L\|_* + \frac{1}{\sqrt{m}} \|S\|_1 \quad \text{subj} \quad L + S = L_0 + S_0,$$

and the minimizer is unique.

See also ... [Chandrakaran, Sanghavi, Willsky, Parillo '11], [Candès, Li, Ma, W. '11], [Hsu, Kakade, Zhang '12], [Agarwal, Wainwright '12],

# Recover low-rank and sparse components

58 images of one person under varying lighting:



$D$

...

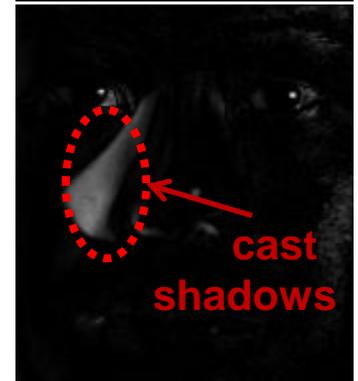
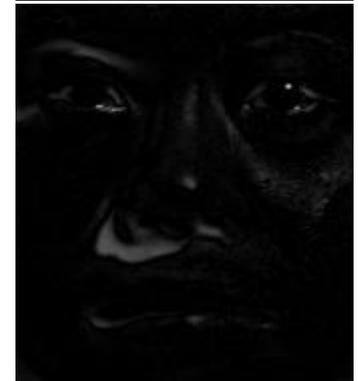
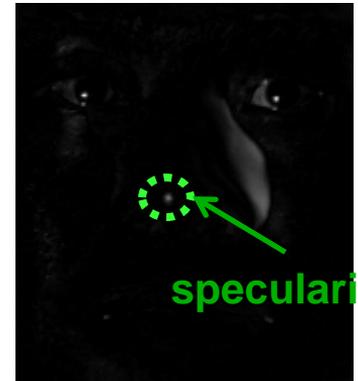
$D$



$L$



$S$



# Cone-preserving dimensionality reduction

A low-rank and sparse decomposition that provably **preserves verification performance?**

Would need to solve

$$\text{minimize } \|\mathbf{L}\|_* + \lambda\|\mathbf{S}\|_1 \quad \text{s.t.} \quad \delta\left(\text{cone}(\mathbf{L} + \mathbf{S}), \text{cone}(\mathbf{A})\right) \leq \gamma$$

... but constraint is very complicated. Relax to:

**Theorem 6 (Zhang, Mu, Kuo, W. '13)** *Let  $(\mathbf{L}_*, \mathbf{S}_*)$  solve*

$$\text{minimize } \|\mathbf{L}\|_* + \lambda\|\mathbf{S}\|_1 \quad \text{s.t.} \quad \begin{bmatrix} \mathbf{I} & \mathbf{L} + \mathbf{S} - \mathbf{A} \\ \mathbf{L}^* + \mathbf{S}^* - \mathbf{A}^* & \gamma' \mathbf{A}^* \mathbf{A} - \boldsymbol{\mu} \end{bmatrix} \succeq \mathbf{0}, \boldsymbol{\mu} \geq \mathbf{0},$$

*with  $\gamma' = \frac{\gamma}{1+\gamma}$ . Then  $\delta\left(\text{cone}(\mathbf{L}_* + \mathbf{S}_*), \text{cone}(\mathbf{A})\right) \leq \gamma$ .*

# Provable Low-Dim Structures of Visual Data

**Guaranteed Illumination Models for Nonconvex Objects**, Zhang, Mu, Kuo, W., Arxiv '13

**Compressive Principal Component Pursuit**, W., Ganesh, Min, Ma, I&I '13

**Towards a Practical Automatic Face Recognition**, Wagner, W., Ganesh, Zhou, Mobahi, Ma, PAMI '12

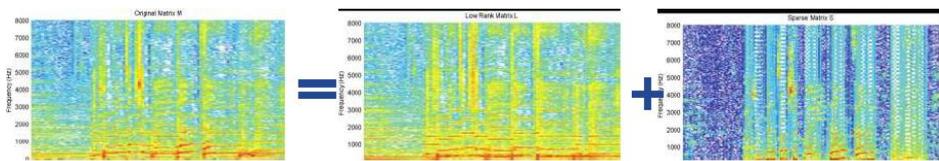
**Robust Principal Component Analysis?** Candes, Li, Ma, W. JACM '11

**Dense Error Correction via  $\ell^1$  Minimization.** W. and Ma, Information Theory '10

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

# Low-rank and sparse: applications

**Audio:** music and voices  
[Huang et. al. '12]



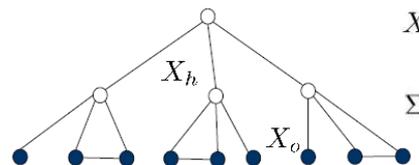
**Video:** foreground/background  
[CLMW '11]



**Alignment:** images online  
[PGWM'11]



**Learning graphical models:**  
[Chandrasekaran et. al. '12]



**Many others ... text (robust topic models), community detection in graphs, subspace segmentation, 3D reconstruction...**

# Generalizations

Upgrade a certificate to a *compressive* certificate. More generally:

**Multiple structure decomposition:**

$$\min \sum_i \lambda_i \|\mathbf{X}_i\|_{\diamond_i} \quad \text{s.t.} \quad \sum_i \mathbf{X}_i = \mathbf{D}$$

**Compressive multiple structure decomposition:**

$$\min \sum_i \lambda_i \|\mathbf{X}_i\|_{\diamond_i} \quad \text{s.t.} \quad \mathcal{P}_Q[\sum_i \mathbf{X}_i] = \mathcal{P}_Q[\mathbf{D}]$$

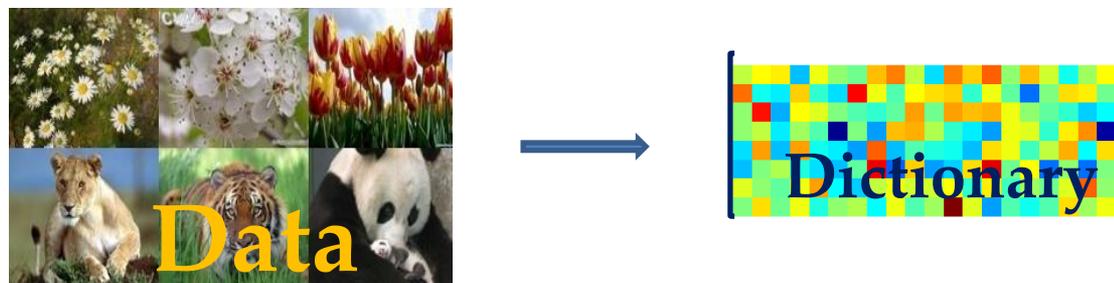
Examples: **PCP** [Chandrakaran, Sanghavi, Willsky, Parillo'11, CLMW '11],  
**outlier pursuit** [Xu+Caramanis+Sanghavi],  
**morphological component analysis** [Bobin et. al.], many more ...

# Learning simple signal models?

$$y = A x_0$$

with  $x_0 \in \mathbb{R}^n$  **sparse** - most of the  $x_0(i)$  are zero.

Good model for many types of imagery data, especially if we can **learn the dictionary**  $A = [a_1 \mid \cdots \mid a_n] \in \mathbb{R}^{m \times n}$

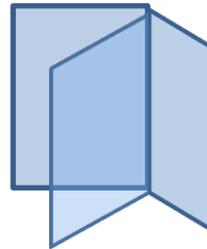


# The model problem

Given  $\mathbf{Y} \approx \mathbf{A}\mathbf{X}$  with  $\mathbf{x}_j$  sparse,  
( $\mathbf{A}$ ,  $\mathbf{X}$ ) unknown, recover  $\mathbf{A}$  and  $\mathbf{X}$ .

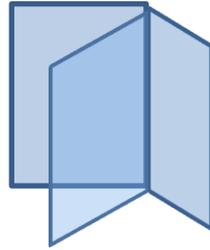
Ambiguities: ( $\mathbf{A}$ ,  $\mathbf{X}$ ) or ( $\mathbf{A}\mathbf{\Pi}\mathbf{\Lambda}$ ,  $\mathbf{\Lambda}^{-1}\mathbf{\Pi}^*\mathbf{X}$ )?

Peculiar geometry:

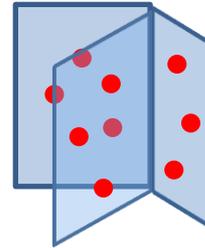


$k$  column subspaces of  $\mathbf{A}$

# When is dictionary learning well-posed?



$\binom{n}{k}$   $k$  column  
subspaces of  $\mathbf{A}$

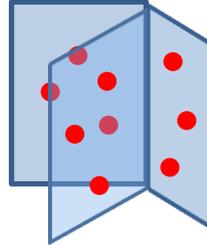


$\mathbf{Y}$   $k+1$  points  
per subspace

Solution is unique:

**Theorem 1** (ess. Aharon et. al. '05) *(sketch)* There exists  $k$  column sparse  $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_p]$ , of size  $p = (k+1)\binom{n}{k}$  such that if we observe  $\mathbf{Y} = \mathbf{A}\mathbf{X}$ ,  $(\mathbf{A}, \mathbf{X})$  is essentially the only  $k$ -column sparse factorization of  $\mathbf{Y}$ .

# When does a learned dictionary generalize?



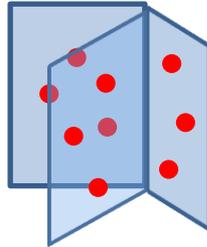
**Theorem 2 (Vainsencher, Mannor and Bruckstein '11)** (sketch) If  $\mathbf{y} \sim_{iid} \mu$  on  $\mathbb{S}^{m-1}$ ,  $p > p_0$ ,  $\lambda > \lambda_0$ , then with prob.  $1 - e^{-t}$  in  $\mathbf{Y}$ ,

$$\mathbb{E}_{\mathbf{y}} \min_{\|\mathbf{x}\|_1 \leq \lambda} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|$$

$$\leq \frac{1.1}{p} \sum_i \|\mathbf{y}_i - \mathbf{A}\mathbf{x}_i\| + \boxed{9 \frac{mn \log(\lambda p) + t}{p}}$$

See also [Maurer and Pontil '10].

# How can we learn a good dictionary?



$$\mathbf{Y} \approx \mathbf{A}\mathbf{X}, \mathbf{X} \text{ sparse.}$$

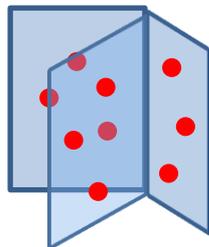
## Alternating directions to minimize sparsity surrogate

[Engan et. al., '99, Aharon et. al. '05, Yaghoobi '10]

$$\min \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 + J(\mathbf{X})$$

**Recently:** Supervised variants [Mairal et. al. '08], structured dictionaries [Rubenstein et. al. '10], highly scalable variants [Mairal et. al. '10] ... and many, many more...

# Is the desired solution a local minimum?



$$Y = AX, X \text{ sparse.}$$

$$\min \|X'\|_1 \quad \text{s.t.} \quad Y = A'X', A' \in \mathcal{A}$$

For square  $A$ , under probabilistic assumptions on  $X$ ,  
 $(A, X)$  is a local minimum whp:

**Theorem 3 (Gribonval + Schnass '10)** *(sketch)* Let  $X_{ij} = \Omega_{ij}V_{ij}$ , with  $\Omega \sim \text{Ber}(\theta)$ ,  $V \sim \mathcal{N}(0, 1)$ . For square, incoherent  $A$ ,  $(A, X)$  is a local minimum of  $\|\cdot\|_1$  with high probability, provided  $p = \Omega(n \log n / \theta)$ .

# Is the desired solution a local minimum?

$$\min \|\mathbf{X}'\|_1 \quad \text{s.t.} \quad \mathbf{Y} = \mathbf{A}'\mathbf{X}', \quad \mathbf{A}' \in \mathcal{A}$$

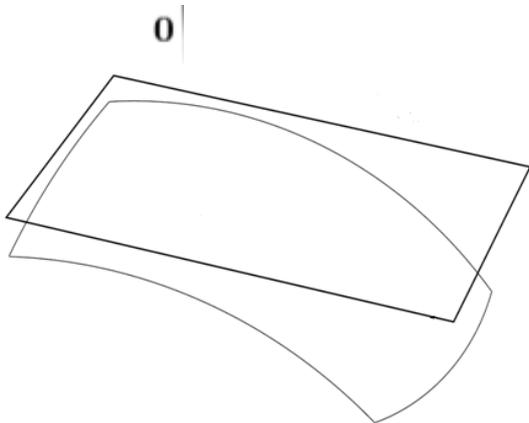
For general  $\mathbf{A}$ , under probabilistic assumptions on  $\mathbf{X}$ ,  
 $(\mathbf{A}, \mathbf{X})$  is a **local minimum whp**:

**Theorem 4** (Geng, W., '11). *Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $k < C/\mu(\mathbf{A})$ , and  $\mathbf{X} \in \mathbb{R}^{n \times p}$  with random  $k$ -sparse support, independent Gaussian nonzeros. Then  $(\mathbf{A}, \mathbf{X})$  is a local minimum of the  $\ell^1$ -norm w.p.  $\geq 1 - \tilde{O}(n^{3/2}k^{1/2}p^{-1/2})$ .*

# Is this obvious?

Maybe ... but surprisingly resistant to analysis ...

$$\min \|\mathbf{X}'\|_1 \quad \text{s.t.} \quad \mathbf{Y} = \mathbf{A}'\mathbf{X}', \quad \mathbf{A}' \in \mathcal{A}$$

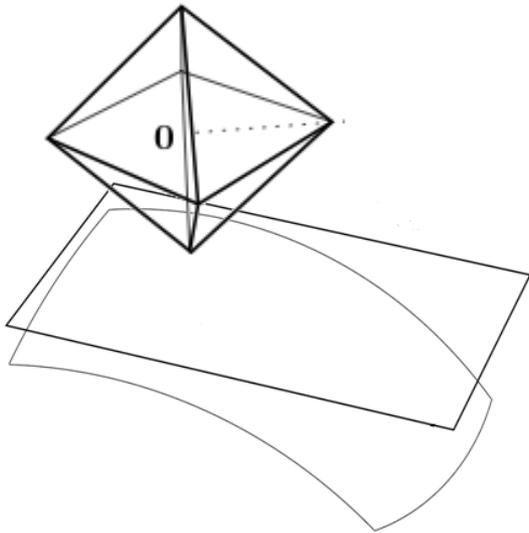


Feasible  $(\mathbf{A}, \mathbf{X})$

# Is this obvious?

Maybe ... but surprisingly resistant to analysis ...

$$\min \|\mathbf{X}'\|_1 \quad \text{s.t.} \quad \mathbf{Y} = \mathbf{A}'\mathbf{X}', \quad \mathbf{A}' \in \mathcal{A}$$

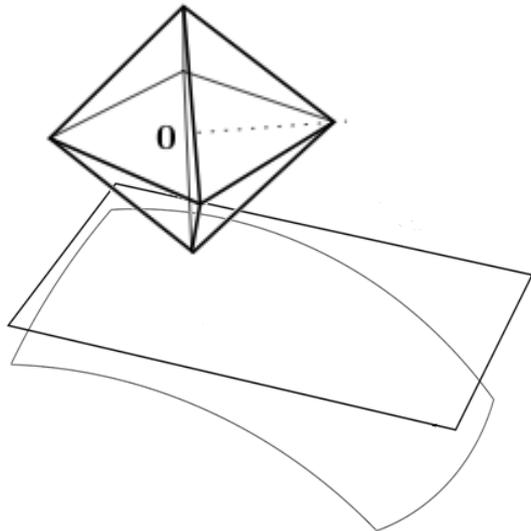


Feasible  $(\mathbf{A}, \mathbf{X})$

# Is this obvious?

Maybe ... but surprisingly resistant to analysis ...

$$\min \|\mathbf{X}'\|_1 \quad \text{s.t.} \quad \mathbf{Y} = \mathbf{A}'\mathbf{X}', \quad \mathbf{A}' \in \mathcal{A}$$



Feasible  $(\mathbf{A}, \mathbf{X})$

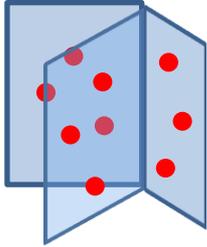
Have to analyze an  $\ell^1$  problem  
over an affine space.

RIP ect., fail here  
ess. sign-permutation ambiguity

Use ideas from **low-rank recovery**  
[Gross '09], [Candes, Li, Ma, W. '12].

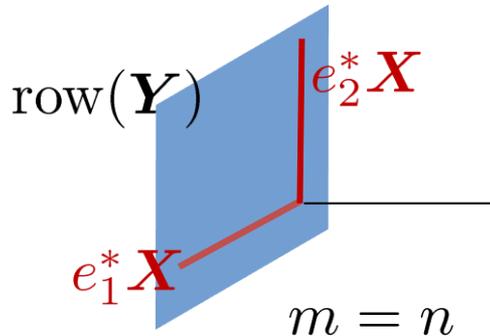
# Can we obtain global results?

In general, not easy... special geometry for square case



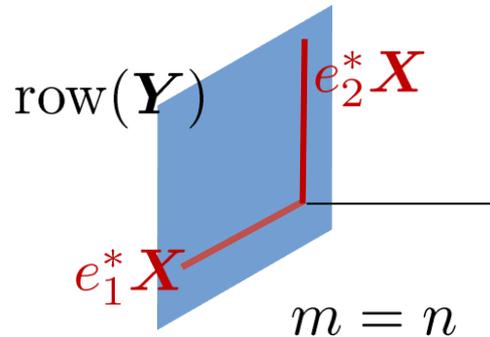
$$Y = AX, X \text{ sparse.}$$

If  $A \in \mathbb{R}^{m \times m}$  nonsingular,  $\text{row}(X) = \text{row}(Y)$ .



Rows of  $X$  are sparse  
vectors in a known subspace.

# Uniqueness – square dictionaries

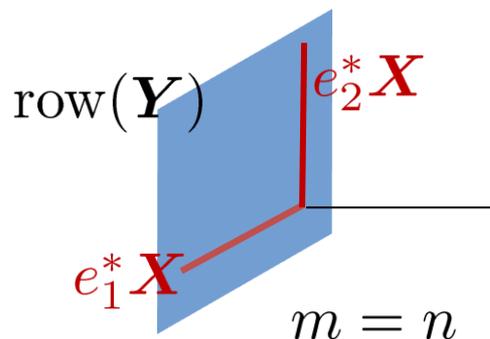


Rows of  $\mathbf{X}$  are sparse vectors in a known subspace.

If  $p > cn \log n$  then whp. rows of  $\mathbf{X}$  are the sparsest vectors in  $\text{row}(\mathbf{Y})$

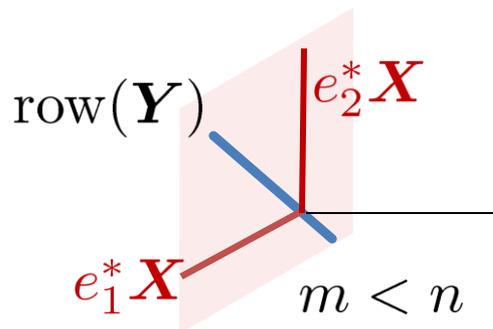
# Uniqueness – square dictionaries

Square:



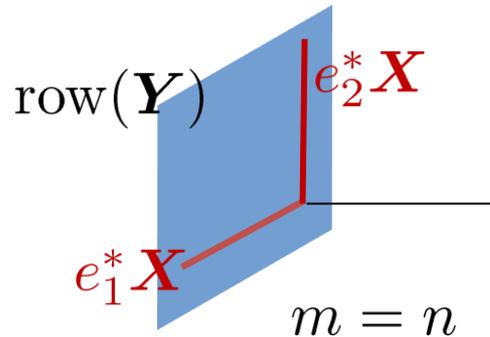
**Theorem [Spielman, Wang, W. '11]:**  
Decomposition essentially unique from  $\Omega(n \log n)$  random observations.

Overcomplete:



**Theorem [Aharon, Elad, Bruckstein '05]:**  
Decomposition is essentially unique from  $(k + 1) \binom{n}{k}$  strategically located observations.

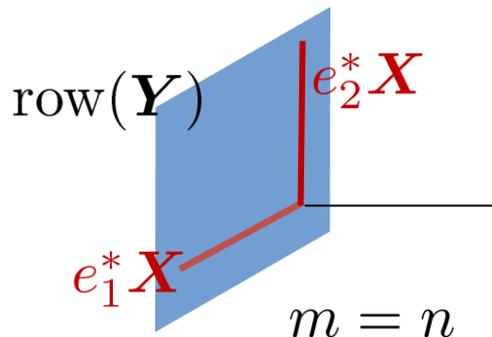
# Algorithms – square dictionaries



Rows of  $\mathbf{X}$  are sparsest  
vectors in  $\text{row}(\mathbf{Y})$

$$\text{minimize } \|\mathbf{w}^* \mathbf{Y}\|_0 \quad \text{subject to } \mathbf{w} \neq 0.$$

# Algorithms – square dictionaries



Rows of  $\mathbf{X}$  are sparsest vectors in  $\text{row}(\mathbf{Y})$

$$\text{minimize } \|\mathbf{w}^* \mathbf{Y}\|_0 \quad \text{subject to } \mathbf{w} \neq 0.$$

$$\text{minimize } \|\mathbf{w}^* \mathbf{Y}\|_1 \quad \text{subject to } \mathbf{r}^* \mathbf{w} = 1.$$

# Algorithms – square dictionaries

$$\text{minimize } \|\mathbf{w}^* \mathbf{Y}\|_1 \quad \text{subject to } \mathbf{r}^* \mathbf{w} = 1.$$

What choice of  $\mathbf{r}$  will make  $\hat{\mathbf{w}}^* \mathbf{Y} = \mathbf{e}_i^* \mathbf{X}$ ?

Change variables  $\mathbf{q} = \mathbf{A}^* \mathbf{w}$ :

$$\text{minimize } \|\mathbf{q}^* \mathbf{X}\|_1 \quad \text{subject to } (\mathbf{A}^{-1} \mathbf{r})^* \mathbf{q} = 1.$$

If  $\mathbf{r} = \mathbf{A} \mathbf{e}_i$  we're golden ...

Don't have this; use  $\mathbf{y}_j = \sum_{i \in I} X_{ij} \mathbf{A} \mathbf{e}_i$

# Algorithms – square dictionaries

**ER-SpUD(SC):** Exact Recovery of Sparsely-Used Dictionaries using single columns of  $Y$  as constraint vectors.

For  $j = 1 \dots p$

Solve  $\min_w \|w^T Y\|_1$  subject to  $(Y e_j)^T w = 1$ , and set  $s_j = w^T Y$ .

**Greedy:** A Greedy Algorithm to Reconstruct  $X$  and  $A$ .

1. **REQUIRE:**  $\mathcal{S} = \{s_1, \dots, s_T\} \subset \mathbb{R}^p$ .

2. For  $i = 1 \dots n$

**REPEAT**

$l \leftarrow \arg \min_{s_l \in \mathcal{S}} \|s_l\|_0$ , breaking ties arbitrarily

$x_i = s_l$

$\mathcal{S} = \mathcal{S} \setminus \{s_l\}$

**UNTIL**  $\text{rank}([x_1, \dots, x_i]) = i$

3. Set  $X = [x_1, \dots, x_n]^T$ , and  $A = Y Y^T (X Y^T)^{-1}$ .

# Algorithms – square dictionaries

**ER-SpUD(DC):** Exact Recovery of Sparsely-Used Dictionaries using the sum of two columns of  $Y$  as constraint vectors.

1. Randomly pair columns of  $Y$  into  $p/2$  groups  $g_i = \{Y e_{i1}, Y e_{i2}\}$ .
2. For  $j = 1 \dots p/2$

Let  $r_j = Y e_{j1} + Y e_{j2}$ , where  $Y e_{j1}, Y e_{j2} \in g_j$ .

Solve  $\min_w \|w^T Y\|_1$  subject to  $r_j^T w = 1$ , and set  $s_j = w^T Y$ .

**Greedy:** A Greedy Algorithm to Reconstruct  $X$  and  $A$ .

1. **REQUIRE:**  $\mathcal{S} = \{s_1, \dots, s_T\} \subset \mathbb{R}^p$ .
2. For  $i = 1 \dots n$

REPEAT

$l \leftarrow \arg \min_{s_l \in \mathcal{S}} \|s_l\|_0$ , breaking ties arbitrarily

$x_i = s_l$

$\mathcal{S} = \mathcal{S} \setminus \{s_l\}$

UNTIL  $\text{rank}([x_1, \dots, x_i]) = i$

3. Set  $X = [x_1, \dots, x_n]^T$ , and  $A = Y Y^T (X Y^T)^{-1}$ .

# Provable solution in the complete case

$$\text{minimize } \|\mathbf{WY}\|_1 \quad \text{subject to } \text{diag}[\mathbf{WY}] = \mathbf{1}.$$

If the expected nonzeros per column is smaller than  $\sqrt{n}$   
the algorithm **succeeds whp**:

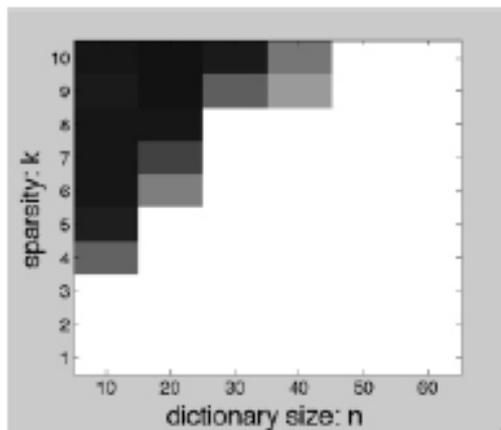
**Theorem 12 (Spielman, Wang, W. '12)** *(sketch)* Let  $\mathbf{X}$  Bernoulli( $\theta$ )–Rademacher or Bernoulli( $\theta$ ) – Gaussian. If  $n > n_0$ ,  $p > c_p n^2 \log^2 n$ , and the nonzero probability satisfies

$$\frac{2}{n} \leq \theta \leq \frac{c}{\sqrt{n}}, \quad (1)$$

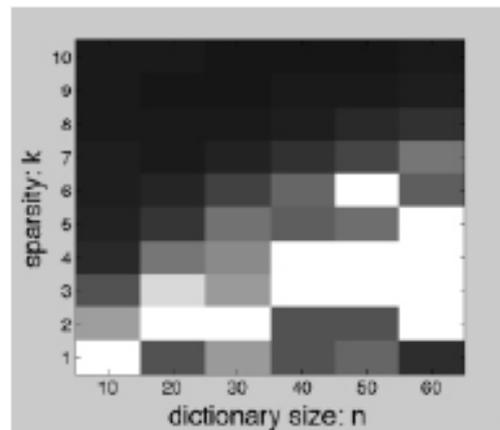
with high probability **ER-SpUD (DC)** recovers all  $n$  rows of  $\mathbf{X}$ .

**Sample requirement**  $p > cn^2 \log^2 n$

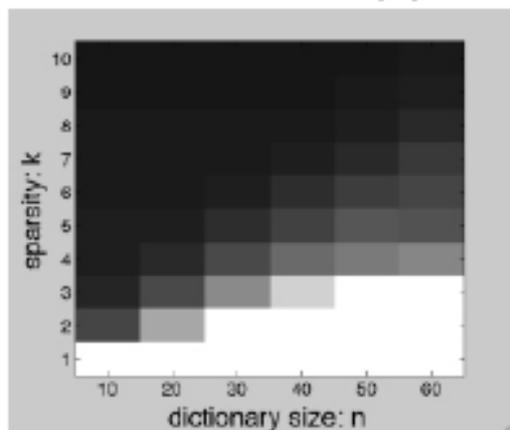
# Does it really work?



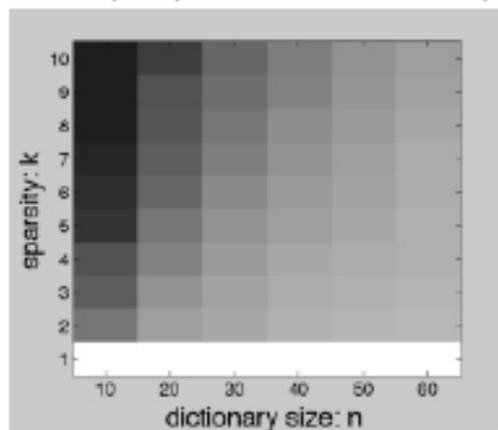
(a) ER-SpUD(SC)



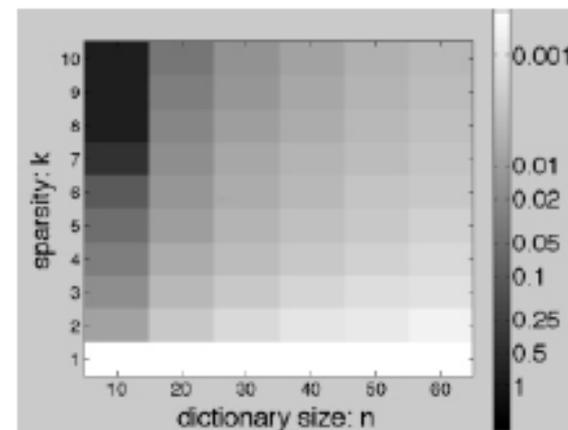
(b) SIV



(c) K-SVD



(d) Online

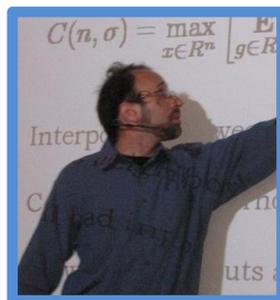
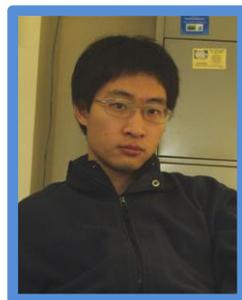


(e) Rel. Newton

Caveat: exact sparse, noiseless setting.

# Provable Representations for Visual Data

Thanks to ...



**Guaranteed Illumination Models for Nonconvex Objects**, Zhang, Mu, Kuo, W., Arxiv '13

**Compressive Principal Component Pursuit**, W., Ganesh, Min, Ma, I&I '13

**Local correctness of  $\ell^1$ -minimization for dictionary learning**, Geng, W., Arxiv '11

**Exact Recovery of Sparsely-Used Dictionaries**, Spielman, Wang, W., COLT '12

**Robust Principal Component Analysis?** Candes, Li, Ma, W. JACM '11

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

**Towards a Practical Automatic Face Recognition ...** Wagner, W., Ganesh, Zhou, Mobahi, Ma, PAMI '12