

Dictionary Learning by ℓ^1 -Minimization

John Wright

Electrical Engineering

Columbia University

Sparse Approximation

Model:

$$y = A x_0$$

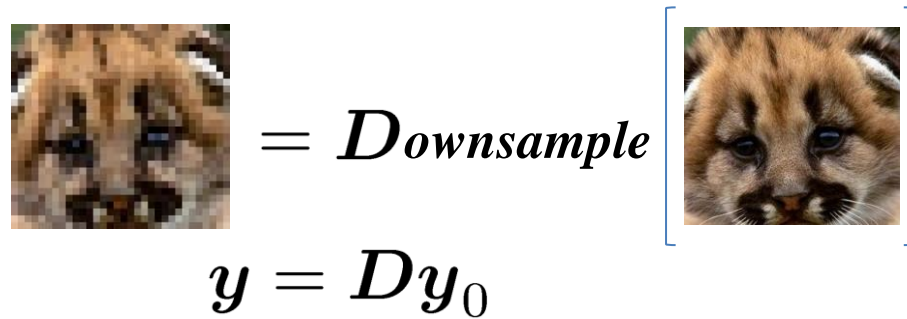
with $x_0 \in \mathbb{R}^n$ **sparse** - most of the $x_0(i)$ are zero.

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



Motivating Applications (biased sample)

Single image superresolution:


$$y = Dy_0$$

reconstruct the original high-resolution image y_0 :

$$\hat{x} \in \arg \min \|x\|_1 \text{ s.t. } y = DAx \quad y_0 \approx A\hat{x}$$

Motivating Applications (biased sample)

High-resolution hyperspectral imaging for cultural heritage:



Ultra high-res RGB camera

Moshe Ben-Ezra
Microsoft Research

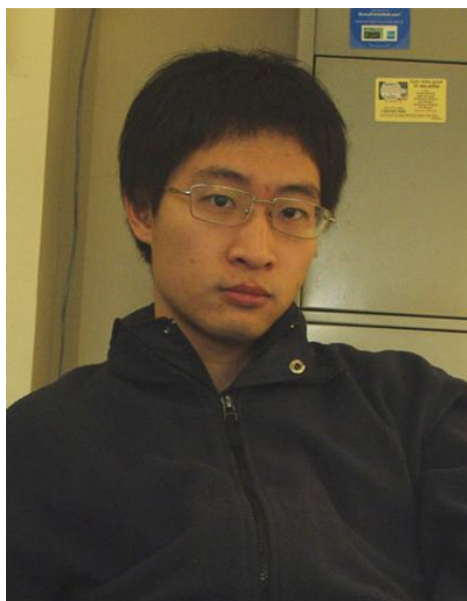


Buddhist Frescos

Dunhuang, China

Can dictionary learning help overcome hardware limitations?

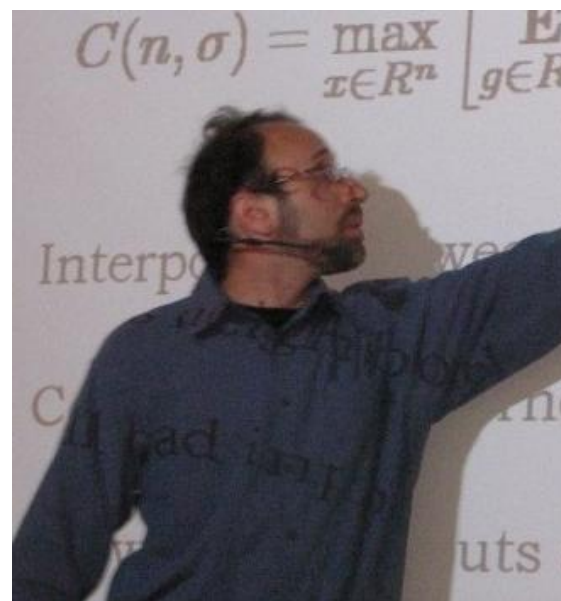
When do dictionary learning algorithms succeed?



Huan Wang (Yale)



Quan Geng (UIUC)



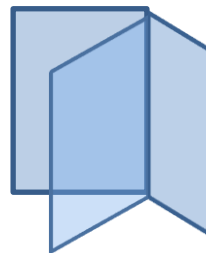
Dan Spielman (Yale)

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} .

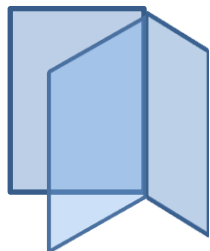
Ambiguities: (\mathbf{A} , \mathbf{X}) or ($\mathbf{A}\mathbf{\Pi}\mathbf{\Lambda}$, $\mathbf{\Lambda}^{-1}\mathbf{\Pi}^*\mathbf{X}$) ?

Peculiar geometry:

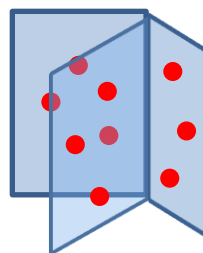


k column subspaces of \mathbf{A}

When is dictionary learning well-posed?



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

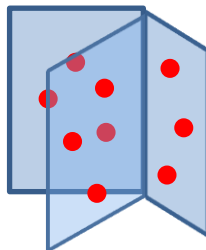


\mathbf{Y} $k+1$ points
per subspace

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?



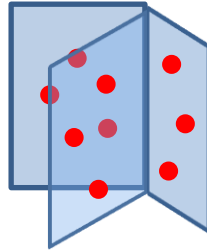
Theorem 2 (Vainsencher, Mannor and Bruckstein '11) *(sketch)* If $\mathbf{y} \sim_{iid} \mu$ on \mathbb{S}^{m-1} , $p > p_0$, $\lambda > \lambda_0$, then with prob. $1 - e^{-t}$ in \mathbf{Y} ,

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

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

See also [Maurer and Pontil '10].

How can we learn a good dictionary?



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

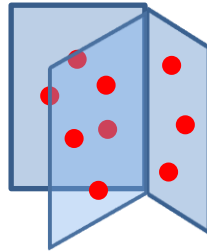
Alternating directions to minimize sparsity surrogate

[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?

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

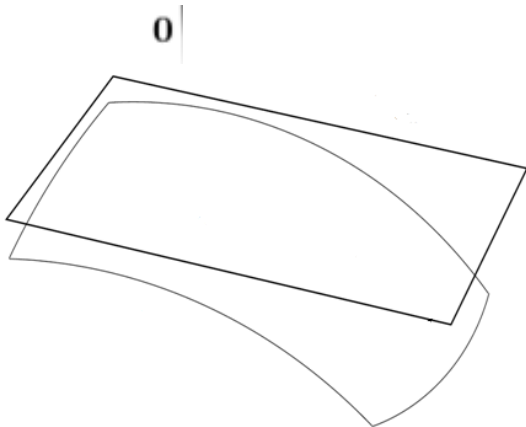
For general \mathbf{A} , under probabilistic assumptions on \mathbf{X} ,
 (\mathbf{A}, \mathbf{X}) is a **local minimum whp**:

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 ℓ^1 -norm $\text{wp} \geq 1 - \tilde{O}(n^{3/2}k^{1/2}p^{-1/2})$.*

Is this obvious?

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

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

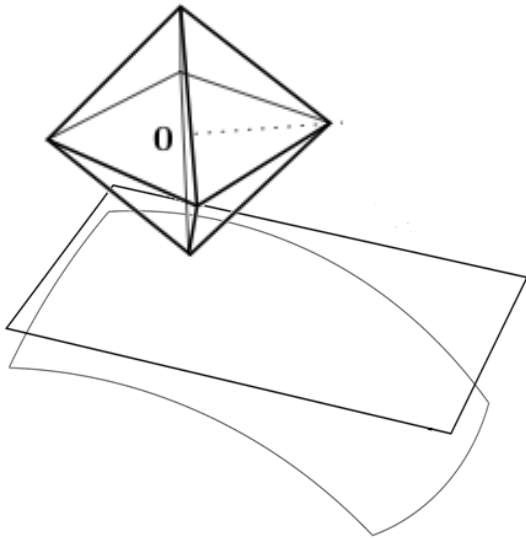


Feasible (\mathbf{A}, \mathbf{X})

Is this obvious?

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

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

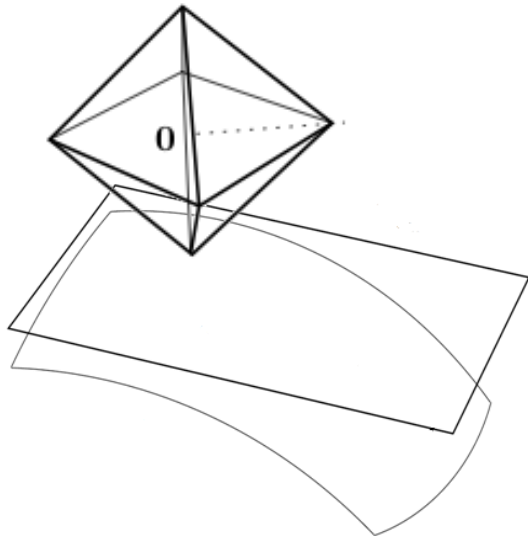


Feasible (\mathbf{A}, \mathbf{X})

Is this obvious?

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

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



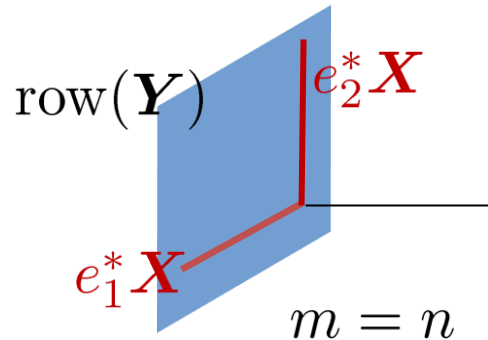
Feasible (\mathbf{A}, \mathbf{X})

Have to analyze an ℓ^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].

Uniqueness – square dictionaries

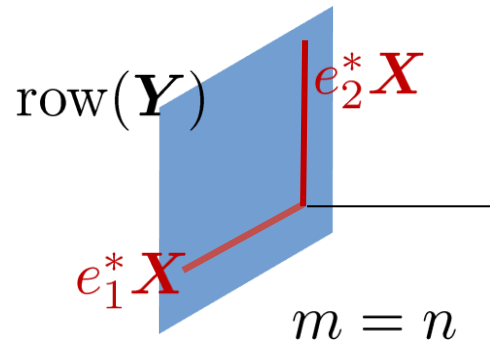


Rows of X are sparse vectors in a known subspace.

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

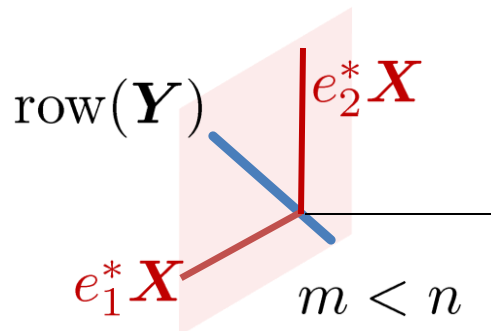
Uniqueness – square dictionaries

Square:



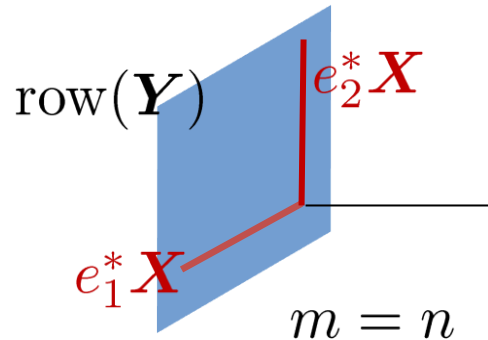
Theorem [Spielman, Wang, W. '11]:
Decomposition essentially unique from $\Omega(n \log n)$ random observations.

Overcomplete:



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

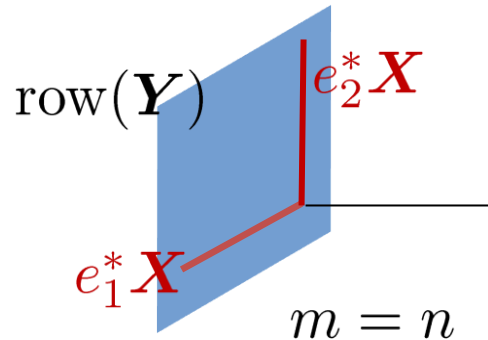
Algorithms – square dictionaries



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 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 $A = Y Y^T (X Y^T)^{-1}$.

Recovery guarantee – square dictionaries

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

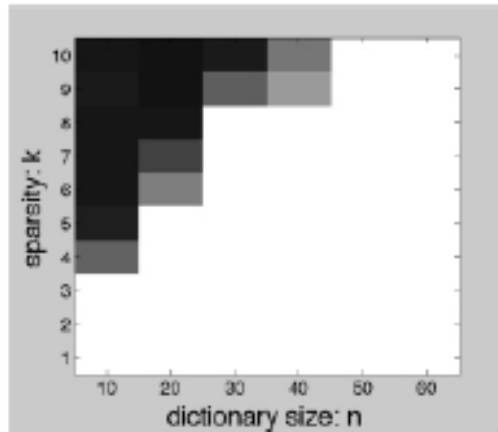
Theorem 5 (Spielman, Wang, W. '12) *(sketch)* Let \mathbf{X} Bernoulli(θ)–Rademacher or Bernoulli(θ) – Gaussian. If $n > n_0$, $p > c_p n^2 \log^2 n$, and the nonzero probability satisfies

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

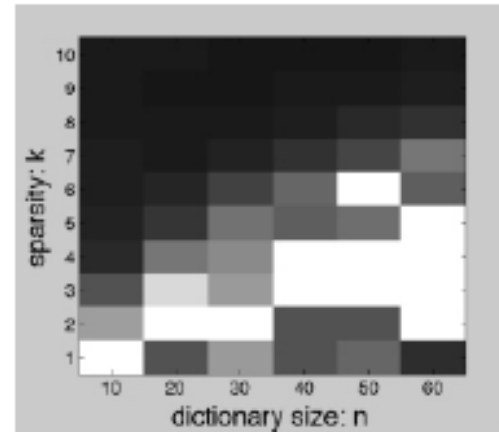
with high probability **ER-SpUD (DC)** recovers all n rows of \mathbf{X} .

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

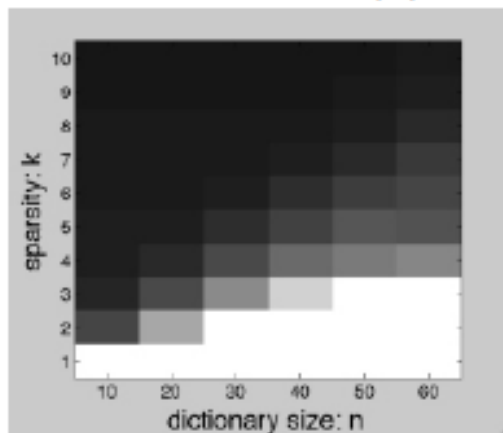
Does it really work?



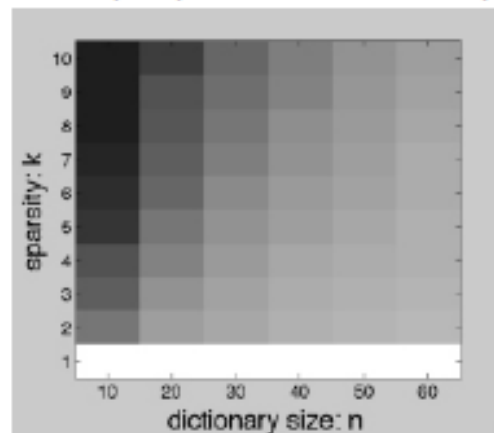
(a) ER-SpUD(SC)



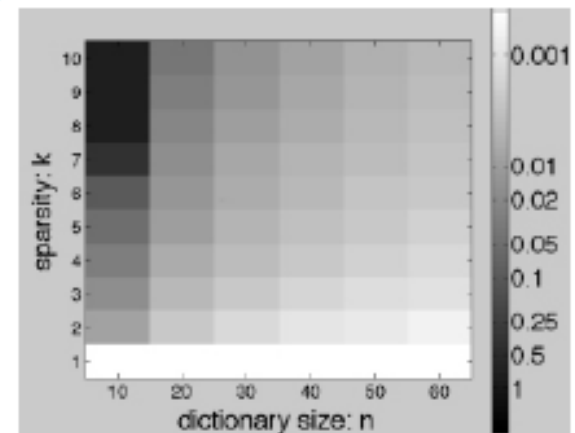
(b) SIV



(c) K-SVD



(d) Online



(e) Rel. Newton

Caveat: exact sparse, noiseless setting.

Good news / bad news ...

If the expected nonzeros per column exceeds $\sqrt{n \log n}$
the algorithm **fails whp**:

Theorem 6 (Spielman, Wang, W. '12) *(sketch)* If n large, $p \geq cn$, and the nonzero probability θ satisfies

$$\theta \geq \sqrt{\frac{\log n}{n}}, \quad (1)$$

then the probability (in \mathbf{X}) that the algorithm correctly recovers one of the rows is at most n^{-C} .

Theory is almost tight in the sparsity level.

For denser \mathbf{X} , think about **different constraints**.

Summary and open questions

Two main mathematical results:

Local recovery in the rectangular case

Exact (global) recovery in the square case

Many open questions:

Past the \sqrt{n} barrier?

Noise tolerance, multiple vectors?

Other coefficient structures?

Dictionary Learning by ℓ^1 Minimization

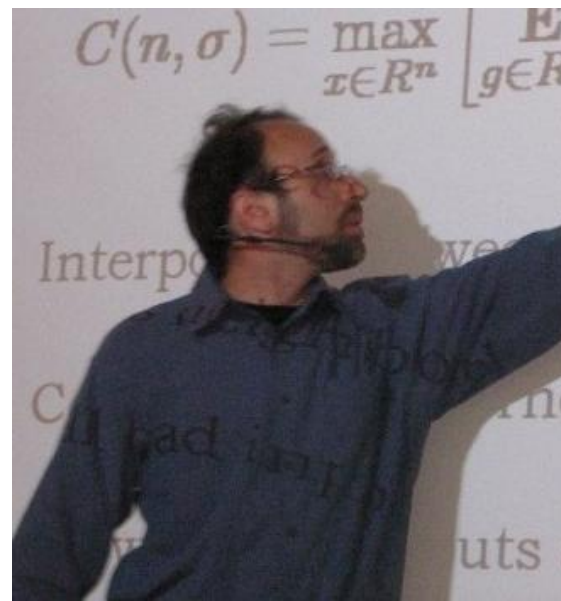
Thanks to ...



Huan Wang (Yale)



Quan Geng (UIUC)



Dan Spielman (Yale)

Local correctness of ℓ^1 -minimization for dictionary learning, Geng, W., Arxiv
Exact recovery of sparse dictionaries, Spielman, Wang, W., COLT '12.