

Dictionary learning for sparse representations

A statistical analysis using L1 minimization

Rémi Gribonval

METISS team (audio signal processing, indexing, source separation)

INRIA, Rennes, France

Karin Schnass

LTS2, EPFL, Switzerland



Journées STAR 2009
IRISA, 22-23 octobre 2009



Outline

1. Preliminaries:

- ✦ sparsity & (overcomplete) dictionaries of atoms
- ✦ blind source separation & dictionary learning

2. Objectives of (theoretical) dictionary learning

3. L1 minimization for dictionary learning

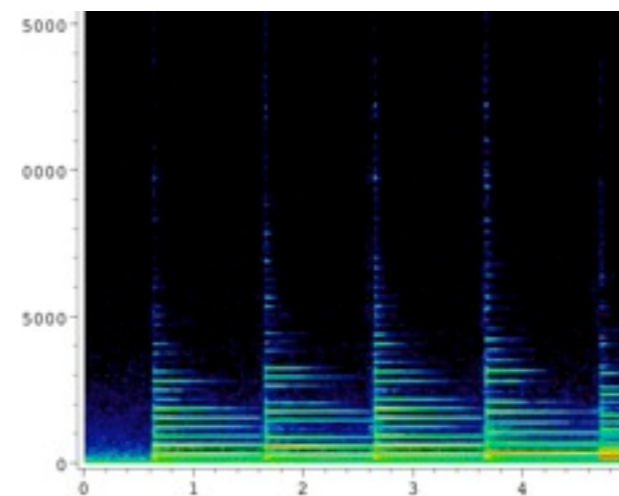
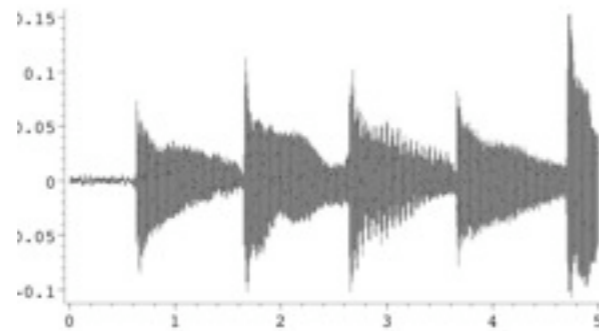
4. Main results

- ✦ geometric “local” identifiability condition
- ✦ random model and finite sample size analysis

5. Discussion, conclusion & challenges

Sparse data representations

- Short Time Fourier Transform (for time-varying harmonic signals)

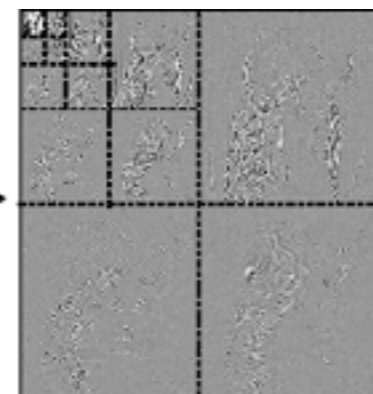


Black = zero

- Wavelet transform (for piecewise smooth images)



ORIGINAL
128, 128, 125, 64, 65,



TRANSFORM COEFFICIENTS
4123, -12.4, -96.7, 45,

Grey = zero

Sparse signal models

- An image / a signal = sum of few atoms \mathbf{a}_k

$$\mathbf{b} = \sum_k x_k \mathbf{a}_k = \mathbf{A}x$$

- ❖ *Dictionary* = collection of atoms = \mathbf{A}
- ❖ *Representation* = coefficient vector = x

- Sparsity of x : enables compression, separation ...
- Sparsity of x ? Only if dictionary \mathbf{A} is “well chosen”
 - ❖ Pre-chosen atoms: wavelets, Gabor, etc.
 - ❖ Learned dictionary = from collection of signals / images

$$\mathbf{b}_n = \mathbf{A}x_n, \quad 1 \leq n \leq N$$

Dictionary learning for sparse representations

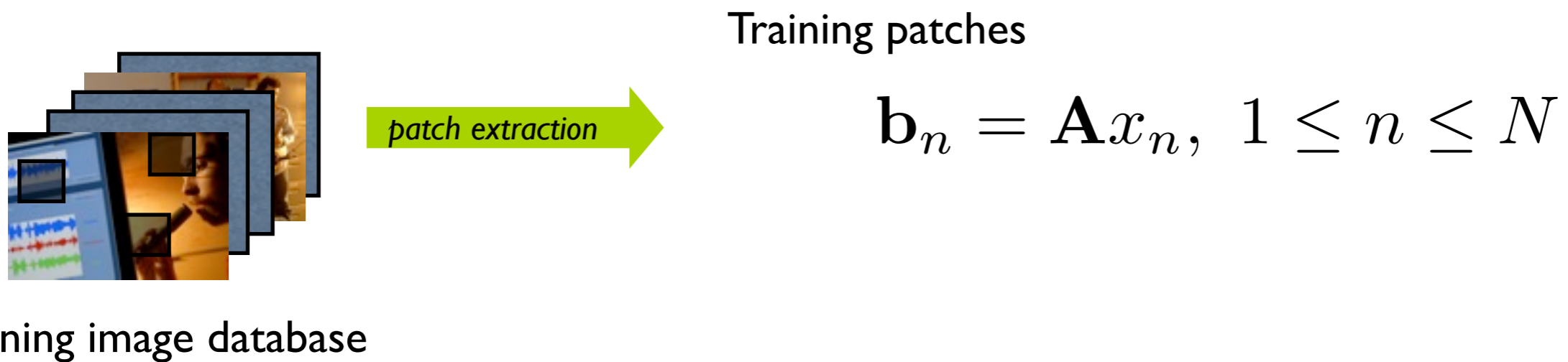
- Sparse modeling : choose a dictionary



Training image database

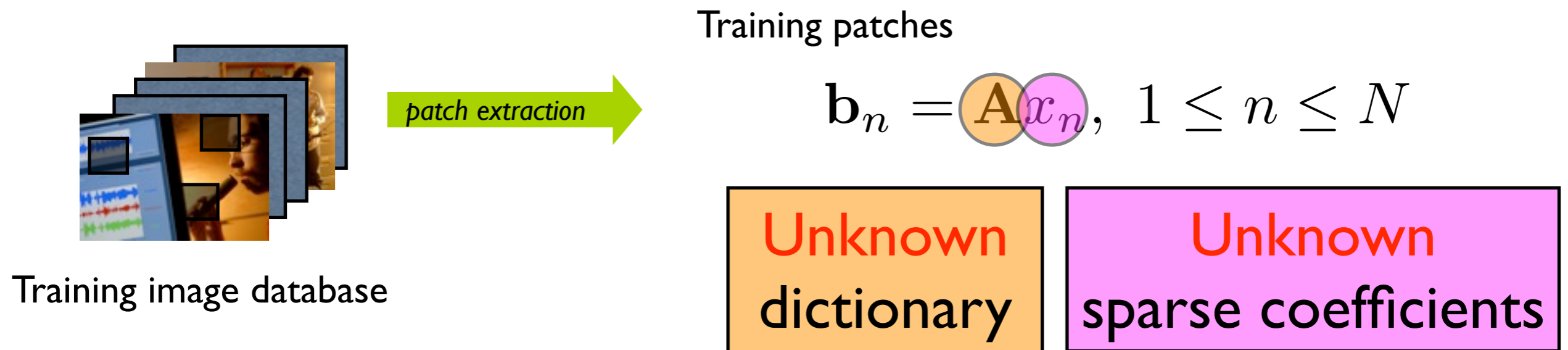
Dictionary learning for sparse representations

- Sparse modeling : choose a dictionary



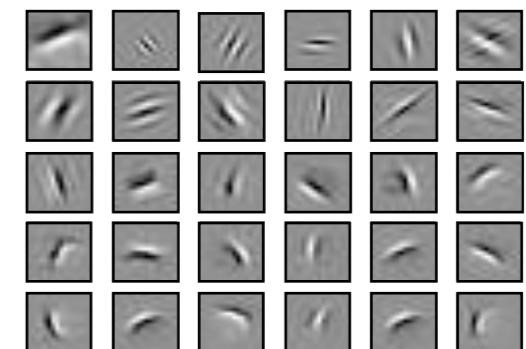
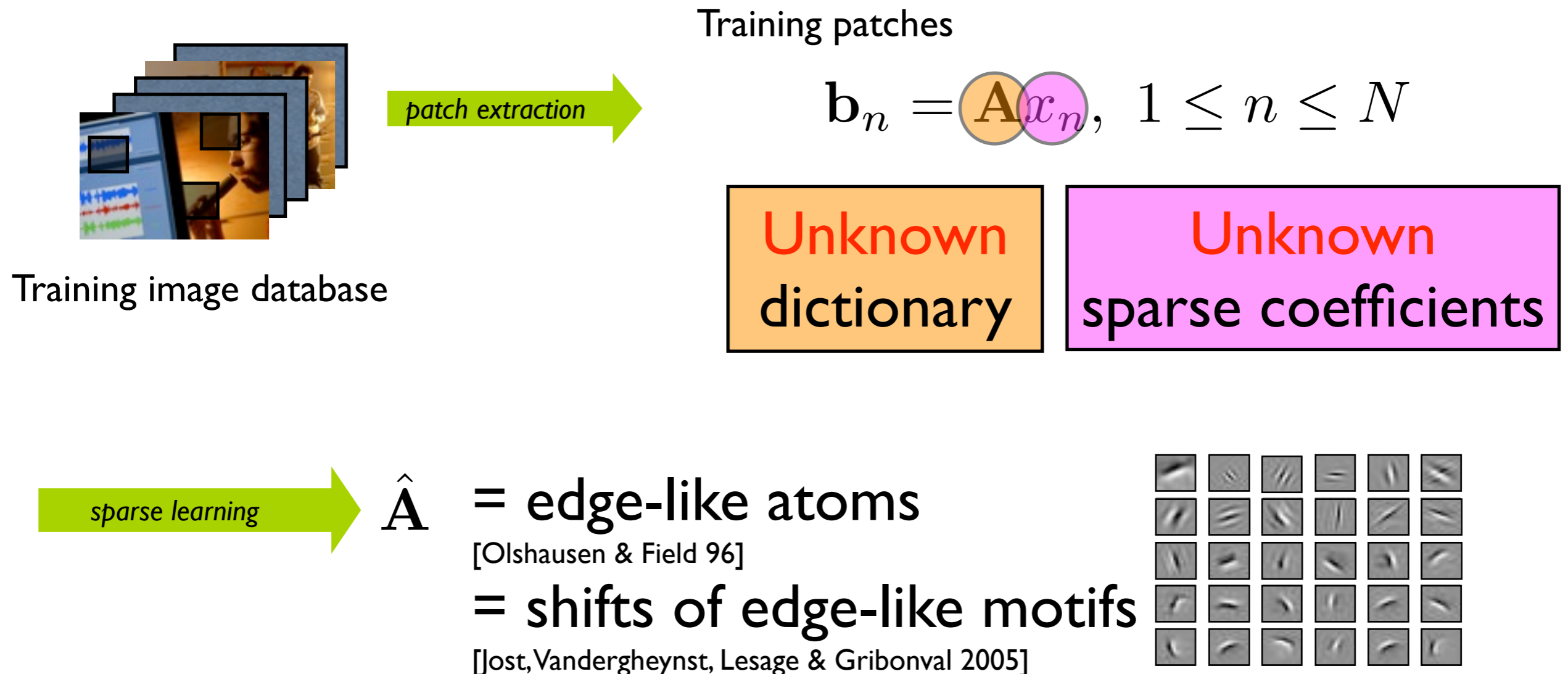
Dictionary learning for sparse representations

- Sparse modeling : choose a dictionary



Dictionary learning for sparse representations

- Sparse modeling : choose a dictionary



Dictionary learning ?

- Problem : estimate a matrix \mathbf{A} given observed samples

$$\mathbf{b}_n = \mathbf{A}x_n, \quad 1 \leq n \leq N$$

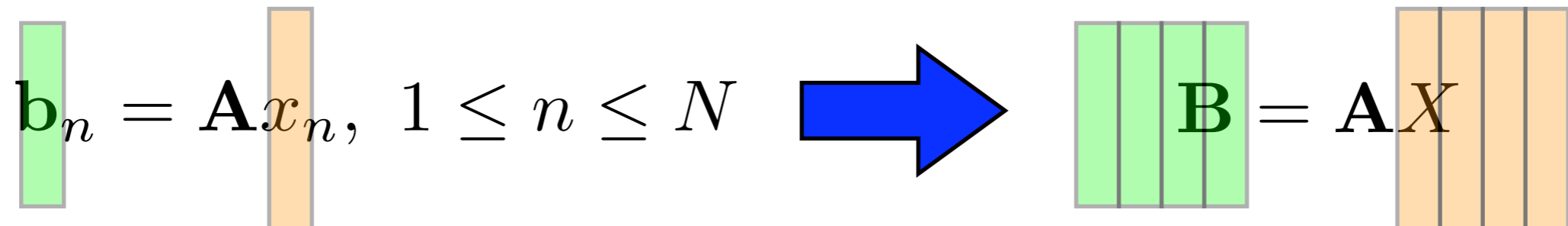
Dictionary learning ?

- Problem : estimate a matrix A given observed samples

$$\mathbf{b}_n = A \mathbf{x}_n, \quad 1 \leq n \leq N$$

Dictionary learning ?

- Problem : estimate a matrix A given observed samples

$$\mathbf{b}_n = \mathbf{A} x_n, \quad 1 \leq n \leq N \quad \Rightarrow \quad \mathbf{B} = \mathbf{A} X$$


Dictionary learning ?

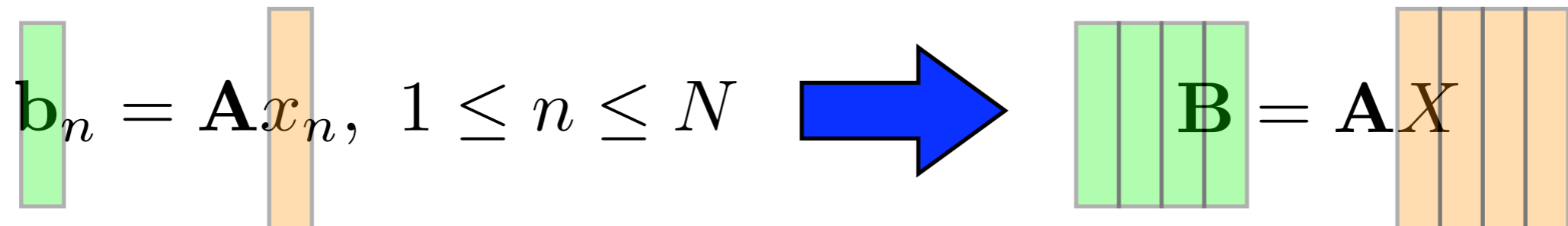
- Problem : estimate a matrix A given observed samples

$$\mathbf{b}_n = \mathbf{A} \mathbf{x}_n, \quad 1 \leq n \leq N \quad \longrightarrow \quad \mathbf{B} = \mathbf{A} \mathbf{X}$$

\mathbf{A} { **Unknown** mixing matrix (blind source separation)
Unknown dictionary (sparse signal approximation)
Unknown channel filter (blind channel estimation) ...

Dictionary learning ?

- Problem : estimate a matrix A given observed samples

$$\mathbf{b}_n = \mathbf{A} \mathbf{x}_n, \quad 1 \leq n \leq N \quad \longrightarrow \quad \mathbf{B} = \mathbf{A} \mathbf{X}$$


A { **Unknown** mixing matrix (blind source separation)
Unknown dictionary (sparse signal approximation)
Unknown channel filter (blind channel estimation) ...

X **Unknown** sources / signal representations / ...

Dictionary learning ?

- Problem : estimate a matrix A given observed samples

$$\mathbf{b}_n = \mathbf{A} \mathbf{x}_n, \quad 1 \leq n \leq N \quad \longrightarrow \quad \mathbf{B} = \mathbf{A} \mathbf{X}$$

A { **Unknown** mixing matrix (blind source separation)
Unknown dictionary (sparse signal approximation)
Unknown channel filter (blind channel estimation) ...

X **Unknown** sources / signal representations / ...

- Fundamentally ill-posed **factorization problem** :
need (weak) model on unknown coefficients X and /
or matrix A

Theoretical dictionary learning

- Problem : estimate a matrix \mathbf{A} given samples

$$\mathbf{b}_n = \mathbf{A}x_n, \quad 1 \leq n \leq N$$

$$\mathbf{B} = \mathbf{A}X$$

| | |
|-----------------|--|
| | ICA (Independent Component Analysis) |
| Model of ... | probability density function $p(X)$ |
| Assumption | Independence $p(X) = \prod_{nk} p(x_n(k))$ |
| Identifiability | Darmois theorem |
| Identification | Contrast functions $\mathbf{A} \sim \hat{\mathbf{W}}^{-1} \quad \hat{\mathbf{W}} := \arg \min_{\mathbf{W}} E_X(f(\mathbf{W}\mathbf{A}X))$ |
| Issues | In practice : finite training sets expectation \longrightarrow sample average |

Theoretical dictionary learning

- Problem : estimate a matrix \mathbf{A} given samples

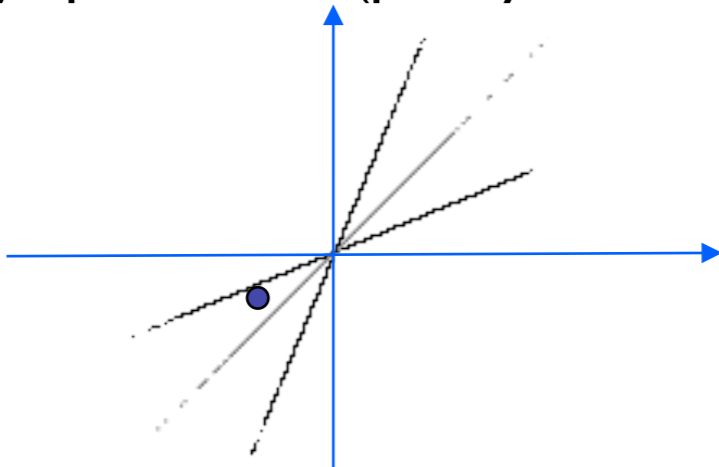
$$\mathbf{b}_n = \mathbf{A}x_n, \quad 1 \leq n \leq N \quad \mathbf{B} = \mathbf{A}X$$

| | ICA (Independent Component Analysis) | SCA (Sparse Component Analysis) |
|-----------------|--|--|
| Model of ... | probability density function $p(X)$ | sample matrix X |
| Assumption | Independence $p(X) = \prod_{nk} p(x_n(k))$ | Sparsity / geometry ★ many zeroes in X ★ x_n and \mathbf{b}_n concentrate around union of low dimensional subspaces |
| Identifiability | Darmois theorem | [Georgiev, Theis & Cichocki 05] [Aharon, Elad & Bruckstein 06] |
| Identification | Contrast functions $\mathbf{A} \sim \hat{\mathbf{W}}^{-1} \quad \hat{\mathbf{W}} := \arg \min_{\mathbf{W}} E_X(f(\mathbf{W}\mathbf{A}X))$ | Combinatorial algorithms |
| Issues | In practice : finite training sets expectation \longrightarrow sample average | Identifiability assumes: ★ highly sparse coefficients ★ (combinatorially ?) many training examples |

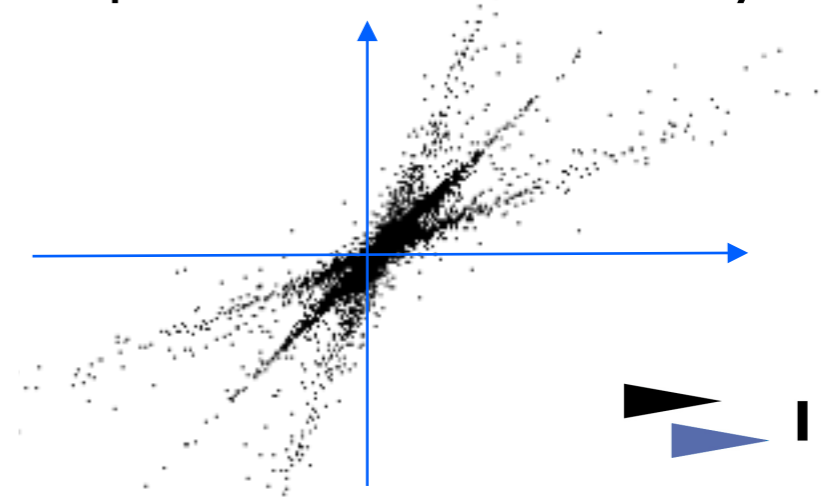
Objectives

- Long term goal :
 - ✦ identifiability conditions on X to recover \mathbf{A} from $\mathbf{B} = \mathbf{A}X$
 - ✦ provably good + efficient identification algorithms
- Focus : exploit sparsity of X
- Desirable features
 - ✦ geometric understanding of identifiability conditions
 - ✦ **robustness** to “weakly-sparse” data
 - ✦ identifiability with **limited number of training samples**
 - ✦ **non-combinatorial** algorithms

Exactly sparse data (purely academic)



Example of real data = “weakly” sparse



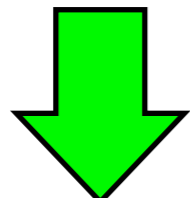
L1 minimization for dictionary learning

Holy grail: provably good + efficient sparse learning

● Sparse representations

- ♦ **Known** matrix \mathbf{A}
- ♦ Data model $\mathbf{b} = \mathbf{A}x_0$

- ♦ Identifiability theorems:
 $\|x_0\|_0 \leq k_1(\mathbf{A})$



$$x_0 = \arg \min_{\mathbf{A}x = \mathbf{b}} \|x\|_1$$

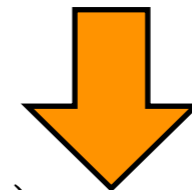
- ♦ Much literature since 2001
(Donoho & Huo, Elad & Bruckstein,
Gribonval & Nielsen, Candès & Romberg
& Tao, Tropp, Donoho & Tanner, ... and
many others)

● Dictionary learning

- ♦ **Unknown** matrix \mathbf{A}_0
- ♦ Data model $\mathbf{B} = \mathbf{A}_0 X_0$

- ♦ Identifiability theorem ?

$$\mathbf{A}_0, X_0 \in \boxed{?}$$

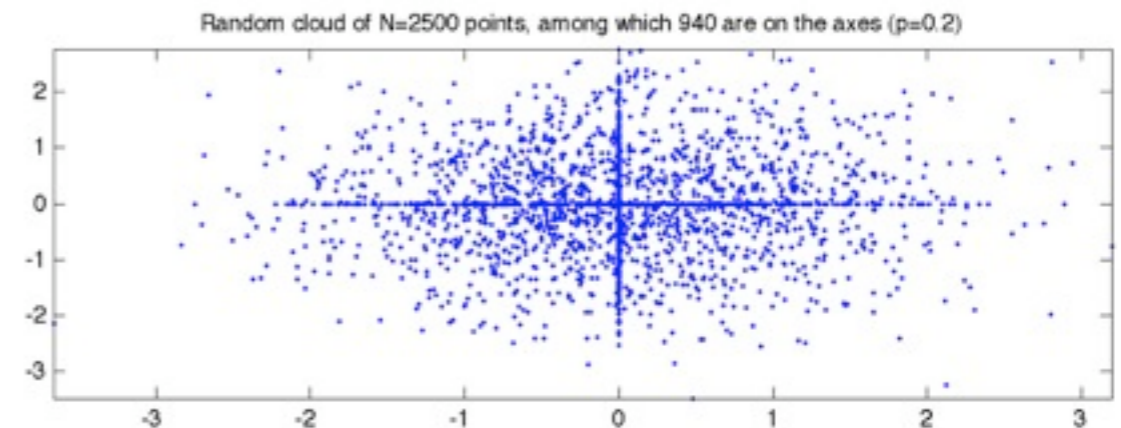


$$(\mathbf{A}_0, X_0) \in \arg \min_{\mathbf{A}X = \mathbf{B}} \|X\|_1$$

- ♦ Most literature on Independent
Component Analysis (ICA),
density model rather than finite
sample size geometric model

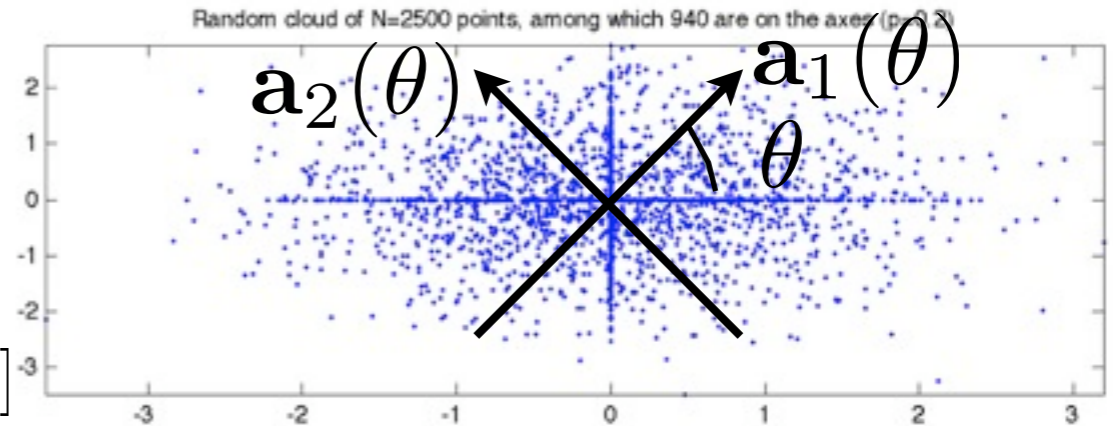
Numerical example

- Cloud of 2500 training samples in \mathbb{R}^2
 - ◆ ~1000 sparse [= on axes]
 - ◆ ~1500 non-sparse



Numerical example

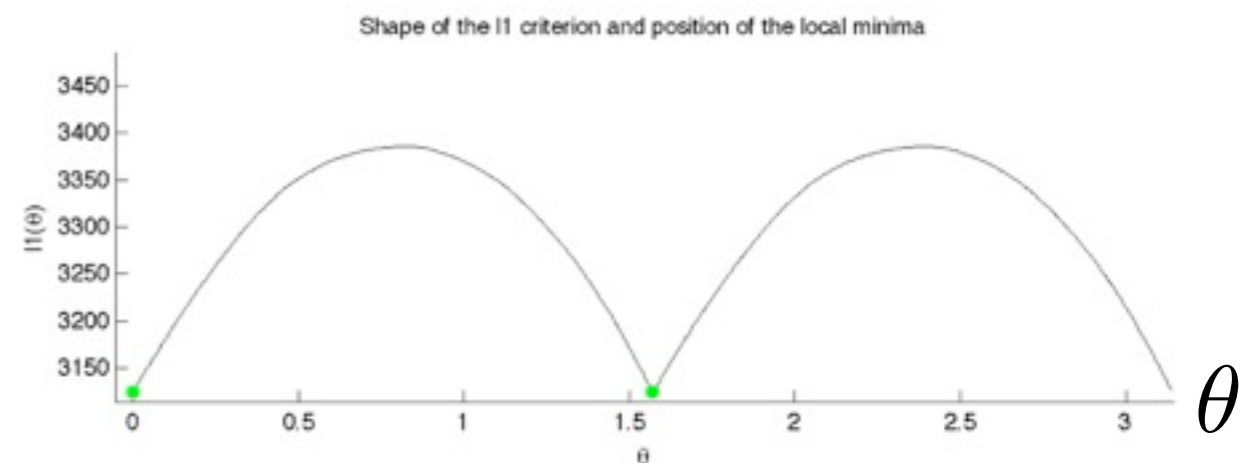
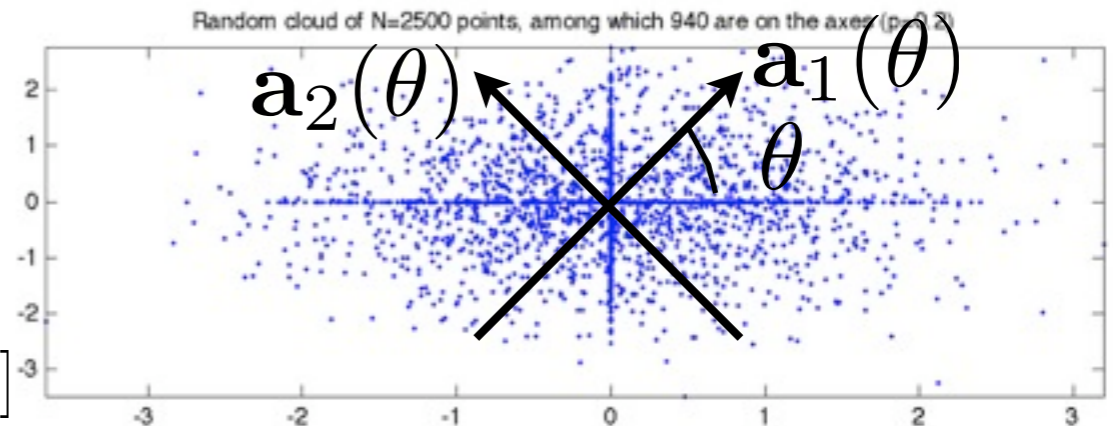
- Cloud of 2500 training samples in \mathbb{R}^2
 - ◆ ~1000 sparse [= on axes]
 - ◆ ~1500 non-sparse
- Orthonormal basis
 - ◆ Angle $\theta \longleftrightarrow A_\theta = [a_1(\theta), a_2(\theta)]$



Numerical example

- Cloud of 2500 training samples in \mathbb{R}^2
 - ♦ ~1000 sparse [= on axes]
 - ♦ ~1500 non-sparse
- Orthonormal basis
 - ♦ Angle $\theta \leftrightarrow A_\theta = [a_1(\theta), a_2(\theta)]$
- LI criterion

$$\|A_\theta^{-1} A_0 X\|_1$$



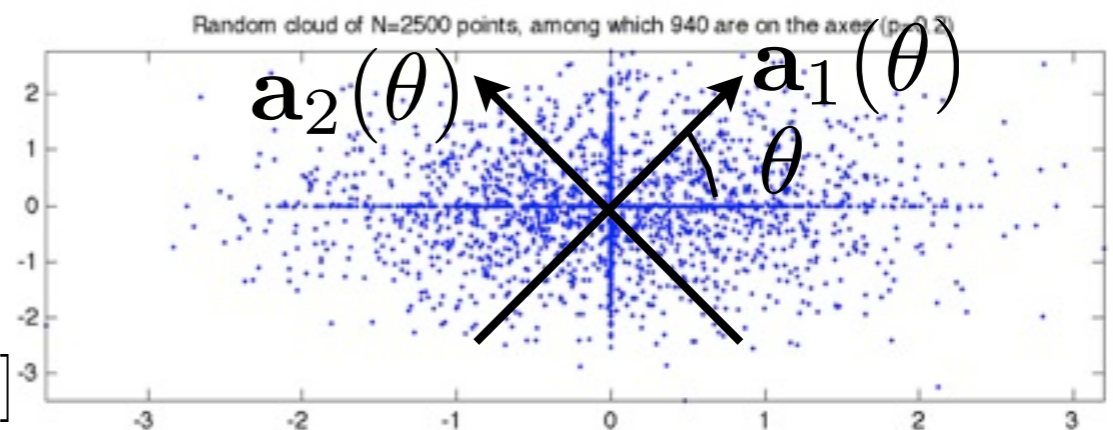
Numerical example

- Cloud of 2500 training samples in \mathbb{R}^2

- ♦ ~1000 sparse [= on axes]
- ♦ ~1500 non-sparse

- Orthonormal basis

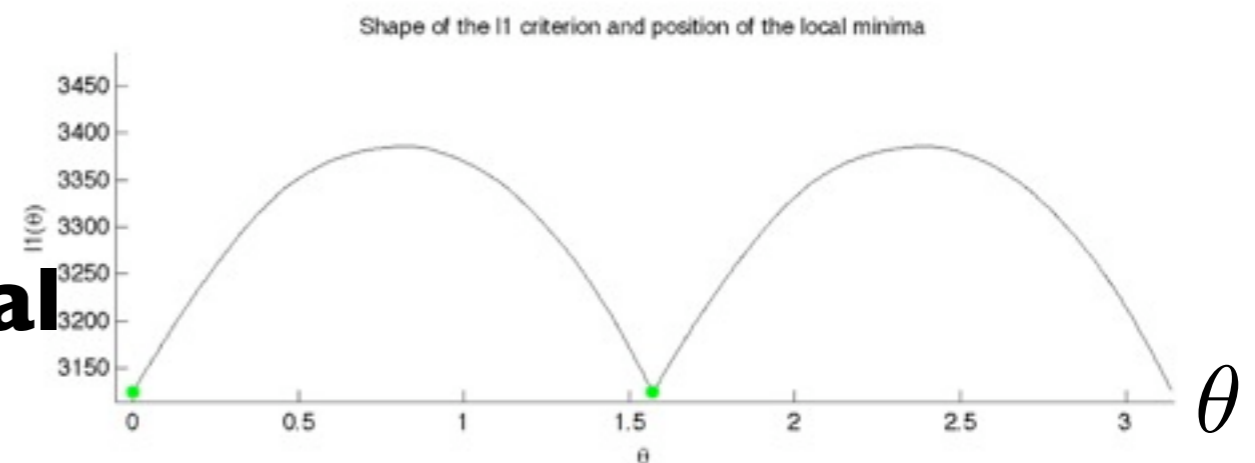
- ♦ Angle $\theta \leftrightarrow A_\theta = [a_1(\theta), a_2(\theta)]$



- LI criterion

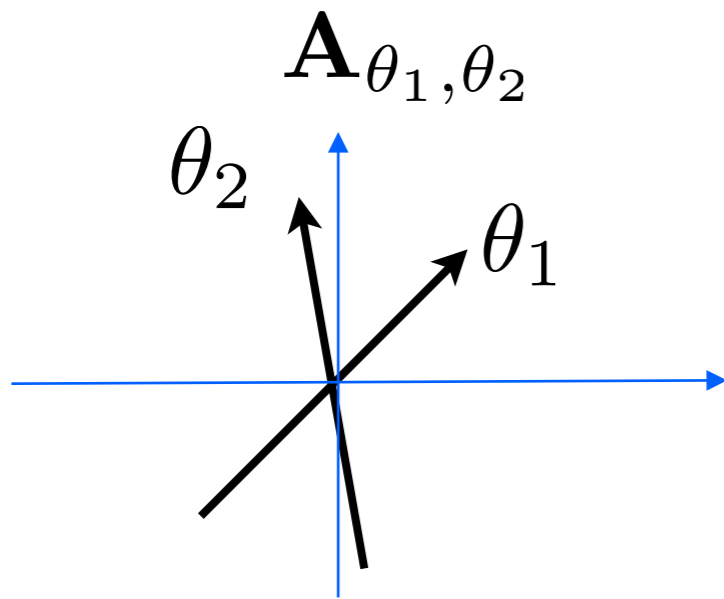
$$\|A_\theta^{-1} A_0 X\|_1$$

- ♦ global optimum=**original**
- ♦ *no other local minimum*

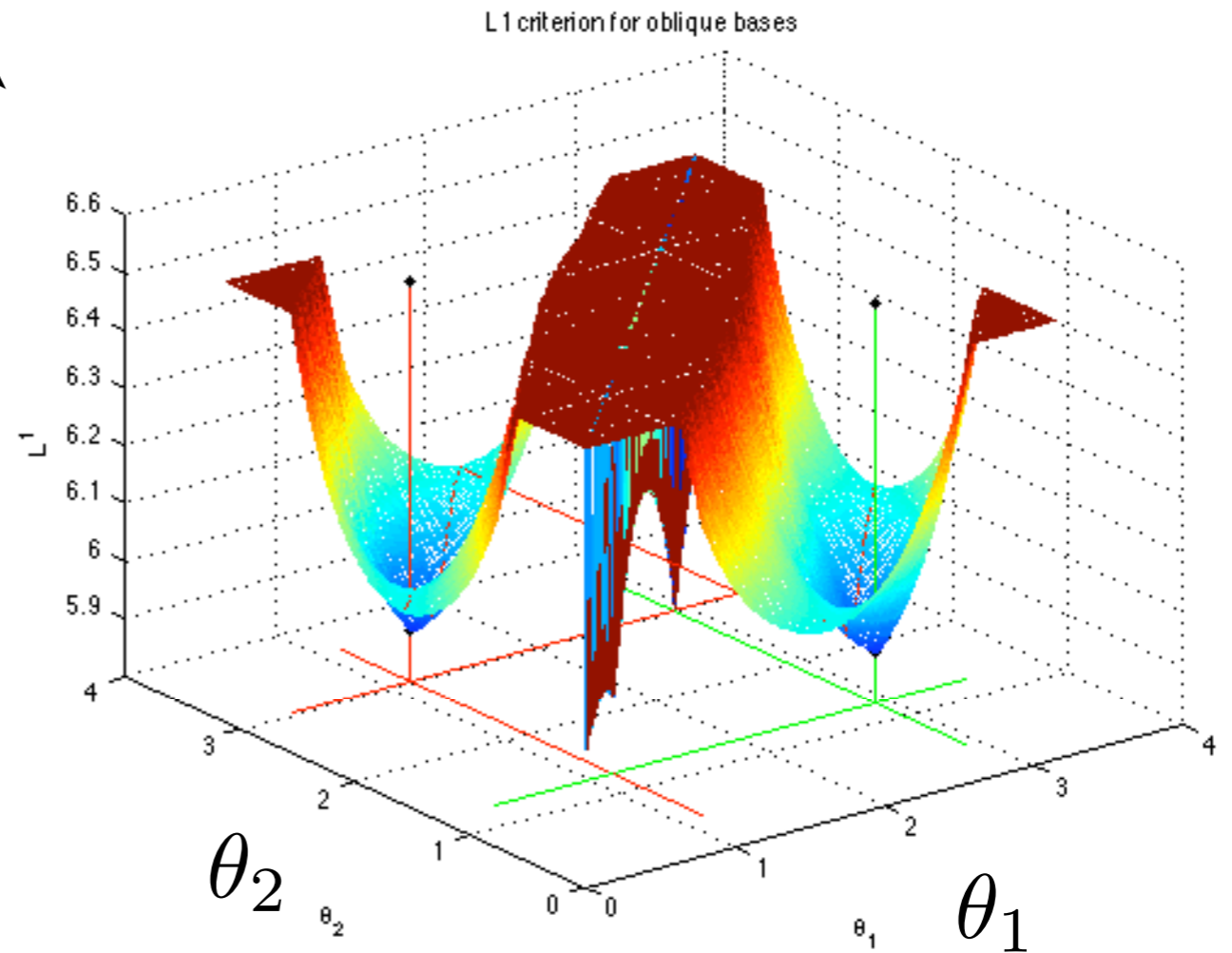


Numerical example

Non orthogonal bases

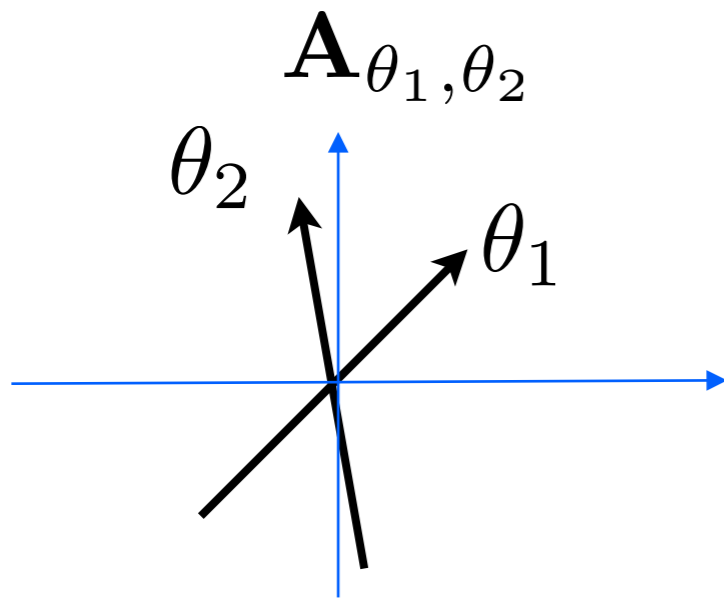


$$\|A_{\theta_1, \theta_2}^{-1} A_0 X\|_1$$

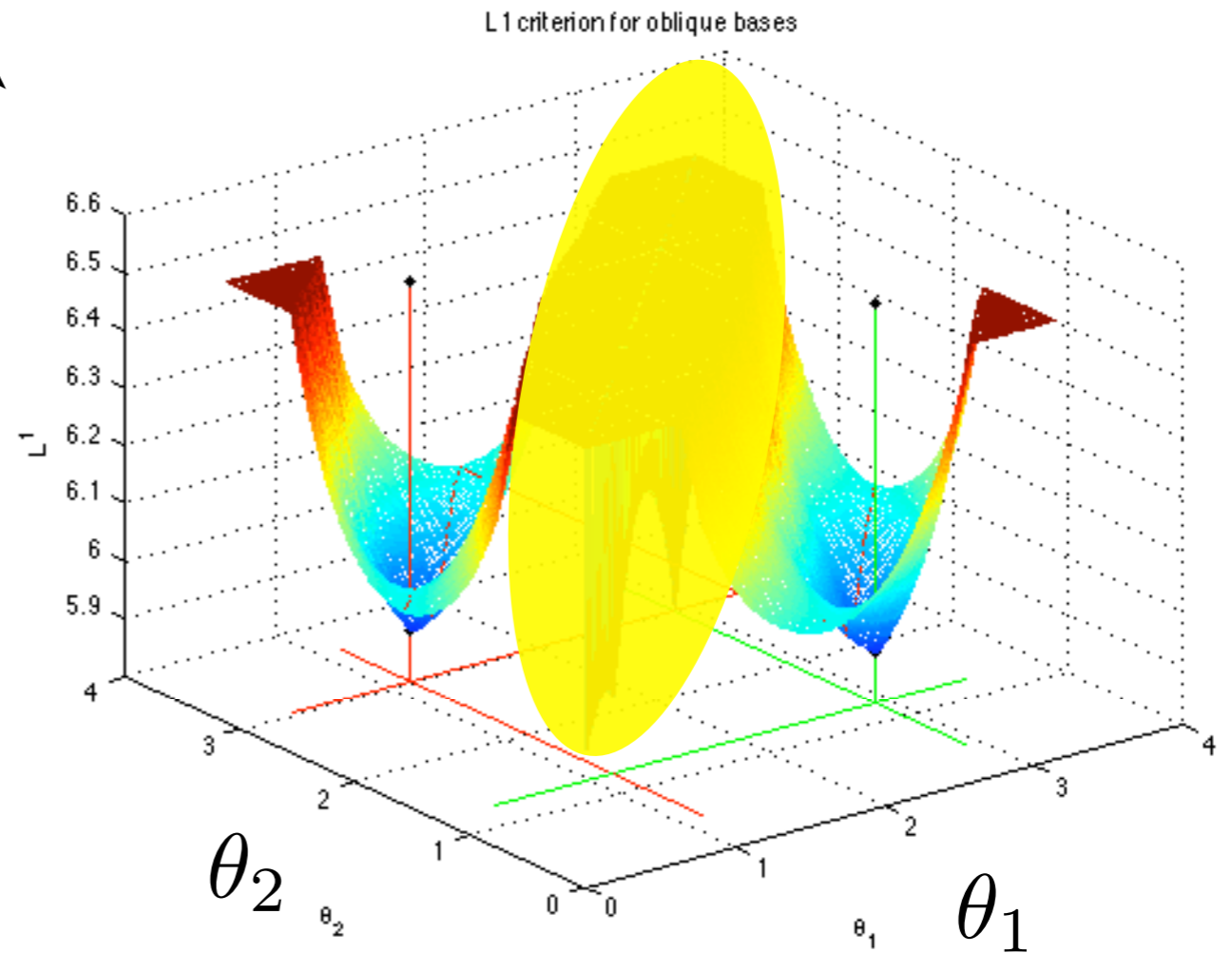


Numerical example

Non orthogonal bases

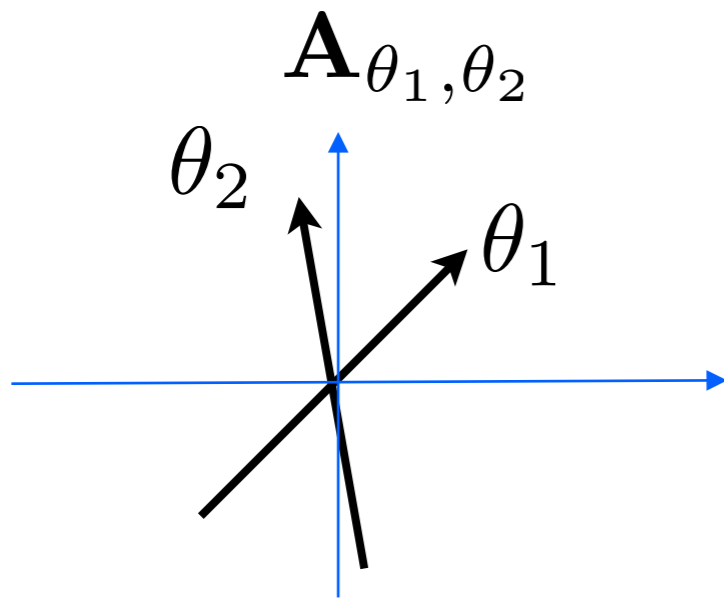


$$\|A_{\theta_1, \theta_2}^{-1} A_0 X\|_1$$

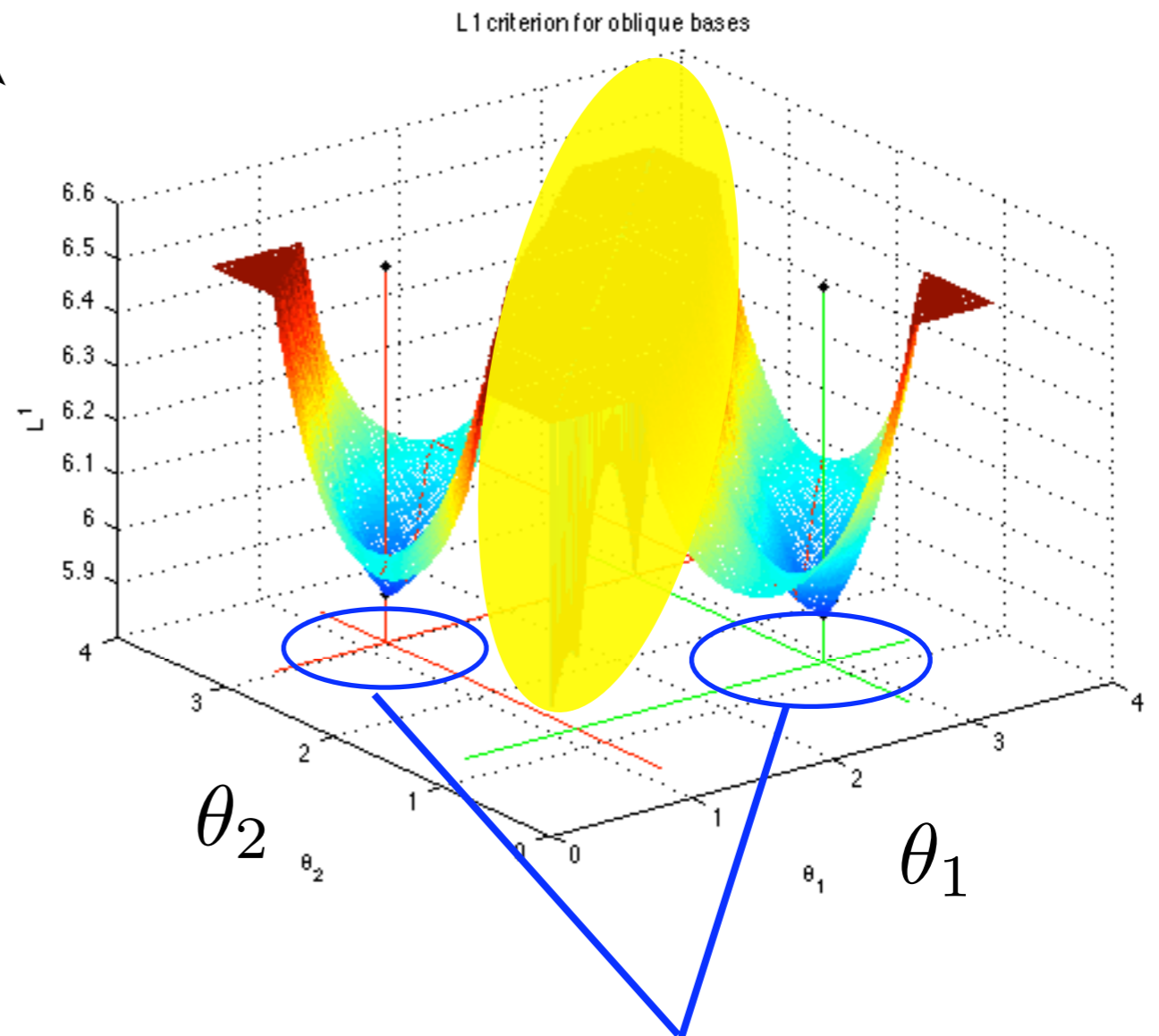


Numerical example

Non orthogonal bases



$$\|A_{\theta_1, \theta_2}^{-1} A_0 X\|_1$$



Empirical observations

- a) Global minima match the original basis
- b) There is no other local minimum.

Theoretical results

1. “Local identifiability” for (non overcomplete) LI dictionary learning

- ♦ algebraic / geometric *characterization of local minima*

2. Probability of identifiability

- ♦ model on X : random, weakly-sparse
- ♦ analysis of identifiability *for (small) finite sample size*

Local identifiability result

- **Two assumptions:**

- ♦ X : for each row k , up to column permutation, has decomposition

$$X = \begin{array}{c|c} & \\ \hline & \\ \hline S_k & 0 \\ \hline X_k & \bar{X}_k \\ \hline \end{array}$$

$\Lambda_k \qquad \bar{\Lambda}_k$

and there exists d_k , $\|d_k\|_\infty < 1$,

$$X_k s_k^T = \bar{X}_k d_k$$

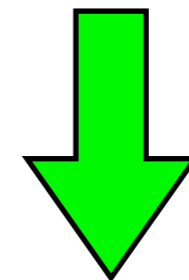
- ♦ A_0 = **basis of sufficiently incoherent unit atoms**

$$\forall k \|a_k\|_2 = 1 \quad \max_{k \neq l} |\langle a_k, a_l \rangle| \ll 1$$

- **Conclusion :**

- ♦ A_0 = **local minimum** of LI among (not necessarily orthonormal) bases

$$\begin{cases} (A', X') \approx (A_0, X) \\ A' X' = A_0 X \end{cases}$$



$$\|X'\|_1 \geq \|X\|_1$$

Trivial example

- **Two assumptions:**

- ♦ X : for each row k , up to column permutation, has decomposition

$$X = \begin{array}{c|c} & \\ \hline & \\ \hline k & \begin{array}{c} S_k \\ 0 \end{array} \\ \hline & \begin{array}{c} X_k \\ \bar{X}_k \end{array} \\ \hline & \end{array}$$

$\Lambda_k \qquad \bar{\Lambda}_k$

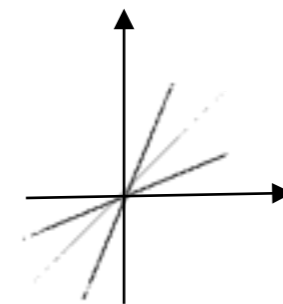
and there exists d_k , $\|d_k\|_\infty < 1$,

$$X_k s_k^T = \bar{X}_k d_k$$

- ♦ $A_0 =$ **basis of sufficiently incoherent unit atoms**

$$\forall k \|a_k\|_2 = 1 \qquad \max_{k \neq l} |\langle a_k, a_l \rangle| \ll 1$$

- If X has at most one nonzero entry per column (at *unknown* positions)



- ♦ *Simply choose*

- How robust is the condition to weakly-sparse outliers ?
- How many samples N does it then typically require ?

Trivial example

- **Two assumptions:**

- ♦ X : for each row k , up to column permutation, has decomposition

$$X = \begin{array}{c|c} & \\ \hline & \\ \hline k & \begin{array}{c} S_k \\ X_k \end{array} \\ \hline \end{array} \begin{array}{c} \\ \\ \hline 0 \\ \hline \bar{X}_k \\ \hline \end{array} \begin{array}{c} \\ \\ \hline =0 \\ \hline \end{array}$$

$\Lambda_k \qquad \bar{\Lambda}_k$

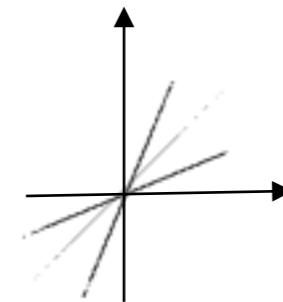
and there exists d_k , $\|d_k\|_\infty < 1$,

$$0 = X_k S_k^T = \bar{X}_k d_k$$

- ♦ $A_0 =$ **basis of sufficiently incoherent unit atoms**

$$\forall k \quad \|a_k\|_2 = 1 \quad \max_{k \neq l} |\langle a_k, a_l \rangle| \ll 1$$

- If X has at most one nonzero entry per column (at *unknown* positions)



- ♦ *Simply choose*

- How robust is the condition to weakly-sparse outliers ?
- How many samples N does it then typically require ?

Trivial example

- **Two assumptions:**

- ♦ X : for each row k , up to column permutation, has decomposition

$$X = \begin{array}{c|c} & \\ \hline k & \\ \hline S_k & 0 \\ \hline X_k & \bar{X}_k \\ \hline \end{array} \quad \begin{array}{c} \\ \\ \\ =0 \\ \end{array}$$

$\Lambda_k \qquad \bar{\Lambda}_k$

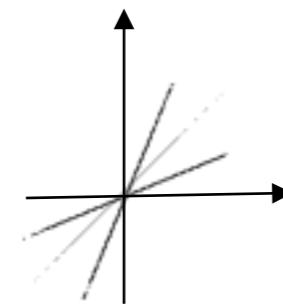
and there exists d_k , $\|d_k\|_\infty < 1$,

$$0 = X_k s_k^T = \bar{X}_k d_k$$

- ♦ A_0 = **basis of sufficiently incoherent unit atoms**

$$\forall k \|a_k\|_2 = 1 \quad \max_{k \neq l} |\langle a_k, a_l \rangle| \ll 1$$

- If X has at most one nonzero entry per column (at *unknown* positions)



- ♦ *Simply choose* $d_k = 0$

Trivial example

- **Two assumptions:**
 - ✦ X : for each row k , up to column permutation, has decomposition

The diagram shows a matrix X being decomposed into two parts, Λ_k and $\bar{\Lambda}_k$. The matrix X is represented as a 2x2 block matrix. The top-left block is labeled S_k and is shaded gray. The top-right block is labeled 0 and is white. The bottom-left block is labeled X_k and is shaded gray. The bottom-right block is labeled \bar{X}_k and is shaded gray. A large red $=0$ is placed between the two columns. Below the matrix, two horizontal lines with vertical tick marks indicate the column partitions, labeled Λ_k and $\bar{\Lambda}_k$ respectively.

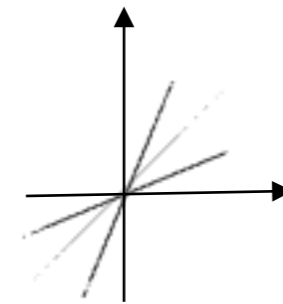
and there exists $d_k, \quad ||d_k||_\infty < 1,$

$$\mathbf{0} = X_k s_k^T = \bar{X}_k d_k$$

- ✦ \mathbf{A}_0 = **basis of sufficiently incoherent unit atoms**

$$\forall k \|a_k\|_2 = 1 \quad \max_{k \neq l} |\langle a_k, a_l \rangle| \ll 1$$

- If X has at most one nonzero entry per column (at *unknown* positions)



- ✦ Simply choose $d_k = 0$

- *How robust is the condition to weakly-sparse outliers ?*

Trivial example

- **Two assumptions:**

- ♦ X : for each row k , up to column permutation, has decomposition

$$X = \begin{array}{c|c} & \\ \hline & \\ \hline k & \begin{array}{c} S_k \\ X_k \end{array} \\ \hline \end{array} \begin{array}{c} \\ \\ \hline 0 \\ \hline \bar{X}_k \\ \hline \end{array} \quad \begin{array}{c} \\ \\ \hline =0 \\ \hline \end{array}$$

$\Lambda_k \qquad \bar{\Lambda}_k$

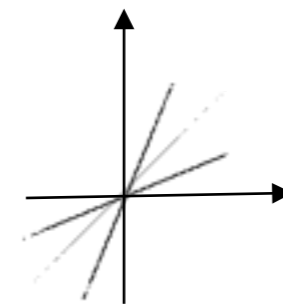
and there exists d_k , $\|d_k\|_\infty < 1$,

$$0 = X_k S_k^T = \bar{X}_k d_k$$

- ♦ $A_0 =$ **basis of sufficiently incoherent unit atoms**

$$\forall k \quad \|a_k\|_2 = 1 \quad \max_{k \neq l} |\langle a_k, a_l \rangle| \ll 1$$

- If X has at most one nonzero entry per column (at *unknown* positions)

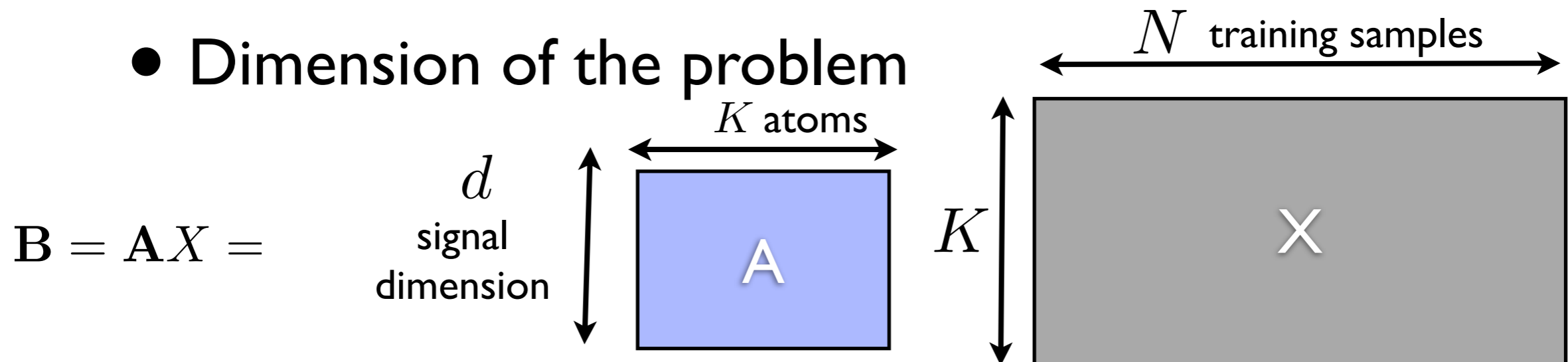


- ♦ *Simply choose* $d_k = 0$

- How robust is the condition to weakly-sparse outliers ?
- How many samples N does it then typically require ?

How many training samples ?

- Dimension of the problem



- General dictionary $K \geq d$, basis $K = d$

- Required number of training samples:

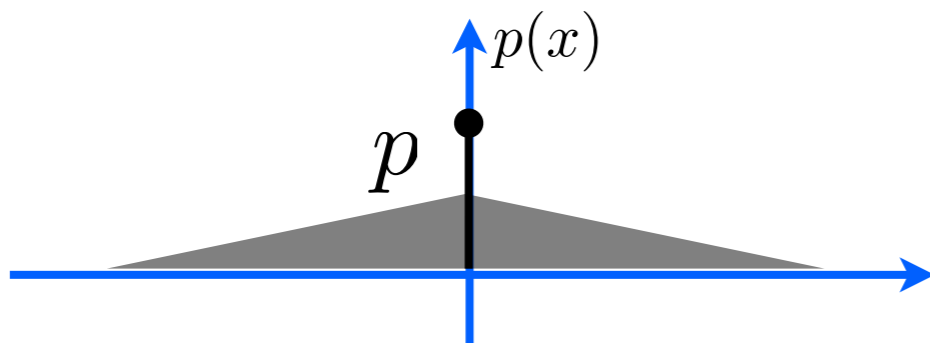
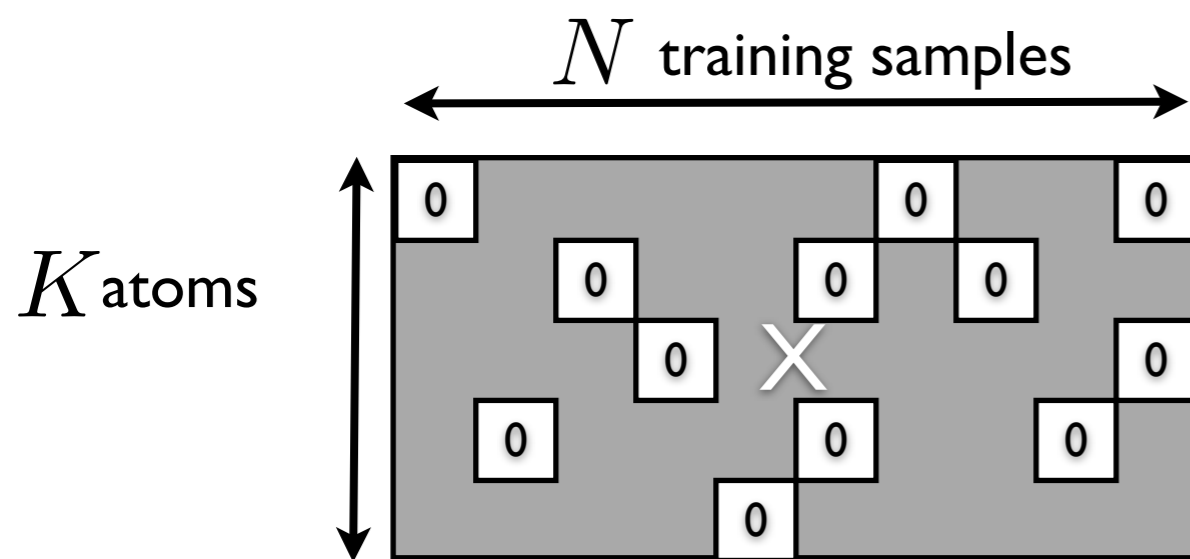
✦ With $N = K$, maximum sparsity achieved for $\hat{\mathbf{A}} = \mathbf{B} \neq \mathbf{A}$ $\hat{\mathbf{X}} = \mathbf{Id}$

1 atom $\hat{\mathbf{a}}_k =$ 1 training sample \mathbf{b}_n

- ✦ Identifiability from $N \leq CK \log K$ samples for all “nice” \mathbf{A} ?
- ✦ Identifiability with weakly-sparse \mathbf{X} ?

Second result : probability of identifiability

- Random model $X = (x_{kn})$
 - ♦ i.i.d. (sub)Gaussian entries in \mathbb{R}^K
 - ♦ a fraction p set to zero at random



- Using concentration of measure :

Probability of failure ...

$$P(\text{😞}) \leq C \exp(aK \log K - bN)$$

Conclusion

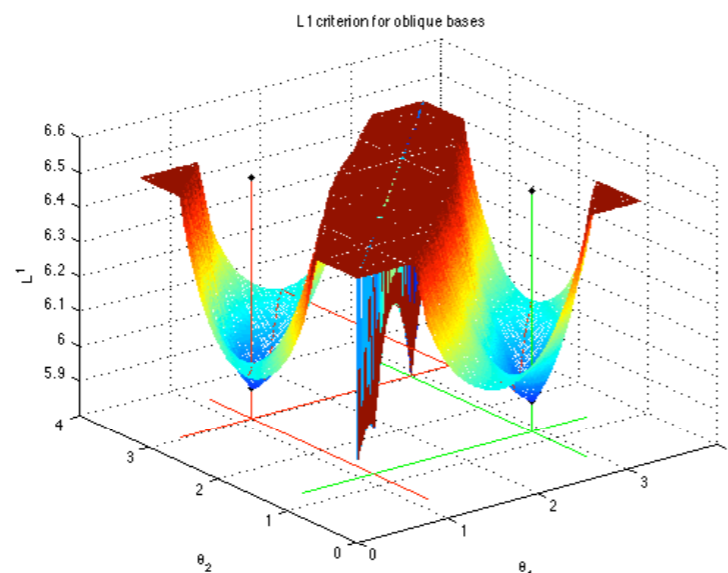
Local identifiability guaranteed with high probability from only “few” training samples:

$$N \geq C(p) \cdot K \log K$$

(almost linear in dimension K , even for small p)

Summary

- L1-minimization for dictionary learning:
 - ✦ Sufficient condition for local identifiability of bases
 - ✦ Condition typically valid
 - ❖ even if only weakly-sparse training samples
 - ❖ even with relatively few training samples (non combinatorial training set)
$$N \geq C(p) \cdot K \log K$$
- Consequence :
 - ✦ ideal convergence of descent algorithms *conditionally on good initialization*
 - ✦ *conjecture : with high probability, no spurious local minima*



Perspectives & challenges

- Main open questions:

- ❖ Probability of spurious local minima
- ❖ Optimization algorithm (LI criterion is nonconvex ...), in progress
- ❖ Stability/robustness to noise / compressible X ?

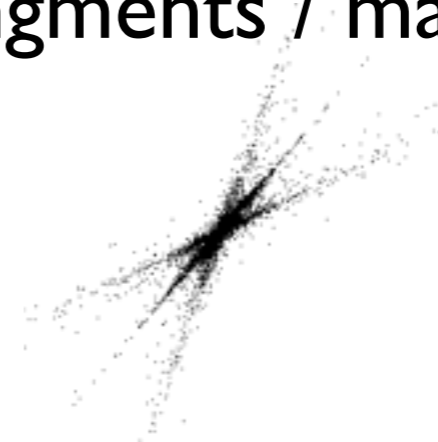
- Extensions:

- ◆ other learning paradigms: efficiency? equivalence?

- ❖ greedy approaches (“deflation”, ongoing work)
- ❖ alternate optimization (MOD, K-SVD, ...)

- ◆ blind sparse deconvolution

- ◆ learning general subspace arrangements / manifolds [cf Yi Ma]

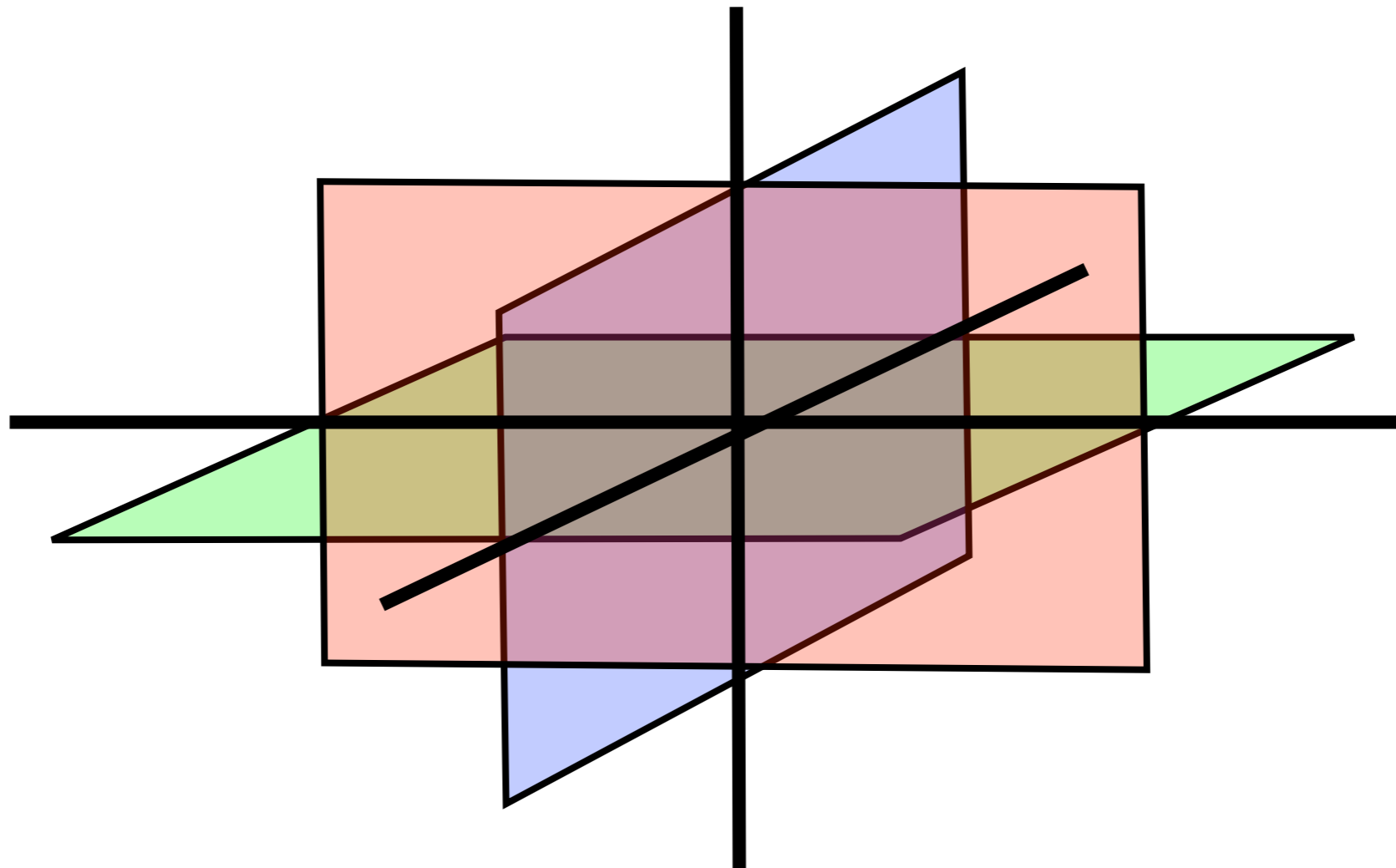


THE END

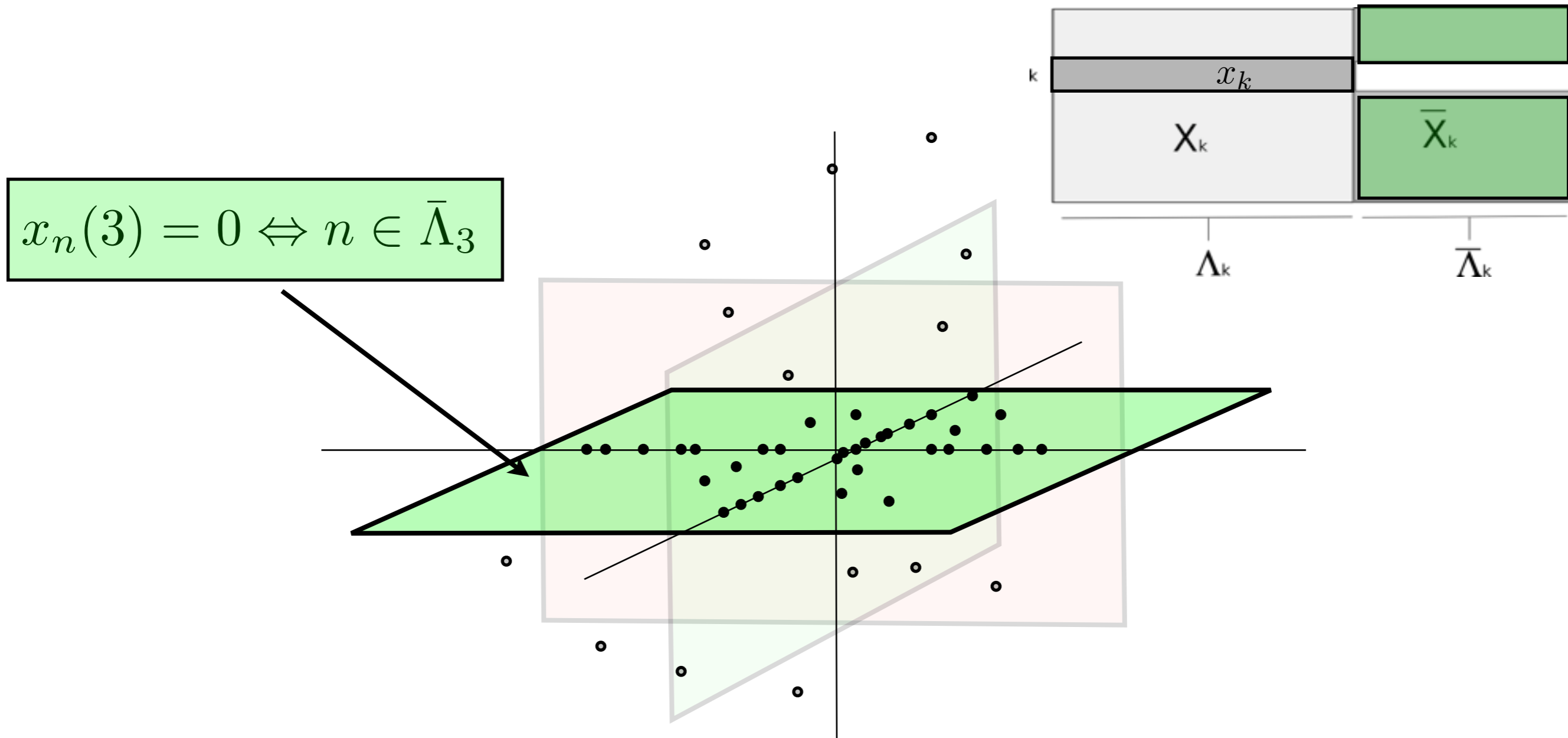
remi.gribonval@inria.fr

Geometric interpretation

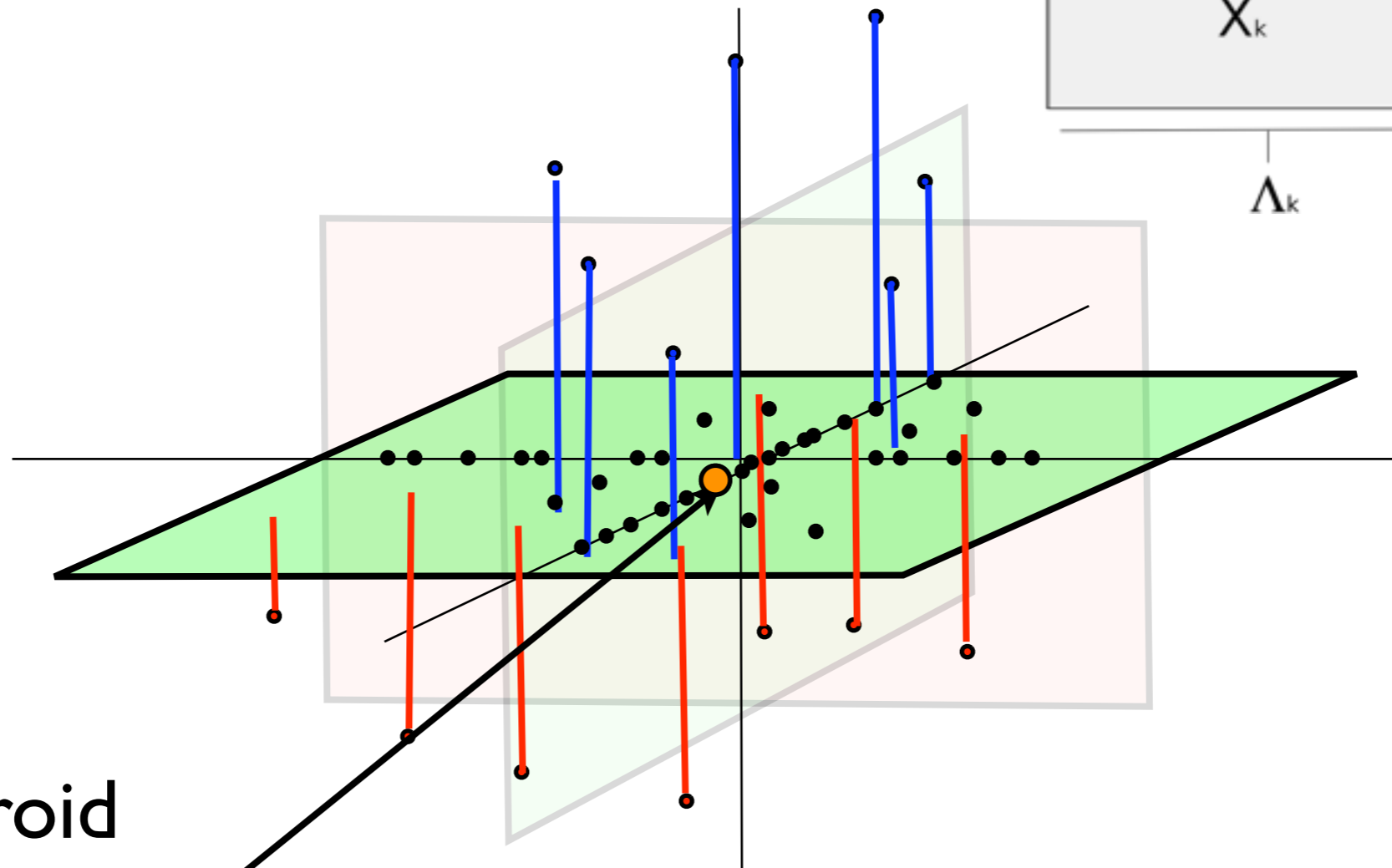
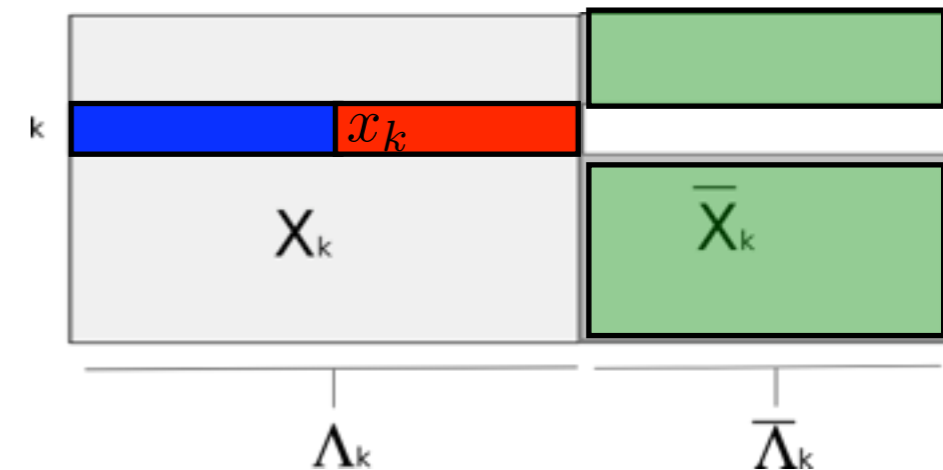
- Many sparse training examples lie on low-dimensional subspaces



Geometric interpretation



Geometric interpretation



~ centroid

$$X_k \cdot \text{sign}(x_k)^T$$

Geometric interpretation

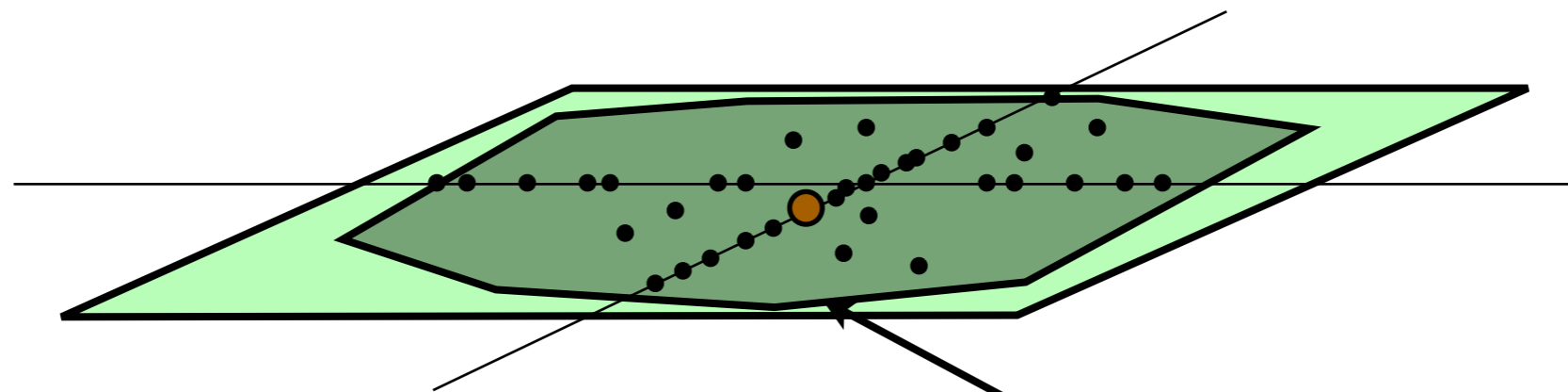
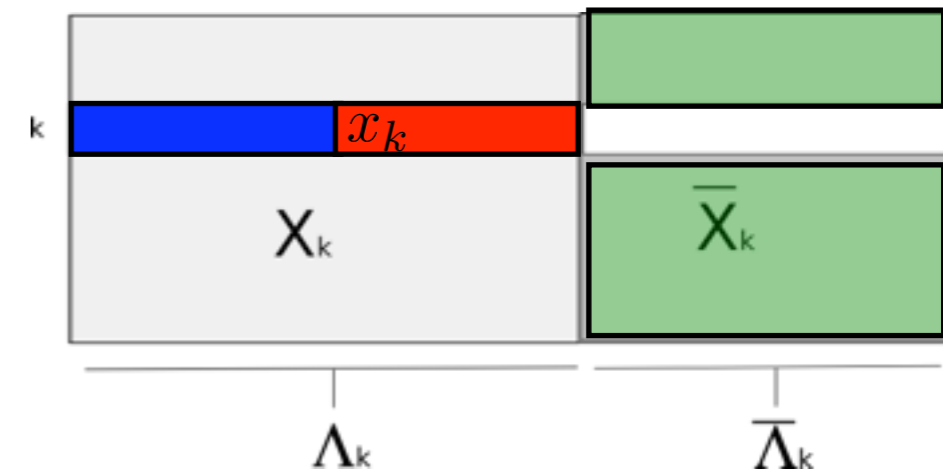


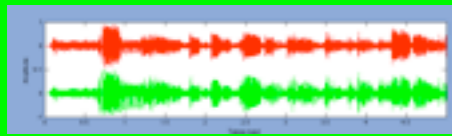
Image of unit hypercube by \bar{X}_3

$$X_k \cdot \text{sign}(x_k)^T = \bar{X}_k d_k, \|d_k\|_\infty < 1$$

Sparse data models ?

- Sparse signals, sparse images ...
 - ✦ challenge = **large-scale algorithms**

Sparse models = Fourier, *lets, ...



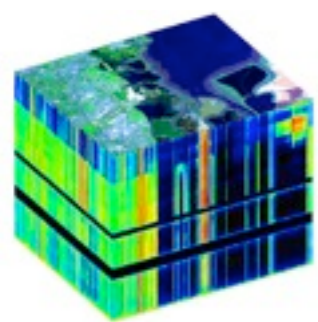
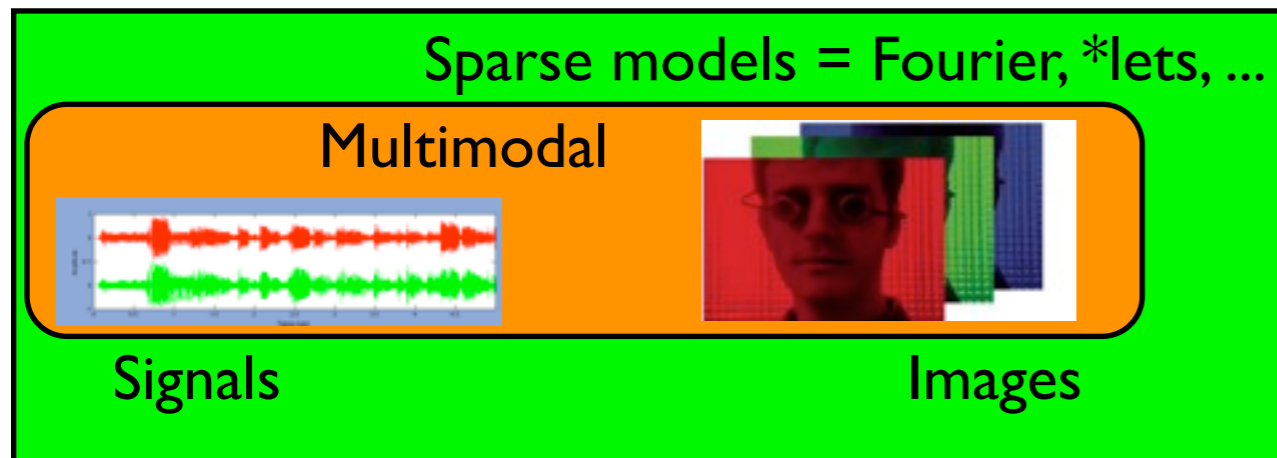
Signals



Images

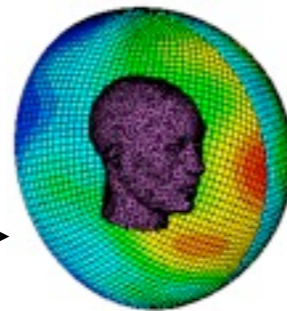
Sparse data models ?

- Sparse signals, sparse images ...
 - ✦ challenge = **large-scale algorithms**

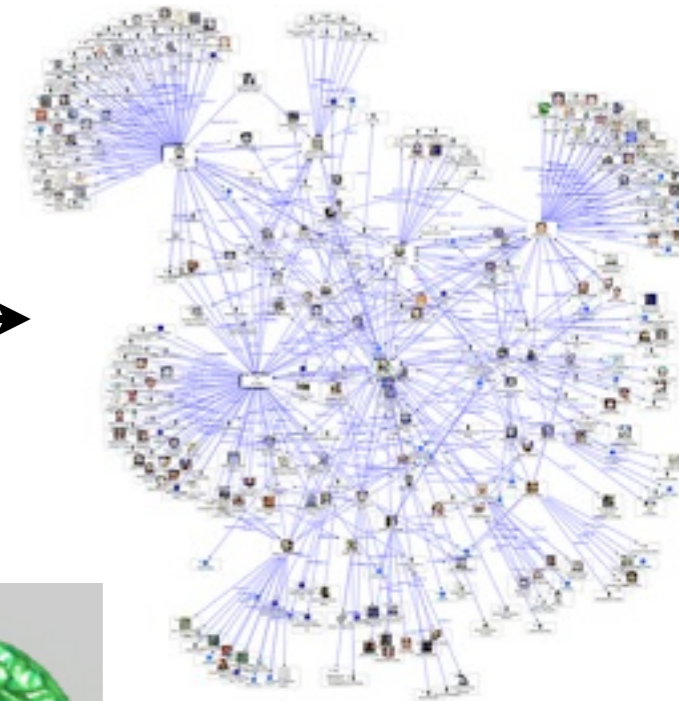
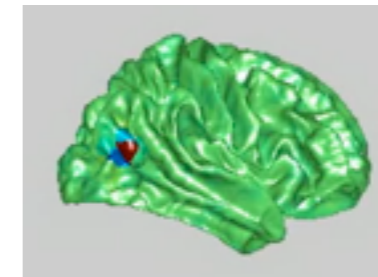


Hyperspectral
Satellite imaging

Spherical geometry
Cosmology, HRTF (3D audio)



Data on graphs
Social networks
Brain connectivity



Vector valued
Diffusion tensor

- “Exotic” or complex data
 - ✦ challenge = **generic models / model building tools**

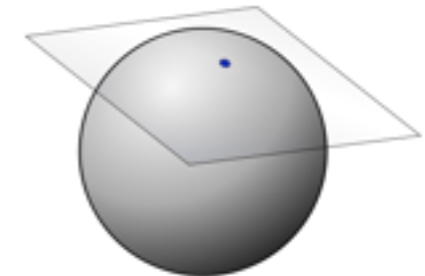
Local identifiability analysis (I)

- Normalization convention needed (as $\mathbf{A}X = \mathbf{B} \Rightarrow (2\mathbf{A})(X/2) = \mathbf{B}$)

$$\mathbf{A} \in \mathcal{A} := \{\|\mathbf{a}_k\|_2 = 1\}$$

- Compatible basis perturbations : tangent plane

$$\mathbf{A} + \delta_{\mathbf{A}} \quad \delta_{\mathbf{A}} \in T_{\mathbf{A}}\mathcal{A}$$



- Coupling between basis and coefficients

$$\mathbf{A}X = \mathbf{B} \quad \Rightarrow \quad \delta_{\mathbf{A}} \cdot X + \mathbf{A} \cdot \delta_X = 0 \quad \Rightarrow \quad \delta_X = -\mathbf{A}^{-1} \delta_{\mathbf{A}} \cdot X$$

Local identifiability analysis (2)

- First order approx. of LI criterion (Λ = support of X)

$$\|X + \delta_X\|_1 - \|X\|_1 \approx \langle \delta_X, \text{sign}(X) \rangle + \|(\delta_X)_{\bar{\Lambda}}\|_1$$

$$\forall \delta_X, |\langle \delta_X, \text{sign}(X) \rangle| < \|(\delta_X)_{\bar{\Lambda}}\|_1 \rightarrow \text{local minimum}$$

- Admissible perturbations (from previous slide)

$$\delta_X = -\mathbf{A}^{-1} \delta_{\mathbf{A}} \cdot X \quad \delta_{\mathbf{A}} \in T_{\mathbf{A}} \mathcal{A}$$

- (...) local minimum iff for zero-diagonal matrices Z

$$\forall Z \neq 0, |\langle Z, X \text{sign}(X)^T - \text{diag}(\|x_k\|_1) \mathbf{A}^T \mathbf{A} \rangle| < \|(ZX)_{\bar{\Lambda}}\|_1$$

- Decoupling between rows: for orthonormal \mathbf{A} ...

$$|\langle z, X \text{sign}(x_k)^T \rangle| < \|(X^T z)_{\bar{\Lambda}_k}\|_1$$

Dictionary learning is *not* about ...

- Channel estimation with **known** pulse sequence

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

- ♦ x = **known** channel input
- ♦ \mathbf{b} = **observed** channel output
- ♦ \mathbf{A} = **unknown** channel response

- [Ex: Pfander, Rauhut & Tanner, "Identification of Matrices having a Sparse Representation", 2008]

$$\mathbf{A} = \sum_k \alpha_k \mathbf{A}_k$$

$$\mathbf{b} = \sum_k \alpha_k (\mathbf{A}_k x)$$

Dictionary learning is *not* about ...

- Channel estimation with **known** sequence

♦ \mathcal{X} = **known** channel

♦ \mathbf{b} = **observed**

♦ \mathbf{A} = **unknown** response

- [Example] Identification of Matrices having a Sparse

$$\mathbf{b} = \sum_k \alpha_k \mathbf{A}_k$$

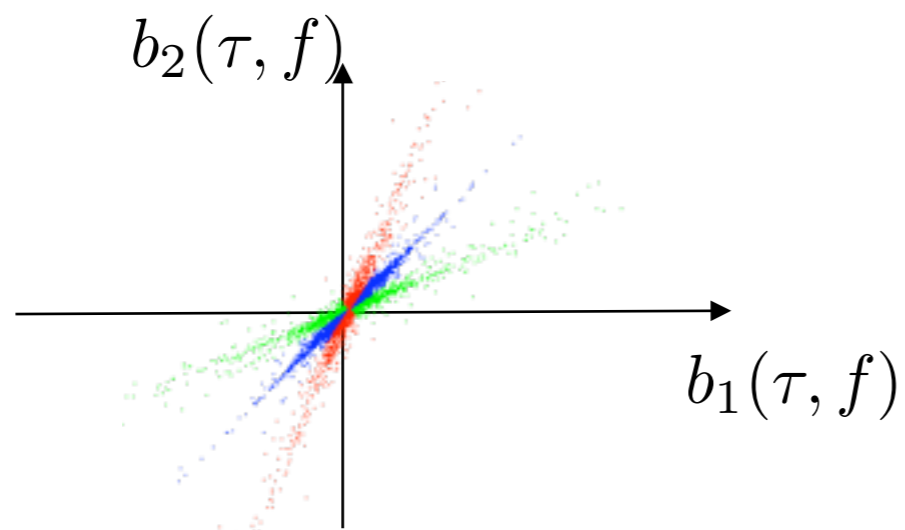
$$\mathbf{b} = \sum_k \alpha_k (\mathbf{A}_k$$

Blind Source Separation

- Mixing model in the time-frequency domain

$$\begin{matrix} b_1(\tau, f) \\ b_2(\tau, f) \end{matrix} \begin{pmatrix} \text{[Spectrogram 1]} \\ \text{[Spectrogram 2]} \end{pmatrix} = \mathbf{A} \mathbf{X}(\tau, f)$$

- And “miraculously” ...



... time-frequency representations of audio signals are (often) **almost disjoint**.

Identifiability of mixing matrix \mathbf{A}
= **geometric properties** of the scatter plot
(concentration along lines)
= thanks to **sparsity** of \mathbf{X}