

# When Are Nonconvex Problems Not Scary?

*A few friendly nonconvex optimization problems*

---

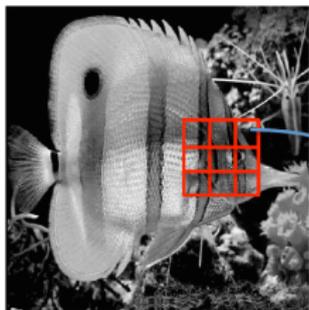
**Ju Sun**

Department of Mathematics

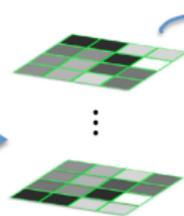
Stanford University

Joint with **Qing Qu** (Columbia U.), **John Wright** (Columbia U.)

# A curious experiment



An image

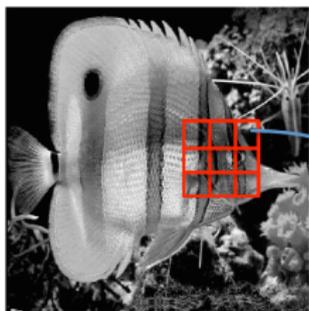


Patches

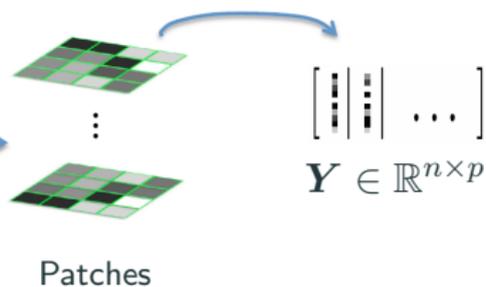
$$\begin{bmatrix} | & | & \dots \\ | & | & \dots \\ | & | & \dots \end{bmatrix}$$

$\mathbf{Y} \in \mathbb{R}^{n \times p}$

# A curious experiment

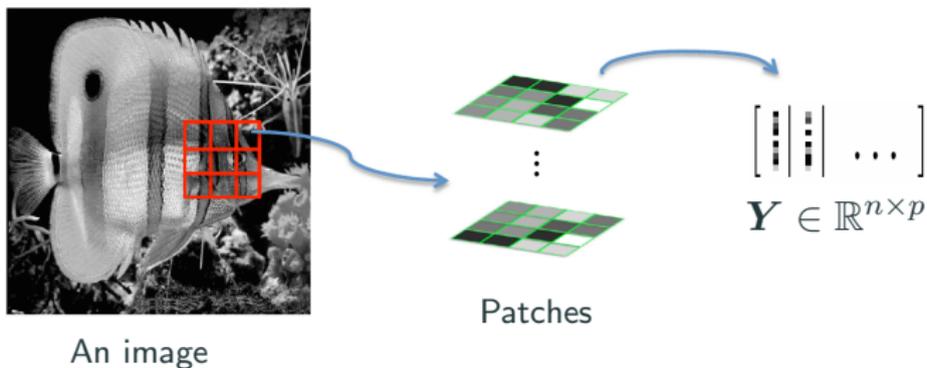


An image



Seek a **concise approximation**:  $Y \approx QX$ , with  $Q \in O_n$  and  $X$  as sparse as possible.

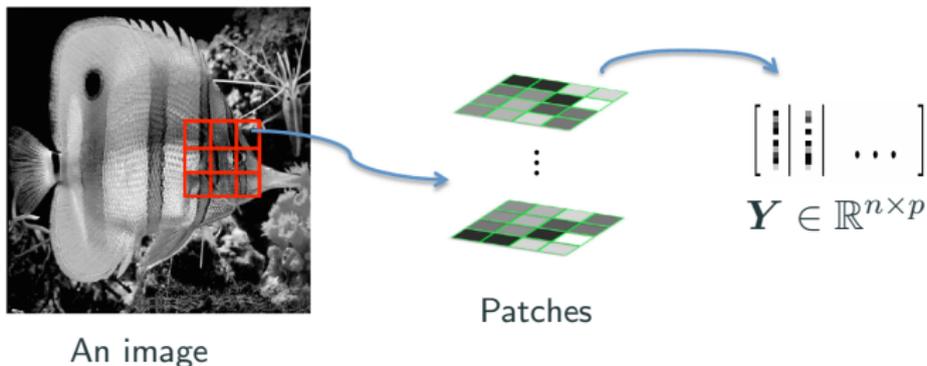
## A curious experiment



Seek a **concise approximation**:  $\mathbf{Y} \approx \mathbf{Q}\mathbf{X}$ , with  $\mathbf{Q} \in O_n$  and  $\mathbf{X}$  as sparse as possible.

... by solving  $\min \frac{1}{2} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1$ , s.t.  $\mathbf{Q} \in O_n$ .

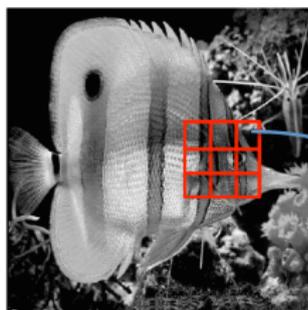
# A curious experiment



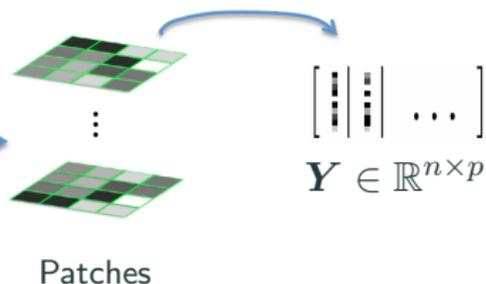
$$\min \quad f(\mathbf{Q}, \mathbf{X}) \doteq \frac{1}{2} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1, \quad \text{s.t.} \quad \mathbf{Q} \in O_n.$$

- The map  $(\mathbf{Q}, \mathbf{X}) \mapsto \mathbf{Q}\mathbf{X}$  is **bilinear**
  - **Combinatorially many global minimizers:**  $(\mathbf{Q}, \mathbf{X})$  or  $(\mathbf{Q}\mathbf{\Pi}, \mathbf{\Pi}^* \mathbf{X})$  ( $2^n n!$  signed permutations  $\mathbf{\Pi}$ )
- Orthogonal group  $O_n$  is a **nonconvex** set

## A curious experiment



An image



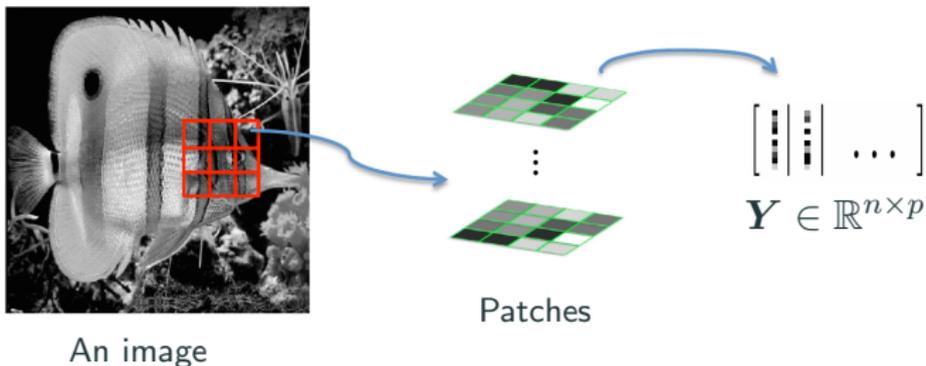
$$\min \quad f(\mathbf{Q}, \mathbf{X}) \doteq \frac{1}{2} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1, \quad \text{s.t. } \mathbf{Q} \in O_n$$

Naive **alternating directions**: starting from a random  $\mathbf{Q}_0 \in O_n$

$$\mathbf{X}_k = \arg \min_{\mathbf{X}} f(\mathbf{Q}_{k-1}, \mathbf{X})$$

$$\mathbf{Q}_k = \arg \min_{\mathbf{Q}} f(\mathbf{Q}, \mathbf{X}_k), \quad \text{s.t. } \mathbf{Q} \in O_n.$$

## A curious experiment



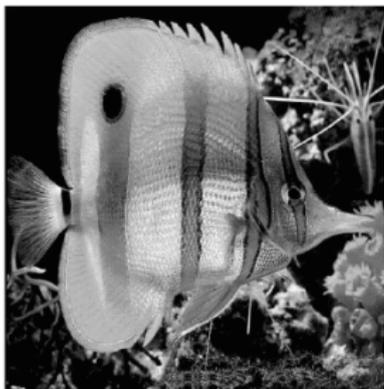
$$\min f(\mathbf{Q}, \mathbf{X}) \doteq \frac{1}{2} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1, \quad \text{s.t. } \mathbf{Q} \in O_n$$

Naive **alternating directions**: starting from a random  $\mathbf{Q}_0 \in O_n$

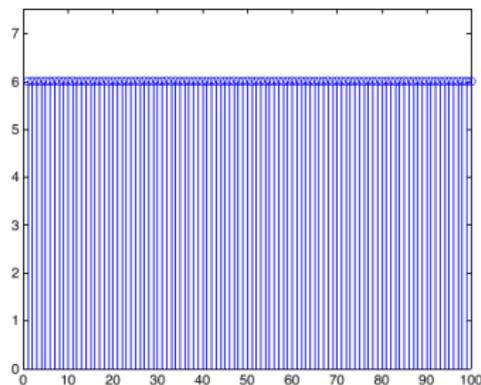
$$\mathbf{X}_k = \mathcal{S}_\lambda [\mathbf{Q}_{k-1}^* \mathbf{Y}]$$

$$\mathbf{Q}_k = \mathbf{U}\mathbf{V}^*, \text{ where } \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* = \text{SVD}(\mathbf{Y}\mathbf{X}_k^*).$$

# A curious experiment



An image



Final  $f(Q_\infty, X_\infty)$ , varying  $Q_0$ .

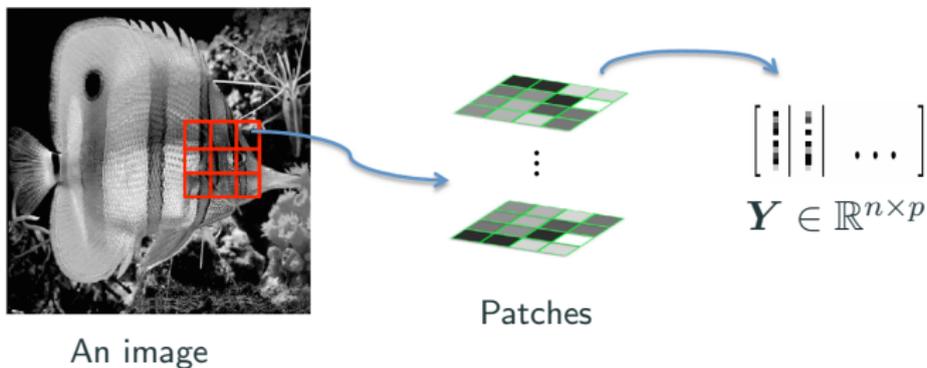
$$\min f(Q, X) \doteq \frac{1}{2} \|Y - QX\|_F^2 + \lambda \|X\|_1, \quad \text{s.t. } Q \in O_n$$

Naive **alternating directions**: starting from a random  $Q_0 \in O_n$

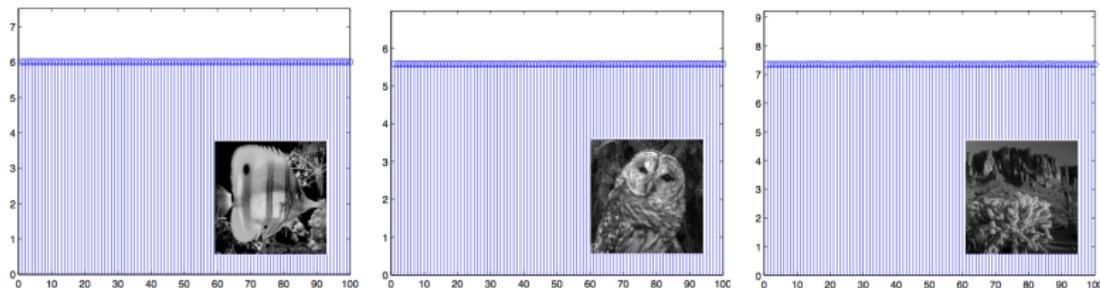
$$X_k = \mathcal{S}_\lambda [Q_{k-1}^* Y]$$

$$Q_k = UV^*, \text{ where } U\Sigma V^* = \text{SVD}(YX^*).$$

# Global solutions to feature learning on real images?



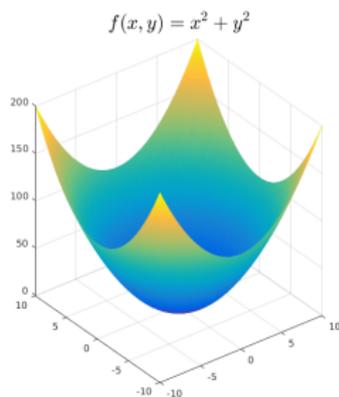
$$\min f(Q, X) \doteq \frac{1}{2} \|Y - QX\|_F^2 + \lambda \|X\|_1, \quad \text{s.t. } Q \in O_n$$



# Nonconvex optimization

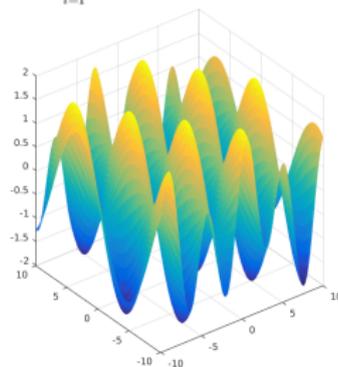
Many problems in modern **signal processing**, **machine learning**, **statistics**, ..., are most naturally formulated as **nonconvex** optimization problems.

$$\begin{aligned} \min f(\mathbf{x}) \\ \text{s. t. } \mathbf{x} \in \mathcal{D}. \end{aligned}$$



“easy”

$$g(x, y) = \sum_{i=1}^2 a_i \sin(b_i x + c_i y) + d_i \cos(e_i x + f_i y)$$



“hard”

**Nonconvex: Even computing a local minimizer is NP-hard!**  
(see, e.g., [Murty and Kabadi, 1987])

# This talk

**In practice:** Heuristic algorithms are often surprisingly successful.

**In theory:** Even computing a local minimizer is NP-hard!

*Which nonconvex optimization problems are easy?*

## Working hypothesis

- Certain nonconvex optimization problems have a **benign structure** when the input data are **large** and **random/generic**.
- This benign structure allows "**initialization-free**" iterative methods to **efficiently** find a **global** minimizer.

$\mathcal{X}$  functions

Examples from practical problems

Sparse (complete) dictionary learning [S., Qu, Wright, '15]

Generalized phase retrieval [S., Qu, Wright, '16]

Other examples in the literature

Numerical optimization methods

Comparison with alternatives

## $\mathcal{X}$ functions

Examples from practical problems

Sparse (complete) dictionary learning [S., Qu, Wright, '15]

Generalized phase retrieval [S., Qu, Wright, '16]

Other examples in the literature

Numerical optimization methods

Comparison with alternatives

# A classical example: the Rayleigh quotient

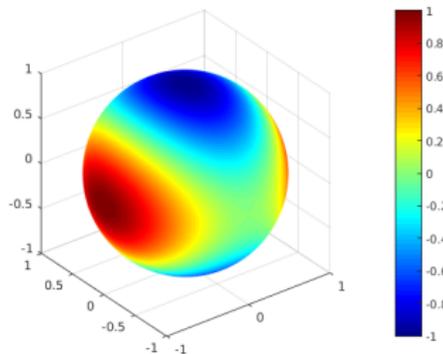
For a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,

$$\min \mathbf{x}^\top \mathbf{A} \mathbf{x} \quad \text{s.t.} \quad \|\mathbf{x}\|_2 = 1.$$

Let  $\mathbf{v}_i$  the eigenvectors of  $\mathbf{A}$ ,  $\lambda_i$  the eigenvalues. Suppose

$$\lambda_1 > \lambda_2 \geq \dots \lambda_{n-1} > \lambda_n.$$

- Only **global minimizers** are  $\pm \mathbf{v}_n$
- Only **global maximizers** are  $\pm \mathbf{v}_1$
- All  $\{\pm \mathbf{v}_i\}$  for  $2 \leq i \leq n - 1$  are **saddle points** with a **directional negative curvature**.



$$\mathbf{A} = \text{diag}(1, 0, -1)$$

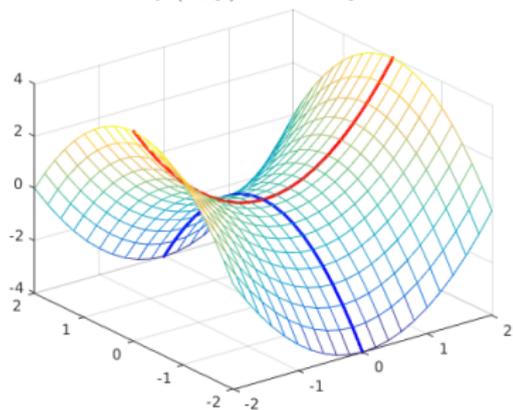
$\mathcal{X}$  functions (qualitative version):

- (P-1) All local minimizers are also global
- (P-2) All saddle points have a directional negative curvature

Thanks to (P-1), focus on finding a local minimizer!

## More on (P-2): Saddle points

$$f(x, y) = x^2 - y^2$$

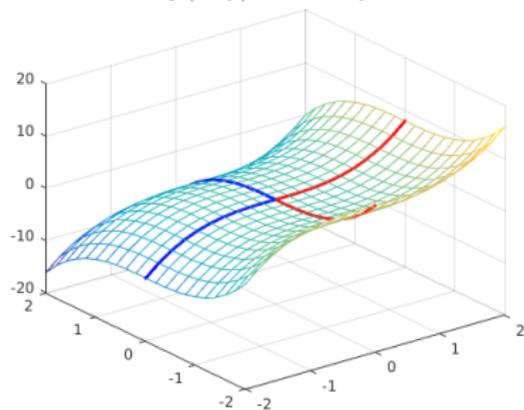


$$\nabla^2 f = \text{diag}(2, -2)$$

**Ridable saddle**

(**strict saddle** [Ge et al., 2015])

$$g(x, y) = x^3 - y^3$$



$$\nabla^2 f = \text{diag}(6x, -6y)$$

local shape determined by  
high-order derivatives around 0

## More on (P-2): Ridable-saddle functions

Consider twice continuously differentiable function  $f : \mathcal{M} \mapsto \mathbb{R}$ , where  $\mathcal{M}$  is a Riemannian manifold.

### (P-2)+

- (P-2A) For all local minimizers,  $\text{Hess } f \succ \mathbf{0}$ , and
  - (P-2B) For all other critical points,  $\lambda_{\min}(\text{Hess } f) < 0$ .
- 
- (P-2A)  $\implies$  local strong convexity around any local minimizer
  - (P-2B)  $\implies$  local directional strict concavity around local maximizers and **saddle points**; particularly, **all saddles are ridable (strict)**.

### Definition

A smooth function  $f : \mathcal{M} \mapsto \mathbb{R}$  is called Morse if  
*all critical points are nondegenerate.*

**All Morse functions are rideable (strict)-saddle functions!**

The Morse functions form an open, dense subset of all smooth functions  $\mathcal{M} \mapsto \mathbb{R}$ .

**A typical/generic function is Morse!**



Marston Morse  
(1892 – 1977)

## More on (P-2): A quantitative definition

**Ridable-saddle (strict-saddle) functions** A function  $f : \mathcal{M} \mapsto \mathbb{R}$  is  $(\alpha, \beta, \gamma, \delta)$ -ridable ( $\alpha, \beta, \gamma, \delta > 0$ ) if any point  $\mathbf{x} \in \mathcal{M}$  obeys **at least one of the following**:

- 1) [**Strong gradient**]  $\|\text{grad } f(\mathbf{x})\| \geq \beta$ ;
- 2) [**Negative curvature**] There exists  $\mathbf{v} \in T_{\mathbf{x}}\mathcal{M}$  with  $\|\mathbf{v}\| = 1$  such that  $\langle \text{Hess } f(\mathbf{x})[\mathbf{v}], \mathbf{v} \rangle \leq -\alpha$ ;
- 3) [**Strong convexity around minimizers**] There exists a local minimizer  $\mathbf{x}_*$  such that  $\|\mathbf{x} - \mathbf{x}_*\| \leq \delta$ , and for all  $\mathbf{y} \in \mathcal{M}$  that is in  $2\delta$  neighborhood of  $\mathbf{x}_*$ ,  $\langle \text{Hess } f(\mathbf{y})[\mathbf{v}], \mathbf{v} \rangle \geq \gamma$  for any  $\mathbf{v} \in T_{\mathbf{y}}\mathcal{M}$  with  $\|\mathbf{v}\| = 1$ .

( $\mathcal{M}$  is a Riemannian manifold,  $T_{\mathbf{x}}\mathcal{M}$  is the tangent space of  $\mathcal{M}$  at point  $\mathbf{x}$ )

$\mathcal{X}$  functions

Examples from practical problems

Sparse (complete) dictionary learning [S., Qu, Wright, '15]

Generalized phase retrieval [S., Qu, Wright, '16]

Other examples in the literature

Numerical optimization methods

Comparison with alternatives

## Benign structure

- (P-1) All local minimizers are also global,
- (P-2A) For all local minimizers,  $\text{Hess } f \succ \mathbf{0}$ , and
- (P-2B) For all other critical points,  $\lambda_{\min}(\text{Hess } f) < 0$ .

... focus on finding a local minimizer

$\mathcal{X}$  functions

Examples from practical problems

Sparse (complete) dictionary learning [S., Qu, Wright, '15]

Generalized phase retrieval [S., Qu, Wright, '16]

Other examples in the literature

Numerical optimization methods

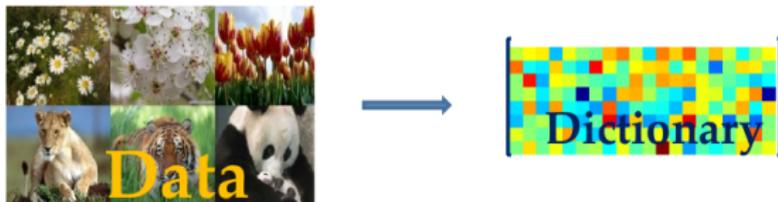
Comparison with alternatives

# Example I: Sparse Dictionary Learning



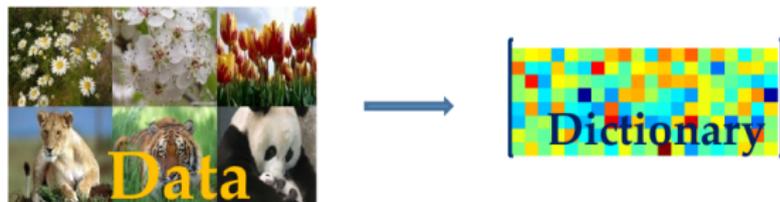
- Algorithmic study initiated in neuroscience [Olshausen and Field, 1996].
- Important algorithmic contributions from many researchers: [Lewicki and Sejnowski, 2000, Engan et al., 1999, Aharon et al., 2006], many others
- Widely used in image processing, visual recognition, compressive signal acquisition, deep architecture for signal classification (see, e.g., [Mairal et al., 2014])

# Dictionary recovery - the complete case



Given  $Y$  generated as  $Y = Q_0 X_0$ , recover  $Q_0$  and  $X_0$ .

# Dictionary recovery - the complete case



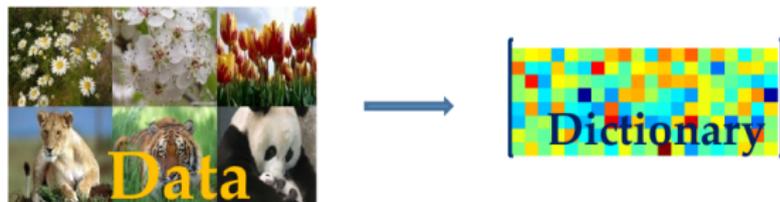
Given  $Y$  generated as  $Y = Q_0 X_0$ , recover  $Q_0$  and  $X_0$ .

## Random Data Model

$Q_0$  complete (square and invertible)

$X_0$  Bernoulli( $\theta$ )-Gaussian:  $X_0 = \Omega \odot G$ ,  $\Omega \sim_{iid} \text{Ber}(\theta)$ ,  $G \sim_{iid} \mathcal{N}(0, 1)$ .

# Dictionary recovery - the complete case



Given  $Y$  generated as  $Y = Q_0 X_0$ , recover  $Q_0$  and  $X_0$ .

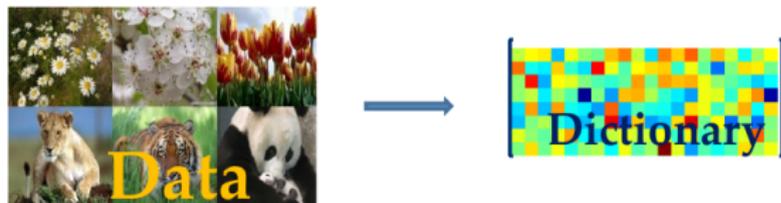
## Random Data Model

$Q_0$  complete (square and invertible)

$X_0$  Bernoulli( $\theta$ )-Gaussian:  $X_0 = \Omega \odot G$ ,  $\Omega \sim_{iid} \text{Ber}(\theta)$ ,  $G \sim_{iid} \mathcal{N}(0, 1)$ .

- $Q_0$  complete  $\implies$   $\boxed{\text{row}(Y) = \text{row}(X_0)}$

# Dictionary recovery - the complete case



Given  $Y$  generated as  $Y = Q_0 X_0$ , recover  $Q_0$  and  $X_0$ .

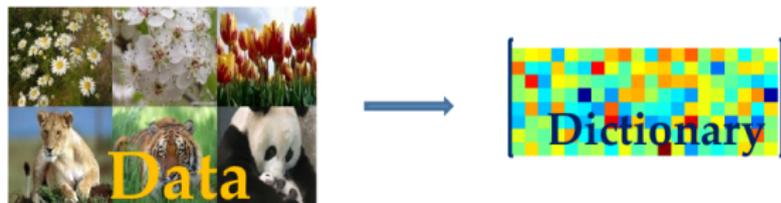
## Random Data Model

$Q_0$  complete (square and invertible)

$X_0$  Bernoulli( $\theta$ )-Gaussian:  $X_0 = \Omega \odot G, \Omega \sim_{iid} \text{Ber}(\theta), G \sim_{iid} \mathcal{N}(0, 1)$ .

- $Q_0$  complete  $\implies \boxed{\text{row}(Y) = \text{row}(X_0)}$
- Rows of  $X_0$  are **sparse** vectors in  $\text{row}(Y)$

# Dictionary recovery - the complete case



Given  $Y$  generated as  $Y = Q_0 X_0$ , recover  $Q_0$  and  $X_0$ .

## Random Data Model

$Q_0$  complete (square and invertible)

$X_0$  Bernoulli( $\theta$ )-Gaussian:  $X_0 = \Omega \odot G, \Omega \sim_{iid} \text{Ber}(\theta), G \sim_{iid} \mathcal{N}(0, 1)$ .

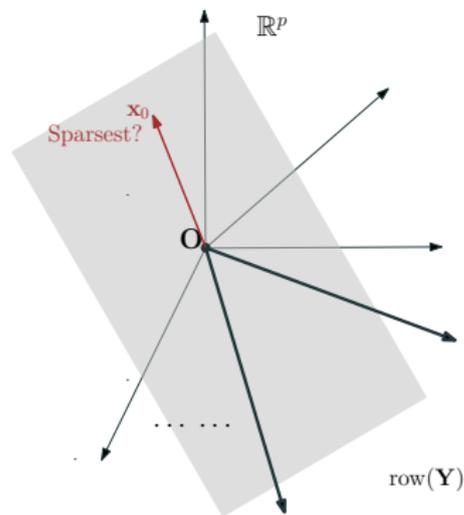
- $Q_0$  complete  $\implies \boxed{\text{row}(Y) = \text{row}(X_0)}$
- Rows of  $X_0$  are **sparse** vectors in  $\text{row}(Y)$
- When  $p \geq \Omega(n \log n)$ , rows of  $X_0$  are **the sparsest** vectors in  $\text{row}(Y)$  [Spielman et al., 2012]

# Dictionary recovery - the complete case

Dictionary recovery: Given  $\mathbf{Y} = \mathbf{Q}_0 \mathbf{X}_0$ , recover  $\mathbf{Q}_0$  and  $\mathbf{X}_0$ .

$\mathbf{Q}_0$  square, invertible:  $\text{row}(\mathbf{Y}) = \text{row}(\mathbf{X}_0)$

**Find the sparsest vectors in  $\text{row}(\mathbf{Y})$ :**



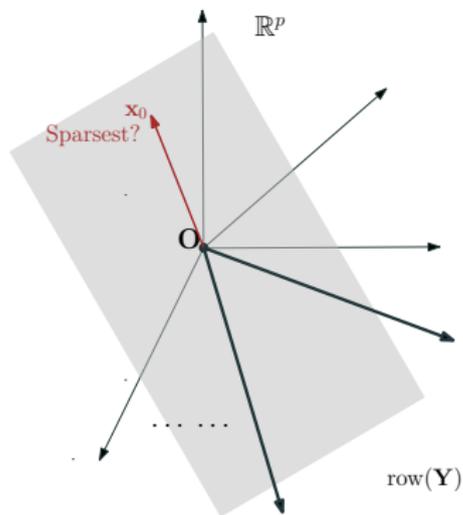
# Dictionary recovery - the complete case

Dictionary recovery: Given  $\mathbf{Y} = \mathbf{Q}_0 \mathbf{X}_0$ , recover  $\mathbf{Q}_0$  and  $\mathbf{X}_0$ .

$\mathbf{Q}_0$  square, invertible:  $\text{row}(\mathbf{Y}) = \text{row}(\mathbf{X}_0)$

**Find the sparsest vectors in  $\text{row}(\mathbf{Y})$ :**

$$\min_{\mathbf{q}} \|\mathbf{q}^* \mathbf{Y}\|_0 \quad \text{s.t.} \quad \mathbf{q} \neq \mathbf{0}.$$



# Dictionary recovery - the complete case

Dictionary recovery: Given  $\mathbf{Y} = \mathbf{Q}_0 \mathbf{X}_0$ , recover  $\mathbf{Q}_0$  and  $\mathbf{X}_0$ .

$\mathbf{Q}_0$  square, invertible:  $\text{row}(\mathbf{Y}) = \text{row}(\mathbf{X}_0)$

**Find the sparsest vectors in  $\text{row}(\mathbf{Y})$ :**

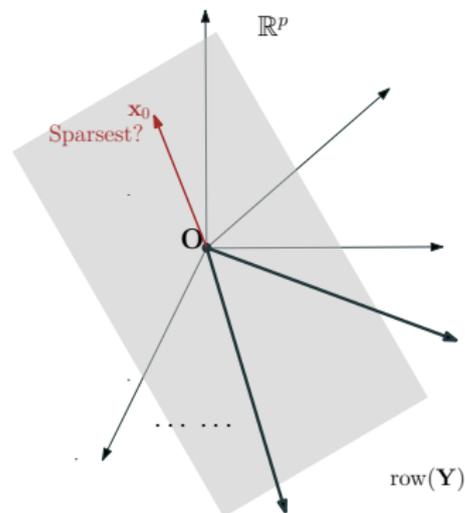
$$\min_{\mathbf{q}} \|\mathbf{q}^* \mathbf{Y}\|_0 \quad \text{s.t.} \quad \mathbf{q} \neq \mathbf{0}.$$

**Nonconvex “relaxation”**

$$\min_{\mathbf{q}} \|\mathbf{q}^* \mathbf{Y}\|_1 \quad \text{s.t.} \quad \|\mathbf{q}\|_2^2 = 1.$$

Many precedents, e.g.,

[Zibulevsky and Pearlmutter, 2001] in blind source separation.



# Towards geometric understanding

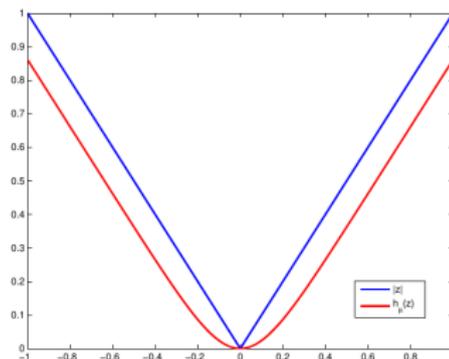
## Model problem

$$\min_{\mathbf{q}} \frac{1}{p} \|\mathbf{q}^* \mathbf{Y}\|_1 = \frac{1}{p} \sum_{i=1}^p |\mathbf{q}^* \mathbf{y}_i| \quad \text{s.t.} \quad \|\mathbf{q}\|_2^2 = 1. \quad \mathbf{Y} \in \mathbb{R}^{n \times p}$$

## Smoothed model problem

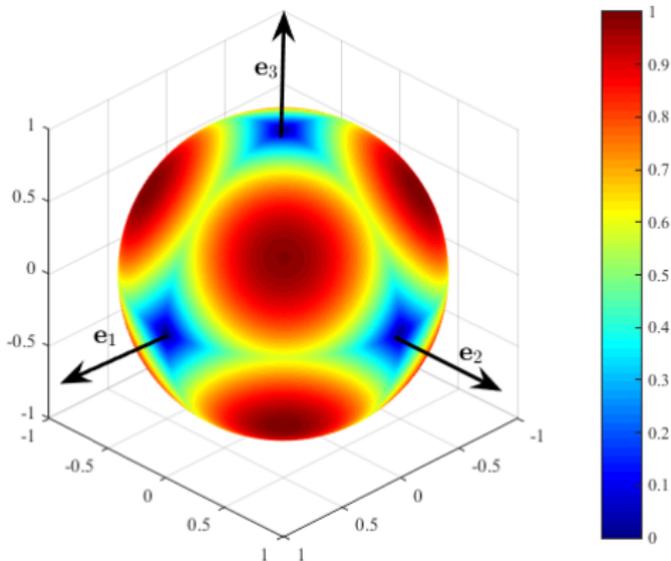
$$\min_{\mathbf{q}} f(\mathbf{q}) \doteq \frac{1}{p} \sum_{i=1}^p h_{\mu}(\mathbf{q}^* \mathbf{y}_i) \quad \text{s.t.} \quad \|\mathbf{q}\|_2^2 = 1. \quad \mathbf{Y} \in \mathbb{R}^{n \times p}$$

$$h_{\mu}(z) = \mu \log \cosh \frac{z}{\mu}$$

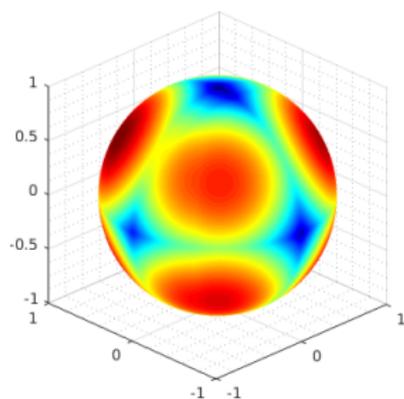


# An $\chi$ function!

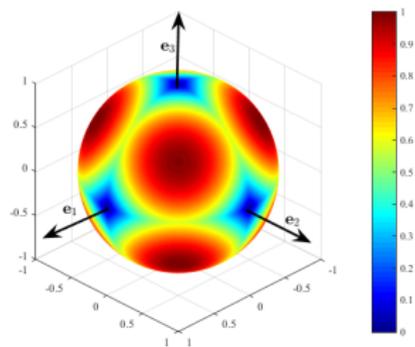
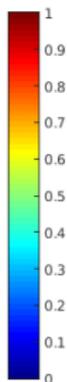
A low-dimensional example ( $n = 3$ ) of the landscape when the target dictionary  $Q_0$  is  $I$  and  $p \rightarrow \infty$



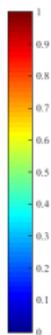
# From finite samples



$$p = 100$$



$$p \rightarrow \infty$$



When  $p \sim n^3$  (suppressing log factors, dependence on  $\mu$ ), the finite sample version is also “nice”.

## The results

$$\min \quad f(\mathbf{q}) \doteq \frac{1}{p} \sum_{i=1}^p h_{\mu}(\mathbf{q}^* \mathbf{y}_i) \quad \text{s.t.} \quad \|\mathbf{q}\|_2^2 = 1. \quad \mathbf{Y} \in \mathbb{R}^{n \times p}$$

### Theorem (Informal, S., Qu, Wright '15)

When  $p$  is reasonably large, and  $\theta \leq 1/3$ , with high probability,

- All local minimizers produce close approximations to rows of  $\mathbf{X}_0$
- $f$  is  $(c\theta, c\theta, c\theta/\mu, \sqrt{2}\mu/7)$ -ridable over  $\mathbb{S}^{n-1}$  for some  $c > 0$

Algorithms later ...

# Comparison with the DL Literature

- **Efficient algorithms** with performance guarantees

[Spielman et al., 2012]  $Q \in \mathbb{R}^{n \times n}$ ,  $\theta = \tilde{O}(1/\sqrt{n})$

[Agarwal et al., 2013b]  $Q \in \mathbb{R}^{m \times n}$  ( $m \leq n$ ),  $\theta = \tilde{O}(1/\sqrt{n})$

[Arora et al., 2013]  $Q \in \mathbb{R}^{m \times n}$  ( $m \leq n$ ),  $\theta = \tilde{O}(1/\sqrt{n})$

[Arora et al., 2015]  $Q \in \mathbb{R}^{m \times n}$  ( $m \leq n$ ),  $\theta = \tilde{O}(1/\sqrt{n})$

- **Quasipolynomial algorithms** with better guarantees

[Arora et al., 2014] different model,  $\theta = O(1/\text{polylog}(n))$

[Barak et al., 2014] sum-of-squares,  $\theta = \tilde{O}(1)$   
polytime for  $\theta = O(n^{-\epsilon})$ .

- Other theoretical work on **local geometry**:

[Gribonval and Schnass, 2010], [Geng and Wright, 2011], [Schnass, 2014], etc

This work: the first polynomial-time algorithm for complete  $Q$  with  $\theta = \Omega(1)$ .

See also refined SOS [Ma et al., 2016].

$\mathcal{X}$  functions

Examples from practical problems

Sparse (complete) dictionary learning [S., Qu, Wright, '15]

Generalized phase retrieval [S., Qu, Wright, '16]

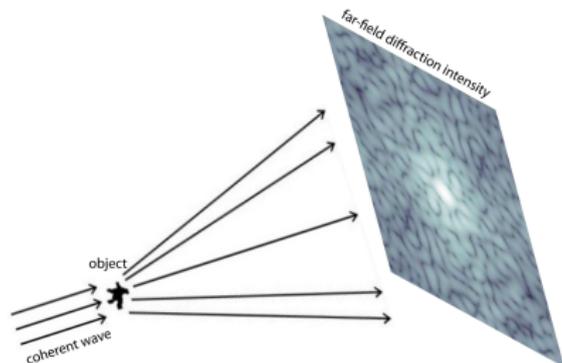
Other examples in the literature

Numerical optimization methods

Comparison with alternatives

## Example II: Generalized phase retrieval

**Phase retrieval:** Given phaseless information of a complex signal, recover the signal



**Applications:** X-ray crystallography, diffraction imaging (left), optics, astronomical imaging, and microscopy

### Coherent Diffraction Imaging<sup>1</sup>

For a complex signal  $x \in \mathbb{C}^n$ , given  $|\mathcal{F}x|$ , recover  $x$ .

<sup>1</sup>Image courtesy of [Shechtman et al., 2015]

# Generalized phase retrieval

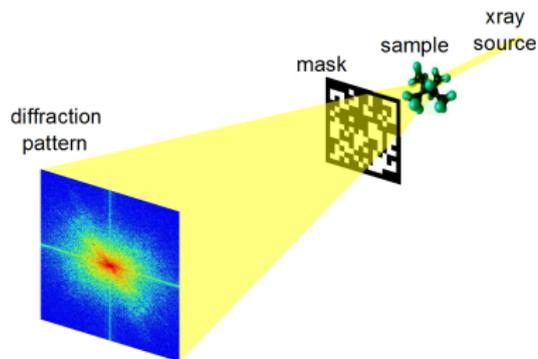
For a complex signal  $x \in \mathbb{C}^n$ , given  $|\mathcal{F}x|$ , recover  $x$ .

## Generalized phase retrieval:

For a complex signal  $x \in \mathbb{C}^n$ , given measurements of the form  $|\mathbf{a}_k^* x|$  for  $k = 1, \dots, m$ , recover  $x$ .

... in practice, generalized measurements by design such as masking, grating, structured illumination, etc

2



## A nonconvex formulation

- Given  $y_k = |\mathbf{a}_k^* \mathbf{x}|$  for  $k = 1, \dots, m$ , recover  $\mathbf{x}$  (**up to a global phase**).
- A natural **nonconvex** formulation (see also [Candès et al., 2015b])

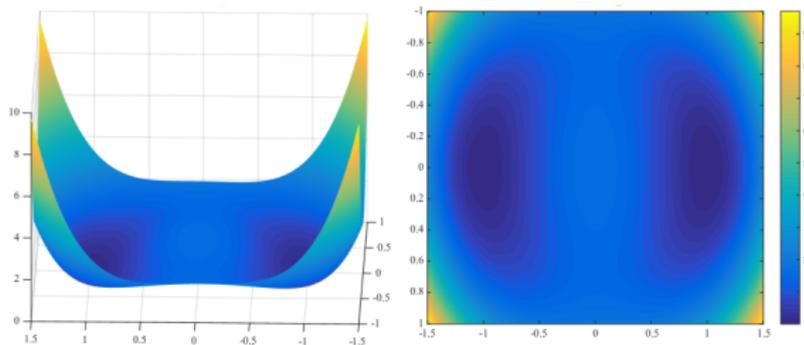
$$\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) \doteq \frac{1}{2m} \sum_{k=1}^m (y_k^2 - |\mathbf{a}_k^* \mathbf{z}|^2)^2.$$

# A nonconvex formulation

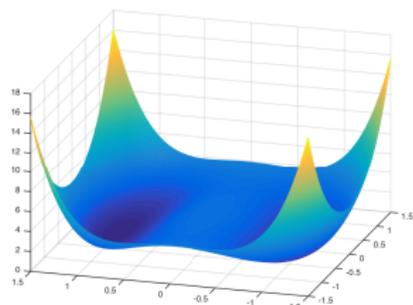
- Given  $y_k = |\mathbf{a}_k^* \mathbf{x}|$  for  $k = 1, \dots, m$ , recover  $\mathbf{x}$  (up to a global phase).
- A natural **nonconvex** formulation (see also [Candès et al., 2015b])

$$\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) \doteq \frac{1}{2m} \sum_{k=1}^m (y_k^2 - |\mathbf{a}_k^* \mathbf{z}|^2)^2.$$

When  $\mathbf{a}_k$ 's are iid standard complex Gaussian vectors and  $m$  large



# The results



$$\min_{z \in \mathbb{C}^n} f(z) \doteq \frac{1}{2m} \sum_{k=1}^m (y_k^2 - |\mathbf{a}_k^* z|^2)^2.$$

## Theorem (Informal, S., Qu, Wright '16)

Let  $\mathbf{a}_k \sim_{\text{iid}} \mathcal{CN}(0, 1)$ . When  $m \geq \Omega(n \log^3(n))$ , w.h.p.,

- All local (and global) minimizers are of the form  $\mathbf{x}e^{i\phi}$ .
- $f$  is  $(c, c/(n \log m), c, c/(n \log m))$ -ridable over  $\mathbb{C}^n$  for some  $c > 0$ .

# Comparison with the literature

- **SDP relaxations and their analysis:**

[Candès et al., 2013a]	SDP relaxation
[Candès et al., 2013b]	Guarantees for $m \sim n \log n$ , adaptive
[Candès and Li, 2014]	Guarantees for $m \sim n$ , non-adaptive
[Candès et al., 2015a]	Coded diffraction patterns
[Waldspurger et al., 2015]	SDP relaxation in phase space

- **Nonconvex methods** (spectral init. + local refinement):

[Netrapalli et al., 2013]	Spectral init. sample splitting
[Candès et al., 2015b]	Spectral init. + gradient descent, $m \sim n \log n$ .
[White et al., 2015]	Spectral init. + gradient descent
[Chen and Candès, 2015]	Spectral init. + truncation, $m \sim n$ .

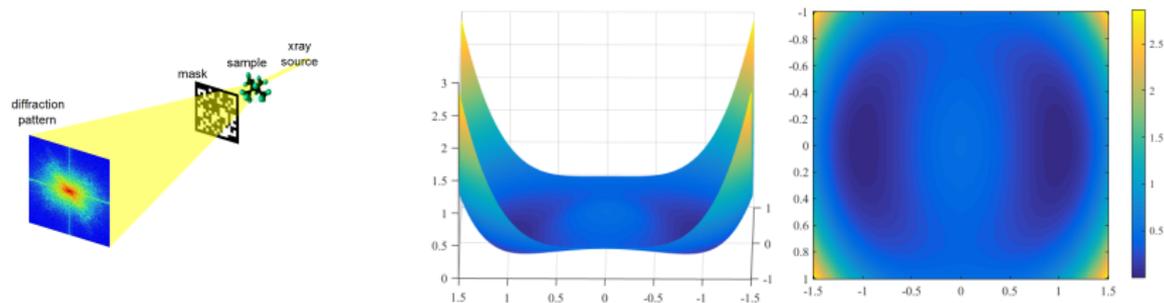
This work: a global characterization of the geometry of the problem.  
Algorithms succeed independent of initialization, when  
 $m \sim n \log^3 n$ .

Later improved to  $m \sim n \log n$ .

# Other measurement models for GPR

## Other measurements

- Coded diffraction model [Candès et al., 2015a]



- Convolutional model:  $y = |a \circledast x|$  [Qu et al., 2017]

$\mathcal{X}$  functions

Examples from practical problems

Sparse (complete) dictionary learning [S., Qu, Wright, '15]

Generalized phase retrieval [S., Qu, Wright, '16]

Other examples in the literature

Numerical optimization methods

Comparison with alternatives

## Example III: Orthogonal tensor decomposition

... generalizes eigen-decomposition of matrices  $M = \sum_{i=1}^r \lambda_i \mathbf{a}_i \otimes \mathbf{a}_i$

**Orthogonally decomposable (OD)**  $d$ -th order tensors

$$\mathcal{T} = \sum_{i=1}^r \lambda_i \mathbf{a}_i^{\otimes d}, \quad \mathbf{a}_i^\top \mathbf{a}_j = \delta_{ij} \quad \forall i, j, (\mathbf{a}_i \in \mathbb{R}^n \quad \forall i)$$

where  $\otimes$  generalizes the usual outer product of vectors.

## Example III: Orthogonal tensor decomposition

... generalizes eigen-decomposition of matrices  $M = \sum_{i=1}^r \lambda_i \mathbf{a}_i \otimes \mathbf{a}_i$

**Orthogonally decomposable (OD)**  $d$ -th order tensors

$$\mathcal{T} = \sum_{i=1}^r \lambda_i \mathbf{a}_i^{\otimes d}, \quad \mathbf{a}_i^\top \mathbf{a}_j = \delta_{ij} \quad \forall i, j, (\mathbf{a}_i \in \mathbb{R}^n \quad \forall i)$$

where  $\otimes$  generalizes the usual outer product of vectors.

**Orthogonal tensor decomposition:** given OD tensor  $\mathcal{T}$ , find the components  $\mathbf{a}_i$ 's (up to sign and permutations).

**Applications:** independent component analysis (ICA), blind source separation, latent variable model learning, etc (see, e.g., [Anandkumar et al., 2014a])

# One component each time

Focus on OD tensors of the form

$$\mathcal{T} = \sum_{i=1}^n \mathbf{1} \cdot \mathbf{a}_i^{\otimes 4}, \quad \mathbf{a}_i^\top \mathbf{a}_j = \delta_{ij} \quad \forall i, j, (\mathbf{a}_i \in \mathbb{R}^n \quad \forall i)$$

# One component each time

Focus on OD tensors of the form

$$\mathcal{T} = \sum_{i=1}^n \mathbf{1} \cdot \mathbf{a}_i^{\otimes 4}, \quad \mathbf{a}_i^\top \mathbf{a}_j = \delta_{ij} \quad \forall i, j, (\mathbf{a}_i \in \mathbb{R}^n \quad \forall i)$$

Consider

$$\min f(\mathbf{u}) \doteq -\mathcal{T}(\mathbf{u}, \mathbf{u}, \mathbf{u}, \mathbf{u}) = -\sum_{i=1}^n (\mathbf{a}_i^\top \mathbf{u})^4 \quad \text{s. t.} \quad \|\mathbf{u}\|_2 = 1$$

[Ge et al., 2015] proved that

- $\pm \mathbf{a}_i$ 's are the only minimizers
- $f$  is  $(7/n, 1/\text{poly}(n), 3, 1/\text{poly}(n))$ -ridable over  $\mathbb{S}^{n-1}$

# All components in one shot

Focus on OD tensors of the form

$$\mathcal{T} = \sum_{i=1}^n \mathbf{1} \cdot \mathbf{a}_i^{\otimes 4}, \quad \mathbf{a}_i^\top \mathbf{a}_j = \delta_{ij} \quad \forall i, j, (\mathbf{a}_i \in \mathbb{R}^n \quad \forall i)$$

Consider the “contrast” formulation

$$\begin{aligned} \min g(\mathbf{u}_1, \dots, \mathbf{u}_n) &\doteq \sum_{i \neq j} \mathcal{T}(\mathbf{u}_i, \mathbf{u}_i, \mathbf{u}_j, \mathbf{u}_j) \\ &= \sum_{i \neq j} \sum_{k=1}^n (\mathbf{a}_k^\top \mathbf{u}_i)^2 (\mathbf{a}_k^\top \mathbf{u}_j)^2, \\ \text{s. t. } &\|\mathbf{u}_i\| = 1 \quad \forall i \in [n] \end{aligned}$$

[Ge et al., 2015] proved that

- All local minimizers of  $g$  are equivalent (i.e., signed permuted) copies of  $[\mathbf{a}_1, \dots, \mathbf{a}_n]$
- $g$  is  $(1/\text{poly}(n), 1/\text{poly}(n), 1, 1/\text{poly}(n))$ -ridable

## Example IV: Phase synchronization

**Synchronization:** recovery from **noisy/incomplete** pairwise relative measurements

- angles/phases – from  $e^{i(\theta_i - \theta_j)} + \Delta_{ij}$ ;
- rotations – from  $\mathbf{R}_i \mathbf{R}_j^{-1} + \Delta_{ij}$ ,  $\mathbf{R}_i, \mathbf{R}_j \in \text{SO}(3)$
- group elements – from  $g_i g_j^{-1} + \Delta_{ij}$  for  $g_i, g_j$  over a compact group  $\mathcal{G}$

**Applications:** signal reconstruction, computer vision (structure from motion, surface reconstruction), cryo-electron microscopy, digital communications, ranking, ... (see, e.g., [Bandeira et al., 2014, Bandeira et al., 2015])

## Example IV: Phase synchronization

**Phase synchronization:** Let  $\mathbf{z} \in \mathbb{C}^n$  and  $|z_1| = \dots = |z_n| = 1$ .  
Given measurements  $C_{ij} = z_i \bar{z}_j + \Delta_{ij}$ , recover  $\mathbf{z}$ .

In matrix form,  $\mathbf{C} = \mathbf{z}\mathbf{z}^* + \mathbf{\Delta}$  and assume  $\mathbf{\Delta}$  Hermitian.

Least-squares formulation:

$$\min_{\mathbf{x} \in \mathbb{C}^n} \|\mathbf{x}\mathbf{x}^* - \mathbf{C}\|_F^2, \quad \text{s. t.} \quad |x_1| = \dots = |x_n| = 1.$$

Equivalent to

$$\min_{\mathbf{x} \in \mathbb{C}^n: |x_1| = \dots = |x_n| = 1} f(\mathbf{u}) \doteq -\mathbf{x}^* \mathbf{C} \mathbf{x}$$

## Quadratic over products of circles

$C = zz^* + \Delta$  and assume  $\Delta$  Hermitian

$$\min_{\mathbf{x} \in \mathbb{C}^N: |x_1| = \dots = |x_n| = 1} f(\mathbf{u}) \doteq -\mathbf{x}^* C \mathbf{x}$$

## Quadratic over products of circles

$C = zz^* + \Delta$  and assume  $\Delta$  Hermitian

$$\min_{\mathbf{x} \in \mathbb{C}^N: |x_1| = \dots = |x_n| = 1} f(\mathbf{u}) \doteq -\mathbf{x}^* C \mathbf{x}$$

[Boumal, 2016] showed when  $\Delta$  is “small”,

*second-order necessary conditions for optimality is also sufficient and the global minimizers recover  $z$ .*

This implies

*all local minimizers are global; all saddles are rideable.*

## Quadratic over products of circles

$C = zz^* + \Delta$  and assume  $\Delta$  Hermitian

$$\min_{\mathbf{x} \in \mathbb{C}^N: |x_1| = \dots = |x_n| = 1} f(\mathbf{u}) \doteq -\mathbf{x}^* C \mathbf{x}$$

[Boumal, 2016] showed when  $\Delta$  is “small” ,

*second-order necessary conditions for optimality is also sufficient and the global minimizers recover  $z$ .*

This implies

*all local minimizers are global; all saddles are rideable.*

Analogous results obtained on synchronization over signs and two-block community detection [Bandeira et al., 2016].

## A partial list of known examples

- (Complete) dictionary learning [Sun et al., 2015]
- Generalized phase retrieval [Sun et al., 2016]
- Orthogonal tensor decomposition [Ge et al., 2015]
- Phase synchronization [Boumal, 2016]
- Low-rank matrix completion/recovery [Ge et al., 2016, Bhojanapalli et al., 2016]
- Shallow networks [Soltanolkotabi et al., 2017]
- Deep networks [Kawaguchi, 2016]
- Blind deconvolution [Zhang et al., 2017]
- **Your examples!**

$\mathcal{X}$  functions

Examples from practical problems

Sparse (complete) dictionary learning [S., Qu, Wright, '15]

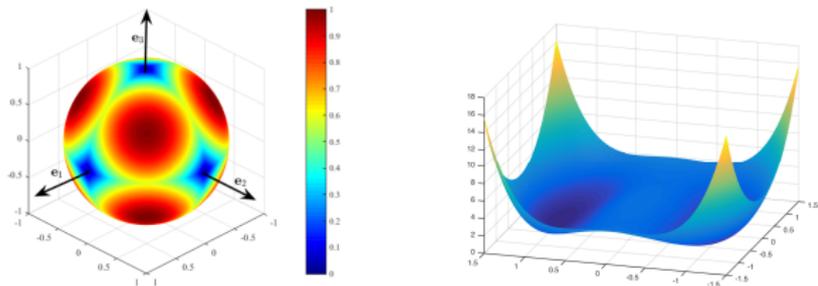
Generalized phase retrieval [S., Qu, Wright, '16]

Other examples in the literature

**Numerical optimization methods**

Comparison with alternatives

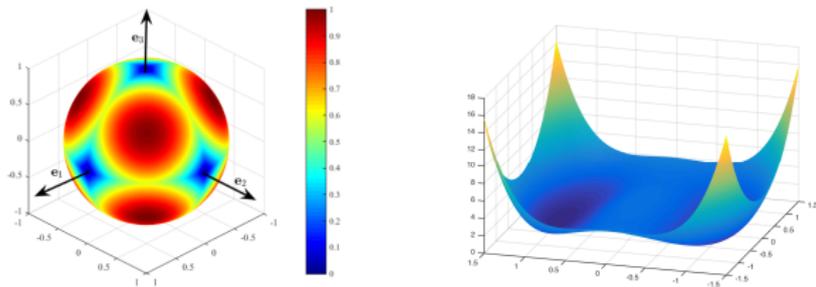
# Benign structure



- (P-1) All local minimizers are also global,
- (P-2A) For all local minimizers,  $\text{Hess } f \succ \mathbf{0}$ , and
- (P-2B) For all other critical points,  $\lambda_{\min}(\text{Hess } f) < 0$ .

... focus on **escaping saddle points** and finding a **local minimizer**.

# Algorithmic possibilities



- **Second-order trust-region method** ([Conn et al., 2000], [Nesterov and Polyak, 2006])
- Curvilinear search [Goldfarb, 1980, Goldfarb et al., 2017]
- Noisy/stochastic gradient descent [Ge et al., 2015, Jin et al., 2017]
- Open: More sophisticated working methods?

## Second-order methods can escape ridable saddles

Taylor expansion at a saddle point  $\mathbf{x}$ :

$$\hat{f}(\boldsymbol{\delta}; \mathbf{x}) = f(\mathbf{x}) + \frac{1}{2} \boldsymbol{\delta}^* \nabla^2 f(\mathbf{x}) \boldsymbol{\delta}.$$

Choosing  $\boldsymbol{\delta} = \mathbf{v}_{\text{neg}}$ , then

$$\hat{f}(\boldsymbol{\delta}; \mathbf{x}) - f(\mathbf{x}) \leq -\frac{1}{2} |\lambda_{\text{neg}}| \|\mathbf{v}_{\text{neg}}\|^2.$$

Guaranteed decrease in  $f$  when **movement is small** such that the **approximation is reasonably good**.

# Trust-region method - Euclidean Space

Generate iterates  $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots$  by

- Forming a second order approximation of the objective  $f(\mathbf{x})$  about  $\mathbf{x}_k$ :

$$\hat{f}(\boldsymbol{\delta}; \mathbf{x}_k) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\delta} \rangle + \frac{1}{2} \boldsymbol{\delta}^* \mathbf{B}_k \boldsymbol{\delta}.$$

and minimizing the approximation within a small radius - the trust region

$$\boldsymbol{\delta}_* \in \arg \min_{\|\boldsymbol{\delta}\| \leq \Delta} \hat{f}(\boldsymbol{\delta}; \mathbf{x}_k) \quad (\text{Trust-region subproblem})$$

- Next iterate is  $\mathbf{x}_{k+1} = \mathbf{x}_k + \boldsymbol{\delta}_*$ .

Can choose  $\mathbf{B}_k = \nabla^2 f(\mathbf{x}^{(k)})$  or an approximation.

## The trust-region subproblem

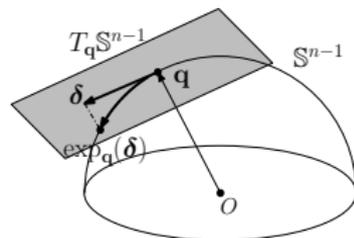
$$\delta_\star \in \arg \min_{\|\delta\| \leq \Delta} \widehat{f}(\delta; \mathbf{x}_k) \quad (\text{Trust-region subproblem})$$

- QCQP, but can be solved in polynomial time by:
  - Root finding [Moré and Sorensen, 1983]
  - SDP relaxation [Rendl and Wolkowicz, 1997].
- In practice, only need an approximate solution (with controllable quality) to ensure convergence.

# Trust-region method - Riemannian Manifold

Local quadratic approximation:

$$\begin{aligned} f(\exp_{\mathbf{q}}(\boldsymbol{\delta})) \\ = \underbrace{f(\mathbf{q}) + \boldsymbol{\delta}^* \text{grad } f(\mathbf{q}) + \frac{1}{2} \boldsymbol{\delta}^* \text{Hess } f(\mathbf{q}) \boldsymbol{\delta}}_{\hat{f}(\boldsymbol{\delta}; \mathbf{q})} + O(\|\boldsymbol{\delta}\|^3). \end{aligned}$$



Basic **Riemannian trust-region method**:

$$\begin{aligned} \boldsymbol{\delta}_* \in \arg \min_{\boldsymbol{\delta} \in T_{\mathbf{q}_k} \mathbb{S}^{n-1}, \|\boldsymbol{\delta}\| \leq \Delta} \hat{f}(\boldsymbol{\delta}; \mathbf{q}_k) \\ \mathbf{q}_{k+1} = \exp_{\mathbf{q}_k}(\boldsymbol{\delta}_*). \end{aligned}$$

More details on Riemannian TRM in [Absil et al., 2007] and [Absil et al., 2009].

# Proof of convergence

- Strong gradient or negative curvature  
     $\implies$  at least a fixed reduction in  $f(\mathbf{x})$  at each iteration
- Strong convexity near a local minimizer  
     $\implies$  quadratic convergence  $\|\mathbf{x}_{k+1} - \mathbf{x}_\star\| \leq c \|\mathbf{x}_k - \mathbf{x}_\star\|^2$ .

## Theorem (Meta..., very informal)

*For (quantitative)  $\mathcal{X}$  functions, starting from an **arbitrary initialization**, the iteration sequence with **sufficiently small trust-region size** converges to a global minimizer in **polynomial number of steps**.*

Worked out examples in [Sun et al., 2015, Sun et al., 2016];

Guarantee of 2-nd order method [Boumal et al., 2016];

Guarantee of 1-st order method [Ge et al., 2015, Jin et al., 2017].

$\mathcal{X}$  functions

Examples from practical problems

Sparse (complete) dictionary learning [S., Qu, Wright, '15]

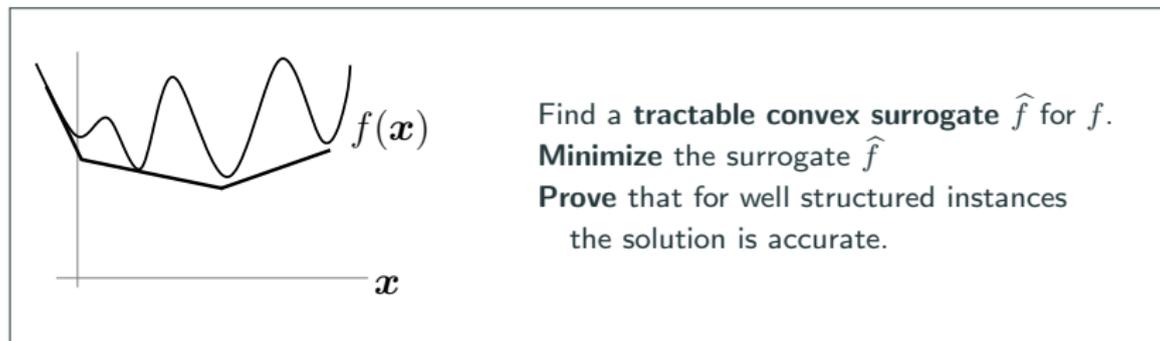
Generalized phase retrieval [S., Qu, Wright, '16]

Other examples in the literature

Numerical optimization methods

Comparison with alternatives

# Convexification – a recipe



**Separates formulations/analysis from algorithms**

**Beautiful mathematical results, substantial applied impact:**

- Examples: sparse recovery, low-rank matrix recovery/completion
- General frameworks:
  - Atomic norms [Chandrasekaran et al., 2012]
  - Submodular sparsity inducers [Bach, 2010]
  - Restricted strong convexity [Negahban et al., 2009]
  - Conic statistical dimensions [Amelunxen et al., 2014], etc.

# But... sometimes the recipe doesn't work

## The natural convex surrogates may be intractable ...

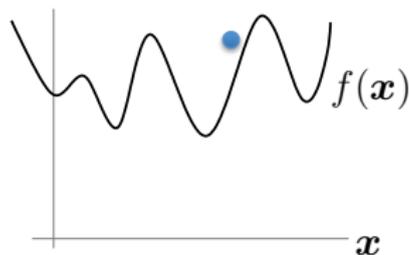
Tensor recovery	[Hillar and Lim, 2013]
Nonnegative low-rank approximation	[Vavasis, 2009]

## ... or may not work as well as we might hope.

Simultaneous structure estimation	[Oymak et al., 2012]
Tensor recovery	[Mu et al., 2014]
Sparse PCA	[Berthet and Rigollet, 2013]
Dictionary learning	[Spielman et al., 2012]

**Substantial and provable gaps between the performance of known convex relaxations and the information theoretic optimum.**

## Prior work: proving nonconvex recovery



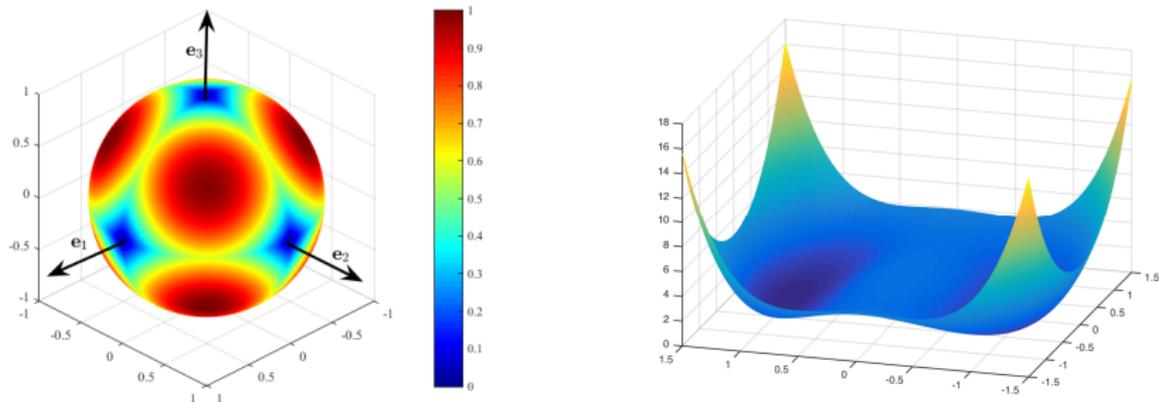
Use problem structure to find a  
**clever initial guess.**

Analyze iteration-by-iteration  
**in the vicinity of the optimum.**

- **Matrix completion/recovery:** [Keshavan et al., 2010], [Jain et al., 2013], [Hardt, 2014], [Hardt and Wootters, 2014], [Netrapalli et al., 2014], [Jain and Netrapalli, 2014], [Sun and Luo, 2014], [Zheng and Lafferty, 2015], [Tu et al., 2015], [Chen and Wainwright, 2015], [Sa et al., 2015], [Wei et al., 2015]. Also [Jain et al., 2010]
- **Dictionary learning:** [Agarwal et al., 2013a], [Arora et al., 2013], [Agarwal et al., 2013b], [Arora et al., 2015]
- **Tensor recovery:** [Jain and Oh, 2014], [Anandkumar et al., 2014c], [Anandkumar et al., 2014b], [Anandkumar et al., 2015]
- **Phase retrieval:** [Netrapalli et al., 2013], [Candès et al., 2015b], [Chen and Candès, 2015], [White et al., 2015]
- **More on the webpage:** <http://sunju.org/research/nonconvex/>

See also [Loh and Wainwright, 2011]

# This work



- We characterize the **geometry**, which is critical to algorithm design whether initialization is used or not
- The geometry effectively allows **arbitrary initialization**

# Thanks to ...



**John Wright**

Columbia



**Qing Qu**

Columbia

**Nonconvex Optimization Meets Symmetry: Examples, Algorithms, and Open problems,**  
S. and Wright, '17 (forthcoming)

**A Geometric Analysis of Phase Retrieval,** S., Qu, Wright, '16

**Complete Dictionary Recovery over the Sphere,** S., Qu, Wright, '15

**When are Nonconvex Optimization Problems Not Scary,** S., Qu, Wright, NIPS  
Workshop, '15

**Finding a Sparse Vector in a Subspace: Linear Sparsity Using Alternating Directions,**  
Qu, S., Wright, '15

**Webpage on provable nonconvex heuristics:**

<http://sunju.org/research/nonconvex/>

# References I

- [Absil et al., 2007] Absil, P.-A., Baker, C. G., and Gallivan, K. A. (2007). **Trust-region methods on Riemannian manifolds.** *Foundations of Computational Mathematics*, 7(3):303–330.
- [Absil et al., 2009] Absil, P.-A., Mahoney, R., and Sepulchre, R. (2009). **Optimization Algorithms on Matrix Manifolds.** Princeton University Press.
- [Agarwal et al., 2013a] Agarwal, A., Anandkumar, A., Jain, P., Netrapalli, P., and Tandon, R. (2013a). **Learning sparsely used overcomplete dictionaries via alternating minimization.** *arXiv preprint arXiv:1310.7991*.
- [Agarwal et al., 2013b] Agarwal, A., Anandkumar, A., and Netrapalli, P. (2013b). **Exact recovery of sparsely used overcomplete dictionaries.** *arXiv preprint arXiv:1309.1952*.
- [Aharon et al., 2006] Aharon, M., Elad, M., and Bruckstein, A. (2006). **K-svd: An algorithm for designing overcomplete dictionaries for sparse representation.** *Trans. Sig. Proc.*, 54(11):4311–4322.
- [Amelunxen et al., 2014] Amelunxen, D., Lotz, M., McCoy, M. B., and Tropp, J. A. (2014). **Living on the edge: Phase transitions in convex programs with random data.** *Information and Inference*, page iau005.
- [Anandkumar et al., 2014a] Anandkumar, A., Ge, R., Hsu, D., Kakade, S. M., and Telgarsky, M. (2014a). **Tensor decompositions for learning latent variable models.** *The Journal of Machine Learning Research*, 15(1):2773–2832.
- [Anandkumar et al., 2014b] Anandkumar, A., Ge, R., and Janzamin, M. (2014b). **Analyzing tensor power method dynamics: Applications to learning overcomplete latent variable models.** *arXiv preprint arXiv:1411.1488*.
- [Anandkumar et al., 2014c] Anandkumar, A., Ge, R., and Janzamin, M. (2014c). **Guaranteed non-orthogonal tensor decomposition via alternating rank-1 updates.** *arXiv preprint arXiv:1402.5180*.
- [Anandkumar et al., 2015] Anandkumar, A., Jain, P., Shi, Y., and Niranjan, U. N. (2015). **Tensor vs matrix methods: Robust tensor decomposition under block sparse perturbations.** *arXiv preprint arXiv:1510.04747*.
- [Arora et al., 2014] Arora, S., Bhaskara, A., Ge, R., and Ma, T. (2014). **More algorithms for provable dictionary learning.** *arXiv preprint arXiv:1401.0579*.

# References II

- [Arora et al., 2015] Arora, S., Ge, R., Ma, T., and Moitra, A. (2015). **Simple, efficient, and neural algorithms for sparse coding.** *arXiv preprint arXiv:1503.00778*.
- [Arora et al., 2013] Arora, S., Ge, R., and Moitra, A. (2013). **New algorithms for learning incoherent and overcomplete dictionaries.** *arXiv preprint arXiv:1308.6273*.
- [Bach, 2010] Bach, F. R. (2010). **Structured sparsity-inducing norms through submodular functions.** In *Advances in Neural Information Processing Systems*, pages 118–126.
- [Bandeira et al., 2014] Bandeira, A. S., Boumal, N., and Singer, A. (2014). **Tightness of the maximum likelihood semidefinite relaxation for angular synchronization.** *arXiv preprint arXiv:1411.3272*.
- [Bandeira et al., 2016] Bandeira, A. S., Boumal, N., and Voroninski, V. (2016). **On the low-rank approach for semidefinite programs arising in synchronization and community detection.** *arXiv preprint arXiv:1602.04426*.
- [Bandeira et al., 2015] Bandeira, A. S., Chen, Y., and Singer, A. (2015). **Non-unique games over compact groups and orientation estimation in cryo-em.** *arXiv preprint arXiv:1505.03840*.
- [Barak et al., 2014] Barak, B., Kelner, J. A., and Steurer, D. (2014). **Dictionary learning and tensor decomposition via the sum-of-squares method.** *arXiv preprint arXiv:1407.1543*.
- [Berthet and Rigollet, 2013] Berthet, Q. and Rigollet, P. (2013). **Complexity theoretic lower bounds for sparse principal component detection.** In *Conference on Learning Theory*.
- [Bhojanapalli et al., 2016] Bhojanapalli, S., Neyshabur, B., and Srebro, N. (2016). **Global optimality of local search for low rank matrix recovery.** *arXiv preprint arXiv:1605.07221*.
- [Boumal, 2016] Boumal, N. (2016). **Nonconvex phase synchronization.** *arXiv preprint arXiv:1601.06114*.
- [Boumal et al., 2016] Boumal, N., Absil, P.-A., and Cartis, C. (2016). **Global rates of convergence for nonconvex optimization on manifolds.** *arXiv preprint arXiv:1605.08101*.
- [Candès et al., 2013a] Candès, E. J., Eldar, Y. C., Strohmer, T., and Voroninski, V. (2013a). **Phase retrieval via matrix completion.** *SIAM Journal on Imaging Sciences*, 6(1).

# References III

- [Candès and Li, 2014] Candès, E. J. and Li, X. (2014). **Solving quadratic equations via phaselift when there are about as many equations as unknowns.** *Foundations of Computational Mathematics*, 14(5):1017–1026.
- [Candès et al., 2015a] Candès, E. J., Li, X., and Soltanolkotabi, M. (2015a). **Phase retrieval from coded diffraction patterns.** *Applied and Computational Harmonic Analysis*, 39(2):277–299.
- [Candès et al., 2015b] Candès, E. J., Li, X., and Soltanolkotabi, M. (2015b). **Phase retrieval via wirtinger flow: Theory and algorithms.** *Information Theory, IEEE Transactions on*, 61(4):1985–2007.
- [Candès et al., 2013b] Candès, E. J., Strohmer, T., and Vershynina, V. (2013b). **Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming.** *Communications on Pure and Applied Mathematics*, 66(8):1241–1274.
- [Chandrasekaran et al., 2012] Chandrasekaran, V., Recht, B., Parrilo, P. A., and Willsky, A. S. (2012). **The convex geometry of linear inverse problems.** *Foundations of Computational mathematics*, 12(6):805–849.
- [Chen and Candès, 2015] Chen, Y. and Candès, E. J. (2015). **Solving random quadratic systems of equations is nearly as easy as solving linear systems.** *arXiv preprint arXiv:1505.05114*.
- [Chen and Wainwright, 2015] Chen, Y. and Wainwright, M. J. (2015). **Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees.** *arXiv preprint arXiv:1509.03025*.
- [Conn et al., 2000] Conn, A. R., Gould, N. I. M., and Toint, P. L. (2000). **Trust-region Methods.** Society for Industrial and Applied Mathematics, Philadelphia, PA, USA.
- [Engan et al., 1999] Engan, K., Aase, S. O., and Hakon Husoy, J. (1999). **Method of optimal directions for frame design.** In *Acoustics, Speech, and Signal Processing, 1999. Proceedings., 1999 IEEE International Conference on*, volume 5, pages 2443–2446. IEEE.
- [Ge et al., 2015] Ge, R., Huang, F., Jin, C., and Yuan, Y. (2015). **Escaping from saddle points—online stochastic gradient for tensor decomposition.** In *Proceedings of The 28th Conference on Learning Theory*, pages 797–842.

# References IV

- [Ge et al., 2016] Ge, R., Lee, J. D., and Ma, T. (2016). **Matrix completion has no spurious local minimum.** *arXiv preprint arXiv:1605.07272*.
- [Geng and Wright, 2011] Geng, Q. and Wright, J. (2011). **On the local correctness of  $\ell^1$ -minimization for dictionary learning.** Submitted to *IEEE Transactions on Information Theory*. Preprint: <http://www.columbia.edu/~jw2966>.
- [Goldfarb, 1980] Goldfarb, D. (1980). **Curvilinear path steplength algorithms for minimization which use directions of negative curvature.** *Mathematical programming*, 18(1):31–40.
- [Goldfarb et al., 2017] Goldfarb, D., Mu, C., Wright, J., and Zhou, C. (2017). **Using negative curvature in solving nonlinear programs.** *arXiv preprint arXiv:1706.00896*.
- [Gribonval and Schnass, 2010] Gribonval, R. and Schnass, K. (2010). **Dictionary identification - sparse matrix-factorization via  $\ell^1$ -minimization.** *IEEE Transactions on Information Theory*, 56(7):3523–3539.
- [Hardt, 2014] Hardt, M. (2014). **Understanding alternating minimization for matrix completion.** In *Foundations of Computer Science (FOCS), 2014 IEEE 55th Annual Symposium on*, pages 651–660. IEEE.
- [Hardt and Wootters, 2014] Hardt, M. and Wootters, M. (2014). **Fast matrix completion without the condition number.** In *Proceedings of The 27th Conference on Learning Theory*, pages 638–678.
- [Hillar and Lim, 2013] Hillar, C. J. and Lim, L.-H. (2013). **Most tensor problems are NP-hard.** *Journal of the ACM (JACM)*, 60(6):45.
- [Jain et al., 2010] Jain, P., Meka, R., and Dhillon, I. S. (2010). **Guaranteed rank minimization via singular value projection.** In *Advances in Neural Information Processing Systems*, pages 937–945.
- [Jain and Netrapalli, 2014] Jain, P. and Netrapalli, P. (2014). **Fast exact matrix completion with finite samples.** *arXiv preprint arXiv:1411.1087*.

# References V

- [Jain et al., 2013] Jain, P., Netrapalli, P., and Sanghavi, S. (2013). **Low-rank matrix completion using alternating minimization.** In *Proceedings of the forty-fifth annual ACM symposium on Theory of Computing*, pages 665–674. ACM.
- [Jain and Oh, 2014] Jain, P. and Oh, S. (2014). **Provable tensor factorization with missing data.** In *Advances in Neural Information Processing Systems*, pages 1431–1439.
- [Jin et al., 2017] Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I. (2017). **How to escape saddle points efficiently.** *arXiv preprint arXiv:1703.00887*.
- [Kawaguchi, 2016] Kawaguchi, K. (2016). **Deep learning without poor local minima.** *arXiv preprint arXiv:1605.07110*.
- [Keshavan et al., 2010] Keshavan, R. H., Montanari, A., and Oh, S. (2010). **Matrix completion from a few entries.** *Information Theory, IEEE Transactions on*, 56(6):2980–2998.
- [Lewicki and Sejnowski, 2000] Lewicki, M. S. and Sejnowski, T. J. (2000). **Learning overcomplete representations.** *Neural computation*, 12(2):337–365.
- [Loh and Wainwright, 2011] Loh, P.-L. and Wainwright, M. J. (2011). **High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity.** In *Advances in Neural Information Processing Systems*, pages 2726–2734.
- [Ma et al., 2016] Ma, T., Shi, J., and Steurer, D. (2016). **Polynomial-time tensor decompositions with sum-of-squares.** In *Foundations of Computer Science (FOCS), 2016 IEEE 57th Annual Symposium on*, pages 438–446. IEEE.
- [Mairal et al., 2014] Mairal, J., Bach, F., and Ponce, J. (2014). **Sparse modeling for image and vision processing.** *Foundations and Trends in Computer Graphics and Vision*, 8(2-3):85–283.
- [Moré and Sorensen, 1983] Moré, J. J. and Sorensen, D. C. (1983). **Computing a trust region step.** *SIAM Journal on Scientific and Statistical Computing*, 4(3):553–572.

# References VI

- [Mu et al., 2014] Mu, C., Huang, B., Wright, J., and Goldfarb, D. (2014). **Square deal: Lower bounds and improved convex relaxations for tensor recovery.** *Journal of Machine Learning Research*, 1:1–48.
- [Murty and Kabadi, 1987] Murty, K. G. and Kabadi, S. N. (1987). **Some NP-complete problems in quadratic and nonlinear programming.** *Mathematical programming*, 39(2):117–129.
- [Negahban et al., 2009] Negahban, S., Yu, B., Wainwright, M. J., and Ravikumar, P. K. (2009). **A unified framework for high-dimensional analysis of  $m$ -estimators with decomposable regularizers.** In *Advances in Neural Information Processing Systems*, pages 1348–1356.
- [Nesterov and Polyak, 2006] Nesterov, Y. and Polyak, B. T. (2006). **Cubic regularization of newton method and its global performance.** *Mathematical Programming*, 108(1):177–205.
- [Netrapalli et al., 2013] Netrapalli, P., Jain, P., and Sanghavi, S. (2013). **Phase retrieval using alternating minimization.** In *Advances in Neural Information Processing Systems*, pages 2796–2804.
- [Netrapalli et al., 2014] Netrapalli, P., Niranjan, U. N., Sanghavi, S., Anandkumar, A., and Jain, P. (2014). **Non-convex robust PCA.** In *Advances in Neural Information Processing Systems*, pages 1107–1115.
- [Olshausen and Field, 1996] Olshausen, B. A. and Field, D. J. (1996). **Emergence of simple-cell receptive field properties by learning a sparse code for natural images.** *Nature*, 381(6583):607–609.
- [Oymak et al., 2012] Oymak, S., Jalali, A., Fazel, M., Eldar, Y. C., and Hassibi, B. (2012). **Simultaneously structured models with application to sparse and low-rank matrices.** *arXiv preprint arXiv:1212.3753*.
- [Qu et al., 2017] Qu, Q., Zhang, Y., Eldar, Y., and Wright, J. (2017). **Convolutional phase retrieval.** In *Advances in Neural Information Processing Systems*.
- [Rendl and Wolkowicz, 1997] Rendl, F. and Wolkowicz, H. (1997). **A semidefinite framework for trust region subproblems with applications to large scale minimization.** *Mathematical Programming*, 77(1):273–299.

# References VII

- [Sa et al., 2015] Sa, C. D., Re, C., and Olukotun, K. (2015). **Global convergence of stochastic gradient descent for some non-convex matrix problems.** In *The 32nd International Conference on Machine Learning*, volume 37, pages 2332–2341.
- [Schnass, 2014] Schnass, K. (2014). **Local identification of overcomplete dictionaries.** *arXiv preprint arXiv:1401.6354*.
- [Shechtman et al., 2015] Shechtman, Y., Eldar, Y. C., Cohen, O., Chapman, H. N., Miao, J., and Segev, M. (2015). **Phase retrieval with application to optical imaging: A contemporary overview.** *Signal Processing Magazine, IEEE*, 32(3):87–109.
- [Soltanolkotabi et al., 2017] Soltanolkotabi, M., Javanmard, A., and Lee, J. D. (2017). **Theoretical insights into the optimization landscape of over-parameterized shallow neural networks.** *arXiv preprint arXiv:1707.04926*.
- [Spielman et al., 2012] Spielman, D. A., Wang, H., and Wright, J. (2012). **Exact recovery of sparsely-used dictionaries.** In *Proceedings of the 25th Annual Conference on Learning Theory*.
- [Sun et al., 2015] Sun, J., Qu, Q., and Wright, J. (2015). **Complete dictionary recovery over the sphere.** *arXiv preprint arXiv:1504.06785*.
- [Sun et al., 2016] Sun, J., Qu, Q., and Wright, J. (2016). **A geometric analysis of phase retrieval.** *arXiv preprint arXiv:1602.06664*.
- [Sun and Luo, 2014] Sun, R. and Luo, Z.-Q. (2014). **Guaranteed matrix completion via non-convex factorization.** *arXiv preprint arXiv:1411.8003*.
- [Tu et al., 2015] Tu, S., Boczar, R., Soltanolkotabi, M., and Recht, B. (2015). **Low-rank solutions of linear matrix equations via procrustes flow.** *arXiv preprint arXiv:1507.03566*.
- [Vavasis, 2009] Vavasis, S. A. (2009). **On the complexity of nonnegative matrix factorization.** *SIAM Journal on Optimization*, 20(3):1364–1377.
- [Waldspurger et al., 2015] Waldspurger, I., d'Aspremont, A., and Mallat, S. (2015). **Phase recovery, maxcut and complex semidefinite programming.** *Mathematical Programming*, 149(1-2):47–81.

# References VIII

- [Wei et al., 2015] Wei, K., Cai, J.-F., Chan, T. F., and Leung, S. (2015). **Guarantees of Riemannian optimization for low rank matrix recovery.** *arXiv preprint arXiv:1511.01562*.
- [White et al., 2015] White, C. D., Ward, R., and Sanghavi, S. (2015). **The local convexity of solving quadratic equations.** *arXiv preprint arXiv:1506.07868*.
- [Zhang et al., 2017] Zhang, Y., Lau, Y., Kuo, H.-w., Cheung, S., Pasupathy, A., and Wright, J. (2017). **On the global geometry of sphere-constrained sparse blind deconvolution.** In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [Zheng and Lafferty, 2015] Zheng, Q. and Lafferty, J. (2015). **A convergent gradient descent algorithm for rank minimization and semidefinite programming from random linear measurements.** *arXiv preprint arXiv:1506.06081*.
- [Zibulevsky and Pearlmutter, 2001] Zibulevsky, M. and Pearlmutter, B. (2001). **Blind source separation by sparse decomposition in a signal dictionary.** *Neural computation*, 13(4):863–882.