

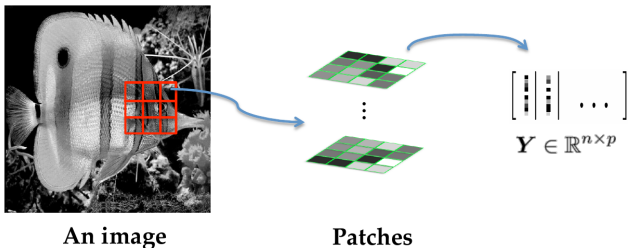
# When Nonconvex Optimization Meets Big Data

Ju Sun

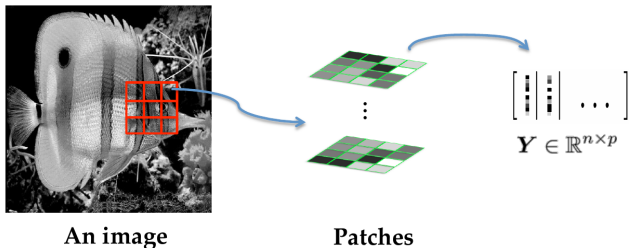
[Advisor: Prof. John Wright]

November 21, 2014

# A curious experiment

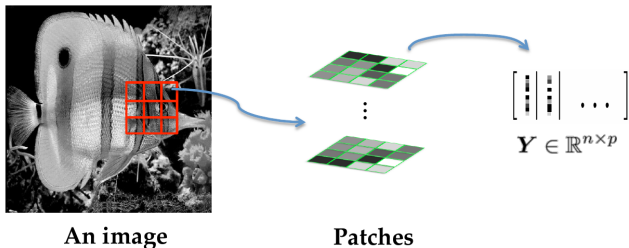


# A curious experiment



Try to learn a **concise approximation**:  $Y \approx QX$ , with  $Q \in O(n)$  and  $X$  sparse.

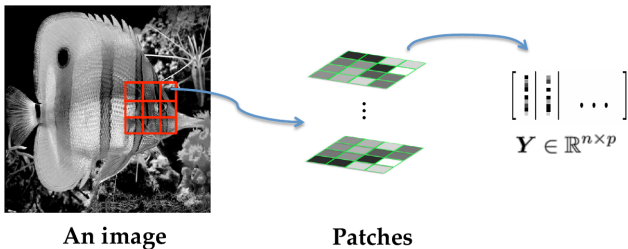
# A curious experiment



Try to learn a **concise approximation**:  $Y \approx QX$ , with  $Q \in O(n)$  and  $X$  sparse.

... by solving  $\min \frac{1}{2} \|Y - QX\|_F^2 + \lambda \|X\|_1$ , s.t.  $Q \in O(n)$ .

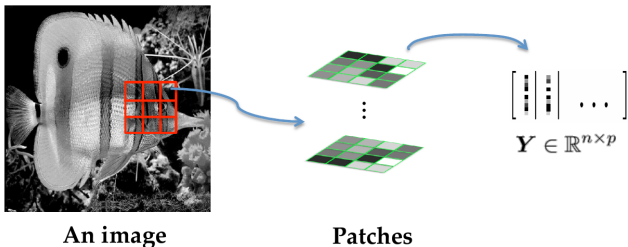
# A curious experiment



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

- Objective is **nonconvex**:  $(Q, X) \mapsto QX$  is bilinear
- Orthogonal group  $O(n)$  is a **nonconvex** set
- **Combinatorially many isolated global minima**:  $(Q, X)$  or  $(Q\Pi, \Pi^* X)$  ( $2^n n!$  many signed permutations  $\Pi$ )

# A curious experiment



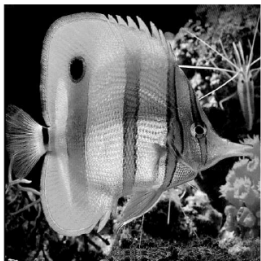
$$\min f(\mathbf{Q}, \mathbf{X}) \doteq \frac{1}{2} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1, \quad \text{s.t. } \mathbf{Q} \in O(n)$$

Apply the naive **alternating directions**: starting from a random  $\mathbf{Q}_0 \in O(n)$

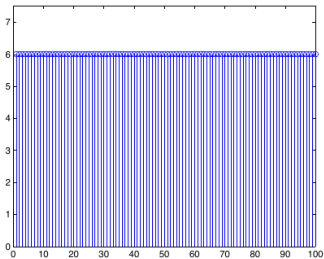
$$\mathbf{X}_k = \arg \min_{\mathbf{X}} f(\mathbf{Q}_{k-1}, \mathbf{X})$$

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

# A curious experiment



An image



Final  $f(\mathbf{Q}_\infty, \mathbf{X}_\infty)$ , varying  $\mathbf{Q}_0$ .

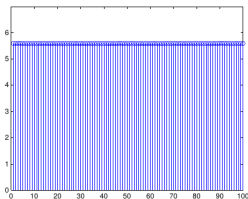
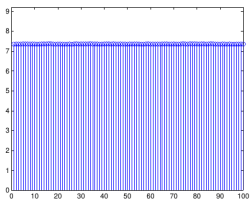
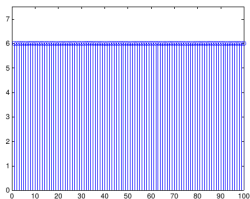
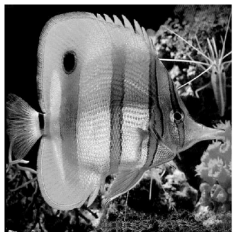
$$\min_{\mathbf{Q}, \mathbf{X}} f(\mathbf{Q}, \mathbf{X}) \doteq \frac{1}{2} \|\mathbf{Y} - \mathbf{Q}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1, \quad \text{s.t. } \mathbf{Q} \in O(n)$$

Apply the naive **alternating directions**: starting from a random  $\mathbf{Q}_0 \in O(n)$

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

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

# What is going on here?

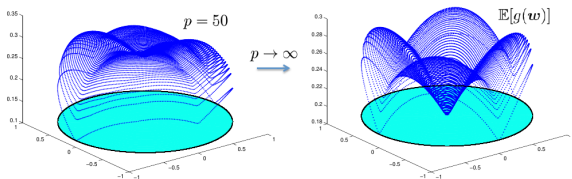


... You can find me and see thousands more!



## Theme of this thesis

- Certain nonconvex optimization problems become tractable when the input data are large and random (generic).
- Geometry of the function landscape provides important clues for algorithm design and analysis.



... starting with *sparse dictionary learning!*

# Sparse dictionary learning



$$Y \approx QX \quad X \in \mathbb{R}^{n \times p} \quad \text{sparse}$$

- Algorithmic study initialized with [Olshausen, Field. '96] in neuroscience.
- Important algorithmic contributions from many researchers: [Lewicki, Sejnowski.'99], [Engan et al. '99], [Aharon, Elad, Bruckstein. '06], many others
- Widely used in image processing, vision, audio, and machine learning



# Dictionary learning - the complete case



$$Y \approx QX \quad X = \Omega \odot V, \quad \Omega \sim \text{Ber}(\theta), \quad V \sim \mathcal{N}(0, 1).$$

# Dictionary learning - the complete case



$$Y \approx QX \quad X = \Omega \odot V, \quad \Omega \sim \text{Ber}(\theta), \quad V \sim \mathcal{N}(0, 1).$$

- Assume  $Q$  is square and nonsingular, then  $\text{row}(Y) = \text{row}(X)$ .
- When  $p \geq \Omega(n \log n)$ , rows of  $X$  are the sparsest vectors in  $\text{row}(Y)$  [Spielman, Wang, Wright. '12]

# Dictionary learning - the complete case



$$Y \approx QX \quad X = \Omega \odot V, \quad \Omega \sim \text{Ber}(\theta), \quad V \sim \mathcal{N}(0, 1).$$

- Assume  $Q$  is square and nonsingular, then  $\text{row}(Y) = \text{row}(X)$ .
- When  $p \geq \Omega (n \log n)$ , rows of  $X$  are the sparsest vectors in  $\text{row}(Y)$  [Spielman, Wang, Wright. '12]

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

if we recover one row (up to scaling) of  $X$ , then we use *deflation* to find the others.

# Dictionary learning: the complete case

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

- Convex relaxation:

$$\min \quad \|\mathbf{q}^* \mathbf{Y}\|_1 \quad \text{s.t. } \mathbf{r}^* \mathbf{q} = 1$$

Provably succeeds when  $\theta n = O(\sqrt{n})$ , provably fails if  $\theta n = \Omega(\sqrt{n \log n})$  [Spielman, Wang, Wright.'12].

- Nonconvex relaxation:

## Model problem

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

many precedents, e.g., [Zibulevsky-Perlmutter, '01] in source separation.

# The model problem

## Model problem

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

- Convex objective function, but nonconvex constraint  $\mathbf{q} \in \mathbb{S}^{n-1}$ .
- If  $p \geq \Omega(n \log n)$ , w.h.p. every global optimizer  $\mathbf{q}_\diamond$  produces  $\mathbf{q}_\diamond^* \mathbf{Y}$  that recovers one row of  $\mathbf{X}$  (up to scaling)

# Towards geometric understanding

## Model problem

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

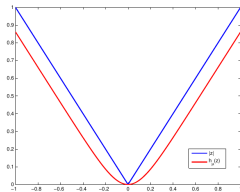
## Slightly modified model problem

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

- Work with a *smooth surrogate* for  $|z|$ :

$$h_\mu(z) = \mu \log \left( \frac{e^{z/\mu} + e^{-z/\mu}}{2} \right)$$

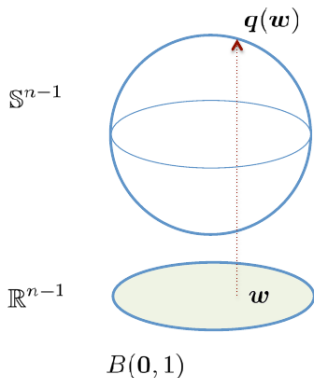
- Recognize the objective as a *normalized sum of independent random variables*  $\rightarrow$  expectation, asymptotically





# Why might this work?

W.l.o.g., assume for analysis that  $Y = X$  (i.e.,  $A = I$ ); we correctly recover a row iff an algorithm produces  $\mathbf{q}_\diamond = \pm \mathbf{e}_i$ ,  $i = 1, \dots, n$



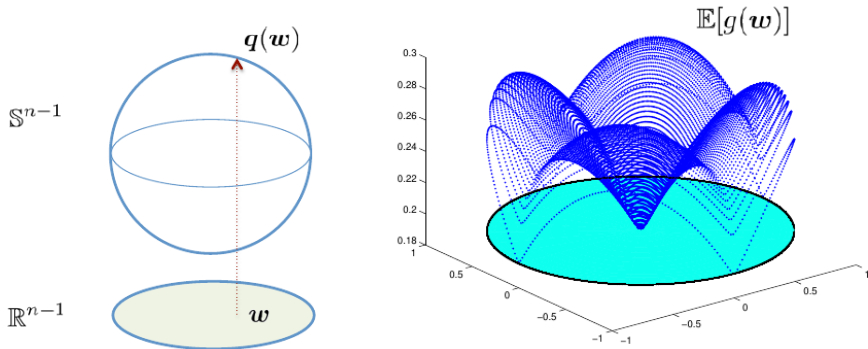
- $f(\mathbf{q}) = \frac{1}{p} \sum_{i=1}^p h_\mu(\mathbf{q}^* \mathbf{x}_i)$
- $\mathbf{q}(\mathbf{w}) = \left( w_1, \dots, w_{n-1}, \sqrt{1 - \|\mathbf{w}\|_2^2} \right)$
- $g(\mathbf{w}) = f(g(\mathbf{w}))$

Target points are

$$\mathbf{w}_\diamond = \mathbf{0}, \pm \mathbf{e}_1, \dots, \pm \mathbf{e}_{n-1}$$

# Why might this work?

W.l.o.g., assume for analysis that  $Y = X$  (i.e.,  $A = I$ ); we correctly recover a row iff an algorithm produces  $q_\diamond = \pm e_i$ ,  $i = 1, \dots, n$

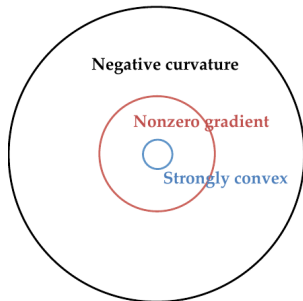


$B(0,1)$

Target points are

$$w_\diamond = \mathbf{0}, \pm e_1, \dots, \pm e_{n-1}$$

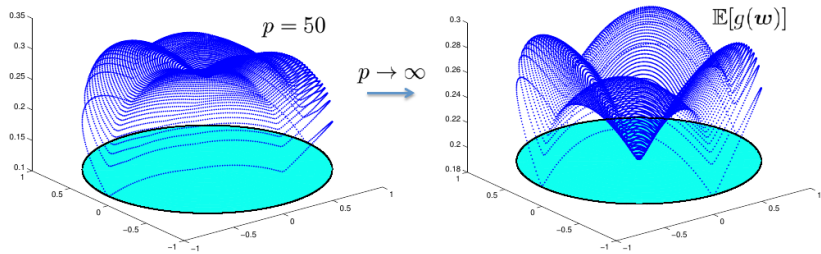
**Lemma:** Suppose  $\theta \in [\frac{1}{n}, \frac{1}{2})$ , and  $\mu < cn^{-5/4}$ . Then



- $\nabla^2 \mathbb{E}[g(\mathbf{w})] \preceq \frac{c'\theta}{\mu} \mathbf{I}$ , for  $\|\mathbf{w}\| \leq \frac{\mu}{4\sqrt{2}}$
- $\frac{\mathbf{w}^* \nabla \mathbb{E}[g(\mathbf{w})]}{\|\mathbf{w}\|} \geq c'\theta$ , for  $\frac{\mu}{4\sqrt{2}} \leq \|\mathbf{w}\| \leq R$
- $\frac{\mathbf{w}^* \nabla \mathbb{E}[g(\mathbf{w})] \mathbf{w}}{\|\mathbf{w}\|^2} \leq -c'\theta$ , for  $\|\mathbf{w}\| \geq R$   
( $R = \Theta(1)$ ).

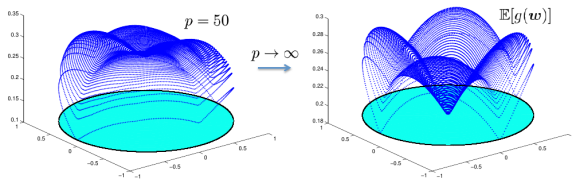
and so, **every local optimizer of  $\mathbb{E}[g(\mathbf{w})]$  is a target point.**

# Convergence in function landscape



When does the finite-sample objective **converge** to the asymptotic one, in optimization sense?  
...informally, is the function geometry “nice” for some large but finite  $p$ ?

# Finite-sample result



- Objective  $g(\mathbf{w}) = \frac{1}{p} \sum_{i=1}^p h_{\mu}(q(\mathbf{w})^* \mathbf{x}_i)$  is a sum of independent RVs.
- The proof follows a typical concentration-expectation path

## Lemma

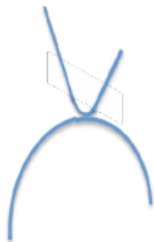
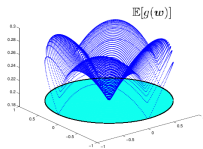
Suppose  $\theta \in [\frac{1}{n}, \frac{1}{2})$ , if  $\mu < cn^{-5/4}$ , and  $p \geq \frac{Cn^3}{\mu^2\theta^2} \log n$ , it holds uniformly w.h.p. that

- $\nabla^2 \mathbb{E}[g(\mathbf{w})] \preceq \frac{c''\theta}{\mu} \mathbf{I}$ , for  $\|\mathbf{w}\| \leq \frac{\mu}{4\sqrt{2}}$
- $\frac{\mathbf{w}^* \nabla \mathbb{E}[g(\mathbf{w})]}{\|\mathbf{w}\|} \geq c''\theta$ , for  $\frac{\mu}{4\sqrt{2}} \leq \|\mathbf{w}\| \leq R$
- $\frac{\mathbf{w}^* \nabla \mathbb{E}[g(\mathbf{w})] \mathbf{w}}{\|\mathbf{w}\|^2} \leq -c''\theta$ , for  $\|\mathbf{w}\| \geq R$  ( $R = \Theta(1)$ ).

... following intuition we build from the geometry:

- Don't know the chart  $q(w)$  ahead of time  $\implies$  work directly on  $q \in \mathbb{S}^{n-1}$ .
- Pull the "niceness" back to the sphere: descent direction in  $w \iff$  descent direction in  $q$  along some curve
- Need to escape saddle points  $\implies$  Use second-order information. Here, via the **trust region method**.

Trust-region on manifolds [Absil, Baker, Gallivan. '07], also [Absil, Mahoney, Sepulchre. '08]

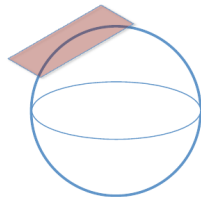


# Trust region method

Consider  $\mathbf{q} \in \mathbb{S}^{n-1}$ ; for  $\boldsymbol{\delta} \perp \mathbf{q}$ , calculus gives

$$\begin{aligned} f(\exp_{\mathbf{q}}(\boldsymbol{\delta})) &= f(\mathbf{q}) + \langle \boldsymbol{\delta}, \nabla f(\mathbf{q}) \rangle + \frac{1}{2} \boldsymbol{\delta}^* (\nabla^2 f(\mathbf{q}) - \langle \mathbf{q}, \nabla f(\mathbf{q}) \rangle) \boldsymbol{\delta} + O(\|\boldsymbol{\delta}\|_2^3) \\ &\doteq \widehat{f}(\mathbf{q}; \boldsymbol{\delta}) + O(\|\boldsymbol{\delta}\|_2^3) \end{aligned}$$

where  $\exp_{\mathbf{q}}(\boldsymbol{\delta}) = \mathbf{q} \cos(\|\boldsymbol{\delta}\|_2) + \frac{\boldsymbol{\delta}}{\|\boldsymbol{\delta}\|_2} \sin(\|\boldsymbol{\delta}\|_2)$ .



Basic Riemannian trust region method:

$$\begin{aligned} \boldsymbol{\delta}_\star &\in \arg \min_{\boldsymbol{\delta} \in T_{\mathbf{q}_k} \mathbb{S}^{n-1}, \|\boldsymbol{\delta}\|_2 \leq \Delta} \widehat{f}(\mathbf{q}_k; \boldsymbol{\delta}) \\ \mathbf{q}_{k+1} &= \exp_{\mathbf{q}_k}(\boldsymbol{\delta}_\star). \end{aligned}$$

The trust region subproblem involves a (possibly nonconvex) quadratic and one quadratic constraint. Solvable in polynomial time by root finding [More+Sorensen '82] or SDP relaxation.

# Pull things together

## Theorem (informal)

Suppose that  $\theta \in [\frac{1}{n}, \frac{1}{2})$ ,  $\mu < cn^{-5/4}$ ,  $\mathbf{Y} = \mathbf{Q}\mathbf{X}$  with  $\mathbf{Q} \in O(n)$ . If we observe  $p \geq \text{poly}(n)$  samples, then applying the trust region method with fixed radius  $\Delta = \frac{1}{\text{poly}(n)}$  for  $T = \text{poly}(n)$  iterations. W.h.p, the algorithm produces a  $\hat{\mathbf{q}}$  such that

$$\|\hat{\mathbf{q}} - \mathbf{q}_\diamond\| \leq C \frac{\mu}{\theta} \sqrt{\frac{n \log p}{p}}$$

for some target solution  $\mathbf{q}_\diamond$  satisfying  $\mathbf{q}_\diamond^* \mathbf{Y} = \pm e_i^* \mathbf{X}$ .

- Using **linear programming rounding + deflation**, one can recover all of  $\mathbf{X}$ , and subsequently  $\mathbf{Q}$ .
- If  $\mathbf{Q}$  is not an orthobasis, apply preconditioning, but need  $p \geq \text{poly}(n, \sigma_{\min}(\mathbf{Q})^{-1})$ .



# Comparison with the Literature

- **Efficient algorithms** with performance guarantees

[Spielman, Wang, Wright,'12]

$$\mathbf{Q} \in \mathbb{R}^{n \times n}, \theta = \tilde{O}(1/\sqrt{n})$$

[Agarwal, Anandkumar, Netrapali,'13]

$$\mathbf{Q} \in \mathbb{R}^{m \times n} (m \leq n), \theta = \tilde{O}(1/\sqrt{n})$$

[Arora, Ge, Moitra,'13]

$$\mathbf{Q} \in \mathbb{R}^{m \times n} (m \leq n), \theta = \tilde{O}(1/\sqrt{n})$$

- **Quasipolynomial algorithms** with better guarantees

[Spielman, Wang, Wright,'12]

$$\mathbf{Q} \in \mathbb{R}^{n \times n}, \theta = \tilde{O}(1/\sqrt{n})$$

[Arora, Bhaskara, Ge, Ma,'14]

different prob. model,  $\theta = O(1/\text{polylog}(n))$

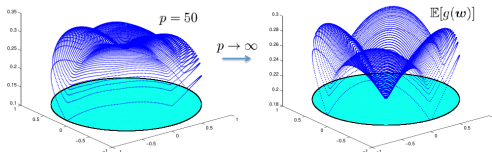
[Barak, Kelner, Steurer,'14]

sum-of-squares,  $\theta = \tilde{O}(1)$

- Other theoretic work on **local geometry**: [Gribonval, Schnass'11], [Geng, Wright, '11], [Schnass'14]

This work: a polynomial algorithm for squared  $\mathbf{Q}$ ,  $\theta = O(1)$ .

# What we have done so far...



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

- Prove as  $p$  becomes **large**, the nonconvex program becomes tractable under our **probabilistic setting**.
- Geometry has guided our analysis and algorithm design.

Related publications:

- Sun, Qu, Wright. Complete dictionary recovery over the sphere. In preparation.
- Qu, Sun, Wright. Finding a sparse vector in a subspace: Linear sparsity using alternating directions. NIPS'14.

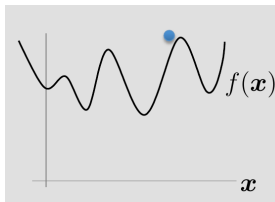
For sparse dictionary learning:

- Streamline the proof and work directly in manifold language
- With directly with the  $\|\cdot\|_1$  function (partial progress)
- Does similar thing happens if we look at structured dictionary directly? (orthogonal group - very likely; tight-frame - likely)
- Algorithm side: how to understand the surprisingly successful alternating direction method in this setting?
- ...

# Proving nonconvex recovery

## Theme of this thesis

- Certain nonconvex optimization problems become tractable when the input data are large and random (generic).
- Geometry of the function landscape provides important clues for algorithm design and analysis.



1. Use problem structure to find a **clever initial guess**.
2. Analyze iteration-by-iteration **in the vicinity of the optimum**.

- **Matrix completion:** [Keshevan, Oh, Montanari.'09], [Jain, Netrapali, Sanghavi.'13], [Hardt'13], [Hardt, Wooters.'14]. Also [Meta, Jain, Dhillon.'09]
- **Dictionary learning:** [Agarwal, Anandkumar, Netrapali.'13], [Arora, Ge, Moitra.'13], [Agarwal, Anandkumar, Jain, Netrapali.'13]
- **Tensor recovery:** [Jain, Oh.'13], [Anandkumar, Ge, Janzamin.'14]
- **Phase retrieval:** [Netrapali, Jain, Sanghavi.'13], [Candes, Li, Soltanokoltabi.'14]

For the analytic strategy:

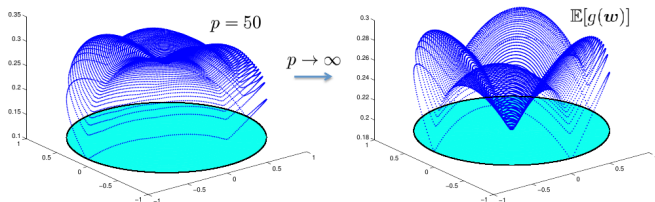
## Generalized Model Problem

$$\min \quad \frac{1}{p} \sum_{i=1}^p f_i(\mathbf{q}) \quad \text{s.t. } \mathbf{q} \in \mathcal{M}.$$

$f_i$ 's are independent, and  $\mathcal{M}$  is some Riemannian manifold ( $\mathbb{S}^{n-1}$ ,  $O(n)$ ,  $\{\mathbf{X} : \text{rank}(\mathbf{X}) = r\}$ , etc)

- Other problems: phase retrieval, matrix/tensor recovery, recovery of signal with simultaneous structures, **blind deconvolution**, etc

THANKS to Prof. John Wright and Mr. Qing Qu.



Questions?