

Taming nonconvexity: from smooth to nonsmooth problems

Ju Sun

Department of Mathematics
Stanford University

SINE Seminar at Coordinated Science Laboratory, UIUC

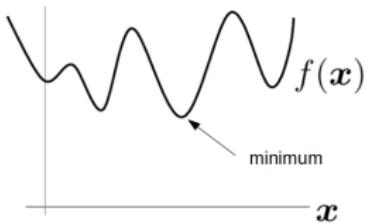
November 5, 2018

Optimization is everywhere

*"Nothing takes place in the world whose meaning is not that of some **maximum or minimum.**"*

$$\min f(x)$$

s. t. $x \in \mathcal{S}$.

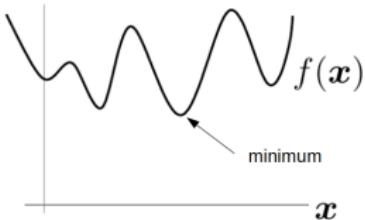


Leonhard Euler

Optimization is everywhere

*"Nothing takes place in the world whose meaning is not that of some **maximum or minimum.**"*

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



Leonhard Euler

Historic heroes

HISTORICAL PERSPECTIVES

300 Years of Optimal Control: From The Brachystochrone to the Maximum Principle

Hector J. Sussmann and Jan C. Willems

Optimal control was born in 1697—300 years ago—in Groningen, a university town in the north of The Netherlands, when Johann Bernoulli, professor of mathematics at the local university from 1695 to 1705, published his solution of the brachystochrone problem. The year before he had challenged his contemporaries to solve this problem. We will tell the story of

the authors. We gladly plead guilty to most of this charge—and state for the record that we are both control theorists, and one of us is a professor at Groningen—asking only that the word “merely” be stricken out. Our biases may of course explain how

Euclid	de Fermat
Newton	Leibniz
Bernoulli's	Euler
Lagrange	Legendre
Gauss	Fourier
Cauchy	Hadamard
	...

Optimization is everywhere

computer vision

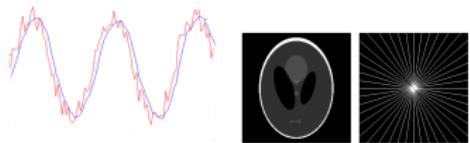


Optimization is everywhere

computer vision



signal processing

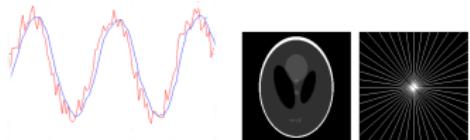


Optimization is everywhere

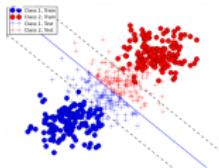
computer vision



signal processing



machine learning



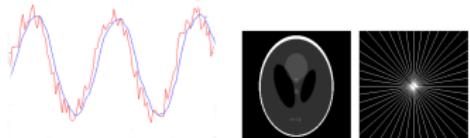
$$P(A|B) \propto P(B|A)P(A)$$

Optimization is everywhere

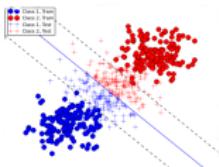
computer vision



signal processing

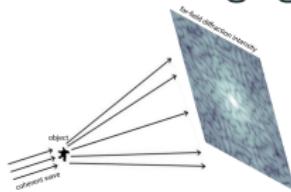


machine learning



$$P(A|B) \propto P(B|A)P(A)$$

scientific imaging

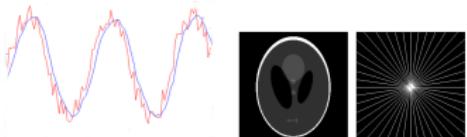


Optimization is everywhere

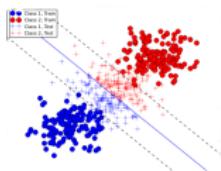
computer vision



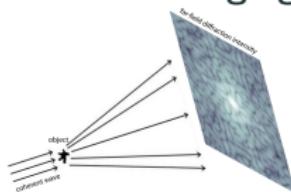
signal processing



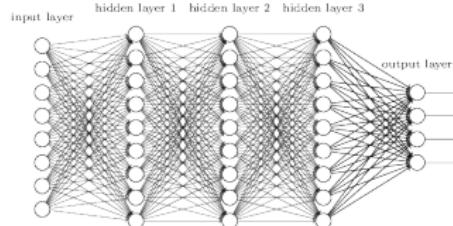
machine learning



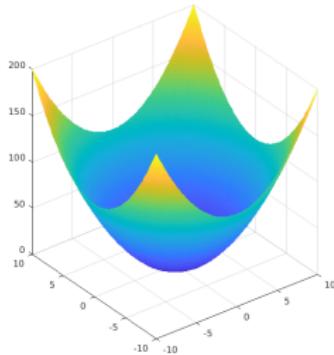
scientific imaging



neural networks

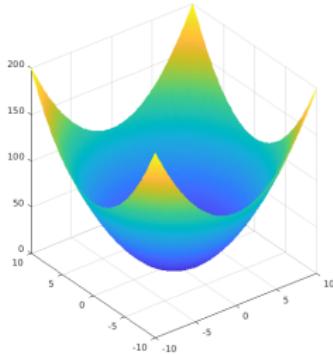


Convex analysis and optimization



All local minimizers are **global!**
(All critical points are **global!**)

Convex analysis and optimization



All local minimizers are **global!**
(All critical points are **global!**)

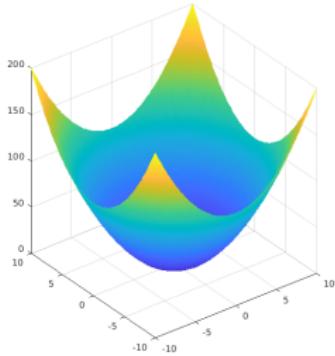
Interior-point method (80's–00's)

– “Most” convex problems can be solved **efficiently**!

Modeling languages (00's–10's)

```
cvx_begin
    variable x(n)
    minimize( norm( A * x - b, 2 ) )
    subject to
        C * x == d
        norm( x, Inf ) <= e
cvx_end
```

Convex analysis and optimization



All local minimizers are **global!**
(All critical points are **global!**)

*...in fact, the great watershed in optimization isn't between linearity and nonlinearity, but **convexity** and **nonconvexity**.*

— R. Tyrrell Rockafellar [Rockafellar, 1993]

Interior-point method (80's–00's)

— “Most” convex problems can be solved **efficiently**!

Modeling languages (00's–10's)

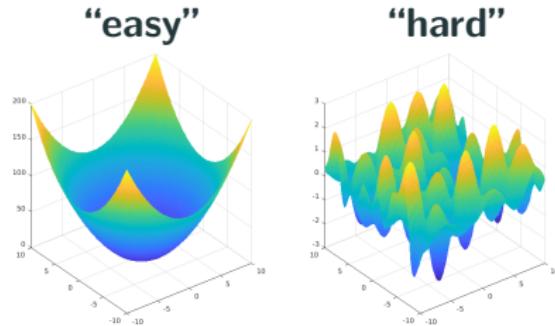
```
minimize ||Ax - b||_2
subject to Cx = d
        ||x||_\infty \leq e
```

```
cvx_begin
    variable x(n)
    minimize( norm( A * x - b, 2 ) )
    subject to
        C * x == d
        norm( x, Inf ) <= e
cvx_end
```



Nonconvex optimization?

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

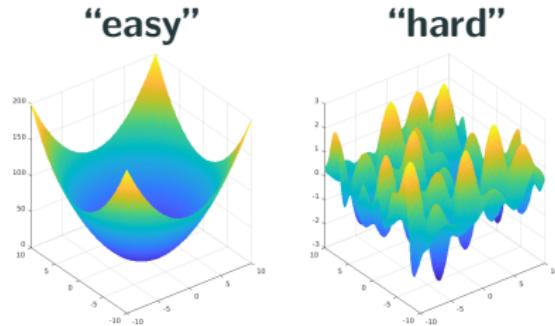


All local mins are global!

Spurious local mins!

Nonconvex optimization?

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

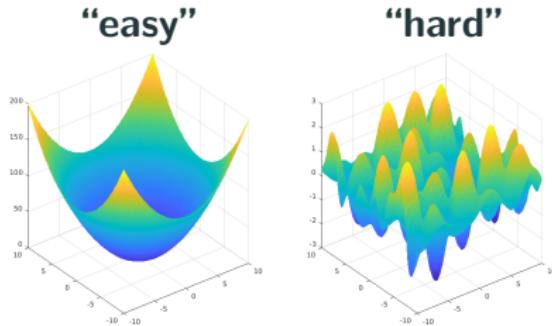


All local mins are global! Spurious local mins!

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

Nonconvex optimization?

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

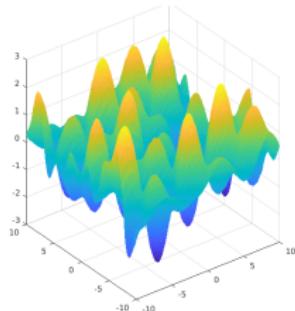


All local mins are global! Spurious local mins!

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

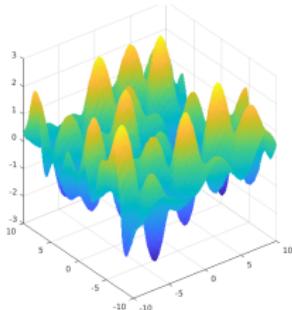
Many problems in modern **signal processing, machine learning, statistics, imaging**, ..., are most naturally formulated as **nonconvex** optimization problems.

Nonscary nonconvex optimization



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

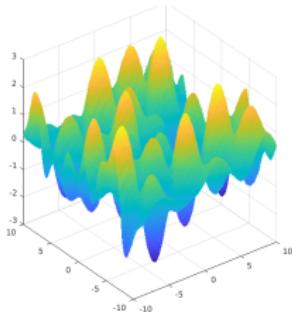
Nonscary nonconvex optimization



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

In practice: Heuristic algorithms are often surprisingly successful.

Nonscary nonconvex optimization

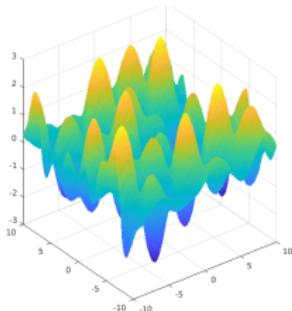


In theory: Even computing a local minimizer is NP-hard!

In practice: Heuristic algorithms are often surprisingly successful.

e.g., training deep neural networks—SGD and variants, plus a few tricks

Nonscary nonconvex optimization



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

In practice: Heuristic algorithms are often surprisingly successful.

e.g., training deep neural networks—SGD and variants, plus a few tricks

Which nonconvex optimization problems are easy?

Running app: sparsifying transform learning

Given data \mathbf{Y} , learn \mathbf{Q} , st $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small

Running app: sparsifying transform learning

Given data \mathbf{Y} , learn \mathbf{Q} , st $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small

The screenshot shows a website for "TRANSFORM LEARNING". The header features the title "TRANSFORM LEARNING" and a sub-section "Sparse Representations at Scale". Below the header is a navigation bar with links: HOME, PROJECTS, PUBLICATIONS, SOFTWARE, and CONTACT. The main content area is titled "Overview". It discusses the use of sparsity in various applications like signal and image processing, machine learning, and medical imaging. It highlights analytical sparsifying transforms like Wavelets and DCT, and data-driven learning of sparse models like synthesis and analysis dictionaries. It also mentions compressed sensing and online learning. The text notes several proposed methods for batch learning of square or overcomplete sparsifying transforms, including double sparsity, union-of-transforms, and filter bank structures. It also covers online learning of sparsifying transforms for big data and real-time applications. The website has demonstrated promising performance in sparse representation, image and video denoising, classification, and compressed sensing (MRI and CT image reconstruction) tasks. Several convergence guarantees are mentioned for transform learning and image reconstruction schemes. The footer contains a "Software" section with a note about available implementations and data for publications.

TRANSFORM LEARNING

Sparse Representations at Scale

HOME PROJECTS PUBLICATIONS SOFTWARE CONTACT

Overview

The sparsity of signals and images in a certain transform domain or dictionary has been exploited in many applications in signal and image processing, machine learning, and medical imaging. Analytical sparsifying transforms such as Wavelets and DCT have been widely used in compression standards. Recently, the data-driven learning of sparse models such as the synthesis dictionary model and the analysis expansion model have been applied in denoising, compressed sensing, etc. Our group's research at the University of Illinois focuses on the data-driven adaptation of the alternative sparsifying transform model, which offers numerous advantages over the synthesis dictionary model.

We have proposed several methods for batch learning of square or overcomplete sparsifying transforms from data. We have also investigated specific structures for these transforms such as double sparsity, union-of-transforms, and filter bank structures, which enable their efficient learning or usage. Apart from batch transform learning, our group has investigated methods for online learning of sparsifying transforms, which are particularly useful for big data or real-time applications. The proposed algorithms for transform learning have been shown to be highly efficient.

Our research has demonstrated promising performance for transform learning methods in sparse representation, image and video denoising, classification, and compressed sensing (MRI and CT image reconstruction) tasks. We also established several convergence guarantees for our transform learning or image reconstruction schemes, which were previously lacking for prior adaptive dictionary-based methods.

This website contains the various manuscripts and theses published by our group on transform learning. It also contains conference posters and oral presentation slides on the various proposed methods.

Software

Software implementations of the various algorithms and data used to generate the results in our publications are available from the [Software](#) tab.

image credit: Professor Yoram Bresler's research website

Running app: sparsifying transform learning

Given data \mathbf{Y} , learn \mathbf{Q} , st $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small

The screenshot shows a website for "TRANSFORM LEARNING". The header features the title "TRANSFORM LEARNING" and a sub-section "Sparse Representations at Scale". Below the header is a navigation bar with links: HOME, PROJECTS, PUBLICATIONS, SOFTWARE, and CONTACT. The main content area is titled "Overview". It discusses the sparsity of signals and images in a certain transform domain or dictionary, mentioning applications in signal and image processing, machine learning, and medical imaging. It highlights analytical sparsifying transforms like Wavelets and DCT, and data-driven learning of sparse models such as synthesis and analysis dictionaries. The text also covers denoising, compressed sensing, etc. A note mentions the adaptation of the alternative sparsifying transform model over the synthesis dictionary model. Another section describes batch learning of square or overcomplete sparsifying transforms from data, mentioning double sparsity, union-of-transforms, and filter bank structures, along with efficient learning and usage. It also covers online learning of sparsifying transforms for big data or real-time applications. A third section discusses promising performance for transform learning methods in sparse representations for image and video denoising, classification, and compressed sensing (MRI and CT image reconstruction tasks). It includes convergence guarantees for transform learning or image reconstruction schemes, previously lacking for prior adaptive dictionary-based methods. The footer contains a "Software" section with a note about software implementations and data used in publications, available from the Software tab.

image credit: Professor Yoram Bresler's research website

- **Dictionary learning:** factor \mathbf{Y} as $\mathbf{Y} \approx \mathbf{AX}$ st \mathbf{X} is sparse

Running app: sparsifying transform learning

Given data \mathbf{Y} , learn \mathbf{Q} , st $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small

The screenshot shows a website for "TRANSFORM LEARNING". The header features the title "TRANSFORM LEARNING" and a sub-section "Sparse Representations at Scale". Below the header is a navigation bar with links: HOME, PROJECTS, PUBLICATIONS, SOFTWARE, and CONTACT. The main content area is titled "Overview". It discusses the use of sparsity in signals and images across various applications like signal and image processing, machine learning, and medical imaging. It highlights analytical sparsifying transforms like Wavelets and DCT, and data-driven learning models. It also mentions compressed sensing and online learning of sparsifying transforms. A note on double sparsity, union-of-transforms, and filter bank structures is included. The "Software" section links to implementations of algorithms and data used in publications.

image credit: Professor Yoram Bresler's research website

- **Dictionary learning:** factor \mathbf{Y} as $\mathbf{Y} \approx \mathbf{AX}$ st \mathbf{X} is sparse
- Apps: image processing, computer vision, computational imaging
[Mairal et al., 2014]

Running app: sparsifying transform learning

Given data \mathbf{Y} , learn \mathbf{Q} , st $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small

The screenshot shows a website for "TRANSFORM LEARNING". The header features the title "TRANSFORM LEARNING" and a sub-section "Sparse Representations at Scale". Below the header is a navigation menu with links to "HOME", "PROJECTS", "PUBLICATIONS", "SOFTWARE", and "CONTACT". The main content area is titled "Overview" and contains several paragraphs of text. It discusses the sparsity of signals and images in a certain transform domain or dictionary, mentioning applications in signal and image processing, machine learning, and medical imaging. It highlights analytical sparsifying transforms like Wavelets and DCT, and more recent data-driven learning of sparse models such as synthesis and analysis dictionaries. The text also covers applications in compressed sensing, denoising, and compressed sensing, emphasizing the group's focus on the data-driven adaptation of alternative sparsifying transform models over synthesis models. A note below states that proposed algorithms for transform learning have been shown to be highly efficient. Another section, "Software", describes software implementations of various algorithms and data used for results in publications, available from the Software tab.

image credit: Professor Yoram Bresler's research website

- **Dictionary learning:** factor \mathbf{Y} as $\mathbf{Y} \approx \mathbf{AX}$ st \mathbf{X} is sparse
- Apps: image processing, computer vision, computational imaging [Mairal et al., 2014]
- Cascaded with nonlinearity [Ravishankar and Wohlberg, 2018]

One element each time

Given \mathbf{Y} , learn **ortho** \mathbf{Q} s.t. $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small.

One element each time

Given \mathbf{Y} , learn **ortho** \mathbf{Q} s.t. $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small.

A naive formulation:

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

One element each time

Given \mathbf{Y} , learn **ortho** \mathbf{Q} s.t. $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small.

A naive formulation:

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

Nonconvex “relaxation”:

$$\min_{\mathbf{q}} \quad f(\mathbf{q}) \doteq \frac{1}{m} \|\mathbf{q}^*\mathbf{Y}\|_1 \quad \text{s.t.} \quad \|\mathbf{q}\|_2^2 = 1.$$

Many precedents, e.g., [Zibulevsky and Pearlmutter, 2001] in blind source separation.

Here, inspired by [Spielman et al., 2012, Sun et al., 2015]

Recovery setting for analysis

Given \mathbf{Y} , learn **ortho** \mathbf{Q} s.t. $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small.

Recovery setting for analysis

Given \mathbf{Y} , learn **ortho** \mathbf{Q} s.t. $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small.

To study possibility of **recovery**, given \mathbf{Q}_0 ortho and \mathbf{X}_0 sparse,

$$\mathbf{Q}_0^* \times \mathbf{Y} = \mathbf{X}_0,$$

recover \mathbf{Q}_0 and \mathbf{X}_0 (**up to signed permutation**) .

Nonconvex “relaxation”:

$$\min f(\mathbf{q}) \doteq \frac{1}{m} \|\mathbf{q}^*\mathbf{Y}\|_1 \quad \text{s.t. } \|\mathbf{q}\|_2^2 = 1.$$

Toward geometric intuition

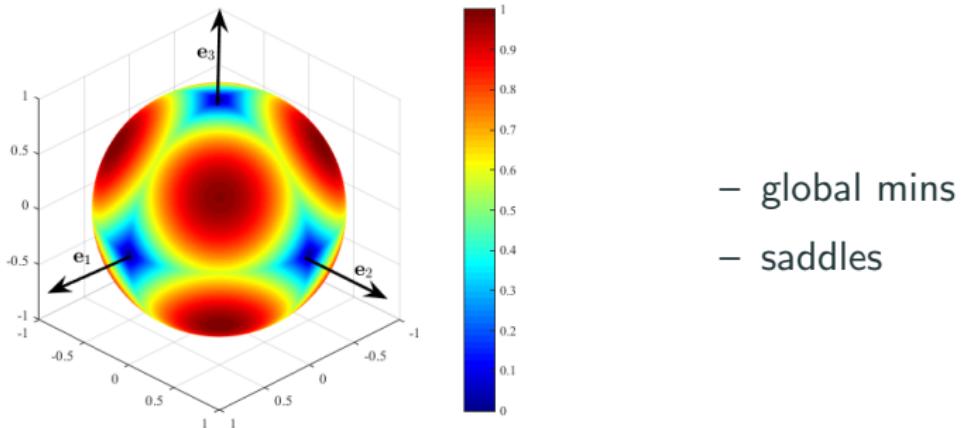
$$\min \quad f(\mathbf{q}) \doteq \frac{1}{m} \|\mathbf{q}^* \mathbf{Y}\|_1 \quad \text{s.t. } \|\mathbf{q}\|_2^2 = 1.$$

A low-dimensional example ($n = 3$) of the landscape when the target transformation $\mathbf{Q}_0 = \mathbf{I}$ and $m \rightarrow \infty$

Toward geometric intuition

$$\min \quad f(\mathbf{q}) \doteq \frac{1}{m} \|\mathbf{q}^* \mathbf{Y}\|_1 \quad \text{s.t. } \|\mathbf{q}\|_2^2 = 1.$$

A low-dimensional example ($n = 3$) of the landscape when the target transformation $\mathbf{Q}_0 = \mathbf{I}$ and $m \rightarrow \infty$

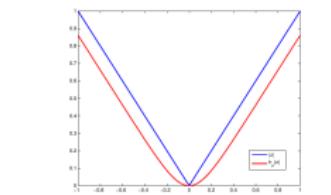


Toward geometric analysis

Smoothed model problem

$$\min f_1(\mathbf{q}) \doteq \frac{1}{m} \sum_{j=1}^m |\mathbf{q}^* \mathbf{y}_j| \text{ s.t. } \|\mathbf{q}\|_2^2 = 1.$$

↓↓ smoothing ↓↓



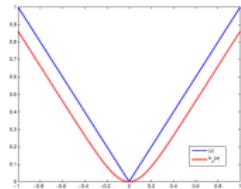
$$h_\mu(z) = \mu \log \cosh \frac{z}{\mu}$$

$$\min f(\mathbf{q}) \doteq \frac{1}{m} \sum_{i=1}^m h_\mu(\mathbf{q}^* \mathbf{y}_i) \text{ s.t. } \|\mathbf{q}\|_2^2 = 1.$$

Toward geometric analysis

Smoothed model problem

$$\min f_1(\mathbf{q}) \doteq \frac{1}{m} \sum_{j=1}^m |\mathbf{q}^* \mathbf{y}_j| \text{ s.t. } \|\mathbf{q}\|_2^2 = 1.$$



↓↓ smoothing ↓↓

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

$$\min f(\mathbf{q}) \doteq \frac{1}{m} \sum_{i=1}^m h_\mu(\mathbf{q}^* \mathbf{y}_i) \text{ s.t. } \|\mathbf{q}\|_2^2 = 1.$$

For analysis: Bernoulli-Gaussian model $\mathbf{X}_0 = \boldsymbol{\Omega}_0 \circ \mathbf{V}_0$,
 $\boldsymbol{\Omega}_0 \sim_{\text{iid}} \text{Ber}(\theta)$, $\mathbf{V}_0 \sim_{\text{iid}} \mathcal{N}(0, 1)$.

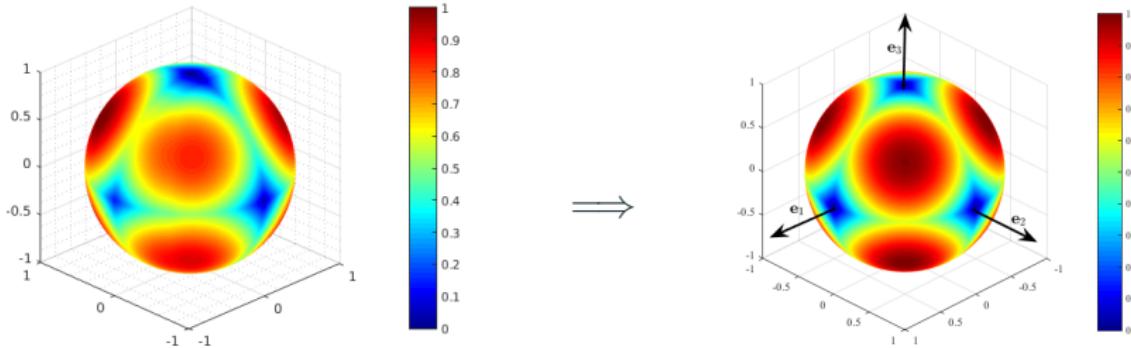
Sparsity parameter θ ; average number of nonzeros per column is θn .

The geometric result

$$\min f(\mathbf{q}) \doteq \frac{1}{m} \sum_{i=1}^m h_\mu (\mathbf{q}^* \mathbf{y}_i) \text{ s.t. } \|\mathbf{q}\|_2^2 = 1.$$

The geometric result

$$\min f(\mathbf{q}) \doteq \frac{1}{m} \sum_{i=1}^m h_\mu (\mathbf{q}^* \mathbf{y}_i) \text{ s.t. } \|\mathbf{q}\|_2^2 = 1.$$



Theorem (Informal, [Sun et al., 2015])

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

All local minimizers are “global”.

Comparison with prior results

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 $\theta = O(n^{-\varepsilon})$.

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 learning complete Q with $\theta = \Omega(1)$.

See also recent refined SOS analysis [Ma et al., 2016a] with similar guarantees.

[Back to the question](#)

Which nonconvex optimization problems are easy?

... two types of partial answers

A1: Problems with nice global landscapes

(P-1) All local minimizers are **global**

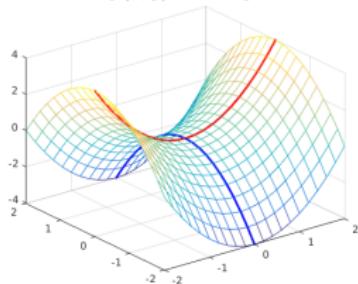
(P-2) All saddle points (and local maximizers) have a **directional negative curvature**, i.e., $\lambda_{\min}(\text{Hess}) < 0$

A1: Problems with nice global landscapes

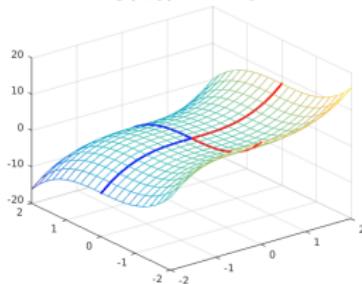
(P-1) All local minimizers are **global**

(P-2) All saddle points (and local maximizers) have a **directional negative curvature**, i.e., $\lambda_{\min}(\text{Hess}) < 0$

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



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



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

Ridable/strict saddle (also

[Ge et al., 2015])

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

locally shaped by high-order derivatives at

A1: Problems with nice global landscapes

All local mins are global, all saddles are strict

A1: Problems with nice global landscapes

All local mins are global, all saddles are strict

Eigenvalue problems (folklore!)

Sparsifying dictionary learning [Sun et al., 2015]

Generalized phase retrieval [Sun et al., 2016]

Orthogonal tensor decomposition [Ge et al., 2015]

Low-rank matrix recovery and completion

[Ge et al., 2016, Bhojanapalli et al., 2016]

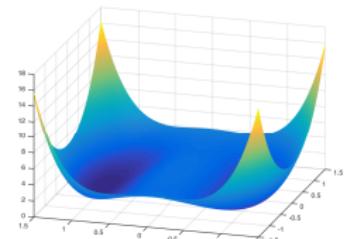
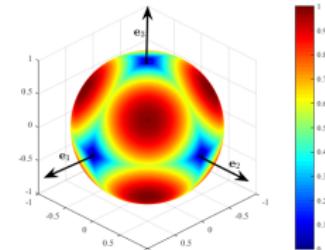
Phase synchronization [Boumal, 2016]

Community detection [Bandeira et al., 2016]

Deep/shallow networks [Kawaguchi, 2016,

Lu and Kawaguchi, 2017, Soltanolkotabi et al., 2017]

Sparse blind deconvolution [Zhang et al., 2017]



A1: Problems with nice global landscapes

All local mins are global, all saddles are strict

Eigenvalue problems (folklore!)

Sparsifying dictionary learning [Sun et al., 2015]

Generalized phase retrieval [Sun et al., 2016]

Orthogonal tensor decomposition [Ge et al., 2015]

Low-rank matrix recovery and completion

[Ge et al., 2016, Bhojanapalli et al., 2016]

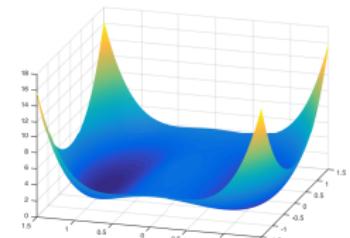
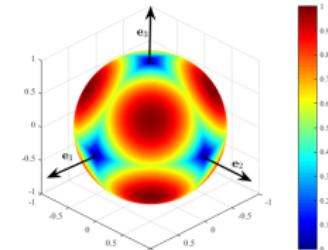
Phase synchronization [Boumal, 2016]

Community detection [Bandeira et al., 2016]

Deep/shallow networks [Kawaguchi, 2016,

Lu and Kawaguchi, 2017, Soltanolkotabi et al., 2017]

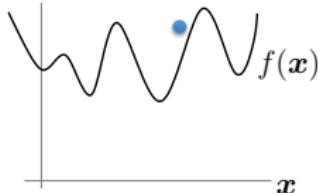
Sparse blind deconvolution [Zhang et al., 2017]



Algorithms: virtually everything reasonable works!

[Conn et al., 2000, Nesterov and Polyak, 2006, Goldfarb, 1980, Jin et al., 2017]

A2: Problems with nice local landscapes



Use problem structure to find a clever (sometimes random) initial guess.

Analyze local search algorithms 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], [Chatterji and Bartlett, 2017], [Gilboa et al., 2018]
- **Tensor recovery:** [Jain and Oh, 2014], [Anandkumar et al., 2014b], [Anandkumar et al., 2014a], [Anandkumar et al., 2015]
- **Phase retrieval:** [Netrapalli et al., 2013], [Candès et al., 2015], [Chen and Candès, 2015], [White et al., 2015], [Wang et al., 2016], [Chen et al., 2018]

Nonscary nonconvex problems

Problems with nice global/local landscapes

- My webpage: <http://sunju.org/research/nonconvex/>
- Sun, Ju and Qu, Qing and Wright, John. **When are nonconvex problems not scary?**. arXiv preprint arXiv:1510.06096 (2015).
- Jain, Prateek and Kar, Purushottam. **Non-convex optimization for machine learning**. Foundations and Trends® in Machine Learning 10.3–4 (2017): 142–336.
- Chen, Yudong and Chi, Yuejie. **Harnessing structures in big data via guaranteed low-rank matrix estimation**. arXiv preprint arXiv:1802.08397 (2018).
- Chi, Yuejie and Lu, Yue M., and Chen, Yuxin. **Nonconvex Optimization Meets Low-Rank Matrix Factorization: An Overview**. arXiv preprint arXiv:1809.09573 (2018).

Common ingredients in analysis

For **smooth** problems,

1st order geometry: ∇f or $v^\top \nabla f$ (directional derivatives)

2nd order geometry: $\nabla^2 f$ or $v^\top \nabla^2 f v$ (directional curvatures)

What about nonsmooth, nonconvex problems?

nonsmooth: may be non-differentiable

Nonsmooth problems are everywhere

Optimization: exact penalty functions

$$\min f(\mathbf{x}) \text{ s. t. } g_i(x) \leq 0, h_j(x) = 0$$

$$\longrightarrow P(\mathbf{x}, c) = f(\mathbf{x}) + c \left(\sum_i g_i(x)_+ + \sum_j |h_j(\mathbf{x})| \right)$$

Nonsmooth problems are everywhere

Optimization: exact penalty functions

$$\min f(\mathbf{x}) \text{ s. t. } g_i(x) \leq 0, h_j(x) = 0$$

$$\longrightarrow P(\mathbf{x}, c) = f(\mathbf{x}) + c \left(\sum_i g_i(x)_+ + \sum_j |h_j(\mathbf{x})| \right)$$

Robust estimation:

$$\min_{\mathbf{x}} \|f(\mathbf{x}) - \mathbf{y}\|_p \quad f \text{ nonlinear}$$

Nonsmooth problems are everywhere

Optimization: exact penalty functions

$$\min f(\mathbf{x}) \text{ s. t. } g_i(x) \leq 0, h_j(x) = 0$$

$$\longrightarrow P(\mathbf{x}, c) = f(\mathbf{x}) + c \left(\sum_i g_i(x)_+ + \sum_j |h_j(\mathbf{x})| \right)$$

Robust estimation:

$$\min_{\mathbf{x}} \|f(\mathbf{x}) - \mathbf{y}\|_p \quad f \text{ nonlinear}$$

Promoting structures:

Sparse phase retrieval

Sparse principal component analysis (SPCA)

Sparse blind deconvolution

Neural networks with nonsmooth activations (e.g., ReLU)

Nonsmooth problems are everywhere

Optimization: exact penalty functions

$$\min f(\mathbf{x}) \text{ s. t. } g_i(x) \leq 0, h_j(x) = 0$$

$$\longrightarrow P(\mathbf{x}, c) = f(\mathbf{x}) + c \left(\sum_i g_i(x)_+ + \sum_j |h_j(\mathbf{x})| \right)$$

Robust estimation:

$$\min_{\mathbf{x}} \|f(\mathbf{x}) - \mathbf{y}\|_p \quad f \text{ nonlinear}$$

Promoting structures:

Sparse phase retrieval

Sparse principal component analysis (SPCA)

Sparse blind deconvolution

Neural networks with nonsmooth activations (e.g., ReLU)

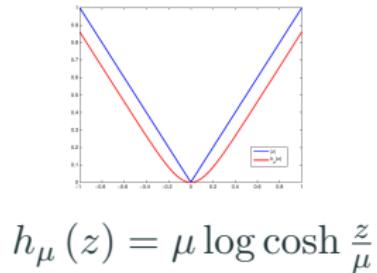
Others [Bagirov et al., 2014, Absil and Hosseini, 2017]

Smoothing?

$$\min \ f_1(\mathbf{q}) \doteq \frac{1}{m} \sum_{j=1}^m |\mathbf{q}^* \mathbf{y}_j| \text{ s.t. } \|\mathbf{q}\|_2^2 = 1.$$

$\downarrow\downarrow$ smoothing $\downarrow\downarrow$

$$\min \ f(\mathbf{q}) \doteq \frac{1}{m} \sum_{i=1}^m h_\mu (\mathbf{q}^* \mathbf{y}_i) \text{ s.t. } \|\mathbf{q}\|_2^2 = 1.$$



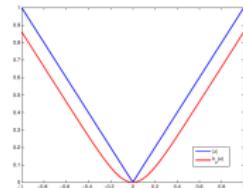
$$h_\mu(z) = \mu \log \cosh \frac{z}{\mu}$$

Smoothing?

$$\min f_1(\mathbf{q}) \doteq \frac{1}{m} \sum_{j=1}^m |\mathbf{q}^* \mathbf{y}_j| \text{ s.t. } \|\mathbf{q}\|_2^2 = 1.$$

$\downarrow\downarrow$ smoothing $\downarrow\downarrow$

$$\min f(\mathbf{q}) \doteq \frac{1}{m} \sum_{i=1}^m h_\mu(\mathbf{q}^* \mathbf{y}_i) \text{ s.t. } \|\mathbf{q}\|_2^2 = 1.$$



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

... at your own risk

Dear Ju Sun,

I regret to inform you that the editorial board has decided that your paper

"Complete dictionary recovery over the sphere"

By Ju Sun, Qing Qu, and John Wright

to Foundations of Computational Mathematics.

which you recently submitted to Foundations of Computational Mathematics cannot be accepted for publication in the journal, based on the statement of the handling editor which I attach below.

The paper will not be sent out for further refereeing.

Sincerely yours,

Handling editor statement: The paper is unusually long (more than 100 pages) for JoFoCM which very rarely publish papers longer than 40-50 pages. The results are strong and there are a number of useful ideas in the paper for further research. I have no doubt that it would be accepted (modulo correctness of the proof, which I did not check in detail) by FOCM if it were of regular length. I do not see a good way to reduce it to FOCM-length without making the paper hard to read. It is long, but well written. That being said, I do not feel it is groundbreaking to overwrite the length restriction, more likely not.

Language for nonsmooth functions?

What about nonsmooth, nonconvex problems?

1st order geometry: ∇f or $v^\top \nabla f$ (directional derivatives)

$\implies ?$

2nd order geometry: $\nabla^2 f$ or $v^\top \nabla^2 f v$ (directional curvatures)

$\implies ?$

Locally Lipschitz (continuous) functions

... functions that are Lipschitz **locally**:

- Continuous convex and concave functions
- Continuously differentiable functions
- Distance function to a set
- Sum of two locally Lipschitz functions: e.g., weakly convex functions ($f(\mathbf{x})$ so that $f(\mathbf{x}) + \rho \|\mathbf{x}\|_2^2$ is convex)
- Products/Quotients of two locally Lipschitz functions
- Compositions of two locally Lipschitz functions: e.g., $h(g(\mathbf{x}))$ with h convex and $g \in \mathcal{C}^1$
- ...

Clarke subdifferentials

We restrict to **finite-dimensional** functions, i.e.,

$f : X \mapsto \mathbb{R}$ with $X \subset \mathbb{R}^n$.

Clarke subdifferentials

We restrict to **finite-dimensional** functions, i.e.,

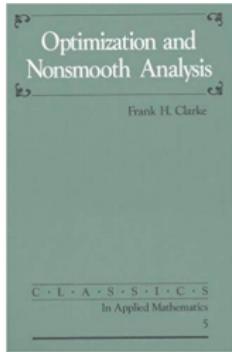
$f : X \mapsto \mathbb{R}$ with $X \subset \mathbb{R}^n$.

Rademacher's theorem: If f is locally Lipschitz, f is differentiable almost everywhere in X .

Clarke subdifferentials

We restrict to **finite-dimensional** functions, i.e.,
 $f : X \mapsto \mathbb{R}$ with $X \subset \mathbb{R}^n$.

Rademacher's theorem: If f is locally Lipschitz, f is differentiable almost everywhere in X .



Definition (Clarke subdifferential [Clarke, 1990])

$$\partial f(x) \doteq \text{conv} \{v : x_k \rightarrow x, \nabla f(x_k) \rightarrow v, f \text{ diff. at } x_k\}$$

... due to **Frank H. Clarke**. Well known in optimal control and economics

Clarke's subdifferential

For locally Lipschitz functions

Definition (Clarke subdifferential [Clarke, 1990])

$$\partial f(\mathbf{x}) \doteq \text{conv} \{ \mathbf{v} : \mathbf{x}_k \rightarrow \mathbf{x}, \nabla f(\mathbf{x}_k) \rightarrow \mathbf{v}, f \text{ diff. at } \mathbf{x}_k \}$$

- $\partial f(\mathbf{x})$ is always nonempty, convex, and compact
- $f \in \mathcal{C}^1$, $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$
- f convex: the usual subdifferential in convex analysis
- Most natural calculus rules hold (under **regularity** conditions, Chap 2 of [Clarke, 1990])
- Optimality: \mathbf{x}_0 is local min $\implies \mathbf{0} \in \partial f(\mathbf{x}_0)$

Language for nonsmooth functions?

What about nonsmooth, nonconvex problems?

Language for nonsmooth functions?

What about nonsmooth, nonconvex problems?

1st order geometry: ∇f or $v^\top \nabla f$ (directional derivatives)
 $\implies \partial f$ or $v^\top \partial f$

Language for nonsmooth functions?

What about nonsmooth, nonconvex problems?

1st order geometry: ∇f or $v^\top \nabla f$ (directional derivatives)
 $\implies \partial f$ or $v^\top \partial f$

2nd order geometry: $\nabla^2 f$ or $v^\top \nabla^2 f v$ (directional curvatures)
 \implies monotonicity of ∂f : f is convex **iff**

$$\langle u_x - u_y, x - y \rangle \geq 0 \quad \forall x, y \text{ and } \forall u_x \in \partial f(x), u_y \in \partial f(y).$$

Back to sparsifying transform learning

Given \mathbf{Y} , learn **ortho** \mathbf{Q} s.t. $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small.

$$\min \quad f(\mathbf{q}) \doteq \frac{1}{m} \|\mathbf{q}^*\mathbf{Y}\|_1 = \frac{1}{m} \sum_i |\mathbf{q}^*\mathbf{y}_i| \quad \text{s.t. } \|\mathbf{q}\|_2^2 = 1.$$

Back to sparsifying transform learning

Given \mathbf{Y} , learn **ortho** \mathbf{Q} s.t. $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small.

$$\min \quad f(\mathbf{q}) \doteq \frac{1}{m} \|\mathbf{q}^*\mathbf{Y}\|_1 = \frac{1}{m} \sum_i |\mathbf{q}^*\mathbf{y}_i| \quad \text{s.t. } \|\mathbf{q}\|_2^2 = 1.$$

Riemannian language: $\partial_R f(\mathbf{q}) = (\mathbf{I} - \mathbf{q}\mathbf{q}^*) \partial f(\mathbf{q})$

[Hosseini and Uschmajew, 2017]

Back to sparsifying transform learning

Given \mathbf{Y} , learn **ortho** \mathbf{Q} s.t. $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small.

$$\min \quad f(\mathbf{q}) \doteq \frac{1}{m} \|\mathbf{q}^*\mathbf{Y}\|_1 = \frac{1}{m} \sum_i |\mathbf{q}^*\mathbf{y}_i| \quad \text{s.t. } \|\mathbf{q}\|_2^2 = 1.$$

Riemannian language: $\partial_R f(\mathbf{q}) = (\mathbf{I} - \mathbf{q}\mathbf{q}^*) \partial f(\mathbf{q})$

[Hosseini and Uschmajew, 2017]

For analysis: Bernoulli-Gaussian model $\mathbf{X}_0 = \boldsymbol{\Omega}_0 \circ \mathbf{V}_0$,
 $\boldsymbol{\Omega}_0 \sim_{\text{iid}} \text{Ber}(\theta)$, $\mathbf{V}_0 \sim_{\text{iid}} \mathcal{N}(0, 1)$. **Sparsity parameter** θ

Back to sparsifying transform learning

Given \mathbf{Y} , learn **ortho** \mathbf{Q} s.t. $\mathbf{Q}^*\mathbf{Y}$ is sparse, i.e., $\|\mathbf{Q}^*\mathbf{Y}\|_0$ is small.

$$\min \quad f(\mathbf{q}) \doteq \frac{1}{m} \|\mathbf{q}^*\mathbf{Y}\|_1 = \frac{1}{m} \sum_i |\mathbf{q}^*\mathbf{y}_i| \quad \text{s.t. } \|\mathbf{q}\|_2^2 = 1.$$

Riemannian language: $\partial_R f(\mathbf{q}) = (\mathbf{I} - \mathbf{q}\mathbf{q}^*) \partial f(\mathbf{q})$

[Hosseini and Uschmajew, 2017]

For analysis: Bernoulli-Gaussian model $\mathbf{X}_0 = \boldsymbol{\Omega}_0 \circ \mathbf{V}_0$,
 $\boldsymbol{\Omega}_0 \sim_{\text{iid}} \text{Ber}(\theta)$, $\mathbf{V}_0 \sim_{\text{iid}} \mathcal{N}(0, 1)$. **Sparsity parameter** θ

When m is large, w.h.p., in a “reasonably large” region of e_n :

$$\inf \langle \partial_R f(\mathbf{q}), \mathbf{q} - e_n \rangle \geq \gamma \|\mathbf{q} - e_n\|$$

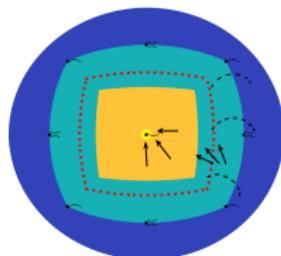


image credit: [Gilboa et al., 2018]

Subgradient descent learns orthogonal dictionaries!

Starting from a $\mathbf{q}^{(0)}$ **uniformly random** on \mathbb{S}^{n-1} , for
 $k = 0, 1, 2, \dots$:

$$\mathbf{q}^{(k+1)} = \frac{\mathbf{q}^{(k)} - \eta^{(k)} \mathbf{v}^{(k)}}{\|\mathbf{q}^{(k)} - \eta^{(k)} \mathbf{v}^{(k)}\|} \quad \text{for any } \mathbf{v} \in \partial_R f(\mathbf{q}^{(k)})$$

Subgradient descent learns orthogonal dictionaries!

Starting from a $\mathbf{q}^{(0)}$ **uniformly random** on \mathbb{S}^{n-1} , for
 $k = 0, 1, 2, \dots$:

$$\mathbf{q}^{(k+1)} = \frac{\mathbf{q}^{(k)} - \eta^{(k)} \mathbf{v}^{(k)}}{\|\mathbf{q}^{(k)} - \eta^{(k)} \mathbf{v}^{(k)}\|} \quad \text{for any } \mathbf{v} \in \partial_R f(\mathbf{q}^{(k)})$$

Ideas:

- Each run finds an e_i with constant probability
- All basis vectors found in $O(n \log n)$ independent runs

Subgradient descent learns orthogonal dictionaries!

Starting from a $\mathbf{q}^{(0)}$ **uniformly random** on \mathbb{S}^{n-1} , for
 $k = 0, 1, 2, \dots$:

$$\mathbf{q}^{(k+1)} = \frac{\mathbf{q}^{(k)} - \eta^{(k)} \mathbf{v}^{(k)}}{\|\mathbf{q}^{(k)} - \eta^{(k)} \mathbf{v}^{(k)}\|} \quad \text{for any } \mathbf{v} \in \partial_R f(\mathbf{q}^{(k)})$$

Ideas:

- Each run finds an e_i with constant probability
- All basis vectors found in $O(n \log n)$ independent runs

Theorem (Informal)

Assume $\theta \in [1/n, 1/2]$. When $m \geq \Omega(\theta^{-2} n^4 \log^3 n)$, whp, the proposed algorithm recovers all basis vectors in polynomial time.

Comparison with the DL literature

Algorithms working in the constant sparsity regime, i.e., $\theta \in \Theta(1)$

- **Convex relaxation based on Sum-of-Squares (SOS):**

[Barak et al., 2015, Ma et al., 2016b, Schramm and Steurer, 2017]

solving huge SDP's or tensor decompositions

Comparison with the DL literature

Algorithms working in the constant sparsity regime, i.e., $\theta \in \Theta(1)$

- **Convex relaxation based on Sum-of-Squares (SOS):**
[Barak et al., 2015, Ma et al., 2016b, Schramm and Steurer, 2017]
solving huge SDP's or tensor decompositions
- **Nonconvex relaxation based on smoothed ℓ_1 :** 2nd order
method [Sun et al., 2015] or 1st order method [Gilboa et al., 2018], still
expensive in computation and involved for analysis

Comparison with the DL literature

Algorithms working in the constant sparsity regime, i.e., $\theta \in \Theta(1)$

- **Convex relaxation based on Sum-of-Squares (SOS):**
[Barak et al., 2015, Ma et al., 2016b, Schramm and Steurer, 2017]
solving huge SDP's or tensor decompositions
- **Nonconvex relaxation based on smoothed ℓ_1 :** 2nd order
method [Sun et al., 2015] or 1st order method [Gilboa et al., 2018], still
expensive in computation and involved for analysis
- **This work: nonconvex relaxation based directly on ℓ_1 :**
lightweight computation and neater analysis — compress the smoothed ℓ_1
analysis by 1/2!

A word on technicalities

Subdifferentials are (convex) sets in general, and randomness in data leads to **random sets**.

- Measure set difference: Hausdorff distance
- Expectation of random sets: selection integrals and support functions [Aubin and Frankowska, 2009, Molchanov, 2013]
- Concentration of Minkowski sum of random sets: support functions and concentration of empirical processes
[Molchanov, 2017, Molchanov, 2013]

A word on technicalities

Subdifferentials are (convex) sets in general, and randomness in data leads to **random sets**.

- Measure set difference: Hausdorff distance
- Expectation of random sets: selection integrals and support functions [Aubin and Frankowska, 2009, Molchanov, 2013]
- Concentration of Minkowski sum of random sets: support functions and concentration of empirical processes
[Molchanov, 2017, Molchanov, 2013]

The $\text{sign}(\cdot)$ function is not Lipschitz in the usual sense

- Careful construction of the ε -net for covering in showing uniform convergence of the subdifferential

Easier ways around nonsmoothness?

Ignore non-differentiable points

- Pathological examples well known
- Performance on “generic” cases not understood

Easier ways around nonsmoothness?

Ignore non-differentiable points

- Pathological examples well known
- Performance on “generic” cases not understood

Smooth out and continue

- Relatively mature for convex problems [Nesterov, 2004]
- Lack in theory for nonconvex problems [Mabahi, 2013]

Easier ways around nonsmoothness?

Ignore non-differentiable points

- Pathological examples well known
- Performance on “generic” cases not understood

Smooth out and continue

- Relatively mature for convex problems [Nesterov, 2004]
- Lack in theory for nonconvex problems [Mobahi, 2013]

Tweak around subdifferential sets

- Intuitive chain rules [Kakade and Lee, 2018]
- Setting a predefined rule—might not be reliable in computation

So far ...

Tame nonconvexity = Live with and understand nonconvexity

So far ...

Tame nonconvexity = Live with and understand nonconvexity

Which nonconvex optimization problems are easy?

- A1: problems with nice **global** landscapes
- A2: problems with nice **local** landscapes

So far ...

Tame nonconvexity = Live with and understand nonconvexity

Which nonconvex optimization problems are easy?

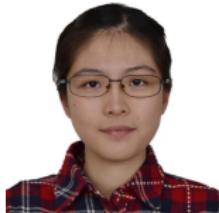
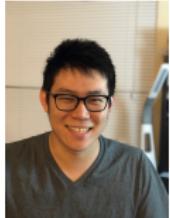
A1: problems with nice **global** landscapes

A2: problems with nice **local** landscapes

What about nonsmooth, nonconvex problems?

we're still picking up the right language...

Thanks to ...



John Wright

Columbia

Qing Qu

NYU

Yu Bai

Stanford

Qijia Jiang

Stanford

E. Candès

Stanford

Block-Reference Coherent Diffraction Imaging, Barmherzig, S., Lane, and Li, '18.

Landscape Analysis of Nonsmooth Functions, S. and Candès, '18.

Subgradient descent learns orthogonal dictionaries, Bai, Jiang, S. and Candès, '18.

Dictionary Learning in Fourier Transform Scanning Tunneling Spectroscopy, Cheung, Shin, Lau, Chen, S., Zhang, Wright, and Pasupathy, '18.

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

My webpage on provable nonconvex heuristics: <http://sunju.org/research/nonconvex/>

Thank you!

References i

- [Absil and Hosseini, 2017] Absil, P. and Hosseini, S. (2017). **A collection of nonsmooth riemannian optimization problems.**
- [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*.
- [Anandkumar et al., 2014a] Anandkumar, A., Ge, R., and Janzamin, M. (2014a). **Analyzing tensor power method dynamics: Applications to learning overcomplete latent variable models.** *arXiv preprint arXiv:1411.1488*.
- [Anandkumar et al., 2014b] Anandkumar, A., Ge, R., and Janzamin, M. (2014b). **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*.

References ii

- [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*.
- [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*.
- [Aubin and Frankowska, 2009] Aubin, J.-P. and Frankowska, H. (2009). **Set-Valued Analysis.** Modern Birkhäuser Classics. Birkhäuser Basel.
- [Bagirov et al., 2014] Bagirov, A., Karmitsa, N., and Mäkelä, M. M. (2014). **Introduction to Nonsmooth Optimization: theory, practice and software.** Springer.
- [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*.
- [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*.

References iii

- [Barak et al., 2015] Barak, B., Kelner, J. A., and Steurer, D. (2015). **Dictionary learning and tensor decomposition via the sum-of-squares method.** In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 143–151. ACM.
- [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*.
- [Candès et al., 2015] Candès, E. J., Li, X., and Soltanolkotabi, M. (2015). **Phase retrieval via wirtinger flow: Theory and algorithms.** *Information Theory, IEEE Transactions on*, 61(4):1985–2007.
- [Chatterji and Bartlett, 2017] Chatterji, N. S. and Bartlett, P. L. (2017). **Alternating minimization for dictionary learning with random initialization.** *arxiv:1711.03634*.
- [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 et al., 2018] Chen, Y., Chi, Y., Fan, J., and Ma, C. (2018). **Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval.** *arXiv:1803.07726*.
- [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*.
- [Clarke, 1990] Clarke, F. H. (1990). **Optimization and nonsmooth analysis, volume 5.** Siam.
- [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.
- [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.
- [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*.

References v

- [Geng and Wright, 2011] Geng, Q. and Wright, J. (2011). **On the local correctness of ℓ^1 -minimization for dictionary learning.** Submitted to *IEEE Transactions on Information Theory*. Preprint: <http://www.columbia.edu/~jw2966>.
- [Gilboa et al., 2018] Gilboa, D., Buchanan, S., and Wright, J. (2018). **Efficient dictionary learning with gradient descent.** *arXiv:1809.10313*.
- [Goldfarb, 1980] Goldfarb, D. (1980). **Curvilinear path steplength algorithms for minimization which use directions of negative curvature.** *Mathematical programming*, 18(1):31–40.
- [Gribonval and Schnass, 2010] Gribonval, R. and Schnass, K. (2010). **Dictionary identification - sparse matrix-factorization via ℓ^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.

- [Hosseini and Uschmajew, 2017] Hosseini, S. and Uschmajew, A. (2017). **A Riemannian gradient sampling algorithm for nonsmooth optimization on manifolds.** *SIAM Journal on Optimization*, 27(1):173–189.
- [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*.
- [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*.
- [Kakade and Lee, 2018] Kakade, S. and Lee, J. D. (2018). **Provably correct automatic subdifferentiation for qualified programs.** *arXiv:1809.08530*.

- [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.
- [Lu and Kawaguchi, 2017] Lu, H. and Kawaguchi, K. (2017). **Depth creates no bad local minima.** *arXiv preprint arXiv:1702.08580*.
- [Ma et al., 2016a] Ma, T., Shi, J., and Steurer, D. (2016a). **Polynomial-time tensor decompositions with sum-of-squares.** In *Foundations of Computer Science (FOCS), 2016 IEEE 57th Annual Symposium on*, pages 438–446. IEEE.
- [Ma et al., 2016b] Ma, T., Shi, J., and Steurer, D. (2016b). **Polynomial-time tensor decompositions with sum-of-squares.**
- [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.
- [Mobahi, 2013] Mobahi, H. (2013). **Optimization by Gaussian smoothing with application to geometric alignment.** PhD thesis, University of Illinois at Urbana-Champaign.

- [Molchanov, 2013] Molchanov, I. (2013). **Foundations of stochastic geometry and theory of random sets.** In *Stochastic Geometry, Spatial Statistics and Random Fields*, pages 1–20. Springer.
- [Molchanov, 2017] Molchanov, I. (2017). **Theory of random sets.** Springer-Verlag London, 2nd edition.
- [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.
- [Nesterov, 2004] Nesterov, Y. (2004). **Smooth minimization of non-smooth functions.** *Mathematical Programming*, 103(1):127–152.
- [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.
- [Ravishankar and Wohlberg, 2018] Ravishankar, S. and Wohlberg, B. (2018). **Learning multi-layer transform models**. *arXiv:1810.08323*.
- [Rockafellar, 1993] Rockafellar, R. T. (1993). **Lagrange multipliers and optimality**. *SIAM review*, 35(2):183–238.
- [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*.
- [Schramm and Steurer, 2017] Schramm, T. and Steurer, D. (2017). **Fast and robust tensor decomposition with applications to dictionary learning**. *arXiv:1706.08672*.
- [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*.

References x

- [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*.
- [Wang et al., 2016] Wang, G., Giannakis, G. B., and Eldar, Y. C. (2016). **Solving systems of random quadratic equations via truncated amplitude flow.** *arxiv:1605.08285*.
- [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.