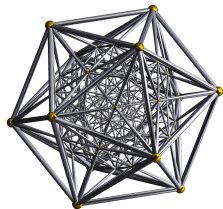


# ICASSP 2022 Short Course One on Low-Dimensional Models for High-Dimensional Data

## Lecture 3: Learning Low-dimensional Models via Nonconvex Optimization

**Sam Buchanan, Yi Ma, Qing Qu  
John Wright, Yuqian Zhang, Zhihui Zhu**

May 25, 2022



# Outline

## ① Introduction & Motivation of Nonconvex Optimization

Motivating Examples

Nonlinearity, Nonconvexity, and Symmetry

## ② Symmetry & Geometry for Nonconvex Problems in Practice

Problems with Rotational Symmetry

Problems with Discrete Symmetry

## ③ Efficient Nonconvex Optimization

Objectives of Nonconvex Optimization

Escaping Saddles

# Example: Low-rank Matrix Completion

We observe:

$$\mathbf{Y} = \mathcal{P}_\Omega \left[ \mathbf{X} \right].$$

Observed ratings = Complete ratings

Users

Items

Observed (Incomplete) Ratings  $\mathbf{Y}$

Complete Ratings  $\mathbf{X}$

## Matrix completion

via **bilinear** low-rank factorization

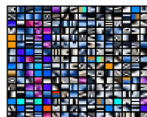
$$\min_{\mathbf{U}, \mathbf{V}} f(\mathbf{U}, \mathbf{V}) = \sum_{(i,j) \in \Omega} [(\mathbf{UV}^*)_{i,j} - \mathbf{Y}_{i,j}]^2 + \underbrace{\frac{\lambda}{2} \|\mathbf{U}\|_F^2 + \frac{\lambda}{2} \|\mathbf{V}\|_F^2}_{\text{reg}(\mathbf{U}, \mathbf{V})}$$

$$\|\mathbf{M}\|_* = \min_{\mathbf{M}=\mathbf{UV}^*} \frac{\lambda}{2} \|\mathbf{U}\|_F^2 + \frac{\lambda}{2} \|\mathbf{V}\|_F^2$$

# Example: Dictionary for Image Representation

Image processing  
(e.g. denoising or super-resolution)  
against a known sparsifying dictionary:

$$I_{\text{noisy}} = \underset{\text{dictionary}}{\mathbf{A}} \times \underset{\text{sparse}}{\mathbf{x}} + \underset{\text{noise}}{\mathbf{z}}. \quad (1)$$

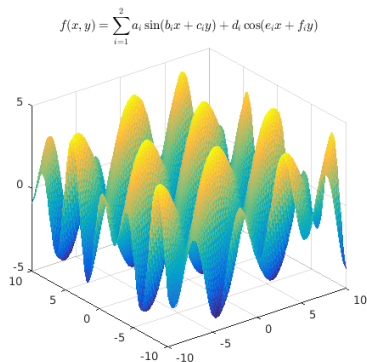
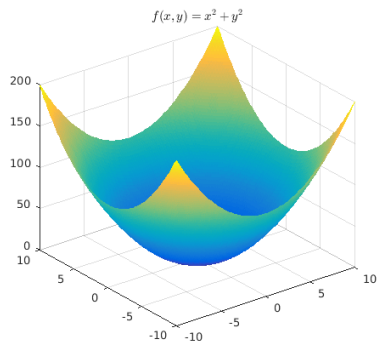


**Dictionary learning:** the motifs or atoms of the dictionary are **unknown**:

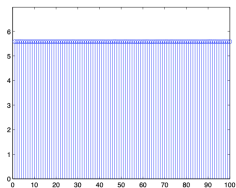
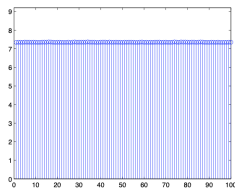
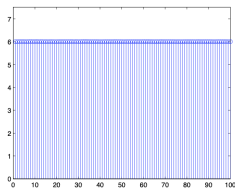
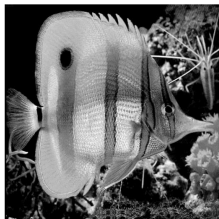
$$\underset{\text{data}}{\mathbf{Y}} = \underset{\text{dictionary}}{\mathbf{A}} \underset{\text{sparse}}{\mathbf{X}}. \quad (2)$$

- Band-limited signals:  $\mathbf{A} = \mathbf{F}$ , the Fourier transform;
- Piecewise smooth signals:  $\mathbf{A} = \mathbf{W}$ , the wavelet transforms;
- Natural images  $\mathbf{A} = ?$  (How to **learn**  $\mathbf{A}$  from the data  $\mathbf{Y}$ ?)

# Convex and Nonconvex Optimization

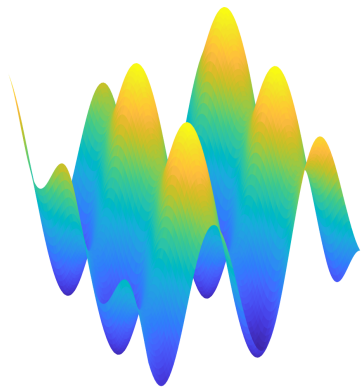


# Dictionary Learning

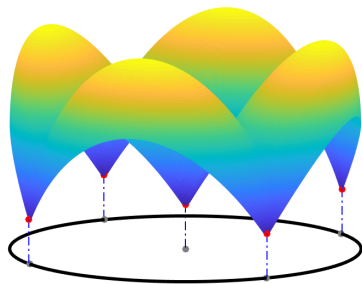


**Recovered solutions always obtain the same objective value.**

# Benign Nonconvex Optimization Landscape



**General Case**



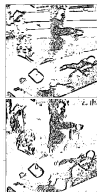
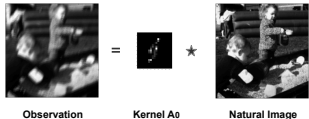
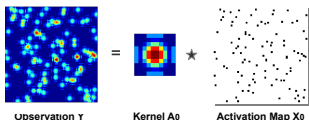
**Structured Case**

# Example: Sparse Blind Deconvolution

**Sparse Blind Deconvolution:**  
the convolutional motif or sparse  
activation signal are **unknown**:

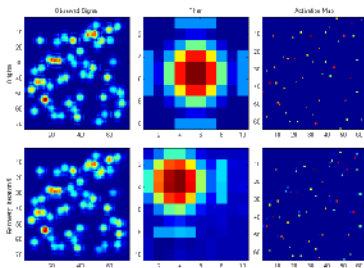
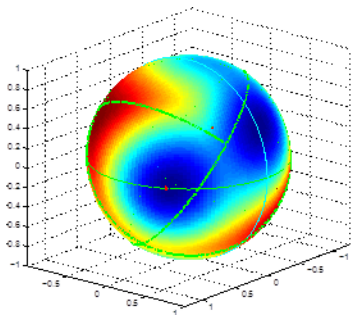
$$\underset{\text{data}}{\mathbf{Y}} = \underset{\text{motif}}{\mathbf{A}} * \underset{\text{sparse}}{\mathbf{X}}. \quad (3)$$

- Scientific signals:  
activation signals are sparse
- Image deblurring:  
natural images are  
sparse in the gradient domain





# Sparse Blind Deconvolution

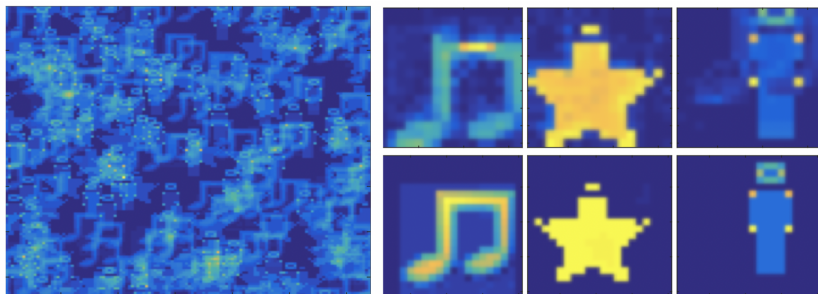


**Recovered solutions are near signed shift-truncations of the ground truth.**

# Convolutional Dictionary learning

$$\mathbf{Y} = \sum_i \mathbf{A}_i * \mathbf{X}_i.$$

data      motif      sparse



Recovered solutions are near signed shift-truncations of the ground truth.

# Challenges of Nonconvex Optimization – Pessimistic Views

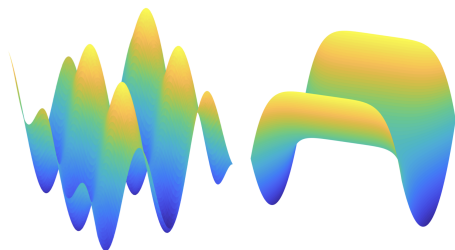
Consider the problem of minimizing a general nonlinear function:

$$\min_z \varphi(z), \quad z \in \mathbb{C}. \quad (4)$$

In **the worst case**, even finding a *local* minimizer can be NP-hard<sup>1</sup>.

Hence typically people seek to work with relatively benign functions with benign guarantees:

- ① convergence to some critical point  $\bar{z}$  such that  $\nabla\varphi(\bar{z}) = \mathbf{0}$ ;
- ② or convergence to some local minimizer  $\nabla^2\varphi(\bar{z}) \succeq \mathbf{0}$ .



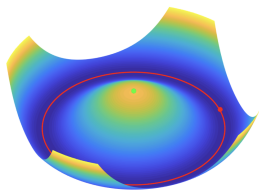
Spurious local minimizers

Flat saddle points

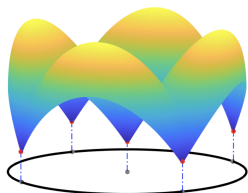
<sup>1</sup>Some NP-complete problems in quadratic and nonlinear programming, K.G Murty and S. N. Kabadi, 1987

## Opportunities – Optimistic Views

However, nonconvex problems that arise from natural physical, geometrical, or statistical origins typically have **nice** structures, in terms of **symmetries!**



Rotational symmetry



Discrete symmetry

The function  $\varphi$  is **invariant** under certain group action:

- for low rank matrix recovery, invariant under a continuous rotation:

$$\varphi((\mathbf{U}\mathbf{\Gamma}, \mathbf{V}\mathbf{\Gamma}^{-1})) = \varphi((\mathbf{U}, \mathbf{V})), \quad \forall \text{invertible } \mathbf{\Gamma}.$$

- for dictionary learning, invariant under signed permutations:

$$\varphi((\mathbf{A}, \mathbf{X})) = \varphi((\mathbf{A}\mathbf{\Pi}, \mathbf{\Pi}^* \mathbf{X})), \quad \forall \mathbf{\Pi} \in \text{SP}(n).$$

## Nonlinearity and Symmetry

Intrinsic ambiguity against the uniqueness of the solution

- **low rank matrix recovery**

$$\mathbf{X} = \mathbf{U}_0 \mathbf{V}_0^T = \mathbf{U}_0 \mathbf{\Gamma} \mathbf{\Gamma}^{-1} \mathbf{V}_0^T$$

for any invertible  $\mathbf{\Gamma}$ .

- **dictionary learning**

$$\mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0 = \mathbf{A}_0 \mathbf{\Pi} \mathbf{\Pi}^* \mathbf{X}_0$$

for any signed permutation  $\mathbf{\Pi}$ .

- **blind deconvolution**

$$\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0 = S_\tau[\mathbf{a}_0] * S_{-\tau}[\mathbf{x}_0]$$

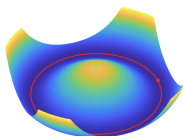
for any signed shift  $\tau$ .

# Optimization under Symmetry

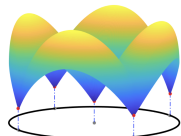
## Definition (Symmetric Function)

Let  $\mathbb{G}$  be a group acting on  $\mathbb{R}^n$ . A function  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^{n'}$  is  $\mathbb{G}$ -symmetric if for all  $\mathbf{z} \in \mathbb{R}^n$ ,  $\mathbf{g} \in \mathbb{G}$ ,  $\varphi(\mathbf{g} \circ \mathbf{z}) = \varphi(\mathbf{z})$ .

Most symmetric objective functions that arise in structure signal recovery **do not** have spurious local minimizers or flat saddles.



Rotational symmetry



Discrete symmetry

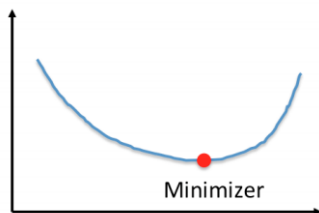
**Slogan 1:** the (only!) local minimizers are symmetric versions of the ground truth.

**Slogan 2:** any local critical point has negative curvature in directions that break symmetry.

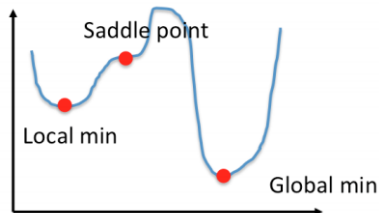
# Basic Calculus

Critical points or stationary points: gradient vanishes

**Convex**



**Non-Convex**

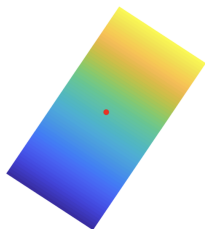


- convex function: critical point = minimizer
- nonconvex function: not all critical point are minimizer

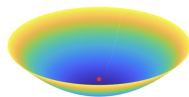
# Basic Calculus

Critical points with non-singular hessian

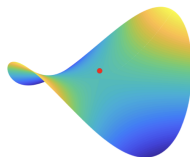
- minimizer: hessian is positive definite
- saddle points: hessian has both positive and negative eigenvalues
- maximizer: hessian is negative definite



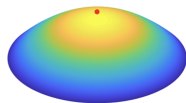
**Noncritical Point** ( $\nabla\varphi \neq \mathbf{0}$ )



**Minimizer**  
 $\nabla^2\varphi > \mathbf{0}$



**Saddle**  
 $\lambda_{\min} \nabla^2\varphi < 0$   
 $\lambda_{\max} \nabla^2\varphi > 0$



**Maximizer**  
 $\nabla^2\varphi < \mathbf{0}$

**Critical Points** ( $\nabla\varphi = \mathbf{0}$ )



# Outline

## ① Introduction & Motivation of Nonconvex Optimization

Motivating Examples

Nonlinearity, Nonconvexity, and Symmetry

## ② Symmetry & Geometry for Nonconvex Problems in Practice

Problems with Rotational Symmetry

Problems with Discrete Symmetry

## ③ Efficient Nonconvex Optimization

Objectives of Nonconvex Optimization

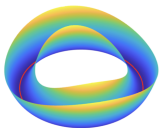
Escaping Saddles

# Problems with Rotational Symmetry

## Nonconvex Problems with Rotational Symmetries

### Eigenspace Computation

Compute the principal subspace of a symmetric matrix.

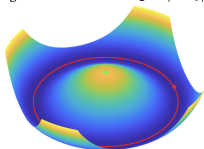


$$\min_{\mathbf{X}^* \mathbf{X} = \mathbf{I}} -\frac{1}{2} \text{trace}[\mathbf{X}^* \mathbf{A} \mathbf{X}].$$

**Symmetry:**  $\mathbf{X} \mapsto \mathbf{X} \mathbf{R}$   
 $\mathbb{G} = O(r)$

### Generalized Phase Retrieval

Recover a complex vector  $\mathbf{x}_0$  from magnitude measurements  $\mathbf{y} = |\mathbf{A} \mathbf{x}_0|$ .

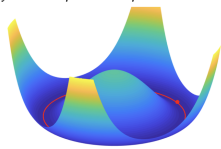


$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y}^2 - |\mathbf{A} \mathbf{x}|^2\|_2^2.$$

**Symmetry:**  $\mathbf{x} \mapsto \mathbf{x} e^{i\phi}$   
 $\mathbb{G} = \mathbb{S}^1 \cong O(2)$

### Matrix Recovery

Recover a low-rank matrix  $\mathbf{X} = \mathbf{U} \mathbf{V}^*$  from incomplete / corrupted observations



$$\min_{\mathbf{U}, \mathbf{V}} \mathcal{L}(\mathbf{Y} - \mathcal{A}[\mathbf{U} \mathbf{V}^*]) + \rho(\mathbf{U}, \mathbf{V}).$$

**Symmetry:**  $(\mathbf{U}, \mathbf{V}) \mapsto (\mathbf{U} \mathbf{\Gamma}, \mathbf{V} \mathbf{\Gamma}^{-*})$   
 $\mathbb{G} = \text{GL}(r)$  or  $\mathbb{G} = O(r)$

# Low rank matrix recovery

Goal: Given  $\mathbf{Y} = \mathcal{A}(\mathbf{X})$ , recover low rank matrix  $\mathbf{X} = \mathbf{U}_0\mathbf{V}_0$

$$\begin{array}{c} \text{Users} \\ \begin{matrix} \text{😊} \\ \text{😊} \\ \vdots \\ \text{😊} \end{matrix} \end{array} \begin{bmatrix} 5 & 3 & \dots & ? \\ ? & 2 & \dots & 4 \\ \vdots & \vdots & \ddots & \vdots \\ 5 & ? & \dots & ? \end{bmatrix} = \mathcal{P}_\Omega \left( \begin{array}{c} \begin{bmatrix} 5 & 3 & \dots & 5 \\ 4 & 2 & \dots & 4 \\ \vdots & \vdots & \ddots & \vdots \\ 5 & 5 & \dots & 3 \end{bmatrix} \\ \text{Complete Ratings } \mathbf{X} \end{array} \right)$$

Items  
Observed (Incomplete) Ratings  $\mathbf{Y}$

- Convex Formulation

$$\min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \|\mathbf{X}\|_\star \quad \text{s.t.} \quad \mathbf{Y} = \mathcal{A}(\mathbf{X})$$

- Nonconvex Formulation

$$\min_{\mathbf{U} \in \mathbb{R}^{m \times r}, \mathbf{V} \in \mathbb{R}^{n \times r}} \|\mathbf{Y} - \mathcal{A}(\mathbf{UV}^T)\|_F^2 + \text{reg}(\mathbf{U}, \mathbf{V})$$

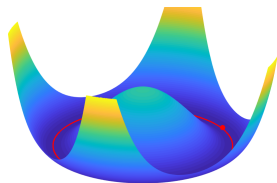
# Low Rank Matrix Recovery

$$\min_{\mathbf{U}, \mathbf{V}} \quad \frac{1}{2} \|\mathbf{Y} - \mathcal{A}(\mathbf{UV}^T)\|_F^2 + \text{reg}(\mathbf{U}, \mathbf{V})$$

**Inherent Symmetry:**

$$\mathbf{X} = \mathbf{U}_0 \mathbf{V}_0^T = \mathbf{U}_0 \mathbf{\Gamma} \mathbf{\Gamma}^{-1} \mathbf{V}_0^T$$

for any invertible  $\mathbf{\Gamma} \in \mathbb{R}^{r \times r}$ .



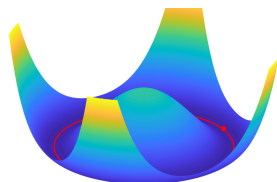
# Low Rank Matrix Recovery

$$\min_{\mathbf{U}, \mathbf{V}} \frac{1}{2} \|\mathbf{Y} - \mathcal{A}(\mathbf{UV}^T)\|_F^2 + \text{reg}(\mathbf{U}, \mathbf{V})$$

**Inherent Symmetry:**

$$\mathbf{X} = \mathbf{U}_0 \mathbf{V}_0^T = \mathbf{U}_0 \mathbf{\Gamma} \mathbf{\Gamma}^{-1} \mathbf{V}_0^T$$

for any invertible  $\mathbf{\Gamma} \in \mathbb{R}^{r \times r}$ .



- Are  $(\mathbf{U}_0 \mathbf{\Gamma}, \mathbf{V}_0 \mathbf{\Gamma}^{-1})$  the only local solutions?
- Does there exist flat stationary points?

# Rank-1 Symmetric Matrix

Simplifications:

- $\mathbf{Y} = \mathcal{A}(\mathbf{X}) = \mathbf{X}$
- $\mathbf{X} = \mathbf{U}_0 \mathbf{U}_0^T$  is symmetric and rank-1

$$\mathbf{X} = \mathbf{u}_0 \mathbf{u}_0^T = (-\mathbf{u}_0)(-\mathbf{u}_0^T)$$

the rotational symmetry is reduced to sign symmetry.

Nonconvex formulation:

$$\min_{\mathbf{u}} \phi(\mathbf{u}) \doteq \frac{1}{4} \|\mathbf{X} - \mathbf{u}\mathbf{u}^T\|_F^2 + \underbrace{\lambda \|\mathbf{u}\|_2^2}_{const}$$

# Rank-1 Symmetric Matrix

$$\min_{\mathbf{u}} \phi(\mathbf{u}) \doteq \frac{1}{4} \|\mathbf{X} - \mathbf{u}\mathbf{u}^T\|_F^2$$

Critical points have zero gradient

$$\begin{aligned}\nabla\phi &= (\mathbf{u}\mathbf{u}^T - \mathbf{X})\mathbf{u} \\ &= \|\mathbf{u}\|_2^2 \mathbf{u} - \mathbf{X}\mathbf{u} \\ &= \mathbf{0}\end{aligned}$$

therefore critical points must be one of the following

- $\mathbf{u} = \pm\mathbf{u}_0$
- $\mathbf{u} = \mathbf{0}$

# Rank-1 Symmetric Matrix

$$\min_{\mathbf{u}} \phi(\mathbf{u}) \doteq \frac{1}{4} \|\mathbf{X} - \mathbf{u}\mathbf{u}^T\|_F^2$$

with the second order derivative

$$\nabla^2 \phi = 2\mathbf{u}\mathbf{u}^T + \|\mathbf{u}\|_2^2 \mathbf{I} - \mathbf{X}.$$



# Rank-1 Symmetric Matrix

$$\min_{\mathbf{u}} \phi(\mathbf{u}) \doteq \frac{1}{4} \|\mathbf{X} - \mathbf{u}\mathbf{u}^T\|_F^2$$

with the second order derivative

$$\nabla^2 \phi = 2\mathbf{u}\mathbf{u}^T + \|\mathbf{u}\|_2^2 \mathbf{I} - \mathbf{X}.$$

Then the critical points can be grouped as

- Local minimizer  $\mathbf{u} = \pm \mathbf{u}_0$  and  $\mathbf{u}\mathbf{u}^T = \mathbf{X}$

$$\nabla^2 \phi = \mathbf{u}\mathbf{u}^T + \|\mathbf{u}\|_2^2 \mathbf{I}.$$

- Maximizer  $\mathbf{u} = \mathbf{0}$

$$\nabla^2 \phi = -\mathbf{X}.$$

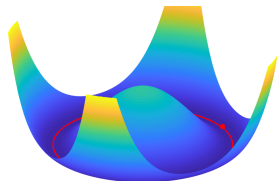
# Low Rank Matrix Recovery

Symmetric low rank matrix

$$\min_{\mathbf{U}} \phi(\mathbf{u}) \doteq \frac{1}{4} \|\mathbf{X} - \mathbf{U}\mathbf{U}^T\|_F^2.$$

General low rank matrix recover

$$\min_{\mathbf{U}, \mathbf{V}} \phi(\mathbf{u}) \doteq \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{V}^T\|_F^2 + \lambda \|\mathbf{U}\|_F^2 + \lambda \|\mathbf{V}\|_F^2.$$



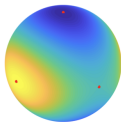
**Local minimizers:** *are ground truth  $\mathbf{U}_0$  and  $\mathbf{V}_0$  up to rotation;*  
**Negative curvature:** *between multiple local minimizers.*

# Problems with Discrete Symmetry

## Nonconvex Problems with Discrete Symmetries

### Eigenvector Computation

Maximize a quadratic form over the sphere.

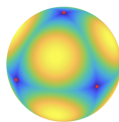


$$\max_{\mathbf{x} \in \mathbb{S}^{n-1}} \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x}.$$

Symmetry:  $\mathbf{x} \mapsto -\mathbf{x}$   
 $G = \{\pm 1\}$

### Dictionary Learning

Approximate a given matrix  $\mathbf{Y}$  as  $\mathbf{Y} \approx \mathbf{A}\mathbf{X}$ , with  $\mathbf{X}$  sparse

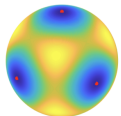


$$\min_{\mathbf{A} \in \mathcal{A}, \mathbf{X}} \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1.$$

Symmetry:  $(\mathbf{A}, \mathbf{X}) \mapsto (\mathbf{A}\Gamma, \mathbf{X}\Gamma^*)$   
 $G = \text{SP}(n)$

### Tensor Decomposition

Determine components  $\mathbf{a}_i$  of an orthogonal decomposable tensor  $\mathbf{T} = \sum_i \mathbf{a}_i \otimes \mathbf{a}_i \otimes \mathbf{a}_i$

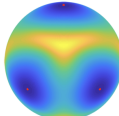


$$\max_{\mathbf{X} \in O(n)} \sum_i T(\mathbf{x}_i, \mathbf{x}_i, \mathbf{x}_i).$$

Symmetry:  $\mathbf{X} \mapsto \mathbf{X}\Gamma$   
 $G = \text{P}(n)$

### Short-and-Sparse Deconvolution

Recover a short  $\mathbf{a}$  and a sparse  $\mathbf{x}$  from their convolution  $\mathbf{y} = \mathbf{a} * \mathbf{x}$ .



$$\min_{\mathbf{a}, \mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{a} * \mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

Symmetry:  $(\mathbf{a}, \mathbf{x}) \mapsto (\alpha s_\tau[\mathbf{a}], \alpha^{-1} s_{-\tau}[\mathbf{x}])$   
 $G = \mathbb{Z}_n \times \mathbb{R}_* \text{ or } G = \mathbb{Z}_n \times \{\pm 1\}$

## Dictionary Learning

Goal: Given dataset  $\mathbf{Y}$ , find the optimal dictionary  $\mathbf{A}$  that renders the sparsest coefficient  $\mathbf{X}$

$$\min_{\mathbf{A}, \mathbf{X}} \|\mathbf{X}\|_1 \quad \text{s.t.} \quad \mathbf{Y} = \mathbf{A}\mathbf{X}.$$

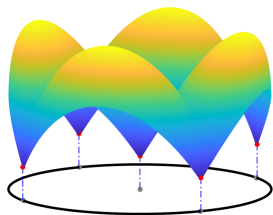
In presence of noise, the optimization problem can be rewritten as

$$\min_{\mathbf{A}, \mathbf{X}} \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1.$$

**Inherent Symmetry:**

$$\mathbf{Y} = \mathbf{A}_0 \mathbf{\Gamma} \mathbf{\Gamma}^* \mathbf{X}_0,$$

for any signed permutation matrix  $\mathbf{\Gamma}$ .



## Orthogonal Dictionary Learning

- Input: matrix  $\mathbf{Y}$  which is the product of an orthogonal matrix  $\mathbf{A}_0$  (called a dictionary) and a sparse matrix  $\mathbf{X}_0$ :

$$\mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0, \quad \mathbf{A}_0 \mathbf{A}_0^* = \mathbf{I}, \mathbf{X}_0 \text{ sparse.}$$

- Optimization Formulation

$$\min_{\mathbf{A}, \mathbf{X}} \|\mathbf{X}\|_1 \quad \text{s.t.} \quad \mathbf{Y} = \mathbf{A}\mathbf{X}, \quad \mathbf{A}\mathbf{A}^* = \mathbf{I}.$$

- Given the optimization constraint,  $\mathbf{X}$  is uniquely defined in terms of  $\mathbf{A}$

$$\mathbf{X} = \mathbf{A}^* \mathbf{A} \mathbf{X} = \mathbf{A}^* \mathbf{Y}.$$

- Equivalent formulation

$$\min_{\mathbf{A} \in \mathcal{O}(n)} \|\mathbf{A}^* \mathbf{Y}\|_1.$$

# Orthogonal Dictionary Learning

Instead of aiming to solve the entire matrix  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$  at once via

$$\min_{\mathbf{A} \in \mathcal{O}(n)} \|\mathbf{A}^* \mathbf{Y}\|_1.$$

A simpler model problem solves for the columns  $\mathbf{a}_i$  one at a time

$$\min_{\|\mathbf{a}\|_2=1} \|\mathbf{a}^* \mathbf{Y}\|_1.$$

More simplifications:

- orthogonal dictionary  $\mathbf{A}_0 = \mathbf{I}$ ;
- sparse coefficients  $\mathbf{X}_0 = \mathbf{I}$

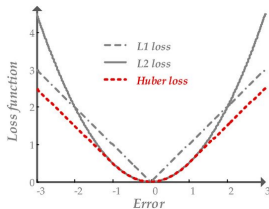
$$\min_{\|\mathbf{a}\|_2=1} \|\mathbf{a}\|_1.$$

# Orthogonal Dictionary Learning

$$\min_{\|\mathbf{a}\|_2=1} \|\mathbf{a}\|_1.$$

To obtain the second order information for stationary points, we use a smoothed  $\ell_1$  penalty — Huber loss

$$h_\lambda(x) = \begin{cases} \lambda|x| - \lambda^2/2 & |x| > \lambda, \\ x^2/2 & |x| \leq \lambda. \end{cases}$$

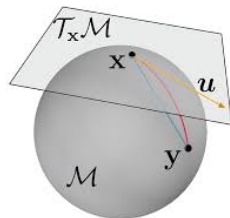


$$\min_{\|\mathbf{a}\|_2=1} \phi(\mathbf{a}) \doteq h_\lambda(\mathbf{a}).$$

# Orthogonal Dictionary Learning — Calculus

$$\min_{\|\mathbf{a}\|_2=1} \phi(\mathbf{a}) = h_\lambda(\mathbf{a}),$$

$$h_\lambda(a_i) = \begin{cases} \lambda|a_i| - \lambda^2/2 & |a_i| > \lambda, \\ a_i^2/2 & |a_i| \leq \lambda. \end{cases}$$



The Euclidean gradient

$$\nabla\phi = \lambda \text{sign}(\mathbf{a}) \circ \mathbf{1}_{|\mathbf{a}|>\lambda} + \mathbf{a} \circ \mathbf{1}_{|\mathbf{a}|\leq\lambda}.$$

With the sphere constraint, a critical point satisfies  $\nabla\phi = \mathbf{0}$  or  $\nabla\phi \propto \mathbf{a}$ .

$$\mathbf{a} \propto \text{sign}(\mathbf{a}).$$



# Orthogonal Dictionary Learning — Calculus

Recall that

$$\nabla\phi = \lambda \text{sign}(\mathbf{a}) \circ \mathbf{1}_{|\mathbf{a}|>\lambda} + \mathbf{a} \circ \mathbf{1}_{|\mathbf{a}|\leq\lambda}$$

has first-order critical points  $\mathbf{a} \propto \text{sign}(\mathbf{a})$ . Denote  $I = \text{supp}(\mathbf{a})$ , then the Riemannian Hessian over the sphere follows

$$\begin{aligned} \text{Hess}[\phi] &= \mathbf{P}_{\mathbf{a}^\perp} \left[ \underbrace{\nabla^2\phi}_{\text{curvature of } \phi} - \underbrace{\langle \nabla\phi, \mathbf{a} \rangle \mathbf{I}}_{\text{curvature of the sphere}} \right] \mathbf{P}_{\mathbf{a}^\perp} \\ &= \mathbf{P}_{\mathbf{a}^\perp} [\mathbf{D}_{\mathbf{1}_{|\mathbf{a}|\leq\lambda}} - \lambda |I| \mathbf{I}] \mathbf{P}_{\mathbf{a}^\perp} \end{aligned}$$

with  $\mathbf{P}_{\mathbf{a}^\perp} = \mathbf{I} - \mathbf{a}\mathbf{a}^T$ . The Hessian exhibits  $|I| - 1$  negative eigenvalues and  $n - |I|$  positive eigenvalues.

# Orthogonal Dictionary Learning — Calculus

- $\mathbf{a} = \pm \mathbf{e}_i$ , then the Hessian is positive definite

$$\text{Hess}[\phi] = \mathbf{P}_{\mathbf{a}^\perp} [(1 - \lambda)\mathbf{I} - \lambda \mathbf{D}_{\mathbf{e}_i}] \mathbf{P}_{\mathbf{a}^\perp} = \mathbf{P}_{\mathbf{a}^\perp} [(1 - \lambda)\mathbf{I}] \mathbf{P}_{\mathbf{a}^\perp}$$

with  $\mathbf{P}_{\mathbf{a}^\perp} = \mathbf{I} - \mathbf{e}_i \mathbf{e}_i^T = \mathbf{I} - \mathbf{D}_{\mathbf{e}_i}$ ;

- $\mathbf{a} = \sum_{i \in I} \pm \frac{1}{\sqrt{|I|}} \mathbf{e}_i$ , there exist negative curvatures along  $\mathbf{e}_i (i \in I)$

$$\text{Hess}[\phi] = \mathbf{P}_{\mathbf{a}^\perp} \left[ (1 - \lambda |I|) \mathbf{D}_{\mathbf{1}_{|a| \leq \lambda}} - \lambda |I| \mathbf{D}_{\mathbf{1}_{|a| > \lambda}} \right] \mathbf{P}_{\mathbf{a}^\perp}.$$

- $\mathbf{a} = \sum_{i \in [n]} \pm \frac{1}{\sqrt{n}} \mathbf{e}_i$ , then  $|I| = n$  and the Hessian is negative definite.

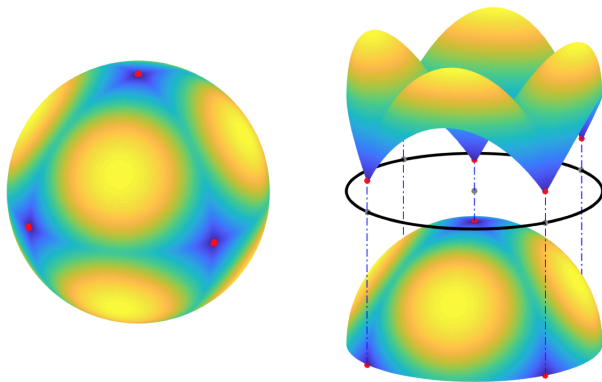
$$\text{Hess}[\phi] = \mathbf{P}_{\mathbf{a}^\perp} [-\lambda n \mathbf{I}] \mathbf{P}_{\mathbf{a}^\perp}$$

with  $\mathbf{P}_{\mathbf{a}^\perp} = \mathbf{I} - \mathbf{a} \mathbf{a}^T = (1 - 1/n) \mathbf{I}$ .

# Orthogonal Dictionary Learning — Geometry

**Local minimizers** are ground truth  $e_i$  or  $-e_i$ .

**Negative curvature** between multiple local minimizers.



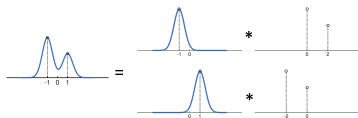
# Short-and-Sparse Blind Deconvolution

Goal: Given convolutional data  $\mathbf{y}$ , find the short signal  $\mathbf{a}$  and the sparse signal  $\mathbf{x}$  such that  $\mathbf{y} = \mathbf{a} * \mathbf{x}$ .

**Inherent Symmetry:**

$$\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0 = \alpha s_l[\mathbf{a}_0] * \frac{1}{\alpha} s_{-l}[\mathbf{x}_0]$$

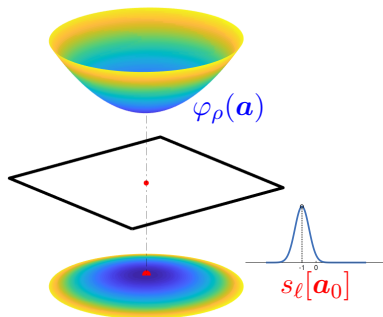
for any shift  $l$  and nonzero scaling.



The practical optimization problem can be written as

$$\min_{\|\mathbf{a}\|_F^2=1, \mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{a} * \mathbf{x}\|_F^2 + \lambda \|\mathbf{x}\|_1.$$

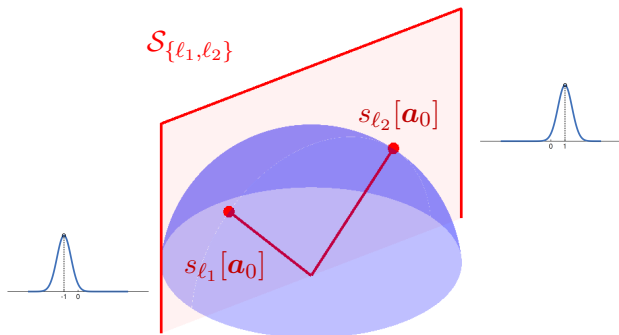
# Objective Function – Near One Shift



$$\mathbb{S}^{p-1} \cap \{\mathbf{a} \in \mathbb{S}^{p-1} \mid \|\mathbf{a} - s_\ell[\mathbf{a}_0]\|_2 \leq r\}$$

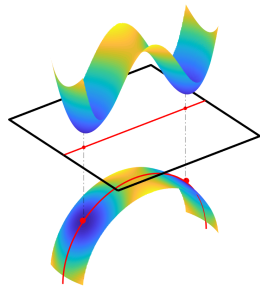
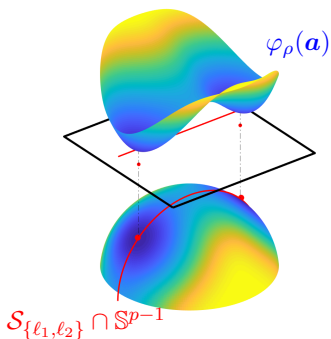
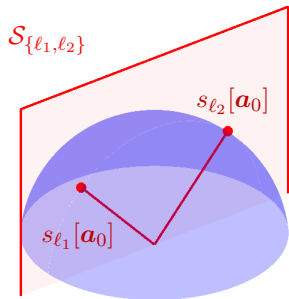
Objective function is **strongly convex** near a shift  $s_\ell[\mathbf{a}_0]$  of the ground truth.

# Objective Function – Linear Span of Two Shifts



$$\text{Subspace } \mathcal{S}_{\{l_1, l_2\}} = \{\alpha_{l_1} s_{l_1}[\mathbf{a}_0] + \alpha_{l_2} s_{l_2}[\mathbf{a}_0] \mid \alpha_{l_1}, \alpha_{l_2} \in \mathbb{R}\}.$$

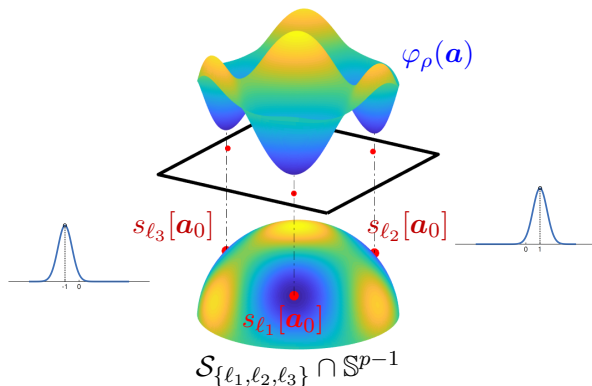
# Objective Function – Linear Span of Two Shifts



**Local minimizers** are near signed shifts  $\pm s_l[\mathbf{a}_0]$ .

**Negative curvature** between two shifts  $s_{l_1}[\mathbf{a}_0]$ ,  $s_{l_2}[\mathbf{a}_0]$ .

# Objective Function – Multiple Shifts

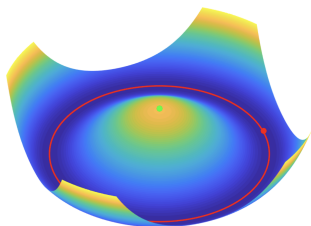


Objective  $\varphi_\rho$  over the linear span  $\mathcal{S}_{l_1, l_2, l_3} = \left\{ \sum_{i=1}^3 \alpha_i s_{l_i}[\mathbf{a}_0] \right\}$   
**Local minimizers** are near signed shifts  $\pm s_{l_i}[\mathbf{a}_0]$ .

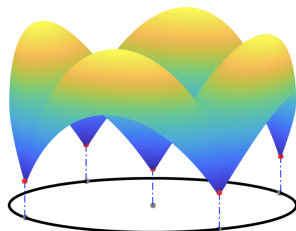


# Symmetry and Nonconvexity

- the (only!) local minimizers are symmetric versions of the ground truth.
- there is negative curvature in directions that break symmetry.



Rotational symmetry



Discrete symmetry

# Outline

- ① Introduction & Motivation of Nonconvex Optimization
  - Motivating Examples
  - Nonlinearity, Nonconvexity, and Symmetry
- ② Symmetry & Geometry for Nonconvex Problems in Practice
  - Problems with Rotational Symmetry
  - Problems with Discrete Symmetry
- ③ Efficient Nonconvex Optimization
  - Objectives of Nonconvex Optimization
  - Escaping Saddles

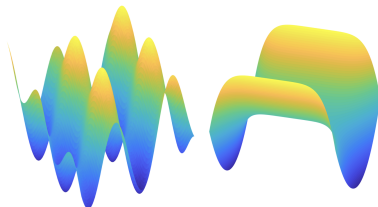
# Nonconvex Optimization

Consider the problem of minimizing a general nonlinear function:

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \mathbf{x} \in \mathcal{C}. \quad (5)$$

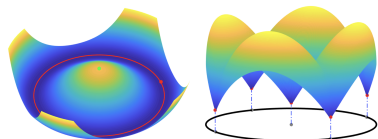
In **the worst case**, even finding a *local* minimizer can be NP-hard<sup>2</sup>.

Nonconvex problems that arise from natural physical, geometrical, or statistical origins typically have **nice** structures, in terms of **symmetries!**



Spurious local minimizers

Flat saddle points



Rotational symmetry

Discrete symmetry

<sup>2</sup>Some NP-complete problems in quadratic and nonlinear programming, K.G Murty and S. N. Kabadi, 1987

## Objectives

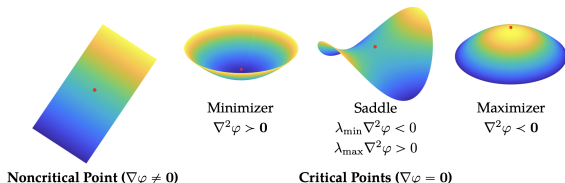
Hence typically people seek to work with relatively benign (gradient/Hessian Lipschitz continuous) functions:

$$\forall \mathbf{x}, \mathbf{y} \quad \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \leq L_1 \|\mathbf{y} - \mathbf{x}\|_2 \quad (6)$$

with benign objectives:

- ① convergence to some critical point  $\mathbf{x}_\star$  such that:  $\nabla f(\mathbf{x}_\star) = \mathbf{0}$ ;
- ② the critical point  $\mathbf{x}_\star$  is second-order stationary:  $\nabla^2 f(\mathbf{x}_\star) \succeq \mathbf{0}$ .

**Example:** in general  $f$  could have irregular second-order stationary points:

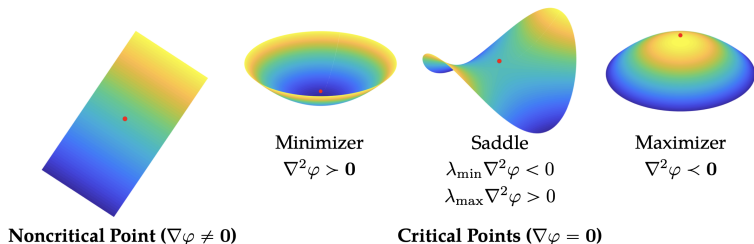


# Objectives

Hence typically people seek to work with relatively benign (gradient/Hessian Lipschitz continuous) functions with benign objectives:

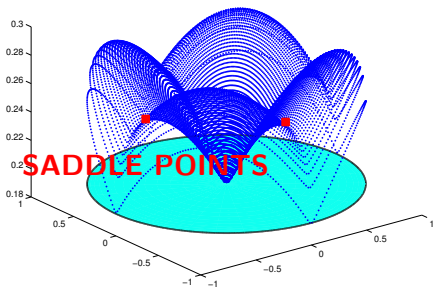
- ① convergence to some critical point  $\mathbf{x}_*$  such that:  $\nabla f(\mathbf{x}_*) = \mathbf{0}$ ;
- ② the critical point  $\mathbf{x}_*$  is second-order stationary:  $\nabla^2 f(\mathbf{x}_*) \succeq \mathbf{0}$ .

**Example:** a function  $\varphi$  with symmetry only has **regular** critical points:



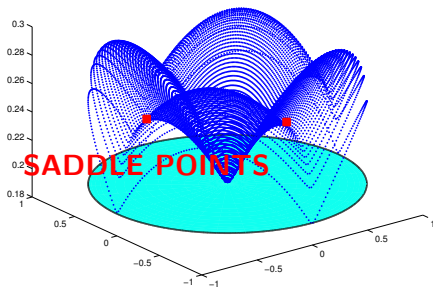
# “Any Reasonable Algorithm” Works

**Key issue:** using negative curvature  
 $\lambda_{\min}(\text{Hess}f) < 0$   
to escape saddles.



# “Any Reasonable Algorithm” Works

**Key issue:** using negative curvature  
 $\lambda_{\min}(\text{Hess}f) < 0$   
to escape saddles.

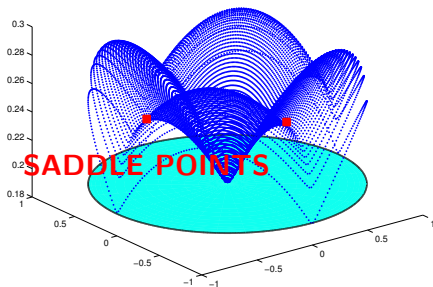


**Efficient (polynomial time) methods:**

- Trust region method, analyses in [Sun, Qu, W., '17]
- Curvilinear search, [Goldfarb, Mu, W., Zhou, '16]
- Noisy (stochastic) gradient descent, [Jin et. al. '17].

# “Any Reasonable Algorithm” Works

**Key issue:** using negative curvature  
 $\lambda_{\min}(\text{Hess}f) < 0$   
 to escape saddles.



## Efficient (polynomial time) methods:

Trust region method, analyses in [Sun, Qu, W., '17]

Curvilinear search, [Goldfarb, Mu, W., Zhou, '16]

Noisy (stochastic) gradient descent, [Jin et. al. '17].

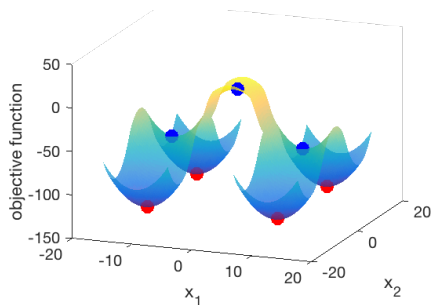
## Randomly initialized gradient descent ....

Obtains a minimizer almost surely [Lee et. al. '16].

Efficient for matrix completion, dictionary learning, ... not efficient in general.

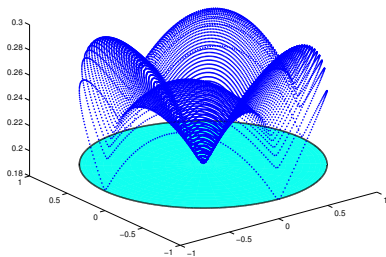


# Worst Case vs. Naturally Occurring Strict Saddle Functions



## Worst Case

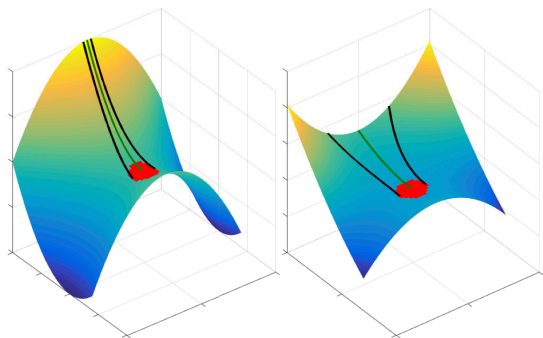
[Du, Jin, Lee, Jordan, Póczos, Singh '17]  
Concentration around stable manifold



## Naturally Occuring

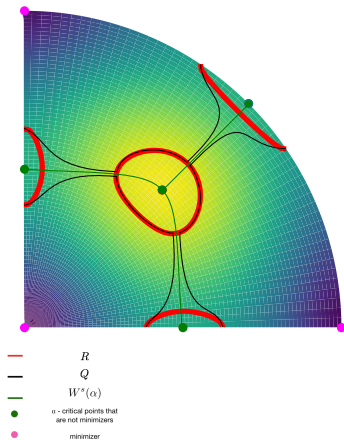
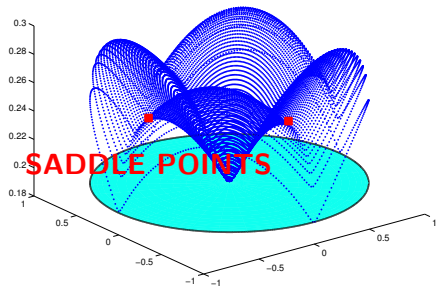
DL, Other sparsification problems  
Dispersion away from stable manifold

# Worst Case vs. Naturally Occurring Strict Saddle Functions

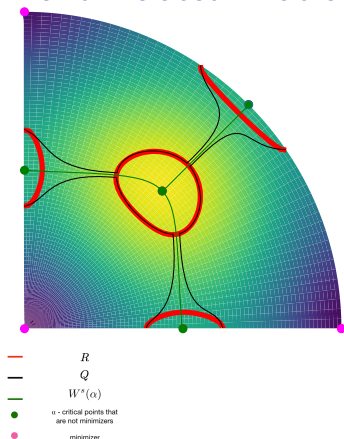
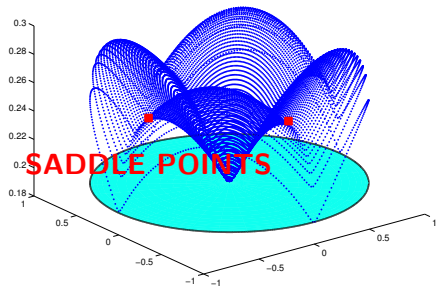


- Red: “slow region” of small gradient around a saddle point.
  - Green: stable manifold associated with the saddle point.
  - Black: points that flow to the slow region.
- Left: global negative curvature normal to the stable manifold
  - Right: positive curvature normal to the stable manifold – randomly initialized gradient descent is more likely to encounter the slow region.




# Gradient Descent Works for DL and Related Problems



# Gradient Descent Works for DL and Related Problems



**Dispersive structure:** Negative curvature  $\perp$  stable manifolds.

W.h.p. in random initialization  $\mathbf{q}^{(0)} \sim \text{uni}(\mathbb{S}^{n-1})$ , **convergence to a neighborhood of a minimizer in polynomial iterations.** [Gilboa,   

# Alternating Descent Method

$$\min_{\mathbf{a} \in \mathbb{S}^{n-1}, \mathbf{x}} \underbrace{\frac{1}{2} \|\mathbf{y} - \mathbf{a} * \mathbf{x}\|_F^2}_{\text{smooth } g} + \lambda \underbrace{\|\mathbf{x}\|_1}_{\text{nonsmooth } h}$$

- Fix  $\mathbf{a}$  and take a **proximal** descent step on  $\mathbf{x}$

$$\mathbf{x}^{(k+1)} \leftarrow \text{prox}_h^{\lambda t} \left( \mathbf{x}^{(k)} - t \nabla g(\mathbf{a}^{(k)}, \mathbf{x}^{(k)}) \right)$$

- Fix  $\mathbf{x}$  and take a **projected** descent step on  $\mathbf{a}$

$$\mathbf{a}^{(k+1)} \leftarrow \mathcal{P}_{\mathbb{S}^{n-1}} \left( \mathbf{a}^{(k)} - t' \text{grad}_g(\mathbf{a}^{(k)}, \mathbf{x}^{(k)}) \right)$$

# Inertial Alternating Descent Method

Accelerating first-order descent with Momentum

- Fix  $\mathbf{a}$  and take an **accelerated proximal** descent step on  $\mathbf{x}$

$$\begin{aligned}\mathbf{w}^{(k)} &= \mathbf{x}^{(k)} + \beta \left( \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \right) \\ \mathbf{x}^{(k+1)} &\leftarrow \text{prox}_h^{\lambda t} \left( \mathbf{x}^{(k)} - t \nabla g(\mathbf{a}^{(k)}, \mathbf{w}^{(k)}) \right)\end{aligned}$$

- Fix  $\mathbf{x}$  and take an **accelerated projected** descent step on  $\mathbf{a}$

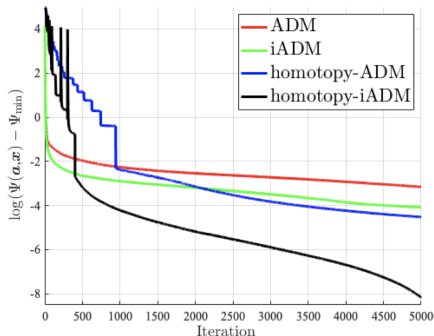
$$\begin{aligned}\mathbf{z}^{(k)} &= \mathbf{a}^{(k)} + \beta \left( \mathbf{a}^{(k)} - \mathbf{a}^{(k-1)} \right) \\ \mathbf{a}^{(k+1)} &\leftarrow \mathcal{P}_{\mathbb{S}^{n-1}} \left( \mathbf{a}^{(k)} - t' \text{grad}_g(\mathbf{z}^{(k)}, \mathbf{x}^{(k)}) \right)\end{aligned}$$

# Convergence Comparison

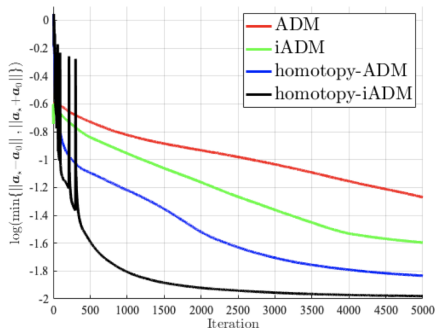
For blind deconvolution problem

$$\min_{\mathbf{a} \in \mathbb{S}^{n-1}, \mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{a} * \mathbf{x}\|_F^2 + \lambda \|\mathbf{x}\|_1$$

(a) function value convergence



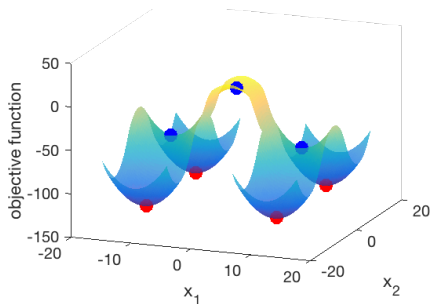
(b) iterate convergence



3

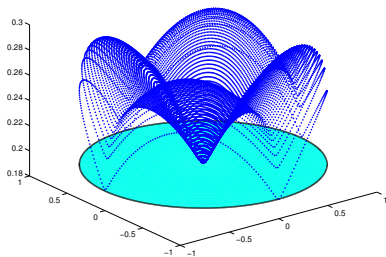
<sup>3</sup>The homotopy counterpart shrinks  $\lambda$  in every iteration.

# Escaping Saddles in Worst Case Problems



## Worst Case

[Du, Jin, Lee, Jordan, Póczos, Singh '17]  
Concentration around stable manifold



## Naturally Occuring

DL, Other sparsification problems  
Dispersion away from stable manifold



# Trust Region Method

**Function class:**  $f$  nonconvex.

**The oracle:** gradient  $\nabla f(\mathbf{x})$ , Hessian  $\nabla^2 f(\mathbf{x})$ , and the trusted region radius  $r$

**Trust region update:**

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \boldsymbol{\delta}$$

with

$$\boldsymbol{\delta} = \arg \min_{\|\boldsymbol{\delta}\| \leq r} f(\mathbf{x}^{(t)}) + \langle \nabla f(\mathbf{x}^{(k)}), \boldsymbol{\delta} \rangle + \frac{1}{2} \boldsymbol{\delta}^T \nabla^2 f(\mathbf{x}^{(k)}) \boldsymbol{\delta}$$

- At any stationary point, the gradient vanishes, and the above optimization problem boils down to the Hessian term;
- At an local solution with positive semi-definite Hessian, the above optimization problem renders  $\boldsymbol{\delta} = \mathbf{0}$ .

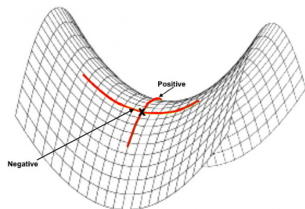
# Gradient Descent with Small Random Noise

**Function class:**  $f$  nonconvex and Lips. continuous.

**The oracle:** gradient  $\nabla f(\mathbf{x})$  and small random noise.

The updates for noisy gradient descent (Langevine dynamics):

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - t_1 \nabla f(\mathbf{x}^{(k)}) + t_2 \mathbf{n},$$



**This avoids computing expensive Hessian.**

# Hybrid Noisy Gradient Descent

**Function class:**  $f$  nonconvex and Lips. continuous.

**The oracle:** gradient  $\nabla f(\mathbf{x})$  and small noise  $\mathbf{n}$ .

**Hybrid noisy gradient descent:**

- **if**  $\|\nabla f(\mathbf{x}_k)\|_2 \geq \epsilon_g$ , **then**  $\mathbf{x}_{k+1} = \mathbf{x}_k - t_1 \nabla f(\mathbf{x}_k)$ ;
- **else**  $\mathbf{x}_k^0 = \mathbf{x}_k$ , and negative curvature descent with noisy gradients:  
**for**  $i = 0, 1, 2, \dots, k_{\max} = O(\log n)$

$$\mathbf{x}_k^{i+1} = \mathbf{x}_k^i - t_1 \nabla f(\mathbf{x}_k^i) + t_2 \mathbf{n}^i,$$

where  $\mathbf{n}^i \sim \mathcal{N}(0, \mathbf{I})$ .

**More saddle-escaping first-order optimization methods in book:**  
**Wright and Ma:** <https://book-wright-ma.github.io>.

# Conclusion and Coming Attractions

For Nonconvex, Sparse and Low-rank problems

- **Benign Geometry:**
  - The only local minimizers are symmetric copies of the ground truth
  - There exist negative curvatures breaking symmetry
- **Efficient Algorithms:**
  - gradient descent algorithms always suffice
  - proximal, projection, acceleration steps can be transferred over

**Next lecture: Exploiting Low-D Structures via Deep Networks.**

**Thank You! Questions?**