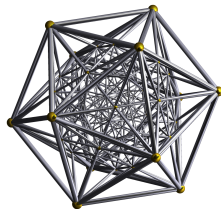ACDL Summer Course 2023

## Lecture 1: Low-Dimensional and Nonconvex Models for Shallow Representation Learning
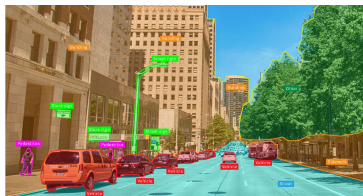
**Qing Qu**

EECS, University of Michigan

June 10th, 2023

# The Success of Deep Learning



**computer vision**

(Credit: Appen. (2019))



**natural language processing**

(Credit: Andrey Suslov (2023))



**gameplay**

(Credit: AlphaGo)



**autonomous driving**
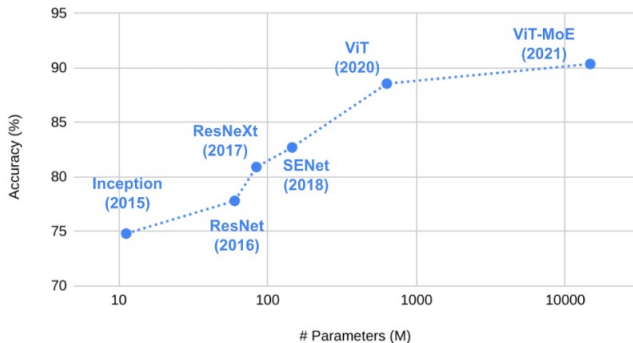
(Credit: Phil Brown (2019))

# The Trend of Large Models...



**Figure**: Accuracy vs. model size for image classification on ImageNet dataset

### ~23 million >> ~1 million
(# Parameters in ResNet-50)     (# Samples in ImageNet)

**In principle, deep network can fit _any_ training labels!**
(*i.e.*, not only clean, but also corrupted labels)

# The Challenges & Opportunities in Large Models...

- **Tremendous cost of computation**
- **Difficult to interpret**
- **Vulnerable to data corruptions**



**Figure**: Accuracy vs. model size for image classification on ImageNet dataset

# The Challenges & Opportunities in Large Models...

- **Tremendous cost of computation**
- **Difficult to interpret**
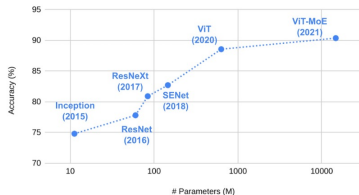- **Vulnerable to data corruptions**



**Figure**: Accuracy vs. model size for image classification on ImageNet dataset

# The Challenges & Opportunities in Large Models...

- **Tremendous cost of computation**
- **Difficult to interpret**
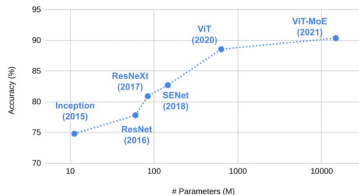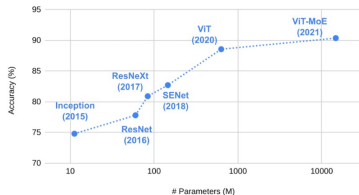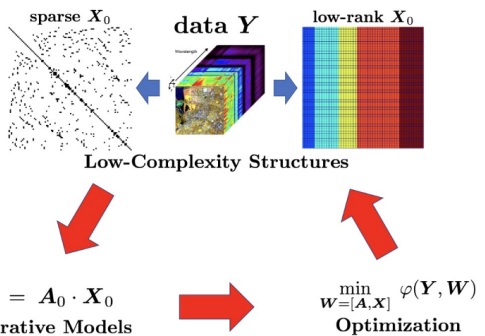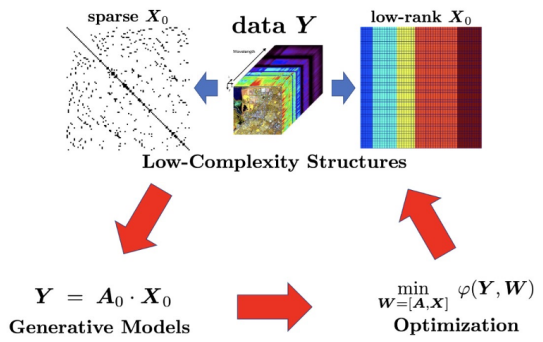- **Vulnerable to data corruptions**



Figure: Accuracy vs. model size for image classification on ImageNet dataset

**Theory and principles behind its success?**

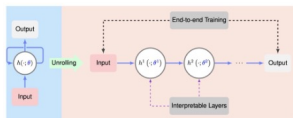# Low-Dimensional Structures Are Largely Ignored...

# Low-Dimensional Structures Are Largely Ignored...



sparse $X_0$   data $Y$   low-rank $X_0$

**Low-Complexity Structures**

$Y = A_0 \cdot X_0$
**Generative Models**

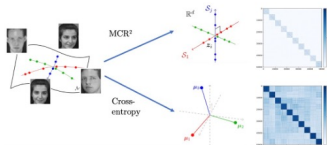$\min_{W=[A,X]} \varphi(Y, W)$
**Optimization**

- **Sparse Recovery**
  [Donoho'06, Candes'08]

- **Low-rank Matrix Recovery**
  [Candes'08, Recht'11,
  Candes'11]

- **(Sparse) Phase Retrieval**
  [Candes'13, Shechtman'15]

- **Super-resolution**
  [Candes'14,
  Fernandez-Granda'16]

- **(Sparse) Blind
  Deconvolution** [Ahmed'14,
  Zhang'17, Kuo'20]

- **(Convolutional) Dictionary
  Learning**[Aharon'06,
  Sun'16, Bristow'13,
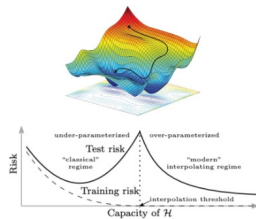  Papyan'17]

# The Emergence of Low-Dim Models in Deep Learning



**Network Architectures**

[Gregor'10, Liu'18, Sulam'18, Papyan'18, Monga'19]

**Representations**

[Pennington'17, Bansal'18, Xiao'18, Wang'20, Ye'20, Qi'20, Han'20, Zhu'21, Fang'21]

**Regularizations & Generalization**

[Neyshabur'17, Mianjy'18, Ulyanov'18, Gidel'19, Arora'19, Belkin'19, Nakkiran'19, Yang'20]

- image credited to Monga et al., Yu et al. & Azizan et al.
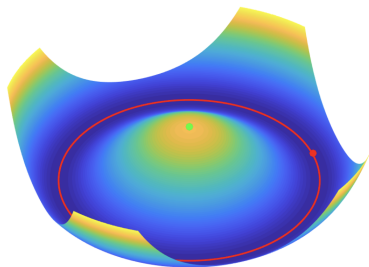
# Outline of Today's Course

**Lec.1** **Low-dimensional Models & Noconvex Optimization (1hrs)**

**Lec.2** **Low-dimensional Representations in Deep Learning I: Neural Collapse (1hrs)**
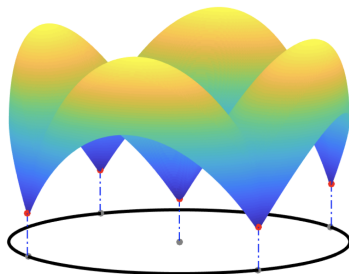
**Lec.3** **Low-dimensional Representations in Deep Learning II: Law-of-Parsimony in GD (1.5hrs)**

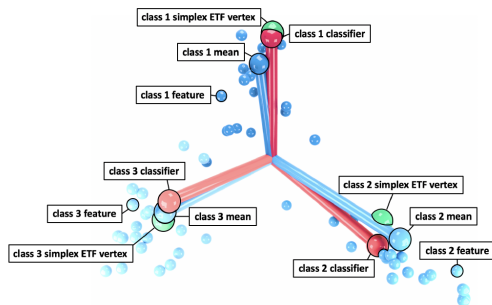**Lec.4** **Low-dimensional Models for Robust Learning (0.5hrs)**

# Outline of Today's Course



**Rotational symmetry**　　　　　　**Discrete symmetry**

**Lec.1 Low-dimensional Models & Noconvex Optimization (1hrs)**
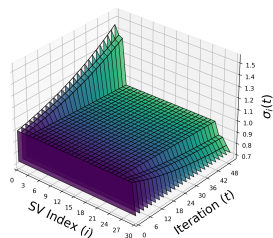
# Outline of Today's Course



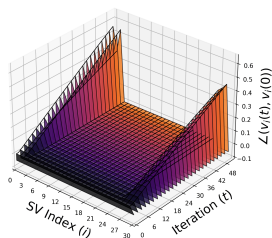**Credit**: Han et al.

Neural Collapse Under MSE Loss: Proximity to and Dynamics on the Central Path. ICLR, 2022.

**Lec.2 Low-dimensional Representations in Deep Learning I: Neural Collapse (1hrs)**
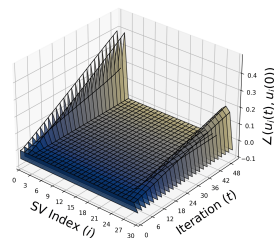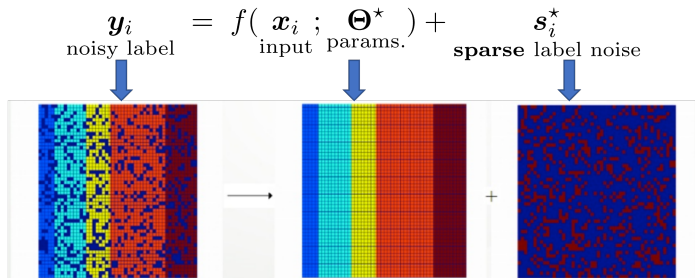
# Outline of Today's Course



Singular Values      Right Singular Vectors      Left Singular Vectors

**Lec.3** **Low-dimensional Representations in Deep Learning II:**
**Law-of-Parsimony in GD (1.5hrs)**

# Outline of Today's Course



$$\underset{\text{noisy label}}{\boldsymbol{y}_i} = f(\underset{\text{input}}{\boldsymbol{x}_i} ; \underset{\text{params.}}{\boldsymbol{\Theta}^{\star}}) + \underset{\textbf{sparse} \text{ label noise}}{\boldsymbol{s}_i^{\star}}$$

**Exact Separation of Sparse Corruption with Incoherence between Data and Noise**

**Lec.4 Low-dimensional Models for Robust Learning (0.5hrs)**

# Outline

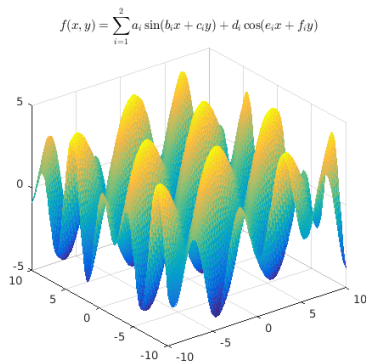# Most of the Machine Learning Problems are Nonconvex...
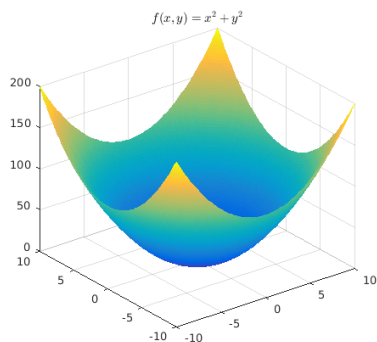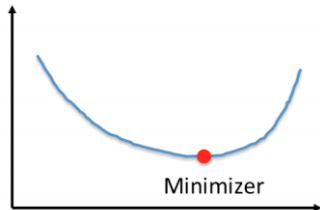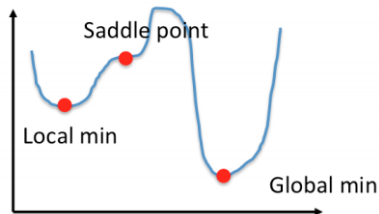


Figure: **Convex vs. Nonconvex Optimization Problems.**

# Basic Calculus

Critical points or stationary points: gradient vanishes



**Convex**          **Non-Convex**

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

# Basic Calculus

Critical points with non-singular hessian

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



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

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

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

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

# Challenges of Nonconvex Optimization – Pessimistic Views

Consider the problem of minimizing a general nonlinear function:

$$\min_{\boldsymbol{z}} \varphi(\boldsymbol{z}), \quad \boldsymbol{z} \in \mathsf{C}. \qquad (1)$$

In **the worst case**, even finding a **local** minimizer can be NP-hard[1].



**Spurious local minimizers**     **Flat saddle points**

---

[1]Some NP-complete problems in quadratic and nonlinear programming, K.G Murty and S. N. Kabadi, 1987

# Challenges of Nonconvex Optimization – Pessimistic Views

Consider the problem of minimizing a general nonlinear function:

$$\min_{\boldsymbol{z}} \varphi(\boldsymbol{z}), \quad \boldsymbol{z} \in \mathsf{C}. \qquad (1)$$

In **the worst case**, even finding a **local** minimizer can be NP-hard[1].

Hence, typically people seek to work with **mild guarantees** for nonconvex problems:



**Spurious local minimizers**          **Flat saddle points**

1. convergence to some **critical point** $\bar{z}$ such that $\nabla\varphi(\bar{z}) = \boldsymbol{0}$;

2. or convergence to some **local minimizer** $\nabla^2\varphi(\bar{z}) \succeq \boldsymbol{0}$.

---

[1]Some NP-complete problems in quadratic and nonlinear programming, K.G Murty and S. N. Kabadi, 1987

# Benign Nonconvex Optimization Landscape



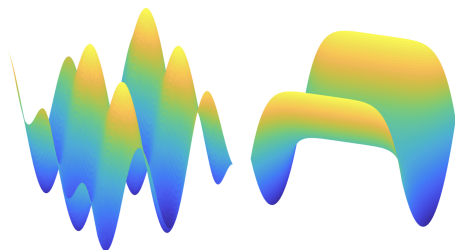**General Case**          **Structured Case**

# Benign Nonconvex Optimization Landscape



General nonconvex problems

Our training problem

**General Case**

**Structured Case**

# Example I: Low-rank Matrix Completion



We observe:

$$\underset{\text{Observed ratings}}{\boldsymbol{Y}} = \mathcal{P}_\Omega \left[ \underset{\text{Complete ratings}}{\boldsymbol{X}} \right].$$

# Example I: Low-rank Matrix Completion



We observe:

$$\underset{\text{Observed ratings}}{\boldsymbol{Y}} = \mathcal{P}_\Omega \left[ \underset{\text{Complete ratings}}{\boldsymbol{X}} \right].$$

**Matrix completion** via nonconvex Burer-Monteiro factorization

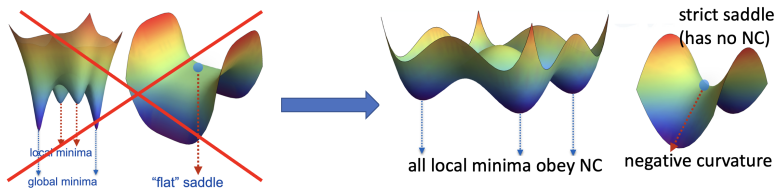$$\min_{\boldsymbol{U}, \boldsymbol{V}} f(\boldsymbol{U}, \boldsymbol{V}) = \sum_{(i,j) \in \Omega} [(\boldsymbol{U}\boldsymbol{V}^*)_{i,j} - \boldsymbol{Y}_{i,j}]^2 + \underbrace{\frac{\lambda}{2}\|\boldsymbol{U}\|_F^2 + \frac{\lambda}{2}\|\boldsymbol{V}\|_F^2}_{\text{reg}(\boldsymbol{U}, \boldsymbol{V})}.$$

# Example II: Dictionary for Image Representation

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



$$I_{\text{noisy}} = \underset{\text{dictionary}}{\boldsymbol{A}} \times \underset{\text{sparse}}{\boldsymbol{x}} + \underset{\text{noise}}{\boldsymbol{z}}. \quad (2)$$

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

$$\underset{\text{data}}{\boldsymbol{Y}} = \underset{\text{dictionary}}{\boldsymbol{A}} \ \underset{\text{sparse}}{\boldsymbol{X}}. \quad (3)$$

# Example II: Dictionary for Image Representation

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



$$I_{\mathsf{noisy}} = \underset{\mathsf{dictionary}}{\boldsymbol{A}} \times \underset{\mathsf{sparse}}{\boldsymbol{x}} + \underset{\mathsf{noise}}{\boldsymbol{z}}. \quad (2)$$

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

$$\underset{\mathsf{data}}{\boldsymbol{Y}} = \underset{\mathsf{dictionary}}{\boldsymbol{A}} \underset{\mathsf{sparse}}{\boldsymbol{X}}. \quad (3)$$

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

# Dictionary Learning



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

# Example: Sparse Blind Deconvolution

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

$$\underset{\text{data}}{\boldsymbol{Y}} = \underset{\text{motif}}{\boldsymbol{A}} * \underset{\text{sparse}}{\boldsymbol{X}}. \qquad (4)$$

- Scientific signals:
  activation signals are sparse

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



Observation Y    Kernel A₀    Activation Map X₀

Observation    Kernel A₀    Natural Image

# Sparse Blind Deconvolution



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

# Convolutional Dictionary learning

$$\underset{\text{data}}{\boldsymbol{Y}} = \sum_i \underset{\text{motif}}{\boldsymbol{A}_i} * \underset{\text{sparse}}{\boldsymbol{X}_i}.$$



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

# Opportunities – Optimistic Views

Nonconvex
problems that arise
in machine learning typically
have **benign** data structures,
in terms of **symmetries**!



Rotational symmetry



Discrete symmetry

# Opportunities – Optimistic Views

Nonconvex
problems that arise
in machine learning typically
have **benign** data structures,
in terms of **symmetries**!



**Rotational symmetry**          **Discrete symmetry**

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

- **low rank matrix recovery**: invariant under a continuous rotation:

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

- **dictionary learning**: invariant under signed permutations:

# Opportunities – Optimistic Views

Nonconvex
problems that arise
in machine learning typically
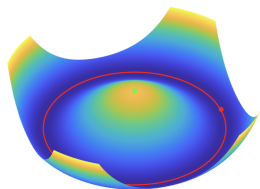have **benign** data structures,
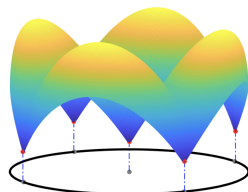in terms of **symmetries**!



**Rotational symmetry**

**Discrete symmetry**

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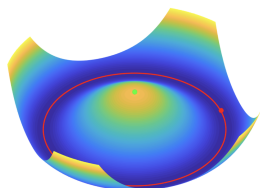under certain group action:

- **low rank matrix recovery**: invariant under a continuous rotation:

$$\varphi((U\Gamma, V\Gamma^{-1})) = \varphi((U, V)), \quad \forall \text{ invertible } \Gamma.$$
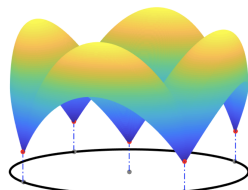
- **dictionary learning**: invariant under signed permutations:

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

## Nonlinearity and Symmetry

Intrinsic ambiguity against the uniqueness of the solution

- **low rank matrix recovery**

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

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

- **dictionary learning**

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

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

- **blind deconvolution**

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

  for any signed shift $\tau$.

# Optimization under Symmetry

## Definition (Symmetric Function)

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

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



Rotational symmetry          Discrete symmetry
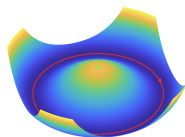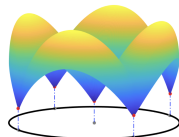
# Optimization under Symmetry

### Definition (Symmetric Function)

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

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



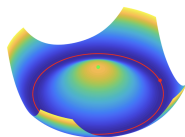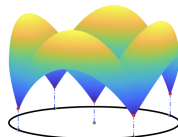Rotational symmetry          Discrete symmetry

**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.

# Outline

# Problems with Rotational Symmetry



**Nonconvex Problems with Rotational Symmetries**

**Eigenspace Computation**
*Compute the principal subspace of a symmetric matrix.*

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

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

**Generalized Phase Retrieval**
*Recover a complex vector $x_0$ from magnitude measurements $y = |Ax_0|$.*

$\min_x \frac{1}{2}\|y^2 - |Ax|^2\|_2^2$.

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

**Matrix Recovery**
*Recover a low-rank matrix $X = UV^*$ from incomplete / corrupted observations*

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

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

# Low Rank Matrix Recovery

**Goal:** Given $Y = \mathcal{A}(X)$, recover low rank matrix $X = U_0 V_0$

# Low Rank Matrix Recovery

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



- **Convex formulation:**

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

# Low Rank Matrix Recovery

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



- **Convex formulation:**

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

- **Nonconvex formulation:**

$$\min_{\boldsymbol{U} \in \mathbb{R}^{m \times r}, \boldsymbol{V} \in \mathbb{R}^{n \times r}} \quad \left\|\boldsymbol{Y} - \mathcal{A}(\boldsymbol{U}\boldsymbol{V}^T)\right\|_F^2 + \mathsf{reg}(\boldsymbol{U}, \boldsymbol{V})$$

# Low Rank Matrix Recovery

$$\min_{\boldsymbol{U}, \boldsymbol{V}} \quad \frac{1}{2} \left\| \boldsymbol{Y} - \mathcal{A}(\boldsymbol{U}\boldsymbol{V}^T) \right\|_F^2 + \mathsf{reg}(\boldsymbol{U}, \boldsymbol{V})$$

**Inherent Symmetry**:

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

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

# Low Rank Matrix Recovery

$$\min_{\boldsymbol{U},\boldsymbol{V}} \quad \frac{1}{2}\left\|\boldsymbol{Y} - \mathcal{A}(\boldsymbol{U}\boldsymbol{V}^T)\right\|_F^2 + \mathsf{reg}(\boldsymbol{U},\boldsymbol{V})$$

**Inherent Symmetry**:

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

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



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

# Simple Setting: Rank-1 Symmetric Matrix

- **Simplifications:**
    - $Y = \mathcal{A}(X) = X$
    - $X = U_0 U_0^T$ is symmetric and rank-1

    $$X = u_0 u_0^T = (-u_0 Q)(-Q^T u_0^T)$$

    the signed rotational symmetry.

# Simple Setting: Rank-1 Symmetric Matrix

- **Simplifications:**
  - $Y = \mathcal{A}(X) = X$
  - $X = U_0 U_0^T$ is symmetric and rank-1

$$X = u_0 u_0^T = (-u_0 Q)(-Q^T u_0^T)$$

  the signed rotational symmetry.

- **Nonconvex formulation:**

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

## Rank-1 Symmetric Matrix

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

- Critical points have zero gradient

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

# Rank-1 Symmetric Matrix

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

- Critical points have zero gradient

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

- Therefore, critical points must be one of the following
    - $\boldsymbol{u} = \pm\boldsymbol{Q}\boldsymbol{u}_0$
    - $\boldsymbol{u} = \boldsymbol{0}$

## Rank-1 Symmetric Matrix

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

with the second-order derivative

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

## Rank-1 Symmetric Matrix

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

with the second-order derivative

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

Then the stationary points can be grouped as

- Local minimizer $\boldsymbol{u} = \pm\boldsymbol{Q}\boldsymbol{u}_0$:

$$\nabla^2 \phi = \boldsymbol{u}\boldsymbol{u}^T + \|\boldsymbol{u}\|_2^2 \, \boldsymbol{I} \succeq \boldsymbol{0}$$

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

$$\nabla^2 \phi = -\boldsymbol{X} < \boldsymbol{0}.$$

# Low Rank Matrix Recovery

- Symmetric low rank matrix recovery:

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

- General low rank matrix recovery:

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

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

# Problems with Discrete Symmetry



**Nonconvex Problems with Discrete Symmetries**

**Eigenvector Computation**
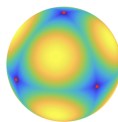
*Maximize a quadratic form over the sphere.*

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

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

**Dictionary Learning**

*Approximate a given matrix $Y$ as $Y \approx AX$, with $X$ sparse*

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

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

**Tensor Decomposition**

*Determine components $a_i$ of an orthogonal decomposable tensor $T = \sum_i a_i \otimes a_i \otimes a_i \otimes a_i$*

$\max_{X \in O(n)} \sum_i T(x_i, x_i, x_i, x_i).$

**Symmetry:** $X \mapsto X\Gamma$
$\mathbb{G} = \mathrm{P}(n)$

**Short-and-Sparse Deconvolution**

*Recover a short $a$ and a sparse $x$ from their convolution $y = a * x$.*

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

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

## Dictionary Learning

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

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

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

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

## Dictionary Learning

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

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

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

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

**Inherent Symmetry**:

$$Y = A_0 \Gamma \Gamma^* X_0,$$

for any signed permutation matrix $\Gamma$.

## Orthogonal Dictionary Learning

- Input: matrix $Y$ which is the product of an orthogonal matrix $A_0$ (called a dictionary) and a sparse matrix $X_0$:

$$Y = A_0 X_0, \quad A_0 A_0^* = I, X_0 \text{ sparse.}$$

- Optimization formulation:

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

## Orthogonal Dictionary Learning

- Input: matrix $Y$ which is the product of an orthogonal matrix $A_0$ (called a dictionary) and a sparse matrix $X_0$:

$$Y = A_0 X_0, \quad A_0 A_0^* = I, X_0 \text{ sparse.}$$

- Optimization formulation:

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

- Given the constraint, $X$ is uniquely defined in terms of $A$

$$X = A^* A X = A^* Y.$$

## Orthogonal Dictionary Learning

- Input: matrix $Y$ which is the product of an orthogonal matrix $A_0$ (called a dictionary) and a sparse matrix $X_0$:

$$Y = A_0 X_0, \quad A_0 A_0^* = I, X_0 \text{ sparse}.$$

- Optimization formulation:

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

- Given the constraint, $X$ is uniquely defined in terms of $A$

$$X = A^* A X = A^* Y.$$

- Equivalent formulation:

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

## Orthogonal Dictionary Learning

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

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

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

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

## Orthogonal Dictionary Learning

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

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

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

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

Stationary Points:

- $\boldsymbol{a} = \pm \boldsymbol{a}_i$, then the Hessian is positive definite
- $\boldsymbol{a} = \sum_{i \in I} \pm \frac{1}{\sqrt{|I|}} \boldsymbol{a}_i$, there exist negative curvatures alone $\boldsymbol{a}_i (i \in I)$

# Orthogonal Dictionary Learning — Geometry

**Local minimizers** are ground truth $a_i$ or $-a_i$.
**Negative curvature** between multiple local minimizers.

# Short-and-Sparse Blind Deconvolution

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

**Inherent Symmetry**:

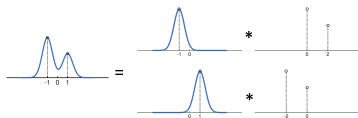$$y = a_0 * x_0 = \alpha s_l[a_0] * \frac{1}{\alpha} s_{-l}[x_0]$$

for any shift $l$ and nonzero scaling.

# Short-and-Sparse Blind Deconvolution

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

**Inherent Symmetry**:

$$y = a_0 * x_0 = \alpha s_l[a_0] * \frac{1}{\alpha} s_{-l}[x_0]$$

for any shift $l$ and nonzero scaling.



The practical optimization problem can be written as

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

# Objective Function – Near One Shift



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

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

# Objective Function – Linear Span of Two Shifts



**Subspace** $\mathcal{S}_{\{\ell_1,\ell_2\}} = \{\alpha_{\ell_1} s_{\ell_1}[\boldsymbol{a}_0] + \alpha_{\ell_2} s_{\ell_2}[\boldsymbol{a}_0] \mid \alpha_{\ell_1}, \alpha_{\ell_2} \in \mathbb{R}\}$.

# Objective Function – Linear Span of Two Shifts



**Local minimizers** are near signed shifts $\pm s_\ell[\boldsymbol{a}_0]$.
**Negative curvature** between two shifts $s_{\ell_1}[\boldsymbol{a}_0]$, $s_{\ell_2}[\boldsymbol{a}_0]$.

# Objective Function – Multiple Shifts



Objective $\varphi_\rho$ over the linear span $\mathcal{S}_{\ell_1,\ell_2,\ell_3} = \{\sum_{i=1}^{3} \alpha_{\ell_i} s_{\ell_i}[\boldsymbol{a}_0]\}$
**Local minimizers** are near signed shifts $\pm s_{\ell_i}[\boldsymbol{a}_0]$.

# Symmetry and Nonconvexity



**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.

# Outline

# Nonconvex Optimization in Generic Setting

Consider the problem of minimizing a **general** nonconvex function:

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}), \quad \boldsymbol{x} \in \mathsf{C}. \qquad (5)$$

In the worst case, even finding a *local* minimizer can be NP-hard[2].

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**

---

[2]Some NP-complete problems in quadratic and nonlinear programming, K.G Murty and S. N. Kabadi, 1987

## Nonconvex Optimization in Generic Setting

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

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

with benign objectives:

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

# Nonconvex Optimization in Generic Setting

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

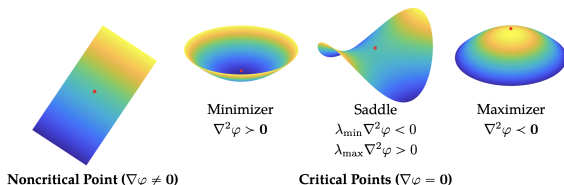$$\forall \boldsymbol{x}, \boldsymbol{y} \quad \|\nabla f(\boldsymbol{y}) - \nabla f(\boldsymbol{x})\|_2 \leq L_1 \|\boldsymbol{y} - \boldsymbol{x}\|_2 \tag{6}$$

with benign objectives:

① convergence to some critical point $\boldsymbol{x}_\star$ such that: $\nabla f(\boldsymbol{x}_\star) = \boldsymbol{0}$;

② the critical point $\boldsymbol{x}_\star$ is second-order stationary: $\nabla^2 f(\boldsymbol{x}_\star) \succeq \boldsymbol{0}$.

**Example:** a function $f$ with symmetry only has **regular** critical points, while general $f$ could have irregular second-order stationary points:



Minimizer
$\nabla^2 \varphi > \boldsymbol{0}$

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

Maximizer
$\nabla^2 \varphi < \boldsymbol{0}$

**Noncritical Point ($\nabla \varphi \neq \boldsymbol{0}$)**          **Critical Points ($\nabla \varphi = \boldsymbol{0}$)**

# Benign Nonconvexity: "Any Reasonable Algorithm" Works

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

# Benign Nonconvexity: "Any Reasonable Algorithm" Works

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



SADDLE POINTS

- **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].

# Benign Nonconvexity: "Any Reasonable Algorithm" Works

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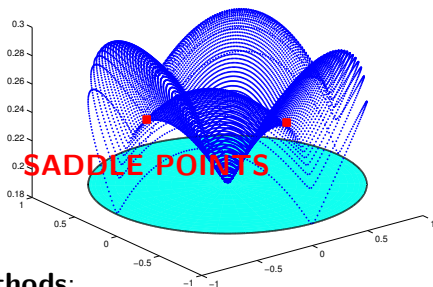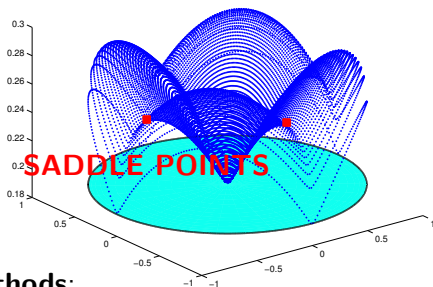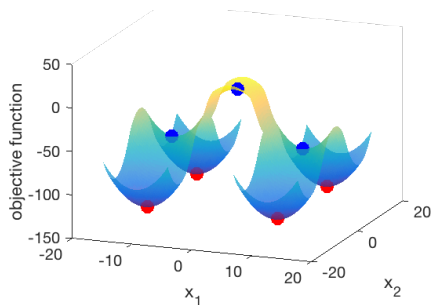$$\lambda_{\min}(\mathrm{Hess}f) < 0$$
to escape saddles.



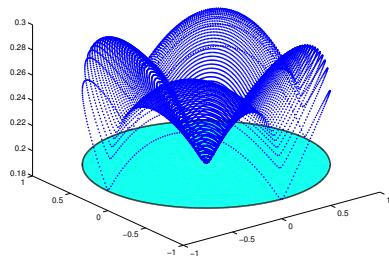- **Efficient (polynomial time) methods**:
  Trust region method, analyses in [Sun, Qu, W., '17]
  Curvilinear search, [Goldfarb, Mu, W., Zhou, '16]
  Noisy (stochastic) gradient descent, [Jin et. al. '17].

- **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
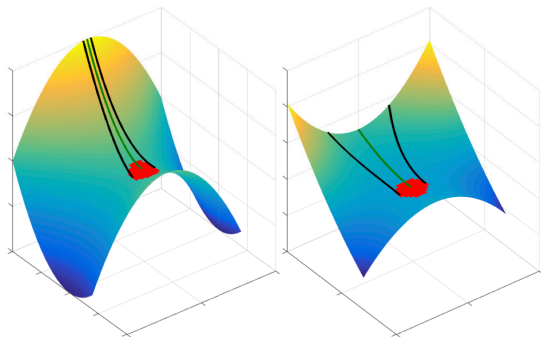


**Worst Case**

[Du, Jin, Lee, Jordan, Poczos, 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



| | |
|---|---|
| — | $R$ |
| — | $Q$ |
| — | $W^s(\alpha)$ |
| ● | $\alpha$ - critical points that are not minimizers |
| ● | minimizer |

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

W.h.p. in random initialization $q^{(0)} \sim \mathrm{uni}(\mathbb{S}^{n-1})$, **convergence to a neighborhood of a minimizer in polynomial iterations.** [Gilboa, Buchanan, W. '18]

# Outline

# References



SCAN ME

1 Zhang Y, Qu Q, Wright J. From symmetry to geometry: Tractable nonconvex problems [J]. arXiv preprint arXiv:2007.06753, 2020.

2 Qu Q, Zhu Z, Li X, et al. Finding the sparsest vectors in a subspace: Theory, algorithms, and applications [J]. arXiv preprint arXiv:2001.06970, 2020.

3 Q. Qu, Y. Zhai, X. Li, Y. Zhang, Z. Zhu, Analysis of optimization landscapes for overcomplete learning, ICLR'20, (oral, top 1.9%)

4 Y. Lau (*), Q. Qu(*), H. Kuo, P. Zhou, Y. Zhang, J. Wright, Short-and-sparse Deconvolution – A Geometric Approach, ICLR'20

# 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

## Thank You! Questions?

# Call for Papers

- IEEE JSTSP Special Issue on Seeking Low-dimensionality in Deep Neural Networks (SLowDNN) Manuscript Due: **Nov. 30, 2023**.

- Conference on Parsimony and Learning (CPAL) January 2024, Hongkong, Manuscript Due: **Aug. 28, 2023**.



SCAN ME



SCAN ME