



EECS 559
Optimization Methods for
SIPML

Lecture 11 – Intro on Nonconvex Problems in SIPML

Instructor: Prof. Qing Qu (qingqu@umich.edu)

Lecture Agenda

- Problem Introduction
- Applications in SIPML

Lecture Agenda

- **Problem Introduction**
 - General Nonconvex Problems
 - Global Optimization?
 - Strict Saddle Problems
- Applications in SIPML

Lecture Agenda

- Problem Introduction
 - General Nonconvex Problems
 - Global Optimization?
 - Strict Saddle Problems
- Applications in SIPML

Nonconvex Optimization

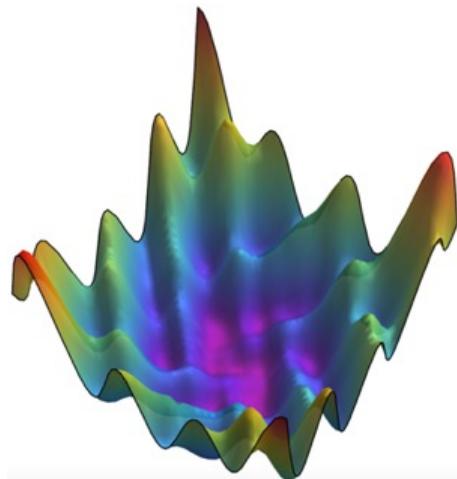
$$\min_{\boldsymbol{x}} f(\boldsymbol{x}), \text{ s.t. } \boldsymbol{x} \in \mathcal{C}.$$

The problem is *nonconvex* if either

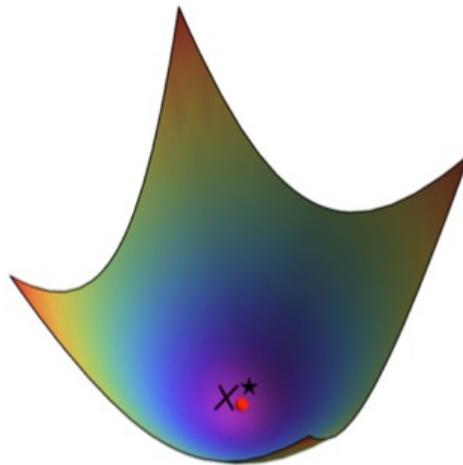
- the objective $f(\boldsymbol{x})$ is nonconvex,
- or the constraint set \mathcal{C} is nonconvex.

Convex vs Nonconvex Problems

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}), \text{ s.t. } \boldsymbol{x} \in \mathbb{R}^n$$

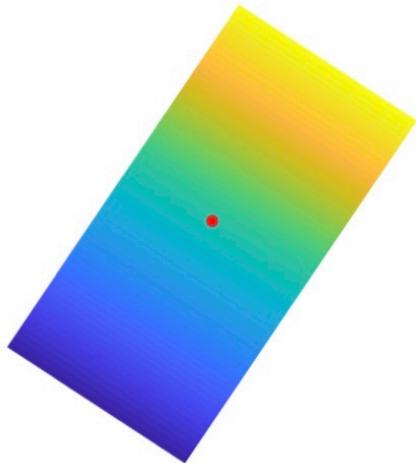


Nonconvex landscape

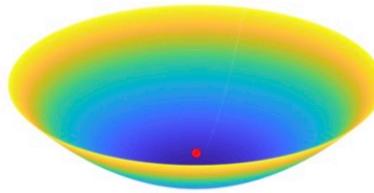


Convex landscape

Classification of Critical Points

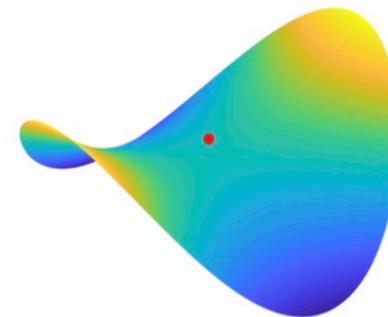


Noncritical Point ($\nabla\varphi \neq \mathbf{0}$)



Minimizer

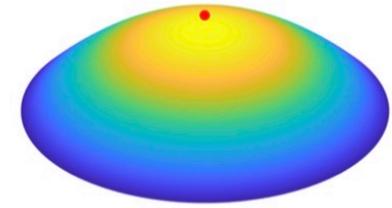
$$\nabla^2\varphi \succ \mathbf{0}$$



Saddle

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

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



Maximizer

$$\nabla^2\varphi \prec \mathbf{0}$$

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

Classification of Critical Points

Definition. (Critical points, local minimizers, & saddle points) Suppose $f : \mathbb{R}^n \mapsto \mathbb{R}$ is continuously differentiable, then

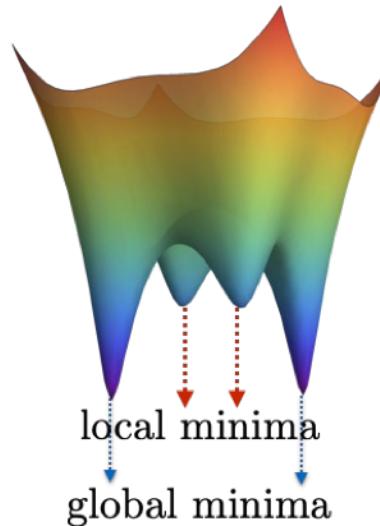
- A point x_* is called a *critical point* of f if $\nabla f(x_*) = \mathbf{0}$;
- A point x_* is a *local minimizer* if $\exists \delta > 0$ such that $f(x_*) \leq f(x)$ for all $x \in \mathcal{B}(x_*, \delta)$;
- A point x_* is a *saddle point* if it is a critical point but neither a local minimizer nor local maximizer.

A local minimizer is a critical point, but not vice versa.

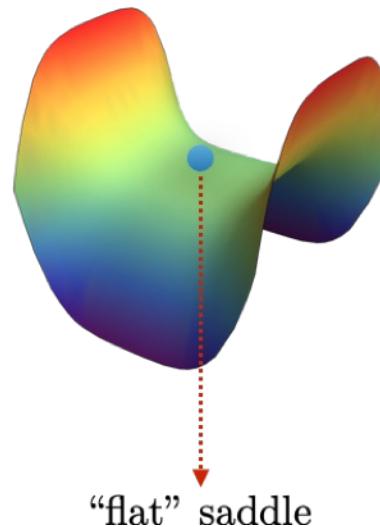
General Nonconvex Problems

$$\min_x f(x), \text{ s.t. } x \in \mathbb{R}^n$$

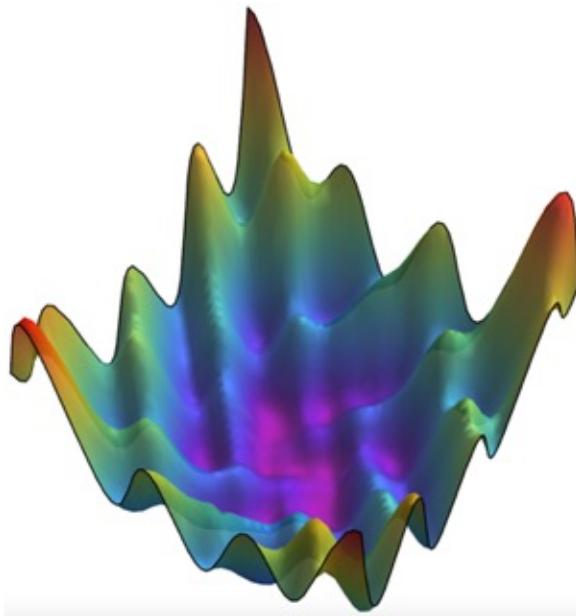
“bad” local minimizers



“flat” saddle points



General Nonconvex Problems



$$\min_{\boldsymbol{x}} f(\boldsymbol{x}), \text{ s.t. } \boldsymbol{x} \in \mathbb{R}^n$$

In the **worst-case**, even finding a local minimizer is NP-hard
(Murty et al. 1987)

Typical Convergence Guarantees

For nonconvex problems, in general we *cannot* hope for efficient global convergence to optimal solutions, but we may have

- guaranteed convergence to critical points (i.e., $\nabla f(\mathbf{x}) = \mathbf{0}$);
- convergence to local minimizers;
- local convergence to global minimizers (with proper initializations).

Convergence to Critical Points

Suppose we are satisfied with any critical point...

This means that our goal is merely to find a point \boldsymbol{x} with

$$\|\nabla f(\boldsymbol{x})\|_2 \leq \varepsilon \quad (\text{called } \varepsilon\text{-approximate critical points})$$

Question: can *vanilla* gradient descent

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \tau_k \cdot \nabla f(\boldsymbol{x}_k)$$

achieve this goal? If so, how fast?

Gradient Descent Converges to Critical Points

Theorem. (Convergence of gradient descent)

Let $f(\mathbf{x})$ be L -smooth function and $\tau_k = \tau \equiv 1/L$.

- In general, gradient descent obeys

$$\min_{0 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|_2 \leq \sqrt{\frac{(f(\mathbf{x}_0) - f(\mathbf{x}_\star))}{k}}.$$

- If $f(\mathbf{x})$ is convex, then gradient descent obeys

$$\min_{k/2 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|_2 \leq \frac{4L \|\mathbf{x}_0 - \mathbf{x}_\star\|}{k}.$$

- The rate cannot be improved for first-order method;
- It does not imply GD converges to critical points.

Lecture Agenda

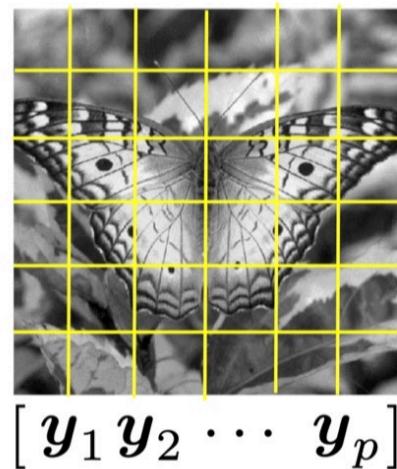
- **Problem Introduction**
 - General Nonconvex Problems
 - **Global Optimization?**
 - Strict Saddle Problems
- Applications in SIPML

Global Optimization of Nonconvex Problems

- Solving general nonconvex problems to global optimality is NP-hard in the worst case;
- However, in practice we often observe many SIPML problems can be efficiently optimized to global solutions.

Example: Dictionary Learning

Given \mathbf{Y} , jointly learn a *compact* dictionary \mathbf{A}_0 and *sparse* \mathbf{X}_0 ?



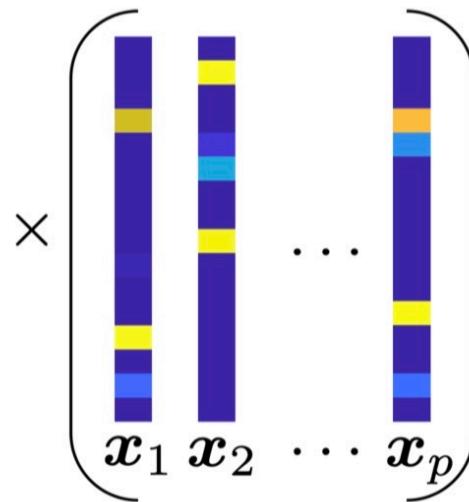
\approx



\mathbf{Y}

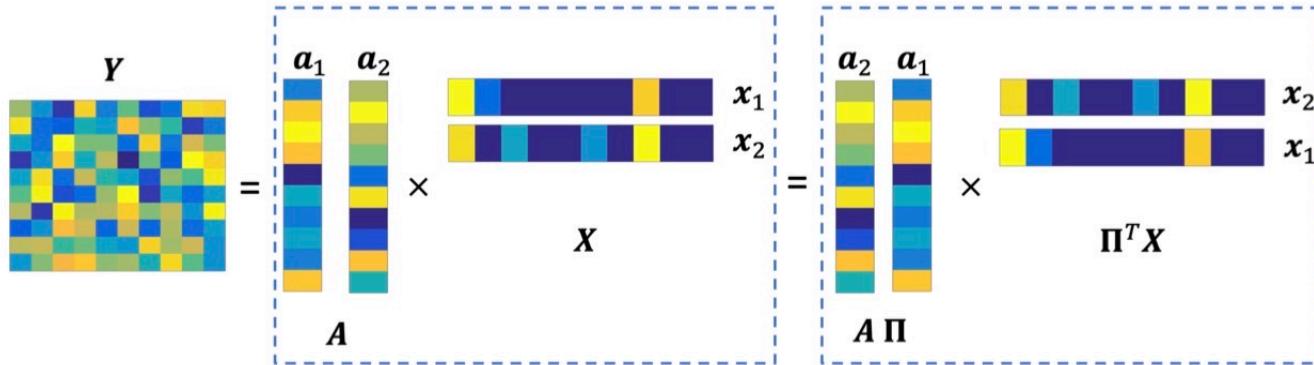
\approx

\mathbf{A}_0



Example: Dictionary Learning

Nonconvexity due to symmetry:

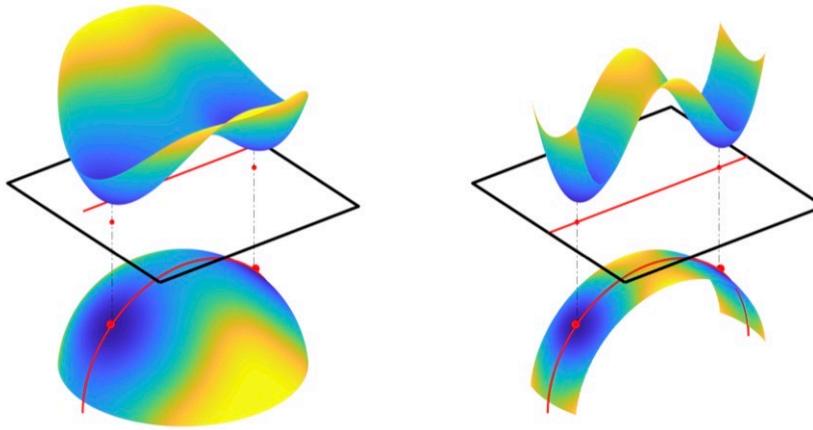


- **Permutation symmetry:** ($2^n n!$ signed permutation Π)

$$Y = A_0 X_0 = (A_0 \Pi) (\Pi^\top X_0)$$

- **Equivalent solution pairs:** $(A_0, X_0) \iff (A_0 \Pi, \Pi^\top X_0)$.

Example: Dictionary Learning



- **Permutation symmetry:** ($2^n n!$ signed permutation Π)
$$Y = A_0 X_0 = (A_0 \Pi) (\Pi^\top X_0)$$
- **Equivalent solution pairs:** $(A_0, X_0) \iff (A_0 \Pi, \Pi^\top X_0)$.

Example: Dictionary Learning

- A natural nonconvex formulation:

$$\min_{\mathbf{A}, \mathbf{X}} f(\mathbf{A}, \mathbf{X}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{AX}\|_F^2 + \lambda \cdot \|\mathbf{X}\|_1, \text{ s.t. } \mathbf{A} \in \mathcal{C}.$$

- A natural algorithmic idea (alternating minimization):

1. optimize \mathbf{X} with \mathbf{A}_k fixed

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

2. optimize \mathbf{A} with \mathbf{X}_{k+1} fixed

$$\mathbf{A}_{k+1} = \arg \min_{\mathbf{A} \in \mathcal{C}} f(\mathbf{A}, \mathbf{X}_{k+1})$$

Example: Dictionary Learning

When the dictionary $\mathbf{A}_0 \in \mathbb{R}^{n \times n}$ is *orthogonal*, i.e.,

$$\mathbf{A}_0 \in \mathcal{O}(n) := \left\{ \mathbf{Q} \in \mathbb{R}^{n \times n} \mid \mathbf{Q}^\top \mathbf{Q} = \mathbf{I} \right\}$$

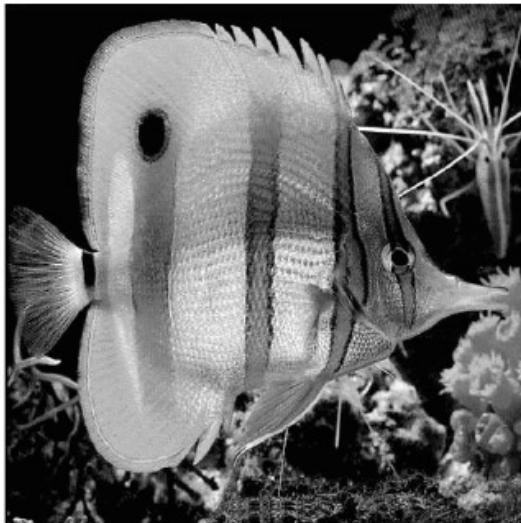
Then we have

$$\mathbf{X}_{k+1} = \arg \min_{\mathbf{X}} f(\mathbf{A}_k, \mathbf{X}) = \text{S}_\lambda [\mathbf{A}_k^\top \mathbf{Y}]$$

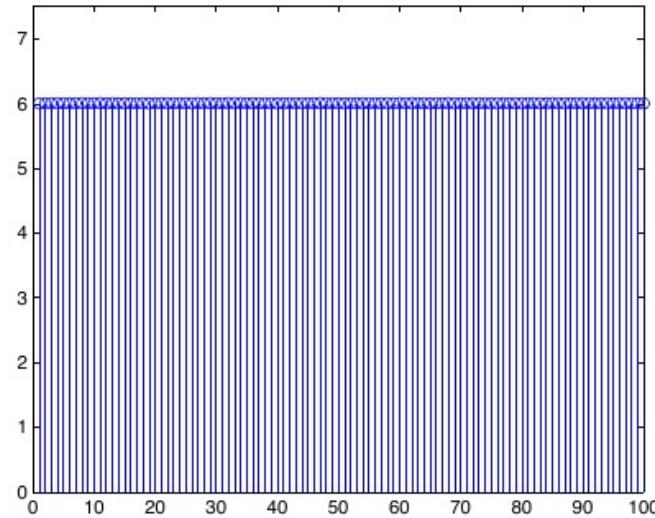
$$\mathbf{A}_{k+1} = \arg \min_{\mathbf{A} \in \mathcal{C}} f(\mathbf{A}, \mathbf{X}_{k+1}) = \mathbf{U} \mathbf{V}^\top,$$

where $\mathbf{U} \Sigma \mathbf{V}^\top = \text{SVD}(\mathbf{Y} \mathbf{X}^\top)$. (Homework)

Example: Dictionary Learning

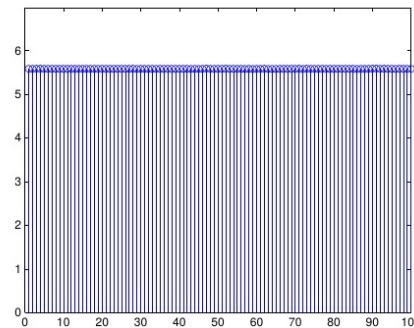
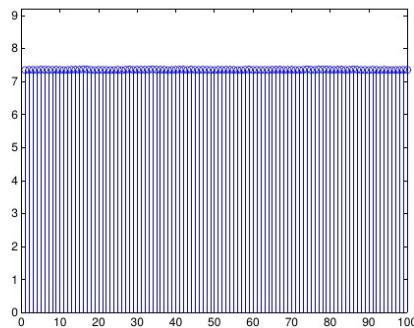
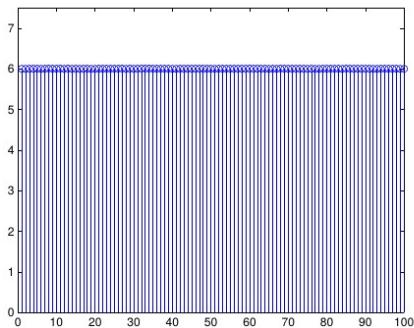
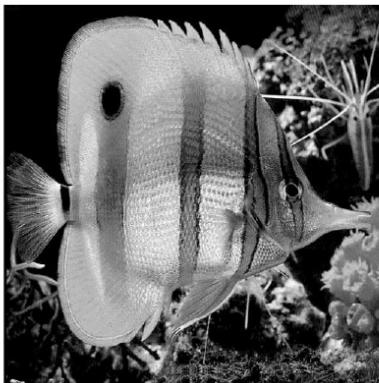


A natural image



$f(A_\infty, X_\infty)$ with random intial.

Example: Dictionary Learning



Further Readings

- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications.* John Wright, Yi Ma. **(Chapter 7)**
- *Complete Dictionary Recovery Over the Sphere I: Overview and the Geometric Picture.* Ju Sun, Qing Qu, and John Wright, IEEE Trans. Info. Theory, 2016.

Lecture Agenda

- **Problem Introduction**
 - General Nonconvex Problems
 - Global Optimization?
 - **Strict Saddle Problems**
- Applications in SIPML

Global Optimization of Nonconvex Problems

- Solving general nonconvex problems to global optimality is NP-hard in the worst case;
- We focus and study on a small class of nonconvex problems, which appears broadly in SIPML, that we can design efficient solvers to global solutions.

A Classical Example: the Rayleigh Quotient

For a given *symmetric* matrix $A^{n \times n}$, solve

$$\min_{\mathbf{x}} -\frac{\mathbf{x}^\top A \mathbf{x}}{\mathbf{x}^\top \mathbf{x}},$$

which is equivalent to

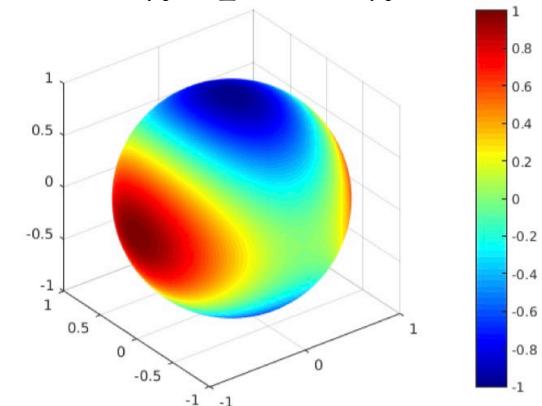
$$\boxed{\min_{\mathbf{x}} f(\mathbf{x}) = -\mathbf{x}^\top A \mathbf{x}, \quad \text{s.t.} \quad \|\mathbf{x}\|_2 = 1.}$$

A Classical Example: the Rayleigh Quotient

$$\min_{\mathbf{x}} f(\mathbf{x}) = -\mathbf{x}^\top \mathbf{A} \mathbf{x}, \quad \text{s.t.} \quad \|\mathbf{x}\|_2 = 1.$$

Let $\mathbf{A} = \mathbf{V}^\top \boldsymbol{\Lambda} \mathbf{V}$, with $\mathbf{V} = [\mathbf{v}_1 \cdots \mathbf{v}_n]$ being the eigenvectors, and $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ with $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_{n-1} > \lambda_n$.

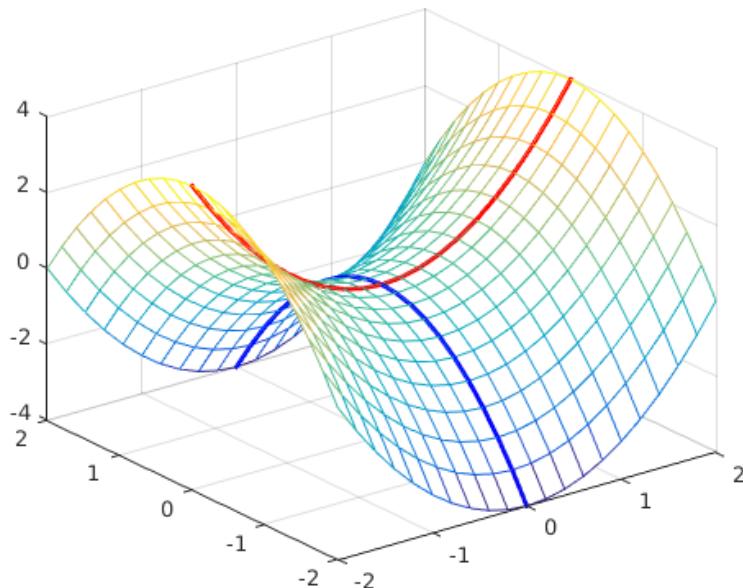
- Only *global minimizers* are $\pm \mathbf{v}_1$;
- Only *global maximizers* are $\pm \mathbf{v}_n$;
- All $\{\pm \mathbf{v}_i\}_{i=1}^n$ for $2 \leq i \leq n-1$ are *saddle points* with *directional negative curvatures*.



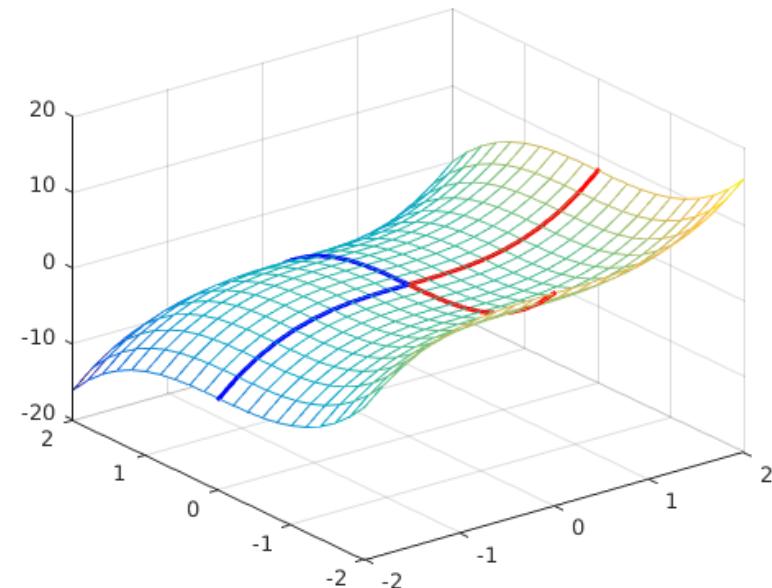
$$\mathbf{A} = \text{diag}(1, 0, -1)$$

Saddle Points with Negative Curvature

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



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



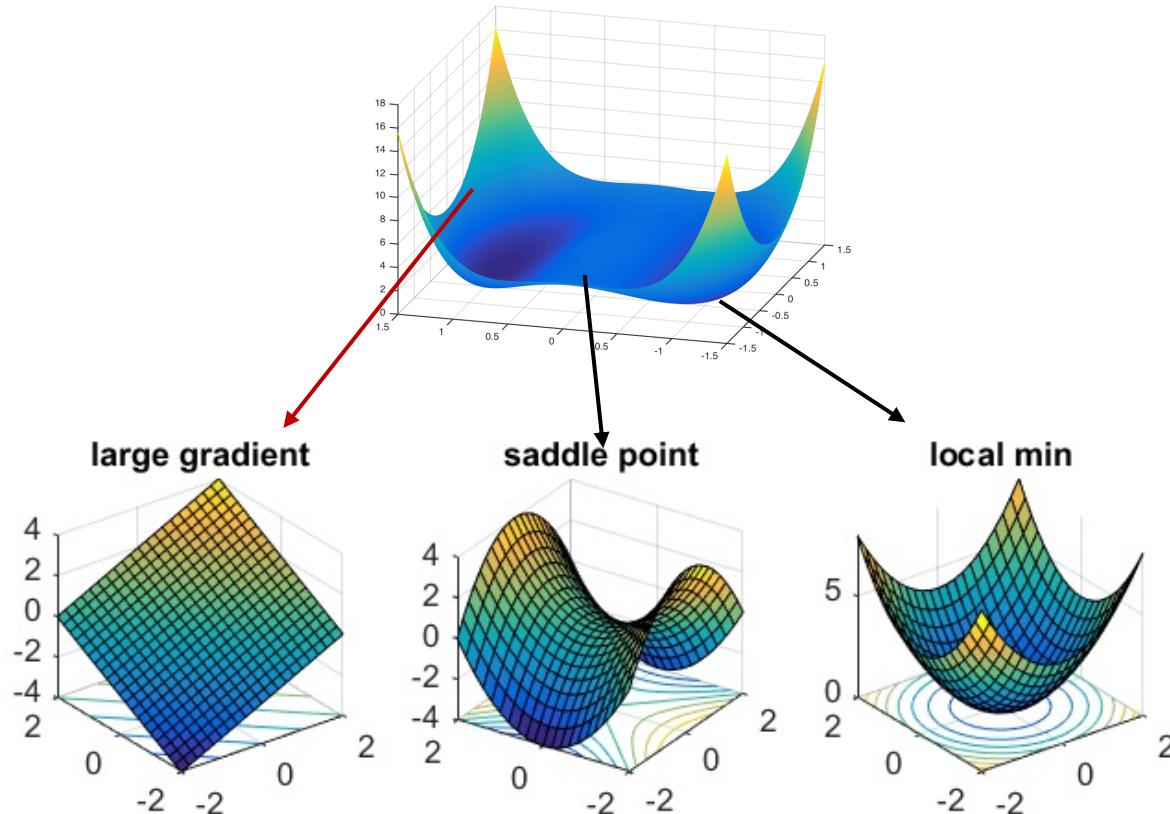
Strict Saddle Function on \mathbb{R}^n

Definition. (Strict Saddle Function in \mathbb{R}^n , Ge et al.'15)

A function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is $(\alpha, \beta, \gamma, \delta)$ -strict saddle, if $\forall x \in \mathbb{R}^n$ obeys *at least one* of the following:

- **[Large gradient]** $\|\nabla f(x)\|_2 \geq \beta$;
- **[Negative curvature]** $\exists v \in \mathbb{S}^{n-1}$, such that
$$v^\top \nabla^2 f(x) v \leq -\alpha;$$
- **[Strong convexity around minimizers]**
 $\exists x_\star$ such that $\|x - x_\star\|_2 \leq \delta$, and for all $y \in \mathcal{B}(x_\star, 2\delta)$, we have $\nabla^2 f(y) \succeq \gamma I$.

Strict Saddle Function on \mathbb{R}^n



Strict Saddle Function on Manifold \mathcal{M}

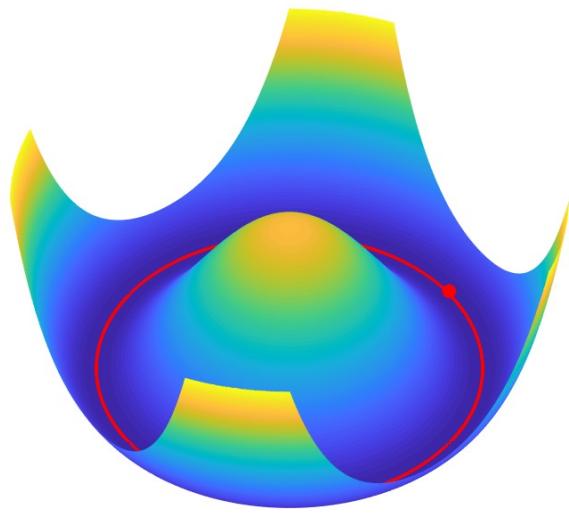
Definition. (Strict Saddle Function in \mathcal{M} , Sun et al.'15)

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

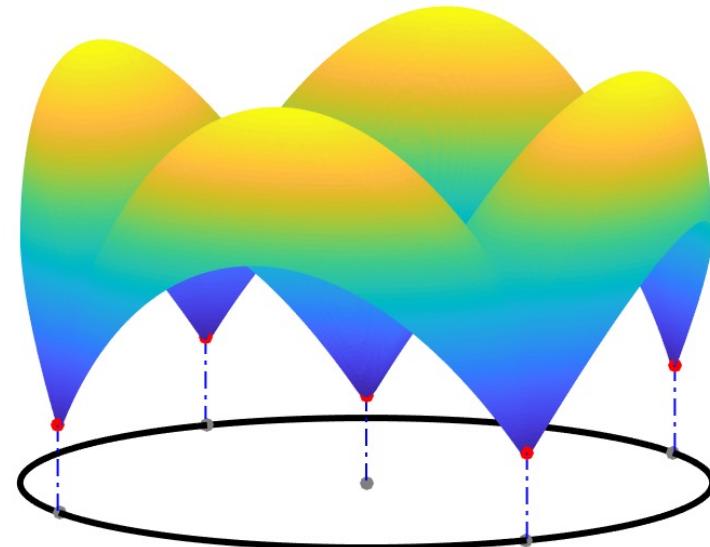
- **[Large gradient]** $\|\text{grad } f(\mathbf{x})\|_2 \geq \beta$;
- **[Negative curvature]** $\exists \mathbf{v} \in T_{\mathbf{x}}\mathcal{M}$ with $\mathbf{v} \in \mathbb{S}^{n-1}$, such that
$$\langle \text{Hess } f(\mathbf{x})[\mathbf{v}], \mathbf{v} \rangle \leq -\alpha$$
 ;
- **[Strong convexity around minimizers]**
 $\exists \mathbf{x}_\star$ such that $\|\mathbf{x} - \mathbf{x}_\star\|_2 \leq \delta$, and for all $\mathbf{y} \in \mathcal{B}(\mathbf{x}_\star, 2\delta) \cap \mathcal{M}$,
we have $\text{Hess } f(\mathbf{y}) \succeq \gamma \mathbf{I}$.

Symmetry Only Creates Equivalent Solutions

Continuous Symmetry

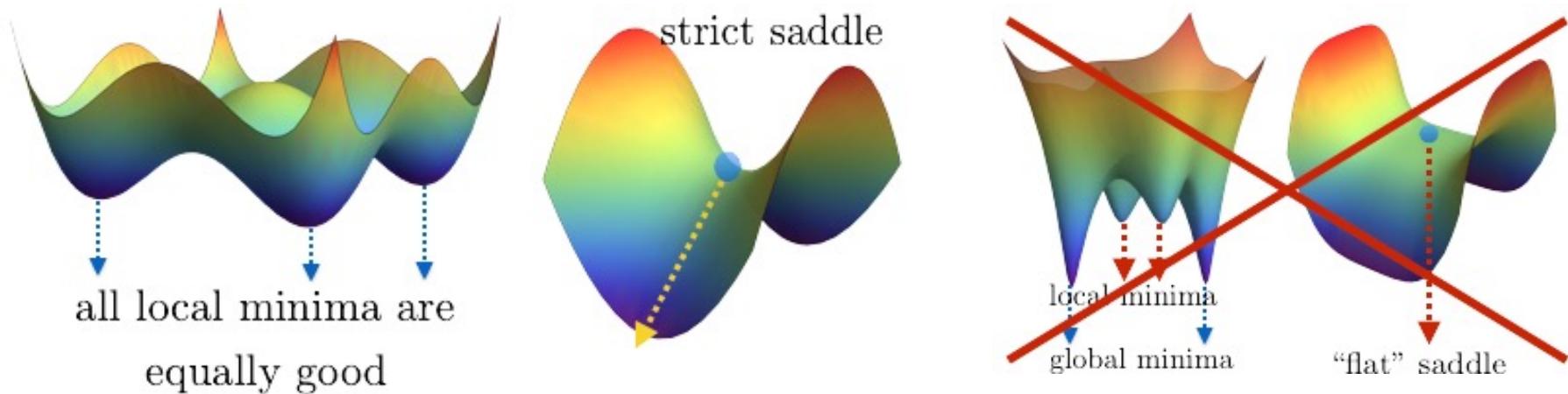


Discrete Symmetry



Symmetry creates *only* equivalent *global* minimizers,
but *not* bad local minima.

Optimizing Nonconvex Problems Globally



Benign nonconvex landscapes enable efficient
global optimization!

Strict Saddle Problems in SIPML

- **Discrete Symmetry**
 - Generalized Phase Retrieval
 - Low-rank Matrix Recovery
 - ...
- **Continuous Symmetry**
 - (Convolutional) Sparse Dictionary Learning
 - Orthogonal Tensor Decomposition
 - Sparse Blind Deconvolution
 - ...

Further Readings

- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications.* John Wright, Yi Ma. (**Chapter 7**)
- *Escaping From Saddle Points --- Online Stochastic Gradient for Tensor Decomposition.* Rong Ge, Furong Huang, Chi Jin, Yang Yuan, COLT'15, 2015.
- *When Are Nonconvex Problems Not Scary?* Ju Sun, Qing Qu, John Wright, NeurIPS Workshop on Nonconvex Optimization, 2015.
- *From Symmetry to Geometry: Tractable Nonconvex Problems.* Yuqian Zhang, Qing Qu, John Wright, 2021.

Lecture Agenda

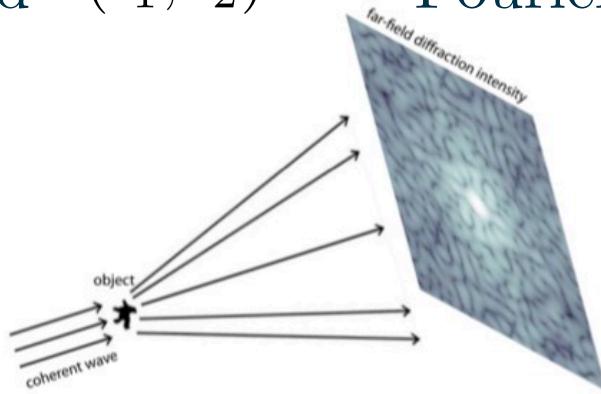
- Problem Introduction
- Applications in SIPML
 - Example I: Generalized Phase Retrieval
 - Example II: Low-Rank Matrix Recovery
 - Example III: Training Deep Neural Networks
 - Example IV: Sparse Dictionary Learning
 - Example V: Sparse Blind Deconvolution

Lecture Agenda

- Problem Introduction
- Applications in SIPML
 - Example I: Generalized Phase Retrieval
 - Example II: Low-Rank Matrix Recovery
 - Example III: Training Deep Neural Networks
 - Example IV: Sparse Dictionary Learning
 - Example V: Sparse Blind Deconvolution

Fourier Phase Retrieval

electric field $x(t_1, t_2) \rightarrow$ Fourier transform $\hat{x}(f_1, f_2)$



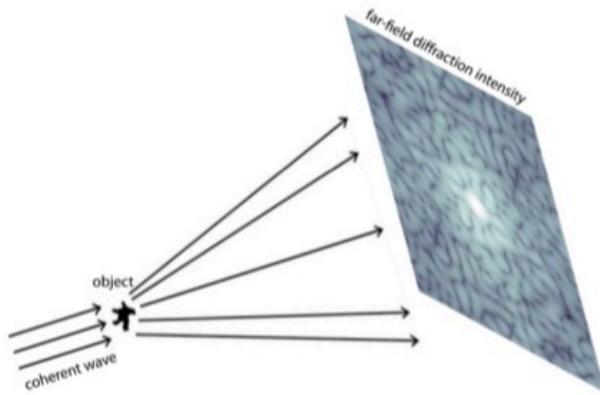
Applications: X-ray crystallography, **diffraction imaging** (left), optics, astronomical imaging, and microscopy

Coherent Diffraction Imaging¹

Fourier phase retrieval: given intensity $y = |\mathcal{F}(x_\star)|$, recover x_\star .

- Fourier phase retrieval is a very challenging ill-posed problem, which is notoriously difficult to solve.

Generalized Phase Retrieval



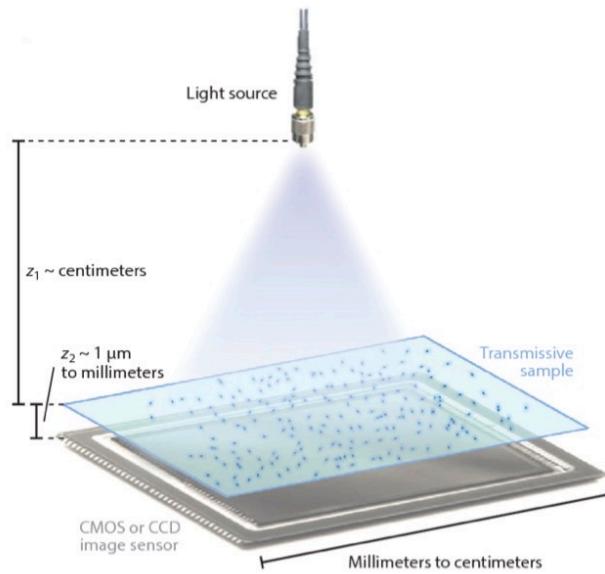
Coherent Diffraction Imaging¹

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

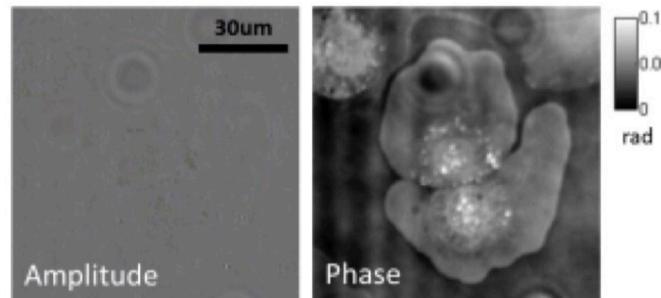
Generalized phase retrieval: given intensity $y = |\mathbf{Ax}_\star|$, recover $\mathbf{x}_\star \in \mathbb{C}^m$.

- The sensing matrix \mathbf{A} can be generic and less structured, making the problem easier to solve.

Applications in Microscopy Imaging

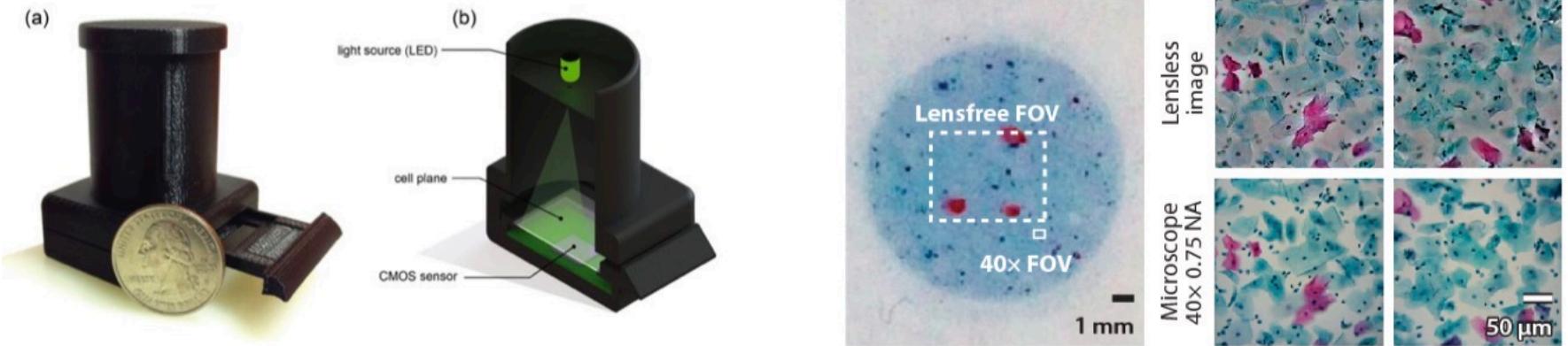


Lensfree microscopy



Phase contrast microscopy:
phase is important¹.

Applications in Microscopy Imaging



Lensfree microscopy imaging via phase retrieval.

Generalized Phase Retrieval

$$\begin{array}{c}
 A \\
 \times \\
 x \\
 = \\
 \hline
 Ax \\
 \rightarrow \\
 y = |Ax|^2
 \end{array}$$

- Solve for $x \in \mathbb{C}^n$ in m quadratic equations

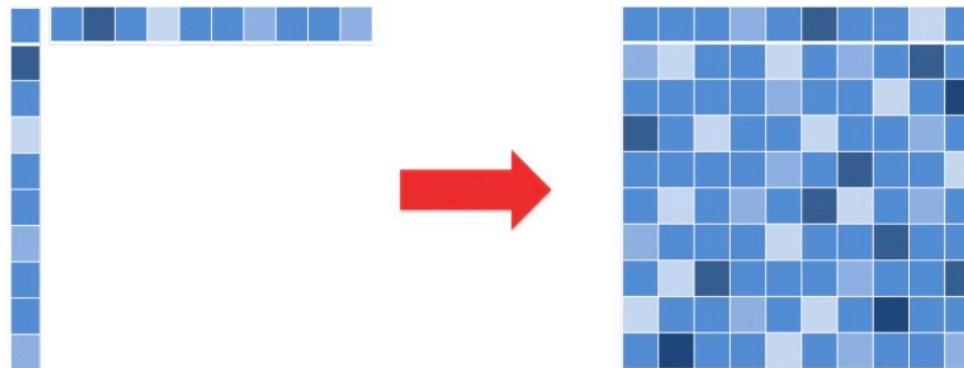
$$y_k = |a_k^\top x|^2, \quad k = 1, \dots, m,$$

or $\mathbf{y} = |\mathbf{Ax}|^2$, where $|\mathbf{z}|^2 := \left[|z_1|^2, \dots, |z_m|^2\right]^\top$.

Convex Relaxation (Lecture 6)

- **Lifting:** introduce $X = xx^*$ to *linearize* the problem

$$y_k = |a_k^* x|^2 = a_k^* \underbrace{(xx^*)}_{X} a_k \implies y_k = \langle a_k a_k^*, X \rangle$$



Convex Relaxation (Lecture 6)

$$\begin{aligned} \text{find } \mathbf{X} \succeq \mathbf{0}, \quad \text{s.t.} \quad y_k &= \langle \mathbf{a}_k \mathbf{a}_k^*, \mathbf{X} \rangle, \quad k = 1, \dots, m \\ \text{rank}(\mathbf{X}) &= 1 \end{aligned}$$

- Convex Relaxation:

$$\min_{\mathbf{X}} \|\mathbf{X}\|_*, \quad \text{s.t.} \quad \mathbf{y} = \mathcal{A}(\mathbf{X}), \quad \mathbf{X} \succeq \mathbf{0}.$$

Needs to optimize a problem of $O(n^2)$ variables instead of $O(n)$, and SDP problems are usually very expensive.

A Natural Nonconvex Formulation

Given $\mathbf{y} \in \mathbb{R}^m$ generated by $\mathbf{y} = |\mathbf{A}\mathbf{x}_\star|$, recover $\mathbf{x}_\star \in \mathbb{C}^n$.

- Optimizing the following *nonconvex* least-squares loss:

$$\min_{\mathbf{x} \in \mathbb{C}^n} f(\mathbf{x}) = \frac{1}{4m} \sum_{i=1}^m \left(y_i^2 - |\mathbf{a}_i^* \mathbf{x}|^2 \right)^2$$

It only involves $O(n)$ variables, but the problem is highly *nonconvex* (**why?**)

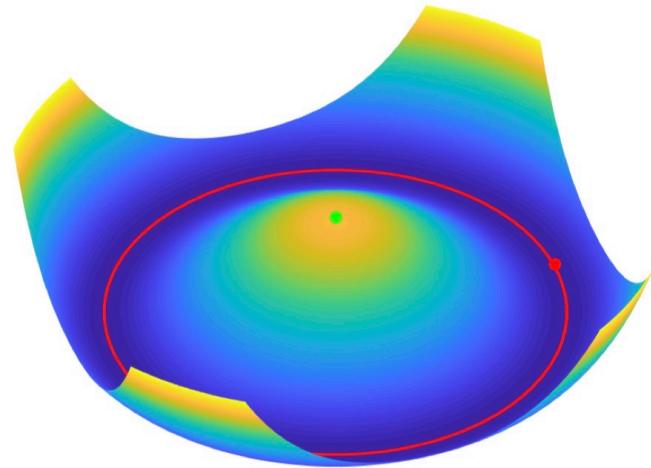
Symmetries in Generalized Phase Retrieval

For any $\phi \in [0, 2\pi)$, we observe

$$\mathbf{y} = |\mathbf{A}\mathbf{x}_\star| = |\mathbf{A}\mathbf{x}_\star e^{i\phi}|,$$

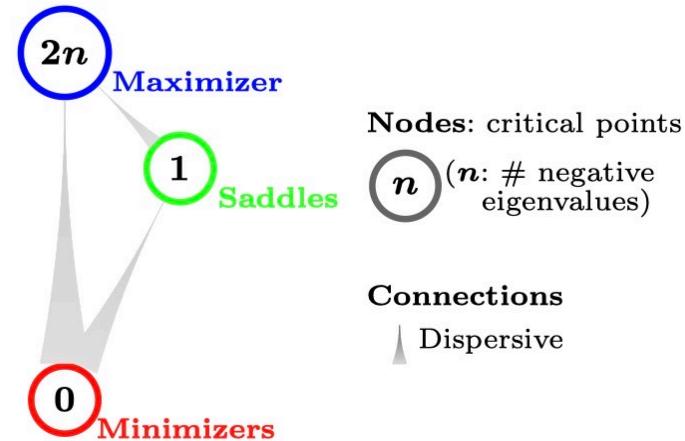
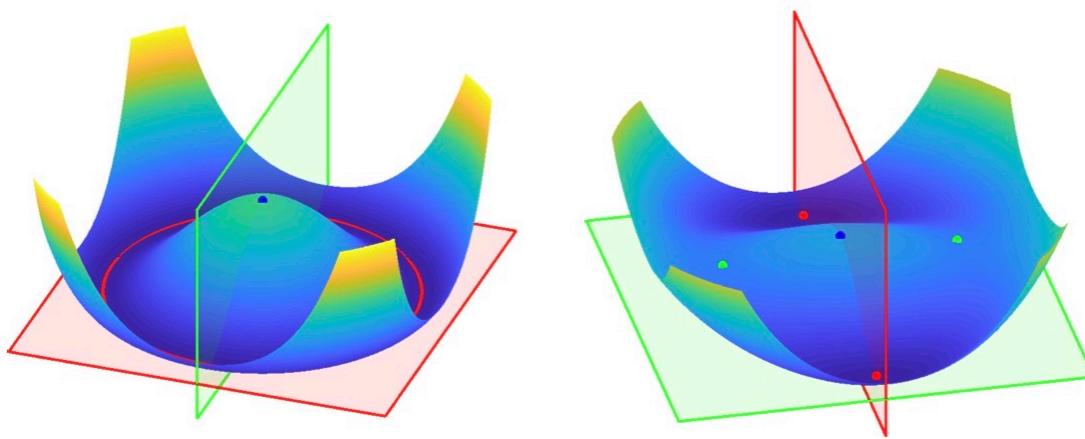
so that

$$\mathcal{X} := \{\mathbf{z} \in \mathbb{C}^n \mid \mathbf{z} = \mathbf{x}_\star e^{i\phi}, \phi \in [0, 2\pi)\}.$$



Plot of the landscape
with a single unknown
 $\mathbf{y} = \mathbf{a}\mathbf{x}_\star.$

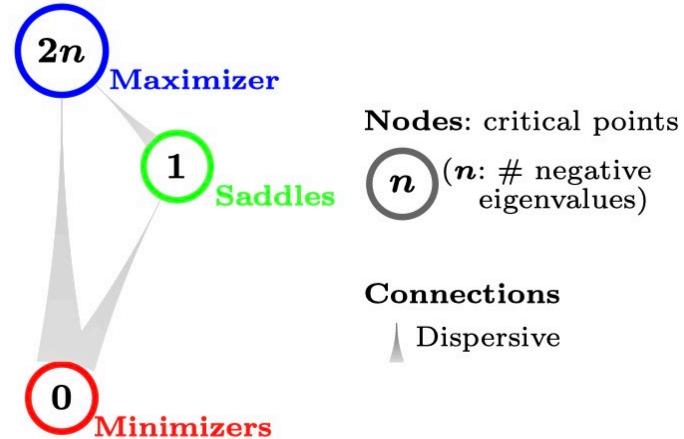
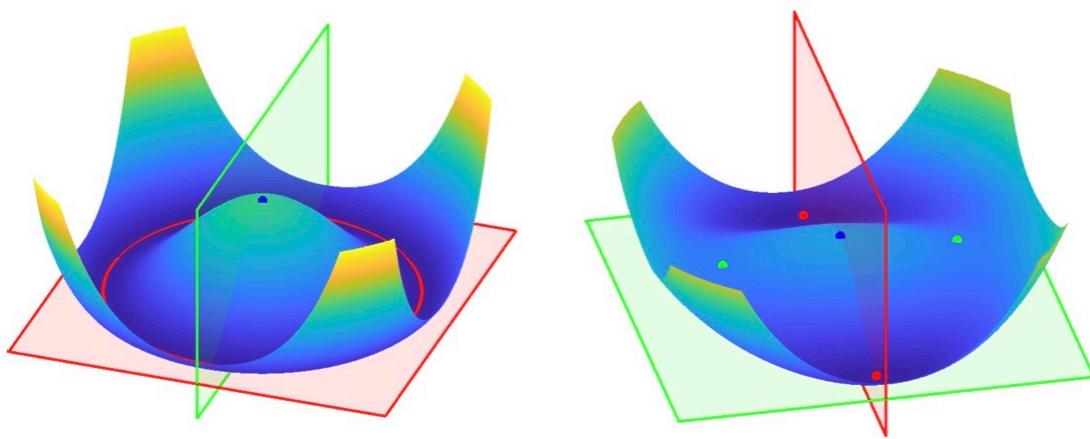
A Natural Nonconvex Formulation



- Optimizing the following *nonconvex* least-squares loss:

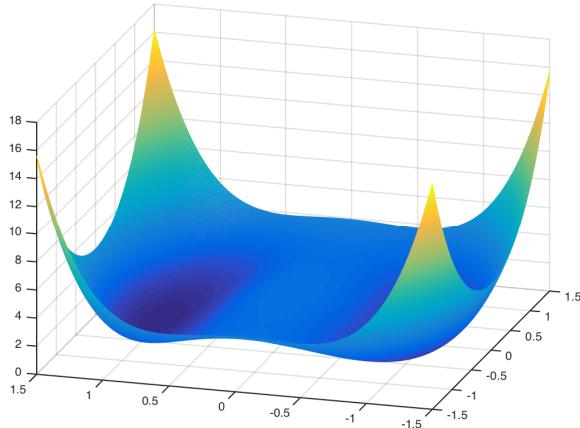
$$\min_{\mathbf{x} \in \mathbb{C}^n} f(\mathbf{x}) = \frac{1}{4m} \sum_{i=1}^m \left(y_i^2 - |\mathbf{a}_i^* \mathbf{x}|^2 \right)^2$$

Generalized Phase Retrieval



- Symmetric copies of the ground truth are minimizers.
- Negative curvature in symmetry-breaking directions.
- Cascade of saddle points.

Generalized Phase Retrieval



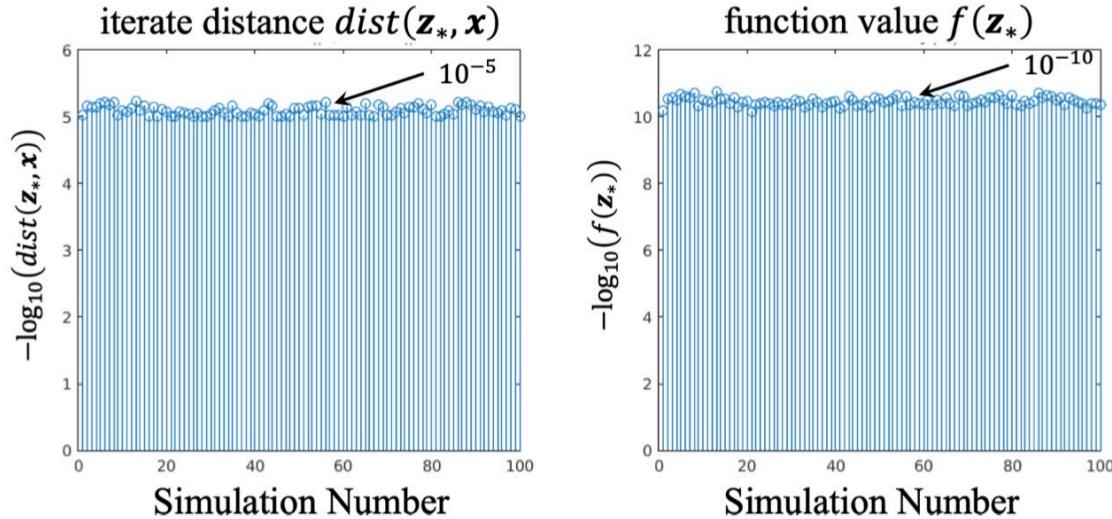
$$\min_{\mathbf{x} \in \mathbb{C}^n} f(\mathbf{x}) = \frac{1}{4m} \sum_{i=1}^m \left(y_i^2 - |\mathbf{a}_i^* \mathbf{x}|^2 \right)^2$$

Theorem. (Informal, Sun et al.'18)

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

- All local (and global) minimizers are of the form $\mathbf{x}e^{i\phi}$, $\phi \in [0, 2\pi)$;
- $f(\mathbf{x})$ is $(c, c/(n \log m), c, c/(n \log m))$ -strict saddle for some $c > 0$.

Generalized Phase Retrieval



Run vanilla gradient descent, with random initializations

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau_k \cdot \nabla f(\mathbf{x}_k)$$

Further Readings

- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications.* John Wright, Yi Ma. **(Chapter 7)**
- *Phase Retrieval via Wirtinger Flow: Theory and Algorithms.* Emmanuel Candes, Xiaodong Li, Mahdi Soltanolkotabi, IEEE Info. Theory, 2014.
- *A Geometric Analysis of Phase Retrieval.* Ju Sun, Qing Qu, John Wright, Foundations of Computational Math, 2018.
- *Phase Retrieval with Application to Optical Imaging: A Contemporary Overview.* Yoav Shechtman, Yonina C. Eldar, Oren Cohen, Henry N. Chapman, Jianwei Miao, Mordechai Segev, IEEE Signal Processing Magazine, 2015.

Lecture Agenda

- Problem Introduction
- Applications in SIPML
 - Example I: Generalized Phase Retrieval
 - Example II: Low-Rank Matrix Recovery
 - Example III: Training Deep Neural Networks
 - Example IV: Sparse Dictionary Learning
 - Example V: Sparse Blind Deconvolution

Low-rank Matrix Recovery

Given $\mathbf{y} = \mathcal{A}(\mathbf{X}_\star)$, recover a rank- r (with $r \ll \min\{n_1, n_2\}$) matrix $\mathbf{X}_\star \in \mathbb{R}^{n_1 \times n_2}$ from $\mathbf{y} \in \mathbb{R}^m$.

- Convex relaxation approaches:

$$\min_{\mathbf{X}} \|\mathbf{X}\|_* \quad \text{s.t.} \quad \mathbf{y} = \mathcal{A}(\mathbf{X}),$$

or in the noisy setting:

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2 + \lambda \|\mathbf{X}\|_*.$$

The nuclear minimization problem with $O(n_1 n_2)$ variables could be very expensive (computing SVD $O(n_1^2 n_2)$).

Low-rank Matrix Recovery

Given $\mathbf{y} = \mathcal{A}(\mathbf{X}_\star)$, recover a rank- r (with $r \ll \min\{n_1, n_2\}$) matrix $\mathbf{X}_\star \in \mathbb{R}^{n_1 \times n_2}$ from $\mathbf{y} \in \mathbb{R}^m$.

- **A nonconvex approach:** let $\mathbf{X} = \mathbf{U}\mathbf{V}^\top$ with $\mathbf{U} \in \mathbb{R}^{n_1 \times r}$ and $\mathbf{V} \in \mathbb{R}^{n_2 \times r}$,

$$\min_{\mathbf{U}, \mathbf{V}} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{V}^\top)\|_2^2 + \rho(\mathbf{U}, \mathbf{V})$$

- The problem is nonconvex due to bilinearity $\mathbf{U}\mathbf{V}^\top$.
- The number of optimization variables $O(\max\{n_1, n_2\}r)$ is much smaller than $O(n_1 n_2)$.

Symmetries in Matrix Recovery

For any *invertible* matrix $\Gamma \in \mathbb{R}^{r \times r}$,

$$\mathbf{X}_\star = \mathbf{U}_\star \mathbf{V}_\star^\top = \mathbf{U}_\star \Gamma \Gamma^{-1} \mathbf{V}_\star^\top = (\mathbf{U}_\star \Gamma) (\mathbf{V}_\star \Gamma^{-\top})^\top.$$

so that the problem possess a *general linear* symmetry

$$(\mathbf{U}_\star, \mathbf{V}_\star) \equiv (\mathbf{U}_\star \Gamma, \mathbf{V}_\star \Gamma^{-\top}), \quad \forall \Gamma \in \text{GL}(r)$$

- The matrix Γ can be *arbitrarily ill-conditioned*, so that \mathbf{U} and \mathbf{V} can be *arbitrarily unbalanced*.
- It can be reduced to a simpler orthogonal symmetry $\mathcal{O}(r)$ by using a penalty such as $\rho(\mathbf{U}, \mathbf{V}) = \|UU^\top - \mathbf{V}\mathbf{V}^\top\|_F^2$
- The landscape is “almost” the same with factorization problem.

Geometry of Symmetric Matrix Factorization

- Study of the nonconvex landscape:

Let us start with a much simpler *symmetric* matrix factorization problem $\mathbf{Y} = \mathbf{X}_\star$, and we aim to factor it as $\mathbf{X}_\star = \mathbf{U}_\star \mathbf{U}_\star^\top$ by

$$\min_{\mathbf{U} \in \mathbb{R}^{n \times r}} \varphi(\mathbf{U}) := \frac{1}{4} \|\mathbf{Y} - \mathbf{U}\mathbf{U}^\top\|_F^2$$

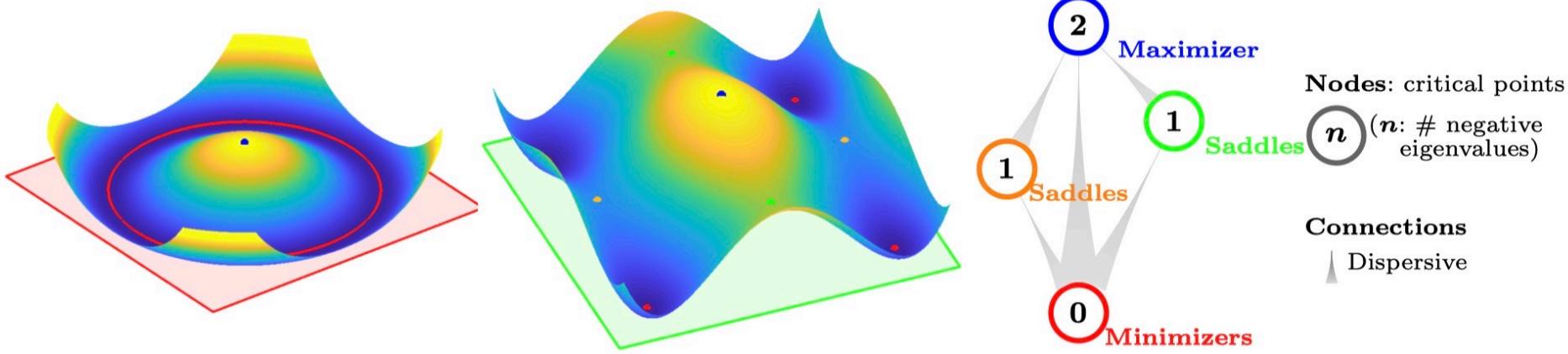
- The problem exhibits a simpler orthogonal symmetry

$$\mathbf{U}\mathbf{U}^\top = \mathbf{U}\boldsymbol{\Gamma}\boldsymbol{\Gamma}^\top\mathbf{U}^\top = (\mathbf{U}\boldsymbol{\Gamma})(\mathbf{U}\boldsymbol{\Gamma})^\top, \quad \forall \boldsymbol{\Gamma} \in \mathcal{O}(r),$$

so that $\varphi(\mathbf{U}) \equiv \varphi(\mathbf{U}\boldsymbol{\Gamma})$.

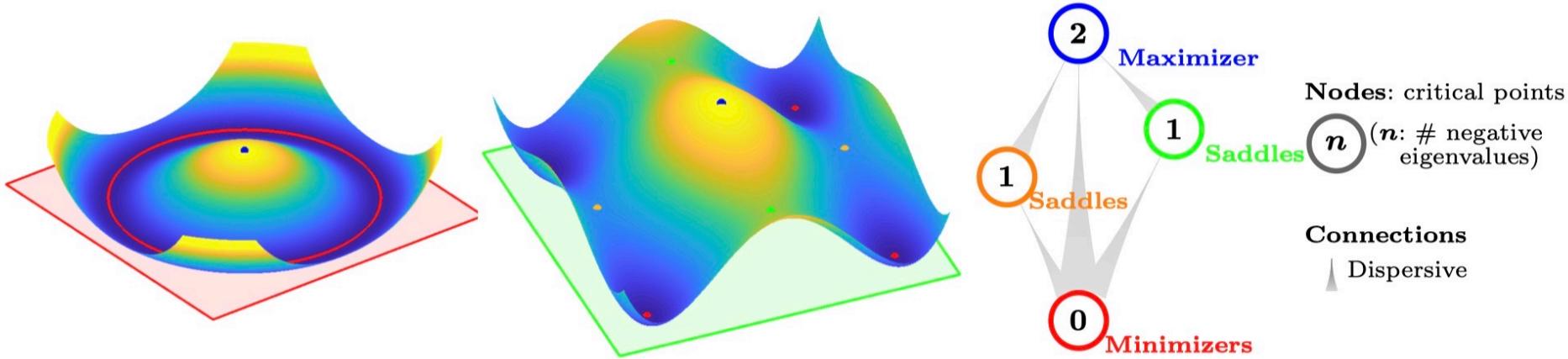
Geometry of Symmetric Matrix Factorization

$$\min_{U \in \mathbb{R}^{n \times r}} \varphi(U) := \frac{1}{4} \|Y - UU^\top\|_F^2$$



X_\star is a rank-2 symmetric matrix, with $\lambda_1 = \frac{3}{4}$ and $\lambda_2 = \frac{1}{2}$.

Geometry of Symmetric Matrix Factorization



- Symmetric copies of the ground truth are minimizers.
- Negative curvature in symmetry-breaking directions.
- Cascade of saddle points.

From Symmetric to Non-symmetric Factorization

This model geometry carries over to *non-symmetric* matrices,

$$\min_{\mathbf{U} \in \mathbb{R}^{n_1 \times r}, \mathbf{V} \in \mathbb{R}^{n_2 \times r}} \frac{1}{2} \left\| \mathbf{Y} - \mathbf{U}\mathbf{V}^\top \right\|_F^2 + \rho(\mathbf{U}, \mathbf{V})$$

by penalization such as $\rho(\mathbf{U}, \mathbf{V}) = \left\| \mathbf{U}\mathbf{U}^\top - \mathbf{V}\mathbf{V}^\top \right\|_F^2$, that we can reduce $\text{GL}(r)$ to $\mathcal{O}(r)$ symmetry.

From Factorization to Matrix Recovery & Completion

Given $\mathbf{y} = \mathcal{A}(\mathbf{X}_\star)$ with $y_i = \langle \mathbf{A}_i, \mathbf{X}_\star \rangle$ ($1 \leq i \leq m$), optimize

$$\min_{\mathbf{U} \in \mathbb{R}^{n_1 \times r}, \mathbf{V} \in \mathbb{R}^{n_2 \times r}} \frac{1}{2} \left\| \mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{V}^\top) \right\|_2^2 + \rho(\mathbf{U}, \mathbf{V})$$

- **Matrix factorization:** $\mathcal{A} = \mathcal{I}$, with \mathcal{I} being an identity mapping.
- **Matrix sensing:** $\mathcal{A}(\cdot)$ satisfies certain rank RIP property:

$$\left| \|\mathcal{A}(\mathbf{X}_r)\|_2^2 - \|\mathbf{X}_r\|_F^2 \right| \leq \delta \|\mathbf{X}_r\|_F^2, \quad \forall \text{rank}(\mathbf{X}_r) \leq r.$$

It can be showed that the objective is strict saddle with benign global landscape.

From Factorization to Matrix Recovery & Completion

Given $\mathbf{y} = \mathcal{A}(\mathbf{X}_\star)$ with $y_i = \langle \mathbf{A}_i, \mathbf{X}_\star \rangle$ ($1 \leq i \leq m$), optimize

$$\min_{\mathbf{U} \in \mathbb{R}^{n_1 \times r}, \mathbf{V} \in \mathbb{R}^{n_2 \times r}} \frac{1}{2} \left\| \mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{V}^\top) \right\|_2^2 + \rho(\mathbf{U}, \mathbf{V})$$

- **Matrix completion:** $\mathcal{A}(\cdot) = \mathcal{P}_\Omega(\cdot)$, where $\mathcal{P}_\Omega(\cdot)$ is not RIP in general, that extra assumption on \mathbf{X}_\star (incoherence) and extra penalization is needed for global benign landscape.

$$\rho_{mc}(\mathbf{U}, \mathbf{V}) = \lambda_1 \sum_{i=1}^{n_1} (\|\mathbf{e}_i^\top \mathbf{U}\|_2 - \alpha_1)_+^4 + \lambda_2 \sum_{i=1}^{n_2} (\|\mathbf{e}_i^\top \mathbf{U}\|_2 - \alpha_2)_+^4.$$

From Factorization to Robust Matrix Recovery

Given $\mathbf{Y} = \mathcal{A}(\mathbf{X}_\star) + \mathbf{S}_\star$ with *low-rank* \mathbf{X}_\star and *sparse* \mathbf{S}_\star ,
optimize

$$\min_{\mathbf{U}, \mathbf{V}, \mathbf{S}} \frac{1}{2} \left\| \mathbf{Y} - \mathcal{A}(\mathbf{U}\mathbf{V}^\top) - \mathbf{S} \right\|_F^2 + \mu \|\mathbf{S}\|_1 + \rho_r(\mathbf{U}, \mathbf{V})$$

- **Robust matrix recovery:** partial minimization w.r.t. \mathbf{S} gives

$$\min_{\mathbf{U}, \mathbf{V}, \mathbf{S}} \frac{1}{2} H_\mu \left(\mathbf{Y} - \mathcal{A}(\mathbf{U}\mathbf{V}^\top) \right) + \rho_r(\mathbf{U}, \mathbf{V}),$$

with $H_\mu(\cdot)$ being Huber loss. Thus, we can perform similar analysis to problems studied previously.

Further Readings

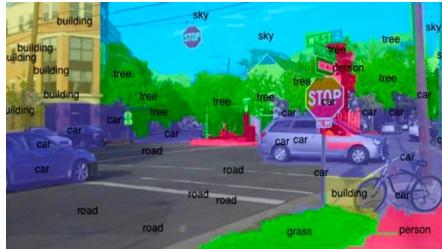
- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications.* John Wright, Yi Ma. **(Chapter 7)**
- *Matrix Completion has No Spurious Local Minimum.* Rong Ge, Jason D. Lee, Tengyu Ma, NeurIPS'16, 2016.
- *No Spurious Local Minima in Nonconvex Low Rank Problems: A Unified Geometric Analysis.* Rong Ge, Chi Jin, Yi Zheng, ICML 2017.

Lecture Agenda

- Problem Introduction
- Applications in SIPML
 - Example I: Generalized Phase Retrieval
 - Example II: Low-Rank Matrix Recovery
 - Example III: Training Deep Neural Networks
 - Example IV: Sparse Dictionary Learning
 - Example V: Sparse Blind Deconvolution

Deep Neural Networks

Deep learning has attained superior performances for many tasks in practice:



Computer vision



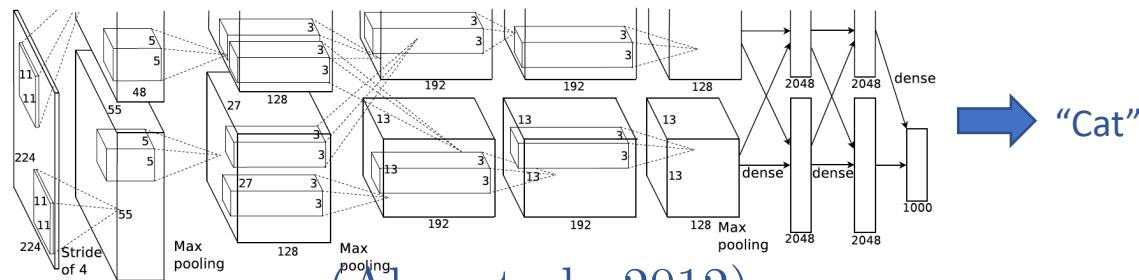
Natural language processing



Gameplay

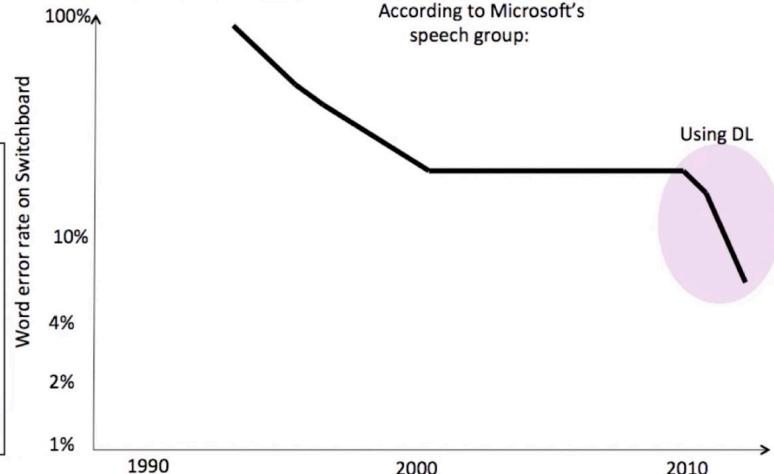
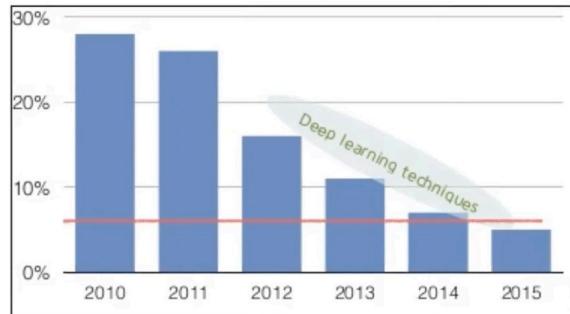


Protein modeling



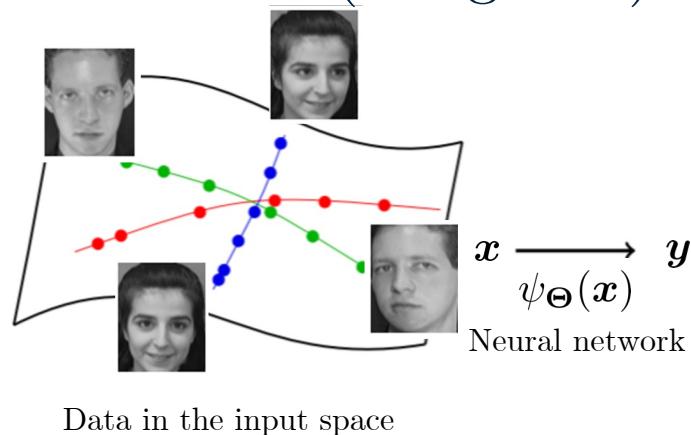
(Alex et al., 2012)

Impact of Deep Learning in ML



Multi-Class Classification Problem

- Learning predictor from **training set** $\{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$
- Example: multi-class classification problem
 - $K = 10$ classes (MNIST, CIFAR10, etc.)
 - $K = 1000$ classes (ImageNet)

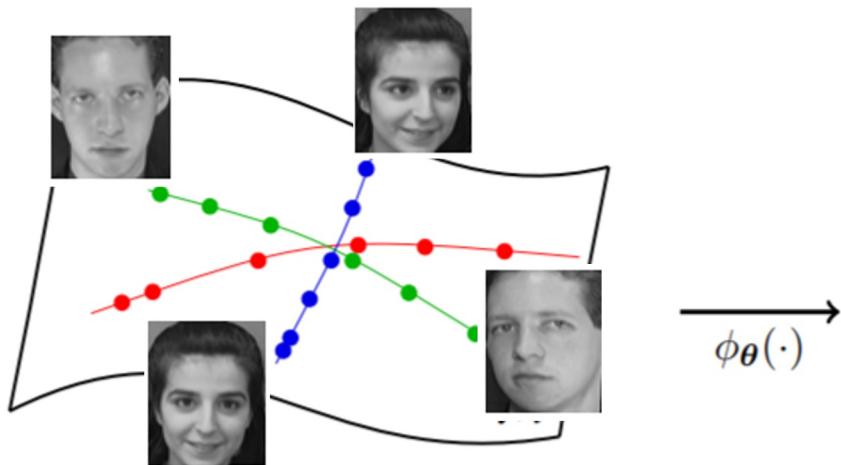


$$\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \cdots \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

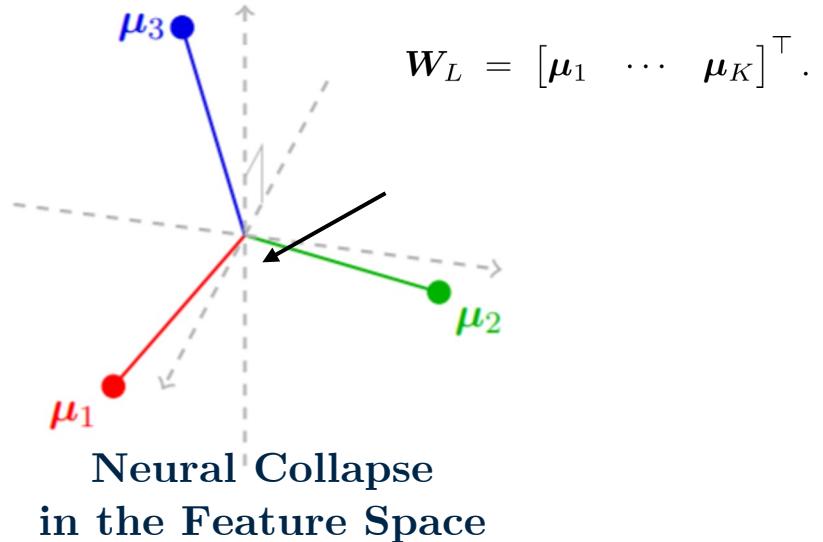
One-hot vectors in \mathbb{R}^K

Deep Neural Network Classifiers

$$\psi_{\Theta}(x) = W_L \underbrace{\sigma(W_{L-1} \cdots \sigma(W_1 x + b_1) + b_{L-1})}_{\text{Last-layer classifier}} + b_L$$
$$\phi_{\theta}(x) =: h \leftarrow \underbrace{\phi_{\theta}(x)}_{\text{Last-layer feature}}$$

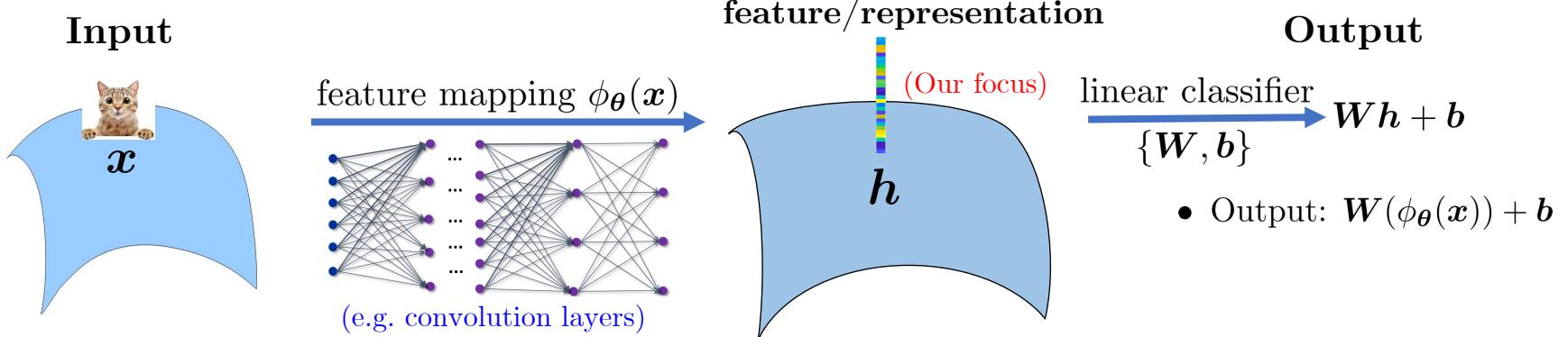


Data in the Input Space



Neural Collapse
in the Feature Space

Deep Neural Network Classifiers



- **Output:** $W \cdot \phi_{\theta}(x) + b$
- **Training problem:** assume balanced dataset with each class n training samples

$$\min_{\theta, W, b} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{CE}(\mathbf{W}(\phi_{\theta}(x_{k,i})) + \mathbf{b}, \mathbf{y}_k) + \lambda \|\mathbf{\theta}, \mathbf{W}, \mathbf{b}\|_F^2$$

Annotations for the equation:

- cross-entropy loss: Points to the cross-entropy loss term \mathcal{L}_{CE} .
- i-th input in the k-th class: Points to the input $x_{k,i}$.
- One-hot vector for the k-th class: Points to the target \mathbf{y}_k .
- Weight decay: Points to the regularization term $\lambda \|\mathbf{\theta}, \mathbf{W}, \mathbf{b}\|_F^2$.

Neural Collapse in Classification

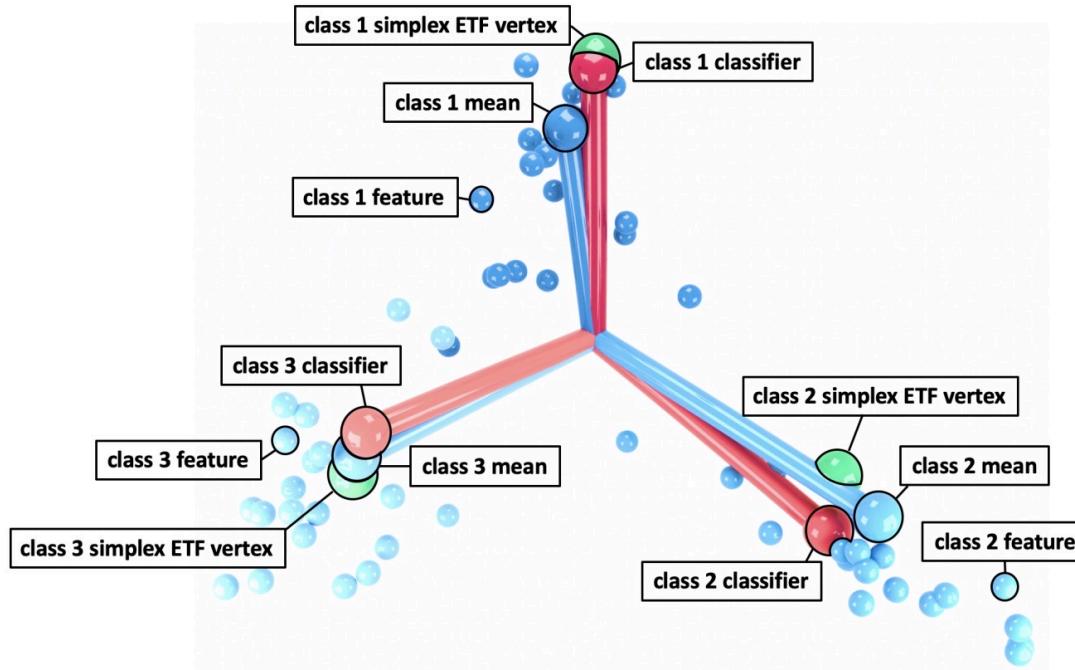


Image credited to Han et al. "Neural Collapse Under MSE Loss: Proximity to and Dynamics on the Central Path"

Neural Collapse in Classification

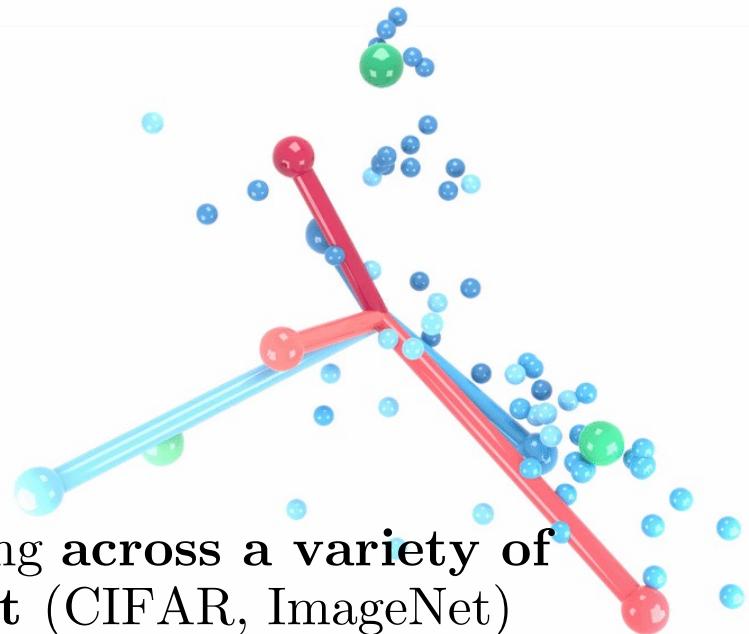
Prevalence of neural collapse during the terminal phase of deep learning training

 Vardan Papyan,  X. Y. Han, and David L. Donoho

[+ See all authors and affiliations](#)

PNAS October 6, 2020 117 (40) 24652-24663; first published September 21, 2020;
<https://doi.org/10.1073/pnas.2015509117>

Contributed by David L. Donoho, August 18, 2020 (sent for review July 22, 2020; reviewed by Helmut Boelschke and Stéphane Mallat)



- Reveals common outcome of network training **across a variety of architectures** (ResNet, VGG) and **dataset** (CIFAR, ImageNet)
- **Precise mathematical structures** within the features and classifier

Neural Collapse: Symmetry and Structures

Balanced training dataset with $n = n_1 = n_2 = \dots = n_K$, and

$$\mathbf{W} := \mathbf{W}_L, \quad \mathbf{H} := [\mathbf{h}_{1,1} \quad \dots \quad \mathbf{h}_{K,n}].$$

Neural Collapse (NC) means that

- 1) *Within-Class Variability Collapse on \mathbf{H} :* features of each class collapse to class-mean with **zero** variability;

$$\mathbf{h}_{k,i} \rightarrow \bar{\mathbf{h}}_k, \quad \forall k \in [K], i \in [n].$$

- 2) *Convergence to Simplex ETF on \mathbf{H} :* the class means are **linearly separable**, and **maximally distant**;

$$\overline{\mathbf{H}}^\top \overline{\mathbf{H}} \sim \mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top, \quad \overline{\mathbf{H}} = [\bar{\mathbf{h}}_1 \quad \dots \quad \bar{\mathbf{h}}_K]$$

Neural Collapse: Symmetry and Structures

Balanced training dataset with $n = n_1 = n_2 = \dots = n_K$, and

$$\mathbf{W} := \mathbf{W}_L, \quad \mathbf{H} := [\mathbf{h}_{1,1} \quad \dots \quad \mathbf{h}_{K,n}].$$

Neural Collapse (NC) means that

- 3) *Convergence to Self-Duality (\mathbf{W}, \mathbf{H}):* the last-layer classifiers are **perfected matched** with the class-means of features.

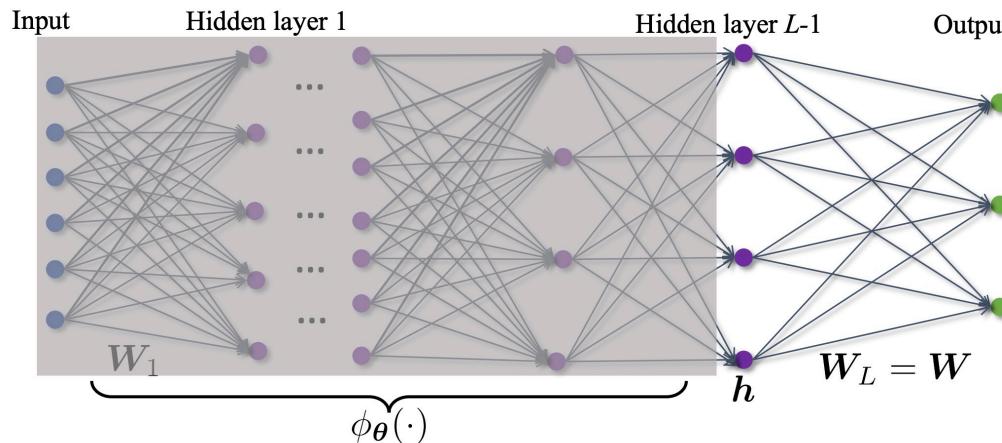
$$\mathbf{w}^k = \beta \bar{\mathbf{h}}_k, \quad \forall k \in [K].$$

- 4) *Simple Decision Rule* via Nearest Class-Center decision.

Simplification: Unconstrained Features

$$\psi_{\Theta}(x) = \underbrace{W_L \sigma(W_{L-1} \cdots \sigma(W_1 x + b_1) + b_{L-1}) + b_L}_{\text{Last-layer classifier}} \quad \phi_{\theta}(x) =: h \leftarrow \underbrace{\phi_{\theta}(x)}_{\text{Last-layer feature}}$$

Treat $H = [h_{1,1} \cdots h_{K,n}]$ as a **free** optimization variable



Simplification: Unconstrained Features

$$\psi_{\Theta}(x) = \underbrace{W_L \sigma(W_{L-1} \cdots \sigma(W_1 x + b_1) + b_{L-1})}_{\text{Last-layer classifier}} + b_L$$

$\phi_{\Theta}(x) =: h \leftarrow$ Last-layer feature

Treat $H = [h_{1,1} \quad \cdots \quad h_{K,n}]$ as a **free** optimization variable

$$\min_{W, H, b} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(W h_{k,i} + b, y_k) + \frac{\lambda_W}{2} \|W\|_F^2 + \frac{\lambda_H}{2} \|H\|_F^2 + \frac{\lambda_b}{2} \|b\|_2^2$$

- **Validity:** Modern network are highly **overparameterized**, that can approximate any point in the feature space [Shaham'18];
- **State-of-the-Art:** also called **Layer-Peeled Model** [Fang'21], existing work [E'20, Lu'20, Mixon'20, Fang'21] **only** studied global optimality conditions.

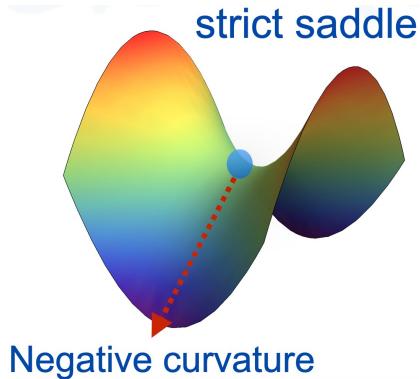
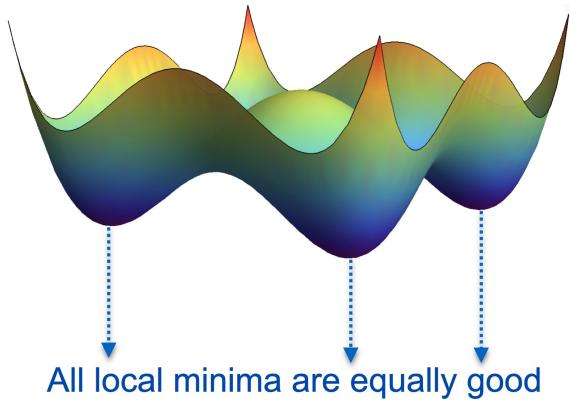
Main Theoretical Results

Theorem (Informal) Consider the nonconvex loss with unconstrained feature model with $K < d$ and balanced data

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} \frac{1}{Kn} \sum_{k=1}^K \sum_{i=1}^n \mathcal{L}_{\text{CE}}(\mathbf{W}\mathbf{h}_{k,i} + \mathbf{b}, \mathbf{y}_k) + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_{\mathbf{b}}}{2} \|\mathbf{b}\|_2^2$$

- (*Global Optimality*) Any global solution $(\mathbf{W}_*, \mathbf{H}_*)$ satisfies the NC properties (1-4).
- (*Benign Global Landscape*) The function has no spurious local minimizer and is a *strict saddle function*, with negative curvature for non-global critical point.

Main Theoretical Results



- Closely relates to **low-rank matrix factorization** problems [Burer et al'03, Bhojanapalli et al'16, Ge et al'16, Zhu et al'18, Li et al'19, Chi et al'19]
 - **Difference in tasks:** classification training vs recovery
 - **Difference in loss functions statistical properties :** cross-entropy vs least-squares; randomness or statistical properties of the sensing matrices
 - **Difference in global solutions.**

Further Readings

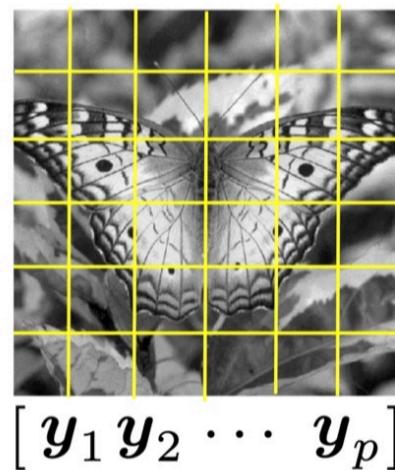
- *Prevalence of Neural Collapse during the terminal phase of deep learning training.* Vardan Papyan, X.Y. Han, David L. Donoho, Proceedings of the National Academy of Sciences, 2020.
- *A geometric analysis of Neural Collapse with unconstrained features.* Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, Qing Qu, NeurIPS, 2021.
- *Exploring deep neural networks via layer-peeled model: Minority collapse in imbalanced training.* Cong Fang, Hangfeng He, Qi Long, Weijie J. Su, Proceedings of the National Academy of Sciences, 2021.

Lecture Agenda

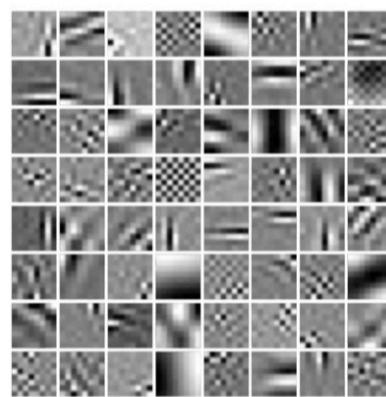
- Problem Introduction
- Applications in SIPML
 - Example I: Generalized Phase Retrieval
 - Example II: Low-Rank Matrix Recovery
 - Example III: Training Deep Neural Networks
 - Example IV: Sparse Dictionary Learning
 - Example V: Sparse Blind Deconvolution

Dictionary Learning

Given \mathbf{Y} , jointly learn a *compact* dictionary \mathbf{A}_0 and *sparse* \mathbf{X}_0 ?



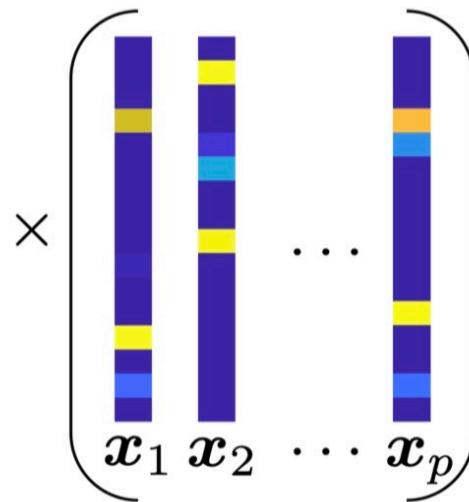
\approx



\mathbf{Y}

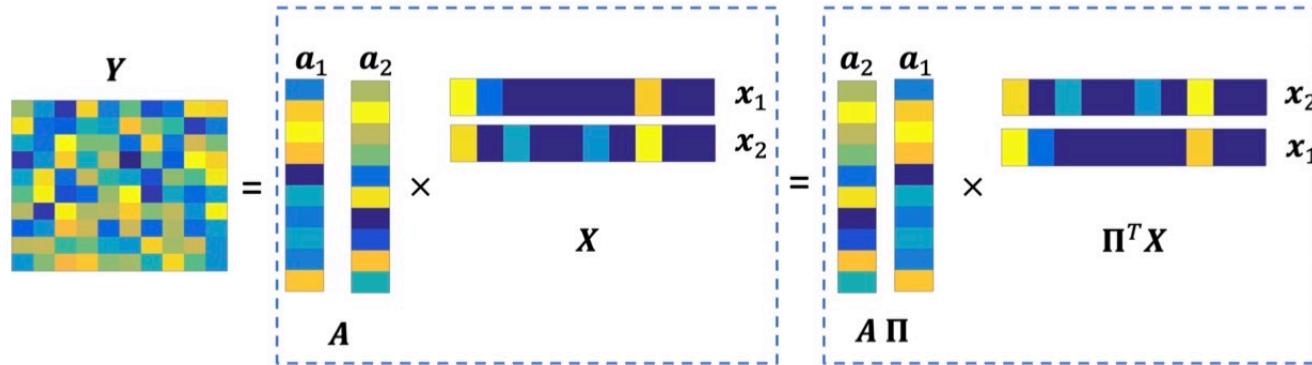
\approx

\mathbf{A}_0



Symmetry in Dictionary Learning

Nonconvexity due to symmetry:



- **Permutation symmetry:** ($2^n n!$ signed permutation Π)

$$Y = A_0 X_0 = (A_0 \Pi) (\Pi^\top X_0)$$

- **Equivalent solution pairs:** $(A_0, X_0) \iff (A_0 \Pi, \Pi^\top X_0)$.

Dictionary Learning – Complete Case

We first study the problem when $\mathbf{A}_0 \in \mathbb{R}^{n \times n}$ is *complete*, that is, \mathbf{A}_0 is square and invertible.

$$\overline{\mathbf{Y}} = (\mathbf{Y}\mathbf{Y}^\top)^{-1/2} \mathbf{Y} \propto \overline{\mathbf{A}}_0 \mathbf{X}_0, \quad \overline{\mathbf{A}}_0 \in O(n).$$

- We find **one column** of $\overline{\mathbf{A}}_0$ via

$$\min_{\mathbf{q}} \varphi(\mathbf{q}^\top \overline{\mathbf{Y}}), \quad \text{s.t.} \quad \mathbf{q} \in \mathbb{S}^{n-1}.$$

- $\varphi(\cdot)$ is a sparsity promoting function, $\varphi(\cdot) = \|\cdot\|_1, \varphi(\cdot) = -\|\cdot\|_4^4$
- assume that sparse $\mathbf{X}_0 \sim \text{Bernoulli-Gaussian}(\theta)$.

Dictionary Learning – Complete Case

We find **one column** of $\overline{\mathbf{A}}_0$ via

$$\boxed{\min_{\mathbf{q}} \varphi(\mathbf{q}^\top \overline{\mathbf{Y}}), \quad \text{s.t.} \quad \mathbf{q} \in \mathbb{S}^{n-1}.}$$

- $\varphi(\cdot)$ is a sparsity promoting function, $\varphi(\cdot) = \|\cdot\|_1$, $\varphi(\cdot) = -\|\cdot\|_4^4$
- assume that sparse $\mathbf{X}_0 \sim \text{Bernoulli-Gaussian}(\theta)$.

High-level intuition: for any $\mathbf{q}_\star = \pm \overline{\mathbf{a}}_{0i}$

$$\mathbb{E}_{\mathbf{X}_0} [\varphi(\mathbf{q}_\star^\top \overline{\mathbf{Y}})] \propto \varphi(\mathbf{q}_\star^\top \overline{\mathbf{A}}_0) = \varphi(\mathbf{e}_i)$$

Dictionary Learning – Complete Case

We find **one column** of $\bar{\mathbf{A}}_0$ via

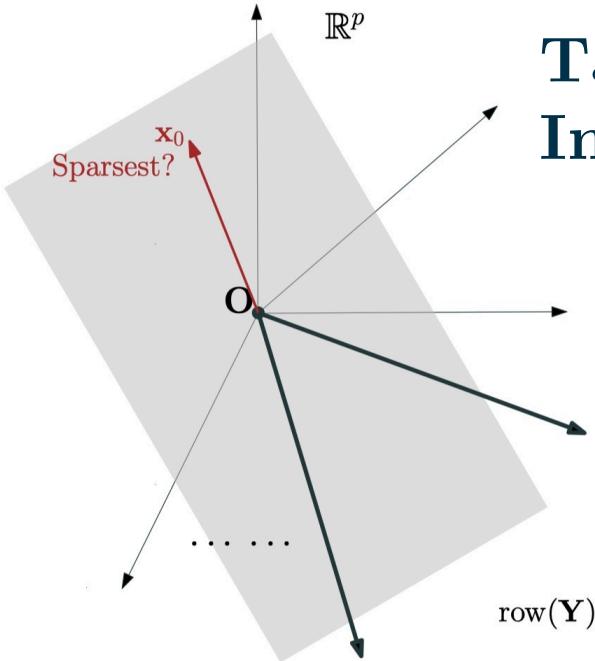
$$\min_{\mathbf{q}} \|\mathbf{q}^\top \bar{\mathbf{Y}}\|_1, \quad \text{s.t.} \quad \mathbf{q} \in \mathbb{S}^{n-1}.$$

High-level intuition: for any $\mathbf{q}_\star = \pm \bar{\mathbf{a}}_{0i}$

$$\mathbb{E}_{\mathbf{X}_0} [\|\mathbf{q}_\star^\top \bar{\mathbf{Y}}\|_1] \propto \|\mathbf{q}_\star^\top \bar{\mathbf{A}}_0\|_1 = \|\mathbf{e}_i\|_1$$

The columns of $\bar{\mathbf{A}}_0$ are equivalent global solutions!

Alternative Interpretations: Finding the Sparsest Vector in a Subspace



Task: Given $\mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0$, recover \mathbf{A}_0 and \mathbf{X}_0

Intuition: \mathbf{A}_0 is square, invertible,

$$\text{row}(\mathbf{Y}) = \text{row}(\mathbf{X}_0)$$

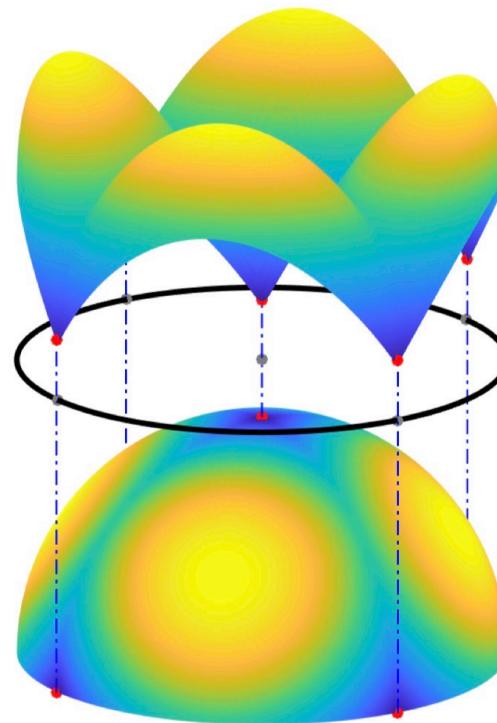
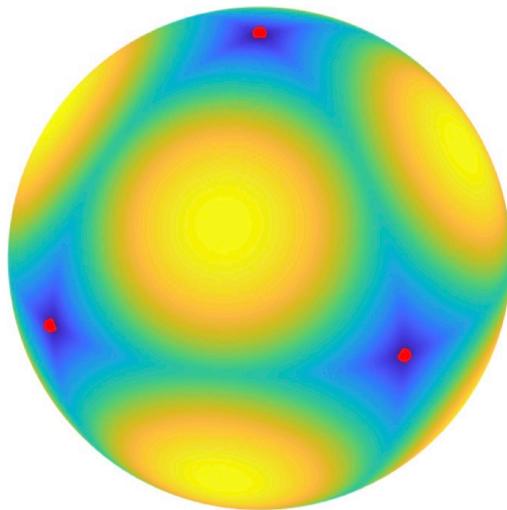
- Find the sparsest vectors in $\text{row}(\mathbf{Y})$:

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

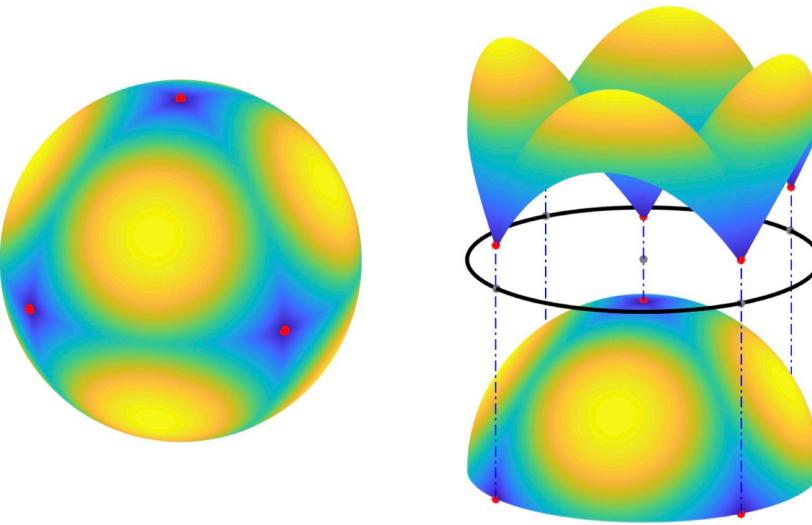
- Nonconvex “relaxation”:

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

Dictionary Learning – Complete Case



Dictionary Learning – Complete Case



- Symmetric copies of the ground truth are minimizers.
- Negative curvature in symmetry-breaking directions.
- Cascade of saddle points.

From Complete Case to Overcomplete Case

We can extend the study to the *overcomplete* case $\mathbf{A}_0 \in \mathbb{R}^{n \times m}$ ($m > n$) under the assumptions that

- Unit norm tight frame:

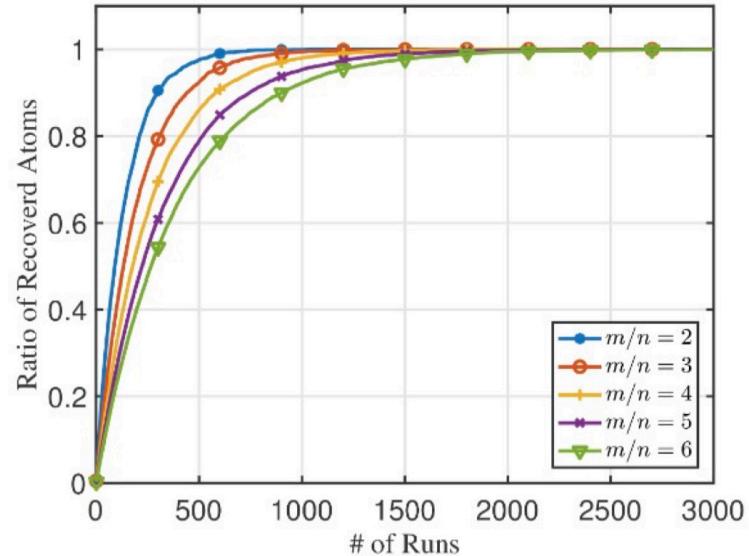
$$\sqrt{\frac{n}{m}} \mathbf{A}_0 \mathbf{A}_0^\top = \mathbf{I}, \quad \|\mathbf{a}_{0i}\|_2 = 1.$$

- Column (near) orthogonal: incoherence

$$\max_{i \neq j} |\langle \mathbf{a}_{0i}, \mathbf{a}_{0j} \rangle| \leq \mu.$$

Recovery of All Dictionary Items

- **Complete case:** this can be done with a deflation process via orthogonal projections
- **Overcomplete case:** we can do repetitive random independent initializations



recover full A_0 via repeated independent trials

Further Readings

- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications.* John Wright, Yi Ma. **(Chapter 7)**
- *Complete Dictionary Recovery Over the Sphere I: Overview and the Geometric Picture.* Ju Sun, Qing Qu, and John Wright, IEEE Trans. Info. Theory, 2016.
- *Analysis of the Optimization Landscapes for Overcomplete Representation Learning.* Qing Qu, Yuexiang Zhai, Xiao Li, Yuqian Zhang, Zhihui Zhu, ICLR'20, 2020.

Lecture Agenda

- Problem Introduction
- Applications in SIPML
 - Example I: Generalized Phase Retrieval
 - Example II: Low-Rank Matrix Recovery
 - Example III: Training Deep Neural Networks
 - Example IV: Sparse Dictionary Learning
 - **Example V: Sparse Blind Deconvolution**

Sparse Blind Deconvolution with Multiple Inputs

Problem: Given multiple \mathbf{y}_i of *circulant* convolutions

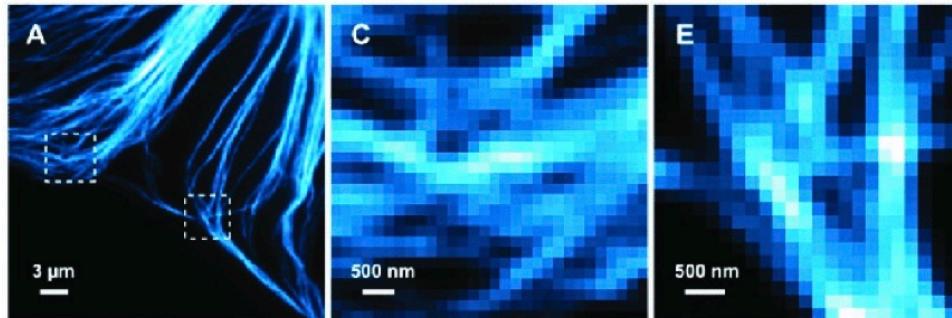
$$\mathbf{y}_i = \mathbf{a} \circledast \mathbf{x}_i, \quad 1 \leq i \leq p$$

jointly learn both \mathbf{a} and *sparse* $\{\mathbf{x}_i\}_{i=1}^p$.

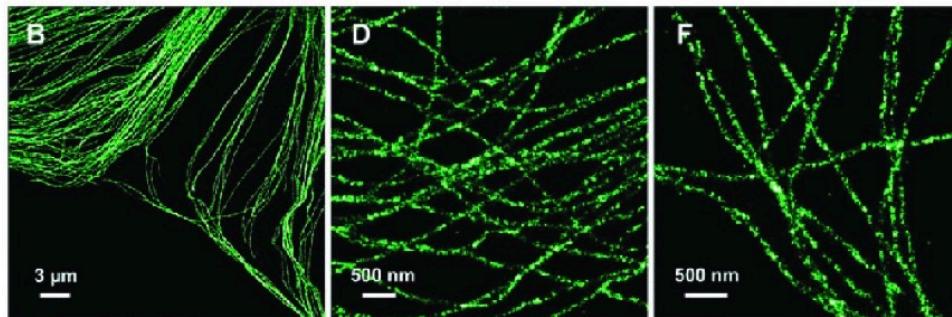
- Here, $\mathbf{y}_i, \mathbf{a}, \mathbf{x}_i \in \mathbb{R}^n$ ($1 \leq i \leq p$).
- Sparse signal \mathbf{x}_i : $\mathbf{x}_i \sim_{i.i.d.} \text{Bernoulli-Gaussian}(\theta)$.

Motivation: Super-resolution Microscopy Imaging

Conventional Fluorescent Optical Microscopy



Stochastic Optical Reconstruction Microscopy (STORM)

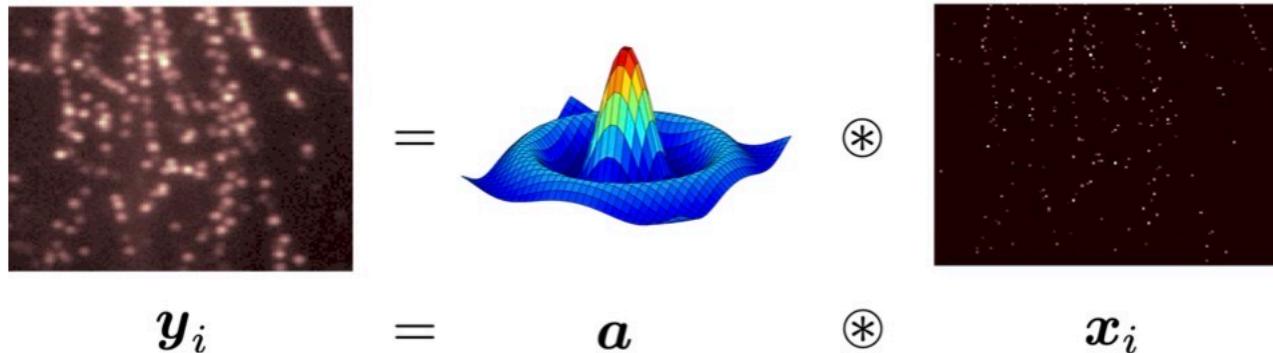


Motivation: Super-resolution Microscopy Imaging

Given multiple \mathbf{y}_i of circulant convolution

$$\mathbf{y}_i = \mathbf{a} \circledast \mathbf{x}_i, \quad 1 \leq i \leq p$$

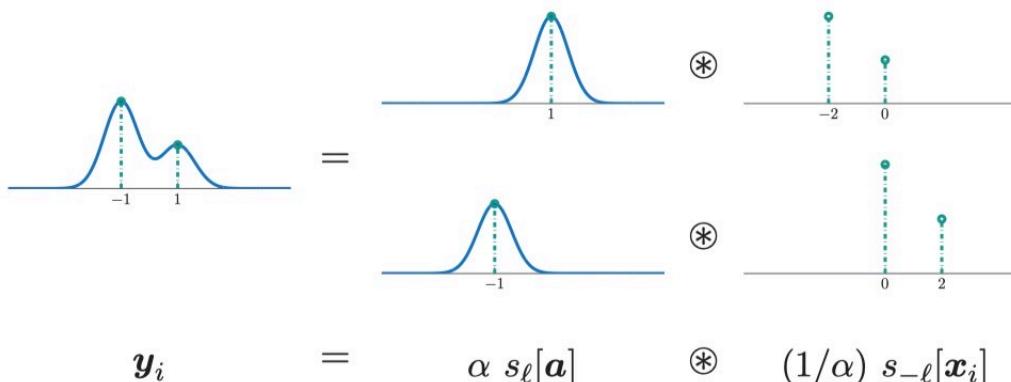
jointly learn both \mathbf{a} and sparse $\{\mathbf{x}_i\}_{i=1}^p$



Symmetry of Sparse Blind Deconvolution

- Scaling symmetry: $\mathbf{y}_i = \mathbf{a} \circledast \mathbf{x}_i = \alpha \mathbf{a} \circledast \alpha^{-1} \mathbf{x}_i$
- Shift symmetry:

$$(\mathbf{a}, \{\mathbf{x}_i\}_{i=1}^p) \iff (\mathbf{s}_\ell[\mathbf{a}], \{\mathbf{s}_{-\ell}[\mathbf{x}_i]\}_{i=1}^p)$$



Symmetry of Sparse Blind Deconvolution

- Scaling symmetry: $\mathbf{y}_i = \mathbf{a} \circledast \mathbf{x}_i = \alpha \mathbf{a} \circledast \alpha^{-1} \mathbf{x}_i$
- Shift symmetry:

$$(\mathbf{a}, \{\mathbf{x}_i\}_{i=1}^p) \iff (\mathbf{s}_\ell[\mathbf{a}], \{\mathbf{s}_{-\ell}[\mathbf{x}_i]\}_{i=1}^p)$$

- can only be solved to a shift ambiguity;
- natural formulations are nonconvex.

Reducing the Problem to Dictionary Learning

- For each $\mathbf{y}_i = \mathbf{a} \circledast \mathbf{x}_i$, equivalently

$$C_{\mathbf{y}_i} = C_{\mathbf{a}} \cdot C_{\mathbf{x}_i}, \quad 1 \leq i \leq p.$$

- Rewrite the problem as

$$\begin{pmatrix} C_{\mathbf{y}_1} & \cdots & C_{\mathbf{y}_p} \end{pmatrix} = C_{\mathbf{a}} \begin{pmatrix} C_{\mathbf{x}_1} & \cdots & C_{\mathbf{x}_p} \end{pmatrix}$$

$$\begin{pmatrix} \text{data } \mathbf{Y} \\ \vdots \\ \text{data } \mathbf{Y} \end{pmatrix} = \text{dictionary} \begin{pmatrix} \text{sparse code } \mathbf{X}_0 \\ \vdots \\ \text{sparse code } \mathbf{X}_0 \end{pmatrix}$$

data \mathbf{Y}

dictionary

sparse code \mathbf{X}_0

Reducing the Problem to Dictionary Learning

Given $\mathbf{Y} = \mathbf{C}_a \mathbf{X}_0$, jointly learn the dictionary \mathbf{C}_a and sparse \mathbf{X}_0 ?

$$\begin{bmatrix} \mathbf{C}_{\mathbf{y}_1} & \cdots & \mathbf{C}_{\mathbf{y}_p} \end{bmatrix} = \mathbf{C}_a \begin{bmatrix} \mathbf{C}_{\mathbf{x}_1} & \cdots & \mathbf{C}_{\mathbf{x}_p} \end{bmatrix}$$

$$\begin{bmatrix} \text{data } \mathbf{Y} \\ \vdots \\ \text{data } \mathbf{Y} \end{bmatrix} = \text{dictionary} \begin{bmatrix} \text{sparse code } \mathbf{X}_0 \\ \vdots \\ \text{sparse code } \mathbf{X}_0 \end{bmatrix}$$

Reducing the Problem to Dictionary Learning

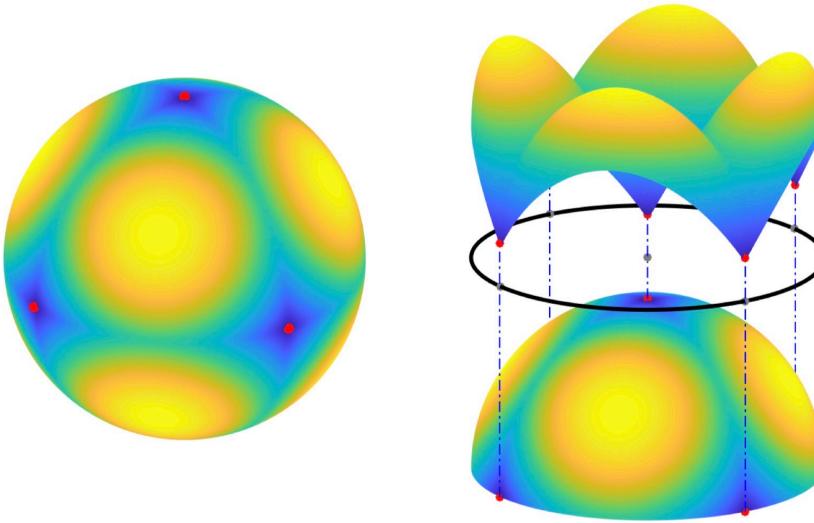
$$\overline{\mathbf{Y}} = (\mathbf{Y}\mathbf{Y}^\top)^{-1/2} \mathbf{Y} \propto \overline{\mathbf{A}}_0 \mathbf{X}_0, \quad \overline{\mathbf{A}}_0 \in O(n).$$

We find **one column** of \mathbf{C}_a via

$$\boxed{\min_{\mathbf{q}} \varphi(\mathbf{q}^\top \overline{\mathbf{Y}}), \quad \text{s.t.} \quad \mathbf{q} \in \mathbb{S}^{n-1}.}$$

- $\varphi(\cdot)$ is a sparsity promoting function, $\varphi(\cdot) = \|\cdot\|_1$, $\varphi(\cdot) = -\|\cdot\|_4^4$
- assume that sparse $\mathbf{X}_0 \sim \text{Bernoulli-Gaussian}(\theta)$.

Benign Global Nonconvex Landscape



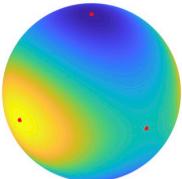
- Symmetric copies of the ground truth are minimizers.
- Negative curvature in symmetry-breaking directions.
- Cascade of saddle points.

Other Applications

Nonconvex Problems with Discrete Symmetries

Eigenvector Computation

Maximize a quadratic form
over the sphere.

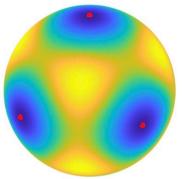


$$\max_{\mathbf{x} \in \mathbb{S}^{n-1}} \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x}.$$

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

Tensor Decomposition

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

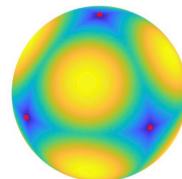


$$\max_{\mathbf{X} \in \mathrm{O}(n)} \sum_i \mathbf{T}(\mathbf{x}_i, \mathbf{x}_i, \mathbf{x}_i, \mathbf{x}_i).$$

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

Dictionary Learning

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

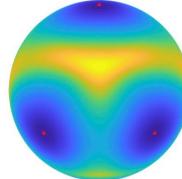


$$\min_{\mathbf{A} \in \mathcal{A}, \mathbf{X}} \frac{1}{2} \|\mathbf{Y} - \mathbf{A}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1.$$

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

Short-and-Sparse Deconvolution

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



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

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

Further Readings

- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications.* John Wright, Yi Ma. **(Chapter 7)**
- *A Nonconvex Approach for Exact and Efficient Multichannel Sparse Blind Deconvolution.* Qing Qu, Xiao Li, Zhihui Zhu, SIAM Imaging Sciences, 2020.
- *Short-and-Sparse Deconvolution -- A Geometric Approach.* Yenson Lau, Qing Qu, Han-Wen Kuo, Pengcheng Zhou, Yuqian Zhang, John Wright, ICLR'20, 2020.

Further Readings

- *Geometry and Symmetry in Short-and-Sparse Deconvolution.* Han-Wen Kuo, Yuqian Zhang, Yenson Lau, John Wright. SIAM Journal on Mathematics of Data Science, 2020.
- *Structured Local Optima in Sparse Blind Deconvolution.* Yuqian Zhang, Han-Wen Kuo, John Wright. IEEE Transaction on Information Theory, 2020.
- *On the Global Geometry of Sphere-Constrained Sparse Blind Deconvolution.* Yuqian Zhang, Yenson Lau, Han-Wen Kuo, Sky Cheung, Abhay Pasupathy, John Wright. EEE Transactions on Pattern Analysis and Machine Intelligence, 2019