

Nonconvex Demixing from Bilinear Measurements

Yuanming Shi



信息科学与技术学院
School of Information Science and Technology

Outline

- **Motivations**

- Blind deconvolution meets blind demixing

- **Two Vignettes:**

- **Implicitly regularized Wirtinger flow**

- ❖ Why nonconvex optimization?

- ❖ Implicitly regularized Wirtinger flow

- **Matrix optimization over manifolds**

- ❖ Why manifold optimization?

- ❖ Riemannian optimization for blind demixing

*Motivations: **Blind deconvolution** meets
blind demixing*



Blind deconvolution

- In many science and engineering problems, the observed signal can be modeled as:

$$z(t) = f(t) * g(t)$$

where $*$ is the convolution operator

- $f(t)$ is a physical signal of interest
- $g(t)$ is the impulse response of the sensory system

- **Applications:** astronomy, neuroscience, image processing, computer vision, wireless communications, microscopy data processing,...
- **Blind deconvolution:** estimate $f(t)$ and $g(t)$ given $z(t)$



Image deblurring

- Blurred images due to camera shake can be modeled as a convolution of the *latent sharp image* and a *kernel* capturing the motion of the camera

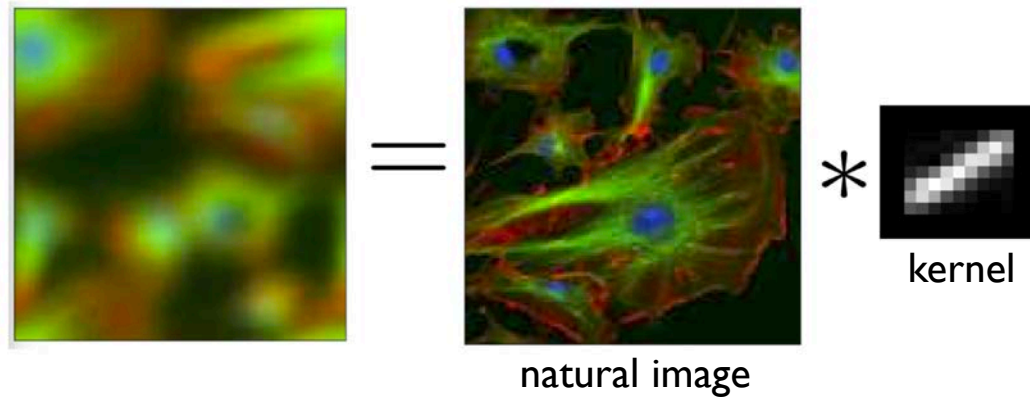
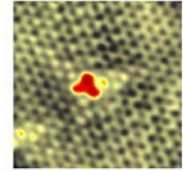
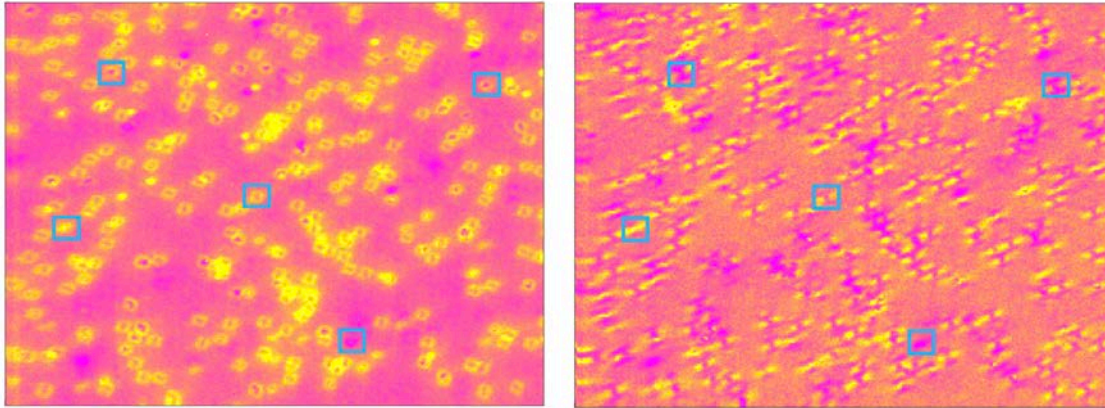


Fig. credit: Chi

How to find the high-resolution image and the blurring kernel simultaneously?

Microscopy data analysis

- **Defects:** the electronic structure of the material is contaminated by randomly and sparsely distributed “defects”



Doped Graphene

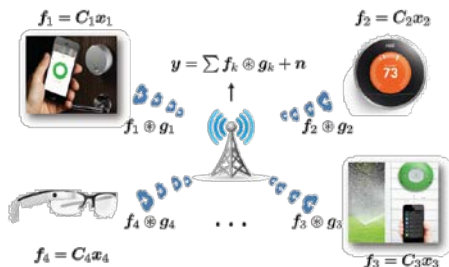
Fig. credit:Wright

How to determine the locations and characteristic signatures of the defects?

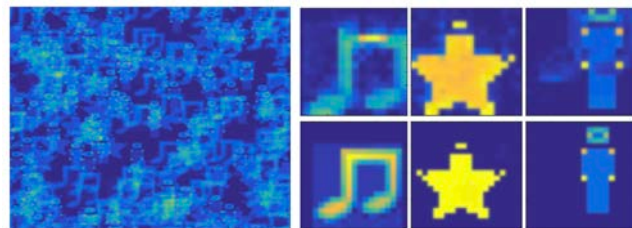
Blind demixing

- The received measurement consists of the sum of all convolved signals

$$z(t) = \sum_{i=1}^S f_i(t) * g_i(t)$$



low-latency communication for IoT



convolutional dictionary learning (multi kernel)

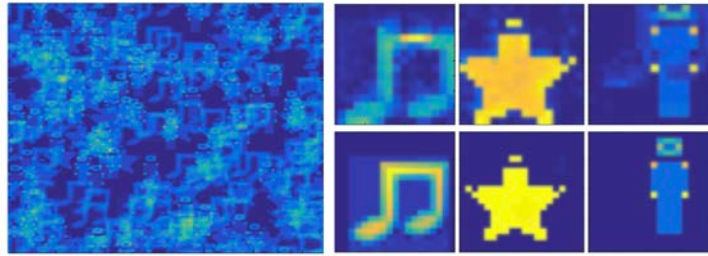
- **Applications:** IoT, dictionary learning, neural spike sorting,...
- **Blind demixing:** estimate $\{f_i(t)\}$ and $\{g_i(t)\}$ given $z(t)$

Convolutional dictionary learning

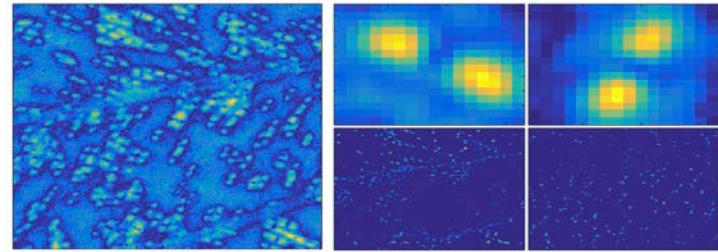
- The observation signal is the superposition of several convolutions

$$z(t) = \sum_{i=1}^s f_i(t) * g_i(t)$$

Fig. credit:Wright



experiment on synthetic image



experiment on microscopy image

How to recover multiple kernels and the corresponding activation signals?

Low-latency communications for IoT

- **Packet structure:** metadata (preamble (PA) and header (H)) and data



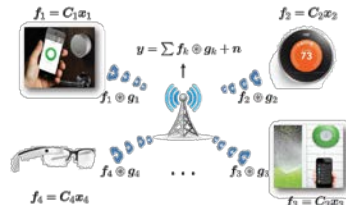
long data packet in current wireless systems



short data packet in IoT

- **Proposal:** transmitters just send overhead-free signals, and the receiver can still extract the information

$$z(t) = \sum_{i=1}^S f_i(t) * g_i(t)$$



How to detect data without channel estimation in multi-user environments?

Demixing from bilinear model?



$$z(t) = \sum_{i=1}^s f_i(t) * g_i(t)$$

Bilinear model

- Translate into the frequency domain...

$$z = \sum_{i=1}^s f_i \odot g_i \in \mathbb{C}^m$$

- Subspace assumptions:** f_i and g_i lie in some known low-dimensional subspaces

$$f_i = A_i x_i^{\natural} \in \mathbb{C}^m \quad g_i = B h_i^{\natural} \in \mathbb{C}^m$$

where $A_i = [a_{i1}, \dots, a_{im}]^* \in \mathbb{C}^{m \times L}$, $B = [b_1, \dots, b_m]^* \in \mathbb{C}^{m \times K}$ and $L, K \ll m$

$$a_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{CN}(0, I) \quad \{b_j\} : \text{partial Fourier basis}$$

- Demixing from bilinear measurements:**

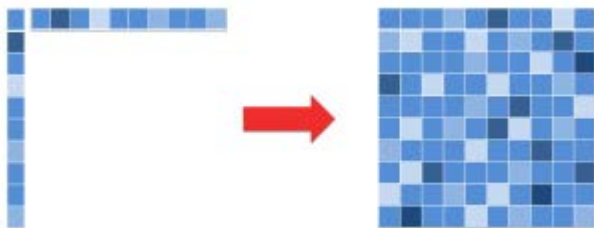
$$\text{find } \{x_i\}, \{h_i\} \quad \text{subject to } z_j = \sum_{i=1}^s b_j^* h_i x_i^* a_{ij}, \quad 1 \leq j \leq m \quad \square$$

An equivalent view: low-rank factorization

- **Lifting:** introduce $M_k^{\natural} = \mathbf{h}_k^{\natural} \mathbf{x}_k^{\natural*}$ to linearize constraints

$$z_j = \sum_{i=1}^s \mathbf{b}_i^* \mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural*} \mathbf{a}_{ij} = \sum_{i=1}^s \mathbf{b}_i^* \underbrace{(\mathbf{h}_i^{\natural} \mathbf{x}_i^{\natural*})}_{M_i^{\natural}} \mathbf{a}_{ij}$$

$M_i^{\natural} \in \mathbb{C}^{K \times L}$



- **Low-rank matrix optimization problem**

$$\begin{aligned} & \text{find} \quad \{M_i\} \\ & \text{subject to} \quad z_j = \sum_{i=1}^s \mathbf{b}_i^* M_i \mathbf{a}_{ij}, \quad j = 1, \dots, m \\ & \quad \text{rank}(M_i) = 1, \quad i = 1, \dots, s, \end{aligned}$$

Convex relaxation

- Ling and Strohmer (TIT'2017) proposed to solve the **nuclear norm minimization problem**:

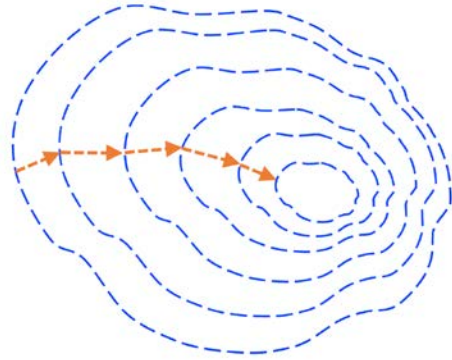
$$\text{minimize } \sum_{k=1}^s \|M_k\|_* \quad \mathbf{a}_{kj} \stackrel{\text{i.i.d.}}{\sim} \mathcal{CN}(\mathbf{0}, \mathbf{I})$$

$$\text{subject to } z_j = \sum_{k=1}^s \mathbf{b}_k^* M_k \mathbf{a}_{kj}, \quad j = 1, \dots, m \quad \{\mathbf{b}_j\}: \text{partial Fourier basis}$$

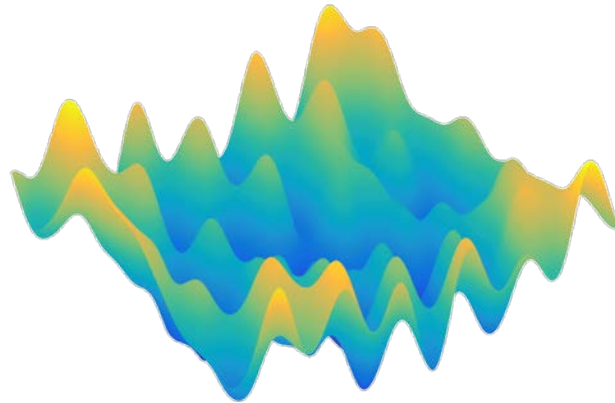
- **Sample-efficient:** $m \gtrsim s^2 \max\{K, L\} \log^2 m$ samples for exact recovery if $\{\mathbf{b}_j\}$ is incoherent w.r.t. $\{\mathbf{h}_k^{\dagger}\}$
- **Computational-expensive:** SDP in the lifting space

Can we solve the nonconvex matrix optimization problem directly?

Vignettes A: *Implicitly regularized Wirtinger flow*



Why nonconvex optimization?

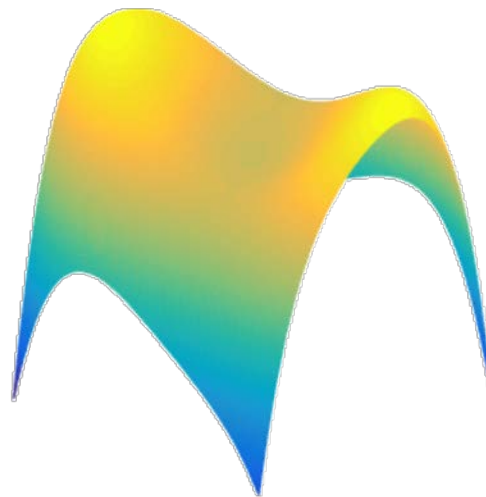


Nonconvex problems are everywhere

- Empirical risk minimization is usually nonconvex

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}; \theta)$$

- low-rank matrix completion
- blind deconvolution/demixing
- dictionary learning
- phase retrieval
- mixture models
- deep learning
- ...



Nonconvex optimization may be super scary

- **Challenges:** saddle points, local optima, bumps,...

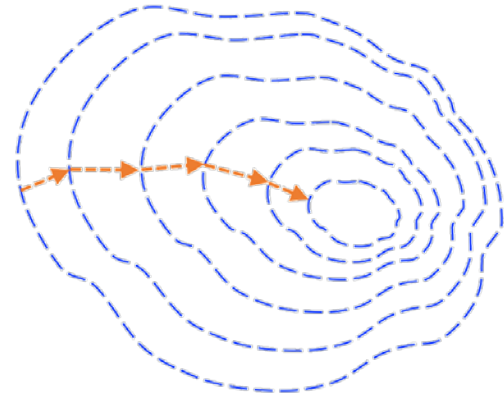
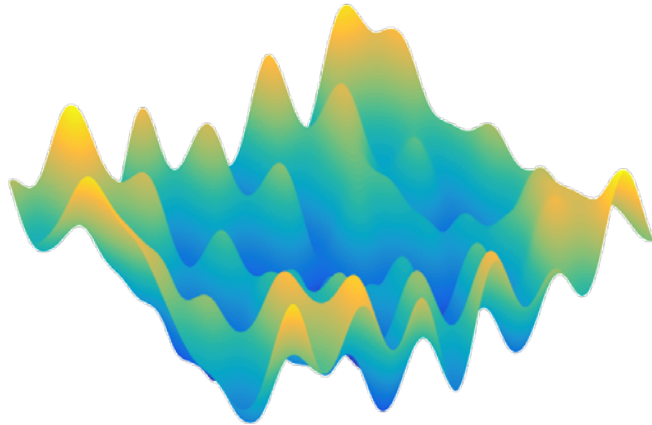


Fig. credit: Chen

- **Fact:** they are usually solved on a daily basis via simple algorithms like (stochastic) gradient descent

Statistical models come to rescue

- **Blessings:** when data are generated by certain statistical models, problems are often much nicer than worst-case instances

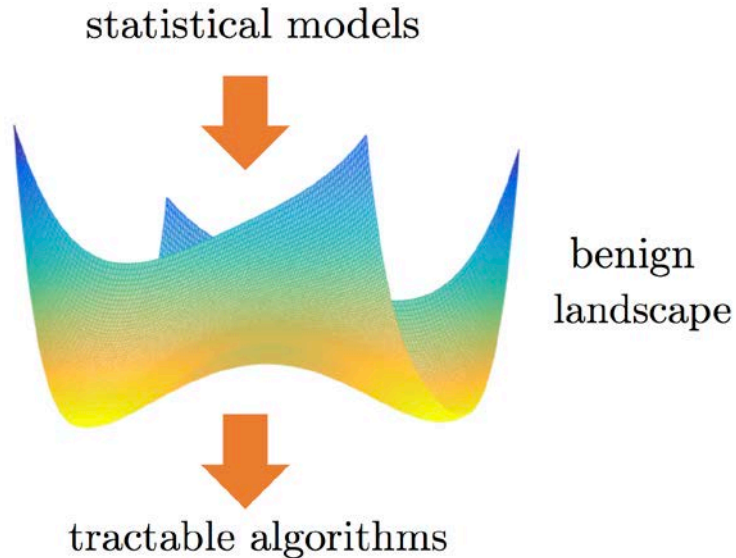
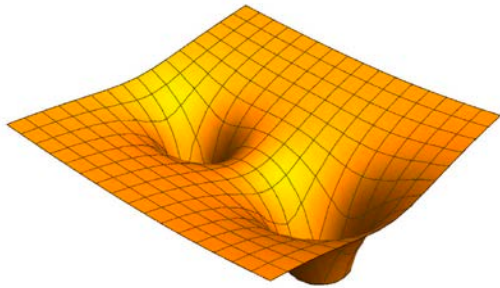


Fig. credit: Chen

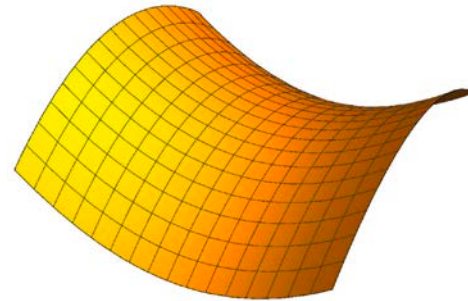
First-order stationary points

- Saddle points and local minima:

$$\lambda_{\min}(\nabla^2 f(\mathbf{z})) \begin{cases} > 0 & \text{local minimum} \\ = 0 & \text{local minimum or saddle point} \\ < 0 & \text{strict saddle point} \end{cases}$$



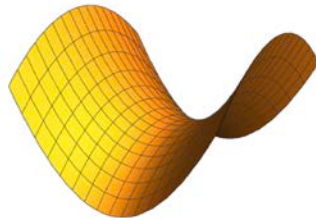
Local minima



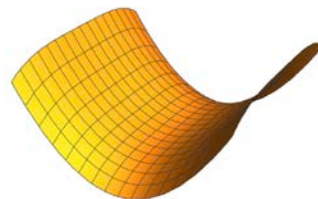
Saddle points/local maxima

First-order stationary points

- **Applications:** PCA, matrix completion, dictionary learning etc.
 - **Local minima:** either all local minima **are** global minima or all local minima **as good as** global minima
 - **Saddle points:** **very poor** compared to global minima; **several** such points



Strict saddle point



Non-strict saddle point

- **Bottomline:** local minima much more desirable than saddle points

How to escape saddle points efficiently?

Statistics meets optimization

- **Proposal:** separation of landscape analysis and generic algorithm design

landscape analysis
(statistics)

all local minima are
global minima

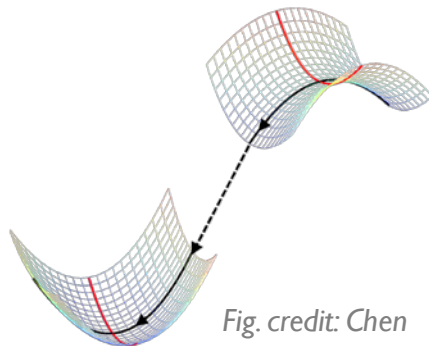
- dictionary learning (Sun et al. '15)
- phase retrieval (Sun et al. '16)
- matrix completion (Ge et al. '16)
- synchronization (Bandeira et al. '16)
- inverting deep neural nets (Hand et al. '17)
- ...



generic algorithms
(optimization)

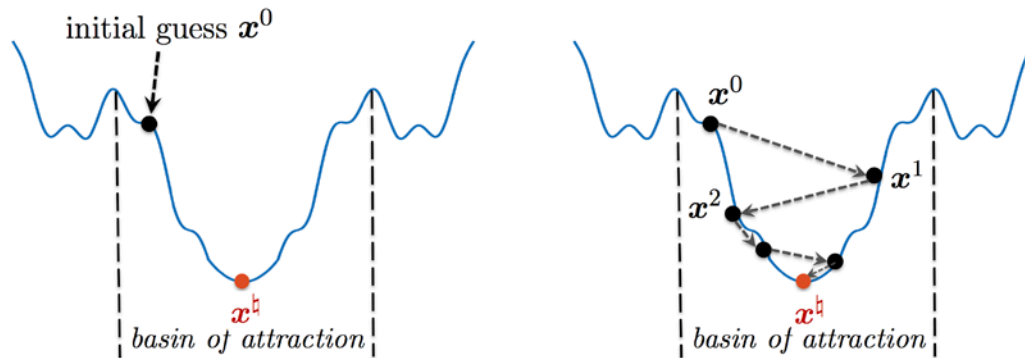
all the saddle points
can be escaped

- gradient descent (Lee et al. '16)
- trust region method (Sun et al. '16)
- perturbed GD (Jin et al. '17)
- cubic regularization (Agarwal et al. '17)
- Natasha (Allen-Zhu '17)
- ...



Issue: conservative computational guarantees for specific problems
(e.g., phase retrieval, blind deconvolution, matrix completion)

Solution: blending landscape and convergence analysis



implicitly regularized Wirtinger flow

A natural least-squares formulation

- **Goal:** demixing from bilinear measurements

$$\text{Given: } y_j = \sum_{i=1}^s \mathbf{b}_j^* \mathbf{h}_i \mathbf{x}_i^* \mathbf{a}_{ij}, \quad 1 \leq j \leq m$$

$$\text{minimize}_{\{\mathbf{h}_k\}, \{\mathbf{x}_k\}} f(\mathbf{h}, \mathbf{x}) := \sum_{j=1}^m \sum_{k=1}^s (\mathbf{b}_j^* \mathbf{h}_k \mathbf{x}_k^* \mathbf{a}_{kj} - y_j)^2$$

- **Pros:** computational-efficient in the natural parameter space
- **Cons:** $f(\cdot)$ is nonconvex: **bilinear** constraint, scaling ambiguity

Wirtinger flow

- Least-square minimization via Wirtinger flow (Candes, Li, Soltanolkotabi '14)

$$\underset{\{\mathbf{h}_k\}, \{\mathbf{x}_k\}}{\text{minimize}} \quad f(\mathbf{h}, \mathbf{x}) := \sum_{j=1}^m \sum_{k=1}^s (\mathbf{b}_j^* \mathbf{h}_k \mathbf{x}_k^* \mathbf{a}_{kj} - y_j)^2$$

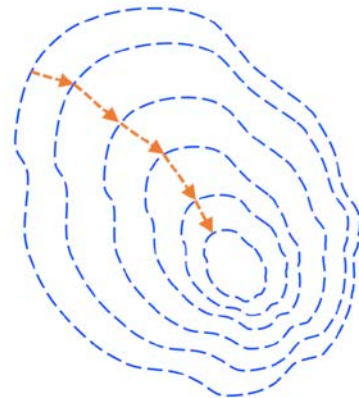
- **Spectral initialization by top eigenvector of**

$$\mathbf{M}_k := \sum_{j=1}^m \mathbf{y}_j \mathbf{b}_j \mathbf{a}_{kj}^*, \quad k = 1, \dots, s$$

- **Gradient iterations**

$$\mathbf{h}_k^{t+1} = \mathbf{h}_k^t - \eta \frac{1}{\|\mathbf{x}_k^t\|_2^2} \nabla_{\mathbf{h}_k} f(\mathbf{h}^t, \mathbf{x}^t)$$

$$\mathbf{x}_k^{t+1} = \mathbf{x}_k^t - \eta \frac{1}{\|\mathbf{h}_k^t\|_2^2} \nabla_{\mathbf{x}_k} f(\mathbf{h}^t, \mathbf{x}^t)$$



Two-stage approach

- **Initialize** within local basin sufficiently close to ground-truth (i.e., strongly convex, no saddle points/ local minima)
- **Iterative refinement** via some iterative optimization algorithms

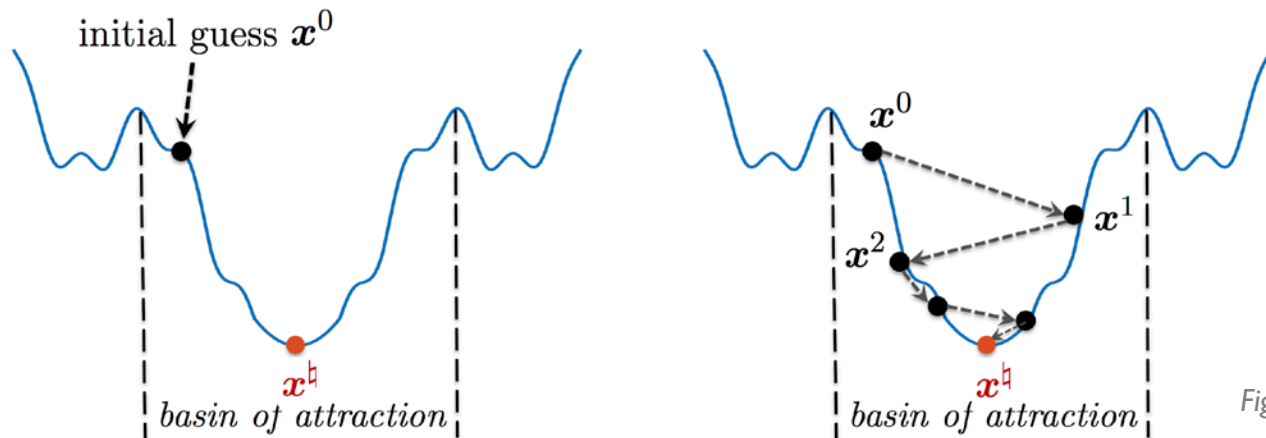
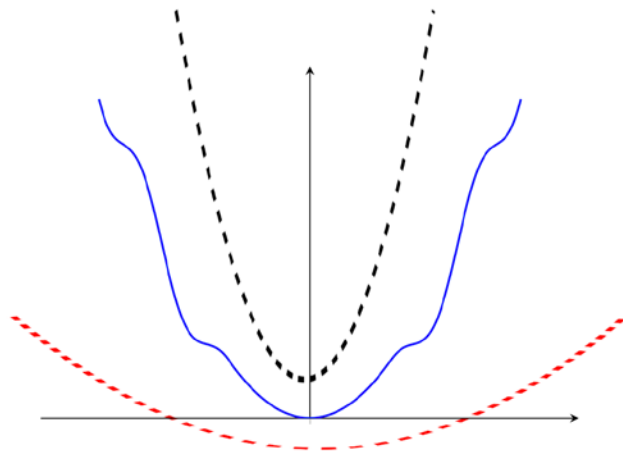


Fig. credit: Chen

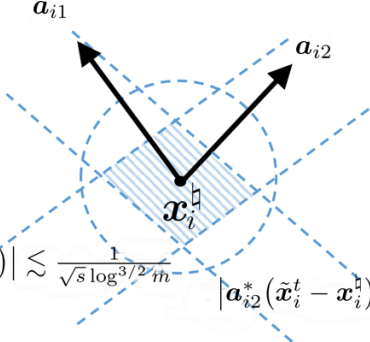
Gradient descent theory

- Two standard conditions that enable **geometric convergence** of GD
 - (local) restricted strong convexity
 - (local) smoothness



Gradient descent theory

- **Question:** which region enjoys both **strong convexity** and **smoothness**?


$$|\mathbf{a}_{i1}^* (\tilde{\mathbf{x}}_i^t - \mathbf{x}_i^h)| \lesssim \frac{1}{\sqrt{s} \log^{3/2} m} \quad |\mathbf{a}_{i2}^* (\tilde{\mathbf{x}}_i^t - \mathbf{x}_i^h)| \lesssim \frac{1}{\sqrt{s} \log^{3/2} m}$$

- \mathbf{x} is not far away from \mathbf{x}^h (**convexity**)
- \mathbf{x} is incoherent w.r.t. sampling vectors (**incoherence region for smoothness**)

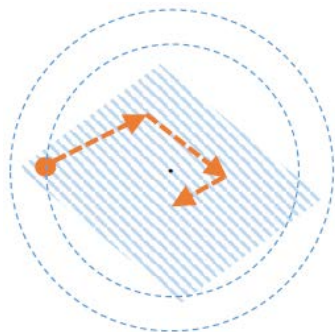
Prior works suggest enforcing **regularization** (e.g., regularized loss [Ling & Strohmer'17]) to promote incoherence

Our finding: WF is implicitly regularized

- WF (GD) implicitly forces iterates to remain **incoherent** with $\{\mathbf{a}_{ij}\}$

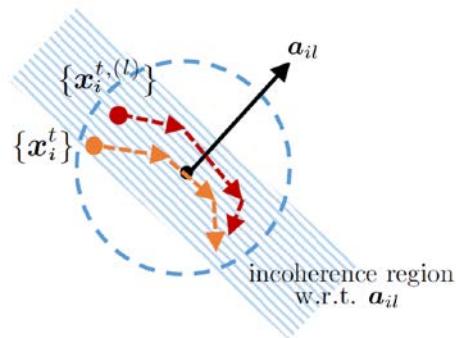
$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{a}_{ij}^* \left(\alpha_i^t \mathbf{x}_i^t - \mathbf{x}_i^{\natural} \right) \right| \lesssim \frac{1}{\sqrt{s} \log^{3/2} m} \|\mathbf{x}_i^{\natural}\|_2$$

- cannot be derived from generic optimization theory
- relies on finer *statistical analysis* for entire trajectory of GD



region of local strong
convexity and smoothness

Key proof idea: leave-one-out analysis



- introduce **leave-one-out** iterates $x_i^{t,(l)}$ by running VWF without l -th sample
- leave-one-out iterate $x_i^{t,(l)}$ is independent of a_{il}
- leave-one-out iterate $x_i^{t,(l)} \approx$ true iterate x_i^t
- x_i^t is nearly independent of (i.e., **nearly orthogonal to**) a_{il}

Theoretical guarantees

- With i.i.d. Gaussian design, WF (**regularization-free**) achieves

- **Incoherence**

$$\max_{1 \leq i \leq s, 1 \leq j \leq m} \left| \mathbf{a}_{ij}^* \left(\alpha_i^t \mathbf{x}_i^t - \mathbf{x}_i^{\natural} \right) \right| \lesssim \frac{1}{\sqrt{s} \log^{3/2} m} \|\mathbf{x}_i^{\natural}\|_2$$

- **Near-linear convergence rate**

$$\text{dist}(\mathbf{z}^t, \mathbf{z}^{\natural}) \lesssim \left(1 - \frac{\eta}{16\kappa} \right)^t \frac{1}{\log^2 m}$$

- **Summary:**

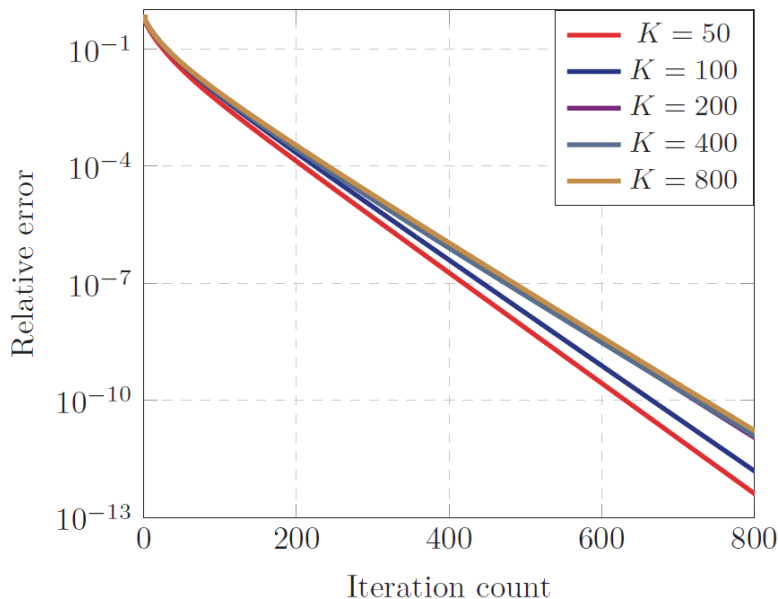
- **Sample size:** $m \gtrsim s^2 \max\{K, L\} \text{poly} \log m$

- **Stepsize:** $\eta \asymp s^{-1}$ vs. $\eta \lesssim (sm)^{-1}$ [Ling & Strohmer'17]

- **Computational complexity:** $\mathcal{O}(s \log \frac{1}{\varepsilon})$ vs. $\mathcal{O}(sm \log \frac{1}{\varepsilon})$ [Ling & Strohmer'17]

Numerical results

- stepsize: $\eta = 0.1$
- number of users: $s = 10$
- sample size: $m = 50K$



linear convergence:

WF attains ε -accuracy within $\mathcal{O}(s \log \frac{1}{\varepsilon})$ iterations

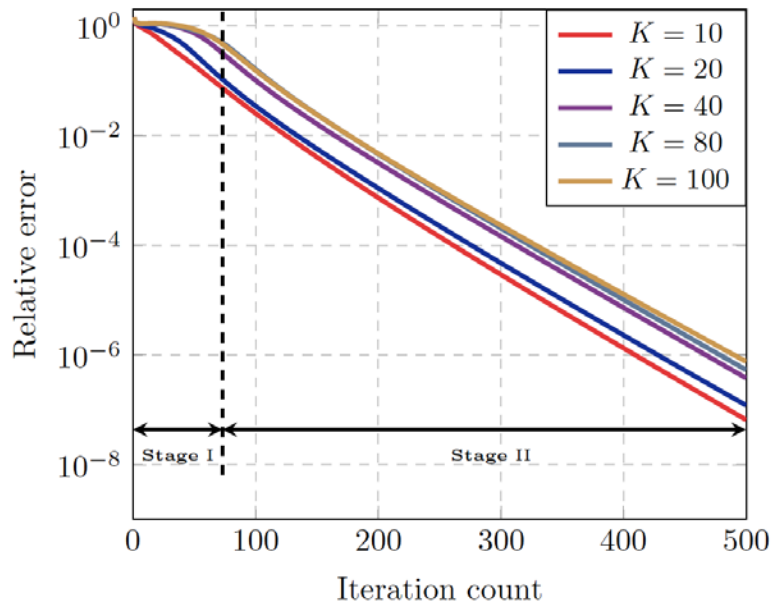
Is carefully-designed initialization necessary?



Numerical results of randomly initialized WF

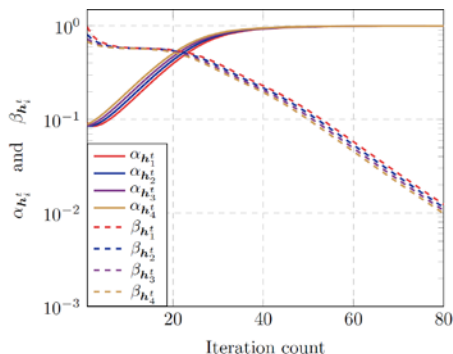
- stepsize: $\eta = 0.1$
- number of users: $s = 10$
- sample size: $m = 50K$
- initial point:

$$\mathbf{h}_i^0 \sim \mathcal{N}(\mathbf{0}, \frac{1}{L} \mathbf{I}_L), \mathbf{x}_i^0 \sim \mathcal{N}(\mathbf{0}, \frac{1}{K} \mathbf{I}_K);$$
$$i = 1, \dots, s, (K = L)$$



Randomly initialized WF enters local basin within $\mathcal{O}(s \log K)$ iterations

Analysis: population dynamics



Population level (infinite sample)

$$\mathbf{h}_i^{t+1} = \mathbf{h}_i^t - \eta \frac{1}{\|\mathbf{x}_i^t\|_2^2} \nabla_{\mathbf{h}_i} F(\mathbf{h}^t, \mathbf{x}^t)$$

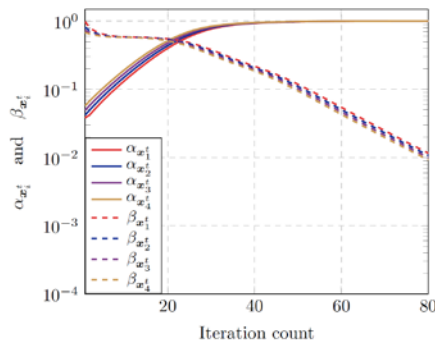
$$\nabla_{\mathbf{h}_i} F(\mathbf{h}, \mathbf{x}) := \mathbb{E}[\nabla_{\mathbf{h}_i} f(\mathbf{h}, \mathbf{x})] = \|\mathbf{x}_i\|_2^2 \mathbf{h}_i - (\mathbf{x}_i^{\mathfrak{h}*} \mathbf{x}_i) \mathbf{h}_i^{\mathfrak{h}};$$

- Signal strength: $\alpha_{\mathbf{h}_i^t} := \langle \mathbf{h}_i^{\mathfrak{h}}, 1/\overline{\omega_i^t} \mathbf{h}_i^t \rangle \|\mathbf{h}_i^{\mathfrak{h}}\|_2$, ω_i^t is the alignment parameter
- Size of residual component: $\beta_{\mathbf{h}_i^t} := \left\| \mathbf{h}_i^t - \langle \mathbf{h}_i^{\mathfrak{h}}, 1/\overline{\omega_i^t} \mathbf{h}_i^t \rangle \mathbf{h}_i^{\mathfrak{h}} \right\|_2$

State evolution

$$\left. \begin{aligned} \alpha_{\mathbf{h}_i^{t+1}} &= (1 - \eta) \alpha_{\mathbf{h}_i^t} + \eta \alpha_{\mathbf{x}_i^t} / (\alpha_{\mathbf{x}_i^t}^2 + \beta_{\mathbf{x}_i^t}^2) \\ \beta_{\mathbf{h}_i^{t+1}} &= (1 - \eta) \beta_{\mathbf{h}_i^t} \end{aligned} \right\} \rightarrow \begin{aligned} T_\gamma &= \mathcal{O}(s \log K) \\ \text{dist}(\mathbf{h}_i^{T_\gamma}, \mathbf{h}_i^{\mathfrak{h}}) &\leq \gamma \end{aligned} \quad \text{local basin}$$

Analysis: population dynamics



Population level (infinite sample)

$$\mathbf{x}_i^{t+1} = \mathbf{x}_i^t - \eta \frac{1}{\|\mathbf{h}_i^t\|_2^2} \nabla_{\mathbf{x}_i} F(\mathbf{h}^t, \mathbf{x}^t)$$

$$\nabla_{\mathbf{x}_i} F(\mathbf{h}, \mathbf{x}) := \mathbb{E}[\nabla_{\mathbf{x}_i} f(\mathbf{h}, \mathbf{x})] = \|\mathbf{h}_i\|_2^2 \mathbf{x}_i - (\mathbf{h}_i^{\natural*} \mathbf{h}_i) \mathbf{x}_i^{\natural}$$

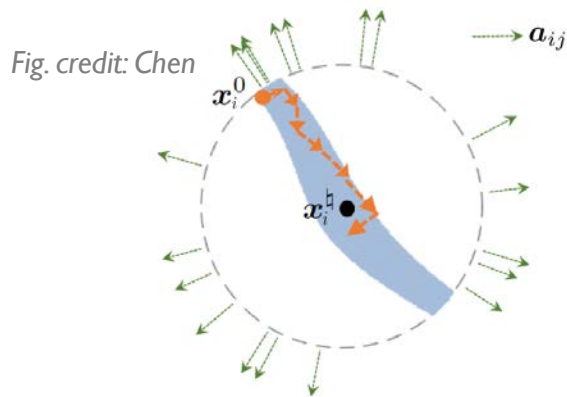
- Signal strength: $\alpha_{\mathbf{x}_i^t} := \langle \mathbf{x}_i^{\natural}, \omega_i^t \mathbf{x}_i^t \rangle \|\mathbf{x}_i^{\natural}\|_2$, ω_i^t is the alignment parameter
- Size of residual component: $\beta_{\mathbf{x}_i^t} := \left\| \mathbf{x}_i^t - \langle \mathbf{x}_i^{\natural}, \omega_i^t \mathbf{x}_i^t \rangle \mathbf{x}_i^{\natural} \right\|_2$

State evolution

$$\left. \begin{aligned} \alpha_{\mathbf{x}_i^{t+1}} &= (1 - \eta) \alpha_{\mathbf{x}_i^t} + \eta \alpha_{\mathbf{h}_i^t} / (\alpha_{\mathbf{h}_i^t}^2 + \beta_{\mathbf{h}_i^t}^2) \\ \beta_{\mathbf{x}_i^{t+1}} &= (1 - \eta) \beta_{\mathbf{x}_i^t} \end{aligned} \right\} \rightarrow \begin{aligned} &T_\gamma = \mathcal{O}(s \log K) \\ &\text{dist}(\mathbf{x}_i^{T_\gamma}, \mathbf{x}_i^{\natural}) \leq \gamma \end{aligned} \quad \text{local basin}$$

Analysis: finite-sample analysis

$$\mathbf{z}_i^{t+1} = \begin{bmatrix} \mathbf{h}_i^{t+1} \\ \mathbf{x}_i^{t+1} \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{h}_i^t - \eta/\|\mathbf{x}_i^t\|_2 \cdot \nabla_{\mathbf{h}_i} F(\mathbf{z}) \\ \mathbf{x}_i^t - \eta/\|\mathbf{x}_i^t\|_2 \cdot \nabla_{\mathbf{x}_i} F(\mathbf{z}) \end{bmatrix}}_{:=\mathbf{m}(\mathbf{z}_i^t)} - \underbrace{\begin{bmatrix} \eta/\|\mathbf{x}_i^t\|_2 \cdot (\nabla_{\mathbf{h}_i} f(\mathbf{z}) - \nabla_{\mathbf{h}_i} F(\mathbf{z})) \\ \eta/\|\mathbf{h}_i^t\|_2 \cdot (\nabla_{\mathbf{x}_i} f(\mathbf{z}) - \nabla_{\mathbf{x}_i} F(\mathbf{z})) \end{bmatrix}}_{:=\mathbf{r}(\mathbf{z}_i^t)}$$



$\mathbf{r}(\mathbf{z}_i^t)$ is well-controlled
in this region

- **Population-level analysis** holds approximately if $\mathbf{r}(\mathbf{z}_i^t) \ll \mathbf{m}(\mathbf{z}_i^t)$
- $\mathbf{r}(\mathbf{z}_i^t)$ is well-controlled if \mathbf{x}_i^t is independent of $\{\mathbf{a}_{ij}\}$
- **Key analysis ingredient:** show \mathbf{x}_i^t is “nearly independent” of each $\{\mathbf{a}_{ij}\}$

Theoretical guarantees

- With i.i.d. Gaussian design, WF with random initialization achieves

$$\text{dist}(\mathbf{z}^t, \mathbf{z}^\natural) \lesssim \gamma \left(1 - \frac{\eta}{16\kappa}\right)^{t-T_\gamma} \|\mathbf{z}^\natural\|_2, \quad t \geq T_\gamma$$

Summary:

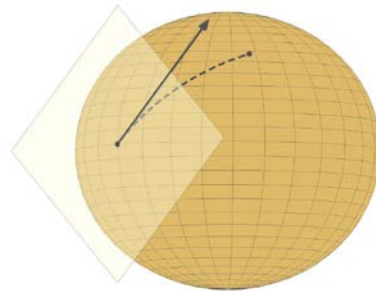
- **Stepsize:** $\eta \asymp s^{-1}$
- **Sample size:** $m \gtrsim s^2 \max\{K, L\} \text{poly log } m$
- **Stage I:** reach local basin $\text{dist}(\mathbf{z}^t, \mathbf{z}^\natural) \leq \gamma$ within $T_\gamma = \mathcal{O}(s \log K)$ iterations
- **Stage II:** linear convergence $\mathcal{O}(s \log \frac{1}{\varepsilon})$
- **Computational complexity:** $\mathcal{O}(s \log K + s \log \frac{1}{\varepsilon})$

Vignettes **B: Matrix optimization** over **manifolds**



Optimization over Riemannian Manifolds (non-Euclidean geometry)

Why manifold optimization?

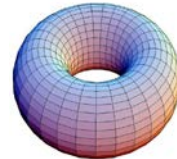
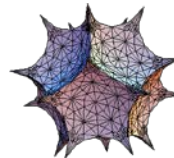
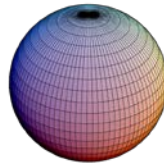


What is manifold optimization?

- Manifold (or manifold-constrained) optimization problem

$$\underset{M \in \mathbb{C}^{m \times n}}{\text{minimize}} \quad f(M) \quad \text{subject to} \quad M \in \mathcal{M}$$

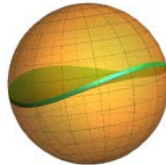
- $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is a **smooth function**
- \mathcal{M} is a **Riemannian manifold**: **spheres**, orthonormal bases (Stiefel), rotations, **positive definite matrices**, **fixed-rank matrices**, Euclidean distance matrices, **semidefinite fixed-rank matrices**, **linear subspaces (Grassmann)**, phases, essential matrices, **fixed-rank tensors**, Euclidean spaces...



Convergence results of manifold optimization

- Convergence guarantees for Riemannian **trust regions**
 - Global convergence to **second-order critical points**
 - Quadratic convergence rate locally
 - Reach ϵ -**second order stationary point** $\|\text{grad}f(z)\| \leq \epsilon$ and $\nabla^2 f(z) \succeq -\epsilon I$ in $\mathcal{O}(1/\epsilon^3)$ iterations under Lipschitz assumptions [Cartis & Absil'16]

Escape *strict* saddle points via finding second-order stationary point



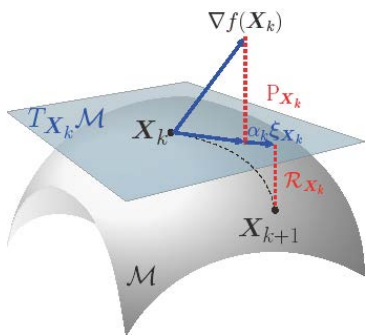
Recent applications of manifold optimization

- **High-dimensional data analysis:** matrix/tensor completion/recovery: [Vandereycken'13], [Boumal-Absil'15], [Kasai-Mishra'16]; phase retrieval: [Sun-Qu-Wright'17]; community detection: [Boumal'16], [Bandeira-Boumal-Voroninski'16],...
- **Machine and deep learning:** Gaussian mixture models: [Hosseini-Sra'15]; dictionary learning: [Sun-Qu-Wright'17]; deep metric learning: [Roy-Mhammedi-Harandi'18],...
- **Wireless transceivers design:** [Shi-Zhang-Letaief'16], [Yu-Shen-Zhang-K. B. Letaief'16], [Shi-Mishra-Chen'17],...

Exploit manifold geometry to address non-convex problems

The power of manifold optimization paradigms

- Generalize Euclidean gradient (Hessian) to *Riemannian gradient (Hessian)*



$$\nabla_{\mathcal{M}} f(\mathbf{X}^{(k)}) = P_{\mathbf{X}^{(k)}}(\nabla f(\mathbf{X}^{(k)}))$$

Riemannian Gradient Euclidean Gradient

$$\mathbf{X}^{(k+1)} = \mathcal{R}_{\mathbf{X}^{(k)}}(-\alpha^{(k)} \nabla_{\mathcal{M}} f(\mathbf{X}^{(k)}))$$

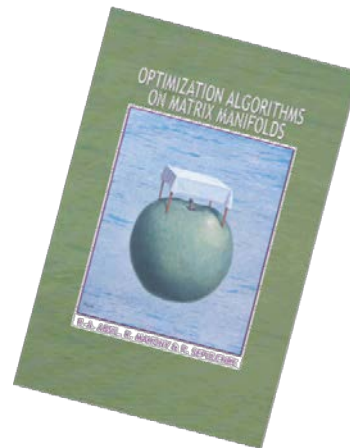
Retraction Operator

- We need Riemannian geometry: 1) linearize search space \mathcal{M} into a **tangent space** $T_{\mathbf{X}}\mathcal{M}$; 2) pick a **metric** on $T_{\mathbf{X}}\mathcal{M}$ to give intrinsic notions of **gradient** and **Hessian**

An excellent book

Optimization algorithms on matrix manifolds

A Matlab toolbox



Manopt

[Home](#)

[Tutorial](#)

[Forum](#)

[About](#)

[Contact](#)

Welcome to Manopt!

A Matlab toolbox for optimization on manifolds

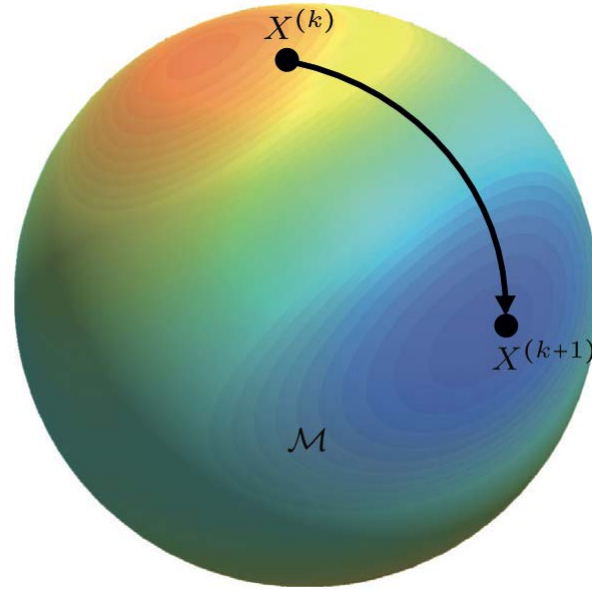
Optimization on manifolds is a powerful paradigm to address nonlinear optimization problems. With Manopt, it is easy to deal with various types of constraints that arise naturally in applications, such as orthonormality or low rank.

[Download](#)

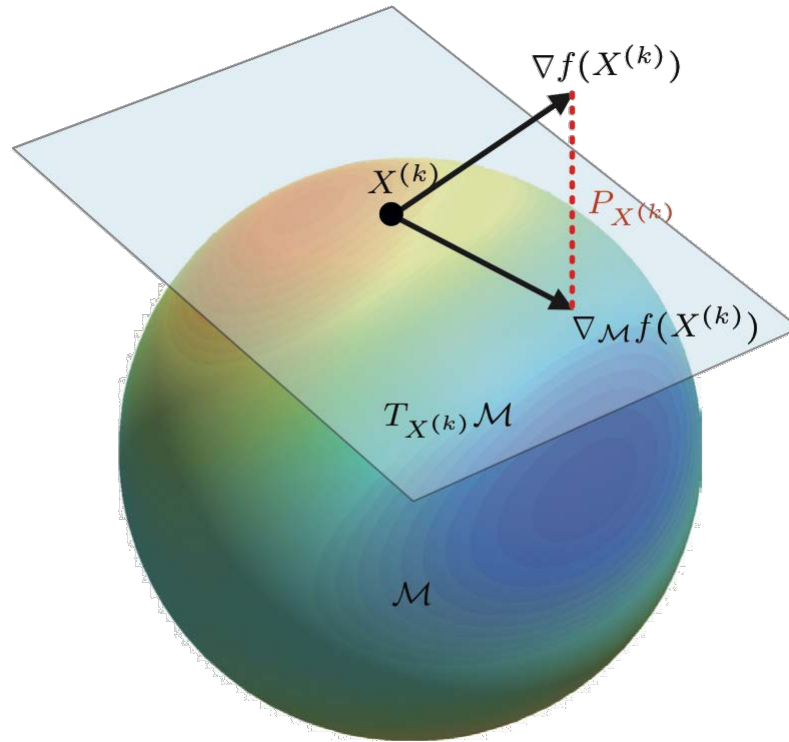
[Get started](#)

*Taking a close look at **gradient descent***

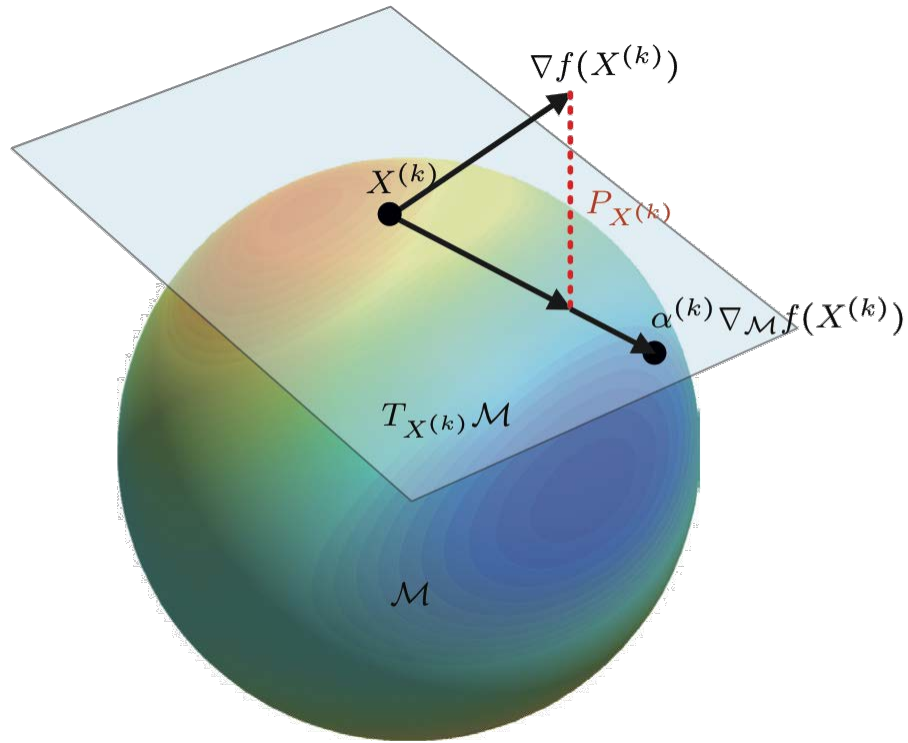
Optimization on the manifold: main idea



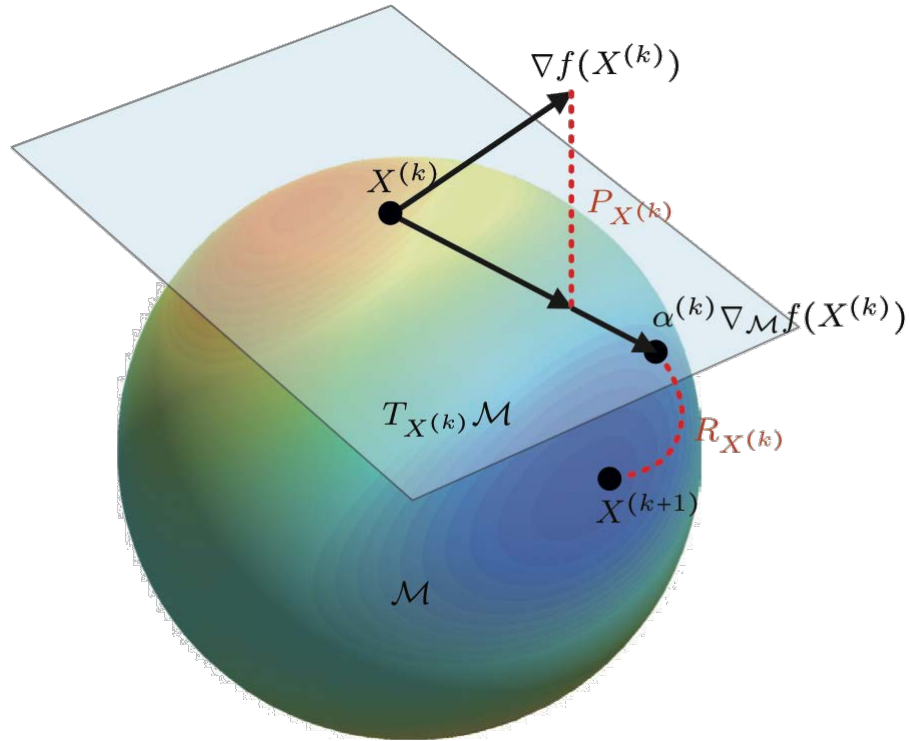
Optimization on the manifold: main idea



Optimization on the manifold: main idea



Optimization on the manifold: main idea



Example: Rayleigh quotient

- Optimization over (sphere) manifold $\mathbb{S}^{n-1} = \{x \in \mathbb{R}^n : x^T x = 1\}$

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x) = -x^T A x \quad \text{subject to} \quad x^T x = 1$$

- The cost function is smooth on \mathbb{S}^{n-1} , symmetric matrix $A \in \mathbb{R}^{n \times n}$

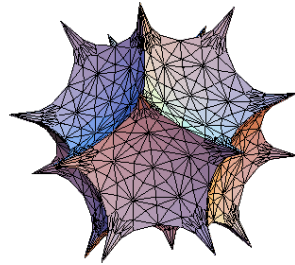
- Step 1: Compute the **Euclidean gradient** in \mathbb{R}^n

$$\nabla f(x) = -2Ax$$

- Step 2: Compute the **Riemannian gradient** on \mathbb{S}^{n-1} via projecting $\nabla f(x)$ to the tangent space using the orthogonal projector $\text{Proj}_x u = (I - xx^T)u$

$$\text{grad} f(x) = \text{Proj}_x \nabla f(x) = -2(I - xx^T)Ax$$

Riemannian optimization for blind demixing



Blind demixing via low-rank optimization

- **Linear mapping:** from bilinear model to linear model

$$y_j = \sum_{i=1}^s \mathbf{b}_j^* \mathbf{h}_i \mathbf{x}_i^* \mathbf{a}_{ij}, \quad 1 \leq j \leq m$$

➤ $\mathbf{b}_j^* \mathbf{h}_i \mathbf{x}_i^* \mathbf{a}_{ij} = \langle \mathbf{b}_j \mathbf{a}_{ij}^*, \mathbf{h}_i \mathbf{x}_i^* \rangle$

➤ $\mathcal{A}_i(\mathbf{X}_i) := \{\langle \mathbf{b}_j \mathbf{a}_{ij}^*, \mathbf{h}_i \mathbf{x}_i^* \rangle\}_{j=1}^L = \{\langle \mathbf{A}_{ij}, \mathbf{X}_k \rangle\}_{j=1}^L, \quad \mathbf{X}_i = \mathbf{h}_i \mathbf{x}_i^*$

- **Proposal:** (non-convex) low-rank optimization problem

$$\begin{aligned} \mathcal{P} : \text{minimize}_{\mathbf{W}_k \in \mathbb{C}^{N \times K}} & \quad \left\| \sum_{k=1}^s \mathcal{A}_k(\mathbf{W}_k) - \mathbf{y} \right\|^2 \\ \text{subject to} & \quad \text{rank}(\mathbf{W}_k) = 1, \quad k = 1, \dots, s, \end{aligned}$$

- Challenges: nonconvex constraints, **complex asymmetric matrices**

Blind demixing via Riemannian optimization

- Handle complex asymmetric matrices

➤ Define linear map $\mathcal{J}_k : \mathbb{S}_+^{(N+K) \times (N+K)} \rightarrow \mathbb{C}^L$ as

$$[\mathcal{J}_k(\mathbf{Y}_k)]_i = \langle \mathbf{J}_{ki}, \mathbf{Y}_k \rangle, \mathbf{Y}_k \in \mathbb{S}_+^{(N+K)(N+K)} \quad \mathbf{J}_{ki} = \begin{bmatrix} \mathbf{0}_{N \times N} & \mathbf{A}_{ki} \\ \mathbf{0}_{K \times N} & \mathbf{0}_{K \times K} \end{bmatrix}$$

- Matrix optimization over the **product manifolds**

$$\text{minimize}_{\mathbf{M}_k \in \mathbb{S}_+^{(N+K)}} \left\| \sum_{k=1}^s \mathcal{J}_k(\mathbf{M}_k) - \mathbf{y} \right\|^2$$

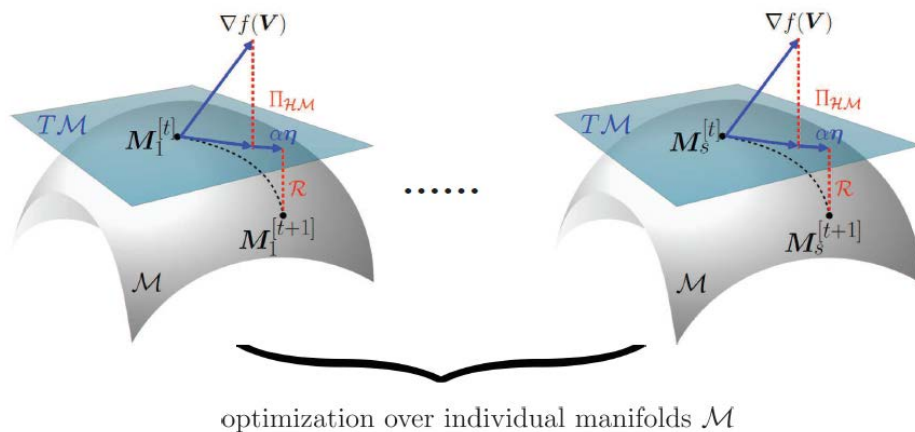
$$\text{subject to } \text{rank}(\mathbf{M}_k) = 1, k = 1, \dots, s$$

➤ **Key observations:** rank-one Hermitian positive semidefinite matrices is a manifold; multiple rank-one constraints construct a manifold

Riemannian optimization over product manifolds

- **Elementwise extension principles**

- The manifold topology of the product manifold is equivalent to the product topology



Element-wise optimization-related ingredients

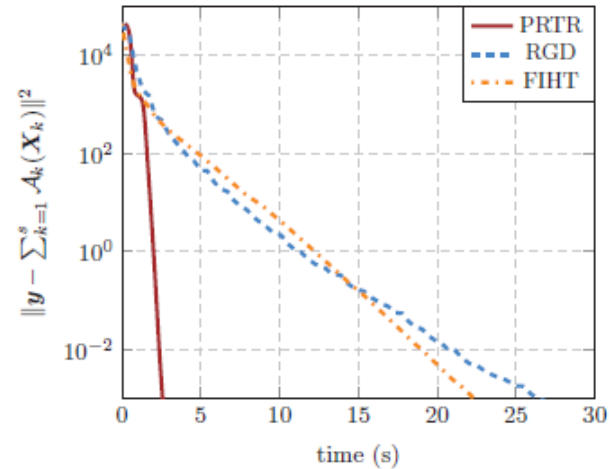
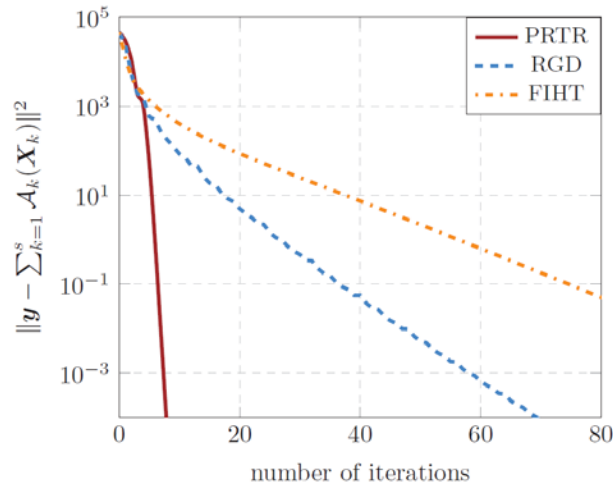
- Riemannian optimization for blind demixing

$$\begin{aligned} & \underset{\mathbf{M}_k \in \mathbb{S}_+^{(N+K)}}{\text{minimize}} \quad \left\| \sum_{k=1}^s \mathcal{J}_k(\mathbf{M}_k) - \mathbf{y} \right\|^2 \\ & \text{subject to} \quad \text{rank}(\mathbf{M}_k) = 1, \quad k = 1, \dots, s \end{aligned}$$

	minimize $\mathbf{w}_k \in \mathcal{M} \left\ \sum_{k=1}^s \mathcal{J}_k(\mathbf{w}_k \mathbf{w}_k^H) - \mathbf{y} \right\ ^2$
Computational space \mathcal{M}	\mathbb{C}_*^{N+K}
Quotient space \mathcal{M}/\sim	$\mathbb{C}_*^{N+K}/\text{SU}(1)$
Riemannian metric $g_{\mathbf{w}_k}$	$g_{\mathbf{w}_k}(\boldsymbol{\zeta}_{\mathbf{w}_k}, \boldsymbol{\eta}_{\mathbf{w}_k}) = \text{Tr}(\boldsymbol{\zeta}_{\mathbf{w}_k}^H \boldsymbol{\eta}_{\mathbf{w}_k} + \boldsymbol{\eta}_{\mathbf{w}_k}^H \boldsymbol{\zeta}_{\mathbf{w}_k})$
Horizontal space $\mathcal{H}_{\mathbf{w}_k} \mathcal{M}$	$\boldsymbol{\eta}_{\mathbf{w}_k} \in \mathbb{C}^{N+K} : \boldsymbol{\eta}_{\mathbf{w}_k}^H \mathbf{w}_k = \mathbf{w}_k^H \boldsymbol{\eta}_{\mathbf{w}_k}$
Horizontal space projection	$\Pi_{\mathcal{H}_{\mathbf{w}_k} \mathcal{M}}(\boldsymbol{\eta}_{\mathbf{w}_k}) = \boldsymbol{\eta}_{\mathbf{w}_k} - a \mathbf{w}_k, \quad a = (\mathbf{w}_k^H \boldsymbol{\eta}_{\mathbf{w}_k} - \boldsymbol{\eta}_{\mathbf{w}_k}^H \mathbf{w}_k) / 2 \mathbf{w}_k^H \mathbf{w}_k$
Riemannian gradient $\text{grad}_{\mathbf{w}_k} f$	$\text{grad}_{\mathbf{w}_k} f = \Pi_{\mathcal{H}_{\mathbf{w}_k} \mathcal{M}}(\frac{1}{2} \nabla_{\mathbf{w}_k} f(\mathbf{v}))$
Riemannian Hessian $\text{Hess}_{\mathbf{w}_k} f[\boldsymbol{\eta}_{\mathbf{w}_k}]$	$\text{Hess}_{\mathbf{w}_k} f[\boldsymbol{\eta}_{\mathbf{w}_k}] = \Pi_{\mathcal{H}_{\mathbf{w}_k} \mathcal{M}}(\frac{1}{2} \nabla_{\mathbf{w}_k}^2 f(\mathbf{v})[\boldsymbol{\eta}_{\mathbf{w}_k}])$
Retraction $\mathcal{R}_{\mathbf{w}_k} : T_{\mathbf{w}_k} \mathcal{M} \rightarrow \mathcal{M}$	$\mathcal{R}_{\mathbf{w}_k}(\boldsymbol{\eta}_{\mathbf{w}_k}) = \mathbf{w}_k + \boldsymbol{\eta}_{\mathbf{w}_k}$

Numerical results

- Optimize over the product of multiple rank-one Hermitian positive semidefinite matrices



Riemannian algorithms: 1) exploit the rank structure in a principled way; 2) develop second-order algorithms systematically; 3) scalable, SVD-free

Concluding remarks

- **Implicitly regularized Wirtinger flow**

- **Implicit regularization:** vanilla gradient descent automatically forces iterates to stay *incoherent*
- Even **simplest** nonconvex methods are remarkably efficient under suitable **statistical** models

- **Matrix optimization over manifolds**

- Exploit the manifold geometry of multiple rank-one Hermitian positive semidefinite matrices
- Develop second-order algorithms systematically: escape saddle points, quadratic convergence rate

- **Future works:** sparse blind demixing, convolutional dictionary learning [Wright, CVPR'17], convolutional neural network [Papayan, et al., SPM'18],...

Reference

- J. Dong and Y. Shi, “Nonconvex demixing from bilinear measurements,” *IEEE Trans. Signal Process.*, vol. 66, no. 19, pp. 5152-5166, Oct., 2018.
- J. Dong, K. Yang, and Y. Shi, “Blind demixing for low-latency communication,” *IEEE Trans. Wireless Commun.*, vol. 18, no. 2, pp. 897-911, Feb., 2019.
- J. Dong, Y. Shi, and Z. Ding, “Blind over-the-air computation and data fusion via provable Wirtinger flow,” <https://arxiv.org/abs/1811.04644>.
- J. Dong and Y. Shi, “Blind Demixing via Wirtinger Flow with Random initialization,” in *Proc. Int. Conf. Artificial Intell. Stat. (AISTATS)*, 2019.

Thanks