

Implicit Regularization in Nonconvex Statistical Estimation



Yuxin Chen

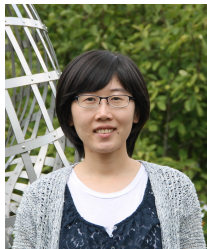
Electrical Engineering, Princeton University



Cong Ma
Princeton ORFE



Kaizheng Wang
Princeton ORFE

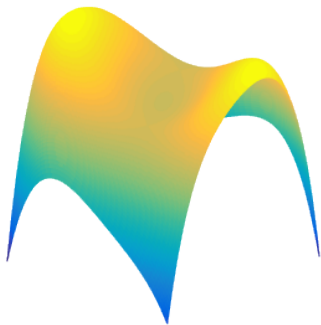


Yuejie Chi
CMU ECE

Nonconvex estimation problems are everywhere

Empirical risk minimization is usually nonconvex

$$\begin{array}{llll} \text{minimize}_x & \ell(\mathbf{x}; \mathbf{y}) & \rightarrow & \text{may be nonconvex} \\ \text{subj. to} & \mathbf{x} \in \mathcal{S} & \rightarrow & \text{may be nonconvex} \end{array}$$

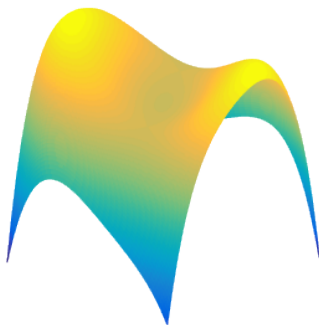


Nonconvex estimation problems are everywhere

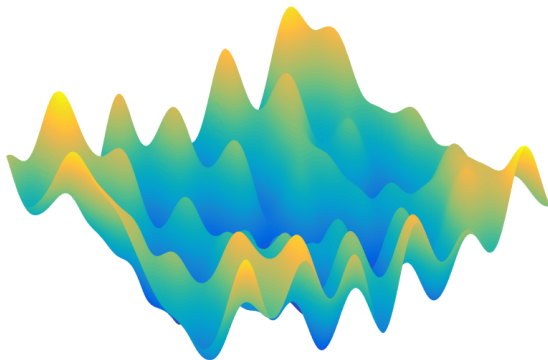
Empirical risk minimization is usually nonconvex

$$\begin{array}{llll} \text{minimize}_x & \ell(\mathbf{x}; \mathbf{y}) & \rightarrow & \text{may be nonconvex} \\ \text{subj. to} & \mathbf{x} \in \mathcal{S} & \rightarrow & \text{may be nonconvex} \end{array}$$

- low-rank matrix completion
- graph clustering
- dictionary learning
- mixture models
- deep learning
- ...



Nonconvex optimization may be super scary



There may be bumps everywhere and exponentially many local optima

e.g. 1-layer neural net (Auer, Herbster, Warmuth '96; Vu '98)

Nonconvex optimization may be super scary



There may be bumps everywhere and exponentially many local optima

e.g. 1-layer neural net (Auer, Herbster, Warmuth '96; Vu '98)

... but is sometimes much nicer than we think

Under certain **statistical models**,
we see benign global geometry: **no spurious local optima**

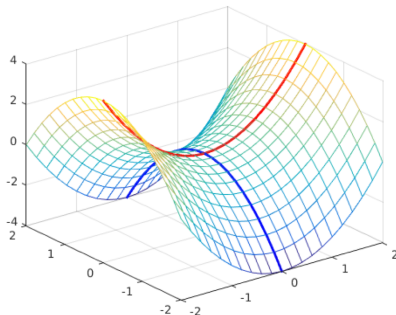
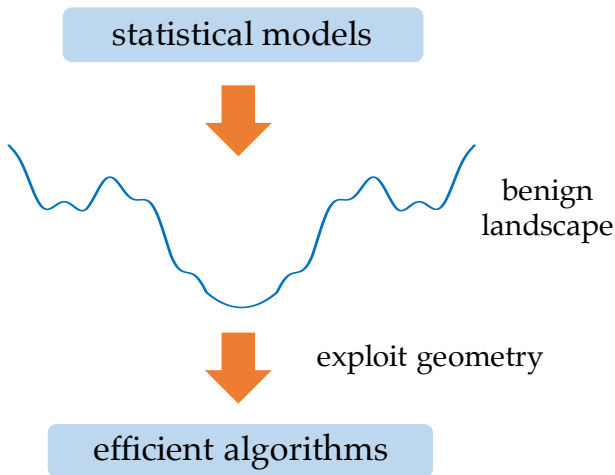
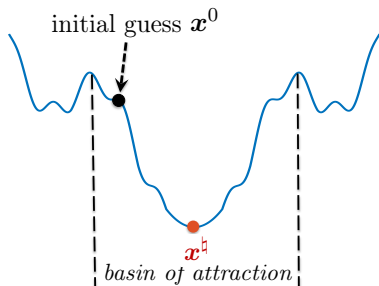


Fig credit: Sun, Qu & Wright

... but is sometimes much nicer than we think

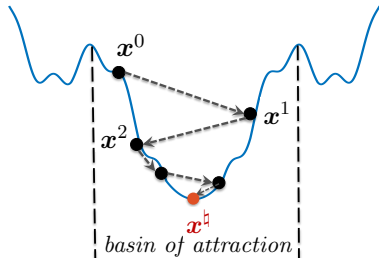
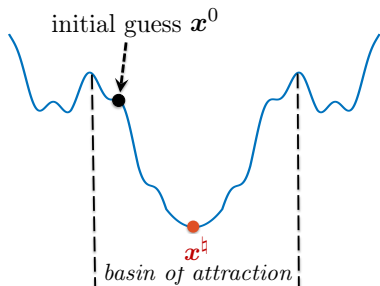


Optimization-based methods: two-stage approach



- Start from an appropriate initial point

Optimization-based methods: two-stage approach



- Start from an appropriate initial point
- Proceed via some iterative optimization algorithms

Roles of regularization

- Prevents overfitting and improves generalization
 - e.g. ℓ_1 penalization, SCAD, nuclear norm penalization, ...

Roles of regularization

- Prevents overfitting and improves generalization
 - e.g. ℓ_1 penalization, SCAD, nuclear norm penalization, ...
- Improves computation by stabilizing search directions
 - e.g. trimming, projection, regularized loss

Roles of regularization

- Prevents overfitting and improves generalization
 - e.g. ℓ_1 penalization, SCAD, nuclear norm penalization, ...
- Improves computation by stabilizing search directions
 - \implies focus of this talk
 - e.g. trimming, projection, regularized loss

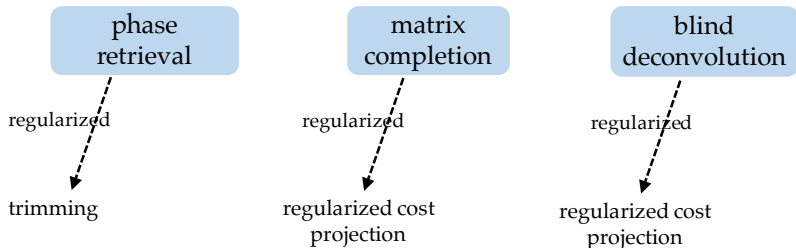
3 representative nonconvex problems

phase
retrieval

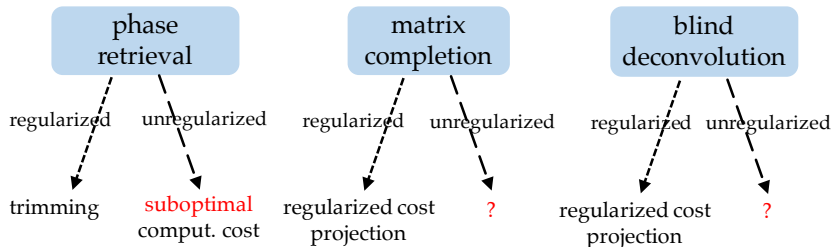
matrix
completion

blind
deconvolution

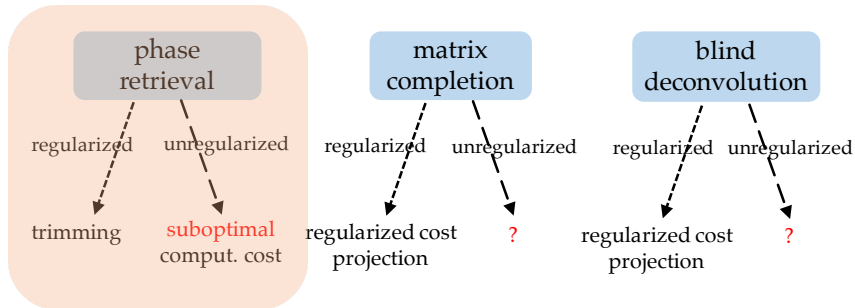
Regularized methods



Regularized vs. **unregularized** methods



Regularized vs. **unregularized** methods



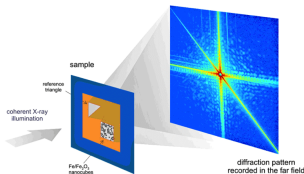
Are unregularized methods suboptimal for nonconvex estimation?

Missing phase problem

Detectors record **intensities** of diffracted rays

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

Fig credit: Stanford SLAC



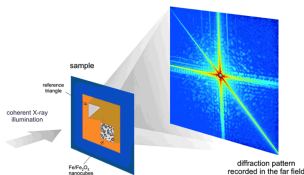
$$\text{intensity of electrical field: } |\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

Missing phase problem

Detectors record **intensities** of diffracted rays

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

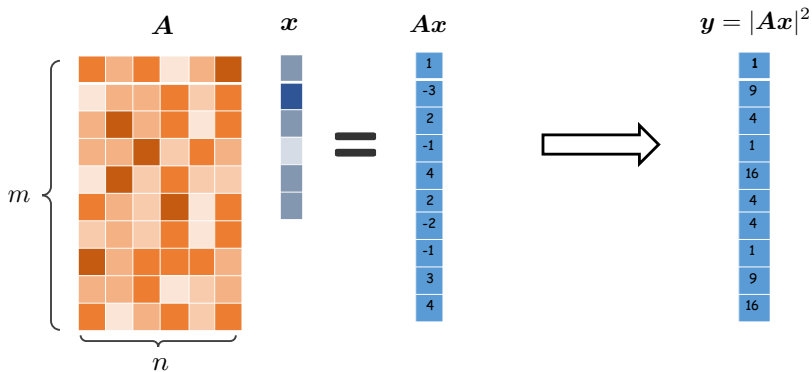
Fig credit: Stanford SLAC



intensity of electrical field: $|\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$

Phase retrieval: recover signal $x(t_1, t_2)$ from intensity $|\hat{x}(f_1, f_2)|^2$

Solving quadratic systems of equations



Recover $\mathbf{x}^\natural \in \mathbb{R}^n$ from m random quadratic measurements

$$y_k = |\mathbf{a}_k^\top \mathbf{x}^\natural|^2, \quad k = 1, \dots, m$$

Assume w.l.o.g. $\|\mathbf{x}^\natural\|_2 = 1$

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

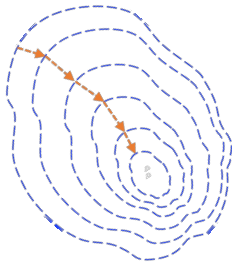
Empirical risk minimization

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\mathbf{a}_k^\top \mathbf{x})^2 - y_k \right]^2$$

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

Empirical risk minimization

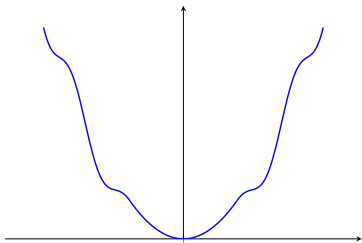
$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\mathbf{a}_k^\top \mathbf{x})^2 - y_k \right]^2$$



- **Initialization by spectral method**
- **Gradient iterations:** for $t = 0, 1, \dots$

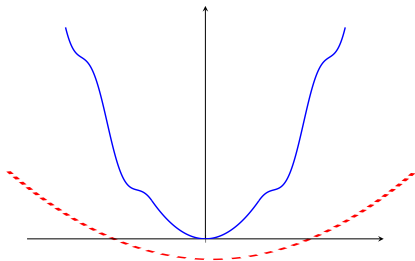
$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t)$$

Gradient descent theory revisited



Two standard conditions that enable geometric convergence of GD

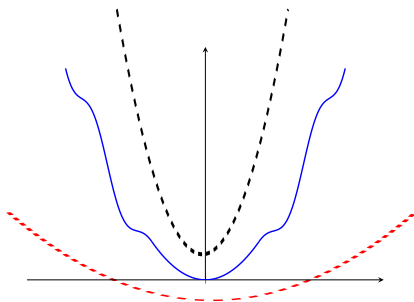
Gradient descent theory revisited



Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity (or regularity condition)

Gradient descent theory revisited



Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity (or regularity condition)
- (local) smoothness

$$\nabla^2 f(\mathbf{x}) \succcurlyeq \mathbf{0} \quad \text{and} \quad \text{is well-conditioned}$$

Gradient descent theory revisited

f is said to be α -strongly convex and β -smooth if

$$\mathbf{0} \preceq \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{x}$$

ℓ_2 **error contraction:** GD with $\eta = 1/\beta$ obeys

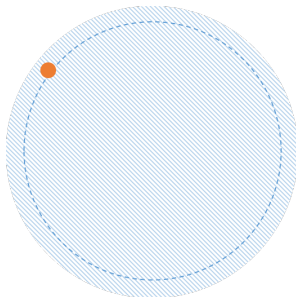
$$\|\mathbf{x}^{t+1} - \mathbf{x}^\natural\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}^\natural\|_2$$

Gradient descent theory revisited

$$\|\mathbf{x}^{t+1} - \mathbf{x}^\natural\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}^t - \mathbf{x}^\natural\|_2$$



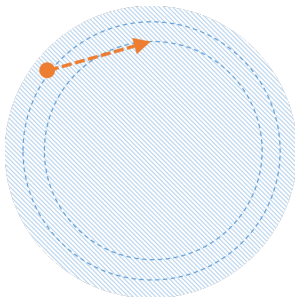
region of local strong convexity + smoothness



Gradient descent theory revisited

$$\|\mathbf{x}^{t+1} - \mathbf{x}^\natural\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}^t - \mathbf{x}^\natural\|_2$$

- region of local strong convexity + smoothness

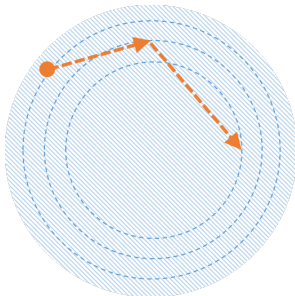


Gradient descent theory revisited

$$\|\mathbf{x}^{t+1} - \mathbf{x}^{\natural}\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}^t - \mathbf{x}^{\natural}\|_2$$



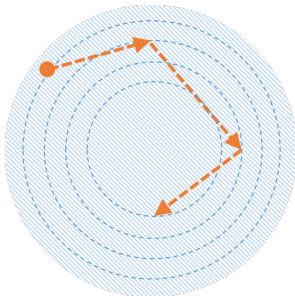
region of local strong convexity + smoothness



Gradient descent theory revisited

$$\|\mathbf{x}^{t+1} - \mathbf{x}^{\natural}\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}^t - \mathbf{x}^{\natural}\|_2$$

- region of local strong convexity + smoothness



Gradient descent theory revisited

$$0 \preceq \alpha I \preceq \nabla^2 f(\mathbf{x}) \preceq \beta I, \quad \forall \mathbf{x}$$

ℓ_2 error contraction: GD with $\eta = 1/\beta$ obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}^\natural\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}^\natural\|_2$$

- Condition number β/α determines rate of convergence

Gradient descent theory revisited

$$0 \preceq \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{x}$$

ℓ_2 **error contraction:** GD with $\eta = 1/\beta$ obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}^\natural\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}^\natural\|_2$$

- Condition number β/α determines rate of convergence
- Attains ε -accuracy within $O(\frac{\beta}{\alpha} \log \frac{1}{\varepsilon})$ iterations

What does this optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

What does this optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

Population level (infinite samples)

$$\mathbb{E}[\nabla^2 f(\mathbf{x})] = 3 \underbrace{\left(\|\mathbf{x}\|_2^2 \mathbf{I} + 2\mathbf{x}\mathbf{x}^\top \right) - \left(\|\mathbf{x}^\natural\|_2^2 \mathbf{I} + 2\mathbf{x}^\natural \mathbf{x}^{\natural\top} \right)}_{\text{locally positive definite and well-conditioned}}$$

Consequence: WF converges within $O(\log \frac{1}{\epsilon})$ iterations if $m \rightarrow \infty$

What does this optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

Finite-sample level ($m \asymp n \log n$)

$$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$$

What does this optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

Finite-sample level ($m \asymp n \log n$)

$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ but ill-conditioned (even locally)
condition number $\asymp n$

What does this optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

Finite-sample level ($m \asymp n \log n$)

$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ but ill-conditioned (even locally)
condition number $\asymp n$

Consequence (Candès et al '14): WF attains ε -accuracy within $O(n \log \frac{1}{\varepsilon})$ iterations if $m \asymp n \log n$

What does this optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

Finite-sample level ($m \asymp n \log n$)

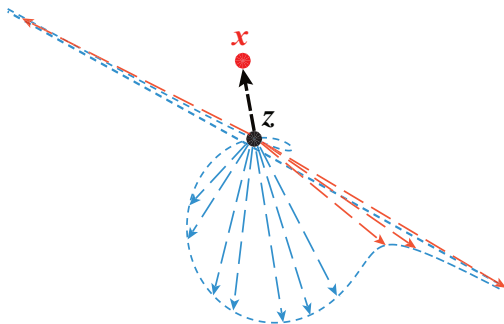
$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ but ill-conditioned (even locally)
condition number $\asymp n$

Consequence (Candès et al '14): WF attains ε -accuracy within $O(n \log \frac{1}{\varepsilon})$ iterations if $m \asymp n \log n$

Too slow ... can we accelerate it?

One solution: truncated WF (Chen, Candès '15)

Regularize / trim gradient components to accelerate convergence



But wait a minute ...

WF converges in $O(n)$ iterations

But wait a minute ...

WF converges in $O(n)$ iterations



Step size taken to be $\eta_t = O(1/n)$

But wait a minute ...

WF converges in $O(n)$ iterations



Step size taken to be $\eta_t = O(1/n)$



This choice is suggested by **generic** optimization theory

But wait a minute ...

WF converges in $O(n)$ iterations



Step size taken to be $\eta_t = O(1/n)$



This choice is suggested by **worst-case** optimization theory

But wait a minute ...

WF converges in $O(n)$ iterations



Step size taken to be $\eta_t = O(1/n)$

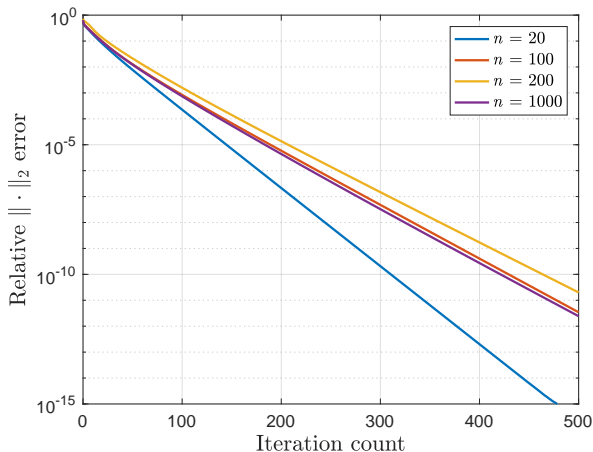


This choice is suggested by **worst-case** optimization theory



Does it capture what really happens?

Numerical surprise with $\eta_t = 0.1$



Vanilla GD (WF) can proceed much more aggressively!

A second look at gradient descent theory

Which region enjoys both strong convexity and smoothness?

$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{k=1}^m \left[3(\mathbf{a}_k^\top \mathbf{x})^2 - (\mathbf{a}_k^\top \mathbf{x})^2 \right] \mathbf{a}_k \mathbf{a}_k^\top$$

A second look at gradient descent theory

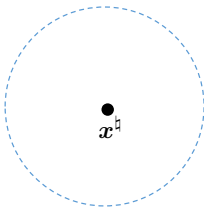
Which region enjoys both strong convexity and smoothness?

$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{k=1}^m \left[3(\mathbf{a}_k^\top \mathbf{x})^2 - (\mathbf{a}_k^\top \mathbf{x})^2 \right] \mathbf{a}_k \mathbf{a}_k^\top$$

- Not smooth if \mathbf{x} and \mathbf{a}_k are too close (coherent)

A second look at gradient descent theory

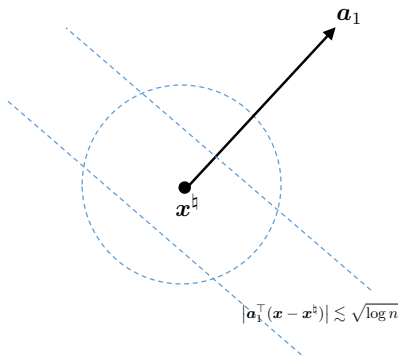
Which region enjoys both strong convexity and smoothness?



- x is not far away from x^h

A second look at gradient descent theory

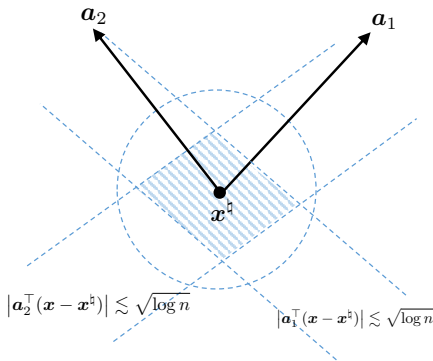
Which region enjoys both strong convexity and smoothness?



- x is not far away from x^{\natural}
- x is incoherent w.r.t. sampling vectors (incoherence region)

A second look at gradient descent theory

Which region enjoys both strong convexity and smoothness?

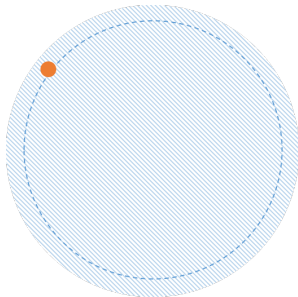


- x is not far away from x^\dagger
- x is incoherent w.r.t. sampling vectors (incoherence region)

A second look at gradient descent theory



region of local strong convexity + smoothness

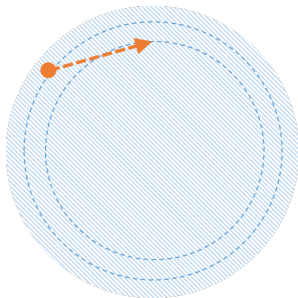


- Prior theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

A second look at gradient descent theory



region of local strong convexity + smoothness

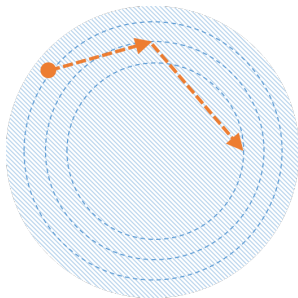


- Prior theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

A second look at gradient descent theory



region of local strong convexity + smoothness

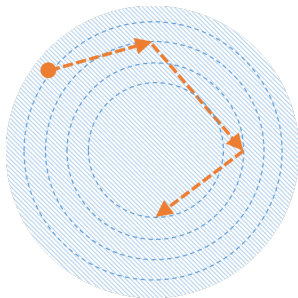


- Prior theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

A second look at gradient descent theory



region of local strong convexity + smoothness

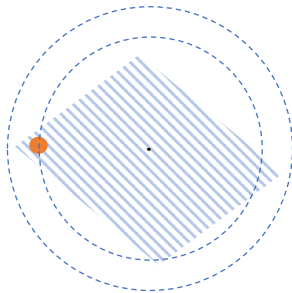
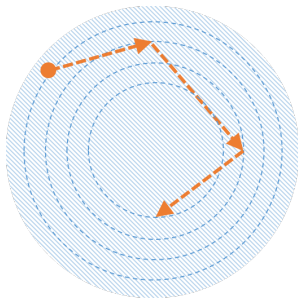


- Prior theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

A second look at gradient descent theory



region of local strong convexity + smoothness

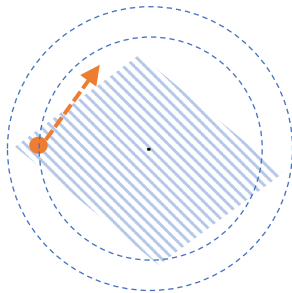
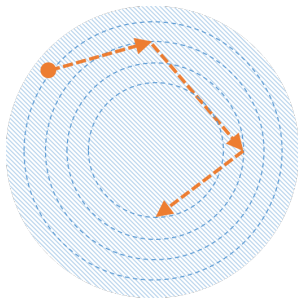


- Prior theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

A second look at gradient descent theory



region of local strong convexity + smoothness

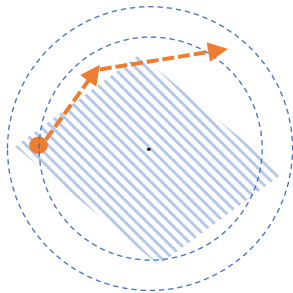
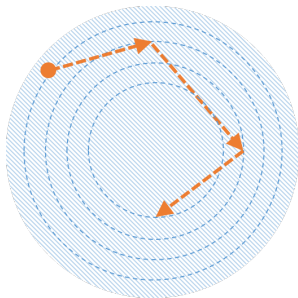


- Prior theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

A second look at gradient descent theory



region of local strong convexity + smoothness

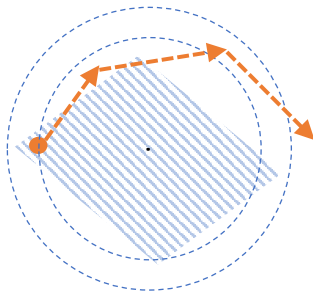
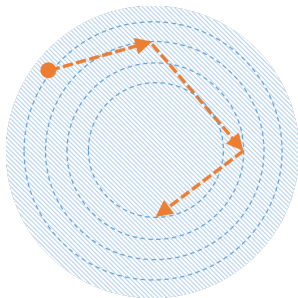


- Prior theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

A second look at gradient descent theory



region of local strong convexity + smoothness

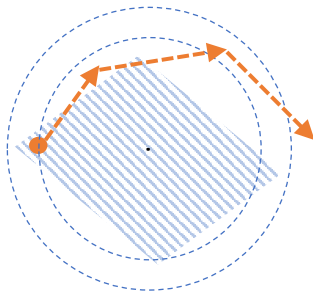
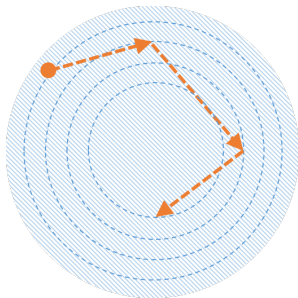


- Prior theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

A second look at gradient descent theory



region of local strong convexity + smoothness

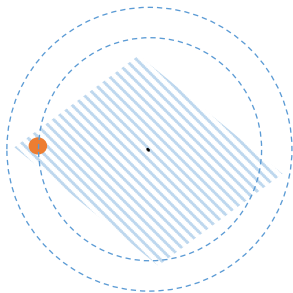


- Prior theory only ensures that iterates remain in ℓ_2 ball but not incoherence region
- *Prior theory enforces regularization to promote incoherence*

Our findings: GD is implicitly regularized



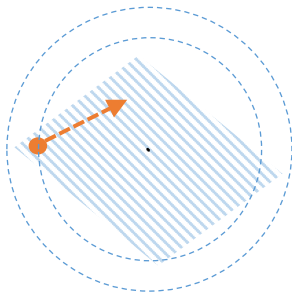
region of local strong convexity + smoothness



Our findings: GD is implicitly regularized



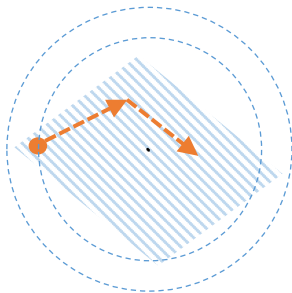
region of local strong convexity + smoothness



Our findings: GD is implicitly regularized



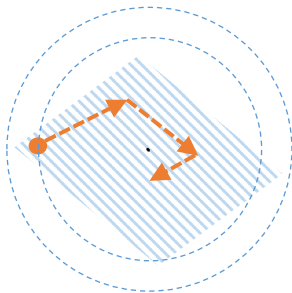
region of local strong convexity + smoothness



Our findings: GD is implicitly regularized



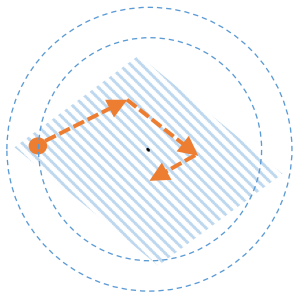
region of local strong convexity + smoothness



Our findings: GD is implicitly regularized



region of local strong convexity + smoothness



GD implicitly forces iterates to remain **incoherent**

Theoretical guarantees

Theorem 1 (Phase retrieval)

Under i.i.d. Gaussian design, WF achieves

- $\max_k |\mathbf{a}_k^\top (\mathbf{x}^t - \mathbf{x}^\natural)| \lesssim \sqrt{\log n} \|\mathbf{x}^\natural\|_2$ (incoherence)

Theoretical guarantees

Theorem 1 (Phase retrieval)

Under i.i.d. Gaussian design, WF achieves

- $\max_k |\mathbf{a}_k^\top (\mathbf{x}^t - \mathbf{x}^\natural)| \lesssim \sqrt{\log n} \|\mathbf{x}^\natural\|_2$ (incoherence)
- $\|\mathbf{x}^t - \mathbf{x}^\natural\|_2 \lesssim (1 - \frac{\eta}{2})^t \|\mathbf{x}^\natural\|_2$ (near-linear convergence)

provided that step size $\eta \asymp \frac{1}{\log n}$ and sample size $m \gtrsim n \log n$.

Theoretical guarantees

Theorem 1 (Phase retrieval)

Under i.i.d. Gaussian design, WF achieves

- $\max_k |\mathbf{a}_k^\top (\mathbf{x}^t - \mathbf{x}^\natural)| \lesssim \sqrt{\log n} \|\mathbf{x}^\natural\|_2$ (incoherence)
- $\|\mathbf{x}^t - \mathbf{x}^\natural\|_2 \lesssim (1 - \frac{\eta}{2})^t \|\mathbf{x}^\natural\|_2$ (near-linear convergence)

provided that step size $\eta \asymp \frac{1}{\log n}$ and sample size $m \gtrsim n \log n$.

- Step size: $\frac{1}{\log n}$ (vs. $\frac{1}{n}$)

Theoretical guarantees

Theorem 1 (Phase retrieval)

Under i.i.d. Gaussian design, WF achieves

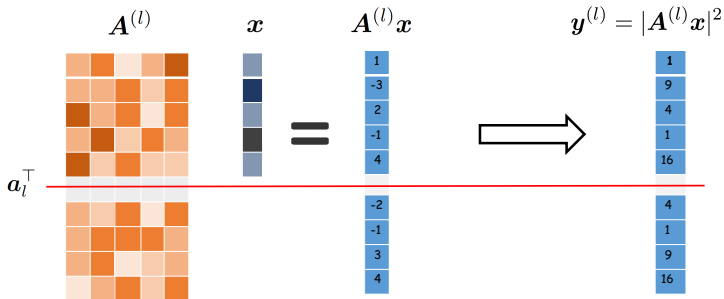
- $\max_k |\mathbf{a}_k^\top (\mathbf{x}^t - \mathbf{x}^\natural)| \lesssim \sqrt{\log n} \|\mathbf{x}^\natural\|_2$ (incoherence)
- $\|\mathbf{x}^t - \mathbf{x}^\natural\|_2 \lesssim (1 - \frac{\eta}{2})^t \|\mathbf{x}^\natural\|_2$ (near-linear convergence)

provided that step size $\eta \asymp \frac{1}{\log n}$ and sample size $m \gtrsim n \log n$.

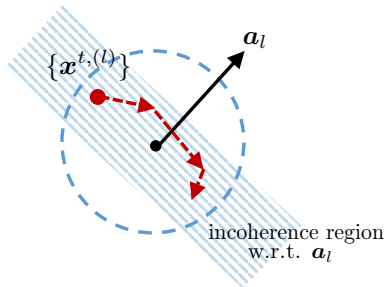
- Step size: $\frac{1}{\log n}$ (vs. $\frac{1}{n}$)
- Computational complexity: $\frac{n}{\log n}$ times faster than existing theory

Key ingredient: leave-one-out analysis

For each $1 \leq l \leq m$, introduce leave-one-out iterates $x^{t,(l)}$ by dropping l th measurement

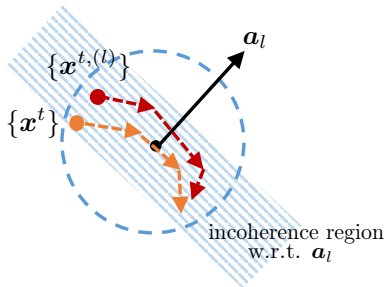


Key ingredient: leave-one-out analysis



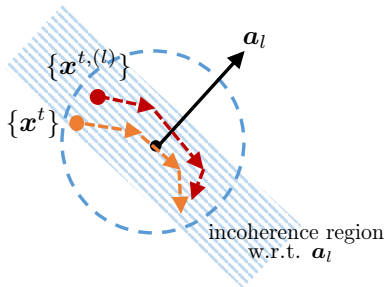
- Leave-one-out iterates $\{x^{t,(l)}\}$ are independent of a_l , and are hence **incoherent** w.r.t. a_l with high prob.

Key ingredient: leave-one-out analysis



- Leave-one-out iterates $\{\mathbf{x}^{t,(l)}\}$ are independent of \mathbf{a}_l , and are hence **incoherent** w.r.t. \mathbf{a}_l with high prob.
- Leave-one-out iterates $\mathbf{x}^{t,(l)} \approx$ true iterates \mathbf{x}^t

Key ingredient: leave-one-out analysis



- Leave-one-out iterates $\{\mathbf{x}^{t,(l)}\}$ are independent of \mathbf{a}_l , and are hence **incoherent** w.r.t. \mathbf{a}_l with high prob.
- Leave-one-out iterates $\mathbf{x}^{t,(l)} \approx$ true iterates \mathbf{x}^t
- $|\mathbf{a}_l^\top (\mathbf{x}^t - \mathbf{x}^{\natural})| \leq |\mathbf{a}_l^\top (\mathbf{x}^{t,(l)} - \mathbf{x}^{\natural})| + |\mathbf{a}_l^\top (\mathbf{x}^t - \mathbf{x}^{t,(l)})|$

This recipe is quite general

Low-rank matrix completion

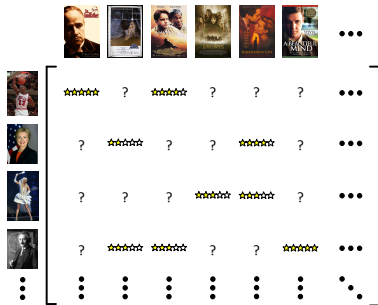


Fig. credit: Candès

Given partial samples Ω of a *low-rank* matrix M , fill in missing entries

Prior art

$$\text{minimize}_{\mathbf{X}} \quad f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left(\mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - M_{j,k} \right)^2$$

Prior art

$$\text{minimize}_{\mathbf{X}} \quad f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left(\mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - M_{j,k} \right)^2$$

Existing theory on gradient descent requires

Prior art

$$\text{minimize}_{\mathbf{X}} \quad f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left(\mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - M_{j,k} \right)^2$$

Existing theory on gradient descent requires

- regularized loss (solve $\min_{\mathbf{X}} f(\mathbf{X}) + R(\mathbf{X})$ instead)
 - e.g. Keshavan, Montanari, Oh '10, Sun, Luo '14, Ge, Lee, Ma '16

Prior art

$$\text{minimize}_{\mathbf{X}} \quad f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left(\mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - M_{j,k} \right)^2$$

Existing theory on gradient descent requires

- regularized loss (solve $\min_{\mathbf{X}} f(\mathbf{X}) + R(\mathbf{X})$ instead)
 - e.g. Keshavan, Montanari, Oh '10, Sun, Luo '14, Ge, Lee, Ma '16
- projection onto set of incoherent matrices
 - e.g. Chen, Wainwright '15, Zheng, Lafferty '16

Theoretical guarantees

Theorem 2 (Matrix completion)

Suppose M is rank- r , incoherent and well-conditioned. *Vanilla gradient descent* (with spectral initialization) achieves ε accuracy

- in $O(\log \frac{1}{\varepsilon})$ iterations

if step size $\eta \lesssim 1/\sigma_{\max}(M)$ and sample size $\gtrsim nr^3 \log^3 n$

Theoretical guarantees

Theorem 2 (Matrix completion)

Suppose M is rank- r , incoherent and well-conditioned. *Vanilla gradient descent* (with spectral initialization) achieves ε accuracy

- in $O(\log \frac{1}{\varepsilon})$ iterations w.r.t. $\|\cdot\|_F$, $\|\cdot\|$, and $\underbrace{\|\cdot\|_{2,\infty}}_{\text{incoherence}}$

if step size $\eta \lesssim 1/\sigma_{\max}(M)$ and sample size $\gtrsim nr^3 \log^3 n$

Theoretical guarantees

Theorem 2 (Matrix completion)

Suppose M is rank- r , incoherent and well-conditioned. *Vanilla gradient descent* (with spectral initialization) achieves ε accuracy

- in $O(\log \frac{1}{\varepsilon})$ iterations w.r.t. $\|\cdot\|_F$, $\|\cdot\|$, and $\underbrace{\|\cdot\|_{2,\infty}}_{\text{incoherence}}$

if step size $\eta \lesssim 1/\sigma_{\max}(M)$ and sample size $\gtrsim nr^3 \log^3 n$

- Byproduct: vanilla GD controls **entrywise error**
— errors are spread out across all entries

Blind deconvolution

image deblurring

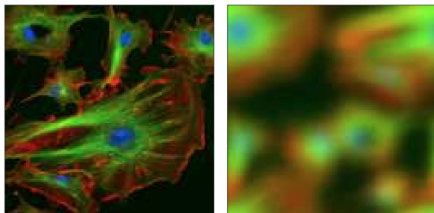


Fig. credit: Romberg

multipath in wireless comm

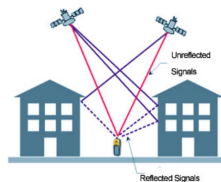


Fig. credit:

EngineeringsALL

Reconstruct two signals from their convolution; equivalently,

$$\text{find } \mathbf{h}, \mathbf{x} \in \mathbb{C}^n \quad \text{s.t.} \quad \mathbf{b}_k^* \mathbf{h} \mathbf{x}^* \mathbf{a}_k = y_k, \quad 1 \leq k \leq m$$

Prior art

$$\text{minimize}_{\mathbf{x}, \mathbf{h}} \quad f(\mathbf{x}, \mathbf{h}) = \sum_{k=1}^m \left| \mathbf{b}_k^* \left(\mathbf{h} \mathbf{x}^* - \mathbf{h}^\dagger \mathbf{x}^{\dagger*} \right) \mathbf{a}_k \right|^2$$

$\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\{\mathbf{b}_k\}$: partial Fourier basis

Prior art

$$\text{minimize}_{\mathbf{x}, \mathbf{h}} \quad f(\mathbf{x}, \mathbf{h}) = \sum_{k=1}^m \left| \mathbf{b}_k^* \left(\mathbf{h} \mathbf{x}^* - \mathbf{h}^\dagger \mathbf{x}^{\dagger*} \right) \mathbf{a}_k \right|^2$$

$\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\{\mathbf{b}_k\}$: partial Fourier basis

Existing theory on gradient descent requires

- regularized loss + projection
 - e.g. Li, Ling, Strohmer, Wei '16, Huang, Hand '17, Ling, Strohmer '17

Prior art

$$\text{minimize}_{\mathbf{x}, \mathbf{h}} \quad f(\mathbf{x}, \mathbf{h}) = \sum_{k=1}^m \left| \mathbf{b}_k^* \left(\mathbf{h} \mathbf{x}^* - \mathbf{h} \mathbf{x} \right) \mathbf{a}_k \right|^2$$

$\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\{\mathbf{b}_k\}$: partial Fourier basis

Existing theory on gradient descent requires

- regularized loss + projection
 - e.g. Li, Ling, Strohmer, Wei '16, Huang, Hand '17, Ling, Strohmer '17
 - requires m iterations even with regularization

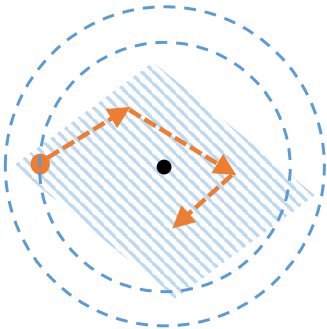
Theoretical guarantees

Theorem 3 (Blind deconvolution)

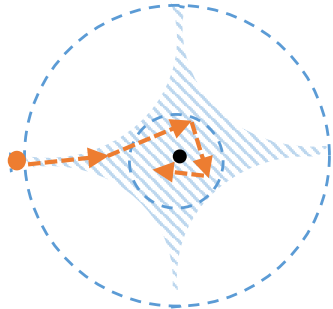
Suppose \mathbf{h}^\natural is incoherent w.r.t. $\{\mathbf{b}_k\}$. *Vanilla gradient descent* (with spectral initialization) achieves ε accuracy in $O(\log \frac{1}{\varepsilon})$ iterations, provided that step size $\eta \lesssim 1$ and sample size $m \gtrsim n \text{poly} \log(m)$.

- Regularization-free
- Converges in $O(\log \frac{1}{\varepsilon})$ iterations (vs. $O(m \log \frac{1}{\varepsilon})$ iterations in prior theory)

Incoherence region in high dimensions



2-dimensional

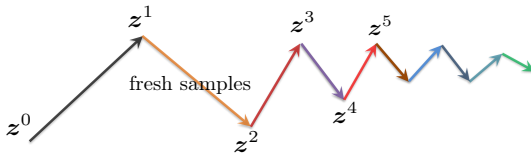


high-dimensional (mental representation)

incoherence region is vanishingly small

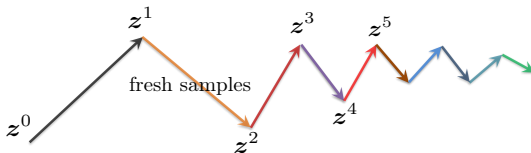
Complicated dependencies across iterations

- Several prior sample-splitting approaches: require **fresh samples** at each iteration; not what we actually run in practice

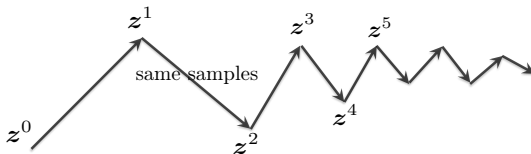


Complicated dependencies across iterations

- Several prior sample-splitting approaches: require **fresh samples** at each iteration; not what we actually run in practice



- This work:** reuses all samples in all iterations



Summary

- **Implicit regularization:** vanilla gradient descent automatically forces iterates to stay *incoherent*

Summary

- **Implicit regularization:** vanilla gradient descent automatically forces iterates to stay *incoherent*
- Enable error controls in a much stronger sense (e.g. *entrywise error control*)

Paper:

“Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution”, Cong Ma, Kaizheng Wang, Yuejie Chi, Yuxin Chen, arXiv:1711.10467