

Refined analysis of local convergence: implicit regularization



Cong Ma

University of Chicago, Winter 2024

Revisit phase retrieval

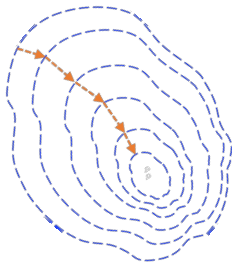
$$\text{given: } y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2, \quad 1 \leq k \leq m$$

↓

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

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

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



- **spectral initialization:** $\mathbf{x}^0 \leftarrow$ leading eigenvector of certain data matrix
- **gradient descent:**

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t), \quad t = 0, 1, \dots$$

First theory of WF

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) := \min\{\|\mathbf{x}^t \pm \mathbf{x}^*\|_2\}$$

Theorem 9.1 (Candès, Li, Soltanolkotabi '14)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{4}\right)^{t/2} \|\mathbf{x}^*\|_2,$$

with high prob., provided that step size $\eta \lesssim 1/n$ and sample size: $m \gtrsim n \log n$.

- Iteration complexity: $O(n \log \frac{1}{\epsilon})$
- Sample complexity: $O(n \log n)$
- Derived based on (worst-case) local geometry

What does 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 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 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 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 $m \asymp n \log n$

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations



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

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations



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



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

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations



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

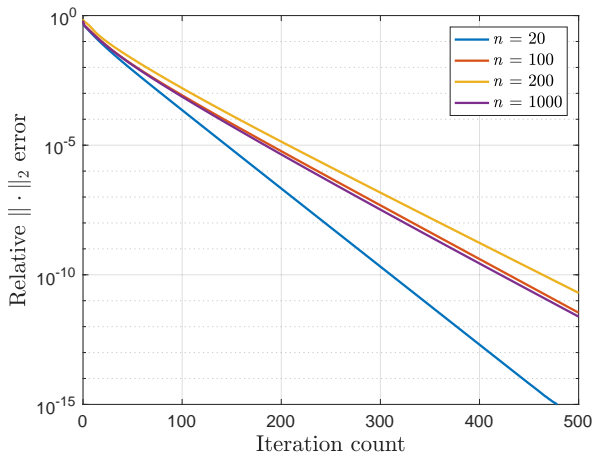


This choice is suggested by **worst-case** optimization theory



Does it capture what really happens?

Numerical efficiency with $\eta_t = 0.1$



Vanilla GD (WF) converges fast for a constant step size!

Improved theory of WF

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) := \min\{\|\mathbf{x}^t \pm \mathbf{x}^*\|_2\}$$

Theorem 9.2 (Ma, Wang, Chi, Chen '17)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{2}\right)^t \|\mathbf{x}^*\|_2$$

with high prob., provided that step size $\eta \asymp 1/\log n$ and sample size $m \gtrsim n \log n$.

- Iteration complexity: $O(n \log \frac{1}{\epsilon}) \searrow O(\log n \log \frac{1}{\epsilon})$
- Sample complexity: $O(n \log n)$
- Derived based on finer analysis of GD trajectory

A second look at gradient descent theory

Which local 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

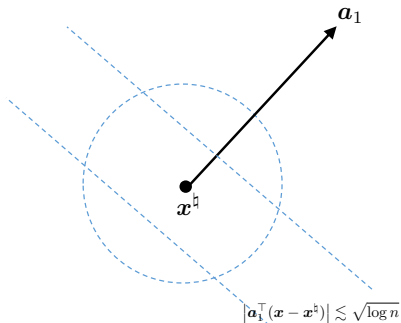
Which local region enjoys both strong convexity and smoothness?

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

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

A second look at gradient descent theory

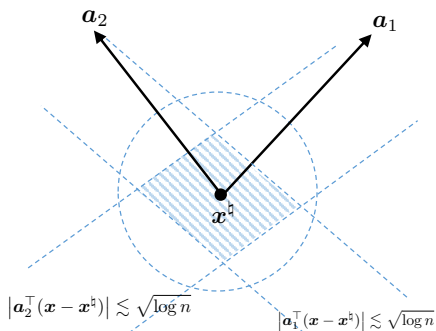
Which local region enjoys both strong convexity and smoothness?



- x is incoherent w.r.t. sampling vectors $\{a_k\}$ (incoherence region)

A second look at gradient descent theory

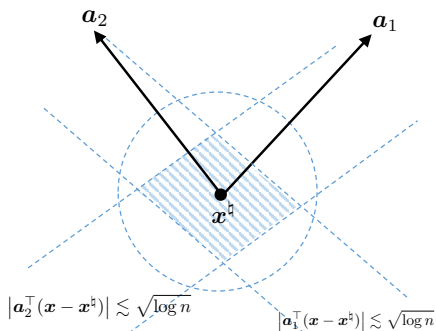
Which local region enjoys both strong convexity and smoothness?



- x is incoherent w.r.t. sampling vectors $\{a_k\}$ (incoherence region)

A second look at gradient descent theory

Which local region enjoys both strong convexity and smoothness?

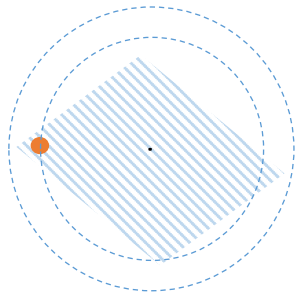


- x is incoherent w.r.t. sampling vectors $\{a_k\}$ (incoherence region)

Prior works suggest enforcing **regularization** (e.g. truncation, projection, regularized loss) to promote incoherence

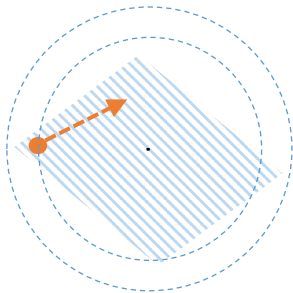
Encouraging message: GD is implicitly regularized

- region of local strong convexity + smoothness



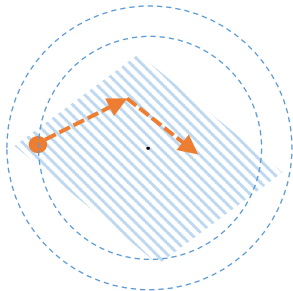
Encouraging message: GD is implicitly regularized

- region of local strong convexity + smoothness



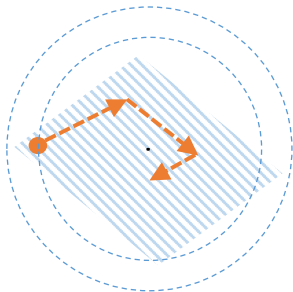
Encouraging message: GD is implicitly regularized

- region of local strong convexity + smoothness



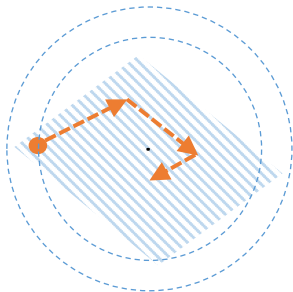
Encouraging message: GD is implicitly regularized

- region of local strong convexity + smoothness



Encouraging message: GD is implicitly regularized

- region of local strong convexity + smoothness



GD implicitly forces iterates to remain **incoherent with** $\{\mathbf{a}_k\}$

$$\max_k |\mathbf{a}_k^\top (\mathbf{x}^t - \mathbf{x}^*)| \lesssim \sqrt{\log n} \|\mathbf{x}^*\|_2, \quad \forall t$$

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

Theoretical guarantees for local refinement stage

Theorem 9.3 (Ma, Wang, Chi, Chen '17)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

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

Theoretical guarantees for local refinement stage

Theorem 9.3 (Ma, Wang, Chi, Chen '17)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

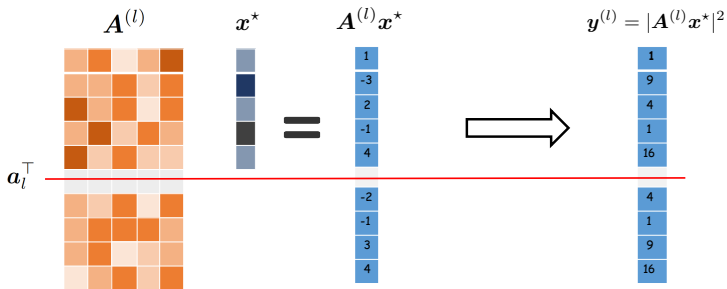
- $\max_k |\mathbf{a}_k^\top \mathbf{x}^t| \lesssim \sqrt{\log n} \|\mathbf{x}^*\|_2$ (incoherence)
- $\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim (1 - \frac{\eta}{2})^t \|\mathbf{x}^*\|_2$ (linear convergence)

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

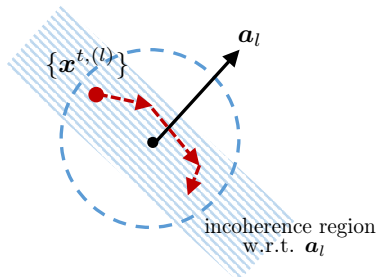
- Attains ε accuracy within $O(\log n \log \frac{1}{\varepsilon})$ iterations

Key proof idea: leave-one-out analysis

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

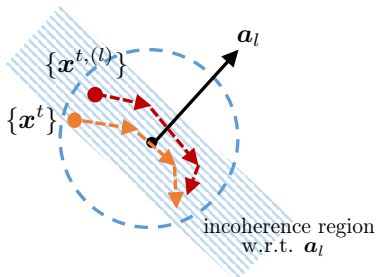


Key proof idea: leave-one-out analysis



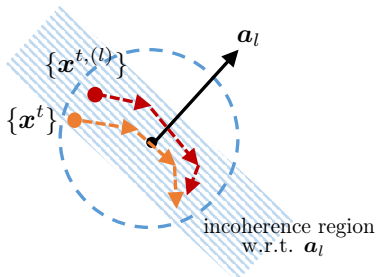
- Leave-one-out iterate $\mathbf{x}^{t,(l)}$ is independent of \mathbf{a}_l

Key proof idea: leave-one-out analysis



- Leave-one-out iterate $x^{t,(l)}$ is independent of \mathbf{a}_l
- Leave-one-out iterate $x^{t,(l)} \approx$ true iterate x^t

Key proof idea: leave-one-out analysis

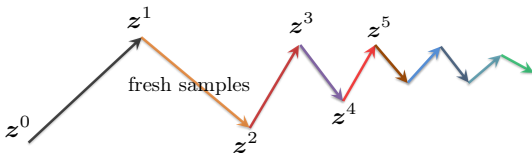


- Leave-one-out iterate $x^{t,(l)}$ is independent of \mathbf{a}_l
- Leave-one-out iterate $x^{t,(l)} \approx$ true iterate x^t

$\implies x^t$ is nearly independent of \mathbf{a}_l
nearly orthogonal to

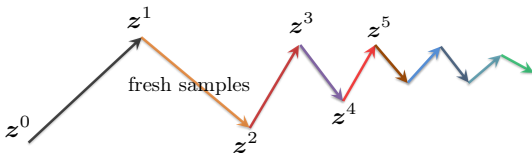
No need of sample splitting

- Several prior works use sample-splitting: require **fresh samples** at each iteration; not practical but helps analysis

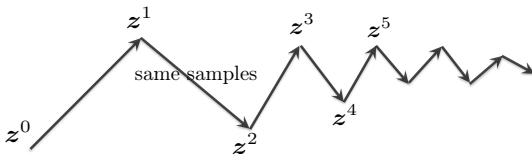


No need of sample splitting

- Several prior works use sample-splitting: require **fresh samples** at each iteration; not practical but helps analysis



- In contrast, we reuse all samples in all iterations



Architecture of the proof

Local geometry

Lemma 9.4

Suppose $m \geq c_0 n \log n$ for some sufficiently large constant $c_0 > 0$.
With high probability,

$$\nabla^2 f(\mathbf{x}) \succeq (1/2) \cdot \mathbf{I}_n$$

holds simultaneously for all $\mathbf{x} \in \mathbb{R}^n$ satisfying $\|\mathbf{x} - \mathbf{x}^*\|_2 \leq 2C_1$; and

$$\nabla^2 f(\mathbf{x}) \preceq (5C_2(10 + C_2) \log n) \cdot \mathbf{I}_n$$

holds simultaneously for all $\mathbf{x} \in \mathbb{R}^n$ obeying

$$\|\mathbf{x} - \mathbf{x}^*\|_2 \leq 2C_1, \tag{9.1a}$$

$$\max_{1 \leq j \leq m} \left| \mathbf{a}_j^\top (\mathbf{x} - \mathbf{x}^*) \right| \leq C_2 \sqrt{\log n}. \tag{9.1b}$$

Error contraction

Lemma 9.5

If \mathbf{x}^t obeys the conditions (9.1), whp. one has

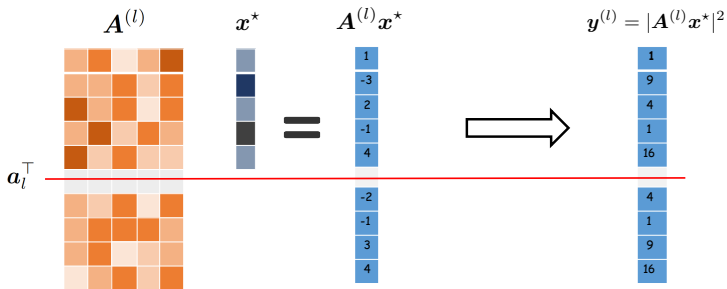
$$\left\| \mathbf{x}^{t+1} - \mathbf{x}^* \right\|_2 \leq (1 - \eta/2) \left\| \mathbf{x}^t - \mathbf{x}^* \right\|_2 \quad (9.2)$$

provided that the step size satisfies $0 < \eta \leq 1 / [5C_2 (10 + C_2) \log n]$.

— how to insure incoherence?

Key proof idea: leave-one-out analysis

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



Induction hypotheses

We aim at proving the following claims using induction

$$\|\mathbf{x}^t - \mathbf{x}^*\|_2 \leq C_1, \quad (9.3a)$$

$$\max_{1 \leq l \leq m} \|\mathbf{x}^t - \mathbf{x}^{t,(l)}\|_2 \leq C_3 \sqrt{\frac{\log n}{n}} \quad (9.3b)$$

$$\max_{1 \leq j \leq m} \left| \mathbf{a}_j^\top (\mathbf{x}^t - \mathbf{x}^*) \right| \leq C_2 \sqrt{\log n}. \quad (9.3c)$$

Proximity between \mathbf{x}^t and $\mathbf{x}^{t,(l)}$

Lemma 9.6

Suppose that the sample size obeys $m \geq Cn \log n$ for some sufficiently large constant $C > 0$ and that the stepsize obeys $0 < \eta < 1/[5C_2(10 + C_2) \log n]$. Then whp., one has

$$\max_{1 \leq l \leq m} \left\| \mathbf{x}^{t+1} - \mathbf{x}^{t+1,(l)} \right\|_2 \leq C_3 \sqrt{\frac{\log n}{n}}. \quad (9.4)$$

Incoherence of leave-one-out iterates

By construction, $\mathbf{x}^{t+1,(l)}$ is statistically independent of the sampling vector \mathbf{a}_l . One thus has

$$\begin{aligned} \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}^{t+1,(l)} - \mathbf{x}^*) \right| &\leq 5\sqrt{\log n} \|\mathbf{x}^{t+1,(l)} - \mathbf{x}^*\|_2 \\ &\stackrel{(i)}{\leq} 5\sqrt{\log n} \left(\|\mathbf{x}^{t+1,(l)} - \mathbf{x}^{t+1}\|_2 + \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2 \right) \\ &\stackrel{(ii)}{\leq} 5\sqrt{\log n} \left(C_3 \sqrt{\frac{\log n}{n}} + C_1 \right) \\ &\leq C_4 \sqrt{\log n} \end{aligned} \tag{9.5}$$

holds for some constant $C_4 \geq 6C_1 > 0$ and n sufficiently large. Here, (i) comes from the triangle inequality, and (ii) arises from the proximity bound (9.4) and the conclusion (9.2).

Combining the bounds

$$\begin{aligned} \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}^{t+1} - \mathbf{x}^*) \right| &\leq \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}^{t+1} - \mathbf{x}^{t+1,(l)}) \right| \\ &\quad + \max_{1 \leq l \leq m} \left| \mathbf{a}_l^\top (\mathbf{x}^{t+1,(l)} - \mathbf{x}^*) \right| \\ &\stackrel{(i)}{\leq} \max_{1 \leq l \leq m} \|\mathbf{a}_l\|_2 \|\mathbf{x}^{t+1} - \mathbf{x}^{t+1,(l)}\|_2 + C_4 \sqrt{\log n} \\ &\stackrel{(ii)}{\leq} \sqrt{6n} \cdot C_3 \sqrt{\frac{\log n}{n}} + C_4 \sqrt{\log n} \leq C_2 \sqrt{\log n} \end{aligned}$$

Another example: Low-rank matrix completion

Low-rank matrix completion

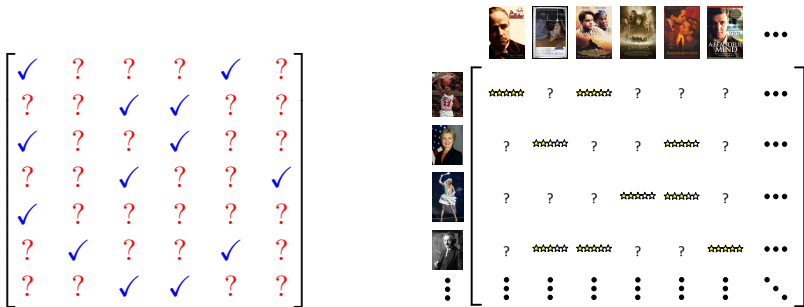
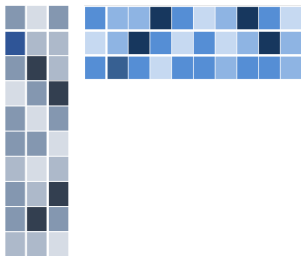


figure credit: Candès

- consider a low-rank matrix $M^* = U^* \Sigma^* U^{*\top}$
- each entry $M_{i,j}^*$ is observed independently with prob. p
- **Goal:** estimate M^*

A natural least-squares loss

Represent low-rank matrix by $\mathbf{X}\mathbf{X}^\top$ with $\underbrace{\mathbf{X} \in \mathbb{R}^{n \times r}}_{\text{low-rank factor}}$



$$\underset{\mathbf{X} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(\mathbf{X}) = \sum_{(i,j) \in \Omega} \left[(\mathbf{X}\mathbf{X}^\top)_{i,j} - M_{i,j}^* \right]^2$$

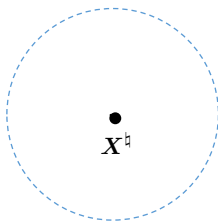
—how does local geometry look like?

Incoherence region

Which region enjoys both restricted strong convexity and smoothness?

Incoherence region

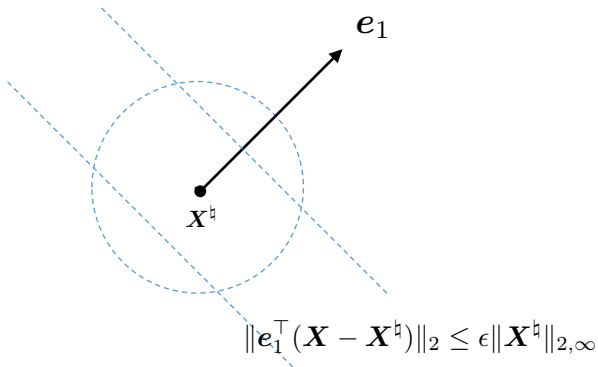
Which region enjoys both restricted strong convexity and smoothness?



- X is not far away from X^h in Euclidean metric

Incoherence region

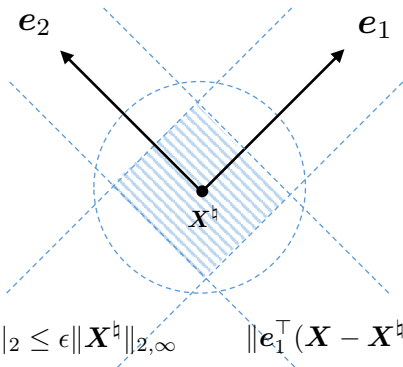
Which region enjoys both restricted strong convexity and smoothness?



- X is not far away from X^b in Euclidean metric
- X is incoherent w.r.t. sampling basis (incoherence region)

Incoherence region

Which region enjoys both restricted strong convexity and smoothness?



$$\|e_2^\top (\mathbf{X} - \mathbf{X}^\dagger)\|_2 \leq \epsilon \|\mathbf{X}^\dagger\|_{2,\infty} \quad \|e_1^\top (\mathbf{X} - \mathbf{X}^\dagger)\|_2 \leq \epsilon \|\mathbf{X}^\dagger\|_{2,\infty}$$

- \mathbf{X} is not far away from \mathbf{X}^\dagger in Euclidean metric
- \mathbf{X} is incoherent w.r.t. sampling basis (incoherence region)

Local geometry of $f(\cdot)$

Lemma 9.7

Suppose that $n^2 p \geq C \kappa^2 \mu r n \log n$ for some sufficiently large constant $C > 0$. Then with high probability, the Hessian $\nabla^2 f(\mathbf{X})$ obeys

$$\begin{aligned} \text{vec}(\mathbf{V})^\top \nabla^2 f(\mathbf{X}) \text{vec}(\mathbf{V}) &\geq \frac{\sigma_{\min}}{2} \|\mathbf{V}\|_F^2 \\ \|\nabla^2 f(\mathbf{X})\| &\leq \frac{5}{2} \sigma_{\max} \end{aligned}$$

for all \mathbf{X} , $\mathbf{V} = \mathbf{Y} \mathbf{H}_Y - \mathbf{X}^*$ s.t.

$$\mathbf{H}_Y := \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{Y} \mathbf{R} - \mathbf{X}^*\|_F,$$

$$\|\mathbf{X} - \mathbf{X}^*\|_{2,\infty} \leq \epsilon \|\mathbf{X}^*\|_{2,\infty},$$

where $\epsilon \ll 1/\sqrt{\kappa^3 \mu r \log^2 n}$.

Restricted local strong convexity

- Due to rotation ambiguity, $f(\cdot)$ cannot be strongly convex along every direction; it is strongly convex along specific directions $V = YH_Y - X^*$
- Instead of ℓ_F ball, $f(X)$ is strongly convex in a local $\ell_{2,\infty}$ ball; X needs to be incoherent in the sense that

$$\|X\|_{2,\infty} \lesssim \sqrt{\frac{\mu r}{n}} \|X^*\|$$

Revisit Incoherence

Definition 9.8

Fix an orthonormal matrix $U^* \in \mathbb{R}^{n \times r}$. Define its incoherence to be

$$\mu(U^*) := \frac{n \|U^*\|_{2,\infty}^2}{r}$$

—recover incoherence of eigenvector when $r = 1$

- For $M^* = U^* \Sigma^* U^{*\top}$, define $\mu(M^*) := \mu(U^*)$

Existing solutions to guarantee incoherence

- regularized loss (solve $\text{minimize}_{\mathbf{X}} f(\mathbf{X}) + R(\mathbf{X})$ instead)
 - e.g. Keshavan, Montanari, Oh '10, Sun, Luo '14, Ge, Lee, Ma '16, Chen, Li '17

Existing solutions to guarantee incoherence

- regularized loss (solve $\text{minimize}_{\mathbf{X}} f(\mathbf{X}) + R(\mathbf{X})$ instead)
 - e.g. Keshavan, Montanari, Oh '10, Sun, Luo '14, Ge, Lee, Ma '16, Chen, Li '17

- projection onto set of incoherent matrices
 - e.g. Chen, Wainwright '15, Zheng, Lafferty '16

Projected gradient descent for matrix completion

- (1) **Projected spectral initialization:** let $U^0 \Sigma^0 U^{0\top}$ be rank- r eigendecomposition of

$$\frac{1}{p} \mathcal{P}_\Omega(\mathbf{Y}).$$

and set $\mathbf{Z}^0 = U^0 (\Sigma^0)^{1/2}$, and incoherence set

$$\mathcal{C} := \{\mathbf{X} \mid \|\mathbf{X}\|_{2,\infty} \leq \sqrt{\frac{2\mu r}{n}} \|\mathbf{Z}^0\|\}$$

let $\mathbf{X}^0 = \mathcal{P}_\mathcal{C}(\mathbf{Z}^0)$

- (2) **Projected gradient descent updates:**

$$\mathbf{X}^{t+1} = \mathcal{P}_\mathcal{C}(\mathbf{X}^t - \eta_t \nabla f(\mathbf{X}^t)), \quad t = 0, 1, \dots$$

Projection operator

Projection onto can be implemented via a row-wise “clipping operation”

$$[\mathcal{P}_C(\mathbf{X})]_{i,\cdot} = \min \left\{ 1, \sqrt{\frac{2\mu r}{n}} \frac{\|\mathbf{Z}^0\|}{\|\mathbf{X}_{i,\cdot}\|_2} \right\} \cdot \mathbf{X}_{i,\cdot}.$$

Performance guarantees

Theorem 9.9

Suppose that $n^2 p \geq c_0 \mu^2 r^2 \kappa^2 n \log n$ for some large constant $c_0 > 0$. With high probability, one has for all $t \geq 0$

$$\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\|_{\text{F}}^2 \leq \left(1 - \frac{c_1}{\mu^2 r^2 \kappa^2}\right)^t \sigma_r(\mathbf{M}^*),$$

provided that step size is chosen as $\eta \asymp \frac{1}{\mu^2 r^2 \kappa \sigma_1(\mathbf{M}^*)}$

Here \mathbf{Q}^t is the optimal alignment matrix between \mathbf{X}^t and \mathbf{X}^*

$$\mathbf{Q}^t := \operatorname{argmin}_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{X}^t \mathbf{R} - \mathbf{X}^*\|_{\text{F}}$$

Regularity condition

Key to prove convergence is the following regularity condition

Lemma 9.10

Suppose that $n^2 p \geq \mu^2 r^2 \kappa^2 n \log n$. Then with high probability, for all $\mathbf{X} \in \mathcal{C}$, and $\|\mathbf{X} - \mathbf{X}^* \mathbf{H}\|_{\text{F}}^2 \leq \frac{1}{16} \sigma_r(\mathbf{M}^*)$ f obeys

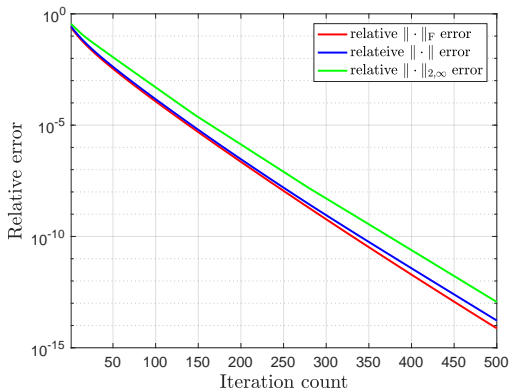
$$\begin{aligned} \langle \nabla f(\mathbf{X}), \mathbf{X} - \mathbf{X}^* \mathbf{H} \rangle &\geq \frac{99}{512} \sigma_r(\mathbf{M}^*) \|\mathbf{X} - \mathbf{X}^* \mathbf{H}\|_{\text{F}}^2 \\ &\quad + \frac{1}{13196 \mu^2 r^2 \kappa \sigma_1(\mathbf{M}^*)} \|\nabla f(\mathbf{X})\|_{\text{F}}^2 \end{aligned}$$

Here \mathbf{H} is optimal alignment matrix

Is regularization necessary for nonconvex matrix completion?

Numerical surprise with unregularized GD

$$n = 1000, r = 10, p = 0.1, \eta = 0.2$$

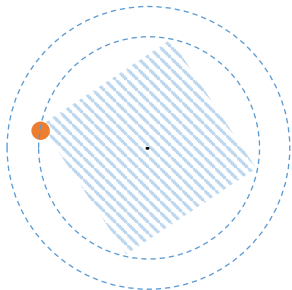


Vanilla GD without regularization converges fast for MC!

Our findings: GD is implicitly regularized

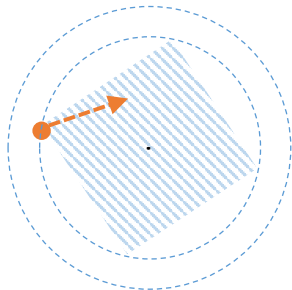


region of local strong convexity + smoothness



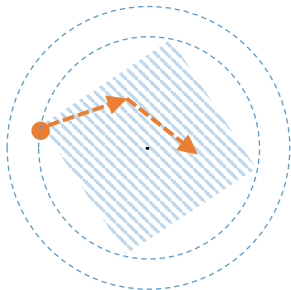
Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



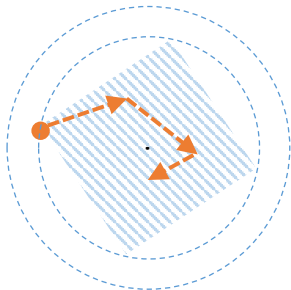
Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



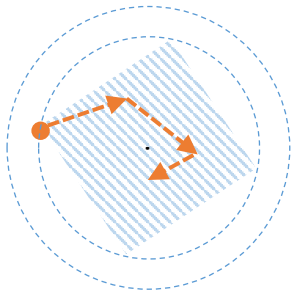
Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



Our findings: GD is implicitly regularized

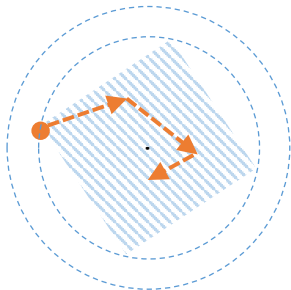
- region of local strong convexity + smoothness



GD implicitly forces iterates to remain **incoherent**

Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



GD implicitly forces iterates to remain **incoherent**

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

Theoretical guarantees

Theorem 9.11 (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 9.11 (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 9.11 (Matrix completion)

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

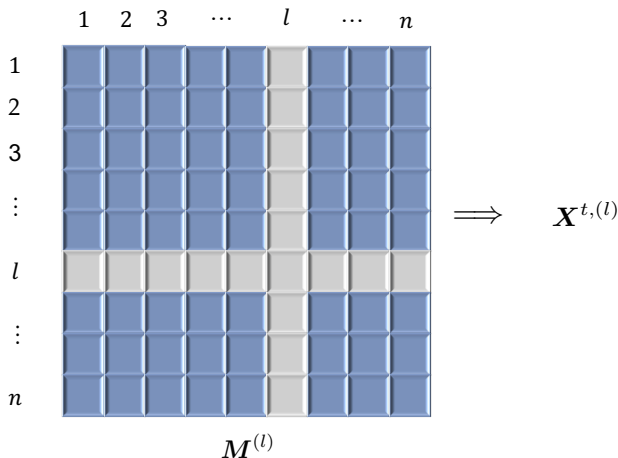
- in $O(\log \frac{1}{\varepsilon})$ iterations w.r.t. $\|\cdot\|_F$, $\|\cdot\|$, and $\underbrace{\|\cdot\|_{2,\infty}}_{\text{incoherence}}$

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

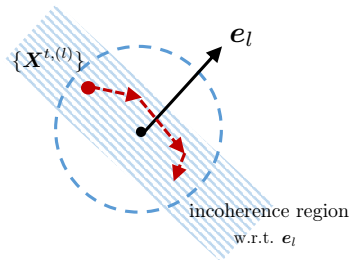
- *Byproduct: vanilla GD controls **entrywise error***
— errors are spread out across all entries

Key ingredient: leave-one-out analysis

For each $1 \leq l \leq n$, introduce leave-one-out iterates $\mathbf{X}^{t,(l)}$ by replacing l th row and column with **true** values

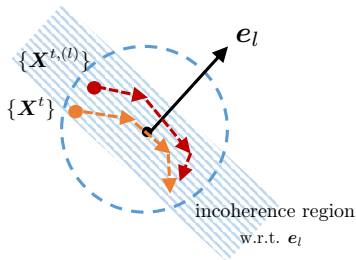


Key ingredient: leave-one-out analysis



- $X^{t,(l)}$ contains more information of l th row of X^\dagger ; indep. of randomness in l th row

Key ingredient: leave-one-out analysis



- $X^{t,(l)}$ contains more information of l th row of X^\dagger ; indep. of randomness in l th row
- Leave-one-out iterates $X^{t,(l)} \approx$ true iterates X^t