

# **Refined analysis of local convergence: implicit regularization**



Cong Ma

University of Chicago, Autumn 2021

# A natural least-squares formulation

---

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

$\Downarrow$

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

# A natural least-squares formulation

---

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

$\Downarrow$

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

- **pros:** often exact as long as sample size is sufficiently large

# A natural least-squares formulation

---

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

$\Downarrow$

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

- **pros:** often exact as long as sample size is sufficiently large
- **cons:**  $f(\cdot)$  is highly nonconvex  
 $\longrightarrow$  *computationally challenging!*

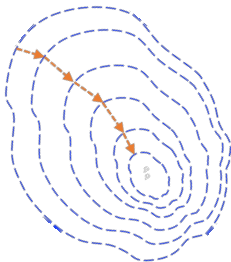
# Wirtinger flow (Candès, Li, Soltanolkotabi '14)

---

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

# Wirtinger flow (Candès, Li, Soltanolkotabi '14)

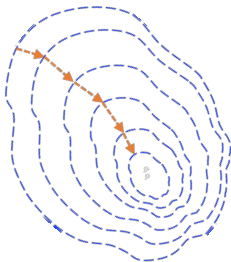
$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m \left[ (\mathbf{a}_k^\top \mathbf{x})^2 - y_k \right]^2$$



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

# 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

---

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

## Theorem 9.1 (Candès, Li, Soltanolkotabi '14)

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

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

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



# First theory of WF

---

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

## Theorem 9.1 (Candès, Li, Soltanolkotabi '14)

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

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

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

- Iteration complexity:  $O(n \log \frac{1}{\epsilon})$

# First theory of WF

---

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

## Theorem 9.1 (Candès, Li, Soltanolkotabi '14)

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

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

*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:  $O(n \log n)$

# First theory of WF

---

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

## Theorem 9.1 (Candès, Li, Soltanolkotabi '14)

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

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

*with high prob., provided that step size and sample size: .*

- 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  $\varepsilon$ -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?

# Improved theory of WF

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

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

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

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

*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

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

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

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

**Consequence (Candès et al '14):** WF attains  $\varepsilon$ -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  $\eta = O(1/n)$

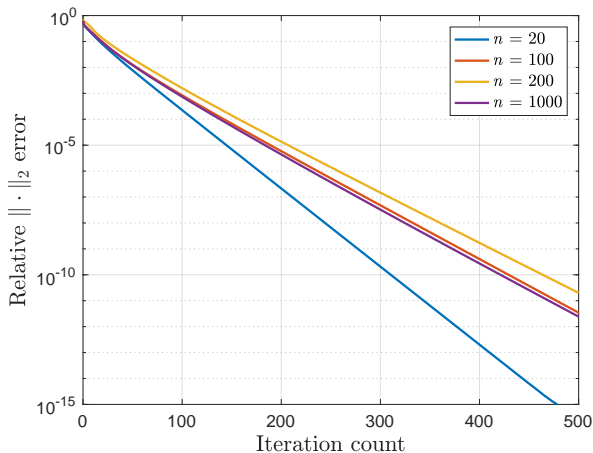


This choice is suggested by **worst-case** optimization theory



Does it capture what really happens?

# Numerical efficiency with $\eta_t = 0.1$



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

## A second look at gradient descent theory

---

Which local region enjoys both strong convexity and smoothness?

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

## A second look at gradient descent theory

---

Which local region enjoys both strong convexity and smoothness?

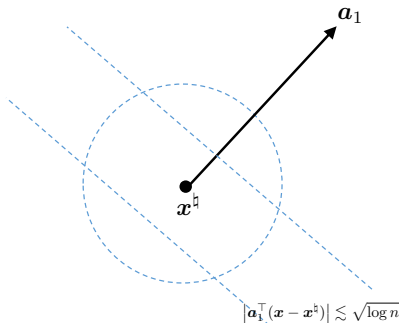
$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{k=1}^m \left[ 3(\mathbf{a}_k^\top \mathbf{x})^2 - (\mathbf{a}_k^\top \mathbf{x}^\star)^2 \right] \mathbf{a}_k \mathbf{a}_k^\top$$

- 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?

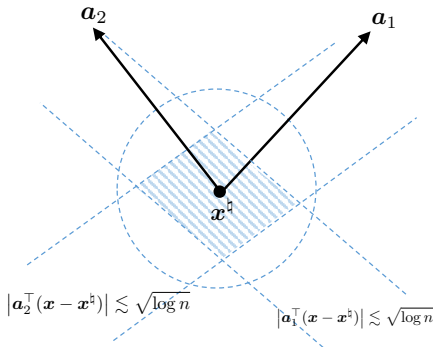


- $x$  is incoherent w.r.t. sampling vectors  $\{a_k\}$  (incoherence region)



## A second look at gradient descent theory

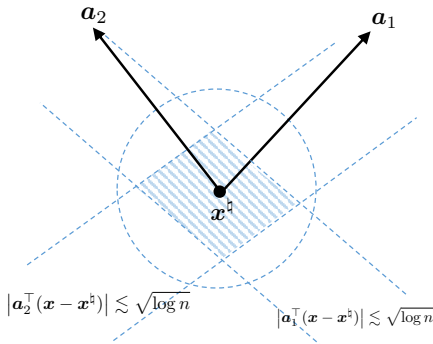
Which local region enjoys both strong convexity and smoothness?



- $x$  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)

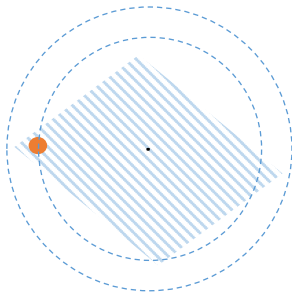
Prior works suggest enforcing **regularization** (e.g. truncation, projection, regularized loss) to promote incoherence

# Encouraging message: GD is implicitly regularized

---



region of local strong convexity + smoothness

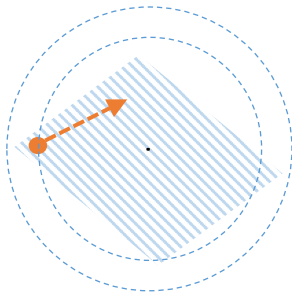


# Encouraging message: GD is implicitly regularized

---



region of local strong convexity + smoothness

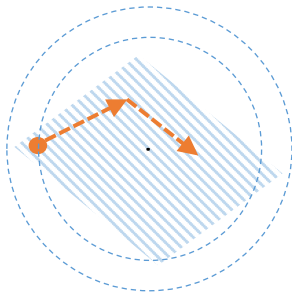


# Encouraging message: GD is implicitly regularized

---



region of local strong convexity + smoothness

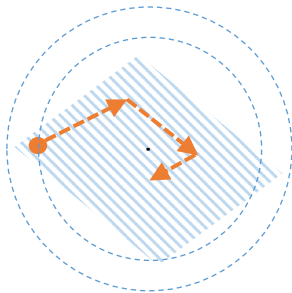


# Encouraging message: GD is implicitly regularized

---



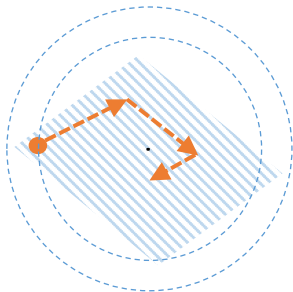
region of local strong convexity + smoothness



# Encouraging message: GD is implicitly regularized



region of local strong convexity + smoothness



GD implicitly forces iterates to remain **incoherent with  $\{a_k\}$**

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

- 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*

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

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

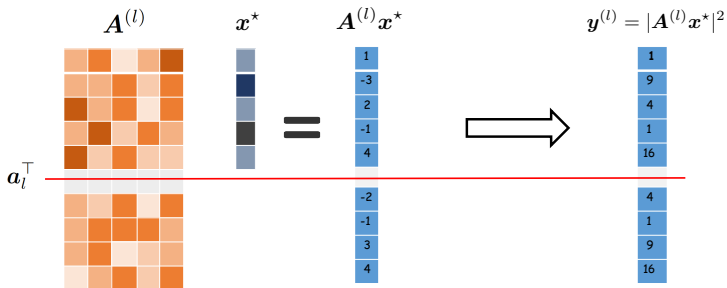
*provided that step size  $\eta \asymp 1/\log n$  and sample size  $m \gtrsim n \log n$ .*

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

# Key proof idea: leave-one-out analysis

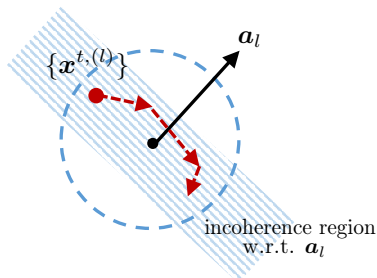
---

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

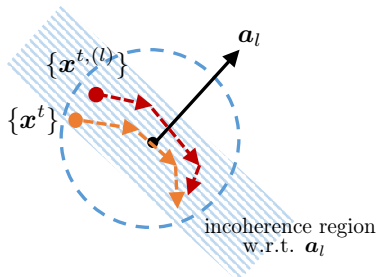
---



- Leave-one-out iterate  $x^{t,(l)}$  is independent of  $a_l$

# Key proof idea: leave-one-out analysis

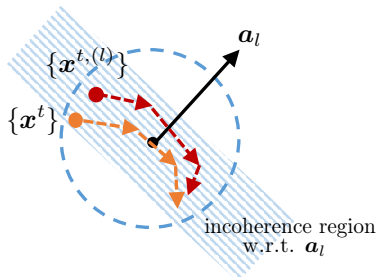
---



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

# Key proof idea: leave-one-out analysis

---



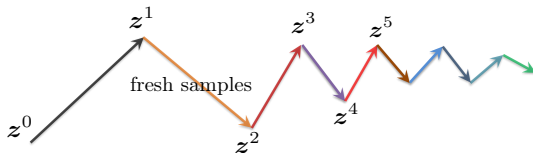
- Leave-one-out iterate  $x^{t,(l)}$  is independent of  $a_l$
- Leave-one-out iterate  $x^{t,(l)} \approx$  true iterate  $x^t$

$\implies x^t$  is nearly independent of  $a_l$   
nearly orthogonal to

# No need of sample splitting

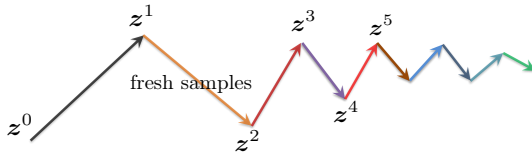
---

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

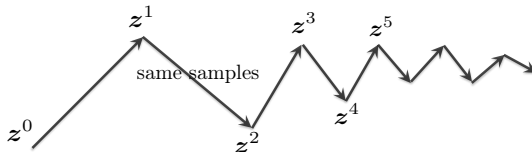


# No need of sample splitting

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



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



## **Low-rank matrix completion**