

# Solving Overparametrized Systems of Random Equations

Andrea Montanari

Stanford University

June 26, 2024

# The optimization puzzle in modern machine learning

- ▶ Empirical Risk Minimization (ERM) is highly non-convex
- ▶ Gradient methods find global optima

## Working hypothesis

*ERM becomes 'easy' if sufficiently overparametrized*

# The optimization puzzle in modern machine learning

- ▶ Empirical Risk Minimization (ERM) is highly non-convex
- ▶ Gradient methods find global optima

## Working hypothesis

*ERM becomes 'easy' if sufficiently overparametrized*

# The optimization puzzle in modern machine learning

## Working hypothesis

*ERM becomes 'easy' if sufficiently overparametrized*

Can we understand this in a simple model?

# Outline

- 1 A simple model with a long history
- 2 Gradient descent: Local analysis
- 3 Hessian descent
- 4 Exact solutions



Eliran Subag  
Weizmann Institute

## A small experiment with a small neural net

## An experiment: 2-Layer ELU network

$$f(\mathbf{z}; \mathbf{W}) = \frac{a}{\sqrt{m}} \sum_{j=1}^m s_j \sigma(\langle \mathbf{w}_j, \mathbf{z} \rangle), \quad \mathbf{z} \in \mathbb{R}^D.$$

$$\sigma(x) = \begin{cases} x & \text{if } x \geq 0, \\ e^x - 1 & \text{if } x < 0. \end{cases}, \quad \|\mathbf{W}\|_{\mathbb{F}}^2 = \sum_{i=1}^m \|\mathbf{w}_i\|_2^2 \leq m.$$



# Empirical Risk Minimization via SGD

$$R_n(\mathbf{W}) := \frac{1}{2n} \sum_{i=1}^n (y_i - f(\mathbf{z}_i; \mathbf{W}))^2,$$

$$\widetilde{\mathbf{W}}^{k+1} = \mathbf{W}^k - \frac{\eta_k}{2} \sum_{i \in B(k)} \nabla_{\mathbf{w}} (y_i - f(\mathbf{z}_i; \mathbf{W}^k))^2,$$

$$\mathbf{W}^{k+1} = \text{Proj}(\widetilde{\mathbf{W}}^{k+1}).$$

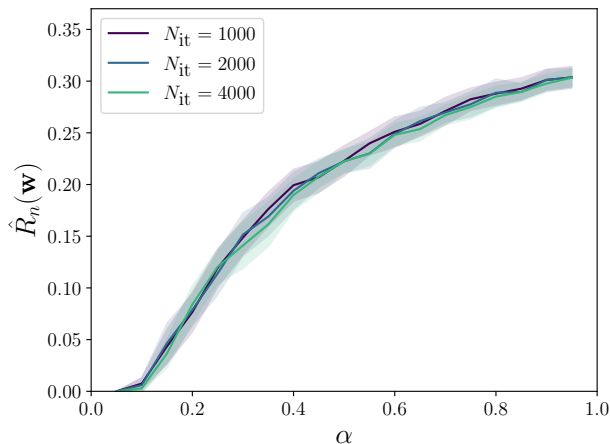
---

$$\eta_k = \text{lr}/(1 + \text{epoch}(k))^{1/2} \quad \mathbf{W}^0 \sim \mathcal{N}(0, \varepsilon^2 \mathbf{I}_{mD}/D), \quad \varepsilon = 0.03$$

## Data distribution

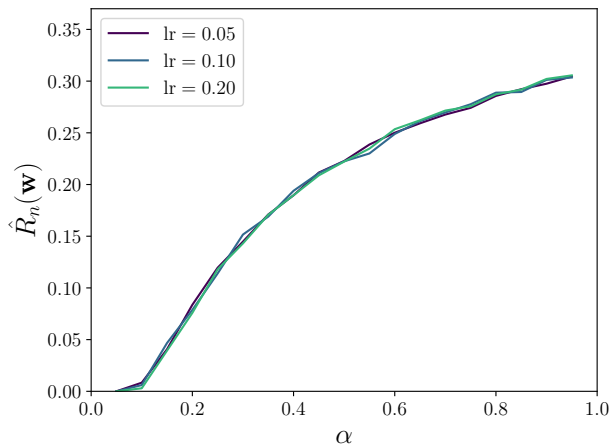
$$(\mathbf{z}_i, y_i) \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_D) \otimes \text{Unif}(\{+1, -1\})$$

## Varying number of epochs; $\alpha = n/mD$



- ▶  $m = D = 20$ ,  $lr = 0.1$
- ▶ Averages over 20 realizations; one std bands

## Varying learning rate; $\alpha = n/mD$



- ▶  $m = D = 20$ ,  $N_{it} = 2,000$
- ▶ Averages over 20 realizations

## A simple model with a long history

# General problem

- 1  $F_1, \dots, F_n : \mathbb{R}^d \rightarrow \mathbb{R}$  i.i.d. random functions.
- 2 Find  $\mathbf{x} \in \mathbb{S}^{d-1}$  such that  $F_i(\mathbf{x}) = 0$  for all  $i \leq n$ .

# Gaussian model

$F_i$  centered, rotation invariant (in distr) Gaussian processes

- ▶ Covariance

$$\mathbb{E}[F_i(\mathbf{x}^1)F_j(\mathbf{x}^2)] = \delta_{ij}\xi(\langle \mathbf{x}^1, \mathbf{x}^2 \rangle).$$

- ▶  $\mathbf{F}(\mathbf{x}) = (F_1(\mathbf{x}), \dots, F_n(\mathbf{x}))^\top$ .

## Cost function and approximate solutions

$$R_n(\mathbf{x}) := \frac{1}{2} \|\mathbf{F}(\mathbf{x})\|_2^2$$

► At a fixed  $\mathbf{x}_0 \in \mathbb{S}^{d-1}$ ,  $R_n(\mathbf{x}_0) = n\xi(1)/2 + o(n)$

► Solutions

$$\text{Sol}_{n,d}(\varepsilon) := \left\{ \mathbf{x} \in \mathbb{S}^{d-1} : \|\mathbf{F}(\mathbf{x})\|_2^2 \leq n\xi(1) \cdot \varepsilon \right\}.$$

►  $n, d \rightarrow \infty$ ,  $n/d \rightarrow \alpha$ .



## Questions

Q1 Do exact solutions exist:  $\text{Sol}_{n,d}(0) \neq \emptyset$ ?

Do approximate solution exist,  $\text{Sol}_{n,d}(\varepsilon) \neq \emptyset$ ?

Q2 Can we find them in polytime?

# Questions

Q1 Do exact solutions exist:  $\text{Sol}_{n,d}(0) \neq \emptyset$ ?

Do approximate solution exist,  $\text{Sol}_{n,d}(\varepsilon) \neq \emptyset$ ?

Q2 Can we find them in polytime?

## History: Classical (complex) setting

$$\mathbf{F} : \mathbb{C}^d \rightarrow \mathbb{C}^n, \text{ homogeneous, } \deg(F_i) = p_i$$

Q1 Bezout's theorem (1779)

For  $n = d - 1$ , deterministically:

$$|\text{Sol}_{n,d}(0)| = \prod_{i=1}^n p_i$$

Q2

- ▶ Smale 17th problem (1993-1998)
- ▶ Positive answer (Lairez, 2020)
- ▶ Homotopy methods

## History: Classical (complex) setting

$$\mathbf{F} : \mathbb{C}^d \rightarrow \mathbb{C}^n, \text{ homogeneous, } \deg(F_i) = p_i$$

**Q1** Bezout's theorem (1779)

For  $n = d - 1$ , deterministically:

$$|\text{Sol}_{n,d}(0)| = \prod_{i=1}^n p_i$$

**Q2**

- ▶ Smale 17th problem (1993-1998)
- ▶ Positive answer (Lairez, 2020)
- ▶ Homotopy methods

## History: Classical (complex) setting

$$\mathbf{F} : \mathbb{C}^d \rightarrow \mathbb{C}^n, \text{ homogeneous, } \deg(F_i) = p_i$$

**Q1** Bezout's theorem (1779)

For  $n = d - 1$ , deterministically:

$$|\text{Sol}_{n,d}(0)| = \prod_{i=1}^n p_i$$

**Q2**

- ▶ **Smale 17th problem** (1993-1998)
- ▶ Positive answer (Lairez, 2020)
- ▶ Homotopy methods

# History: Real setting

Q1 Homogeneous case,  $n = d - 1$ .

- ▶ Subag 2022: With high probability

$$|\text{Sol}_{n,d}(0)| = (1 + o(1)) \prod_{i=1}^n \sqrt{p_i}$$

Q1 Non-homogeneous case,  $n/d \rightarrow \alpha \in (0, \infty)$ .

- ▶ Next slide

Q2 The rest of this talk

## History: Real setting

Q1 Homogeneous case,  $n = d - 1$ .

- ▶ Subag 2022: With high probability

$$|\text{Sol}_{n,d}(0)| = (1 + o(1)) \prod_{i=1}^n \sqrt{p_i}$$

Q1 Non-homogeneous case,  $n/d \rightarrow \alpha \in (0, \infty)$ .

- ▶ Next slide

Q2 The rest of this talk

## History: Real setting

**Q1** Homogeneous case,  $n = d - 1$ .

- ▶ Subag 2022: With high probability

$$|\text{Sol}_{n,d}(0)| = (1 + o(1)) \prod_{i=1}^n \sqrt{p_i}$$

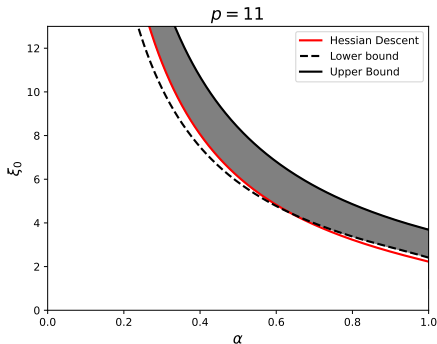
**Q1** Non-homogeneous case,  $n/d \rightarrow \alpha \in (0, \infty)$ .

- ▶ Next slide

**Q2** The rest of this talk



Q1:  $\xi(\mathbf{q}) = \xi_0 + \mathbf{q}^{11}$  (polynomial of degree 11)



- ▶ Above gray region,  $\alpha > \alpha_{\text{UB}}(\xi_0)$ :  $\text{Sol}_{n,d}(\varepsilon) = \emptyset$
- ▶ Below gray region,  $\alpha < \alpha_{\text{LB}}(\xi_0)$ :  $\text{Sol}_{n,d}(0) = \emptyset$

---

See paper for formal statements.

## Gradient descent: Local analysis

# Projected Gradient Descent

## Gradient descent

$$\begin{aligned} \mathbf{z}^{k+1} &= \mathbf{x}^k - \eta \mathbf{P}_{\mathcal{T}, \mathbf{x}^k} \nabla R_n(\mathbf{x}^k), \\ \mathbf{x}^{k+1} &= \frac{\mathbf{z}^{k+1}}{\|\mathbf{z}^{k+1}\|_2}. \end{aligned}$$

## Projected gradient flow

$$\dot{\mathbf{x}}(t) = -\mathbf{P}_{\mathcal{T}, \mathbf{x}(t)} \nabla R_n(\mathbf{x}(t)).$$

Difficult to analyze sharply!

# Projected Gradient Descent

## Gradient descent

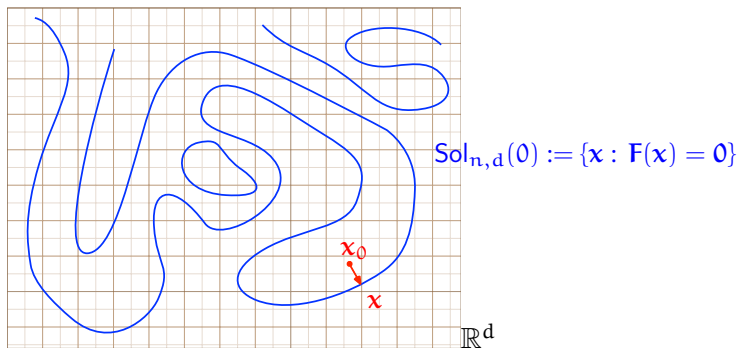
$$\begin{aligned} \mathbf{z}^{k+1} &= \mathbf{x}^k - \eta \mathbf{P}_{\mathcal{T}, \mathbf{x}^k} \nabla R_n(\mathbf{x}^k), \\ \mathbf{x}^{k+1} &= \frac{\mathbf{z}^{k+1}}{\|\mathbf{z}^{k+1}\|_2}. \end{aligned}$$

## Projected gradient flow

$$\dot{\mathbf{x}}(t) = -\mathbf{P}_{\mathcal{T}, \mathbf{x}(t)} \nabla R_n(\mathbf{x}(t)).$$

Difficult to analyze sharply!

# Local analysis: Taylor expand around initialization



---

**State of the art in ML Theory:** Jacot, Gabriel, Hongler, 2018; Du, Zhai, Póczos, Singh 2018; Allen-Zhu, Li, Song 2018; Chizat, Bach, 2019; Arora, Du, Hu, Li, Salakhutdinov, Wang, 2019; Oymak, Soltanolkotabi, 2019; ...

## Local analysis

$$\underline{\alpha}_{\text{GD}}(\xi) := \frac{c_0 \xi'(1)^2}{\xi''(1)\xi(1)(\log(\xi'''(1)/\xi''(1)) \vee 1)}.$$

Theorem (M, Subag, 2023)

If  $\alpha < \underline{\alpha}_{\text{GD}}(\xi)$ , and  $\eta < 1/(C_1 d)$ , then whp for all  $k \geq 1$ ,

$$\|\mathbf{F}(\mathbf{x}^k)\|_2^2 \leq 2n\xi(1) e^{-c_2(\sqrt{d}-\sqrt{n})^2(\eta k)}.$$

Special case:  $\xi(q) = \xi_0 + q^p$

$$\alpha_{\text{GD}}(\xi_0, p) \asymp \frac{1}{\xi_0 \log p}.$$

To be compared with

$$\alpha_{\text{LB}}(\xi_0, p) = \frac{\log p}{\xi_0} \cdot (1 + o_p(1)).$$

## Hessian descent



## Problem with GD

- ▶ Need  $\|\nabla R_n(\mathbf{x}^k)\|_2 \geq \varepsilon$ .
- ▶ Cannot be true uniformly.
- ▶  $\Rightarrow$  Local analysis :(

## Problem with GD

- ▶ Need  $\|\nabla R_n(\mathbf{x}^k)\|_2 \geq \varepsilon$ .
- ▶ Cannot be true uniformly.
- ▶  $\Rightarrow$  Local analysis :(

## Problem with GD

- ▶ Need  $\|\nabla R_n(\mathbf{x}^k)\|_2 \geq \varepsilon$ .
- ▶ Cannot be true uniformly.
- ▶  $\Rightarrow$  Local analysis :(

## Problem with GD

- ▶ Need  $\|\nabla R_n(\mathbf{x}^k)\|_2 \geq \varepsilon$ .
- ▶ Cannot be true uniformly.
- ▶  $\Rightarrow$  Local analysis :(

*Idea:* Use Hessian

## Problem with Hessian Descent

$$\mathbf{D}^2\mathcal{R}_n(\mathbf{x}) = \nabla^2\mathcal{R}_n(\mathbf{x})|_{\mathbb{T}(\mathbf{x})} - \langle \mathbf{x}, \nabla\mathcal{R}_n(\mathbf{x}) \rangle \mathbf{I}$$

- ▶  $\mathbf{D}^2 =$  Riemannian Hessian
- ▶  $\nabla^2 =$  Euclidean Hessian
- ▶  $\mathbb{T}(\mathbf{x}) =$  Tangent space at  $\mathbf{x}$
- ▶ Problem:  $\langle \mathbf{x}, \nabla\mathcal{R}_n(\mathbf{x}) \rangle$  does not concentrate

## Problem with Hessian Descent

$$\mathbf{D}^2\mathcal{R}_n(\mathbf{x}) = \nabla^2\mathcal{R}_n(\mathbf{x})|_{\mathbb{T}(\mathbf{x})} - \langle \mathbf{x}, \nabla\mathcal{R}_n(\mathbf{x}) \rangle \mathbf{I}$$

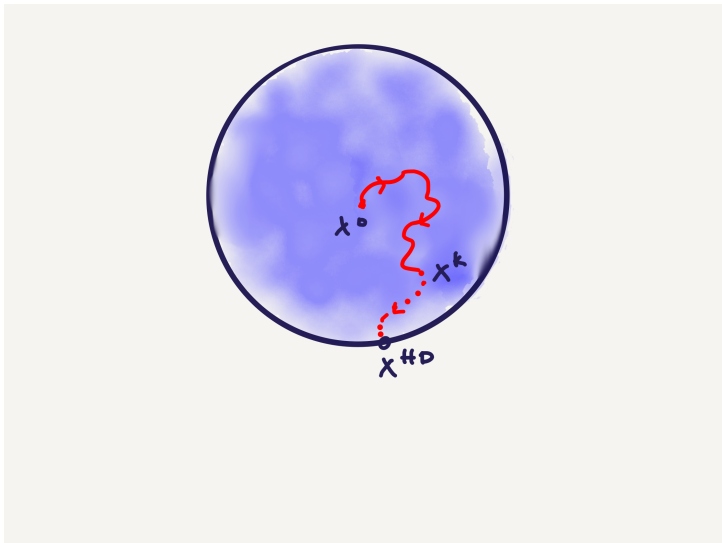
- ▶  $\mathbf{D}^2 =$  Riemannian Hessian
- ▶  $\nabla^2 =$  Euclidean Hessian
- ▶  $\mathbb{T}(\mathbf{x}) =$  Tangent space at  $\mathbf{x}$
- ▶ Problem:  $\langle \mathbf{x}, \nabla\mathcal{R}_n(\mathbf{x}) \rangle$  does not concentrate

## Problem with Hessian Descent

$$\mathbf{D}^2\mathcal{R}_n(\mathbf{x}) = \nabla^2\mathcal{R}_n(\mathbf{x})|_{\mathcal{T}(\mathbf{x})} - \langle \mathbf{x}, \nabla\mathcal{R}_n(\mathbf{x}) \rangle \mathbf{I}$$

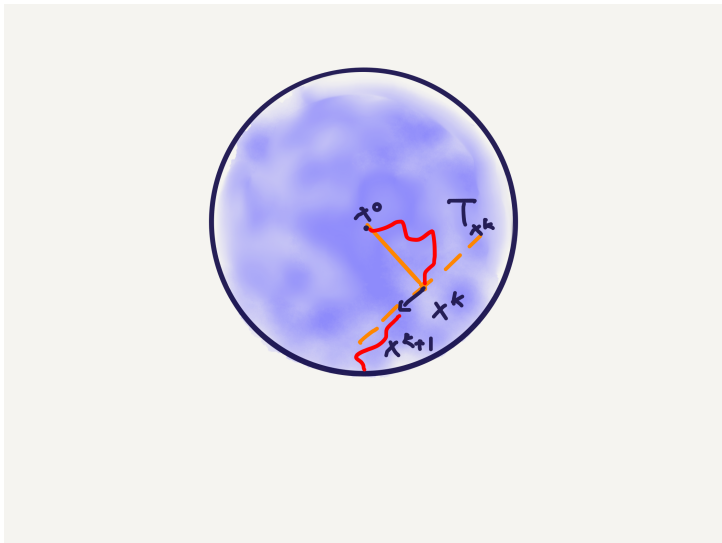
*Idea:* Relax sphere constraint

## Algorithm: Sketch





## Algorithm: Orthogonal steps



## Algorithm: Simplified

$\mathbf{v}_{\min}(\mathbf{A})$  := eigenvector associated to smallest eigenvalue of  $\mathbf{A}$

---

---

```
Initialize  $\mathbf{x}^1 \sim \sqrt{\delta} \cdot \text{Unif}(\mathbb{S}^{d-1})$ ;  
for  $k \in \{1, \dots, K := 1/\delta - 1\}$  do  
    Compute  $\mathbf{v}(\mathbf{x}^k) = \mathbf{v}_{\min}(\nabla^2 R_n(\mathbf{x}^k)|_{T_{\mathbf{x}^k}})$ ;  
     $s_k := \text{sign}(\langle \mathbf{v}(\mathbf{x}^k), \nabla R_n(\mathbf{x}^k) \rangle)$ ;  
     $\mathbf{x}^{k+1} = \mathbf{x}^k - s_k \sqrt{\delta} \mathbf{v}(\mathbf{x}^k)$ ;  
end  
return  $\mathbf{x}^{\text{HD}} = \mathbf{x}^K$ ;
```

---

## Full algorithm

---

---

Initialize  $\mathbf{x}^1 \sim \sqrt{\delta} \cdot \text{Unif}(\mathbb{S}^{d-1})$ ;

**for**  $k \in \{1, \dots, K := 1/\delta - 1\}$  **do**

    Compute  $\mathbf{v} = \mathbf{v}(\mathbf{x}^k) \in \mathbb{T}_{\mathbf{x}^k}$  such that  $\|\mathbf{v}\|_2 = 1$  and

$$\langle \mathbf{v}, \nabla^2 \mathcal{R}_n(\mathbf{x}^k) \mathbf{v} \rangle \leq \lambda_{\min}(\nabla^2 \mathcal{R}_n(\mathbf{x}^k)|_{\mathbb{T}, \mathbf{x}^k}) + d\delta;$$

$s_k := \text{sign}(\langle \mathbf{v}(\mathbf{x}^k), \nabla \mathcal{R}_n(\mathbf{x}^k) \rangle)$ ;

$\mathbf{x}^{k+1} = \mathbf{x}^k - s_k \sqrt{\delta} \mathbf{v}(\mathbf{x}^k)$ ;

**end**

**return**  $\mathbf{x}^{\text{HD}} = \mathbf{x}^K$ ;

---

# Analysis

## Theorem (M, Subag, 2023)

For  $\alpha \in (0, 1)$ ,  $\mathbf{a}, \mathbf{b} \in \mathbb{R}_{\geq 0}$ , define

$$z_*(\alpha, \mathbf{a}, \mathbf{b}) := \inf_{m>0} \left\{ \frac{1}{m} - \frac{\alpha \mathbf{b}}{1 + \mathbf{b}m} + \alpha^2 m \right\}.$$

Let  $u(\cdot; \alpha, \xi) : [0, 1] \rightarrow \mathbb{R}$  be the unique solution of the ODE

$$\frac{du}{dt}(t) = -\frac{1}{2\alpha} z_*(\alpha; \sqrt{2\alpha u(t)\xi''(t)}, \xi'(t)), \quad u(0) = \frac{1}{2}\xi(0).$$

Then whp

$$\frac{1}{n} R_n(\mathbf{x}^{\text{HD}}) \leq u(1; \alpha, \xi) + C_0 \delta.$$

---

Recall  $R_n(\mathbf{x}) := \|\mathbf{F}(\mathbf{x}^{\text{HD}})\|_2^2/2$ .

# Where does this come from?

## Hessian

$$\nabla^2 R_n(\mathbf{x}) = \sum_{\ell=1}^n F_{\ell}(\mathbf{x}) \nabla^2 F_{\ell}(\mathbf{x}) + \mathbf{D}F(\mathbf{x})^{\top} \mathbf{D}F(\mathbf{x})$$

Distribution as  $\mathbf{x}$ ,  $\|\mathbf{x}\|^2 = q$

$$\nabla^2 R_n(\mathbf{x})|_{\Gamma(\mathbf{x})} = \sqrt{R_n(\mathbf{x}) \xi''(q)} \mathbf{W} + \xi'(q) \mathbf{Z}^{\top} \mathbf{Z},$$

$$(\mathbf{W}, \mathbf{Z}) \sim \text{GOE}(d-1) \otimes \text{GOE}(n, d-1)$$

Decrease in value

$$\lim_{n, d \rightarrow \infty} \frac{1}{d} \lambda_{\min}(\mathbf{A}_{n, d}) = -z_*(\alpha, a, b) := - \inf_{m > 0} \frac{1}{m} - \frac{\alpha b}{1 + bm} + a^2 m.$$

Where does this come from?

**Hessian**

$$\nabla^2 R_n(\mathbf{x}) = \sum_{\ell=1}^n F_{\ell}(\mathbf{x}) \nabla^2 F_{\ell}(\mathbf{x}) + \mathbf{D}\mathbf{F}(\mathbf{x})^T \mathbf{D}\mathbf{F}(\mathbf{x})$$

Distribution as  $\mathbf{x}$ ,  $\|\mathbf{x}\|^2 = q$

$$\nabla^2 R_n(\mathbf{x})|_{\Gamma(\mathbf{x})} = \sqrt{R_n(\mathbf{x}) \xi''(q)} \mathbf{W} + \xi'(q) \mathbf{Z}^T \mathbf{Z},$$

$$(\mathbf{W}, \mathbf{Z}) \sim \text{GOE}(d-1) \otimes \text{GOE}(n, d-1)$$

Decrease in value

$$\lim_{n, d \rightarrow \infty} \frac{1}{d} \lambda_{\min}(\mathbf{A}_{n,d}) = -z_*(\alpha, a, b) := -\inf_{m>0} \frac{1}{m} - \frac{\alpha b}{1 + bm} + a^2 m.$$

Where does this come from?

**Hessian**

$$\nabla^2 R_n(\mathbf{x}) = \sum_{\ell=1}^n F_{\ell}(\mathbf{x}) \nabla^2 F_{\ell}(\mathbf{x}) + \mathbf{D}F(\mathbf{x})^T \mathbf{D}F(\mathbf{x})$$

**Distribution as  $\mathbf{x}$ ,  $\|\mathbf{x}\|^2 = q$**

$$\nabla^2 R_n(\mathbf{x})|_{\Gamma(\mathbf{x})} = \sqrt{R_n(\mathbf{x}) \xi''(q)} \mathbf{W} + \xi'(q) \mathbf{Z}^T \mathbf{Z},$$

$$(\mathbf{W}, \mathbf{Z}) \sim \text{GOE}(d-1) \otimes \text{GOE}(n, d-1)$$

Decrease in value

$$\lim_{n, d \rightarrow \infty} \frac{1}{d} \lambda_{\min}(\mathbf{A}_{n,d}) = -z_*(\alpha, a, b) := -\inf_{m>0} \frac{1}{m} - \frac{\alpha b}{1 + bm} + a^2 m.$$

Where does this come from?

**Hessian**

$$\nabla^2 R_n(\mathbf{x}) = \sum_{\ell=1}^n F_{\ell}(\mathbf{x}) \nabla^2 F_{\ell}(\mathbf{x}) + \mathbf{D}\mathbf{F}(\mathbf{x})^{\top} \mathbf{D}\mathbf{F}(\mathbf{x})$$

**Distribution as  $\mathbf{x}$ ,  $\|\mathbf{x}\|^2 = q$**

$$\nabla^2 R_n(\mathbf{x})|_{\Gamma(\mathbf{x})} = \sqrt{R_n(\mathbf{x}) \xi''(q)} \mathbf{W} + \xi'(q) \mathbf{Z}^{\top} \mathbf{Z},$$

$$(\mathbf{W}, \mathbf{Z}) \sim \text{GOE}(d-1) \otimes \text{GOE}(n, d-1)$$

**Decrease in value**

$$\lim_{n, d \rightarrow \infty} \frac{1}{d} \lambda_{\min}(\mathbf{A}_{n, d}) = -z_*(\alpha, a, b) := -\inf_{m>0} \frac{1}{m} - \frac{\alpha b}{1 + bm} + a^2 m.$$



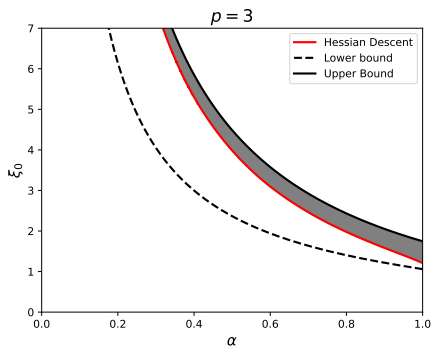
Special case:  $\xi(q) = \xi_0 + q^p$

$$\frac{4(p-1)}{p\xi_0 + 4(p-1)} \leq \alpha_{\text{HD}}(\xi_0, p) \leq \frac{4(p-1)}{p\xi_0}.$$

To be compared with

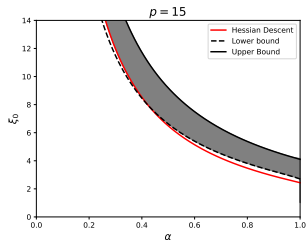
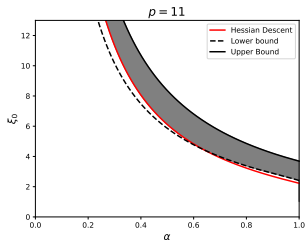
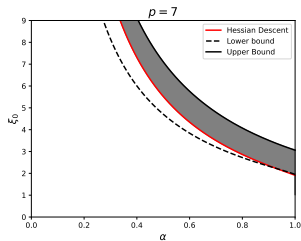
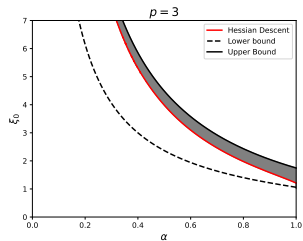
$$\alpha_{\text{GD}}(\xi_0, p) \asymp \frac{1}{\xi_0 \log p},$$
$$\alpha_{\text{LB}}(\xi_0, p) = \frac{\log p}{\xi_0} \cdot (1 + o_p(1)).$$

# Phase diagram ( $\xi(q) = \xi_0 + q^3$ )



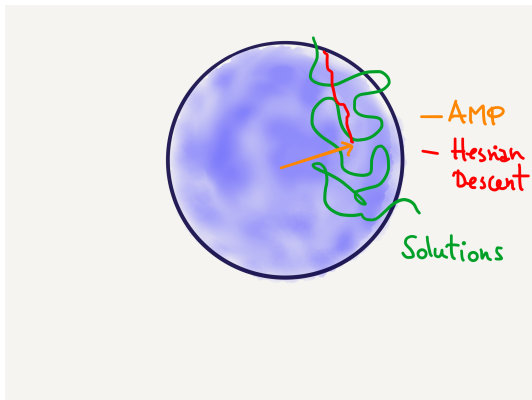
- ▶ Above gray region,  $\alpha > \alpha_{\text{UB}}(\xi_0)$ :  $\text{Sol}_{n,d}(\varepsilon) = \emptyset$
- ▶ Below gray region,  $\alpha < \alpha_{\text{LB}}(\xi_0)$ :  $\text{Sol}_{n,d}(0) = \emptyset$
- ▶ Red line:  $\alpha_{\text{HD}}(\xi_0, p)$

# Phase diagram ( $\xi(q) = \xi_0 + q^p$ )



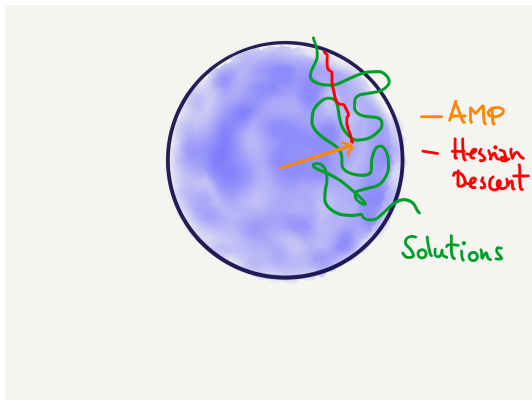
## Is HD optimal (among polytime algs)?

- ▶ **No!** Suboptimal when  $\mathbf{F}(\mathbf{x})$  has degree-1 term
- ▶  $\text{Sol}_{n,d}(0)$  not centered at  $\mathbf{0}$
- ▶ See paper for the fix/general algorithm
- ▶ Conjectured to be optimal among 'stable algorithms'



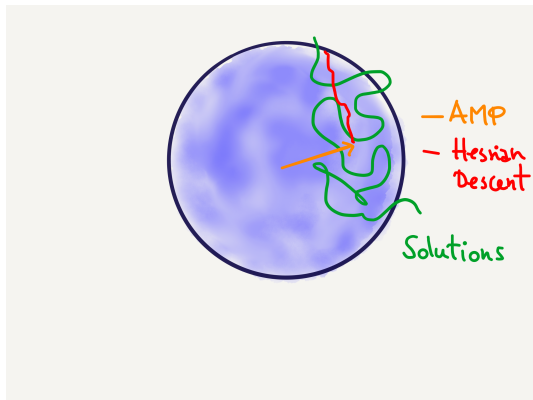
## Is HD optimal (among polytime algs)?

- ▶ **No!** Suboptimal when  $\mathbf{F}(\mathbf{x})$  has degree-1 term
- ▶  $\text{Sol}_{n,d}(0)$  not centered at  $\mathbf{0}$
- ▶ See paper for the fix/general algorithm
- ▶ Conjectured to be optimal among 'stable algorithms'



## Is HD optimal (among polytime algs)?

- ▶ **No!** Suboptimal when  $\mathbf{F}(\mathbf{x})$  has degree-1 term
- ▶  $\text{Sol}_{n,d}(0)$  not centered at  $\mathbf{0}$
- ▶ See paper for the fix/general algorithm
- ▶ Conjectured to be optimal among 'stable algorithms'



What about exact solutions?

## What is an exact solution?

Definition (Shub, Smale, 1993)

$\mathbf{x}_*$  is an *approximate solution* of  $\mathbf{F}(\mathbf{x}) = \mathbf{0}$  if letting  $(\mathbf{x}^k)_{k \geq 0}$  be Newton iterates with  $\mathbf{x}^0 = \mathbf{x}_*$ , then, for all  $k$

$$\|\mathbf{F}(\mathbf{x}^k)\| \leq \|\mathbf{F}(\mathbf{x}^0)\| \cdot \exp \{ -c \cdot 2^k \}.$$

Smale 17th problem over the reals:

*Can we find approximate solutions in polytime?*



## What is an exact solution?

Definition (Shub, Smale, 1993)

$\mathbf{x}_*$  is an *approximate solution* of  $\mathbf{F}(\mathbf{x}) = \mathbf{0}$  if letting  $(\mathbf{x}^k)_{k \geq 0}$  be Newton iterates with  $\mathbf{x}^0 = \mathbf{x}_*$ , then, for all  $k$

$$\|\mathbf{F}(\mathbf{x}^k)\| \leq \|\mathbf{F}(\mathbf{x}^0)\| \cdot \exp \{ -c \cdot 2^k \}.$$

**Smale 17th problem over the reals:**

*Can we find approximate solutions in polytime?*

# What is an exact solution?

## Theorem (M, Subag, 2024)

*Assume  $F_i$  homogeneous, arbitrary (possibly different) degrees. Then there exists a deterministic polytime algorithm such that, if*

$$n \leq d - C\sqrt{d \log d},$$

*then it return a an approximate solution, with high probability wrt  $\mathbf{F}$ .*

## Conclusion #1

- ▶ Random systems of nonlinear equations
- ▶ Rich computational/probabilistic structure
- ▶ Quantitative comparison with neural nets lanscape?

## Conclusion #1

- ▶ Random systems of nonlinear equations
- ▶ Rich computational/probabilistic structure
- ▶ Quantitative comparison with neural nets lanscape?

## Conclusion #2



It is an honor to celebrate Andrew!  
Thanks!

## Epilogue: Revisiting the original experiment

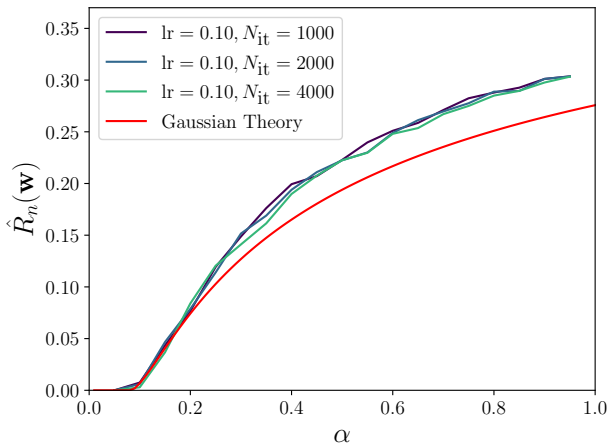
# Empirical Risk Minimization

$$f(\mathbf{z}; \mathbf{W}) = \frac{\alpha}{\sqrt{m}} \sum_{j=1}^m s_j \sigma(\langle \mathbf{w}_j, \mathbf{z} \rangle), \quad \mathbf{z} \in \mathbb{R}^D.$$

$$R_n(\mathbf{W}) := \frac{1}{2n} \sum_{i=1}^n (y_i - f(\mathbf{z}_i; \mathbf{W}))^2,$$

$$\|\mathbf{W}\|_F^2 \leq m.$$

# Experiments vs Gaussian Theory: $\alpha = 1$

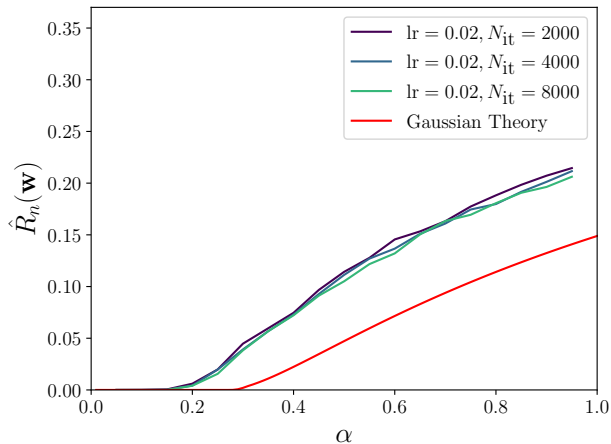


---

Red: Approx matching covariance



## Experiments vs Gaussian Theory: $\alpha = 2$



# Experiments vs Gaussian Theory: $\alpha = 5$

