# From fitting ellipsoids to random points, to learning in large neural networks

*Antoine Maillard*
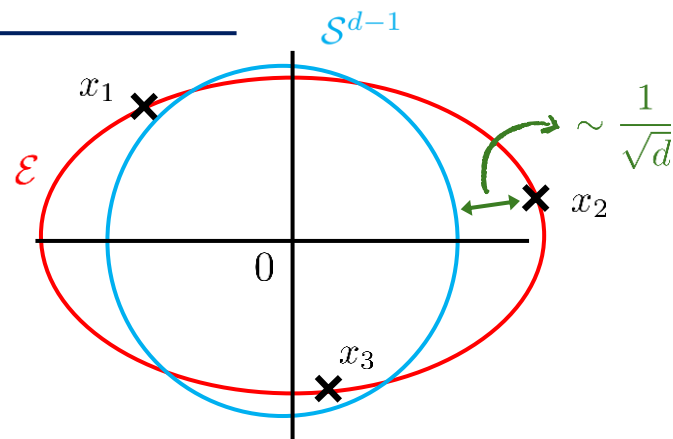
**ETH** *zürich*

Roccella – September 4th 2024

1

Part I: Fitting ellipsoids to random points

# Fitting ellipsoids to random points

$$x_1, \cdots, x_n \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathrm{I}_d/d)$$

$$n, d \to \infty$$

*Does $\mathcal{E}$ exist ?*



$\mathcal{S}^{d-1}$

$x_1$

$\sim \dfrac{1}{\sqrt{d}}$

$\mathcal{E}$

$x_2$

$0$

$x_3$

## Ellipsoid Fitting Property

$$\mathbb{P}[\exists S \in \mathbb{R}^{d \times d} : S \succeq 0 \text{ and } x_i^\top S x_i = 1 \text{ for all } i \in [n]]$$

Principal axes of $\mathcal{E}$ $\Longleftrightarrow$ Eigenspaces of $S$

$$r_i(\mathcal{E}) = \lambda_i(S)^{-1/2}$$

## Some motivations

❖ Low-rank matrix decomposition
  Saunderson & al '12 ; '13 ; '13

❖ Independent Components Analysis
  Podosinnikova & al '19

❖ Discrepancy of random matrices
  Potechin & al '22

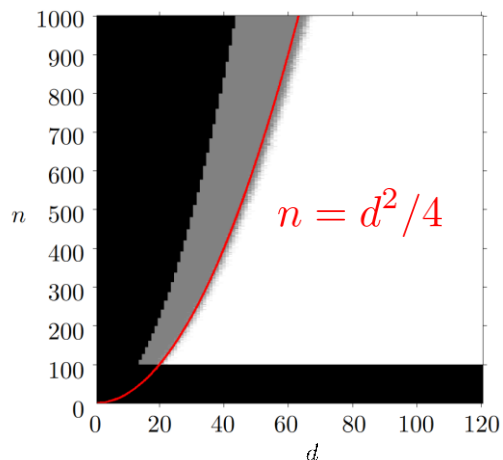❖ Neural networks with quadratic activations
  **More on that later !**

# The ellipsoid fitting conjecture

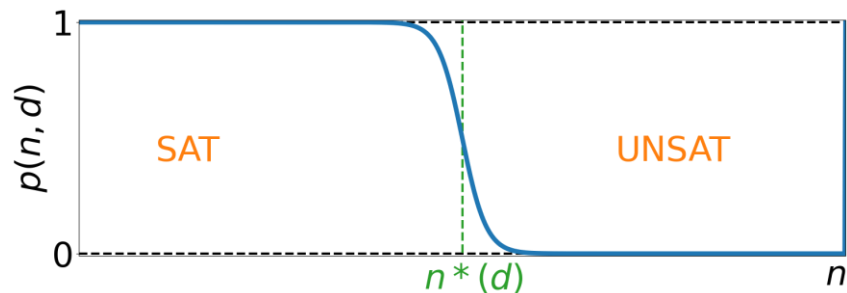Ellipsoid fitting is a **semidefinite program**

⬇

Convex problem + efficient solvers

$p(n, d) = \mathbb{P}[\text{An ellipsoid fit exists}]$

■ : No simulation    ■ : No solutions    □ : Solutions exist



$n = d^2/4$

Saunderson, James, et al. *SIAM Journal on Matrix Analysis and Applications* 2012
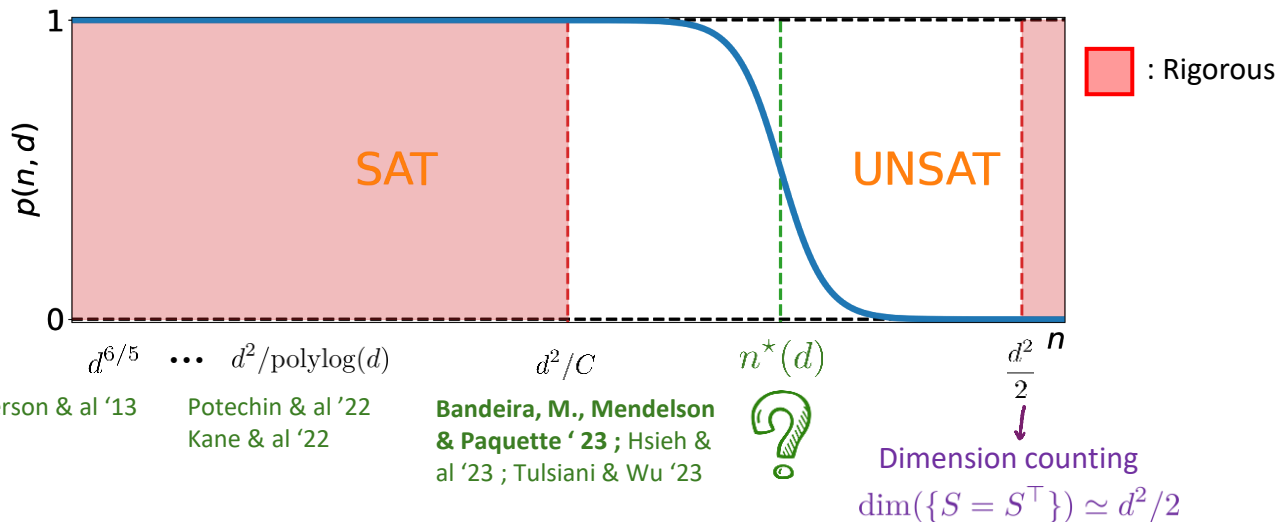


SAT          UNSAT

$n*(d)$

$n$

**Open conjecture**

$$\lim_{d \to \infty} \frac{n^\star(d)}{d^2} = \frac{1}{4}$$

# The ellipsoid fitting conjecture: what is known

**Conjecture** $\boxed{\lim_{d\to\infty} \frac{n^\star(d)}{d^2} = \frac{1}{4}}$



☐ : Rigorous

Progress on **lower bounds**

$d^{6/5}$ ••• $d^2/\text{polylog}(d)$

Saunderson & al '13

Potechin & al '22
Kane & al '22

$d^2/C$

**Bandeira, M., Mendelson & Paquette ' 23 ;** Hsieh & al '23 ; Tulsiani & Wu '23

$n^\star(d)$

$\frac{d^2}{2}$

Dimension counting
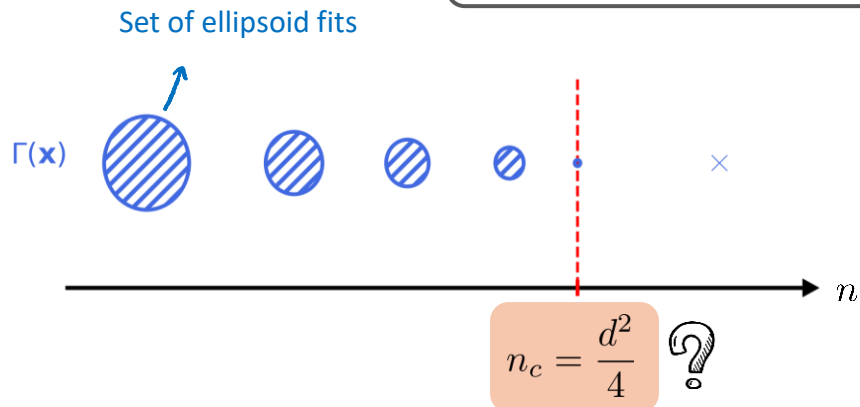$\dim(\{S = S^\top\}) \simeq d^2/2$

**This talk**

We see EFP as a **Random Constraint Satisfaction Problem**

$$\begin{cases} S \succeq 0 \qquad \longrightarrow \text{"spectral" constraint} \\ \{x_i^\top S x_i = 1\}_{i=1}^n \end{cases}$$

**"disordered" model**

# Statistical physics of ellipsoid fitting

**Ellipsoid Fitting Property**

$$\mathbb{P}[\exists S \in \mathbb{R}^{d \times d} : S \succeq 0 \text{ and } x_i^\top S x_i = 1 \text{ for all } i \in [n]]$$

Set of ellipsoid fits

$\Gamma(\mathbf{x})$

$n$

$$n_c = \frac{d^2}{4}$$

Volume of solutions / **"Partition function"**

$$\text{supp}(P_0) \subseteq \mathcal{S}_d^+$$

$$\mathcal{Z} := \int P_0(\mathrm{d}S) \prod_{i=1}^n \delta(x_i^\top S x_i - 1)$$

**Replica calculation** of $\mathbb{E} \log \mathcal{Z}$

- **Analytical derivation** of the threshold $n_c = d^2/4$
- **Shape** (spectrum) of typical ellipsoid fits
- Extensions to **non–Gaussian** vectors

[**M.** & Kunisky '23]

Connections to "HCIZ" integrals in random matrix theory [Matytsin '94 ; Guionnet&al'02]

# Mathematical physics for ellipsoid fitting  [M. & Bandeira '23]

**Two-steps proof**

**I:** • "**Gaussian universality**" lemma : $\dfrac{1}{n}\log\mathcal{Z} \simeq \dfrac{1}{n}\log\mathcal{Z}_G$

[Goldt & al '22, Montanari & Saeed '22, Hu & Lu '22, …]

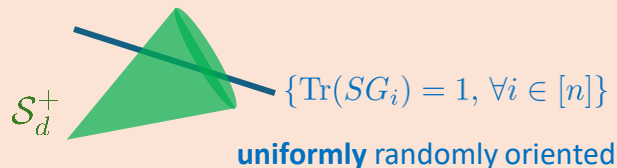$x_i^\top S x_i \longrightarrow \mathrm{Tr}(SG_i)$ ← **Gaussian matrix**

**CRITICAL**

$$\mathcal{Z} := \int P_0(\mathrm{d}S)\prod_{i=1}^{n}\delta(x_i^\top S x_i - 1) \longrightarrow \mathcal{Z}_G := \int P_0(\mathrm{d}S)\prod_{i=1}^{n}\delta(\mathrm{Tr}(SG_i) - 1)$$

**II:** • **Random convex geometry** tools for $\mathcal{Z}_G$

Extensions of Gordon's **min-max theorem**
[Gordon '88, Amelunxen & al '14]

$\mathcal{S}_d^+$

$\{\mathrm{Tr}(SG_i) = 1, \forall i \in [n]\}$

**uniformly** randomly oriented

**Theorem:** The problem associated to $\mathcal{Z}_G$ is

- SAT (whp) if $n \le (1-\varepsilon)\omega(\mathcal{S}_d^+)^2$

- UNSAT (whp) if $n \ge (1+\varepsilon)\omega(\mathcal{S}_d^+)^2$ ←

**Gaussian width**

$$\omega(\mathcal{S}_d^+) := \mathbb{E}\max_{\substack{S\succeq 0\\ \|S\|_F=1}}\mathrm{Tr}[GS]$$

$$\omega(\mathcal{S}_d^+) \sim_{d\to\infty} \frac{d}{2} \implies \boxed{n^\star(\mathcal{Z}_G) \sim \frac{d^2}{4}}$$

# Mathematical physics for ellipsoid fitting  [**M.** & Bandeira '23]

**I:** "**Gaussian universality**" lemma  ➕  **II:** **Random convex geometry** tools

## Theorem

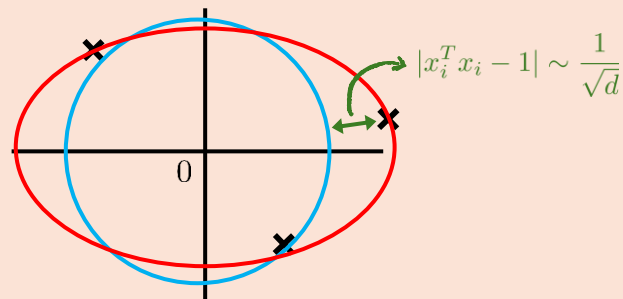$$\mathbf{EFP}_{\varepsilon, M} : \exists S \in \mathbb{R}^{d \times d} : S \succeq 0, \|S\| \leq M \text{ and } \frac{1}{n}\sum_{i=1}^{n} |x_i^T S x_i| \leq \frac{\varepsilon}{\sqrt{d}}$$

$$\mathbf{EFP} = \mathbf{EFP}_{0, \infty}$$

$$n/d^2 \to \alpha \begin{cases} \alpha < {}^1\!/_4 \quad \exists M_\alpha : \forall \varepsilon > 0, \ \mathbb{P}[\mathbf{EFP}_{\varepsilon, M_\alpha}] \to_{d \to \infty} 1 \\ \\ \alpha > {}^1\!/_4 \quad \exists \varepsilon_\alpha : \forall M > 0, \ \mathbb{P}[\mathbf{EFP}_{\varepsilon_\alpha, M}] \to_{d \to \infty} 0 \end{cases}$$



$$|x_i^T x_i - 1| \sim \frac{1}{\sqrt{d}}$$

➢ Strengthen proof to **exact** ellipsoid fitting ?

➢ Extension to **other high-dimensional semidefinite programs** ?

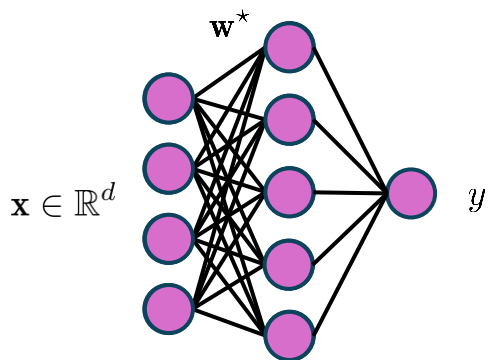➢ What does it have to do with **learning in neural networks** ??

Part II : Learning in neural networks

# A two-layer neural network with quadratic activation [**M.**, Troiani, Martin, Krzakala, Zdeborová '24]

## Teacher network

$\mathbf{w}^\star$

$\mathbf{x} \in \mathbb{R}^d$

$y$

$$y_i = f_{\mathbf{W}^*}(\mathbf{x}_i) := \frac{1}{m} \sum_{k=1}^{m} \left[ \frac{1}{\sqrt{d}} (\mathbf{w}_k^*)^T \cdot \mathbf{x}_i \right]^2$$

$\sim \mathcal{N}(0, \mathrm{I}_d)$

$\mathbf{w}_k^\star \sim \mathcal{N}(0, \mathrm{I}_d)$

**High-dimensional limit**

$$d \to \infty \,;\, m = \Theta(d)$$

**Learning from data**

$$\{(\mathbf{x}_1, y_1), \cdots, (\mathbf{x}_n, y_n)\} \implies \boxed{\mathbf{W}^\star}$$

**Bayes-optimal generalization error**

$$\mathcal{E}_{\text{gen.}} := \mathbb{E}_{\mathbf{W}^\star, \{\mathbf{x}_i\}} \min_{\hat{y}(\{y_i, \mathbf{x}_i\})} \mathbb{E}_{\mathbf{x}_{\text{test}}} [(\hat{y}(\mathbf{x}_{\text{test}}) - f_{\mathbf{W}^\star}(\mathbf{x}_{\text{test}}))^2]$$

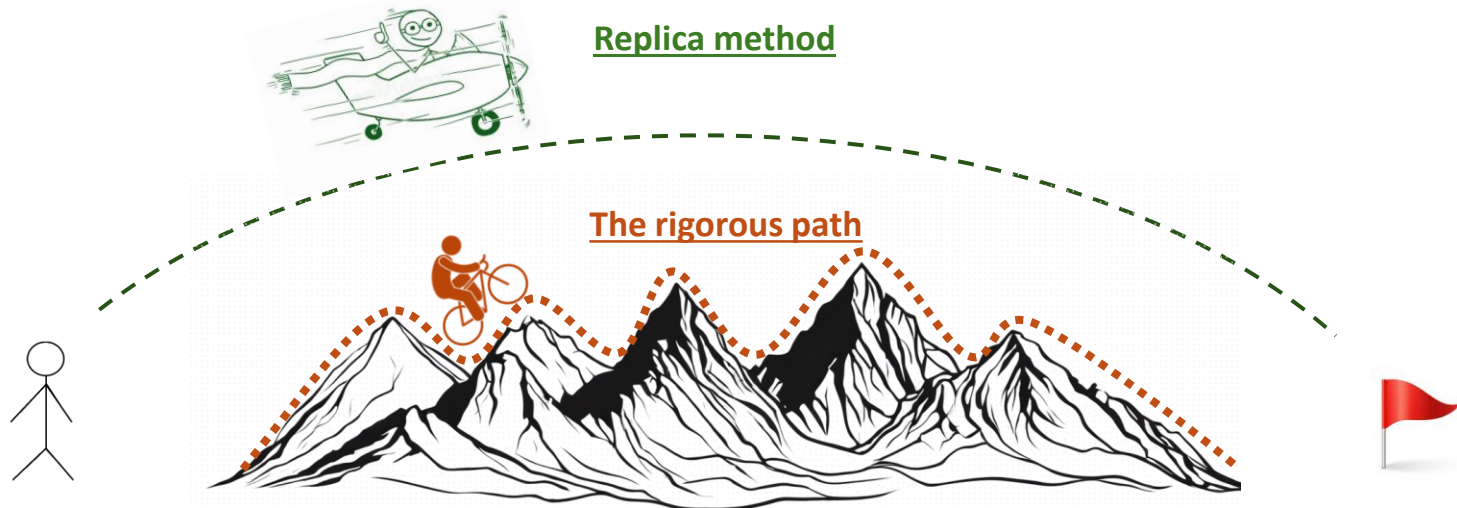- If $n = \mathcal{O}(d)$, the optimal error can be reached by **linear regression**...  Cui&al '23

- But there are $\Theta(d^2)$ weights to learn...

  What happens for $n = \Theta(d^2)$
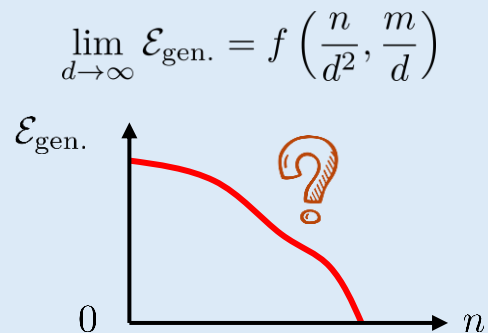
# All roads lead to Rome



Replica method

The rigorous path

$$m = \Theta(d) \qquad n = \Theta(d^2)$$

$$\left\{ y_i = \frac{1}{m} \sum_{k=1}^{m} \left[ \frac{1}{\sqrt{d}} (\mathbf{w}_k^*)^T \cdot \mathbf{x}_i \right]^2 \right\}_{i=1}^{n}$$

$$\lim_{d \to \infty} \mathcal{E}_{\text{gen.}} = f\left( \frac{n}{d^2}, \frac{m}{d} \right)$$

$\mathcal{E}_{\text{gen.}}$

$0$     $n$

# Taking the long road

**Step 0:**
$$y = \frac{1}{m} \sum_{k=1}^{m} \left[ \frac{1}{\sqrt{d}} (\mathbf{w}_k^*)^T \cdot \mathbf{x} \right]^2 = \frac{1}{d} \mathbf{x}^\top \mathbf{S}^\star \mathbf{x} = \mathrm{Tr}[\mathbf{S}^\star \mathbf{\Phi}]$$

$$\mathbf{S}^\star := \frac{1}{m} \sum_{k=1}^{m} \mathbf{w}_k^\star (\mathbf{w}_k^\star)^\top \sim \mathcal{W}_{m,d} \qquad \mathbf{\Phi} := \frac{1}{d} \mathbf{x}\mathbf{x}^\top$$

Can be generalized to **noisy pre-activations**
$$\mathbf{w}_k^\star \cdot \mathbf{x} \rightarrow \mathbf{w}_k^\star \cdot \mathbf{x} + \sqrt{\Delta}\xi_k$$

$$y \sim P_{\mathrm{out}}\left(\cdot | \mathrm{Tr}[\mathbf{S}^\star \mathbf{\Phi}]\right)$$

**Goal:** $\{y_i, \mathbf{x}_i\}_{i=1}^n \implies \hat{\mathbf{S}}_{\mathrm{opt.}} = \arg\min \mathcal{E}_{\mathrm{gen.}}(\hat{\mathbf{w}}_k) = \arg\min \|\hat{\mathbf{S}} - \mathbf{S}^\star\|_F^2 \simeq$ **planted "ellipsoid fitting-like" problem**

**Step 1 : "Gaussian universality"**     $n = \Theta(d^2)$     ⚠ Same scaling regime as ellipsoid fitting !

**Universality** of Bayes-optimal generalization error

Leverages our ellipsoid fitting analysis [**M.** & Bandeira '23]

**CRITICAL**

$$\min \mathcal{E}_{\mathrm{gen.}}(\hat{\mathbf{w}}_k) = \min \|\hat{\mathbf{S}} - \mathbf{S}^\star\|_F^2 \quad = \quad \min \widetilde{\mathcal{E}}_{\mathrm{gen.}}(\hat{\mathbf{S}}) = \min \|\hat{\mathbf{S}} - \mathbf{S}^\star\|_F^2 \times (1 + o(1))$$

from $\{y_i \sim P_{\mathrm{out}}\left(\cdot | \mathrm{Tr}[\mathbf{S}^\star \mathbf{\Phi}_i]\right)\}_{i=1}^n$     from $\{\tilde{y}_i \sim P_{\mathrm{out}}(\cdot | \mathrm{Tr}[\mathbf{S}^\star \mathbf{G}_i])\}_{i=1}^n$

**Gaussian** matrix

# Taking the long road

**Step 2 :** $\{\tilde{y}_i \sim P_{\text{out}}(\cdot | \text{Tr}[\mathbf{S}^\star \mathbf{G}_i])\}_{i=1}^n$   Just a (generalized) **linear model on** $\mathbf{S}^\star$ , with…

➤ Gaussian data $\mathbf{G} := \begin{pmatrix} \text{flatt}(\mathbf{G}_1) \\ \vdots \\ \text{flatt}(\mathbf{G}_n) \end{pmatrix}$ ➕ ➤ Wishart prior $\mathbf{S}^\star \sim \mathcal{W}_{m,d}$

Generalization of
[Barbier & al '19]

"Replica-symmetric" formula for $\widetilde{\mathcal{E}}_{\text{gen.}}$

Involves … → Scalar estimation problem involving $P_{\text{out}}$ ✅

**Step 3 :** **Denoising** problem : $\mathbf{Y} = \sqrt{\lambda}\mathbf{S}^\star + \mathbf{Z}$ ⟶ $\mathbf{S}^\star$ ❓ ✅
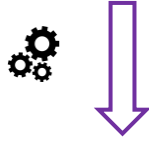
Gaussian (GOE) matrix

[Bun & al '16 ; **M.**, Krzakala & al '22 ; Pourkamali & al '23 ; Semerjian '24 ; …]

❑ The optimal estimator is **spectral :** $\mathbf{Y} = \mathbf{O}\mathbf{D}\mathbf{O}^\top$ ⟹ $\hat{\mathbf{S}}(\mathbf{Y}) = \mathbf{O} f_{\text{opt.}}(\mathbf{D})\mathbf{O}^\top$

❑ Analytical expressions for $f_{\text{opt.}}$ and the **asymptotic MMSE** $\lim_{d \to \infty} \|\hat{\mathbf{S}}(\mathbf{Y}) - \mathbf{S}^\star\|_F^2$

# Taking the long road

$$\left\{ y_i = \frac{1}{m} \sum_{k=1}^{m} \left[ \frac{1}{\sqrt{d}} (\mathbf{w}_k^*)^T \cdot \mathbf{x}_i + \sqrt{\Delta} \xi_k \right]^2 \right\}_{i=1}^{n}$$

$$d \to \infty \qquad \begin{aligned} m &= \kappa d \\ n &= \alpha d^2 \end{aligned}$$

**Combining all steps...**

$$\lim_{d \to \infty} \mathcal{E}_{\text{gen.}} = 2\kappa\alpha\zeta - \Delta(2 + \Delta)$$

$\zeta$ solves the self-consistent equation
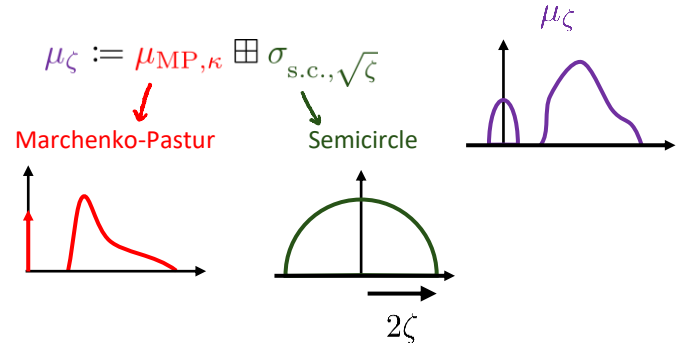
$$(1 - 2\alpha) + \frac{\Delta(2 + \Delta)}{\kappa\zeta} = \frac{4\pi^2\zeta}{3} \int \mu_\zeta(y)^3 \, \mathrm{d}y$$

$$\mu_\zeta := \mu_{\mathrm{MP},\kappa} \boxplus \sigma_{\mathrm{s.c.}, \sqrt{\zeta}}$$

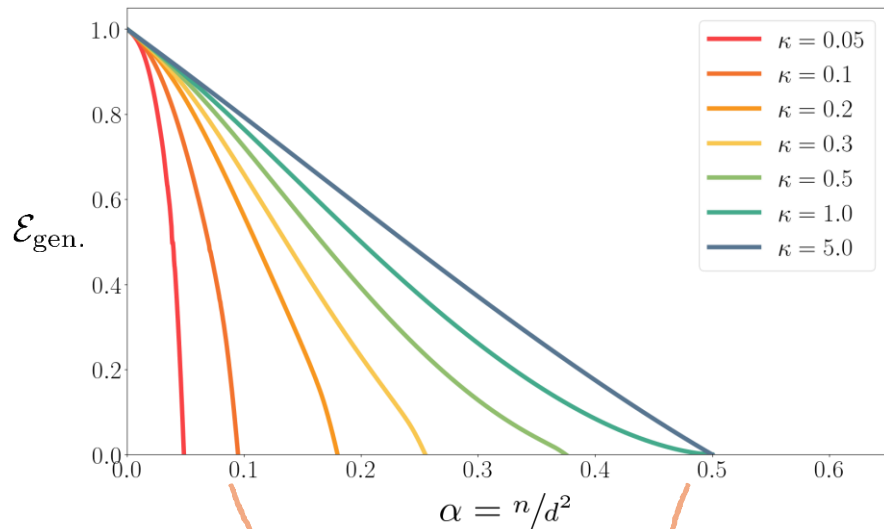Marchenko-Pastur       Semicircle

$\mu_\zeta$

$2\zeta$

➤ **Easy-to-evaluate formula** for the Bayes-optimal generalization error

➤ Not a fully rigorous theorem, some mathematical subtleties in **Steps 1 and 2**

# Optimal generalization error

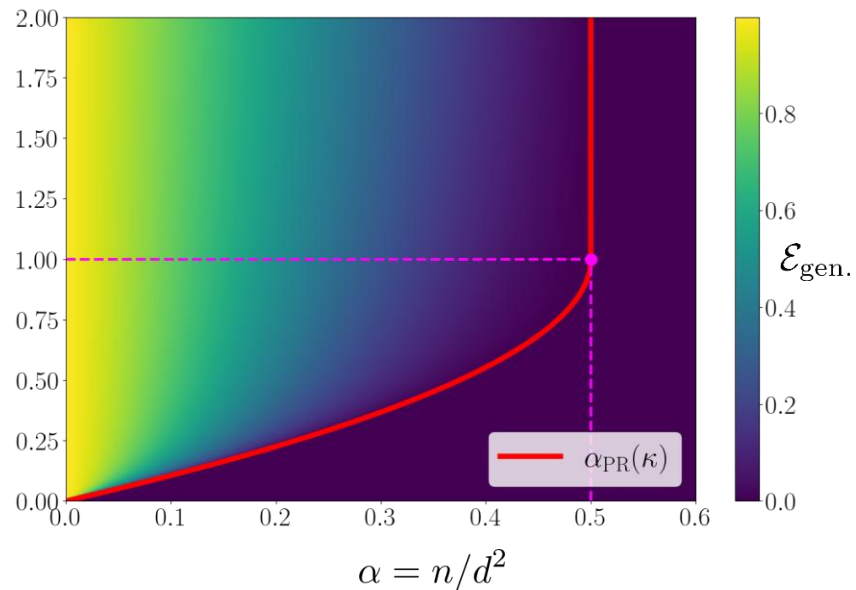Intensive width $\kappa = {}^m/_d$ ;   Sample complexity $\alpha = {}^n/_{d^2}$



Noiseless setting : $\Delta = 0$

$$\alpha_{\mathrm{PR}}(\kappa) = \min\left(\kappa - \frac{\kappa^2}{2}, \frac{1}{2}\right)$$

**Perfect recovery threshold**

Matches a naïve "counting argument"   $\mathrm{DOF}[\{\mathbf{S} : \mathbf{S} = \mathbf{S}^\top \text{ and } \mathrm{rk}(\mathbf{S}) \leq \kappa d\}] \simeq \alpha_{\mathrm{PR}}(\kappa) d^2$

# The GAMP-RIE algorithm

$$y_i \sim P_{\mathrm{out}}\left(\cdot \,|\, \mathrm{Tr}[\mathbf{S}^\star \boldsymbol{\Phi}_i]\right)$$

$$\boldsymbol{\Phi}_i := (\mathbf{x}_i \mathbf{x}_i^\mathsf{T} - \mathrm{I}_d)/\sqrt{d}$$

$$\mathbf{S}^\star \sim \mathcal{W}_{m,d} \quad \text{(Wishart)}$$

$P_{\mathrm{out}}$ : Noise channel

**Informal hypothesis**

Universality $\boldsymbol{\Phi}_i \Rightarrow \mathbf{G}_i$ also holds "at the level of algorithms"

MSE-optimal algorithm

Generalized linear model w. Gaussian data $\longrightarrow$ **Generalized Approximate Message Passing (GAMP)** [Donoho&al '09 ; Rangan '11 ; …]

Each GAMP iteration solves
$$\mathbf{Y} = \sqrt{\lambda}\mathbf{S}^\star + \mathbf{Z} \longrightarrow \boxed{\mathbf{S}^\star} \;?$$
$\longrightarrow$ **Rotationally-Invariant Estimator (RIE)**
$$\mathbf{Y} = \mathbf{O}\mathbf{D}\mathbf{O}^\top \Rightarrow \hat{\mathbf{S}}(\mathbf{Y}) = \mathbf{O}f_{\mathrm{opt.}}(\mathbf{D})\mathbf{O}^\top$$
[Bun & al '16 ; …]

Known "optimal shrinkage" function

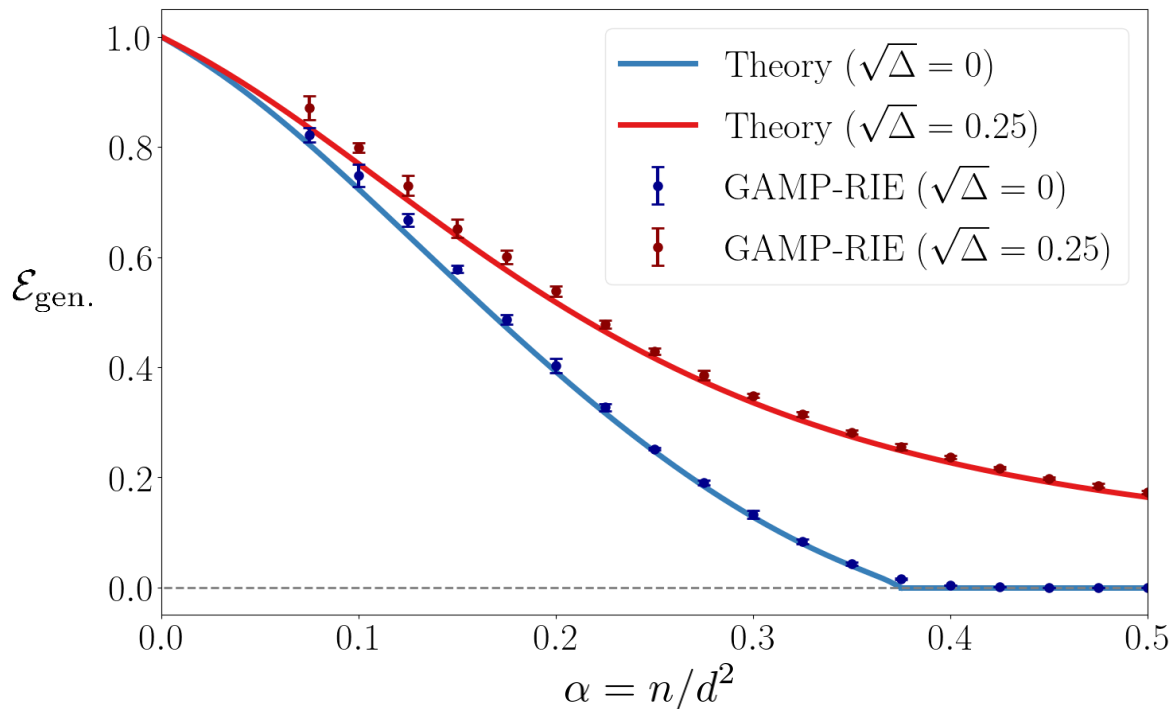**An explicit easy-to-implement polynomial-time algorithm**

GAMP ⸺ RIE

# Absence of hard phase

Intensive width $\kappa = {}^m/d$ ;   Sample complexity $\alpha = {}^n/d^2$

$\kappa = 0.5$
$d = 200$

$$y_i = f_{\mathbf{W}^*}(\mathbf{x}_i) := \frac{1}{m} \sum_{k=1}^{m} \left[ \frac{1}{\sqrt{d}} (\mathbf{w}_k^*)^T \cdot \mathbf{x}_i + \sqrt{\Delta}\xi_{k,i} \right]^2$$



**No computational-to-statistical gap / hard phase**

For $m = \mathcal{O}(1)$ there is a hard phase !

cf. [Barbier & al '19] for $m = 1$

# Gradient descent : noiseless case

Intensive width $\kappa = m/d$ ;  Sample complexity $\alpha = n/d^2$

$$\mathcal{L}(\mathbf{W}) := \frac{1}{n} \sum_{i=1}^{n} \left( y_i - \tilde{f}_{\mathbf{W}}(\mathbf{x}_i) \right)^2 \text{, where } \tilde{f}_{\mathbf{W}}(\mathbf{x}) := \frac{1}{m} \sum_{k=1}^{m} \left[ \frac{1}{\sqrt{d}} (\mathbf{w}_k)^T \cdot \mathbf{x} \right]^2$$
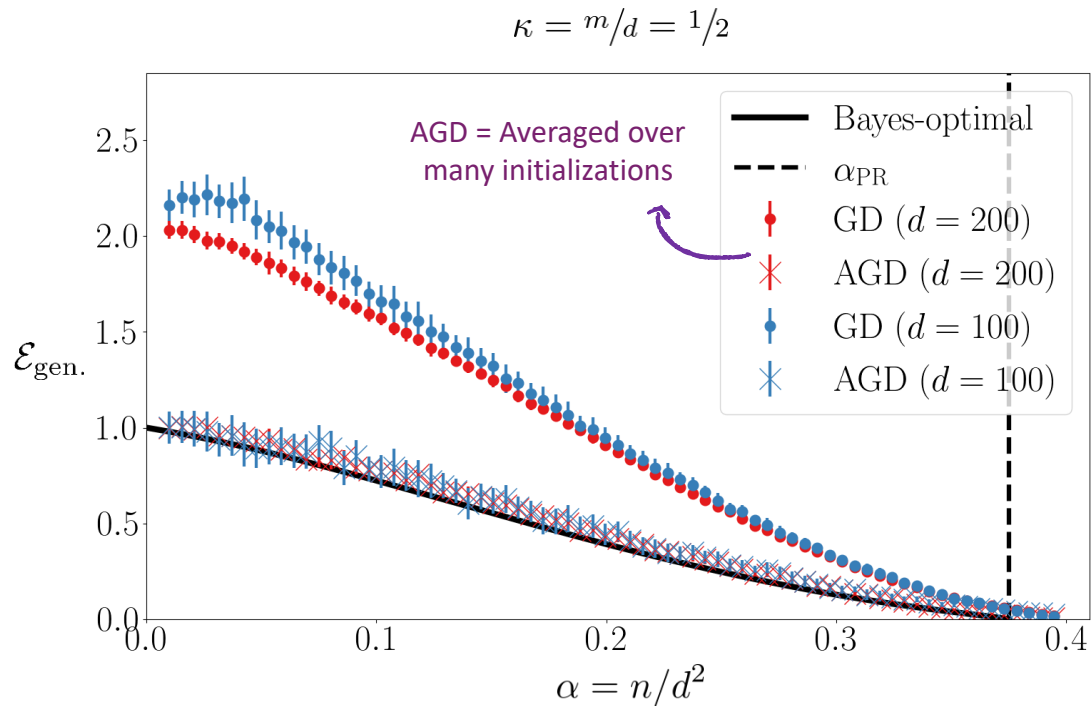
For **any** $\kappa$, AGD seems to reach the Bayes-optimal MMSE

➤ For $\kappa \geq 1$ ($m \geq d$), the problem is **convex**

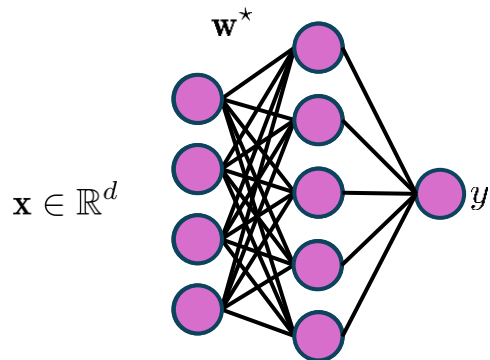over $\mathbf{S} := (1/m) \sum_{k=1}^{m} \mathbf{w}_k \mathbf{w}_k^\top$

The landscape of $\mathcal{L}(\mathbf{W})$ trivializes in this regime
[Du & Lee '18 ; Soltanolkotabi & al '18 ; Venturi & al '19]

➤ For $\kappa < 1$, **non-convex problem**. Still,

naïve GD reaches optimal error !

No longer true for noisy pre-activations.



$\kappa = m/d = 1/2$

AGD = Averaged over many initializations

Legend:
— Bayes-optimal
-- $\alpha_{\mathrm{PR}}$
GD ($d = 200$)
AGD ($d = 200$)
GD ($d = 100$)
AGD ($d = 100$)

$\mathcal{E}_{\mathrm{gen.}}$

$\alpha = n/d^2$

# Summary

$\mathbf{w}^\star$

$\mathbf{x} \in \mathbb{R}^d$

$y$

1. Analytical formula for the **Bayes-optimal generalization error.**

2. Optimal algorithm (GAMP-RIE), **no computational-statistical gap.**

3. (Averaged) **Gradient descent seems to sample from the posterior for noiseless pre-activations**, even in the non-convex regime $\kappa < 1$!
   Not true for noisy case

$$\left\{ y_i = f_{\mathbf{W}^*}(\mathbf{x}_i) := \frac{1}{m} \sum_{k=1}^{m} \left[ \frac{1}{\sqrt{d}} (\mathbf{w}_k^*)^T \cdot \mathbf{x}_i + \sqrt{\Delta} \xi_{k,i} \right]^2 \right\}_{i=1}^{n}$$

$\sim \mathcal{N}(0, \mathrm{I}_d)$

$\mathbf{w}_k^\star \sim \mathcal{N}(0, \mathrm{I}_d)$

$n = \alpha d^2 \; ; \;\; m = \kappa d$

**WHAT NEXT?**

❖ **Other activations ?** (beyond quadratic) **Other architectures ?**

❖ Theoretical analysis of GD properties ?

❖ …

**THANK YOU !**