

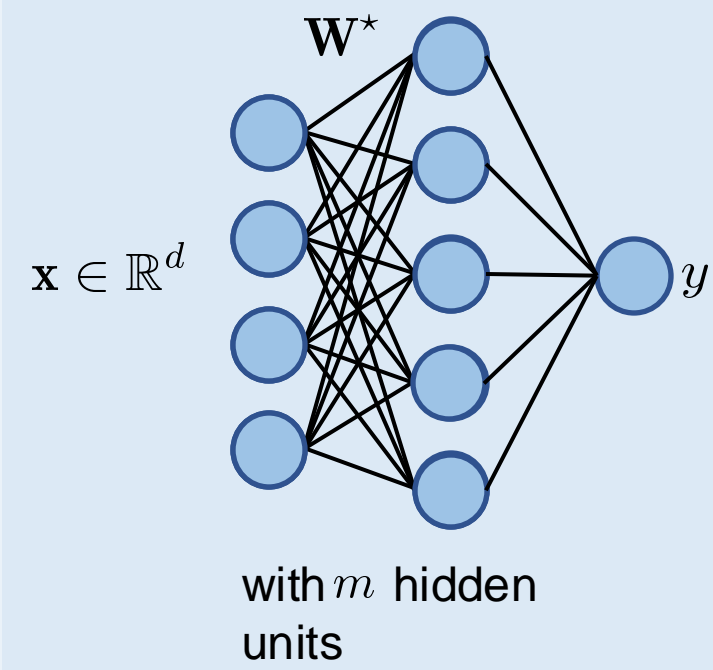


Bayes-optimal learning of an extensive-width neural network from quadratically many samples

Antoine Maillard¹, Emanuele Troiani², Simon Martin³, Florent Krzakala², Lenka Zdeborová²
[1] ETH Zürich [2] EPFL [3] ENS Paris



Learning in large neural networks



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

$\sim \mathcal{N}(0, I_d)$ $\mathbf{w}_k^* \sim \mathcal{N}(0, I_d)$

Learning from data
 $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \Rightarrow \mathbf{W}^*$?

High-dimensional limit
 $d \rightarrow \infty; m = \Theta(d)$

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

- The scaling $m = \Theta(d)$ is an interesting one to study wide networks:
 - Scaling $m = \mathcal{O}(1)$: multi-index phase retrieval [5]
 - Scaling $m \gg d$: linear regression reaches the optimal error
- With $n = \mathcal{O}(d)$ samples, no information on \mathbf{w}^* can be retrieved [1]
- But there are $\Theta(d^2)$ weights to learn
Information theory hints at the scaling $n = \Theta(d^2)$

Optimal generalization error

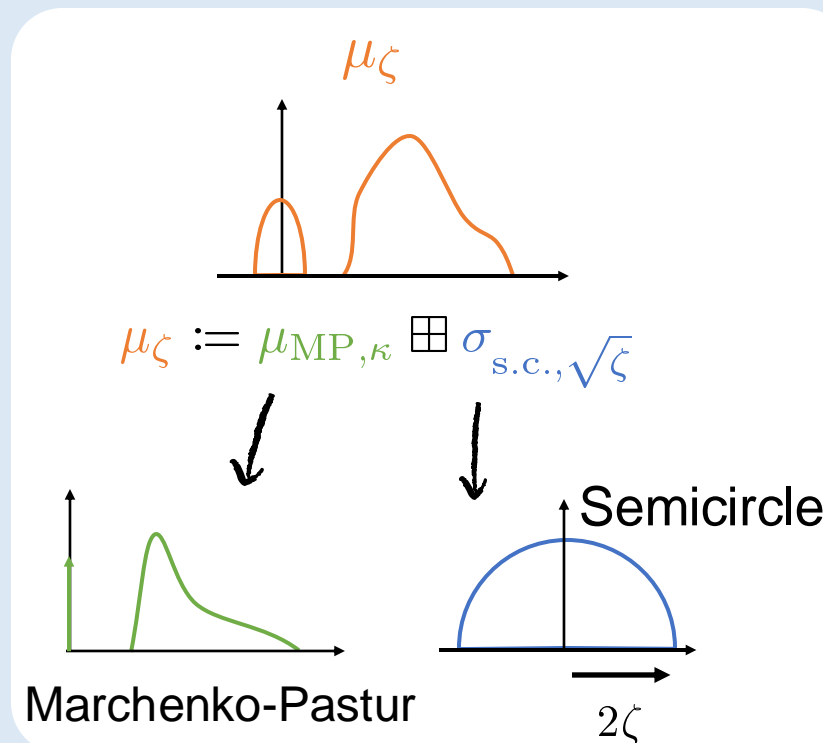
$d \rightarrow \infty$ $m = \kappa d$ $n = \alpha d^2$

$$\left\{ y_i = \frac{1}{m} \sum_{k=1}^m \left[\frac{1}{\sqrt{d}} (\mathbf{w}_k^*)^T \cdot \mathbf{x}_i + \sqrt{\Delta} \xi_k \right]^2 \right\}_{i=1}^n \Rightarrow \lim_{d \rightarrow \infty} \mathcal{E}_{\text{gen.}} = 2\kappa\alpha\zeta - \Delta(2 + \Delta)$$

ζ solves the self-consistent equation

$$(1 - 2\alpha) + \frac{\Delta(2 + \Delta)}{\kappa\zeta} = \frac{4\pi^2\zeta}{3} \int \mu_\zeta(y)^3 dy$$

➤ **Easy-to-evaluate formula** for the Optimal generalization error



Open questions

- Analyze **other activations** (beyond quadratic)
Related to extensive rank tensor denoising
- Account for the **structure** of the data
- Theoretical** analysis of GD properties

Sketch of the analysis

1 A planted linear matrix model

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

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

Wishart prior

Can be generalized to **noisy pre-activations**

$$\mathbf{w}_k^* \cdot \mathbf{x} \rightarrow \mathbf{w}_k^* \cdot \mathbf{x} + \sqrt{\Delta} \xi_k \iff y \sim P_{\text{out}}(\cdot | \text{Tr}[\mathbf{S}^* \Phi])$$

Universality of optimal generalization error [3]

$$\min_{\hat{\mathbf{w}}_k} \mathcal{E}_{\text{gen.}}(\hat{\mathbf{w}}_k) = \min \|\hat{\mathbf{S}} - \mathbf{S}^*\|_F^2 \quad \approx \quad \min_{\hat{\mathbf{S}}} \tilde{\mathcal{E}}_{\text{gen.}}(\hat{\mathbf{S}}) = \min \|\hat{\mathbf{S}} - \mathbf{S}^*\|_F^2$$

from $\{y_i \sim P_{\text{out}}(\cdot | \text{Tr}[\mathbf{S}^* \Phi_i])\}_{i=1}^n$ from $\{\tilde{y}_i \sim P_{\text{out}}(\cdot | \text{Tr}[\mathbf{S}^* \mathbf{G}_i])\}_{i=1}^n$

Gaussian matrix

Just a (generalized) **linear model on \mathbf{S}^***

2 Mapping to GLM: scalar estimation part

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

Formula for $\tilde{\mathcal{E}}_{\text{gen.}}$ (generalization of [4]) It involves **Scalar estimation problem involving P_{out}** and **Matrix denoising on the Wishart prior**

3 Mapping to GLM: matrix denoising part

Denoising problem $\mathbf{Y} = \sqrt{\lambda} \mathbf{S}^* + \mathbf{Z} \rightarrow \mathbf{S}^*$?

Gaussian (GOE) matrix

□ The optimal estimator is **spectral** [3]: $\mathbf{Y} = \mathbf{O} \mathbf{D} \mathbf{O}^T \iff \hat{\mathbf{S}}(\mathbf{Y}) = \mathbf{O} f_{\text{opt.}}(\mathbf{D}) \mathbf{O}^T$

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

Optimal algorithm: GAMP-RIE

A provably optimal for large d , easy-to-implement polynomial-time algorithm

GAMP
RIE

$$\omega_i^t = \frac{\mathbf{x}_i^T \hat{\mathbf{S}} \mathbf{x}_i}{d} - g_{\text{out}}(y_i, \omega_i^{t-1}, V^{t-1}) V^t$$

$$A^t = \frac{2\alpha}{n} \sum_{i=1}^n g_{\text{out}}(y_i, \omega_i^t, V^t)^2$$

$$\mathbf{R}^t = \hat{\mathbf{S}} + \frac{1}{d^{3/2} A^t} \sum_{i=1}^n g_{\text{out}}(y_i, \omega_i^t, V^t) (\mathbf{x}_i \mathbf{x}_i^T - \mathbf{I})$$

$$\hat{\mathbf{S}}^{t+1} = f_{\text{RIE}} \left(\mathbf{R}^t, \frac{1}{2A^t} \right) \quad V^{t+1} = 2f_{\text{RIE}} \left(\frac{1}{2A^t} \right)$$

Generalized linear model w. Gaussian data with non-separable prior

Generalized Approximate Message Passing (GAMP) [2]

Each GAMP iteration solves $\mathbf{Y} = \sqrt{\lambda} \mathbf{S}^* + \mathbf{Z} \rightarrow \mathbf{S}^*$?

Rotationally-Invariant Estimator (RIE) [3]
 $\mathbf{Y} = \mathbf{O} \mathbf{D} \mathbf{O}^T \iff \hat{\mathbf{S}}(\mathbf{Y}) = f_{\text{RIE}}(\mathbf{Y}) = \mathbf{O} f_{\text{opt.}}(\mathbf{D}) \mathbf{O}^T$

Generalization error curves

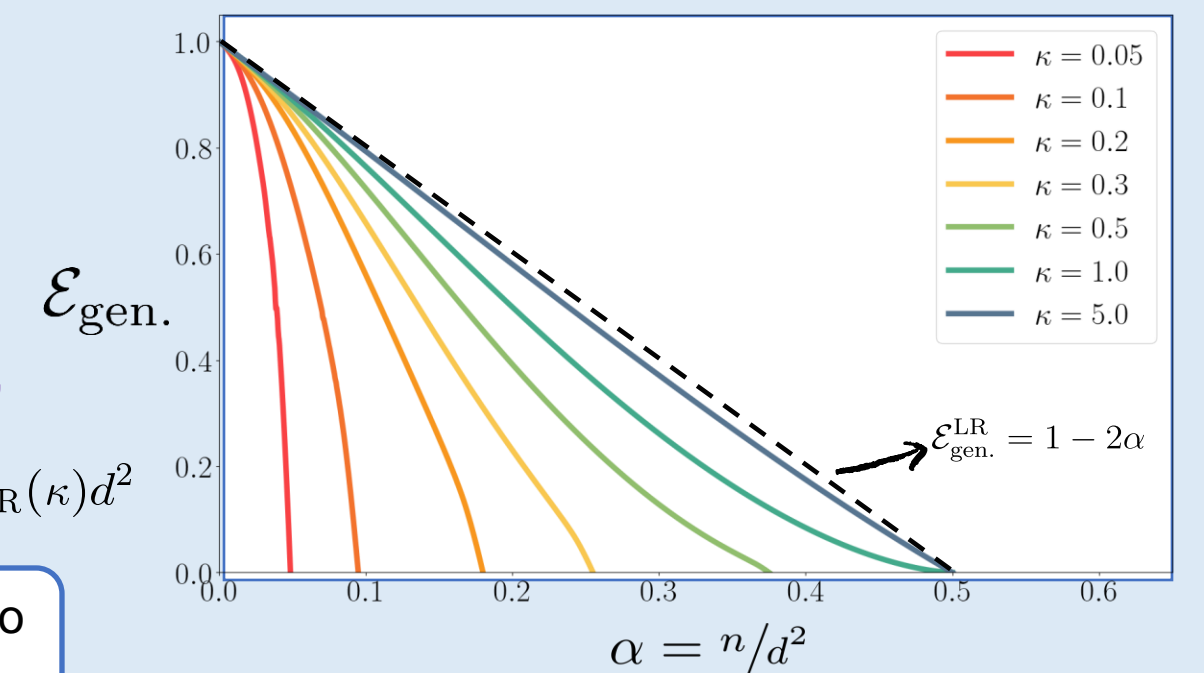
Noiseless setting : $\Delta = 0$

Perfect recovery threshold

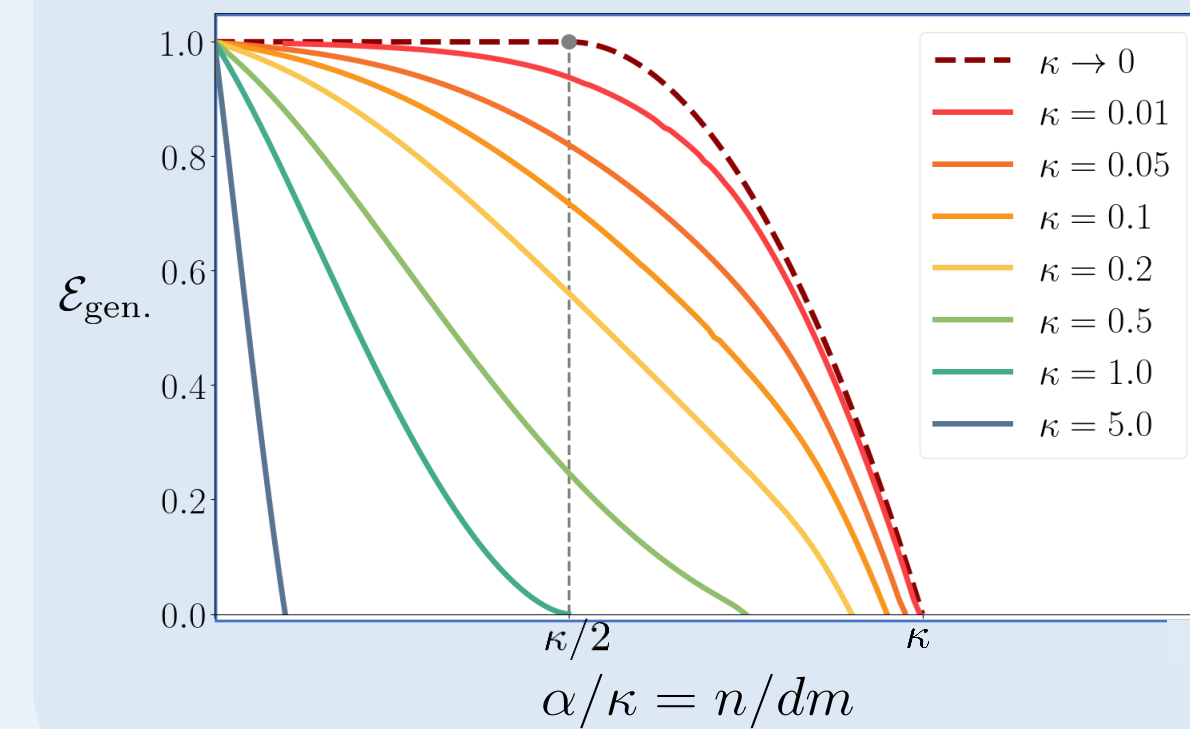
$$\alpha_{\text{PR}}(\kappa) = \min \left(\kappa - \frac{\kappa^2}{2}, \frac{1}{2} \right)$$

Matches a naïve “counting argument”

$$\text{DOF}[\{\mathbf{S} : \mathbf{S} = \mathbf{S}^T \text{ and } \text{rk}(\mathbf{S}) \leq \kappa d\}] \simeq \alpha_{\text{PR}}(\kappa) d^2$$



No computational to statistical gap



Significantly lower generalization error than linear regression ($\kappa \rightarrow \infty$):

$$\mathcal{E}_{\text{gen.}}^{\text{LR}} = 1 - 2\alpha$$

Smooth transition from hidden layer size $m = \mathcal{O}(1)$ to $m = \Theta(d)$ [5]

Gradient descent performance

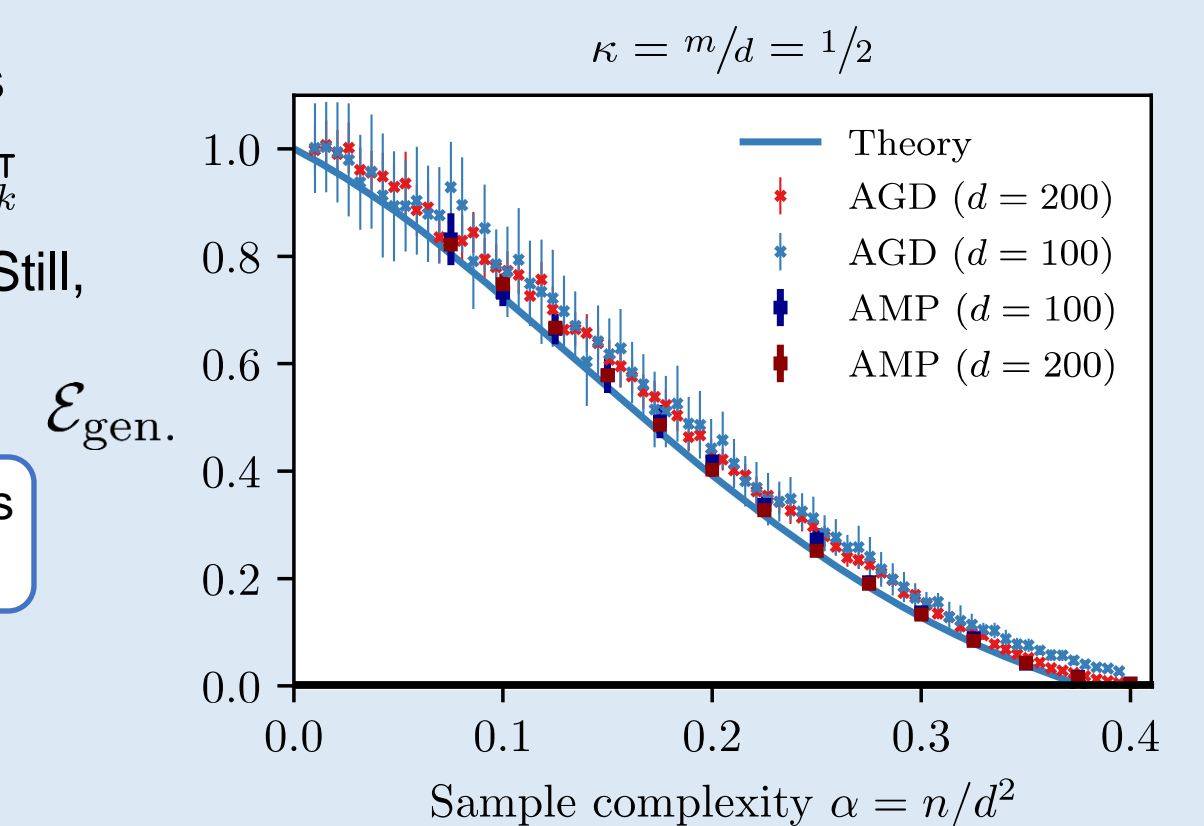
Noiseless setting : $\Delta = 0$

$$\mathcal{L}(\mathbf{W}) := \frac{1}{n} \sum_{i=1}^n (y_i - \tilde{f}_{\mathbf{W}}(\mathbf{x}_i))^2, \text{ where } \tilde{f}_{\mathbf{W}}(\mathbf{x}) := \frac{1}{m} \sum_{k=1}^m \left[\frac{1}{\sqrt{d}} (\mathbf{w}_k)^T \cdot \mathbf{x} \right]^2$$

➤ For $\kappa \geq 1$ ($m \geq d$), the problem is **convex** over $\mathbf{S} := (1/m) \sum_{k=1}^m \mathbf{w}_k \mathbf{w}_k^T$

➤ For $\kappa < 1$, **non-convex problem**. Still, naïve GD reaches optimal error !

For any κ , (averaged) GD seems to reach the optimal MMSE



References

- H. Cui, F. Krzakala, L. Zdeborová (2013). “Bayes-optimal learning of deep random networks of extensive-width” ICML 2013 [Arxiv 2302.00375]
- D. L. Donoho, A. Maleki, A. Montanari (2009). “Message-passing algorithms for compressed sensing” Proceeding of the National Academy of Sciences [Arxiv 0907.3574]
- J. Bun, R. Allez, J. P. Bouchaud, M. Potters (2016). “Rotational invariant estimators for general noisy matrices” IEEE Transactions on Information Theory [Arxiv 1502.06736]
- A. Maillard, A. Bandeira (2023). “Exact threshold for approximate ellipsoid fitting of random points” [Arxiv 2310.05787]
- E. Troiani, Y. Dandl, L. DeFilippis, L. Zdeborová, B. Loureiro, F. Krzakala (2024). “Fundamental computational limits of weak learnability in high-dimensional multi-index models” [Arxiv 2405.15480]