

Learning from long sequences of high-dimensional tokens, and extensive-width neural networks

Antoine Maillard

The Inria logo is a stylized, cursive red script that reads "Inria". It is positioned to the right of the author's name.

Les Houches – February 6th 2025

References

- ❖ Exact threshold for approximate ellipsoid fitting of random points ([arXiv:2310.05787](#))
- ❖ Bayes-optimal learning of an extensive-width neural network from quadratically many samples ([arXiv:2408.03733](#), *NeurIPS '24*)
- ❖ Bilinear Sequence Regression: A Model for Learning from Long Sequences of High-dimensional Tokens ([arXiv:2410.03733](#))



Afonso Bandeira



Lenka Zdeborová



Florent Krzakala



Emanuele Troiani



Simon Martin



Vittorio Erba

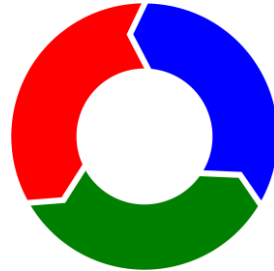


Luca Biggio

Solvable models of learning

High-dimensional statistical / modelling learning

Data Sequences of
(tokens) $\sim p$




Model
 $y = f(\theta, x)$
/ ...

Algorithm $\hat{\theta} = \arg \min_{\theta} \mathcal{R}(\{x_i, y_i\}_{i=1}^n, \theta)$
(and variants)

Why is it advantageous to present the data as long sequences of high-dimensional tokens ?

Objectives:

- Build **solvable models** that retain the **key ingredients** 
- A **rigorous** theory describing learning in these models

What is the **simplest** exactly solvable model that exhibits advantages in learning from long sequences of high-dimensional tokens ?



Bilinear Sequence Regression (BSR)

Erba, Troiani, Biggio, M. & Zdeborová '24

What is the **simplest** exactly solvable model that exhibits advantages in learning from long sequences of high-dimensional tokens ?

The simplest regression model for **vectorized data** $\mathbf{x} \in \mathbb{R}^d$ is **(generalized) linear regression** $y = \varphi(\mathbf{x} \cdot \mathbf{w})$

$\mathbf{X} \in \mathbb{R}^{L \times d}$: sequence of L d -dimensional tokens

L : length of the sequences

d : embedding dimension of tokens

$L, d \gg 1$

r rank/width of the model

The **most basic** regression model is

$$y = f_{\mathbf{M}}(\mathbf{X}) = \varphi(\text{Tr}[\mathbf{X}\mathbf{M}^T])$$



Bilinear Sequence Regression (BSR)

$$f_{\mathbf{U}, \mathbf{V}}(\mathbf{X}) = \varphi(\text{Tr}[\mathbf{X}\mathbf{U}\mathbf{V}^T]) = \varphi\left(\sum_{a,i=1}^{L,d} X_{ai} \sum_{\gamma=1}^r U_{i\gamma} V_{\gamma a}\right)$$

\mathbf{U} : token space

\mathbf{V} : sequence space

Also called **Matrix Sensing** Recht&al '10, Gunasekar&al '17, ...

Like linear regression is a base model for non-sequential data, BSR is a toy base model for sequences of tokens

Setting

Teacher-student

$$y_\mu = \text{Tr}[\mathbf{X}_\mu \mathbf{U}^* (\mathbf{V}^*)^\top] = \sum_{a,i=1}^{L,d} X_{ai}^\mu \sum_{\gamma=1}^r U_{i\gamma}^* V_{\gamma a}^*$$

“Generic”
 Ex: $\sim \mathcal{N}(0, \mathbf{I}_{L \times d})$

i.i.d. $\sim \mathcal{N}(0, 1)$

Learning from data

$$\mathcal{D} = \{(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_n, y_n)\} \Rightarrow (\mathbf{U}^*, \mathbf{V}^*) ?$$

Scalings

High-dimensional

$$L, d \rightarrow \infty; L = \Theta(d)$$

Long sequences of **high-dimensional** tokens

Extensive-width

$$r = \Theta(d)$$

Related to **extensive-width** neural networks

Talks of Guilhem & Jean

Number of samples

$$n = \Theta[r(L + d)]$$

For **non-trivial** gen. error

Low-rank/width

$$r = \Theta(1)$$

[Schülke&al ‘16]

Optimal gen. error, and optimal algorithm (AMP)

Objectives

$$L = \Theta(d); r = \Theta(d) \quad n = \Theta[r(L + d)]$$

$$y_\mu = \text{Tr}[\mathbf{X}_\mu \mathbf{U}^* (\mathbf{V}^*)^\top] = \sum_{a,i=1}^{L,d} X_{ai}^\mu \sum_{\gamma=1}^r U_{i\gamma}^* V_{\gamma a}^*$$

↑ "Generic"
↑ ↑ i.i.d. $\mathcal{N}(0, 1)$

Ex: $\sim \mathcal{N}(0, \mathbf{I}_{L \times d})$

Teacher-student model

Training dataset $\mathcal{D} = \{\mathbf{X}^\mu, y^\mu\}_{\mu=1}^n \Rightarrow \mathbf{U}^*, \mathbf{V}^*$?

Objectives

Bayes-optimal generalization error

$$\mathcal{E}_{\text{gen.}} := \mathbb{E}_{\mathbf{U}^*, \mathbf{V}^*, \mathcal{D}} \min_{\hat{y}(\mathcal{D})} \mathbb{E}_{\mathbf{X}_{\text{test}}} [(\hat{y}(\mathbf{X}_{\text{test}}) - f_{\mathbf{U}^*, \mathbf{V}^*}(\mathbf{X}_{\text{test}}))^2]$$

- Sharp thresholds** ? Phase transitions ?
- Efficiently achievable ? **Hard phases** ?

Comparison to **minimal nuclear-norm** estimator

$$\min\{\|\mathbf{S}\|_{\text{NN}} : \mathbf{S} = \mathbf{U}\mathbf{V}^\top, y_\mu = \text{Tr}[\mathbf{X}_\mu \mathbf{S}]\}$$

Recht&al '10

Comparison to **vectorized data** ? (linear regression)

$$\mathbf{S}^* = \mathbf{U}^* (\mathbf{V}^*)^\top \Rightarrow \mathbf{S}^* \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

- Performance of **GD-based algorithms** ?
Bhojanapalli&al '16

$$\mathcal{L}(\mathbf{U}, \mathbf{V}) = \sum_{\mu=1}^n (y_\mu - \text{Tr}[\mathbf{X}_\mu \mathbf{U}\mathbf{V}^\top])^2$$

- Is GD an **implicit nuclear-norm minimizer** ?
Gunasekar&al '17

Result I : Bayes-optimal error

$$\mathcal{E}_{\text{gen.}} := \mathbb{E}_{\mathbf{U}^*, \mathbf{V}^*, \mathcal{D}} \min_{\hat{y}(\mathcal{D})} \mathbb{E}_{\mathbf{X}_{\text{test}}} [(\hat{y}(\mathbf{X}_{\text{test}}) - f_{\mathbf{U}^*, \mathbf{V}^*}(\mathbf{X}_{\text{test}}))^2]$$



$$\lim_{d \rightarrow \infty} \mathcal{E}_{\text{gen.}} = \alpha \zeta$$



- ❑ **Easy-to-evaluate formula** for the Bayes-optimal generalization error
- ❑ Derivation via the **replica method** of stat. physics
- ❑ Paper with **proofs** in preparation
[Xu, M., Krzakala, Zdeborová '25]

$$\left\{ y_\mu = \sum_{a,i=1}^{L,d} X_{ai}^\mu \sum_{\gamma=1}^r U_{i\gamma}^* V_{\gamma a}^* \right\}_{\mu=1}^n$$

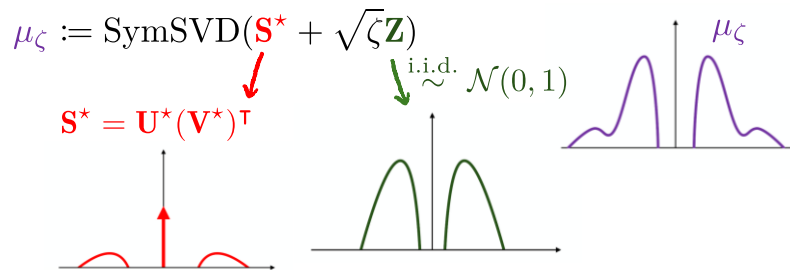
Scaling regime

$$\rho = \frac{r}{d} \quad \beta = \frac{L}{d} \quad \alpha = \frac{n}{dL}$$

ζ solves the self-consistent equation

$$(1 - \alpha) = \frac{2\zeta}{\beta^{3/2}} \int dy \mu_\zeta(y) \left[\frac{(\beta - 1)^2}{2x^2} + \frac{2\pi^2}{3} \mu_\zeta(y)^2 \right]$$

Analytical form using **free probability** tools



Derivation: a symmetric variant to BSR

M., Troiani, Martin, Krzakala & Zdeborová '24

Bilinear Sequence Regression (BSR)

$$f_{\mathbf{U}, \mathbf{V}}(\mathbf{X}) = \varphi(\text{Tr}[\mathbf{XUV}^T]) = \varphi\left(\sum_{a,i=1}^{L,d} X_{ai} \sum_{\gamma=1}^r U_{i\gamma} V_{\gamma a}\right)$$

A symmetric variant

$$\begin{matrix} L = d \\ \mathbf{X} = \mathbf{X}^T \end{matrix} + f_{\mathbf{S}}(\mathbf{X}) = \varphi(\text{Tr}[\mathbf{XS}])$$

Example

“Extensive-width sign retrieval”



Extensive-width 2-layer NN with a quadratic activation

$$\mathbf{X} = \mathbf{x}\mathbf{x}^T \text{ rank-1}$$

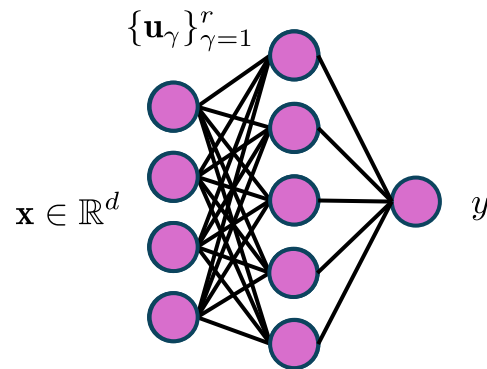
$$\mathbf{S} = \sum_{\gamma=1}^r \mathbf{u}_{\gamma} \mathbf{u}_{\gamma}^T \succeq 0$$

High-dimensional limit

$$d \rightarrow \infty; r = \Theta(d); n = \Theta(d^2)$$

Same objectives as BSR

Cf Jean's talk on Tuesday



$$y_{\mu} = \mathbf{x}_{\mu}^T \mathbf{S}^* \mathbf{x}_{\mu} = \sum_{\gamma=1}^r [(\mathbf{u}_{\gamma}^*)^T \cdot \mathbf{x}_{\mu}]^2$$

↑ $\sim \mathcal{N}(0, \mathbf{I}_d)$
↑ $\mathbf{u}_{\gamma}^* \sim \mathcal{N}(0, \mathbf{I}_d)$

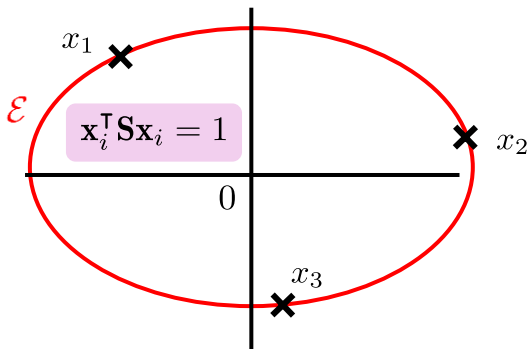
Proof ideas (1)

$$y_\mu = \varphi(\text{Tr}[\mathbf{X}_\mu \mathbf{S}^*]) \quad \mathbf{S}^* = \mathbf{U}^*(\mathbf{U}^*)^\top \sim \text{Wishart matrix}$$

Step 1: “Gaussian universality”

Universality of Bayes-optimal generalization error

Developed for the analysis of ellipsoid fitting [M. & Bandeira '23]



The “generic” assumption on $\mathbb{P}(\mathbf{X}_\mu)$

Gaussian matrix



Uniform Central Limit Theorem for 1-dimensional projections $\text{Tr}[\mathbf{X}\mathbf{S}]$

$$\sup_{\mathbf{S}} |\mathbb{E}\phi(\text{Tr}[\mathbf{X}\mathbf{S}]) - \mathbb{E}\phi(\text{Tr}[\mathbf{G}\mathbf{S}])| \rightarrow 0$$

Satisfied e.g. by $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$, $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$ (extensive-width sign retrieval)



$$\begin{aligned} \min \mathcal{E}_{\text{gen.}}(\hat{\mathbf{w}}_k) = \min \|\mathbf{y}(\hat{\mathbf{S}}) - \mathbf{y}(\mathbf{S}^*)\|^2 &= \min \tilde{\mathcal{E}}_{\text{gen.}}(\hat{\mathbf{S}}) = \min \|\tilde{\mathbf{y}}(\hat{\mathbf{S}}) - \tilde{\mathbf{y}}(\mathbf{S}^*)\|^2 \times (1 + o(1)) \\ \text{from } \{y_\mu = \varphi(\text{Tr}[\mathbf{S}^* \mathbf{X}_\mu])\}_{\mu=1}^n &\quad \text{from } \{\tilde{y}_\mu = \varphi(\text{Tr}[\mathbf{S}^* \mathbf{G}_\mu])\}_{\mu=1}^n \end{aligned}$$

Gaussian matrix



To compute $\mathcal{E}_{\text{gen.}}$ we pretend the data is a Gaussian matrix

Proof ideas (2)

A (generalized) **linear model on S^*** , with...

Step 2 : $\{\tilde{y}_\mu = \varphi(\text{Tr}[S^* G_\mu])\}_{\mu=1}^n$

➤ Gaussian data $G := \begin{pmatrix} \text{flatt}(G_1) \\ \vdots \\ \text{flatt}(G_n) \end{pmatrix}$

+

➤ Wishart prior
 $S^* = U^*(U^*)^\top \sim \mathcal{W}_{r,d}$

[Barbier & al '19]:
i.i.d. priors



Generalizing [Barbier
& al '19]

An explicit “replica-symmetric” formula for $\tilde{\mathcal{E}}_{\text{gen.}}$

Involves ...

Scalar estimation problem involving φ



Step 3 :

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

Guilhem's talk

Gaussian (GOE) matrix

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

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

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

Combining all steps...



$\lim_{d \rightarrow \infty} \mathcal{E}_{\text{gen.}} = \dots$



Result II : The GAMP-RIE algorithm

$$y_\mu \sim \varphi(\cdot | \text{Tr}[\mathbf{S}^* \mathbf{X}_\mu])$$

Example

$$\mathbf{X}_\mu = \mathbf{x}_\mu \mathbf{x}_\mu^\top$$

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

Informal hypothesis

Universality $\mathbf{X}_\mu \Rightarrow \mathbf{G}_\mu$ also holds “at the level of algorithms”

MSE-optimal algorithm

Generalized linear model
w. Gaussian data



Generalized Approximate Message Passing (GAMP)

[Donoho&al '09 ;
Rangan '11 ; ...]

Guilhem's talk

Each GAMP iteration solves

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



Rotationally-Invariant Estimator (RIE)

$$\mathbf{Y} = \mathbf{O} \mathbf{D} \mathbf{O}^\top \Leftrightarrow \hat{\mathbf{S}}(\mathbf{Y}) = \mathbf{O} f_{\text{opt.}}(\mathbf{D}) \mathbf{O}^\top$$

[Bun & al '16 ; ...]

Known “optimal shrinkage” function

An explicit easy-to-implement polynomial-time algorithm

GAMP

RIE

+

Similar algorithm in non-symmetric
model (**Bilinear Sequence Regression**)

Non-symmetric denoising in [Troiani, Erba,
Krzakala, M. & Zdeborová '22]

Objectives

Objectives

Bayes-optimal generalization error

$$\mathcal{E}_{\text{gen.}} := \mathbb{E}_{\mathbf{U}^*, \mathbf{V}^*, \mathcal{D}} \min_{\hat{y}(\mathcal{D})} \mathbb{E}_{\mathbf{X}_{\text{test}}} [(\hat{y}(\mathbf{X}_{\text{test}}) - f_{\mathbf{U}^*, \mathbf{V}^*}(\mathbf{X}_{\text{test}}))^2]$$

- ❑ **Sharp thresholds** ? Phase transitions ?
- ❑ Efficiently achievable ? **Hard phases** ?

Comparison to **minimal nuclear-norm** estimator

$$\min\{\|\mathbf{S}\|_{\text{NN}} : \mathbf{S} = \mathbf{U}\mathbf{V}^T, y_\mu = \text{Tr}[\mathbf{X}_\mu \mathbf{S}]\}$$

Recht&al '10

Comparison to **linear regression** on **vectorized data** ?

- ❑ Performance of **GD-based algorithms** ?
Bhojanapalli&al '16

$$\mathcal{L}(\mathbf{U}, \mathbf{V}) = \sum_{\mu=1}^n (y_\mu - \text{Tr}[\mathbf{X}_\mu \mathbf{U}\mathbf{V}^T])^2$$

- ❑ Is GD an **implicit nuclear-norm minimizer** ?
Gunasekar&al '17

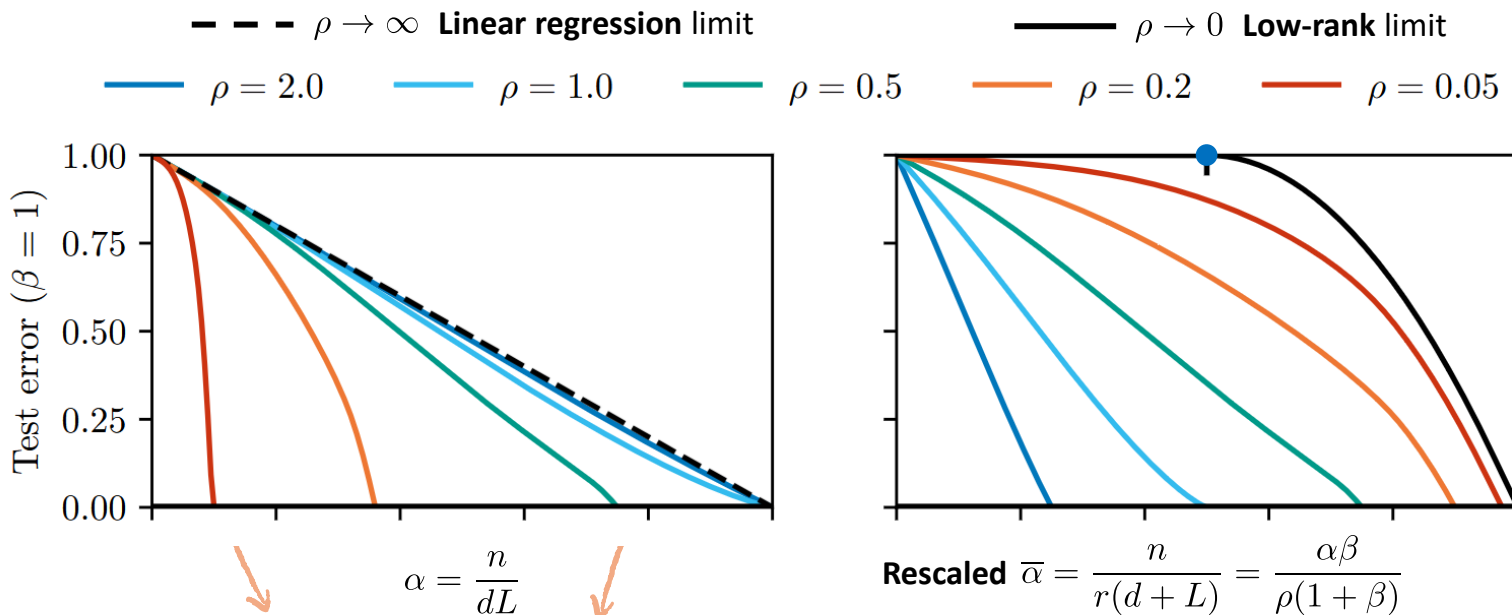


**Bilinear Sequence Regression (BSR) /
Extensive-rank matrix sensing**

Bayes-optimal error

$$\left\{ y_\mu = \sum_{a,i=1}^{L,d} X_{ai}^\mu \sum_{\gamma=1}^r U_{i\gamma}^* V_{\gamma a}^* \right\}_{\mu=1}^n$$

$$\rho = \frac{r}{d} \quad \beta = \frac{L}{d} \quad \alpha = \frac{n}{dL}$$



Perfect recovery threshold

$$\alpha_{\text{PR}} = \begin{cases} \frac{\rho}{\beta}(1+\beta-\rho) & 0 < \rho < 1, \\ 1 & \rho \geq 1 \end{cases}$$

$$\bar{\alpha}_{\text{WR}} \stackrel{\rho \rightarrow 0}{=} \frac{\sqrt{\beta}}{1+\beta}$$

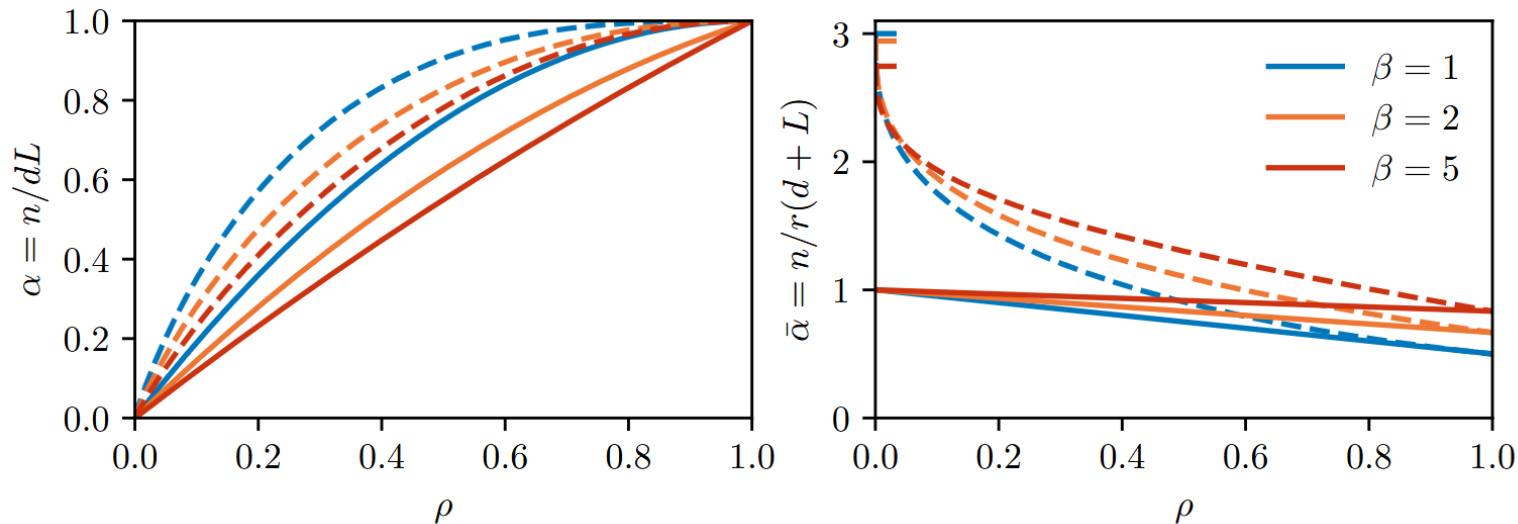
Matches a naïve “counting argument” $\text{DOF}[S = UV^T] \simeq \alpha_{\text{PR}} \cdot d^2$

- **Weak recovery threshold** in the low-rank limit.
- Matches the $r = \mathcal{O}_d(1)$, then $r \rightarrow \infty$ limiting curve.

Minimal nuclear-norm estimator

$$\left\{ y_\mu = \sum_{a,i=1}^{L,d} X_{ai}^\mu \sum_{\gamma=1}^r U_{i\gamma}^* V_{\gamma a}^* \right\}_{\mu=1}^n$$

$$\rho = \frac{r}{d} \quad \beta = \frac{L}{d} \quad \alpha = \frac{n}{dL}$$



— Bayes-optimal **perfect recovery** α_{PR}

- - - Perfect recovery of the **minimal nuclear norm estimator** $\min\{\|\mathbf{S}\|_{\text{NN}} : \mathbf{S} = \mathbf{UV}^T, y_\mu = \text{Tr}[\mathbf{X}_\mu \mathbf{S}]\}$

[Donoho, Gavish & Montanari '13]

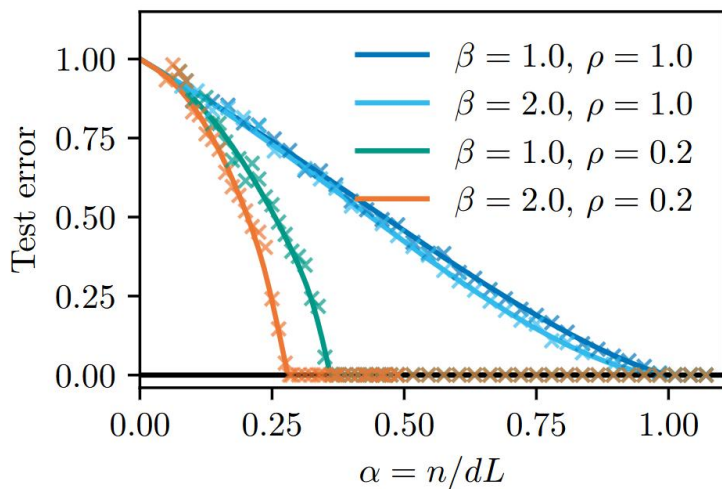


Suboptimality of the min-nuclear
norm estimator

Absence of hard phase

$$\left\{ y_\mu = \sum_{a,i=1}^{L,d} X_{ai}^\mu \sum_{\gamma=1}^r U_{i\gamma}^* V_{\gamma a}^* \right\}_{\mu=1}^n \quad \rho = \frac{r}{d} \quad \beta = \frac{L}{d} \quad \alpha = \frac{n}{dL}$$

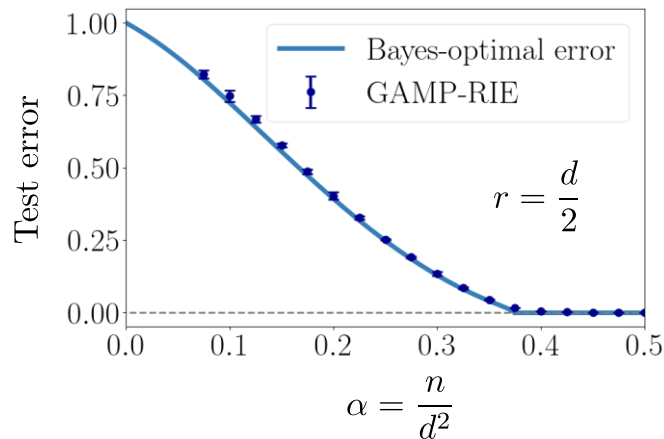
Bilinear Sequence Regression



× GAMP-RIE Algorithm

— Bayes-optimal error

Extensive-width sign retrieval



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

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



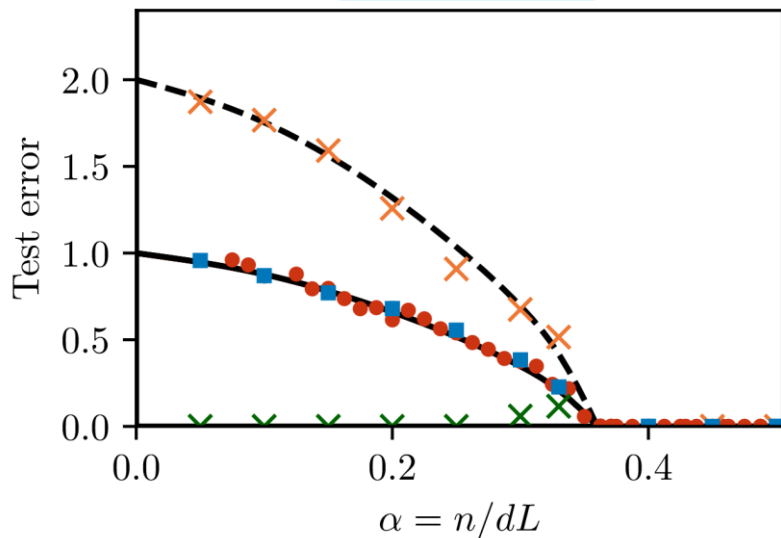
No computational-to-statistical gap /
hard phase in extensive-width regime

Empirical performance of gradient descent

$$\rho = \frac{r}{d} \quad \beta = \frac{L}{d} \quad \alpha = \frac{n}{dL}$$

$$\mathcal{L}(\mathbf{U}, \mathbf{V}) = \sum_{\mu=1}^n (y_{\mu} - \text{Tr}[\mathbf{X}_{\mu} \mathbf{U} \mathbf{V}^{\top}])^2$$

$$\beta = 1 \quad \rho = 0.2$$



- GAMP-RIE
- Error of the Bayes-optimal (BO) estimator $\mathbb{E}[\mathbf{U}, \mathbf{V}|\mathcal{D}]$
- - - $2 \times \text{BO} = \text{Error of posterior sampler } \sim \mathbb{P}(\mathbf{U}, \mathbf{V}|\mathcal{D})$
- * Gradient descent
- * Gradient descent (final training loss value)
- **Averaged** Gradient descent
Averaged over many initializations

- $d = 100$ for GAMP/GD runs
- Random initialization
- Cross-validated learning rate



For **any** ρ , AGD seems to reach the Bayes-optimal MMSE

Despite **non-convexity** of the problem!



- Similar results in **extensive-width phase retrieval**.
- No longer true when adding noise to y_{μ} .

Does GD do implicit nuclear norm regularization ?

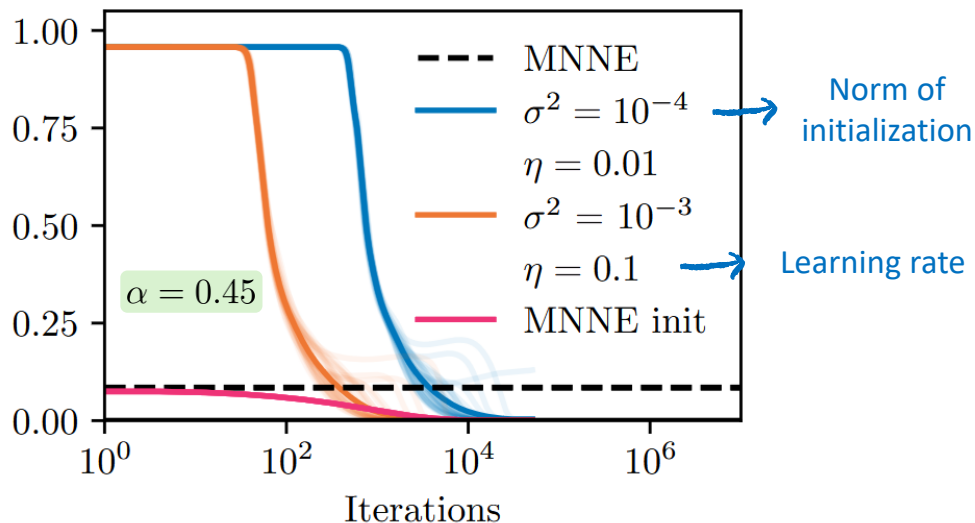
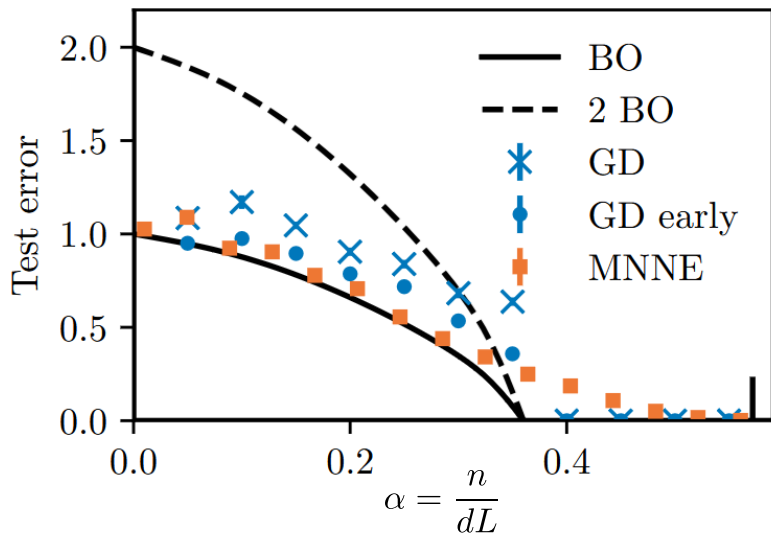
$$\left\{ y_\mu = \sum_{a,i=1}^{L,d} X_{ai}^\mu \sum_{\gamma=1}^r U_{i\gamma}^* V_{\gamma a}^* \right\}_{\mu=1}^n$$

Suggestion: GD with small initialization has an implicit bias towards the minimal nuclear norm estimator (MNNE)

Gunasekar, Woodworth, Bhojanapalli, Neyshabur & Srebro '17

$$\min\{\|\mathbf{S}\|_{\text{NN}} : \mathbf{S} = \mathbf{UV}^\top, y_\mu = \text{Tr}[\mathbf{X}_\mu \mathbf{S}]\}$$

$$\beta = 1 \quad \rho = 0.2$$



GD does not reach the MNNE

What about
overparametrized settings

$$(r^* < r)$$



Summary

Bilinear Sequence Regression (BSR)

Most basic model for learning from **long** sequences of **high-dimensional** tokens

$$y_\mu = \text{Tr}[\mathbf{X}_\mu \mathbf{U}^* (\mathbf{V}^*)^\top] = \sum_{a,i=1}^{L,d} X_{ai}^\mu \sum_{\gamma=1}^r U_{i\gamma}^* V_{\gamma a}^*$$

$$L = \Theta(d); r = \Theta(d) \quad n = \Theta[r(L + d)]$$

THANK YOU !

1. Analytical formula for the **Bayes-optimal generalization error**.
2. Optimal algorithm (GAMP-RIE), **no computational-statistical gap**.
3. Gap between BO error and **linear regression and MNNE**
4. (Averaged) **Gradient descent seems to sample from the posterior in the noiseless setting**, despite non-convexity !



- ❖ Theoretical analysis of GD properties / implicit regularization
- ❖ For extensive-width 2-layer NNs: beyond **quadratic activation** ? Cf Jean's talk on Tuesday
- ❖ Overparametrization ($r^* < r$) ?
- ❖ Correlations between tokens ?