

FUNDAMENTAL LIMITS OF HIGH-DIMENSIONAL ESTIMATION

A STROLL BETWEEN STATISTICAL PHYSICS, PROBABILITY, AND RANDOM
MATRIX THEORY

Antoine Maillard

Under the supervision of Florent Krzakala



Département
de Physique
—
École normale
supérieure



PSL 



PhD defense – August 30th 2021

WHAT IS STATISTICAL INFERENCE ?

Input data

Φ_μ

Model

X^*



Output data

y_μ

$$\Phi = \{\Phi_1, \dots, \Phi_m\} \quad Y = \{y_1, \dots, y_m\}$$

“Signal”

X^*



- Supervised learning in “teacher-student” neural networks
- Signal processing
- Phase retrieval
- Matrix factorization
- Quantitative finance, particle physics, evolutionary biology...

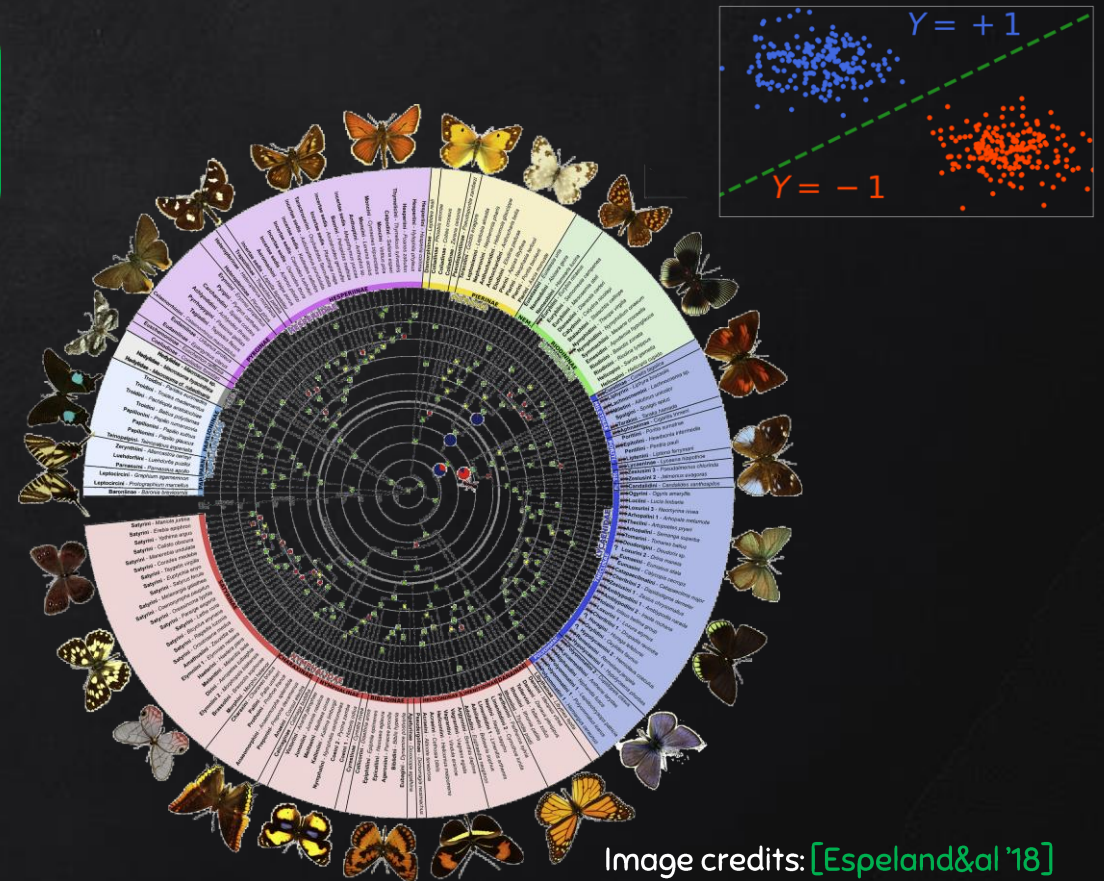
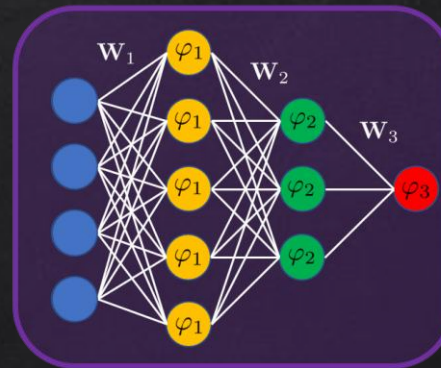
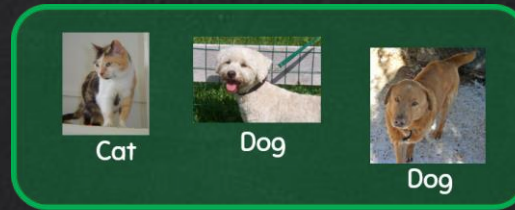


Image credits: [Espeland&al '18]

STATISTICS IN HIGH DIMENSION

$$\left\{ \Phi_\mu \in \mathbb{R}^p \xrightarrow{\text{Model}} \begin{matrix} \text{Model} \\ \mathbf{X}^* \in \mathbb{R}^n \end{matrix} \xrightarrow{\text{gears}} y_\mu \right\}_{\mu=1, \dots, m}$$

Data deluge

Theoretical revolution of the 2000s

Gigantic databases and
explosion of computing power.



High-dimensional statistics

[Donoho, AMS Lectures 2000:
*High-Dimensional Data
Analysis: The Curses and
Blessings of Dimensionality*]

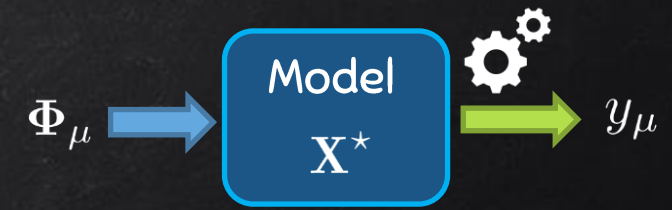
“Modern machine learning”: GoogLeNet [Szegedy&al '15]: $n \simeq 5 \times 10^6$ and $m \simeq 10^6$.

“High-dimensional” limit

Number of parameters $n \rightarrow \infty$
+
Number of data $m \rightarrow \infty$

In this presentation: $m/n \rightarrow \alpha > 0$ (sampling ratio).

BAYESIAN FORMALISM



Posterior distribution

Observation channel,
or likelihood

Prior distribution

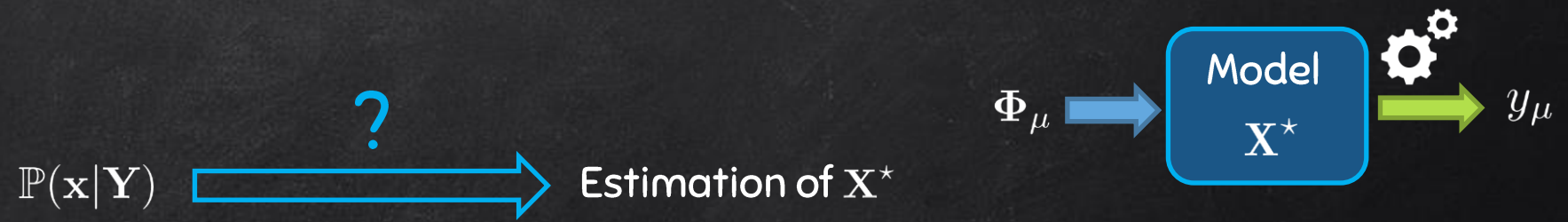
$$\mathbb{P}(X^* = \mathbf{x} | Y, \Phi) = \frac{\mathbb{P}(Y | X^* = \mathbf{x}, \Phi) \times \mathbb{P}(X^* = \mathbf{x})}{\mathbb{P}(Y | \Phi)}$$



Most of this talk: the prior and the observation channel are known to the statistician.

Bayes-optimal setting

ESTIMATORS



Bayesian estimators:

➤ *Maximum A Posteriori* $\hat{\mathbf{X}}_{\text{MAP}} \equiv \arg \max_{\mathbf{x}} \mathbb{P}(\mathbf{x}|\mathbf{Y})$

$$\hat{\mathbf{X}}_{\text{MMSE}} = \int d\mathbf{x} \mathbb{P}(\mathbf{x}|\mathbf{Y}) \mathbf{x} = \langle \mathbf{x} \rangle_{\mathbf{Y}}$$

➤ *Minimal Mean Squared Error* $\hat{\mathbf{X}}_{\text{MMSE}} \equiv \arg \min_{\mathbf{x}} \left\{ \mathbb{E}_{\mathbf{Y}} \int d\mathbf{x}' \mathbb{P}(\mathbf{x}'|\mathbf{Y}) \|\mathbf{x} - \mathbf{x}'\|^2 \right\}$

Many other types of estimators exist, such as the *Empirical Risk Minimizer* $\hat{\mathbf{X}} \equiv \arg \min_{\mathbf{x}} \sum_{\mu=1}^m L(x_\mu, y_\mu)$

“Loss” function

This presentation: we mainly focus on **MMSE estimation** and **empirical risk minimization**.

IMPORTANT EXAMPLE – GENERALIZED LINEAR MODELS

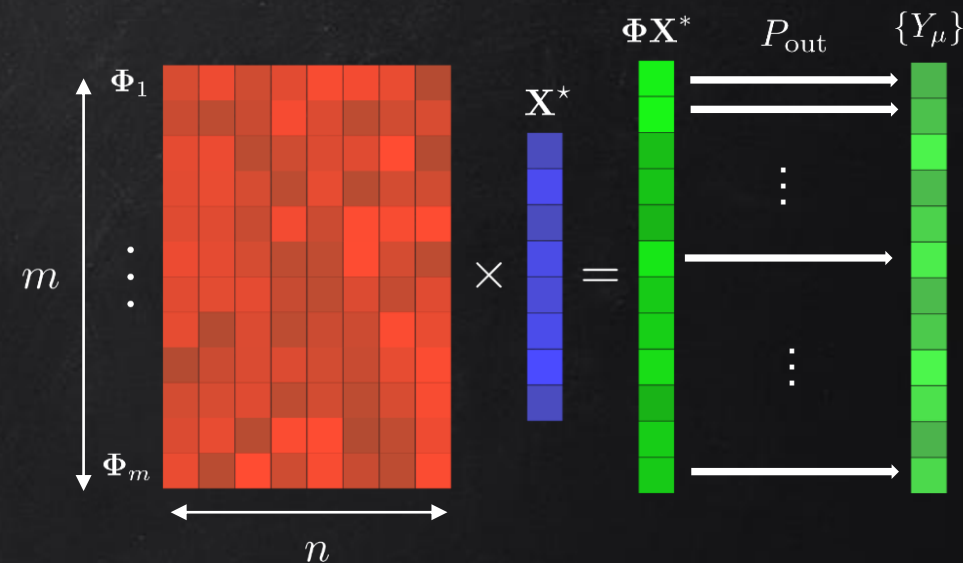
Goal: Recover $\mathbf{X}^* \in \mathbb{R}^n$ from $\{\Phi_\mu, Y_\mu\}_{\mu=1}^m$:

Observations $Y_\mu \in \mathbb{R}$

$$Y_\mu \sim P_{\text{out}}\left(\cdot \mid \frac{1}{\sqrt{n}} \sum_{i=1}^n \Phi_{\mu i} X_i^*\right) \quad \mu \in \{1, \dots, m\}$$

Probabilistic
channel (noise)

Sensing matrix



Many examples: compressed sensing, perceptron learning, phase retrieval, ...

$$\mathbb{P}(\mathbf{x} | \mathbf{Y}, \Phi) = \frac{\mathbb{P}(\mathbf{x}) \mathbb{P}(\mathbf{Y} | \mathbf{x}, \Phi)}{\mathbb{P}(\mathbf{Y} | \Phi)} = \frac{1}{\mathcal{Z}_n(\mathbf{Y}, \Phi)} \prod_{i=1}^n P_X(x_i) \prod_{\mu=1}^m P_{\text{out}}\left[Y_\mu \mid \frac{1}{\sqrt{n}} (\Phi \mathbf{x})_\mu\right].$$

Prior knowledge on the signal

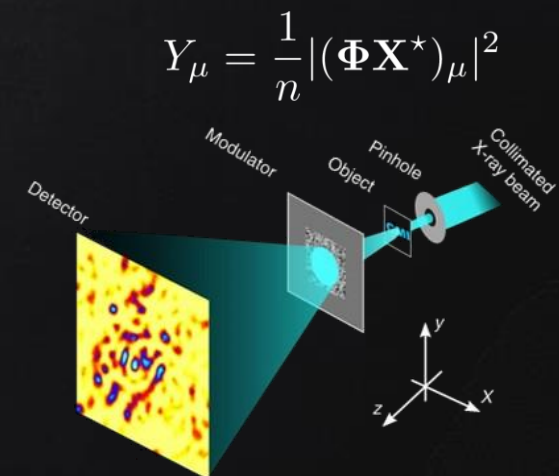


Image credits: [Zhang&al16]

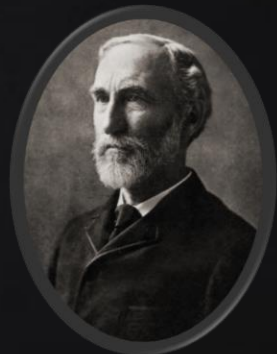
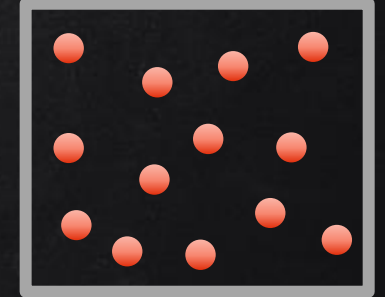
WHERE ARE THE PHYSICS?

$$\mathcal{H} = \frac{m}{2} \sum_{i=1}^n v_i^2$$

“Statistical mechanics 101”

Consider n particles with position x_i , with distribution $P_X(x)$, interacting via the **Hamiltonian** $\mathcal{H}(\mathbf{x})$, at temperature $T = \eta^{-1}$.

Gibbs-Boltzmann probability: $\mathbb{P}_\eta(\mathbf{x}) = \frac{1}{\mathcal{Z}_n(\eta)} e^{-\eta \mathcal{H}(\mathbf{x})} \prod_{i=1}^n P_X(x_i)$



$$\text{GLM: } \mathbb{P}(\mathbf{x}|\mathbf{Y}, \Phi) = \frac{1}{\mathcal{Z}_n(\mathbf{Y}, \Phi)} \prod_{i=1}^n P_X(x_i) \prod_{\mu=1}^m P_{\text{out}} \left[Y_\mu \middle| \frac{1}{\sqrt{n}} (\Phi \mathbf{x})_\mu \right].$$



Statistical physics “disordered” model, with Hamiltonian $\mathcal{H}(\mathbf{x}) = - \sum_{\mu=1}^m \ln P_{\text{out}} \left[Y_\mu \middle| \frac{1}{\sqrt{n}} (\Phi \mathbf{x})_\mu \right] \quad (T = 1)$



Spin glasses

General connection for many statistical models

[Hopfield '82; Mézard&Parisi '85; Gardner&Derrida '89; Anderson '89; Mézard&Montanari '09; ...]

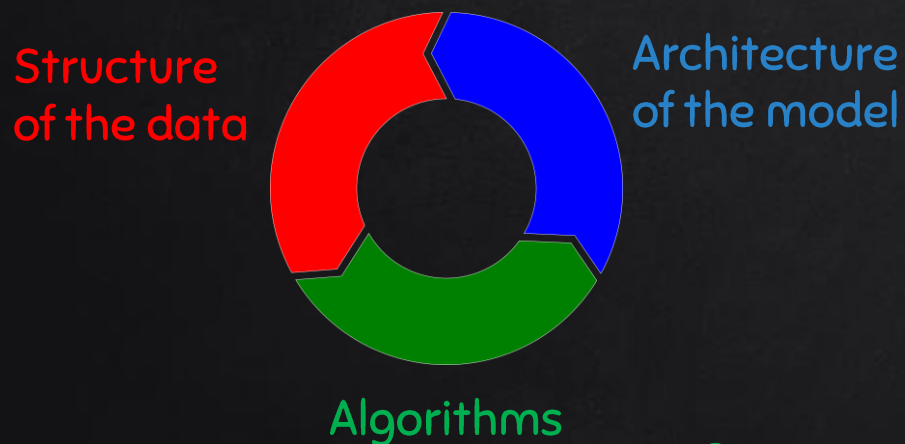


WHERE DO WE GO FROM HERE?

Deep and detailed connection

- | | | |
|---|---|--|
| ➤ Bayesian estimation problems | ↔ | ➤ Finite-temperature statistical physics |
| ➤ Empirical risk minimization | ↔ | ➤ (Zero-temperature) energy landscape minimization |
| ➤ Posterior distribution | ↔ | ➤ Gibbs-Boltzmann distribution |
| ➤ High-dimensional limit | ↔ | ➤ Thermodynamic limit |
| ➤ Randomness of the observations (noise, ...) | ↔ | ➤ Disordered systems, “spin glasses” |

Theory of machine learning / inference



Our “statistical physics-inspired” approach allows to study each of these pieces!

MAIN PHD PROJECTS

❖ Revisiting high-temperature expansions

- High-temperature expansions and message passing algorithms. *J.Stat.Mech.* 2019.
- Towards exact solution of extensive-rank matrix factorization. *In preparation.*

Approximation schemes and algorithms



II

✓ Exacts in high dimension

High-temperature expansions + Diagrammatics and random matrix theory

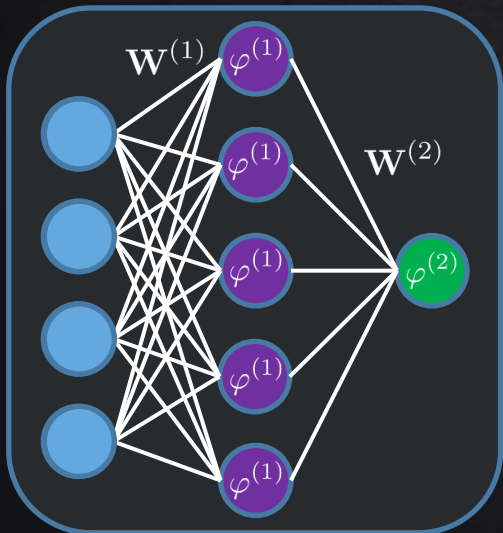
→ I + II



Extensive-rank matrix factorization $\mathbf{Y} = \mathbf{U}\mathbf{V}^T + \mathbf{Z}$

MAIN PHD PROJECTS

- ❖ Optimal estimation in high-dimensional problems
 - The mutual information in random linear estimation beyond iid matrices. *ISIT 2018*.
 - Computational-to-statistical gaps in learning a two-layers neural network. *NeurIPS 2018 & J.Stat.Mech. 2019*.
 - The spiked matrix model with generative priors. *IEEE Trans. Inf. Theory 2020 & NeurIPS 2019*
 - Phase retrieval in high dimensions: statistical and computational phase transitions. *NeurIPS 2020*.
 - Construction of optimal spectral methods in phase retrieval. *MSML 2021*.



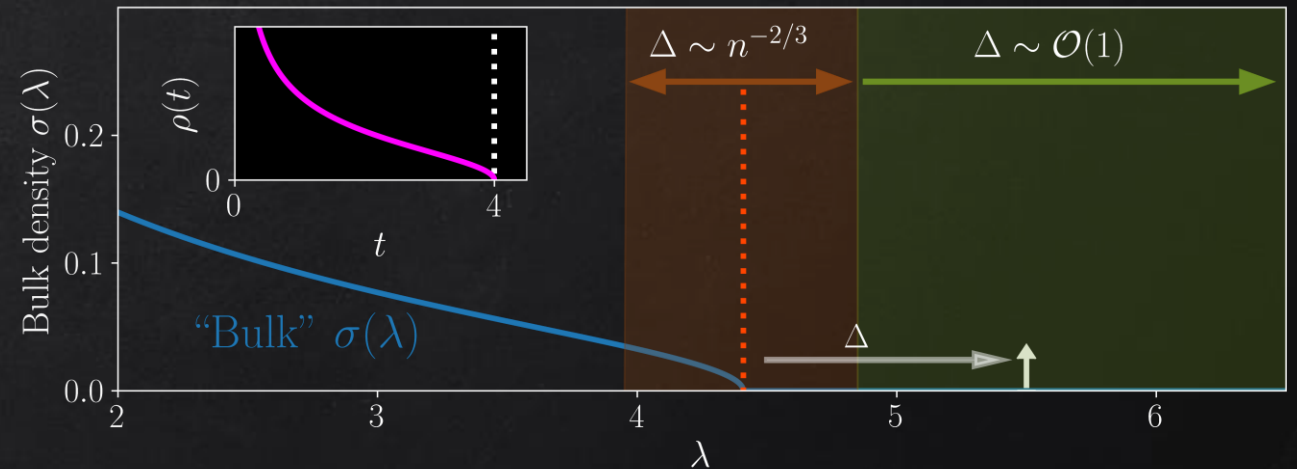
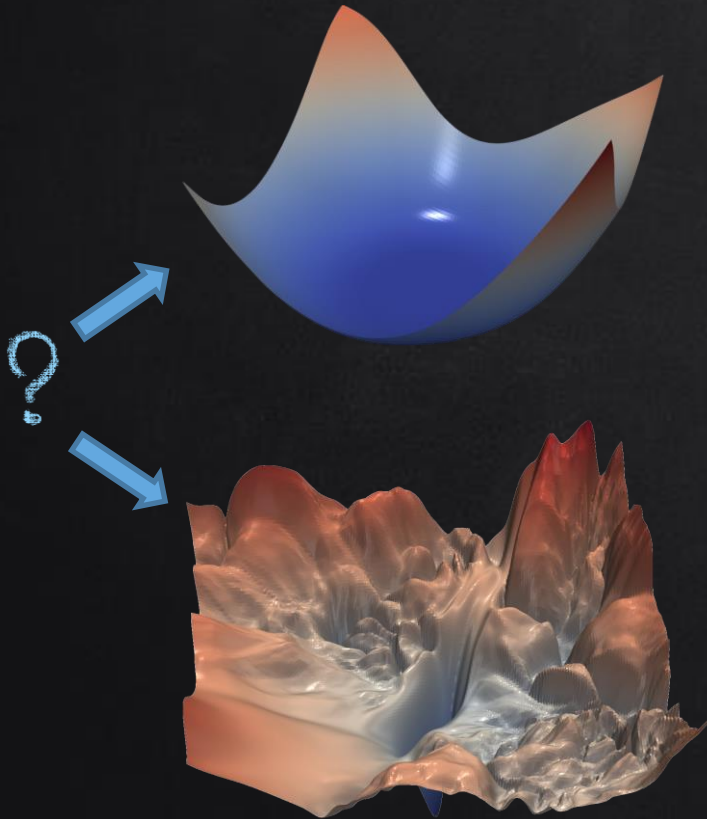
Committee machine

$$Y_\mu = \frac{1}{n} |(\Phi \mathbf{X}^*)_\mu|^2$$



MAIN PHD PROJECTS

- ❖ Towards a topological approach to high-dimensional optimization
 - Landscape complexity for the empirical risk of generalized linear models. *MSML 2020*.
 - Large deviations of extreme eigenvalues of generalized sample covariance matrices. *EPL 2021*.



$$\mathbf{M} = \frac{1}{m} \sum_{\mu=1}^m \rho_{\mu} \mathbf{z}_{\mu} \mathbf{z}_{\mu}^{\dagger}$$

MAIN PHD PROJECTS

❖ Revisiting high-temperature expansions

- High-temperature expansions and message passing algorithms. *J.Stat.Mech.* 2019.
- Towards exact solution of extensive-rank matrix factorization. *In preparation.*

❖ Optimal estimation in high-dimensional problems

- The mutual information in random linear estimation beyond iid matrices. *ISIT 2018.*
- Computational-to-statistical gaps in learning a two-layers neural network. *NeurIPS 2018 & J.Stat.Mech.* 2019.
- I { • The spiked matrix model with generative priors. *IEEE Trans. Inf. Theory* 2020 & *NeurIPS* 2019
- Phase retrieval in high dimensions: statistical and computational phase transitions. *NeurIPS* 2020.
- Construction of optimal spectral methods in phase retrieval. *MSML* 2021.

❖ Towards a topological approach to high-dimensional optimization

- II { • Landscape complexity for the empirical risk of generalized linear models. *MSML* 2020.
- Large deviations of extreme eigenvalues of generalized sample covariance matrices. *EPL* 2021.

I

EXPLOITING DATA STRUCTURE IN SPIKED MATRIX ESTIMATION

Spiked Wigner model [Johnstone '01]

$$\mathbf{Y} = \frac{1}{\sqrt{p}} \mathbf{v}^* (\mathbf{v}^*)^\top + \sqrt{\Delta} \boldsymbol{\xi} \in \mathbb{R}^{p \times p}$$

$$\mathbf{v}^* \sim P_v$$

$$\rho_v \equiv \lim_{p \rightarrow \infty} \frac{1}{p} \mathbb{E}_{P_v} \|\mathbf{v}\|^2$$

$$\begin{cases} \xi_{ij} = \xi_{ji} \\ \xi_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1 + \delta_{ij}) \end{cases}$$

Spiked Wishart model

$$\mathbf{Y} = \frac{1}{\sqrt{p}} \mathbf{u}^* (\mathbf{v}^*)^\top + \sqrt{\Delta} \boldsymbol{\xi} \in \mathbb{R}^{n \times p}$$

$$\mathbf{u}^* \sim P_u$$

$$\mathbf{v}^* \sim P_v$$

$$\xi_{\mu i} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

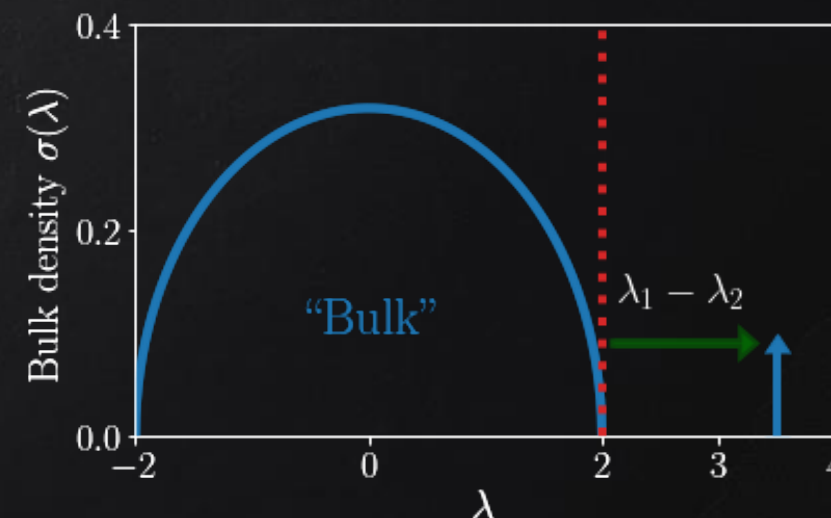
❖ PCA: the dominant eigenvector of \mathbf{Y} .

Optimal for unstructured signal $P_v = \mathcal{N}(0, 1)$.

❖ Leverage prior knowledge on the structure of the signal to improve recovery? \rightarrow “Dimensionality reduction”



“BBP” transition



$$\Delta / \rho_v^2 \approx 1$$

[Edwards&Jones '76; Baik, Ben Arous&Péché '04]

DIMENSIONALITY REDUCTION: SYNTHETIC MODELS

$$\mathbf{Y} = \frac{1}{\sqrt{p}} \mathbf{v}^* (\mathbf{v}^*)^\top + \sqrt{\Delta} \boldsymbol{\xi} \in \mathbb{R}^{p \times p}$$

Sparsity

➤ Natural representation, e.g.:

Images \Rightarrow Wavelet
Sound \Rightarrow Fourier

➤ Efficient algorithms: LASSO, compressed sensing,...

But...

- ❖ Large algorithmically hard phases
- ❖ Impossible to “beat” the BBP transition of PCA.

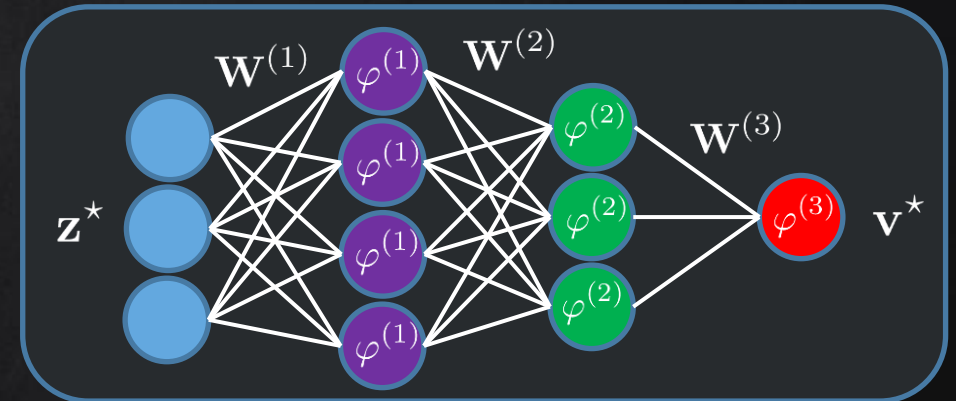
Generative prior

Random weights

Unstructured latent variable

$$\mathbf{v}^* = \varphi^{(L)} \left(\frac{1}{\sqrt{k_L}} \mathbf{W}^{(L)} \dots \varphi^{(1)} \left(\frac{1}{\sqrt{k}} \mathbf{W}^{(1)} \mathbf{z}^* \right) \right)$$

Structured signal



In the limit $p \rightarrow \infty$, for a given $\Delta > 0$, what is the optimal recovery:

- **Information-theoretically** (in exponential time) ?
- With which tractable (**polynomial-time**) **algorithm** ?
- **Cheap** (e.g. **spectral**) **algorithms** that outperform PCA ?

THE REPLICA-SYMMETRIC FORMULA

$$\mathbb{P}(\mathbf{v}|\mathbf{Y}) = \frac{1}{\mathcal{Z}_p(\mathbf{Y})} P_v(\mathbf{v}) \prod_{i < j} e^{-\frac{1}{2\Delta} \left(Y_{ij} - \frac{v_i v_j}{\sqrt{p}} \right)^2}$$

Theorem (informal)

$$\xi \sim \mathcal{N}(0, \mathbf{I}_p)$$

$$\lim_{p \rightarrow \infty} -\frac{1}{p} \mathbb{E}_{\mathbf{Y}} \ln \mathcal{Z}_p(\mathbf{Y}) = \inf_{q \in (0, \rho_v)} f_{\text{RS}}(\Delta, q), \text{ with } f_{\text{RS}}(\Delta, q) = \frac{q(\rho_v - q)}{2\Delta} + \lim_{p \rightarrow \infty} \frac{1}{p} \mathbb{I}(\mathbf{v}; \mathbf{v} + \sqrt{\Delta/q_v} \xi)$$

Moreover:

$$\text{MMSE}_v(\Delta) = \rho_v - \arg \min_{q \in [0, \rho_v]} f_{\text{RS}}(\Delta, q)$$

Mutual information

Information-theoretic MMSE.

How to:

- Derive the result using the non-rigorous replica method [Mézard, Parisi & Virasoro '87]...
- Prove the result (not the method!) using interpolation techniques [Guerra '03 ; Talagrand '06 ; Barbier&al '19]...

Similar results & strategy for: two-layers neural networks [Aubin, A.M.&al '18], compressed sensing with non-i.i.d. matrices [Barbier, A.M.&al'18], phase retrieval with rotationally-invariant matrices [A.M.&al '20], ...

APPLICATION: SINGLE-LAYER GENERATIVE PRIOR

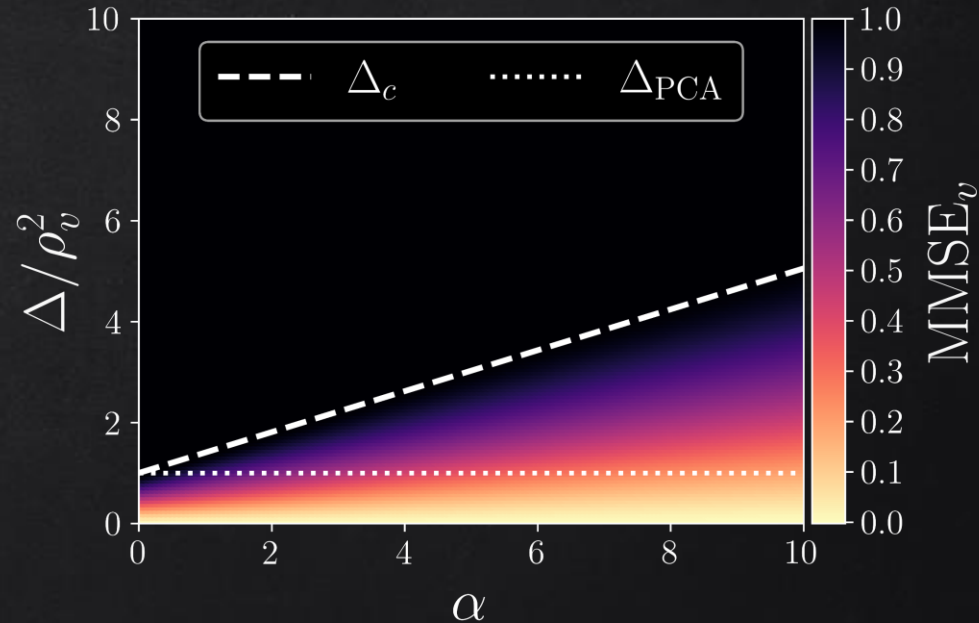
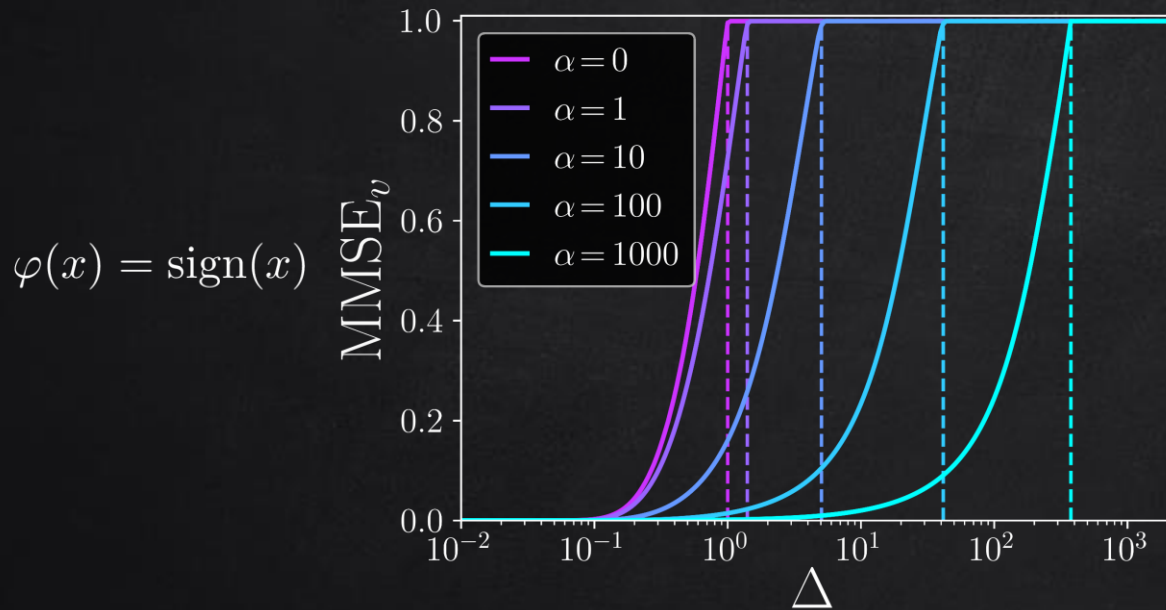
$$\mathbf{v}^* \sim \varphi\left(\frac{1}{\sqrt{k}}\mathbf{W}\mathbf{z}^*\right)$$

$$\begin{cases} W_{il} & \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1) \\ \mathbf{z}^* & \sim P_z \\ \alpha & \equiv p/k \end{cases} \quad \leftarrow \text{Unstructured (i.i.d.) prior}$$



$$f_{\text{RS}}(\Delta, q) = \frac{q(\rho_v - q)}{2\Delta} + \frac{1}{\alpha} \min_{q_z, \hat{q}_z} \left[\frac{q_z \hat{q}_z}{2} - \Psi_z(\hat{q}_z) - \alpha \Psi_{\text{out}}\left(\frac{q}{\Delta}, q_z\right) \right]$$

Tedious, but completely scalar potential!



Weak-recovery: $\Delta_c \equiv \inf\{\Delta : \text{MMSE}_v(\Delta) = 1\}$

❖ Sparse PCA: $\Delta_c = 1$

❖ 1-layer generative prior: $\Delta_c = 1 + \frac{4}{\pi^2}\alpha$

ALGORITHMIC LIMITS

Can we algorithmically (i.e. in polynomial time) achieve the optimal MSE?

Measure of algorithm reconstruction by the **overlaps**

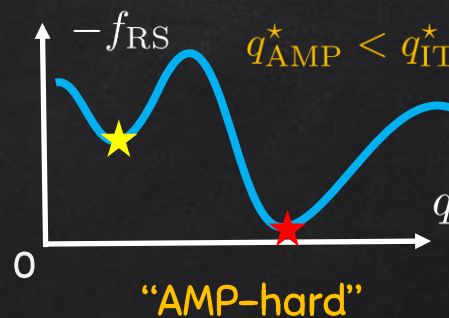
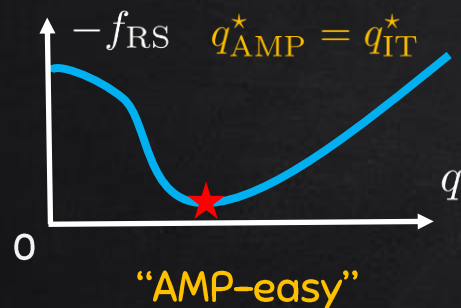
$$q \equiv \lim_{p \rightarrow \infty} \frac{1}{p} \mathbb{E}[\mathbf{v}^\top \mathbf{v}^*]$$

$$q_z \equiv \lim_{k \rightarrow \infty} \frac{1}{k} \mathbb{E}[\mathbf{z}^\top \mathbf{z}^*]$$



$$q^{t+1} = 2\partial_q \Psi_{\text{out}}\left(\frac{q}{\Delta}, q_z\right); \quad q_z^{t+1} = 2\partial_{\hat{q}_z} \Psi_z(\hat{q}_z^t); \quad \hat{q}_z^{t+1} = 2\alpha\partial_{q_z} \Psi_{\text{out}}\left(\frac{q^t}{\Delta}, q_z^t\right)$$

State Evolution (SE) equations: “Fixed point algorithm” on f_{RS} !



Tested settings: single/multi layer and $\varphi \in \{\text{linear, sign, ReLU}\}$.

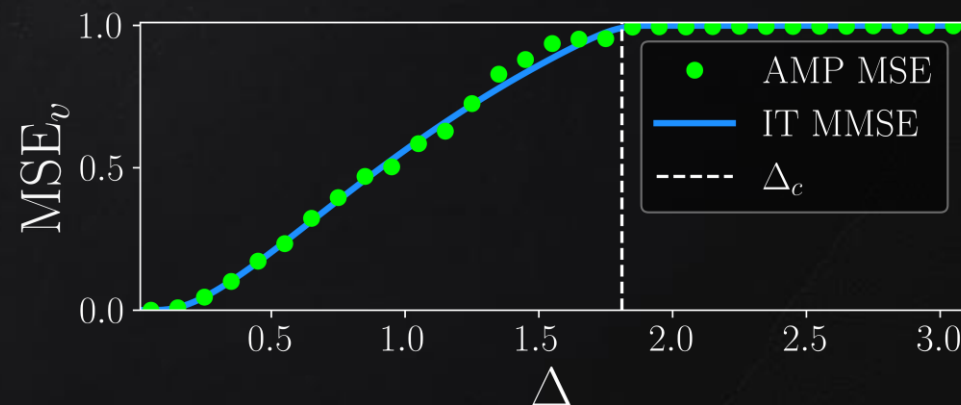


No algorithmically hard phase: very different from sparse PCA!

Approximate Message Passing (AMP)

- 1: **Input:** $Y \in \mathbb{R}^{p \times p}$ and $W \in \mathbb{R}^{p \times k}$;
- 2: **Initialize with:** $\hat{\mathbf{v}}^{t=1} = \mathcal{N}(\mathbf{0}, \sigma^2 I_p)$, $\hat{\mathbf{z}}^{t=1} = \mathcal{N}(\mathbf{0}, \sigma^2 I_k)$, and $\hat{\mathbf{c}}_v^{t=1} = I_p$, $\hat{\mathbf{c}}_z^{t=1} = I_k$, $t = 1$.
- 3: **repeat**
- 4: **Spiked layer denoising:**
- 5: $\mathbf{B}_v^t = \frac{1}{\Delta\sqrt{p}}\hat{\mathbf{v}}^t - \frac{1}{\Delta}\left(\frac{I_p^\top \hat{\mathbf{c}}_v^t}{p}\right)\hat{\mathbf{v}}^{t-1}$ and $A_v^t = \frac{1}{\Delta p}(\|\hat{\mathbf{v}}^t\|_2)^2 I_p$.
- 6: **Generative layer denoising:**
- 7: $V^t = \frac{1}{k}(I_k^\top \hat{\mathbf{c}}_z^t)I_p$, $\omega^t = \frac{1}{\sqrt{k}}W\hat{\mathbf{z}}^t - V^t\mathbf{g}^{t-1}$
- 8: $\mathbf{g}^t = f_{\text{out}}(\mathbf{B}_v^t, A_v^t, \omega^t, V^t)$
- 9: $\Lambda^t = \frac{1}{k}\|\mathbf{g}^t\|_2^2 I_k$ and $\gamma^t = \frac{1}{\sqrt{k}}W^\top \mathbf{g}^t + \Lambda^t \hat{\mathbf{z}}^t$.
- 10: **Marginals estimation:**
- 11: $\hat{\mathbf{v}}^{t+1} = f_v(\mathbf{B}_v^t, A_v^t, \omega^t, V^t)$ and $\hat{\mathbf{c}}_v^{t+1} = \partial_B f_v(\mathbf{B}_v^t, A_v^t, \omega^t, V^t)$,
- 12: $\hat{\mathbf{z}}^{t+1} = f_z(\gamma^t, \Lambda^t)$ and $\hat{\mathbf{c}}_z^{t+1} = \partial_\gamma f_z(\gamma^t, \Lambda^t)$,
- 13: $t = t + 1$.
- 14: **until** Convergence.
- 15: **Output:** $\hat{\mathbf{v}}, \hat{\mathbf{z}}$.

Iteration of the TAP equations of stat. phys.
Optimal among general first-order methods
[Celentano&al '20]

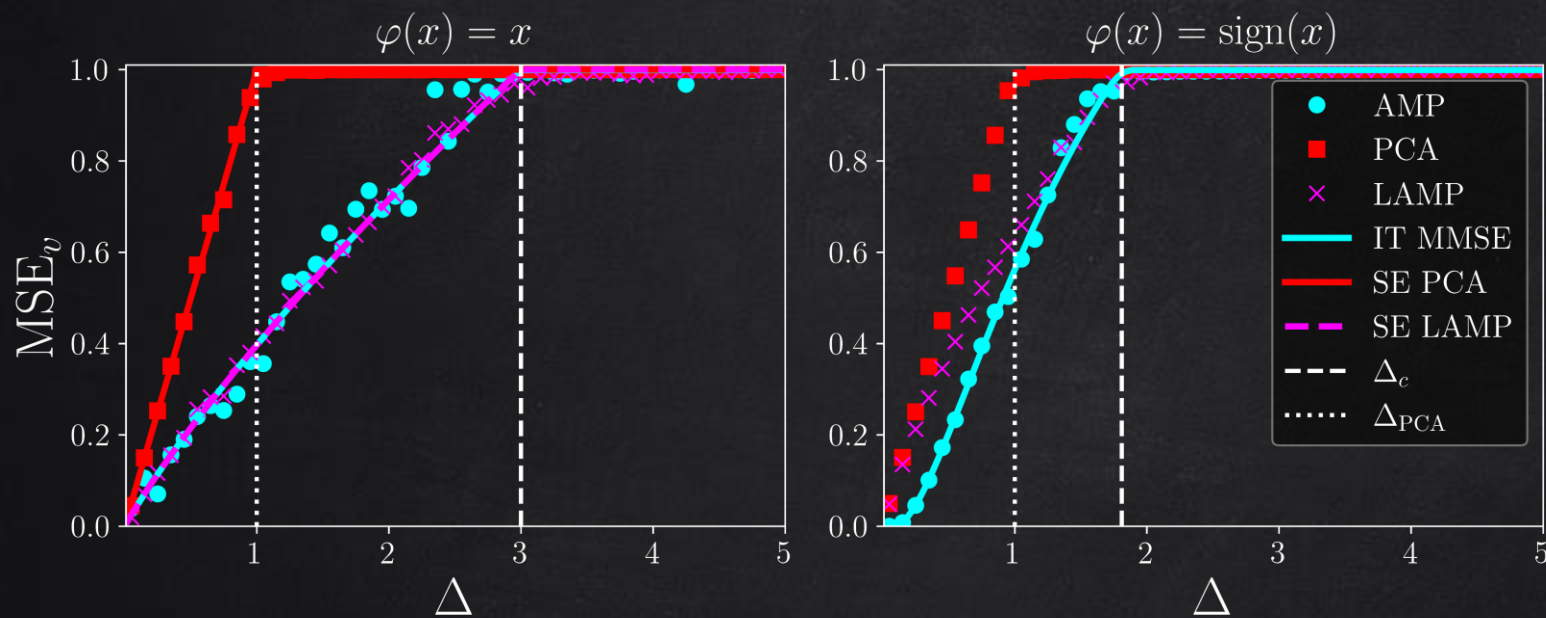


SPECTRAL ALGORITHMS

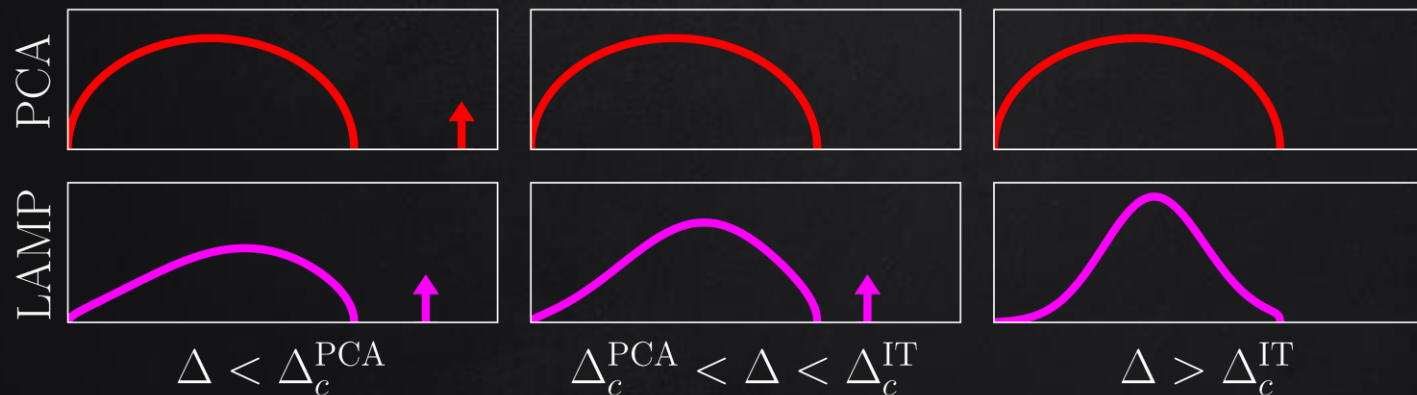
Linearized Approximate Message-Passing (LAMP)

Symmetries → “Trivial” fixed point →

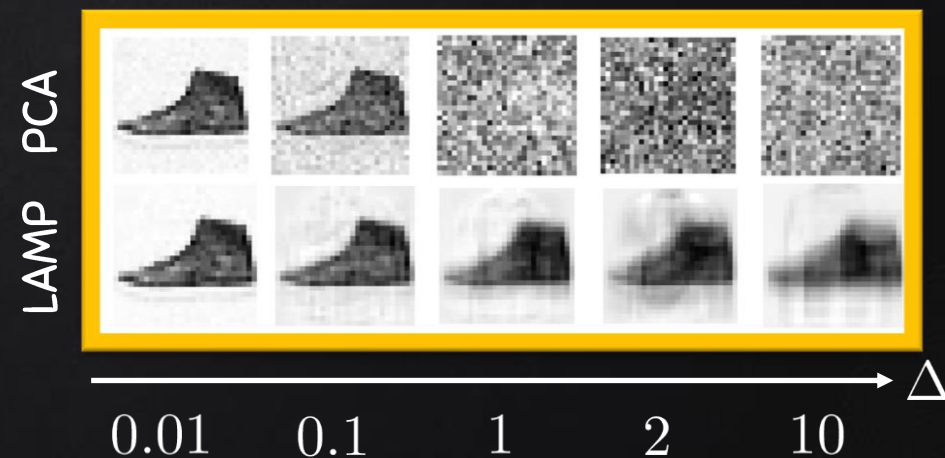
$$\text{Leading eigenvector of } \Gamma_p \equiv \frac{1}{\Delta} \mathbf{K}_p \left[\frac{\mathbf{Y}}{\sqrt{p}} - \mathbf{I}_p \right], \text{ with } \mathbf{K}_p \equiv \frac{1}{k} \mathbb{E}[\mathbf{v}\mathbf{v}^\top].$$



LAMP “beats” the BBP transition of PCA!



Realistic data $\mathbf{K}_p \simeq \frac{1}{n} \sum_{\alpha=1}^n \mathbf{v}^\alpha (\mathbf{v}^\alpha)^\top$



RANDOM MATRIX ANALYSIS

Linear case $\varphi(x) = x$

$$\mathbf{\Gamma}_p = \frac{1}{\Delta} \frac{\mathbf{W}\mathbf{W}^\top}{k} \left[\frac{\mathbf{Y}}{\sqrt{p}} - \mathbf{I}_p \right]$$

μ : asymptotic spectral density of Γ_p ,
with λ_{\max} the right edge of its support.

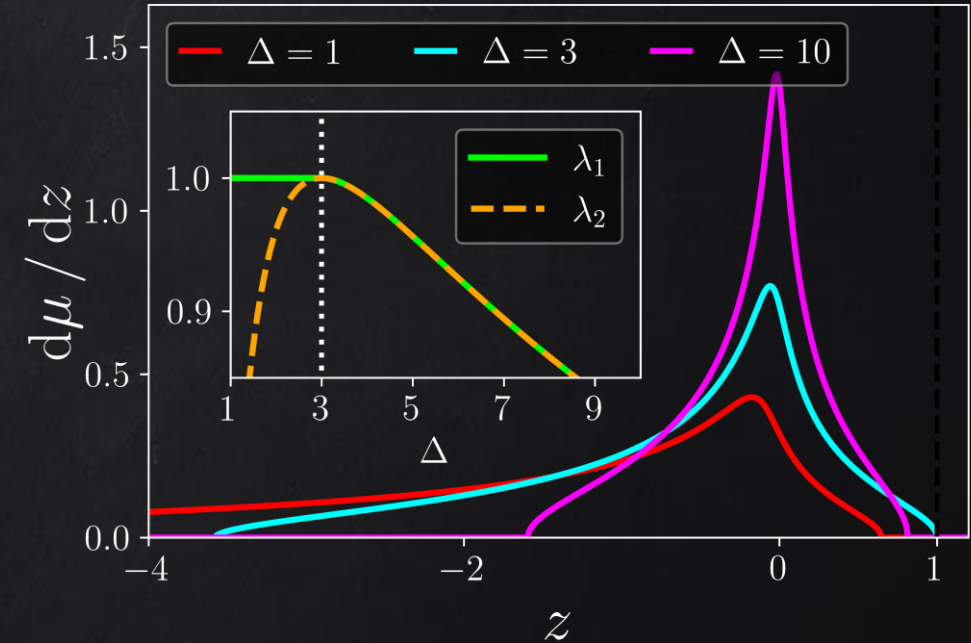
Let $\Delta_c(\alpha) \equiv 1 + \alpha$. We denote $\lambda_1 \geq \lambda_2$ the leading eigenvalues of Γ_p , with normalized eigenvectors $\mathbf{v}_1, \mathbf{v}_2$. Then:

➤ For $\Delta > \Delta_c(\alpha)$, $\lambda_1 \xrightarrow[p \rightarrow \infty]{\text{a.s.}} \lambda_{\max}$ and $\lambda_2 \xrightarrow[p \rightarrow \infty]{\text{a.s.}} \lambda_{\max}$, and $\lambda_{\max} < 1$.

➤ For $\Delta < \Delta_c(\alpha)$, $\lambda_1 \xrightarrow[p \rightarrow \infty]{\text{a.s.}} 1$ and $\lambda_2 \xrightarrow[p \rightarrow \infty]{\text{a.s.}} \lambda_{\max}$, and $\lambda_{\max} < 1$.

Moreover, if $\epsilon(\Delta) \equiv \lim_{p \rightarrow \infty} \frac{1}{p} |\mathbf{v}_1^T \mathbf{v}^*|$, then
$$\begin{cases} \epsilon(\Delta) = 0 & \text{if } \Delta > \Delta_c(\alpha), \\ \epsilon(\Delta) > 0 & \text{if } \Delta < \Delta_c(\alpha). \end{cases}$$

$$\alpha = 2 \Rightarrow \Delta_c = 1 + \alpha = 3$$



- **Main difficulty**: correlation of \mathbf{W} and $\mathbf{Y} = \frac{\mathbf{W}(\mathbf{z}\mathbf{z}^\top)\mathbf{W}^\top}{\sqrt{kp}} + \sqrt{\Delta}\xi$. We use a **cavity computation**, generalizing the classical arguments of [Baik, Ben Arous & Pécché '04].
- Similar results in the **spiked Wishart model**.
- A RMT analysis of **non-linear activations** is still lacking !

SUMMARY ON THE SPIKED MATRIX MODEL

Sparse priors

- ❖ Large hard phases for sparse signals $\rho \ll 1$.
[Deshpande&al '14, Lesieur&al '15]
- ❖ IT weak recovery: $\Delta_c^{\text{IT}} > 1$. But no algorithm can beat the PCA threshold $\Delta_c^{\text{PCA}} = 1$.

Generative priors

- ❖ No algorithmically hard phase, AMP achieves the IT MMSE.
- ❖ Spectral L-AMP outperforms PCA and achieves optimal weak-recovery at $\Delta_c^{\text{LAMP}} = \Delta_c^{\text{AMP}} > 1$.
- ❖ Rigorous RMT analysis of LAMP's performance in the linear case.

Generative priors lead to algorithmically better-behaved problems than sparsity!

Active line of research on the influence of the data structure

- Similar analysis followed in the group, e.g. [Aubin&al '20] for phase retrieval.
- [Goldt&al '19 ; Goldt&al '20]: "hidden manifold" model: theoretical and empirical evidence that many conclusions transfer to **trained** (non-random) **generative priors**.
-

II TOPOLOGY OF HIGH-DIMENSIONAL LANDSCAPES

Squared loss of a noiseless GLM

$$L_2(\mathbf{x}) = \frac{1}{2m} \sum_{\mu=1}^m \left[\varphi(\Phi_{\mu} \cdot \mathbf{x}) - \varphi(\Phi_{\mu} \cdot \mathbf{X}^*) \right]^2$$

$$\begin{aligned} \mathbf{X}^* &\in \mathbb{R}^n, \|\mathbf{X}^*\|^2 = 1 \\ \mathbf{x} &\in \mathbb{R}^n, \|\mathbf{x}\|^2 = 1 \end{aligned}$$

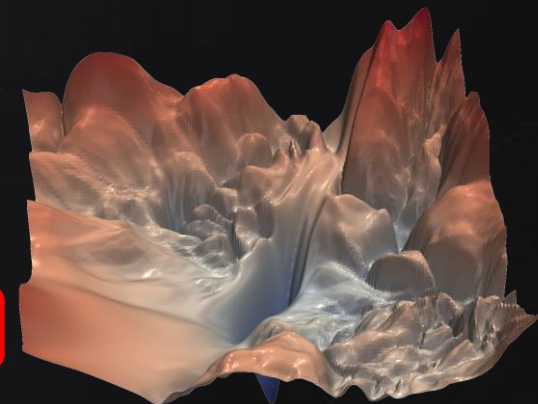
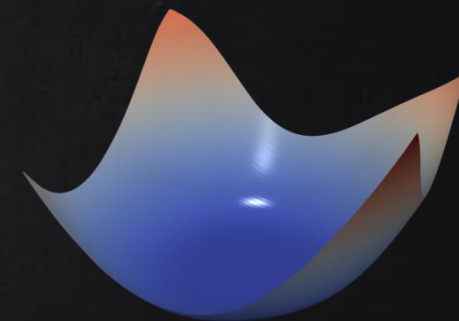
i.i.d. Gaussian data

In many nonconvex problems: one can provably find regimes in which the optimization landscape is 'easy' (matrix decomposition, tensor factorization, neural nets...) [Soudry&al '16; Ge&al '16; Ge&Ma '17; ...]

In practice, local optimization algorithms work far beyond these regimes !

WHY ?

- The bounds on the simplicity of the landscape are not tight enough ?
- Optimization algorithms work in the "hard" regime (i.e. many spurious minima) ?
- Analyze the topological transition, and characterize the 'hard' regime ?

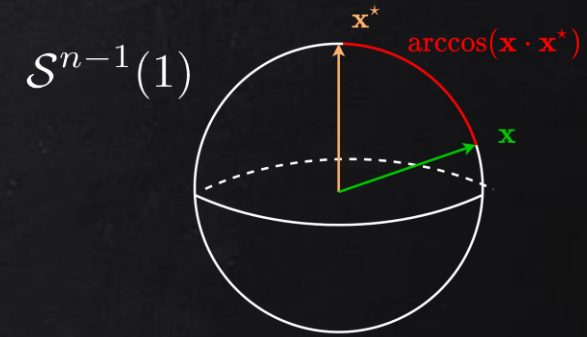


[Li&al '17]

“COMPLEX” LANDSCAPES

$$L_2(\mathbf{x}) = \frac{1}{2m} \sum_{\mu=1}^m \left[\varphi(\Phi_{\mu} \cdot \mathbf{x}) - \varphi(\Phi_{\mu} \cdot \mathbf{X}^*) \right]^2$$

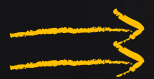
- High-dimensional limit $n, m \rightarrow \infty$ with $\alpha = m/n = \Theta(1)$.
- Count critical points with fixed “energy” $L_2(\mathbf{x})$ and **overlap** $q = \mathbf{X}^* \cdot \mathbf{x}$?



$$\text{Crit}_{\star}(B, Q) \equiv \sum_{\mathbf{x}: \text{grad } L_2(\mathbf{x})=0} \mathbb{1}\{L_2(\mathbf{x}) \in B, \mathbf{x} \cdot \mathbf{X}^* \in Q\}$$

Random variable
(randomness of the data)

- Typically of size $e^{\Theta(n)}$.
- Strongly fluctuating!



- Mean value: **annealed complexity** $\Sigma_{\star}^{(\text{an.})}(B, Q) \equiv \lim_{n \rightarrow \infty} \frac{1}{n} \ln \mathbb{E} \text{Crit}_{\star}(B, Q)$
- Typical value: **quenched complexity** $\Sigma_{\star}^{(\text{qu.})}(B, Q) \equiv \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} \ln \text{Crit}_{\star}(B, Q)$

⚠ Quenched \neq Annealed!
(Few exceptions [Subag '17])

Complexity at a fixed **index** $\text{Crit}_k(B, Q) \equiv \sum_{\mathbf{x}: \text{grad } L_2(\mathbf{x})=0} \mathbb{1}\{\text{i[Hess } L_2(\mathbf{x})] = k, L_2(\mathbf{x}) \in B, \mathbf{x} \cdot \mathbf{X}^* \in Q\}$

OUR MAIN TOOL: THE KAC-RICE FORMULA

Theorem (Kac-Rice)

e.g. [Adler&Taylor '09]

$f : \mathcal{S}^{n-1}(1) \rightarrow \mathbb{R}$ is a smooth random function that is a.s. Morse. Then:

$$\mathbb{E} \text{Crit}_k(f) = \int_{\mathcal{S}^{n-1}(1)} \sigma(d\mathbf{x}) \varphi_{\text{grad } f(\mathbf{x})}(0) \times \mathbb{E} \left[\left| \det \text{Hess } f(\mathbf{x}) \right| \mid \text{grad } f(\mathbf{x}) = 0; i(\text{Hess } f(\mathbf{x})) = k \right]$$

Density of the (random) gradient taken in 0

- Random differential geometry \longrightarrow Random matrix theory.
- Many possible refinements: fix the value of $f(\mathbf{x})$, higher-order moments, ...

- Takeaways:
- **Distribution of $\{\text{Hess } f(\mathbf{x}) \mid \text{grad } f(\mathbf{x}) = 0\}$:** intractable for “generic” functions!
 - To compute $\mathbb{E} \text{Crit}_k(f)$: we need the **large deviations of the k-th largest eigenvalue of the Hessian.**

\longrightarrow Applications limited to Gaussian random functions

Pure p-spin and variants

$$f(\mathbf{x}) = \sum_{i_1, \dots, i_p} \underbrace{J_{i_1, \dots, i_p}}_{\mathcal{N}(0,1)} x_{i_1} \cdots x_{i_p}$$

[Bray&Moore '80, Crisanti&al '95, Fyodorov&al '07, Auffinger&al '13, Ros&al '19,]

MAIN RESULTS

$$L_2(\mathbf{x}) = \frac{1}{2m} \sum_{\mu=1}^m \left[\varphi(\Phi_\mu \cdot \mathbf{x}) - \varphi(\Phi_\mu \cdot \mathbf{X}^*) \right]^2 \xrightarrow{\text{simplification}} \boxed{L_1(\mathbf{x}) = \frac{1}{m} \sum_{\mu=1}^m \varphi(\Phi_\mu \cdot \mathbf{x})} \quad \text{Theorem \& proof transfer to } L_2(\mathbf{x}).$$

First exact high-dimensional result obtained with Kac-Rice for non-Gaussian functions!

Theorem

$$\Sigma_\star^{(\text{an})}(B) = \frac{1 + \ln \alpha}{2} + \sup_{\substack{\nu \in \mathcal{M}_1^+(\mathbb{R}) \\ \int \nu(dt) \varphi(t) \in B}} \left[-\frac{1}{2} \ln \left\{ \int \nu(dt) \varphi'(t)^2 \right\} - \underbrace{\alpha H(\nu | \mathcal{N}(0, 1))}_{\text{Relative entropy}} + \kappa_\alpha(\nu) \right]$$

Involved function: related to the logarithmic potential of the asymptotic spectral measure of \mathbf{zDz}^\top / m if $\mathbf{z} \in \mathbb{R}^{n \times m}$ is a Gaussian i.i.d. matrix and $D_\mu = \varphi''(y_\mu)$ with $y_\mu \stackrel{\text{i.i.d.}}{\sim} \nu$.

- Term $\kappa_\alpha(\nu) \implies$ Analytically very hard variational problem!
- We derive a **closed formula for the quenched complexity**, using the heuristic replica method [Parisi&al '87, Ros&al '19].
- Generic result, applies to **mixture of Gaussians, binary classification, ...**

SKETCH OF PROOF

$$L_1(\mathbf{x}) = \frac{1}{m} \sum_{\mu=1}^m \varphi(\Phi_{\mu} \cdot \mathbf{x})$$

Kac-Rice formula \implies Hessian conditioned by zero gradient.

Main idea: Condition everything by the i.i.d. Gaussian random variables $y_{\mu} \equiv \Phi_{\mu} \cdot \mathbf{x}$

$\mathbb{E}_{\Phi}[\cdots] = \mathbb{E}_{\mathbf{y}}\mathbb{E}[\cdots | \mathbf{y}]$ Under this conditional distribution, and under the gradient being zero:

$$\text{Hess } L_1(\mathbf{x}) \stackrel{\text{d}}{=} \frac{1}{m} \sum_{\mu=1}^m \varphi''(y_{\mu}) \mathbf{z}_{\mu} \mathbf{z}_{\mu}^{\top} + t(\mathbf{y}) \mathbb{1}_n (+ \text{finite rank term})$$

\mathbf{z}_{μ} : i.i.d. standard Gaussian vectors

“Generalized” version of a sample covariance matrix

- We prove fast enough concentration of $\ln |\det \text{Hess}|$ as a function of $\{y_{\mu}\}_{\mu=1}^m$: $\mathbb{E} |\det \text{Hess}| \simeq e^{\mathbb{E} \ln |\det \text{Hess}|}$
- The expectation only depends on the empirical distribution $\nu_{\mathbf{y}} \equiv (1/m) \sum_{\mu=1}^m \delta_{y_{\mu}}$: $\mathbb{E} \ln |\det \text{Hess}| \simeq m \kappa_{\alpha}(\nu_{\mathbf{y}})$
- We use Sanov's theorem: the law of $\nu_{\mathbf{y}}$ satisfies large deviations with rate function: $I(\nu) = \alpha H(\nu | \mathcal{N}(0, 1))$
- Kac-Rice formula and Varadhan's lemma $\implies \Sigma_{\star}^{(\text{an})} = \sup_{\nu \in \mathcal{M}_1^{+}(\mathbb{R})} [\kappa_{\alpha}(\nu) + \underline{G(\nu)} - \alpha H(\nu | \mathcal{N}(0, 1))]$ \blacksquare

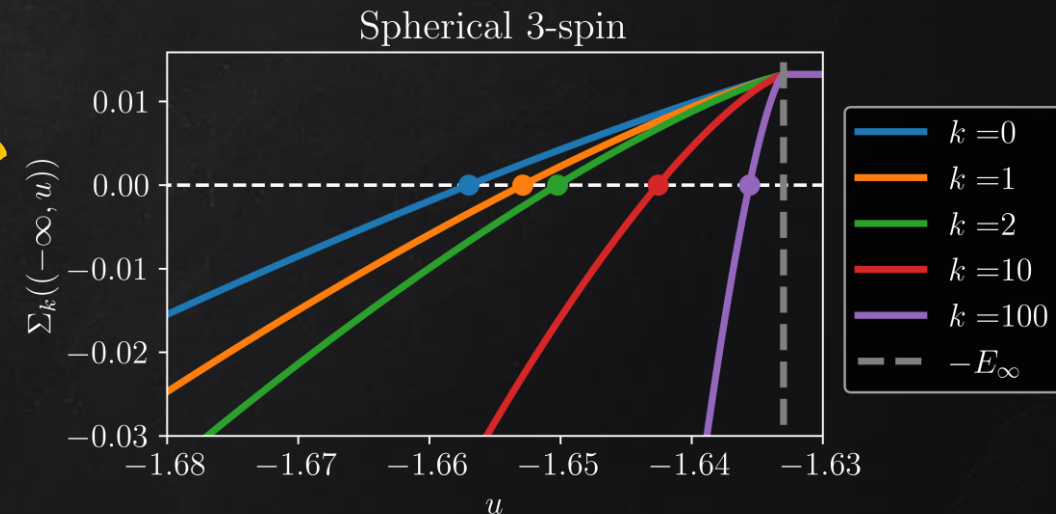
(Gradient density term in Kac-Rice)

SUMMARY & OUTLOOK

- ❖ First exact high-dimensional result obtained with Kac–Rice for non-Gaussian functions !
- ❖ Generalizes to other models: mixture of two Gaussians, binary linear classification... → Neural networks ?

Physical discussion is lacking. Many problems ahead:

- Numerically solve the variational problem? Sign of the complexity given α, φ ? ...
- Count local minima ? We need a LDP for $\lambda_{\min}(\text{Hess } L(\mathbf{x}))$...
→ Obtained in [A.M., EPL 2021] ! To be continued...
“Tilting” of the measure [Biroli&Guionnet’20, Belinschi&al’20, Guionnet&al’20, Husson’20, Augeri&al’21]...

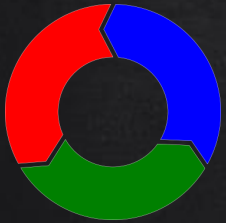


- Discrete systems ?
 - ❖ TAP approach ?
 - ❖ Recent algorithmic progress on F-RSB spin glasses [Subag’21, Montanari’21, El Alaoui & al’20].

CONCLUDING REMARKS

Theory of inference/learning

Data
structure



Architecture
of the model

Algorithms



Diversified toolbox

Statistical physics

Replica, message-passing,
Plefka expansions, DMFT...

Probabilistic methods

Guerra interpolation,
concentration identities...

Topological approach

Kac-Rice, large deviations,
random matrix theory...

To name a few...

Many exciting
questions, e.g.:

- Interplay in more involved learning models ?
- What if we do not know how the data was generated ?

Realistic deep networks ?



- Classical statistical physics approach.
- Extensive-rank HCIZ ? [Matytsin '93]; [Guionnet&Zeitouni '02]
- Plefka expansion ? [Parisi&Potters '95]; [A.M.&al '19], ...

The "extensive-rank problem"

One challenge among many....

Conclusion



STAY
CALM

AND

CONTINUE
TESTING

TESTING

Or

$$Y = UV^T + Z$$