



PSL



EPFL



CONSTRUCTION OF OPTIMAL SPECTRAL METHODS IN PHASE RETRIEVAL

[arXiv:2012.04524](#) ; [MSML'21](#)

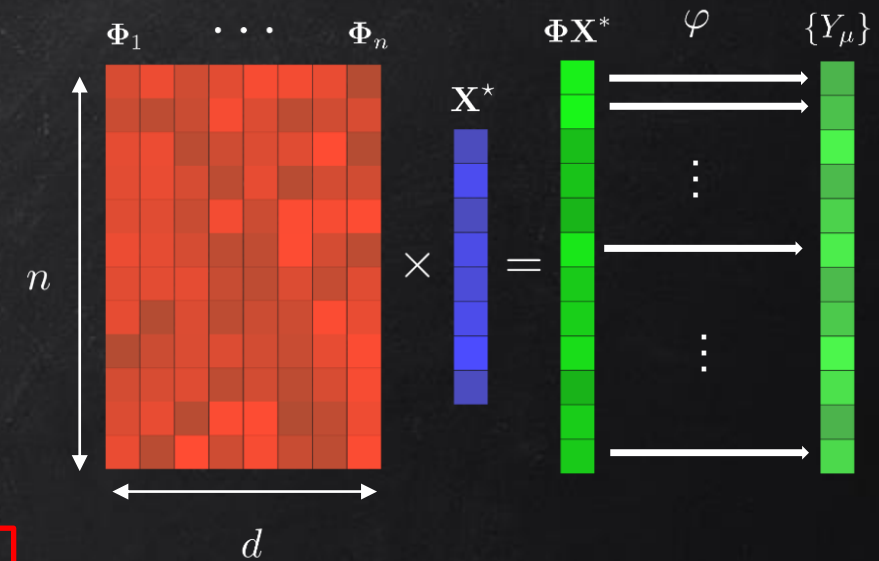
Antoine Maillard

Joint work with Florent Krzakala, Yue M. Lu & Lenka Zdeborová

Random Matrix Theory and Networks – June 16th 2021

INFERENCE IN HIGH DIMENSIONS

Goal: Recover a d -dimensional signal \mathbf{X}^* from n data points $\{\Phi_\mu, Y_\mu\}_{\mu=1}^n$ generated as:



Generalized Linear Model (GLM)

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

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

(Probabilistic) channel
with possible noise.

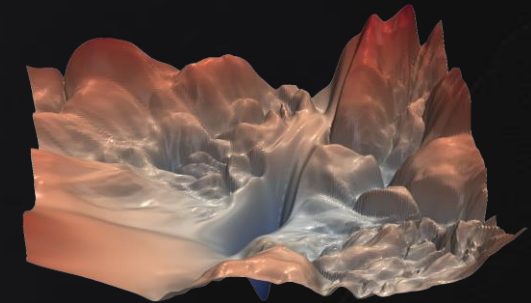
Sensing matrix (real/complex)

Signal (real/complex), d -dimensional

Real / Complex
 $\beta = 1$ $\beta = 2$

In general, one needs to solve a **highly non-convex** optimization problem in high dimensions.

Example (empirical risk minimization – square loss): $\hat{\mathbf{X}} = \operatorname{argmin}_{\mathbf{x}} \sum_{\mu=1}^n \left(Y_\mu - \varphi(\mathbf{x} \cdot \Phi_\mu / \sqrt{n}) \right)^2$



PHASE RETRIEVAL

$$Y_\mu \sim P_{\text{out}} \left(\cdot \middle| \frac{1}{\sqrt{d}} \sum_{i=1}^d \Phi_{\mu i} X_i^* \right)$$

In **phase retrieval**, one measures the modulus $P_{\text{out}}(y|z) = P_{\text{out}}(y||z|)$, e.g. noiseless $Y_\mu = \frac{1}{d} |(\Phi \mathbf{X}^*)_\mu|^2$; Poisson-noise $Y_\mu \sim \text{Pois}(\Lambda |(\Phi \mathbf{X}^*)_\mu|^2 / d)$.

- Classical problem of learning with a non-convex landscape.
- Arises in **signal processing, statistical estimation, optics, X-ray crystallography, astronomy, microscopy**... : optical detectors lose information on the phase of the signals.

How to solve this problem efficiently in high dimensions ? $n, d \rightarrow \infty$

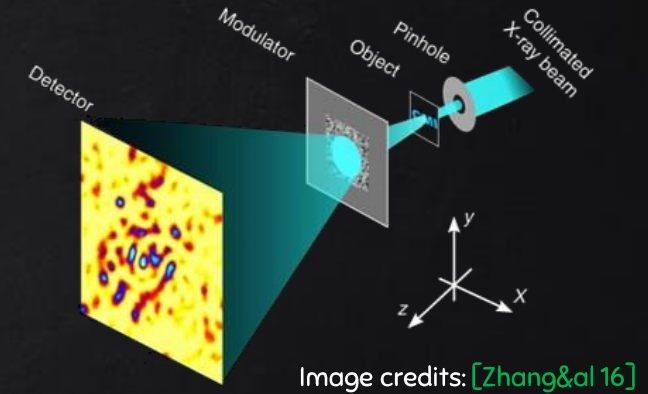


Image credits: [Zhang&al 16]

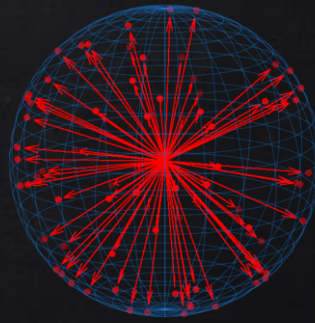
- SDP relaxations [Candès&al '15a&b, Waldspurger&al '15, Goldstein&al '18, ...]
- Non-convex optimization procedures [Netrapalli&al '15, Candès&al '15c, ...]
- Approximate Message-Passing [Barbier&al '19, A.M.&al '20]

Computationally heavy /
Need informed initialization

Spectral methods

[Mondelli&al '18, Luo&al '18, Dudeja&al '19, ...]

Our model: The matrix Φ is **right-orthogonally (unitarily) invariant**, i.e. delocalized right-eigenvectors: $\forall \mathbf{U}, \Phi \stackrel{d}{=} \Phi \mathbf{U}$
 The bulk of eigenvalues of $\Phi^\dagger \Phi / d$ converges to a distribution $\nu(x)$, as $n, d \rightarrow \infty$ with $n/d \rightarrow \alpha > 0$.



Examples: Gaussian matrices, product of Gaussians, random column-orthogonal/unitary, any $\Phi \equiv \mathbf{U} \mathbf{S} \mathbf{V}^\dagger$ with $S_i^2 \stackrel{\text{i.i.d.}}{\sim} \nu$.

WEAK-RECOVERY THRESHOLD IN PHASE RETRIEVAL [A.M&al'20]

What is the minimal number of measurements $\alpha = n/d$ necessary to beat a random guess in polynomial time?

➡ If the signal \mathbf{X}^* has norm $\frac{1}{n} \|\mathbf{X}^*\|^2 = \rho$, this threshold $\alpha_{\text{WR, Algo}}$ is the only solution to:

$$\alpha = \frac{\langle \lambda \rangle_\nu^2}{\langle \lambda^2 \rangle_\nu} \left[1 + \left\{ \int_{\mathbb{R}} dy \frac{\left(\int_{\mathbb{K}} \mathcal{D}_\beta z (|z|^2 - 1) P_{\text{out}} \left[y \left| \sqrt{\frac{\rho \langle \lambda \rangle_\nu}{\alpha}} z \right| \right] \right)^2}{\int_{\mathbb{K}} \mathcal{D}_\beta z P_{\text{out}} \left[y \left| \sqrt{\frac{\rho \langle \lambda \rangle_\nu}{\alpha}} z \right| \right]} \right\}^{-1} \right]$$



This is an implicit equation.

Example: Noiseless phase retrieval:

$$\alpha = \left(1 + \frac{\beta}{2} \right) \frac{\langle \lambda \rangle_\nu^2}{\langle \lambda^2 \rangle_\nu}$$

- Gaussian matrices: $\alpha_{\text{WR, Algo}} = \frac{\beta}{2}$ [Barbier&al '19, Mondelli &al '19, Luo&al '19]
- Random column-orthogonal/unitary matrices: $\alpha_{\text{WR, Algo}} = 1 + \frac{\beta}{2}$ [Dudeja&al '20] for $\beta = 2$

CONSTRUCTION OF SPECTRAL METHODS

Given a generic phase retrieval problem, we want to design **spectral methods** such that:

- Weak-recovery is achieved for all $\alpha > \alpha_{\text{WR,Algo}}$.
- For all α the method is **optimal among all possible spectral methods** in terms of estimation error:
d

$$\text{MSE} \equiv \frac{1}{d\rho} \|\mathbf{X}^* - \hat{\mathbf{X}}_{\text{spectral}}\|^2$$

This talk: Three different strategies, related to the **statistical-physics** approach to high-dimensional inference.

- Method I: Naïve generalization of what is known for Gaussian matrices.
- Method II: Linearization of **message-passing** algorithms.
- Method III: **Bethe Hessian** analysis from the Thouless–Anderson–Palmer [TAP77] free energy.

METHOD I: NAIVE APPROACH

$$y_\mu \sim P_{\text{out}}\left(\cdot \mid \frac{1}{\sqrt{d}} \sum_{i=1}^d \Phi_{\mu i} X_i^*\right)$$

Most previous works reduced to methods of the type

$$\mathbf{M}(\mathcal{T}) \equiv \frac{1}{d} \sum_{\mu=1}^n \mathcal{T}(y_\mu) \Phi_\mu \Phi_\mu^\dagger$$

[Chen&al '15, Wang&al '16, Zhang&al '17, Mondelli&al '19, Luo&al '19, Dudgeon&al '19]

Idea: For large sample size $n \gg d$ we expect $\mathbf{M} \simeq \mathbb{E}[\mathbf{M}] = a\mathbf{I}_n + b\mathbf{X}^*(\mathbf{X}^*)^\dagger$.

For Gaussian matrices Φ the optimal method in this class is given by

$$\mathcal{T}_{\text{Gaussian}}^*(y) \equiv \frac{\partial_\omega g_{\text{out}}(y_\mu, 0, \rho)}{1 + \rho \partial_\omega g_{\text{out}}(y_\mu, 0, \rho)}$$

$$\partial_\omega g_{\text{out}}(y_\mu, 0, \sigma^2) = -\frac{1}{\sigma^2} + \frac{1}{\sigma^4} \frac{\int_{\mathbb{K}} dx e^{-\frac{\beta}{2\sigma^2}|x|^2} |x|^2 P_{\text{out}}(y_\mu|x)}{\int_{\mathbb{K}} dx e^{-\frac{\beta}{2\sigma^2}|x|^2} P_{\text{out}}(y_\mu|x)}$$

($\mathbb{K} = \mathbb{R}, \mathbb{C}$)

- In **noiseless phase retrieval** one has $\mathcal{T}_{\text{Gaussian}}^*(y) = 1 - 1/y$.
- It achieves **weak recovery** at the optimal threshold $\alpha = \alpha_{\text{WR, Algo}}$.
- Optimal also in terms of achieved Mean Squared Error among the class of $\mathbf{M}(\mathcal{T})$.
- We can naively use it for all matrices: $\mathbf{M}_{\text{naive}} \equiv \mathbf{M}(\mathcal{T}_{\text{Gaussian}}^*)$.

METHOD II: LINEARIZATION OF MESSAGE-PASSING

- [Schniter&al '16, A.M.&al '20]: For GLMs with rotationally-invariant matrices, the best-known polynomial-time algorithm (in terms of estimation error) is given by *Generalized Vector Approximate Message-Passing* (G-VAMP).
- But G-VAMP is computationally very expensive \longrightarrow Construct a *spectral method* from the G-VAMP iterations.

Symmetry of the phase retrieval problem $P_{\text{out}}(y|z) = P_{\text{out}}(y||z|)$ \longrightarrow G-VAMP has a trivial fixed point at $\hat{\mathbf{x}} = 0$.

Linearize G-VAMP around this point.

Linearized Approximate Message-Passing (LAMP) spectral method.

$$\mathbf{M}_{\text{LAMP}} \equiv \frac{\rho\langle\lambda\rangle_\nu}{\alpha} \left(\frac{\alpha}{\langle\lambda\rangle_\nu} \frac{\Phi\Phi^\dagger}{d} - \mathbf{I}_n \right) \text{Diag}(\{\partial_\omega g_{\text{out}}(y_\mu, 0, \rho\langle\lambda\rangle_\nu/\alpha)\}) \longrightarrow \hat{\mathbf{x}} \equiv \frac{\Phi^\dagger \text{Diag}(\{\partial_\omega g_{\text{out}}(y_\mu, 0, \rho\langle\lambda\rangle_\nu/\alpha)\}) \hat{\mathbf{u}}}{\|\Phi^\dagger \text{Diag}(\{\partial_\omega g_{\text{out}}(y_\mu, 0, \rho\langle\lambda\rangle_\nu/\alpha)\}) \hat{\mathbf{u}}\|} \sqrt{d\rho}.$$

\mathbf{M}_{LAMP} is a $n \times n$ non-Hermitian matrix (complex spectrum).

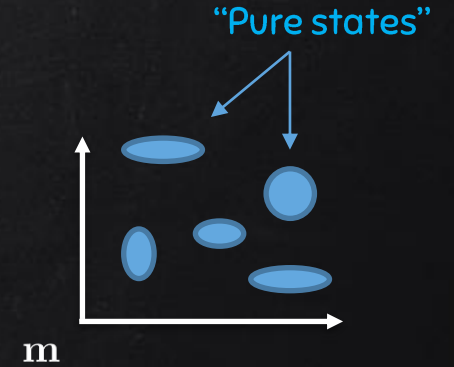
$\hat{\mathbf{u}}$: top eigenvector of \mathbf{M}_{LAMP} .

Similar approaches have been applied in community detection [Krzakala&al '13], phase retrieval with unitary matrices [Ma&al '21] and spiked matrix estimation [Aubin&al '20].

METHOD III: TAP LANDSCAPE AND BETHE HESSIAN

Thouless-Anderson-Palmer approach [TAP77]

- The posterior measure of $\mathbf{x}|\mathbf{Y}$ (the *Gibbs measure*) decomposes along **pure states**.
- These pure states can be found by “tilting” the measure, imposing $m_i = \langle x_i \rangle$ and $\sigma_i^2 = \text{Var}(x_i)$:
They are the **maxima of the free entropy** of this constrained measure, as a function of (\mathbf{m}, σ) .



- TAP free entropy for **rotationally-invariant generalized linear models** derived in [A.M.&al '19], generalizing [Parisi&Potters '95]:
Involved but explicit!

$$f_{\text{TAP}}(\mathbf{m}) = \sup_{\sigma \geq 0} \sup_{\mathbf{g} \in \mathbb{K}^n} \sup_{\substack{b \geq 0 \\ r \geq 0}} \text{extr}_{\omega \in \mathbb{K}^n} \text{extr}_{\substack{\lambda \in \mathbb{K}^d \\ \gamma \geq 0}} \left[\frac{\beta}{d} \sum_{i=1}^d \lambda_i \cdot m_i + \frac{\beta \gamma}{2d} (d\sigma^2 + \sum_{i=1}^d |m_i|^2) - \frac{\beta}{d} \sum_{\mu=1}^n \omega_\mu \cdot g_\mu - \frac{\beta b}{2d} \left(\sum_{\mu=1}^n |g_\mu|^2 - \alpha d r \right) + \frac{1}{d} \sum_{i=1}^d \ln \int_{\mathbb{K}} P_0(dx) e^{-\frac{\beta \gamma}{2} |x|^2 - \beta \lambda_i \cdot x} \right. \\ \left. + \frac{\alpha}{n} \sum_{\mu=1}^n \ln \int_{\mathbb{K}} \frac{dh}{\left(\frac{2\pi b}{\beta} \right)^{\beta/2}} P_{\text{out}}(y_\mu | h) e^{-\frac{\beta |h - \omega_\mu|^2}{2b}} + \frac{\beta}{d} \sum_{i=1}^d \sum_{\mu=1}^n g_\mu \cdot \left(\frac{\Phi_{\mu i}}{\sqrt{d}} m_i \right) + \beta F(\sigma^2, r) \right].$$

$$F(x, y) \equiv \inf_{\zeta_x, \zeta_y > 0} \left[\frac{\zeta_x x}{2} + \frac{\alpha \zeta_y y}{2} - \frac{\alpha - 1}{2} \ln \zeta_y - \frac{1}{2} \langle \ln(\zeta_x \zeta_y + \lambda) \rangle_\nu \right] - \frac{1}{2} \ln x - \frac{\alpha}{2} \ln y - \frac{1 + \alpha}{2}.$$

Weak-recovery impossible $\alpha < \alpha_{\text{WR}}$

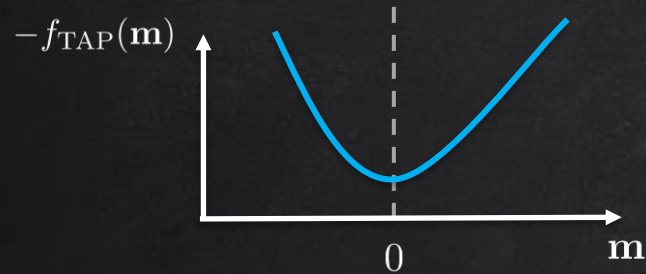
Global maximum of f_{TAP} in $\mathbf{m} = 0$: the uninformative “paramagnetic” point.

Weak-recovery possible $\alpha > \alpha_{\text{WR}}$

$\mathbf{m} = 0$ is an unstable stationary point of f_{TAP} , which has a global maximum in $\mathbf{m} \neq 0$ (optimal estimator).

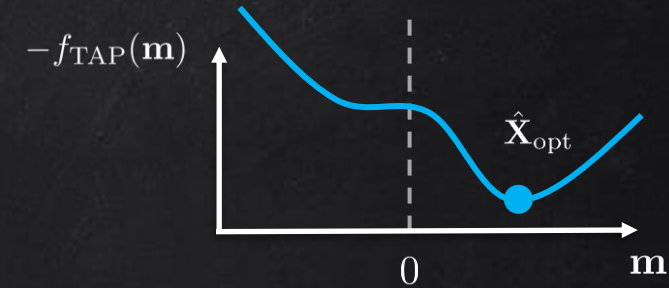
Weak-recovery impossible $\alpha < \alpha_{\text{WR}}$

Global maximum of f_{TAP} in $\mathbf{m} = 0$: the uninformative “paramagnetic” point.



Weak-recovery possible $\alpha > \alpha_{\text{WR}}$

$\mathbf{m} = 0$ is an unstable stationary point of f_{TAP} , which has a global maximum in $\mathbf{m} \neq 0$ (optimal estimator).



A spectral method can only use the physical information available in the uninformative point $\mathbf{m} = 0$.

Compute the Hessian of f_{TAP} at the paramagnetic point.

Constructive derivation of a spectral method that is conjectured to be optimal.

TAP – Bethe Hessian spectral method.

$$\mathbf{M}_{\text{TAP}} \equiv -d \nabla^2 f_{\text{TAP}}(\mathbf{m} = 0) = -\frac{1}{\rho} \mathbf{I}_d + \frac{1}{d} \sum_{\mu=1}^n \frac{\partial_{\omega} g_{\text{out}}(y_{\mu}, 0, \rho \langle \lambda \rangle_{\nu} / \alpha)}{1 + \frac{\rho \langle \lambda \rangle_{\nu}}{\alpha} \partial_{\omega} g_{\text{out}}(y_{\mu}, 0, \rho \langle \lambda \rangle_{\nu} / \alpha)} \Phi_{\mu} \Phi_{\mu}^{\dagger}$$

Similar to previous strategies in community detection.
[Saade&al'14]

RELATIONS BETWEEN THE METHODS AND THE MARGINALITY PUZZLE

Analytical study of the three spectral methods



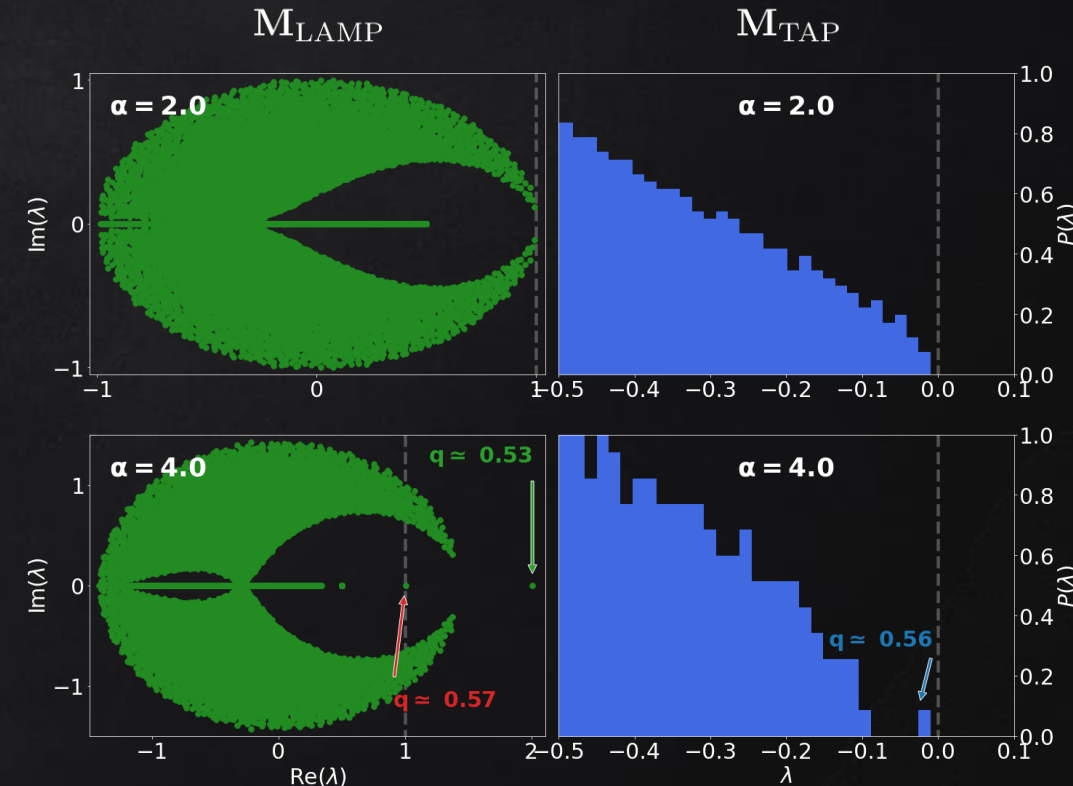
- M_{TAP} is the “naïve” generalization of Method I.
 - M_{TAP} and M_{LAMP} : transition at the optimal $\alpha_{\text{WR, Algo}}$.
- The fixed points of the TAP free entropy are in exact correspondence with the ones of G-VAMP [A.M.&al '19].
 - **Marginal stability:** When recovery is possible, the largest eigenvalue of M_{TAP} concentrates on 0, and is in **exact correspondence** with an eigenvalue of M_{LAMP} that concentrates on 1.
 - **Instability of M_{LAMP} :** No other eigenvector of M_{TAP} achieves non-trivial performance, while the dominant eigenvalue of M_{LAMP} is **another non-trivial estimator that is suboptimal**.

Puzzling disparity between methods that should be equivalent.

- Complex Gaussian Φ
- Poisson channel ($\Lambda = 1$)

$$P_{\text{out}}(y|z) = e^{-\Lambda|z|^2} \sum_{k=0}^{\infty} \delta(y - k) \frac{\Lambda^k |z|^{2k}}{k!} \longrightarrow \alpha_{\text{WR, Algo}} = 2$$

- We denote the overlap $q = \frac{1}{d} \sum_{i=1}^d X_i^* \hat{x}_i$



OPTIMAL SPECTRAL METHOD

$$\mathbf{M}(\mathcal{T}) \equiv \frac{1}{d} \sum_{\mu=1}^n \mathcal{T}(y_{\mu}) \Phi_{\mu} \Phi_{\mu}^{\dagger}$$

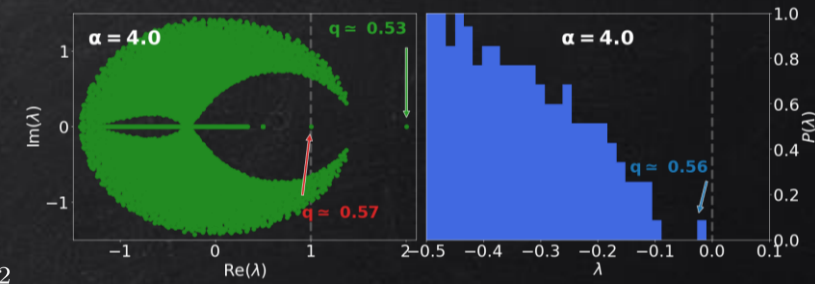
From the Bethe Hessian analysis

Main conjecture: For any right-orthogonally invariant sensing matrix, the optimal spectral method (in terms of weak-recovery threshold and achieved error) belongs to the class of matrices $\mathbf{M}(\mathcal{T})$ and is attained in:

$$\mathcal{T}^*(y) = \frac{\partial_{\omega} g_{\text{out}}(y_{\mu}, 0, \rho \langle \lambda \rangle_{\nu} / \alpha)}{1 + \frac{\rho \langle \lambda \rangle_{\nu}}{\alpha} \partial_{\omega} g_{\text{out}}(y_{\mu}, 0, \rho \langle \lambda \rangle_{\nu} / \alpha)}$$

- We did not assume anything on the form of the method, yet the optimal spectral method we constructed is in the class of $\mathbf{M}(\mathcal{T})$ matrices: we **confirm the validity of the restriction of previous works on spectral methods!**
- The optimal spectral method does not depend on the **spectrum of the sensing matrix (apart from a global scaling)**, nor on the **sampling ratio α !**
 - ➡ Very different from the optimal algorithms! **[A.M.&al '20]**
 - ➡ Consequences for practitioners: **one only needs to know the observation channel** to construct the method!

SPECTRAL METHODS PERFORMANCE



Noiseless complex phase retrieval $Y_\mu = \frac{1}{d} |\Phi X^*|^2$

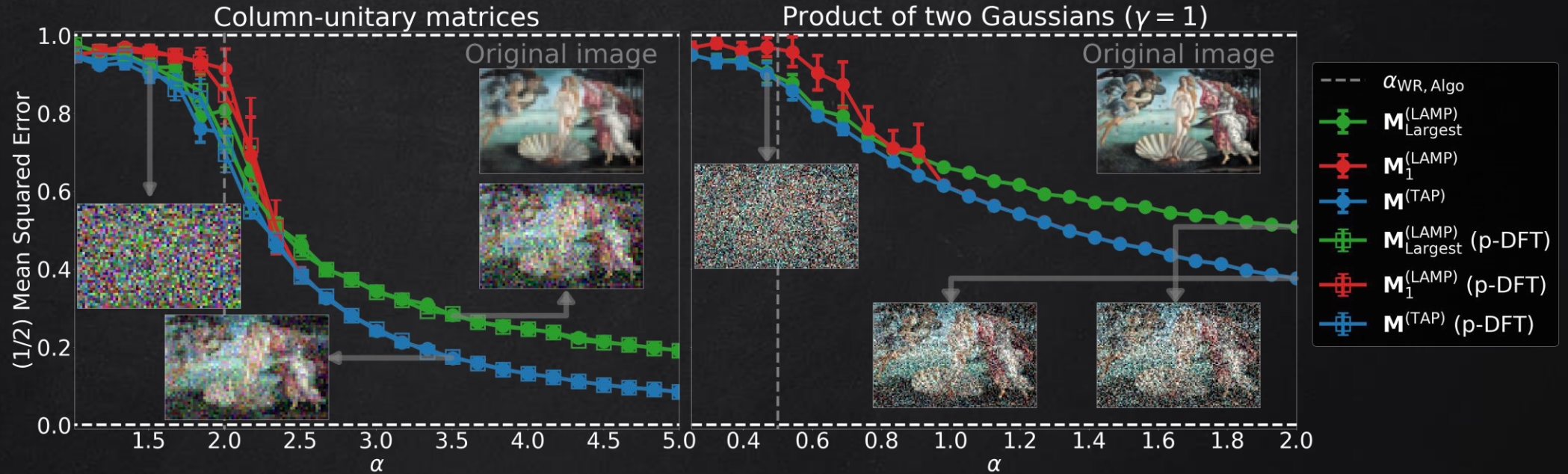


Image: 1280x820, reduced by a factor 20 (left) or 10 (right).

- All three methods share the **same weak-recovery threshold**, in agreement with the best polynomial-time algorithm.
- $\hat{x}_{LAMP}(\lambda = 1) \iff \hat{x}_{TAP}$, achieving the best overlap. Otherwise $\hat{x}_{LAMP}(\lambda_{max})$ is suboptimal in terms of MSE.
- Our theory stays valid for **matrices with controlled structure**. (partial DFT \equiv randomly subsampled DFT)

SPECTRAL INITIALIZATION IN NOISELESS PHASE RETRIEVAL

We combine the optimal spectral method \mathbf{M}_{TAP} with gradient descent on the square loss $L(\mathbf{x}) \equiv \frac{1}{2n} \sum_{\mu=1}^n \left\{ \left| \frac{(\Phi \mathbf{x})_{\mu}}{\sqrt{d}} \right|^2 - \left| \frac{(\Phi \mathbf{X}^*)_{\mu}}{\sqrt{d}} \right|^2 \right\}^2$.




- **Noiseless** phase retrieval with **randomly subsampled DFT** sensing matrix. [A.M&al '20]
 - Already **perfect recovery** at $\alpha \in (3, 4)$. For partial-DFT, perfect recovery with the best polynomial-time algorithm is $\alpha_{\text{PR}} \simeq 2, 3$
- ➡ Very competitive while computationally cheap!

CONCLUSION AND PERSPECTIVES

Main contributions

- Constructive derivation of a **conjecturally optimal spectral method** in generic phase retrieval problems, in a framework that encompasses real/complex variables and a wide variety of sensing matrices.
- Our results apply to **randomly subsampled DFT** matrices and to **real image** (i.e. structured signal) recovery.
- We use two fundamentally equivalent approaches – **message-passing linearization** and **Bethe Hessian analysis** – that yield the same optimal performance, associated with a **marginal stability of the linear dynamics**.

Theory far from complete

- The “marginality vs instability” puzzle: In M_{LAMP} the optimal method is “hidden” inside the bulk and marginally stable, while the dominant eigenvalue is unstable and suboptimal.
- What if we do not know how the data was generated ? What become of the thresholds and of the performance of the spectral estimators ?  A first observation: marginality disappears when using a mismatched channel distribution.

THANK YOU !