

# Graph-based approximate message passing iterations

Cedric Gerbelot

Courant Institute of Mathematical Sciences, NYU

---

Joint work with Raphaël Berthier (EPFL), arXiv 2109.11905

To appear in *Information and Inference : A Journal of the IMA*

Cargese 2023

*Statistical physics and machine learning back together again*

# The starting point : probabilistic inference

A hidden process generates  $\mathbf{w} \in \mathbb{R}^d$  with large  $d$

$$\mathbf{w} \sim p_{\mathbf{w}}(\mathbf{w})$$

# The starting point : probabilistic inference

A hidden process generates  $\mathbf{w} \in \mathbb{R}^d$  with large  $d$

$$\mathbf{w} \sim p_{\mathbf{w}}(\mathbf{w})$$

We observe  $\mathbf{y} \in \mathbb{R}^n$  s.t.

$$\mathbf{y} \sim p_{\mathbf{y}}(\mathbf{y}|\mathbf{w})$$

Estimate  $\mathbf{w}$  ?

# The starting point : probabilistic inference

A hidden process generates  $\mathbf{w} \in \mathbb{R}^d$  with large  $d$

$$\mathbf{w} \sim p_{\mathbf{w}}(\mathbf{w})$$

We observe  $\mathbf{y} \in \mathbb{R}^n$  s.t.

$$\mathbf{y} \sim p_{\mathbf{y}}(\mathbf{y}|\mathbf{w})$$

Estimate  $\mathbf{w}$  ?

MMSE estimator :  $\hat{\mathbf{w}} = \mathbb{E}[\mathbf{w}|\mathbf{y}]$ , i.e.

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} p_{\mathbf{w}}(\mathbf{w}) p_{\mathbf{y}}(\mathbf{y}|\mathbf{w}) d\mu(\mathbf{w})$$

# The starting point : probabilistic inference

A hidden process generates  $\mathbf{w} \in \mathbb{R}^d$  with large  $d$

$$\mathbf{w} \sim p_{\mathbf{w}}(\mathbf{w})$$

We observe  $\mathbf{y} \in \mathbb{R}^n$  s.t.

$$\mathbf{y} \sim p_{\mathbf{y}}(\mathbf{y}|\mathbf{w})$$

Estimate  $\mathbf{w}$  ?

MMSE estimator :  $\hat{\mathbf{w}} = \mathbb{E}[\mathbf{w}|\mathbf{y}]$ , i.e.

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} p_{\mathbf{w}}(\mathbf{w}) p_{\mathbf{y}}(\mathbf{y}|\mathbf{w}) d\mu(\mathbf{w})$$

**Problem : This is a high-dimensional integral**

Typically  $p_{\mathbf{w}} \propto \exp(-\beta f(\mathbf{w}))$  and  $p_{\mathbf{y}} \propto \exp(-\beta g(\mathbf{w}, \mathbf{y}))$

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta (g(\mathbf{w}, \mathbf{y}) + f(\mathbf{w}))) d\mu(\mathbf{w})$$

Typically  $p_{\mathbf{w}} \propto \exp(-\beta f(\mathbf{w}))$  and  $p_{\mathbf{y}} \propto \exp(-\beta g(\mathbf{w}, \mathbf{y}))$

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta (g(\mathbf{w}, \mathbf{y}) + f(\mathbf{w}))) d\mu(\mathbf{w})$$

- Equilibrium Gibbs measure
- Hamiltonian  $\mathcal{H}(\mathbf{w}, \mathbf{y}) = g(\mathbf{w}, \mathbf{y}) + f(\mathbf{w})$
- $\beta$  is the inverse temperature

## Link with statistical physics : disordered systems

- distributions involve a dense, random interaction matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , here i.i.d.  $\mathcal{N}(0, \frac{1}{d})$  elements
- $\mathbf{y}$  can come from another random model, i.e.

$$\mathbf{y} \sim \mathbf{p}_{0,\mathbf{y}}(\mathbf{y}|\mathbf{X}, \mathbf{w}_0, \epsilon), \mathbf{w}_0 \sim p_{\mathbf{w}_0}(\mathbf{w}_0)$$

- estimate the generative model with postulated densities  $g, f$

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta (g(\mathbf{X}\mathbf{w}, \mathbf{y}) + f(\mathbf{w}))) d\mu(\mathbf{w})$$



## Link with statistical physics : disordered systems

- distributions involve a dense, random interaction matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , here i.i.d.  $\mathcal{N}(0, \frac{1}{d})$  elements
- $\mathbf{y}$  can come from another random model, i.e.

$$\mathbf{y} \sim \mathbf{p}_{0,\mathbf{y}}(\mathbf{y}|\mathbf{X}, \mathbf{w}_0, \epsilon), \mathbf{w}_0 \sim p_{\mathbf{w}_0}(\mathbf{w}_0)$$

- estimate the generative model with postulated densities  $g, f$

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta (g(\mathbf{X}\mathbf{w}, \mathbf{y}) + f(\mathbf{w}))) d\mu(\mathbf{w})$$

Typical setup

$$P(\mathbf{w}|\mathbf{y}, \mathbf{X}) = \frac{1}{\mathcal{Z}} \prod_{\mu=1}^n p_z(y_\mu|z_\mu) \prod_{i=1}^d p_w(w_i) \quad \text{where} \quad z_\mu = \sum_{i=1}^d X_{\mu i} w_i$$

# Back to machine learning and inference

Recall

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta (g(\mathbf{X}\mathbf{w}, \mathbf{y}) + f(\mathbf{w}))) d\mu(\mathbf{x})$$

For  $\beta \rightarrow +\infty$ , Laplace's method gives

$$\hat{\mathbf{w}} \xrightarrow{\beta \rightarrow +\infty} \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} g(\mathbf{X}\mathbf{w}, \mathbf{y}) + f(\mathbf{w})$$

**Empirical risk minimization with  $n$  samples in  $\mathbb{R}^d$**

Examples : LASSO, logistic regression, etc ...

# Back to machine learning and inference

Recall

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta (g(\mathbf{X}\mathbf{w}, \mathbf{y}) + f(\mathbf{w}))) d\mu(\mathbf{x})$$

For  $\beta \rightarrow +\infty$ , Laplace's method gives

$$\hat{\mathbf{w}} \xrightarrow{\beta \rightarrow +\infty} \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} g(\mathbf{X}\mathbf{w}, \mathbf{y}) + f(\mathbf{w})$$

**Empirical risk minimization with  $n$  samples in  $\mathbb{R}^d$**

Examples : LASSO, logistic regression, etc ...

**Goal : single letter formulas for the properties of  $\hat{\mathbf{w}}$  when  $n, d \rightarrow \infty$   
with aspect ratio  $\alpha \in (0, \infty)$**

## Belief propagation (BP) and AMP iterations

Consider the LASSO problem :  $\mathbf{X} \in \mathbb{R}^{n \times d}$  i.i.d.  $\mathcal{N}(0, 1/d)$ .

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda_1 \|\mathbf{w}\|_1 \right\}$$

where  $\mathbf{y} = \mathbf{X}\mathbf{w}_0 + \boldsymbol{\epsilon}_0$

# Belief propagation (BP) and AMP iterations

Consider the LASSO problem :  $\mathbf{X} \in \mathbb{R}^{n \times d}$  i.i.d.  $\mathcal{N}(0, 1/d)$ .

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda_1 \|\mathbf{w}\|_1 \right\}$$

where  $\mathbf{y} = \mathbf{X}\mathbf{w}_0 + \boldsymbol{\epsilon}_0$

- **Can we solve this optimization problem ?**  
(existing methods : subgradient, proximal point)
- **Theoretical guarantees of  $\mathbf{w}^*$  ?**

# Belief propagation (BP) and AMP iterations

Consider the LASSO problem :  $\mathbf{X} \in \mathbb{R}^{n \times d}$  i.i.d.  $\mathcal{N}(0, 1/d)$ .

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda_1 \|\mathbf{w}\|_1 \right\}$$

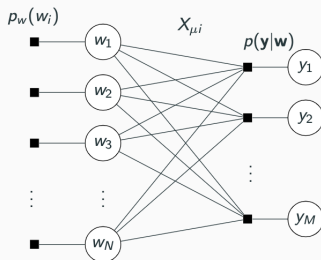
where  $\mathbf{y} = \mathbf{X}\mathbf{w}_0 + \boldsymbol{\epsilon}_0$

- **Can we solve this optimization problem ?**  
(existing methods : subgradient, proximal point)
- **Theoretical guarantees of  $\mathbf{w}^*$  ?**

**Can do both at the same time**

# Belief propagation (BP) and AMP iterations

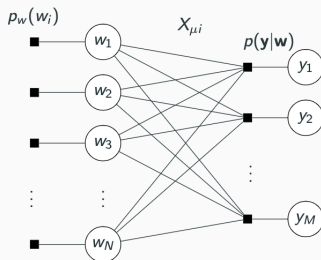
Factor graph representation of LASSO Gibbs measure



Factor graph for  $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$

# Belief propagation (BP) and AMP iterations

Factor graph representation of LASSO Gibbs measure



Factor graph for  $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$

Relaxation of BP equations + concentration for  $n, d \rightarrow +\infty$



TAP equations [Mézard, Parisi & Virasoro '87]

AMP iteration [Donoho et al. '09]



# AMP for the LASSO

LASSO problem.  $\mathbf{X} \in \mathbb{R}^{n \times d}$  i.i.d.  $\mathcal{N}(0, 1/d)$ .

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda_1 \|\mathbf{w}\|_1 \right\}$$

$$\mathbf{y} = \mathbf{X}\mathbf{w}_0 + \boldsymbol{\epsilon}_0$$

# AMP for the LASSO

LASSO problem.  $\mathbf{X} \in \mathbb{R}^{n \times d}$  i.i.d.  $\mathcal{N}(0, 1/d)$ .

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda_1 \|\mathbf{w}\|_1 \right\}$$
$$\mathbf{y} = \mathbf{X}\mathbf{w}_0 + \epsilon_0$$

**Approximate message-passing for the LASSO** [Donoho et al. '09, Bayati & Montanari '11]

$$\mathbf{z}^t = \mathbf{y} - \mathbf{X}\hat{\mathbf{w}}^t + \frac{1}{\alpha} \mathbf{z}^{t-1} \langle \eta'(\hat{\mathbf{w}}^{t-1} + \frac{1}{\alpha} \mathbf{X}^T \mathbf{z}^{t-1}, \theta^{t-1}) \rangle$$
$$\hat{\mathbf{w}}^{t+1} = \eta(\hat{\mathbf{w}}^t + \frac{1}{\alpha} \mathbf{X}^T \mathbf{z}^t, \theta^t)$$

- $\eta$  is the soft-thresholding operator (proximal of  $\ell_1$ )
- $\theta^t$  is a tunable parameter

**AMP** for LASSO

$$\begin{aligned}\mathbf{z}^t &= y - \mathbf{X}\hat{\mathbf{w}}^t + \frac{1}{\alpha}\mathbf{z}^{t-1}\langle\eta'(\hat{\mathbf{w}}^{t-1} + \frac{1}{\alpha}\mathbf{X}^T\mathbf{z}^{t-1}, \theta^{t-1})\rangle \\ \hat{\mathbf{w}}^{t+1} &= \eta(\hat{\mathbf{w}}^t + \frac{1}{\alpha}\mathbf{X}^T\mathbf{z}^t, \theta^t)\end{aligned}$$

## AMP for LASSO

$$\begin{aligned}\mathbf{z}^t &= y - \mathbf{X}\hat{\mathbf{w}}^t + \frac{1}{\alpha}\mathbf{z}^{t-1}\langle\eta'(\hat{\mathbf{w}}^{t-1} + \frac{1}{\alpha}\mathbf{X}^T\mathbf{z}^{t-1}, \theta^{t-1})\rangle \\ \hat{\mathbf{w}}^{t+1} &= \eta(\hat{\mathbf{w}}^t + \frac{1}{\alpha}\mathbf{X}^T\mathbf{z}^t, \theta^t)\end{aligned}$$

**State evolution for  $n, d \rightarrow \infty$ ,  $Z \sim \mathcal{N}(0, 1)$**  [Donoho et al. '09, Bayati & Montanari '11]

$$\begin{aligned}V &= \mathbb{E}_{\mathbf{z}, w_0} \{ [\eta'(W_0 + \sqrt{\frac{\Delta_0 + E}{\alpha}}Z; \theta(V))]^2 \} \\ E &= \mathbb{E}_{\mathbf{z}, w_0} \{ [\eta(W_0 + \sqrt{\frac{\Delta_0 + E}{\alpha}}Z; \theta(V)) - W_0]^2 \}\end{aligned}$$

Same result as the replica computation [Krzakala et al. '12]

# AMP and state evolution equations

## Generic AMP iteration

- $\mathbf{A} \in \mathbb{R}^{n \times d}$  Gaussian random matrix with i.i.d. entries  $A_{ij} \sim \mathcal{N}(0, \frac{1}{d})$
- $e_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,  $g_t : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , pseudo-Lipschitz functions

$$\begin{aligned}\mathbf{u}^{t+1} &= \mathbf{A}^\top g_t(\mathbf{v}^t) - d_t e_t(\mathbf{u}^t) \\ \mathbf{v}^t &= \mathbf{A} e_t(\mathbf{u}^t) - b_t g_{t-1}(\mathbf{v}^{t-1})\end{aligned}$$

where

$$\begin{aligned}d_t &= \frac{1}{m} \operatorname{div}(g_t(\mathbf{v}^t)) \\ b_t &= \frac{1}{m} \operatorname{div}(e_t(\mathbf{u}^t))\end{aligned}$$

## AMP and state evolution equations

$$\begin{aligned}\mathbf{u}^{t+1} &= \mathbf{A}^\top g_t(\mathbf{v}^t) - d_t e_t(\mathbf{u}^t) \\ \mathbf{v}^t &= \mathbf{A} e_t(\mathbf{u}^t) - b_t g_{t-1}(\mathbf{v}^{t-1})\end{aligned}$$

Define the recursion

$$\begin{aligned}\tau_{t+1} &= \mathbb{E} [g_t^2(Z_\sigma^t)] & Z_\sigma^t &\sim \mathcal{N}(0, \sigma_t) \\ \sigma_t &= \mathbb{E} [e_t^2(Z_\tau^t)] & Z_\tau^t &\sim \mathcal{N}(0, \tau_t)\end{aligned}$$

# AMP and state evolution equations

$$\begin{aligned}\mathbf{u}^{t+1} &= \mathbf{A}^\top g_t(\mathbf{v}^t) - d_t e_t(\mathbf{u}^t) \\ \mathbf{v}^t &= \mathbf{A} e_t(\mathbf{u}^t) - b_t g_{t-1}(\mathbf{v}^{t-1})\end{aligned}$$

Define the recursion

$$\begin{aligned}\tau_{t+1} &= \mathbb{E} [g_t^2(Z_\sigma^t)] & Z_\sigma^t &\sim \mathcal{N}(0, \sigma_t) \\ \sigma_t &= \mathbb{E} [e_t^2(Z_\tau^t)] & Z_\tau^t &\sim \mathcal{N}(0, \tau_t)\end{aligned}$$

## Theorem (Bayati & Montanari 2011)

For any pseudo-Lipschitz function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$ ,

$$\begin{aligned}\frac{1}{d} \sum_{i=1}^d \phi(u_i^t) &\xrightarrow[n, d \rightarrow \infty]{a.s.} \mathbb{E} [\phi(Z_\tau^t)] \\ \frac{1}{n} \sum_{i=1}^n \phi(v_i^t) &\xrightarrow[n, d \rightarrow \infty]{a.s.} \mathbb{E} [\phi(Z_\sigma^t)]\end{aligned}$$

# AMP and state evolution equations

Recall inference problem :

- recover  $\mathbf{w}_0$  from observations  $\mathbf{y} = \phi(\mathbf{A}\mathbf{w}_0)$
- prior  $p_{\mathbf{w}_0}$

$$\begin{aligned}\mathbf{u}^{t+1} &= \mathbf{A}^\top g_t(\mathbf{v}^t) - d_t e_t(\mathbf{u}^t) \\ \mathbf{v}^t &= \mathbf{A} e_t(\mathbf{u}^t) - b_t g_{t-1}(\mathbf{v}^{t-1})\end{aligned}$$



# AMP and state evolution equations

Recall inference problem :

- recover  $\mathbf{w}_0$  from observations  $\mathbf{y} = \phi(\mathbf{A}\mathbf{w}_0)$
- prior  $p_{\mathbf{w}_0}$

$$\begin{aligned}\mathbf{u}^{t+1} &= \mathbf{A}^\top g_t(\mathbf{v}^t) - d_t e_t(\mathbf{u}^t) \\ \mathbf{v}^t &= \mathbf{A} e_t(\mathbf{u}^t) - b_t g_{t-1}(\mathbf{v}^{t-1})\end{aligned}$$

- function  $e_t$  estimates  $\mathbf{w}_0$  : denoising  $p_{\mathbf{w}_0}$  blurred with additive Gaussian noise
- function  $g_t$  estimates  $\mathbf{A}\mathbf{w}_0$  : same for  $p(\mathbf{y}|\mathbf{A}\mathbf{w}_0)$

# Composing AMP iterations

**Multilayer AMP** [Mézard '17, Manoel et al. '17] :

Random neural networks, generative models, structured matrices ...

Used in other works [Gabrié et al. '19, Aubin et al. '20]

Recover  $\mathbf{w}_0 \in \mathbb{R}^{N_1}$ , prior  $p_{\mathbf{w}_0}$ ,  $\mathbf{A}_l \in \mathbb{R}^{N_l \times N_{l-1}}$

$$\mathbf{y} = \phi_L(\mathbf{A}_L \phi_{L-1}(\mathbf{A}_{L-1}(\dots \phi_1(\mathbf{A}_1 \mathbf{w}_0))))$$

Standard AMP iteration at each layer  $l$

# Composing AMP iterations

**Multilayer AMP** [Mézard '17, Manoel et al. '17] :

Random neural networks, generative models, structured matrices ...

Used in other works [Gabrié et al. '19, Aubin et al. '20]

Recover  $\mathbf{w}_0 \in \mathbb{R}^{N_1}$ , prior  $p_{\mathbf{w}_0}$ ,  $\mathbf{A}_l \in \mathbb{R}^{N_l \times N_{l-1}}$

$$\mathbf{y} = \phi_L(\mathbf{A}_L \phi_{L-1}(\mathbf{A}_{L-1}(\dots \phi_1(\mathbf{A}_1 \mathbf{w}_0))))$$

Standard AMP iteration at each layer  $l$

$$\begin{aligned}\mathbf{u}_l^{t+1} &= \mathbf{A}_l^\top g_{l,t}(\mathbf{v}_l^t) - d_{l,t} e_{l,t}(\mathbf{u}_l^t) \\ \mathbf{v}_l^t &= \mathbf{A}_l e_{l,t}(\mathbf{u}_l^t) - b_{l,t} g_{l,t-1}(\mathbf{v}_l^{t-1})\end{aligned}$$

Layerwise state evolution equations

# Composing AMP iterations

**Spiked matrix with generative prior** [Aubin et al. '19]

Recover  $\mathbf{v}_0$  from

$$\mathbf{Y} = \sqrt{\frac{\lambda}{N}} \mathbf{v}_0 \mathbf{v}_0^\top + \mathbf{W} \quad \mathbf{W} \sim \text{GOE}(N)$$

with  $\mathbf{v}_0$  generated from

$$\mathbf{v}_0 = \phi_L(\mathbf{A}_L \phi_{L-1}(\mathbf{A}_{L-1}(\dots \phi_1(\mathbf{A}_1 \mathbf{w}_0))))$$

# Composing AMP iterations

**Spiked matrix with generative prior** [Aubin et al. '19]

Recover  $\mathbf{v}_0$  from

$$\mathbf{Y} = \sqrt{\frac{\lambda}{N}} \mathbf{v}_0 \mathbf{v}_0^\top + \mathbf{W} \quad \mathbf{W} \sim \text{GOE}(N)$$

with  $\mathbf{v}_0$  generated from

$$\mathbf{v}_0 = \phi_L(\mathbf{A}_L \phi_{L-1}(\mathbf{A}_{L-1}(\dots \phi_1(\mathbf{A}_1 \mathbf{w}_0))))$$

AMP for the spike part

$$\mathbf{x}^{t+1} = \mathbf{W} f^t(\mathbf{x}^t) - b^t f^{t-1}(\mathbf{x}^{t-1})$$

Combined with MLAMP, also admits layerwise state evolution equations.

# Composition preserves state evolution property

State evolution property is preserved when composing AMP iterations



Common structure in AMP iterations ? Can we use this structure to give a modular proof of SE equations ?

State evolution property is preserved when composing AMP iterations



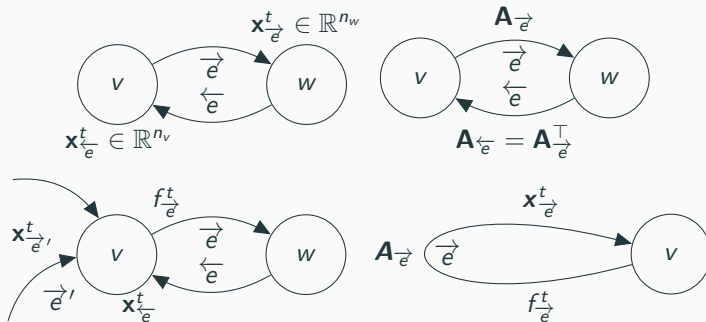
Common structure in AMP iterations ? Can we use this structure to give a modular proof of SE equations ?

## Proposed Solution

- indexation of AMP iterations on an oriented graph
- modular proof of SE equations based on this graph

# Graph-based AMP iterations : the oriented graph

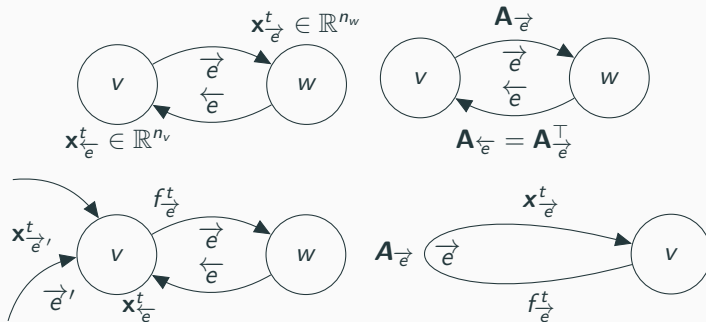
Variables  $\mathbf{x}_{\vec{e}}^t \in \mathbb{R}^{n_e}$ , random matrices  $\mathbf{A}_{\vec{e}}$ , non-linearities  $f_{\vec{e}}^t$





# Graph-based AMP iterations : the oriented graph

Variables  $\mathbf{x}_{\vec{e}}^t$ , random matrices  $\mathbf{A}_{\vec{e}}$ , non-linearities  $f_{\vec{e}}^t$



Arbitrary composition of this structure

## Graph-based AMP : the iteration

The graph-based AMP iteration reads

$$\begin{aligned}\mathbf{x}_{\vec{e}}^{t+1} &= \mathbf{A}_{\vec{e}} \mathbf{m}_{\vec{e}}^t - b_{\vec{e}}^t \mathbf{m}_{\vec{e}}^{t-1}, \\ \mathbf{m}_{\vec{e}}^t &= f_{\vec{e}}^t \left( (\mathbf{x}_{\vec{e}'}^t)_{\vec{e}': \vec{e}' \rightarrow \vec{e}} \right),\end{aligned}$$

where  $b_{\vec{e}}^t$  is the *Onsager term*

$$b_{\vec{e}}^t = \frac{1}{N} \text{Tr} \frac{\partial f_{\vec{e}}^t}{\partial \mathbf{x}_{\vec{e}}^t} \left( (\mathbf{x}_{\vec{e}'}^t)_{\vec{e}': \vec{e}' \rightarrow \vec{e}} \right) \in \mathbb{R}.$$

# Graph-based AMP : the iteration

The graph-based AMP iteration reads

$$\begin{aligned}\mathbf{x}_{\vec{e}}^{t+1} &= \mathbf{A}_{\vec{e}} \mathbf{m}_{\vec{e}}^t - b_{\vec{e}}^t \mathbf{m}_{\vec{e}}^{t-1}, \\ \mathbf{m}_{\vec{e}}^t &= f_{\vec{e}}^t \left( (\mathbf{x}_{\vec{e}'}^t)_{\vec{e}': \vec{e}' \rightarrow \vec{e}} \right),\end{aligned}$$

where  $b_{\vec{e}}^t$  is the *Onsager term*

$$b_{\vec{e}}^t = \frac{1}{N} \text{Tr} \frac{\partial f_{\vec{e}}^t}{\partial \mathbf{x}_{\vec{e}}^t} \left( (\mathbf{x}_{\vec{e}'}^t)_{\vec{e}': \vec{e}' \rightarrow \vec{e}} \right) \in \mathbb{R}.$$

**Theorem (informal):**

- any Graph-based AMP iterations admits rigorous SE
- equations can be deduced from the graph

# Graph-based AMP : the state evolution equations

## Definition (State evolution iterates)

Define independently for each  $\vec{e} \in \vec{E}$ ,  $\mathbf{Z}_{\vec{e}}^0 = \mathbf{x}_{\vec{e}}^0$  and  $(\mathbf{Z}_{\vec{e}}^1, \dots, \mathbf{Z}_{\vec{e}}^t)$  a centered Gaussian random vector of covariance  $(\kappa_{\vec{e}}^{r,s})_{r,s \leq t} \otimes I_{n_w}$ . We then define new state evolution iterates

$$\kappa_{\vec{e}}^{t+1,s+1} = \lim_{n \rightarrow \infty} \frac{1}{N} \mathbb{E} \left[ \langle f_{\vec{e}}^s \left( (\mathbf{Z}_{\vec{e}'}^s)_{\vec{e}': \vec{e}' \rightarrow \vec{e}} \right), f_{\vec{e}}^t \left( (\mathbf{Z}_{\vec{e}'}^t)_{\vec{e}': \vec{e}' \rightarrow \vec{e}} \right) \rangle \right]$$

for all  $s \in \{1, \dots, t\}$ ,  $\vec{e} \in \vec{E}$ .

## Theorem (Gerbelot & Berthier '21)

*Under regularity assumptions, for any sequence of uniformly (in  $n$ ) pseudo-Lipschitz function  $\Phi : \mathbb{R}^{(t+1)N} \rightarrow \mathbb{R}$ ,*

$$\Phi \left( (\mathbf{x}_{\vec{e}}^s)_{0 \leq s \leq t, \vec{e} \in \vec{E}} \right) \stackrel{P}{\simeq} \mathbb{E} \left[ \Phi \left( (\mathbf{Z}_{\vec{e}}^s)_{0 \leq s \leq t, \vec{e} \in \vec{E}} \right) \right]$$

# Graph-based AMP : recovering known examples

**Symmetric AMP** : spiked matrix recovery, SK model [Rangan et al. '12, Javanmard & Montanari '12, Deshpande & Montanari '14, Bolthausen '14]

Recover  $\mathbf{v}_0 \in \mathbb{R}^N$  from:

$$\mathbf{Y} = \sqrt{\frac{\lambda}{N}} \mathbf{v}_0 \mathbf{v}_0^\top + \mathbf{W}$$

where the noise matrix  $\mathbf{W} \in GOE(N)$ .

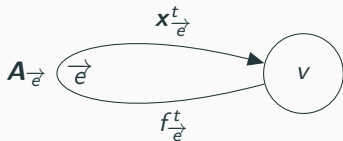
# Graph-based AMP : recovering known examples

**Symmetric AMP** : spiked matrix recovery, SK model [Rangan et al. '12, Javanmard & Montanari '12, Deshpande & Montanari '14, Bolthausen '14]

Recover  $\mathbf{v}_0 \in \mathbb{R}^N$  from:

$$\mathbf{Y} = \sqrt{\frac{\lambda}{N}} \mathbf{v}_0 \mathbf{v}_0^\top + \mathbf{W}$$

where the noise matrix  $\mathbf{W} \in GOE(N)$ .



$$\begin{aligned} \mathbf{x}_{\vec{e}}^{t+1} &= \mathbf{A}_{\vec{e}} \mathbf{m}_{\vec{e}}^t - b_{\vec{e}}^t \mathbf{m}_{\vec{e}}^{t-1}, \\ \mathbf{m}_{\vec{e}}^t &= f_{\vec{e}}^t(\mathbf{x}_{\vec{e}}^t), \end{aligned}$$

# Graph-based AMP : recovering known examples

**Asymmetric AMP** : LASSO, GLM, ... [Donoho et al. '09, Bayati & Montanari '11, Rangan '11]

Recover  $\mathbf{x}_0 \in \mathbb{R}^N$  from

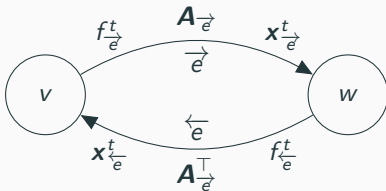
$$\mathbf{y} = \phi(\mathbf{A}\mathbf{x}_0, \epsilon_0)$$

# Graph-based AMP : recovering known examples

**Asymmetric AMP** : LASSO, GLM, ... [Donoho et al. '09, Bayati & Montanari '11, Rangan '11]

Recover  $\mathbf{x}_0 \in \mathbb{R}^N$  from

$$\mathbf{y} = \phi(\mathbf{A}\mathbf{x}_0, \epsilon_0)$$



$$\mathbf{x}_{\vec{e}}^{t+1} = \mathbf{A}_{\vec{e}} \mathbf{m}_{\vec{e}}^t - b_{\vec{e}}^t \mathbf{m}_{\vec{e}}^{t-1},$$

$$\mathbf{m}_{\vec{e}}^t = f_{\vec{e}}^t(\mathbf{x}_{\vec{e}}^t),$$

$$\mathbf{x}_{\overleftarrow{e}}^{t+1} = \mathbf{A}_{\overleftarrow{e}}^T \mathbf{m}_{\overleftarrow{e}}^t - b_{\overleftarrow{e}}^t \mathbf{m}_{\overleftarrow{e}}^{t-1},$$

$$\mathbf{m}_{\overleftarrow{e}}^t = f_{\overleftarrow{e}}^t(\mathbf{x}_{\overleftarrow{e}}^t).$$



# Graph-based AMP iterations : proving heuristic SE equations

**Multilayer AMP** [Manoel et al. '17] :

Random neural networks, generative models, structured matrices ...

Used in other works [Gabrie et al. '19, Aubin et al. '20]

Recover  $\mathbf{x}_0 \in \mathbb{R}^{N_1}$

$$\mathbf{y} = \phi_L(\mathbf{A}_L \phi_{L-1}(\mathbf{A}_{L-1}(\dots \phi_1(\mathbf{A}_1 \mathbf{x}_0))))$$

# Graph-based AMP iterations : proving heuristic SE equations

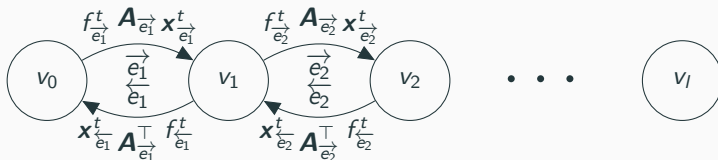
**Multilayer AMP** [Manoel et al. '17] :

Random neural networks, generative models, structured matrices ...

Used in other works [Gabrie et al. '19, Aubin et al. '20]

Recover  $\mathbf{x}_0 \in \mathbb{R}^{N_1}$

$$\mathbf{y} = \phi_L(\mathbf{A}_L \phi_{L-1}(\mathbf{A}_{L-1}(\dots \phi_1(\mathbf{A}_1 \mathbf{x}_0))))$$



# Graph-based AMP : proving heuristic SE equations

**Spiked matrix with generative prior** [Aubin et al. '19]

Recover  $\mathbf{v}_0$  from

$$\mathbf{Y} = \sqrt{\frac{\lambda}{N}} \mathbf{v}_0 \mathbf{v}_0^\top + \mathbf{A}_0$$

with  $\mathbf{v}_0$  generated from

$$\mathbf{v}_0 = \phi_L(\mathbf{A}_L \phi_{L-1}(\mathbf{A}_{L-1}(\dots \phi_1(\mathbf{A}_1 \mathbf{w}_0))))$$

# Graph-based AMP : proving heuristic SE equations

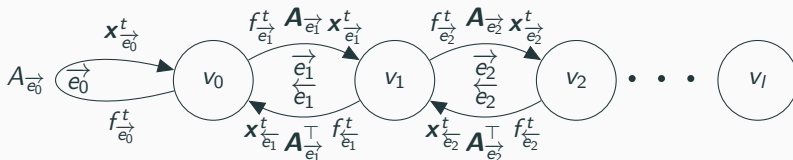
Spiked matrix with generative prior [Aubin et al. '19]

Recover  $\mathbf{v}_0$  from

$$\mathbf{Y} = \sqrt{\frac{\lambda}{N}} \mathbf{v}_0 \mathbf{v}_0^\top + \mathbf{A}_0$$

with  $\mathbf{v}_0$  generated from

$$\mathbf{v}_0 = \phi_L(\mathbf{A}_L \phi_{L-1}(\mathbf{A}_{L-1}(\dots \phi_1(\mathbf{A}_1 \mathbf{w}_0))))$$



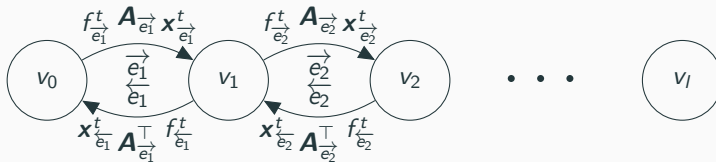
# Graph-based AMP iterations : proving new SE equations

Convolutional multilayer generalized linear estimation [Gerbelot et al. 22]

Recover  $\mathbf{x}_0 \in \mathbb{R}^{N_1}$  from

$$\mathbf{y} = \phi_L(\mathbf{A}_L \phi_{L-1}(\mathbf{A}_{L-1}(\dots \phi_1(\mathbf{A}_1 \mathbf{x}_0))))$$

Random sparse circulant  $\mathbf{A}_l \in \mathbb{R}^{N_{l+1} \times N_l}$



# Proof of Graph-AMP theorem

$\mathbf{A} \sim \text{GOE}(N)$ , then

$$\mathbf{x}^{t+1} = \mathbf{A}\mathbf{m}^t - b_t\mathbf{m}^{t-1}$$

$$\mathbf{m}^t = f_t(\mathbf{x}^t)$$

with initialization at  $\mathbf{x}^0$  and Onsager correction

$$b_t = \text{div} [f_t(\mathbf{x}^t)]$$

# Proof of Graph-AMP theorem

$\mathbf{A} \sim \text{GOE}(N)$ , then

$$\mathbf{x}^{t+1} = \mathbf{A}\mathbf{m}^t - b_t\mathbf{m}^{t-1}$$

$$\mathbf{m}^t = f_t(\mathbf{x}^t)$$

with initialization at  $\mathbf{x}^0$  and Onsager correction

$$b_t = \text{div} [f_t(\mathbf{x}^t)]$$

State evolution : for any  $t$ ,  $\mathbf{x}^t$  behaves as  $\mathbf{Z}^t \sim \mathcal{N}(0, \kappa_{t,t} \mathbf{I}_N)$

$$\text{where } \kappa_{t+1} = \mathbb{E} [(f^t(z^t))^2], \quad z^t \sim \mathcal{N}(0, \kappa_t)$$

Proof idea due to E. Bolthausen '09, '14

## Sketch of proof : Bolthausen conditioning

$$\begin{aligned}\mathbf{x}^{t+1} &= \mathbf{A}\mathbf{m}^t - b_t\mathbf{m}^{t-1} \\ \mathbf{m}^t &= f_t(\mathbf{x}^t)\end{aligned}$$

Define the  $\sigma$ -algebra  $\mathfrak{S}_t = \sigma(\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^t)$ . We then have :

$$\mathbf{x}^{t+1}|_{\mathfrak{S}_t} = \mathbf{A}|_{\mathfrak{S}_t}\mathbf{m}^t - b_t\mathbf{m}^{t-1}$$



## Sketch of proof : Bolthausen conditioning

$$\begin{aligned}\mathbf{x}^{t+1} &= \mathbf{A}\mathbf{m}^t - b_t\mathbf{m}^{t-1} \\ \mathbf{m}^t &= f_t(\mathbf{x}^t)\end{aligned}$$

Define the  $\sigma$ -algebra  $\mathfrak{G}_t = \sigma(\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^t)$ . We then have :

$$\mathbf{x}^{t+1}|_{\mathfrak{G}_t} = \mathbf{A}|_{\mathfrak{G}_t}\mathbf{m}^t - b_t\mathbf{m}^{t-1}$$

Gaussian conditioning lemma

$$\begin{aligned}\mathbf{A}|_{\mathfrak{G}_t} &= \mathbb{E}[\mathbf{A}|\mathfrak{G}_t] + \mathcal{P}_t(\mathbf{A}) \\ &= \mathbf{A} - \mathbf{P}_{\mathbf{M}_{t-1}}^\perp \mathbf{A} \mathbf{P}_{\mathbf{M}_{t-1}}^\perp + \mathbf{P}_{\mathbf{M}_{t-1}}^\perp \tilde{\mathbf{A}} \mathbf{P}_{\mathbf{M}_{t-1}}^\perp\end{aligned}$$

where  $\mathbf{M}_{t-1} = [m^0 | \dots | m^{t-1}]$  and  $\tilde{\mathbf{A}}$  is an independent copy of  $\mathbf{A}$ .

## Sketch of proof : Bolthausen conditioning

A bit of algebra leads to

$$\mathbf{x}_{|\mathfrak{S}_t}^{t+1} = \underbrace{\mathbf{X}_{t-1}\alpha_t + \mathbf{P}_{\mathbf{M}_{t-1}}^\perp \tilde{\mathbf{A}} \mathbf{P}_{\mathbf{M}_{t-1}}^\perp \mathbf{m}^t}_{\text{Part 1}} + \underbrace{[0|\mathbf{M}_{t-2}] \mathbf{B}_t \alpha_t + \mathbf{P}_{\mathbf{M}_{t-1}} \mathbf{A} \mathbf{m}_\perp^t - b_t \mathbf{m}^{t-1}}_{\text{Part 2}}$$

## Sketch of proof : Bolthausen conditioning

A bit of algebra leads to

$$\mathbf{x}_{|\mathfrak{S}_t}^{t+1} = \underbrace{\mathbf{X}_{t-1}\alpha_t + \mathbf{P}_{\mathbf{M}_{t-1}}^\perp \tilde{\mathbf{A}} \mathbf{P}_{\mathbf{M}_{t-1}}^\perp \mathbf{m}^t}_{\text{Part 1}} + \underbrace{[0|\mathbf{M}_{t-2}] \mathbf{B}_t \alpha_t + \mathbf{P}_{\mathbf{M}_{t-1}} \mathbf{A} \mathbf{m}_\perp^t - b_t \mathbf{m}^{t-1}}_{\text{Part 2}}$$

- Part 1 concentrates (induction+Gaussian concentration)
- Part 2 goes to zero w.h.p. as  $N \rightarrow \infty$
- Onsager correction  $b_t$  cancels the bothersome part

## Sketch of proof : embedding of the graph

Embed the graph into a large, matrix valued, non-separable iteration of the form

$$\mathbf{X}^{t+1} = \mathbf{A}\mathbf{M}^t - \mathbf{M}^{t-1}\mathbf{b}_t^\top.$$

## Sketch of proof : embedding of the graph

Embed the graph into a large, matrix valued, non-separable iteration of the form

$$\mathbf{X}^{t+1} = \mathbf{A}\mathbf{M}^t - \mathbf{M}^{t-1}\mathbf{b}_t^\top.$$

where

$$\begin{aligned} \mathbf{A}\mathbf{M}^t &= \begin{pmatrix} \mathbf{A}_{\vec{e}_1} & & & & \\ & \ddots & & & \\ & & \mathbf{A}_{\vec{e}_l} & & \\ & & & * & \mathbf{A}_{\vec{e}_{l+1}} \\ & & & \mathbf{A}_{\vec{e}_{l+1}}^* & \\ & & & & \ddots \end{pmatrix} \begin{pmatrix} f_{\vec{e}_1}^t(\cdot) & & & & \\ & \ddots & & & \\ & & f_{\vec{e}_l}^t(\cdot) & & \\ & & & 0 & f_{\vec{e}_{l+1}}^t(\cdot) \\ & & & f_{\vec{e}_{l+1}}^t(\cdot) & 0 \\ & & & & \ddots \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{A}_{\vec{e}_1} f_{\vec{e}_1}^t \left( (x_{\vec{e}}^t)_{\vec{e}:\vec{e} \rightarrow \vec{e}_1} \right) & & & & \\ & \ddots & & & \\ & & \mathbf{A}_{\vec{e}_l} f_{\vec{e}_l}^t(\cdot) & & * \\ & & & \mathbf{A}_{\vec{e}_{l+1}} f_{\vec{e}_{l+1}}^t(\cdot) & \\ & & * & & \mathbf{A}_{\vec{e}_{l+1}} f_{\vec{e}_{l+1}}^t(\cdot) \\ & & & & \ddots \end{pmatrix}. \end{aligned}$$

and,

$$\begin{aligned}
 M^{t-1} \mathbf{b}_t &= \begin{pmatrix} f_{\vec{e}_1}^{t-1}(\cdot) & & & & \\ & \ddots & & & \\ & & f_{\vec{e}_l}^{t-1}(\cdot) & & \\ & & & 0 & f_{\vec{e}_{l+1}}^{t-1}(\cdot) \\ & & & f_{\vec{e}_{l+1}}^{t-1}(\cdot) & 0 \\ & & \ddots & & \ddots \end{pmatrix} \begin{pmatrix} \mathbf{b}_{\vec{e}_1}^t \\ & \ddots & & & \\ & & \mathbf{b}_{\vec{e}_l}^t & & \\ & & & 0 & \mathbf{b}_{\vec{e}_{l+1}}^t \\ & & & \mathbf{b}_{\vec{e}_{l+1}}^t & 0 \\ & & \ddots & & \ddots \end{pmatrix} \\
 &= \begin{pmatrix} \mathbf{b}_{\vec{e}_1}^t f_{\vec{e}_1}^{t-1} \left( (x_{\vec{e}}^t)_{\vec{e}:\vec{e} \rightarrow \vec{e}_1} \right) & & & & \\ & \ddots & & & \\ & & \mathbf{b}_{\vec{e}_l}^t f_{\vec{e}_l}^{t-1}(\cdot) & & \\ & & & 0 & \mathbf{b}_{\vec{e}_{l+1}}^t f_{\vec{e}_{l+1}}^{t-1}(\cdot) \\ & & & \mathbf{b}_{\vec{e}_{l+1}}^t f_{\vec{e}_{l+1}}^{t-1}(\cdot) & \\ & & & & \mathbf{b}_{\vec{e}_{l+1}}^t f_{\vec{e}_{l+1}}^{t-1}(\cdot) \\ & & & & \ddots \end{pmatrix}.
 \end{aligned}$$

## Low rank perturbations

Define the matrix

$$\hat{\mathbf{A}} = \mathbf{A} + \frac{1}{N} \mathbf{V}_0 \mathbf{V}_0^\top \in \mathbb{R}^{N \times N},$$

# Low rank perturbations

Define the matrix

$$\hat{\mathbf{A}} = \mathbf{A} + \frac{1}{N} \mathbf{V}_0 \mathbf{V}_0^\top \in \mathbb{R}^{N \times N},$$

AMP iteration

$$\mathbf{X}^{t+1} = \hat{\mathbf{A}} \mathbf{M}^t - \mathbf{M}^{t-1} (\mathbf{b}^t)^\top \in \mathbb{R}^{N \times q},$$

$$\mathbf{M}^t = f^t(\mathbf{X}^t) \in \mathbb{R}^{N \times q},$$

$$\mathbf{b}^t = \frac{1}{N} \sum_{i=1}^N \frac{\partial f_i^t}{\partial \mathbf{x}_i}(\mathbf{x}^t) \in \mathbb{R}^{q \times q}.$$



# Low rank perturbations

Define the matrix

$$\hat{\mathbf{A}} = \mathbf{A} + \frac{1}{N} \mathbf{V}_0 \mathbf{V}_0^\top \in \mathbb{R}^{N \times N},$$

AMP iteration

$$\mathbf{X}^{t+1} = \hat{\mathbf{A}} \mathbf{M}^t - \mathbf{M}^{t-1} (\mathbf{b}^t)^\top \in \mathbb{R}^{N \times q},$$

$$\mathbf{M}^t = f^t(\mathbf{X}^t) \in \mathbb{R}^{N \times q},$$

$$\mathbf{b}^t = \frac{1}{N} \sum_{i=1}^N \frac{\partial f_i^t}{\partial \mathbf{X}_i}(\mathbf{X}^t) \in \mathbb{R}^{q \times q}.$$

State evolution recursion, initialized with  $\boldsymbol{\mu}_0 = \mathbf{0}_{q \times q}$ ,

$$\boldsymbol{\kappa}^{1,1} = \lim_{N \rightarrow \infty} \frac{1}{N} f^0(\mathbf{V}_0 \boldsymbol{\mu}_0 + \mathbf{X}^0)^\top f^0(\mathbf{V}_0 \boldsymbol{\mu}_0 + \mathbf{X}^0)$$

$$\boldsymbol{\mu}^{s+1} = \lim_{N \rightarrow +\infty} \frac{1}{N} \mathbb{E} [(\mathbf{V}_0)^\top f^s(\mathbf{V}_0 \boldsymbol{\mu}^s + \mathbf{Z}^s)]$$

$$\boldsymbol{\kappa}^{t+1,s+1} = \boldsymbol{\kappa}^{s+1,t+1} = \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E} [f^s(\mathbf{V}_0 \boldsymbol{\mu}^s + \mathbf{Z}^s)^\top f^t(\mathbf{V}_0 \boldsymbol{\mu}^t + \mathbf{Z}^t)].$$

## Sketch of proof

Similar proof with a single node in [Deshpande, Abbe & Montanari '17].  
Define

$$\begin{aligned}\mathbf{S}^{t+1} &= \mathbf{A}\tilde{\mathbf{M}}^t - \tilde{\mathbf{M}}^{t-1}(\tilde{\mathbf{b}}^t)^\top && \in \mathbb{R}^{N \times q}, \\ \tilde{\mathbf{M}}^t &= f^t(\mathbf{V}_0\boldsymbol{\mu}^t + \mathbf{S}^t) && \in \mathbb{R}^{N \times q}, \\ \tilde{\mathbf{b}}^t &= \frac{1}{N} \sum_{i=1}^N \frac{\partial f_i^t}{\partial \mathbf{S}_i}(\mathbf{V}_0\boldsymbol{\mu}^t + \mathbf{S}^t) && \in \mathbb{R}^{q \times q}.\end{aligned}$$

Has the structure required by main theorem. Then prove

$$\forall t \in \mathbb{N} \quad \frac{1}{\sqrt{N}} \|\mathbf{X}^t - \mathbf{S}^t - \mathbf{V}_0\boldsymbol{\mu}^t\|_F \xrightarrow[N \rightarrow \infty]{P} 0$$

by induction.

## Dependence on linear observations

$$\mathbf{X}^{t+1} = \mathbf{A}\mathbf{M}^t - \mathbf{M}^{t-1}(\mathbf{b}^t)^\top \quad \in \mathbb{R}^{N \times q},$$

$$\mathbf{M}^t = f^t(\varphi(\mathbf{A}\mathbf{W}_0), \mathbf{X}^t) \quad \in \mathbb{R}^{N \times q},$$

$$\mathbf{b}^t = \frac{1}{N} \sum_{i=1}^N \frac{\partial f_i^t}{\partial \mathbf{X}_i}(\varphi(\mathbf{A}\mathbf{W}_0), \mathbf{X}^t) \quad \in \mathbb{R}^{q \times q}.$$

# Dependence on linear observations

$$\mathbf{X}^{t+1} = \mathbf{A}\mathbf{M}^t - \mathbf{M}^{t-1}(\mathbf{b}^t)^\top \in \mathbb{R}^{N \times q},$$

$$\mathbf{M}^t = f^t(\varphi(\mathbf{A}\mathbf{W}_0), \mathbf{X}^t) \in \mathbb{R}^{N \times q},$$

$$\mathbf{b}^t = \frac{1}{N} \sum_{i=1}^N \frac{\partial f_i^t}{\partial \mathbf{X}_i}(\varphi(\mathbf{A}\mathbf{W}_0), \mathbf{X}^t) \in \mathbb{R}^{q \times q}.$$

State evolution

$$\boldsymbol{\nu}^{t+1} = \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E} \left[ \mathbf{W}_0^\top f^t \left( \varphi(\mathbf{Z}\mathbf{W}_0), \mathbf{Z}\mathbf{W}_0 \rho_{\mathbf{W}_0}^{-1} \boldsymbol{\nu}^t + \mathbf{W}_0 \hat{\boldsymbol{\nu}}^t + \mathbf{Z}^t \right) \right]$$

$$\hat{\boldsymbol{\nu}}^{t+1} = \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E} \left[ \sum_{i=1}^N \frac{\partial f_i^t}{\partial \mathbf{Z}\mathbf{W}_{0,i}, \varphi} \left( \varphi(\mathbf{Z}\mathbf{W}_0), \mathbf{Z}\mathbf{W}_0 \rho_{\mathbf{W}_0}^{-1} \boldsymbol{\nu}^t + \mathbf{W}_0 \hat{\boldsymbol{\nu}}^t + \mathbf{Z}^t \right) \right]$$

$$\boldsymbol{\kappa}^{t+1,s+1} =$$

$$\lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E} \left[ \left( f^s \left( \varphi(\mathbf{Z}\mathbf{W}_0), \mathbf{Z}\mathbf{W}_0 \rho_{\mathbf{W}_0}^{-1} \boldsymbol{\nu}^s + \mathbf{W}_0 \hat{\boldsymbol{\nu}}^s + \mathbf{Z}^s \right) - \mathbf{W}_0 \rho_{\mathbf{W}_0}^{-1} \boldsymbol{\nu}^{s+1} \right)^\top \right. \\ \left. \left( f^t \left( \varphi(\mathbf{Z}\mathbf{W}_0), \mathbf{Z}\mathbf{W}_0 \rho_{\mathbf{W}_0}^{-1} \boldsymbol{\nu}^t + \mathbf{W}_0 \hat{\boldsymbol{\nu}}^t + \mathbf{Z}^t \right) - \mathbf{W}_0 \rho_{\mathbf{W}_0}^{-1} \boldsymbol{\nu}^{t+1} \right) \right]$$

- finite size rates [Rush & Venkataramanan '18]
- subGaussian universality [Bayati, Lelarge & Montanari '15]
- rotationally invariant matrices [Rangan, Schniter & Fletcher '16],  
semi-random universality [Dudeja, Lu & Sen '22]

# An application of non-separable non-linearities

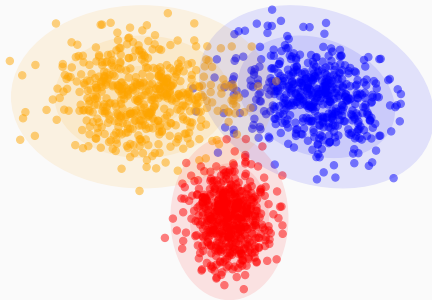
Classifying a high dimensional, non-isotropic Gaussian mixture

joint work with Loureiro. B., Sicuro. G., Pacco. A., Krzakala. F., Zdeborová.  
L. *Neurips 2021*

# Classifying Gaussian Mixtures with Convex GLM

Data and teacher

$$\mathbf{x} \in \mathbb{R}^d, \mathbf{y} \in \mathbb{R}^K \quad P(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^K y_k \rho_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k),$$



$K=3, d=2$

# Classifying Gaussian Mixtures with Convex GLM

Data and teacher

$$\mathbf{x} \in \mathbb{R}^d, \mathbf{y} \in \mathbb{R}^K \quad P(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^K y_k \rho_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k),$$

Student

$$\mathbf{W}^* \in \min_{\mathbf{W} \in \mathbb{R}^{d \times K}} L(\mathbf{Y}, \mathbf{XW}) + \frac{\lambda}{2} \|\mathbf{W}\|_2^2$$

Learn  $K$  separating hyperplanes, i.e. a matrix  $\mathbf{W} \in \mathbb{R}^{d \times K}$

Examples : ridge regression, softmax with cross-entropy, ...



# What are the difficulties ?

---

**Learning a matrix** : how are the different hyperplanes correlated/linked by the learning process ?

# What are the difficulties ?

---

**Learning a matrix** : how are the different hyperplanes correlated/linked by the learning process ?

**Different covariances** : effect of each cluster cannot be characterized with the same quantities

# What are the difficulties ?

**Learning a matrix** : how are the different hyperplanes correlated/linked by the learning process ?

**Different covariances** : effect of each cluster cannot be characterized with the same quantities

**Convex Gaussian Comparison Inequalities break down beyond least-squares**

[Thrampoulidis et al. 20] (identity covariances)

# What are the difficulties ?

**Learning a matrix** : how are the different hyperplanes correlated/linked by the learning process ?

**Different covariances** : effect of each cluster cannot be characterized with the same quantities

**Convex Gaussian Comparison Inequalities break down beyond least-squares**

[Thrampoulidis et al. 20] (identity covariances)

**Solve it with an AMP**

## Sketch of proof

Target :

$$\mathbf{W}^* \in \min_{\mathbf{W} \in \mathbb{R}^{d \times K}} L(\mathbf{Y}, \mathbf{X}\mathbf{W}) + r(\mathbf{W}) \quad (1)$$

Tool :

$$\begin{aligned} \mathbf{u}^{t+1} &= \mathbf{Z}^\top \mathbf{h}_t(\mathbf{v}^t) - \mathbf{e}_t(\mathbf{u}^t) \langle \mathbf{h}'_t \rangle^\top \\ \mathbf{v}^t &= \mathbf{Z} \mathbf{e}_t(\mathbf{u}^t) - \mathbf{h}_{t-1}(\mathbf{v}^{t-1}) \langle \mathbf{e}'_t \rangle^\top \end{aligned} \quad (2)$$

# Sketch of proof

Target :

$$\mathbf{W}^* \in \min_{\mathbf{W} \in \mathbb{R}^{d \times K}} L(\mathbf{Y}, \mathbf{X}\mathbf{W}) + r(\mathbf{W}) \quad (1)$$

Tool :

$$\begin{aligned} \mathbf{u}^{t+1} &= \mathbf{Z}^\top \mathbf{h}_t(\mathbf{v}^t) - \mathbf{e}_t(\mathbf{u}^t) \langle \mathbf{h}'_t \rangle^\top \\ \mathbf{v}^t &= \mathbf{Z} \mathbf{e}_t(\mathbf{u}^t) - \mathbf{h}_{t-1}(\mathbf{v}^{t-1}) \langle \mathbf{e}'_t \rangle^\top \end{aligned} \quad (2)$$

Instructions:

- design  $\mathbf{h}_t, \mathbf{e}_t$  s.t. fixed point of (2) matches opt. cond. of (1)
- find a converging trajectory (convexity helps)
- use state evolution equations (fixed point)

Proof idea [Bayati and Montanari '11], see also [Feng et al. '22]

## Designing the AMP : a quick look

Often designed from a factor graph

**The factor graph for generic multiclass GLM is not obvious ...**

## Designing the AMP : a quick look

Often designed from a factor graph

**The factor graph for generic multiclass GLM is not obvious ...**

Reformulate the optimality condition

$$\mathbf{X}^\top \partial L(\mathbf{Y}, \mathbf{X} \mathbf{W}^*) + \partial r(\mathbf{W}^*) = 0$$



## Designing the AMP : a quick look

Often designed from a factor graph

**The factor graph for generic multiclass GLM is not obvious ...**

Reformulate the optimality condition

$$\mathbf{X}^\top \partial L(\mathbf{Y}, \mathbf{X} \mathbf{W}^*) + \partial r(\mathbf{W}^*) = 0$$

Match it with the fixed point

$$(Id + \mathbf{e}(\bullet) \langle \mathbf{h}' \rangle)(\mathbf{u}) = \mathbf{Z}^\top \mathbf{h}(\mathbf{v})$$

$$(Id + \mathbf{h}(\bullet) \langle \mathbf{e}' \rangle)(\mathbf{v}) = \mathbf{Z} \mathbf{e}(\mathbf{u})$$

Non-separable, block structure gradient

$$\mathbf{z}^\top \begin{bmatrix} \partial \tilde{L}_1(\mathbf{z}_1 \tilde{\mathbf{w}}_1) & & & \\ & \partial \tilde{L}_2(\mathbf{z}_2 \tilde{\mathbf{w}}_2) & (0) & \\ & (0) & \ddots & \\ & & & \partial \tilde{L}_K(\mathbf{z}_K \tilde{\mathbf{w}}_K) \end{bmatrix} + \begin{bmatrix} \partial \tilde{r}(\tilde{\mathbf{w}})_1 & & & \\ & \partial \tilde{r}(\tilde{\mathbf{w}})_2 & (0) & \\ & (0) & \ddots & \\ & & & \partial \tilde{r}(\tilde{\mathbf{w}})_K \end{bmatrix}$$

Non-separable, block structure gradient

$$\mathbf{z}^\top \begin{bmatrix} \partial \tilde{L}_1(\mathbf{z}_1 \tilde{\mathbf{w}}_1) & & & \\ & \partial \tilde{L}_2(\mathbf{z}_2 \tilde{\mathbf{w}}_2) & (0) & \\ & (0) & \ddots & \\ & & & \partial \tilde{L}_K(\mathbf{z}_K \tilde{\mathbf{w}}_K) \end{bmatrix} + \begin{bmatrix} \partial \tilde{r}(\tilde{\mathbf{w}})_1 & & & \\ & \partial \tilde{r}(\tilde{\mathbf{w}})_2 & (0) & \\ & (0) & \ddots & \\ & & & \partial \tilde{r}(\tilde{\mathbf{w}})_K \end{bmatrix}$$

**Spatially coupled, non-separable, matrix-valued, two-layer AMP**

## Converging trajectory

Prove  $\lim_{t \rightarrow \infty} \lim_{d \rightarrow \infty} \frac{1}{d} \|\mathbf{u}^t - \mathbf{u}^{t+1}\|_F^2 = 0$

Use SE equations and well chosen initialization

$$\lim_{d \rightarrow \infty} \frac{1}{d} \|\mathbf{u}^t - \mathbf{u}^{t+1}\|_F^2 = \mathbf{C} - \mathbf{C}_{t,t+1}$$

$\mathbf{C}_{t,t+1} \stackrel{d \rightarrow +\infty}{=} \frac{1}{d} (\mathbf{u}^t)^\top \mathbf{u}^{t+1} = \text{expectation over SE fields, } \mathbf{C}_{t,t} = \mathbf{C}$

# Converging trajectory

Prove  $\lim_{t \rightarrow \infty} \lim_{d \rightarrow \infty} \frac{1}{d} \|\mathbf{u}^t - \mathbf{u}^{t+1}\|_F^2 = 0$

Use SE equations and well chosen initialization

$$\lim_{d \rightarrow \infty} \frac{1}{d} \|\mathbf{u}^t - \mathbf{u}^{t+1}\|_F^2 = \mathbf{C} - \mathbf{C}_{t,t+1}$$

$\mathbf{C}_{t,t+1} \stackrel{d \rightarrow +\infty}{=} \frac{1}{d} (\mathbf{u}^t)^\top \mathbf{u}^{t+1} = \text{expectation over SE fields, } \mathbf{C}_{t,t} = \mathbf{C}$

Prescribes an iteration  $\mathbf{C}_{t,t+1} = \mathcal{O}(\mathbf{C}_{t,t-1})$

In the strongly convex case, can prove  $\mathcal{O}$  is a contraction



Trajectory converges

# Main result (informal)

## Theorem [LSGPKZ21]

Fixed-point of self-consistent equations

$$\begin{cases} \mathbf{Q}_k = \frac{1}{d} \mathbb{E}_{\Xi} [\mathbf{G} \boldsymbol{\Sigma}_k \mathbf{G}^{\top}] \\ \mathbf{M}_k = \frac{1}{\sqrt{d}} \mathbb{E}_{\Xi} [\mathbf{G} \boldsymbol{\mu}_k] \\ \mathbf{V}_k = \frac{1}{d} \mathbb{E}_{\Xi} \left[ \left( \mathbf{G} \odot \left( \hat{\mathbf{Q}}_k \otimes \boldsymbol{\Sigma}_k \right)^{-\frac{1}{2}} \odot (\mathbf{I}_K \otimes \boldsymbol{\Sigma}_k) \right) \Xi_k^{\top} \right] \end{cases} \quad \begin{cases} \hat{\mathbf{Q}}_k = \alpha \rho_k \mathbb{E}_{\xi} [\mathbf{f}_k \mathbf{f}_k^{\top}] \\ \hat{\mathbf{V}}_k = -\alpha \rho_k \mathbf{Q}_k^{-\frac{1}{2}} \mathbb{E}_{\xi} [\mathbf{f}_k \boldsymbol{\xi}^{\top}] \\ \hat{\mathbf{m}}_k = \alpha \rho_k \mathbb{E}_{\xi} [\mathbf{f}_k] \end{cases}$$

where  $\mathbf{G} = \mathbf{A}^{\frac{1}{2}} \odot \text{Prox}_{r(\mathbf{A}^{\frac{1}{2}} \odot \bullet)}(\mathbf{A}^{\frac{1}{2}} \odot \mathbf{B})$ ,  $\mathbf{A}^{-1} \equiv \sum_k \hat{\mathbf{V}}_k \otimes \boldsymbol{\Sigma}_k$ ,  $\mathbf{B} \equiv \sum_k \left( \boldsymbol{\mu}_k \hat{\mathbf{m}}_k^{\top} + \Xi_k \odot \sqrt{\hat{\mathbf{Q}}_k \otimes \boldsymbol{\Sigma}_k} \right)$

$\mathbf{f}_k \equiv \mathbf{V}_k^{-1}(\mathbf{h}_k - \boldsymbol{\omega}_k)$ ,  $\mathbf{h}_k = \mathbf{V}_k^{1/2} \text{Prox}_{\ell(\mathbf{e}_k, \mathbf{V}_k^{1/2} \bullet)}(\mathbf{V}_k^{-1/2} \boldsymbol{\omega}_k)$ ,  $\boldsymbol{\omega}_k \equiv \mathbf{M}_k + \mathbf{b} + \mathbf{Q}_k^{1/2} \boldsymbol{\xi}_k$

# Main result (informal)

## Theorem [LSGPKZ21]

Fixed-point of self-consistent equations

$$\begin{cases} \mathbf{Q}_k = \frac{1}{d} \mathbb{E}_{\Xi} [\mathbf{G} \boldsymbol{\Sigma}_k \mathbf{G}^{\top}] \\ \mathbf{M}_k = \frac{1}{\sqrt{d}} \mathbb{E}_{\Xi} [\mathbf{G} \boldsymbol{\mu}_k] \\ \mathbf{V}_k = \frac{1}{d} \mathbb{E}_{\Xi} \left[ \left( \mathbf{G} \odot \left( \hat{\mathbf{Q}}_k \otimes \boldsymbol{\Sigma}_k \right)^{-\frac{1}{2}} \odot (\mathbf{I}_K \otimes \boldsymbol{\Sigma}_k) \right) \Xi_k^{\top} \right] \end{cases} \quad \begin{cases} \hat{\mathbf{Q}}_k = \alpha \rho_k \mathbb{E}_{\xi} [\mathbf{f}_k \mathbf{f}_k^{\top}] \\ \hat{\mathbf{V}}_k = -\alpha \rho_k \mathbf{Q}_k^{-\frac{1}{2}} \mathbb{E}_{\xi} [\mathbf{f}_k \xi^{\top}] \\ \hat{\mathbf{m}}_k = \alpha \rho_k \mathbb{E}_{\xi} [\mathbf{f}_k] \end{cases}$$

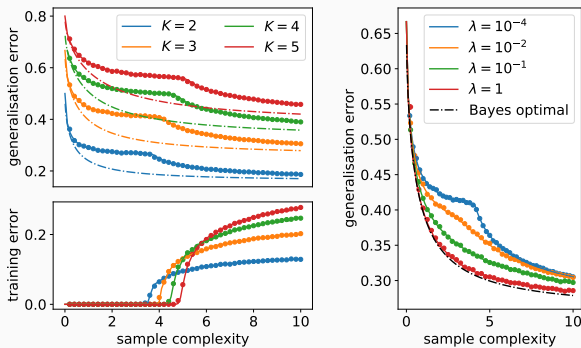
where  $\mathbf{G} = \mathbf{A}^{\frac{1}{2}} \odot \text{Prox}_{r(\mathbf{A}^{\frac{1}{2}} \odot \mathbf{B})}(\mathbf{A}^{\frac{1}{2}} \odot \mathbf{B})$ ,  $\mathbf{A}^{-1} \equiv \sum_k \hat{\mathbf{V}}_k \otimes \boldsymbol{\Sigma}_k$ ,  $\mathbf{B} \equiv \sum_k \left( \mu_k \hat{\mathbf{m}}_k^{\top} + \Xi_k \odot \sqrt{\hat{\mathbf{Q}}_k \otimes \boldsymbol{\Sigma}_k} \right)$

$\mathbf{f}_k \equiv \mathbf{V}_k^{-1}(\mathbf{h}_k - \boldsymbol{\omega}_k)$ ,  $\mathbf{h}_k = \mathbf{V}_k^{1/2} \text{Prox}_{\ell(\mathbf{e}_k, \mathbf{V}_k^{1/2} \bullet)}(\mathbf{V}_k^{-1/2} \boldsymbol{\omega}_k)$ ,  $\boldsymbol{\omega}_k \equiv \mathbf{M}_k + \mathbf{b} + \mathbf{Q}_k^{1/2} \xi_k$

Training and generalization for  $n, d \rightarrow \infty$  :

$$\epsilon_t = 1 - \sum_{k=1}^K \rho_k \mathbb{E}_{\xi} [\hat{y}_k(\mathbf{h}_k)], \quad \epsilon_g = 1 - \sum_{k=1}^K \rho_k \mathbb{E}_{\xi} [\hat{y}_k(\boldsymbol{\omega}_k)].$$

# Examples : synthetic random design problems



Ridge penalized logistic regression on K Gaussian clusters,  $\Sigma_k = \Delta Id$ . (Left) Sample complexity (Right) Regularization

Related works : [T. Cover '69] [E. Gardner, B. Derrida '89] [E.J. Candès, P. Sur '20] [F. Mignacco, F. Krzakala, Y. Lu, P. Urbani, L. Zdeborova '20] [C. Thrampoulidis, S. Oymak, M. Soltanolkotabi '20]



**Thank you**