Statistical physics of learning : a mathematical perspective

Cedric Gerbelot

ENS & INRIA, Paris, France April 7, 2022

An introductory example : binary classification on MNIST



Figure 1: Binary classification on Mnist/Fashion-Mnist, odd vs even, ℓ_2 regularized logistic regression

An introductory example : binary classification on MNIST



Figure 1: Binary classification on Mnist/Fashion-Mnist, odd vs even, ℓ_2 regularized logistic regression

Exact theory ? So what's the answer ? [B. Loureiro, G. Sicuro, C. Gerbelot, A. Pacco, F. Krzakala, L. Zdeborova '21]

An introductory example : binary classification on MNIST

Theorem 1 (Concentration properties of the estimator). Let $\xi_{k \in [K]} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ be collection of *K*-dimensional standard normal vectors independent of other quantities. Let also be $\{\Xi_k\}$ a set of *K* matrices, $\Xi_k \in \mathbb{R}^{K \times d}$, with i.i.d. standard normal entries, independent of other quantities. Under the set of assumptions (A1–A5), for any pseudo-Lispchitz functions of finite order $\phi_1 : \mathbb{R}^{K \times d} \rightarrow \mathbb{R}$, $\phi_2 : \mathbb{R}^{K \times n} \rightarrow \mathbb{R}$, the estimator W^* and the matrix $Z^* = \frac{1}{\sqrt{d}}W^*X$ verify:

$$\phi_1(\boldsymbol{W}^{\star}) \xrightarrow{P} \mathbb{E}_{\boldsymbol{\Xi}} \left[\phi_1(\boldsymbol{G}) \right], \qquad \qquad \phi_2(\boldsymbol{Z}^{\star}) \xrightarrow{P} \mathbb{E}_{\boldsymbol{\xi}} \left[\phi_2(\boldsymbol{H}) \right], \qquad \qquad (6)$$

where we have introduced the proximal for the loss:

$$\boldsymbol{h}_{k} = \boldsymbol{V}_{k}^{1/2} \operatorname{Prox}_{\ell(\boldsymbol{e}_{k}, \boldsymbol{V}_{k}^{1/2} \bullet)}(\boldsymbol{V}_{k}^{-1/2} \boldsymbol{\omega}_{k}) \in \mathbb{R}^{K}, \qquad \boldsymbol{\omega}_{k} \equiv \boldsymbol{m}_{k} + \boldsymbol{b} + \boldsymbol{Q}_{k}^{1/2} \boldsymbol{\xi}_{k},$$
(7)

and $\mathbf{H} \in \mathbb{R}^{K \times n}$ is obtained by concatenating each \mathbf{h}_k , $\rho_k n$ times. We have also introduced the matrix proximal $\mathbf{G} \in \mathbb{R}^{K \times d}$:

$$\boldsymbol{G} = \mathbf{A}^{\frac{1}{2}} \odot \operatorname{Prox}_{r(\mathbf{A}^{\frac{1}{2}} \odot \bullet)}(\mathbf{A}^{\frac{1}{2}} \odot \boldsymbol{B}), \quad \mathbf{A}^{-1} \equiv \sum_{k} \hat{V}_{k} \otimes \boldsymbol{\Sigma}_{k}, \ \boldsymbol{B} \equiv \sum_{k} \left(\boldsymbol{\mu}_{k} \hat{\boldsymbol{m}}_{k}^{\top} + \boldsymbol{\Xi}_{k} \odot \sqrt{\hat{\boldsymbol{Q}}_{k} \otimes \boldsymbol{\Sigma}_{k}} \right).$$

The collection of parameters $(Q_k, m_k, V_k, \hat{Q}_k, \hat{m}_k, \hat{V}_k)_{k \in [K]}$ is given by the fixed point of the following self-consistent equations:

$$\begin{cases} Q_{k} = \frac{1}{q} \mathbb{E}_{\Xi}[G\Sigma_{k}G^{\top}] \\ m_{k} = \frac{1}{\sqrt{d}} \mathbb{E}_{\Xi}[G\mu_{k}] \\ V_{k} = \frac{1}{d} \mathbb{E}_{\Xi}\left[\left(G \circ \left(\hat{Q}_{k} \otimes \Sigma_{k} \right)^{-\frac{1}{2}} \odot (\mathbf{I}_{K} \otimes \Sigma_{k}) \right) \Xi_{k}^{\top} \right] \end{cases} \begin{cases} \hat{Q}_{k} = \alpha \rho_{k} \mathbb{E}_{\xi} \left[f_{k} f_{k}^{\top} \right] \\ \hat{V}_{k} = -\alpha \rho_{k} Q_{k}^{-\frac{1}{2}} \mathbb{E}_{\xi} \left[f_{k} \xi^{\top} \right] \\ \hat{m}_{k} = \alpha \rho_{k} \mathbb{E}_{\xi} \left[f_{k} \xi^{\top} \right] \end{cases} \tag{8}$$
where $f_{k} \equiv V_{k}^{-1}(h_{k} - \omega_{k})$, and the vector b^{\star} is such that $\sum_{k} \rho_{k} \mathbb{E}_{\xi} \left[V_{k} f_{k}^{k} \right] = 0$ holds.

Figure 2: Main theorem from [B. Loureiro, G. Sicuro, C. Gerbelot, A. Pacco, F. Krzakala, L. Zdeborova '21]

What's the motivation for these formulas ? Why are they useful ? How are they obtained ?

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

$$\mathbf{w} \sim p_{\mathbf{w}}(\mathbf{w})$$
 (1)

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

$$\mathbf{w} \sim \rho_{\mathbf{w}}(\mathbf{w})$$
 (1)

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

$$\mathbf{y} \sim p_{\mathbf{y}} \left(\mathbf{y} | \mathbf{w} \right) \tag{2}$$

Estimate w ?

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

$$\mathbf{w} \sim p_{\mathbf{w}}(\mathbf{w})$$
 (1)

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

$$\mathbf{y} \sim p_{\mathbf{y}} \left(\mathbf{y} | \mathbf{w} \right) \tag{2}$$

Estimate w ?

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

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} p_{\mathbf{w}}(\mathbf{w}) p_{\mathbf{y}}(\mathbf{y}|\mathbf{w}) \, d\mu(\mathbf{w}) \tag{3}$$

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

$$\mathbf{w} \sim p_{\mathbf{w}}(\mathbf{w})$$
 (1)

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

$$\mathbf{y} \sim p_{\mathbf{y}} \left(\mathbf{y} | \mathbf{w} \right) \tag{2}$$

Estimate w ?

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

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} p_{\mathbf{w}}(\mathbf{w}) p_{\mathbf{y}}(\mathbf{y}|\mathbf{w}) \, d\mu(\mathbf{w}) \tag{3}$$

Problem : This is a high-dimensional integral !

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

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta \left(g(\mathbf{w}, \mathbf{y}) + f(\mathbf{w})\right)\right) d\mu(\mathbf{w}) \tag{4}$$

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

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta \left(g(\mathbf{w}, \mathbf{y}) + f(\mathbf{w})\right)\right) d\mu(\mathbf{w})$$
(4)

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

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

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta \left(g(\mathbf{w}, \mathbf{y}) + f(\mathbf{w})\right)\right) d\mu(\mathbf{w}) \tag{4}$$

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

Stat. phys. toolbox deals with this problem

Link with statistical physics : disordered systems

- distributions involve a dense, random interaction matrix $\mathbf{X} \in \mathbb{R}^{n imes d}$
- y can come from another stochastic model, i.e.

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

• estimate the generative model with postulated densities g, f

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta \left(g(\mathbf{X}\mathbf{w}, \mathbf{y}) + f(\mathbf{w})\right)) d\mu(\mathbf{w})$$
 (5)

Link with statistical physics : disordered systems

- distributions involve a dense, random interaction matrix $\mathbf{X} \in \mathbb{R}^{n imes d}$
- y can come from another stochastic model, i.e.

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

• estimate the generative model with postulated densities g, f

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta \left(g(\mathbf{X}\mathbf{w}, \mathbf{y}) + f(\mathbf{w})\right)\right) d\mu(\mathbf{w})$$
(5)

- Disordered equilibrium Boltzmann measure
- toolbox : replica/cavity method, belief-propagation, etc...
- focus on **X** with i.i.d. $\mathcal{N}(0,1)$ elements

Central limit theorem: If $\{w_i\}_{i=1}^d$ i.i.d. in L^2 , then $\frac{1}{\sqrt{d}} \left(\sum_{i=1}^d w_i - d\mathbb{E}[w] \right) \xrightarrow[d \to +\infty]{in \text{ law}} \mathcal{N}(0, \sigma^2)$ (6)

Central limit theorem: If $\{w_i\}_{i=1}^d$ i.i.d. in L^2 , then $\frac{1}{\sqrt{d}} \left(\sum_{i=1}^d w_i - d\mathbb{E}[w] \right) \xrightarrow[d \to +\infty]{in \text{ law}} \mathcal{N}(0, \sigma^2)$ (6)

Concentration of measure :

$$\Phi(\mathbf{w}) \xrightarrow{P} \mathbb{E}\left[\Phi(\mathbf{w})\right] = \mathbb{E}\left[\phi(w)\right]$$
(7)

for sufficiently regular Φ and fast decaying p_w .

Central limit theorem:
If
$$\{w_i\}_{i=1}^d$$
 i.i.d. in L^2 , then

$$\frac{1}{\sqrt{d}} \left(\sum_{i=1}^d w_i - d\mathbb{E}[w] \right) \xrightarrow[d \to +\infty]{in law} \mathcal{N}(0, \sigma^2)$$
(6)

Concentration of measure :

$$\Phi(\mathbf{w}) \xrightarrow[d \to \infty]{P} \mathbb{E}\left[\Phi(\mathbf{w})\right] = \mathbb{E}\left[\phi(w)\right]$$
(7)

for sufficiently regular Φ and fast decaying p_w .

Want the same thing for $\hat{\boldsymbol{w}}$

Central limit theorem: If $\{w_i\}_{i=1}^d$ i.i.d. in L^2 , then $1 \left(\frac{d}{d}\right)$ in law

$$\frac{1}{\sqrt{d}} \left(\sum_{i=1}^{d} w_i - d\mathbb{E}\left[w \right] \right) \xrightarrow[d \to +\infty]{\text{in law}} \mathcal{N}(0, \sigma^2)$$
(6)

Concentration of measure :

$$\Phi(\mathbf{w}) \xrightarrow{P} \mathbb{E}\left[\Phi(\mathbf{w})\right] = \mathbb{E}\left[\phi(w)\right]$$
(7)

for sufficiently regular Φ and fast decaying p_w .

Want the same thing for $\hat{\boldsymbol{w}}$

Problem : ŵ is not i.i.d.! Strong coupling notably due to X

Finding a decoupled measure

Stat.phys. toolbox (and ensuing math !!) \Downarrow Product measure of simple components describing $\hat{\mathbf{w}}$ for large n, d

Stat.phys. toolbox (and ensuing math !!) \downarrow Product measure of simple components describing $\hat{\mathbf{w}}$ for large n, d

Theorems then take the form

$$\begin{split} \Phi(\hat{\mathbf{w}}) & \xrightarrow{P} \mathbb{E} \left[\Phi\left(\Omega(\mathbf{Z}), \{q_s\}_{s \in \mathcal{S}} \right) \right] \\ \Phi(\hat{\mathbf{w}}) & \xrightarrow{P} \mathbb{E} \left[\phi\left(\omega(z), \{q_s\}_{s \in \mathcal{S}} \right) \right] & \text{if separable problem, 1D integral} \end{split}$$

Stat.phys. toolbox (and ensuing math !!) \downarrow Product measure of simple components describing $\hat{\mathbf{w}}$ for large n, d

Theorems then take the form

$$\begin{split} \Phi(\hat{\mathbf{w}}) &\xrightarrow{P} \mathbb{E} \left[\Phi\left(\Omega(\mathbf{Z}), \{q_s\}_{s \in \mathcal{S}}\right) \right] \\ \Phi(\hat{\mathbf{w}}) &\xrightarrow{P} \mathbb{E} \left[\phi\left(\omega(z), \{q_s\}_{s \in \mathcal{S}}\right) \right] & \text{if separable problem, 1D integral} \end{split}$$

- $\mathbf{Z} \in \mathbb{R}^d$ i.i.d. $\mathcal{N}(0,1)$
- Ω explicit, bounded variation function
- q_s are low-dimensional paramaters
- given by explicit, self-consistent equations

Recall

$$\hat{\mathbf{w}} = \frac{1}{\mathcal{Z}(\mathbf{y})} \int \mathbf{w} \exp(-\beta \left(g(\mathbf{X}\mathbf{w}, \mathbf{y}) + f(\mathbf{w})\right)) d\mu(\mathbf{x})$$
(8)

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

$$\hat{\mathbf{w}} \xrightarrow[\beta \to +\infty]{} \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} g(\mathbf{X}\mathbf{w}, \mathbf{y}) + f(\mathbf{w})$$
 (9)

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

Examples : LASSO, logistic regression, etc ...

From the decoupled measure : training, generalization error, phenomenology, etc ...

- typical case
- benchmark, random design problems
- exact solutions
- strong assumptions (i.i.d. Gaussian !)

- typical case
- benchmark, random design problems
- exact solutions
- strong assumptions (i.i.d. Gaussian !)

How realistic are the stat. phys. benchmarks ? What can we do to make them more realistic ? Observe "teacher" generative model

 $oldsymbol{y} = f_0(oldsymbol{X}oldsymbol{w}_0) \in \mathbb{R}^n, \quad oldsymbol{w}_0 \in \mathbb{R}^d \quad oldsymbol{X} \in \mathbb{R}^{n imes d} ext{ i.i.d. } \mathcal{N}(0,1)$

Teacher Student Generalized Linear Model

Observe "teacher" generative model

 $m{y} = f_0(m{X}m{w}_0) \in \mathbb{R}^n, \quad m{w}_0 \in \mathbb{R}^d \quad m{X} \in \mathbb{R}^{n imes d} \text{ i.i.d. } \mathcal{N}(0,1)$

Learn with "student"

$$oldsymbol{w}^{\star} \in \operatorname*{argmin}_{oldsymbol{w} \in \mathbb{R}^d} L(oldsymbol{y}, oldsymbol{X} oldsymbol{w}) + r(oldsymbol{w})$$

- L, r are a convex loss and penalty
- $n, d \rightarrow +\infty$ with fixed ratio

Teacher Student Generalized Linear Model

Observe "teacher" generative model

 $m{y} = f_0(m{X}m{w}_0) \in \mathbb{R}^n, \quad m{w}_0 \in \mathbb{R}^d \quad m{X} \in \mathbb{R}^{n imes d} \text{ i.i.d. } \mathcal{N}(0,1)$

Learn with "student"

$$\boldsymbol{w}^{\star} \in \operatorname*{argmin}_{\boldsymbol{w} \in \mathbb{R}^{d}} L(\boldsymbol{y}, \boldsymbol{X} \boldsymbol{w}) + r(\boldsymbol{w})$$

- L, r are a convex loss and penalty
- $n, d \rightarrow +\infty$ with fixed ratio

Goal : statistical properties of w^*

Teacher Student Generalized Linear Model

Observe "teacher" generative model

 $m{y} = f_0(m{X}m{w}_0) \in \mathbb{R}^n, \quad m{w}_0 \in \mathbb{R}^d \quad m{X} \in \mathbb{R}^{n imes d} \text{ i.i.d. } \mathcal{N}(0,1)$

Learn with "student"

$$\boldsymbol{w}^{\star} \in \operatorname*{argmin}_{\boldsymbol{w} \in \mathbb{R}^{d}} L(\boldsymbol{y}, \boldsymbol{X} \boldsymbol{w}) + r(\boldsymbol{w})$$

- L, r are a convex loss and penalty
- $n, d \rightarrow +\infty$ with fixed ratio

Goal : statistical properties of w^*

Beyond i.i.d. assumption : introduce correlation

Introducing Correlation : a Block Covariance Model

Teacher and student with different feature spaces

Block covariate model proposed in [B. Loureiro, **CG**, H. Cui, S. Goldt, M. Mézard, F. Krzakala, L. Zdeborova '21]

$$\begin{bmatrix} \boldsymbol{u} \\ \boldsymbol{v} \end{bmatrix} \in \mathbb{R}^{p+d} \sim \mathcal{N}\left(0, \begin{bmatrix} \Psi & \Phi \\ \Phi^{\top} & \Omega \end{bmatrix}\right) \quad y^{\mu} = f_0\left(\frac{1}{\sqrt{p}} \boldsymbol{w}_0^{\top} \boldsymbol{u}^{\mu}\right),$$
$$\boldsymbol{w}^{\star} = \operatorname{argmin}_{\boldsymbol{w} \in \mathbb{R}^d} \left[\sum_{\mu=1}^n l\left(\frac{\boldsymbol{w}^{\top} \boldsymbol{v}^{\mu}}{\sqrt{d}}, y^{\mu}\right) + r(\boldsymbol{w})\right]$$

Many works: [E. Dobriban, S. Wager '15][PL. Bartlett, PM. Long, G. Lugosi, A. Tsigler '19][T. Hastie, A. Montanari, S. Rosset, RJ. Tibshirani '19][M. Celentano, A. Montanari, Y. Wei '20]

Solution to Block Covariance model

Theorem (informal)[B. Loureiro, C.Gerbelot, H. Cui, S. Goldt, M. Mézard, F. Krzakala, L. Zdeborova '21]

Unique fixed point of self-consistent equations

$$\begin{cases} V = \mathbb{E}_{(\omega,\bar{\theta})\sim\mu} \left[\frac{\omega}{\lambda+\bar{V}\omega} \right] \\ m = \frac{\hat{m}}{\sqrt{\gamma}} \mathbb{E}_{(\omega,\bar{\theta})\sim\mu} \left[\frac{\bar{\theta}^2}{\lambda+\bar{V}\omega} \right] \\ q = \mathbb{E}_{(\omega,\bar{\theta})\sim\mu} \left[\frac{\hat{m}^2 \bar{\theta}^2 \omega + \hat{q}\omega^2}{(\lambda+\bar{V}\omega)^2} \right] \end{cases}, \quad \begin{cases} \hat{V} = \frac{\alpha}{V} (1 - \mathbb{E}_{s,h\sim\mathcal{N}(0,1)} [z'(V,m,q)]) \\ \hat{m} = \frac{1}{\sqrt{\rho\gamma}} \frac{\alpha}{V} \mathbb{E}_{s,h\sim\mathcal{N}(0,1)} \left[sz(V,m,q) - \frac{m}{\sqrt{\rho}} z'(V,m,q) \right] \\ \hat{q} = \frac{\alpha}{V^2} \mathbb{E}_{s,h\sim\mathcal{N}(0,1)} \left[\left(\frac{m}{\sqrt{\rho}} s + \sqrt{q - \frac{m^2}{\rho}} h - z(V,m,q) \right)^2 \right] \end{cases}$$

where $z(V, m, q) = \text{prox}_{Vl(., f_0(\sqrt{\rho}s))}(\rho^{-1/2}ms + \sqrt{q - \rho^{-1}m^2}h)$

Solution to Block Covariance model

Theorem (informal)[B. Loureiro, C.Gerbelot, H. Cui, S. Goldt, M. Mézard, F. Krzakala, L. Zdeborova '21]

Unique fixed point of self-consistent equations

$$\begin{cases} V = \mathbb{E}_{(\omega,\bar{\theta})\sim\mu} \left[\frac{\omega}{\lambda+\bar{V}\omega} \right] \\ m = \frac{\hat{m}}{\sqrt{\gamma}} \mathbb{E}_{(\omega,\bar{\theta})\sim\mu} \left[\frac{\bar{\theta}^2}{\lambda+\bar{V}\omega} \right] \\ q = \mathbb{E}_{(\omega,\bar{\theta})\sim\mu} \left[\frac{\hat{m}^2\bar{\theta}^2\omega+\hat{q}\omega^2}{(\lambda+\bar{V}\omega)^2} \right] \end{cases}, \quad \begin{cases} \hat{V} = \frac{\alpha}{V} (1 - \mathbb{E}_{s,h\sim\mathcal{N}(0,1)}[z'(V,m,q)]) \\ \hat{m} = \frac{1}{\sqrt{\rho\gamma}} \frac{\alpha}{V} \mathbb{E}_{s,h\sim\mathcal{N}(0,1)} \left[sz(V,m,q) - \frac{m}{\sqrt{\rho}} z'(V,m,q) \right] \\ \hat{q} = \frac{\alpha}{V^2} \mathbb{E}_{s,h\sim\mathcal{N}(0,1)} \left[\left(\frac{m}{\sqrt{\rho}} s + \sqrt{q - \frac{m^2}{\rho}} h - z(V,m,q) \right)^2 \right] \end{cases}$$

where
$$z(V, m, q) = \text{prox}_{Vl(., f_0(\sqrt{\rho}s))}(\rho^{-1/2}ms + \sqrt{q - \rho^{-1}m^2}h)$$

 $\textit{n},\textit{p},\textit{d} \rightarrow \infty,$ training and generalization error :

$$\begin{split} & \mathcal{E}_{\text{train.}}(\hat{\mathbf{w}}) \xrightarrow{P} \mathbb{E}_{s,h \sim \mathcal{N}(0,1)} \left[I\left(\operatorname{prox}_{V^{\star}I(.,f_{0}(\sqrt{\rho}s))}\left(\frac{m^{\star}}{\sqrt{\rho}}s + \sqrt{q^{\star} - \frac{m^{\star 2}}{\rho}}h\right), f_{0}(\sqrt{\rho}s) \right) \right] \\ & \mathcal{E}_{\text{gen.}}(\hat{\mathbf{w}}) \xrightarrow{P} \mathbb{E}_{(\nu,\lambda)} \left[\hat{g}\left(\hat{f}(\lambda), f_{0}(\nu)\right) \right] \end{split}$$

How well does it work ?

Ridge regression works well ...



Figure 3: (Left) Ridge regression on real data. (Right) Logistic regression with real and synthetic (GAN) data

... but classification is more problematic

How well does it work ?

Ridge regression works well ...



Figure 3: (Left) Ridge regression on real data. (Right) Logistic regression with real and synthetic (GAN) data

... but classification is more problematic

Need for another realistic benchmark problem

Study classification of k-Gaussian mixture with convex GLM

Study classification of k-Gaussian mixture with convex GLM

- Benchmark problem in ML, universal approximation, ...
- many scenarios described by Gaussian mixtures (GANs, 'Neural collapse', ...)

[M. Seddik, C. Louart, M. Tamaazousti, R. Couillet, '20][V. Papyan, X. Han, D.Donoho, '20]

Classifying Gaussian Mixtures with Convex GLM

Data and teacher

$$oldsymbol{x} \in \mathbb{R}^{d}, oldsymbol{y} \in \mathbb{R}^{K} \quad P(oldsymbol{x},oldsymbol{y}) = \sum_{k=1}^{K} y_k
ho_k \mathcal{N}\left(oldsymbol{x} \mid oldsymbol{\mu}_k, oldsymbol{\Sigma}_k
ight),$$



Figure 4: K=3, d=2

Data and teacher

$$oldsymbol{x} \in \mathbb{R}^{d}, oldsymbol{y} \in \mathbb{R}^{K} \quad P(oldsymbol{x},oldsymbol{y}) = \sum_{k=1}^{K} y_k
ho_k \mathcal{N}\left(oldsymbol{x} \mid oldsymbol{\mu}_k, oldsymbol{\Sigma}_k
ight),$$

Student

$$\boldsymbol{W}^{\star} \in \min_{\boldsymbol{W} \in \mathbb{R}^{d imes K}} L(\boldsymbol{Y}, \boldsymbol{X} \boldsymbol{W}) + r(\boldsymbol{W})$$

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

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

Main result (informal)

Theorem [B. Loureiro, G. Sicuro, CG, A. Pacco, F. Krzakala, L. Zdeborova '21]

Fixed-point of self-consistent equations

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

where
$$\mathbf{G} = \mathbf{A}^{\frac{1}{2}} \odot \operatorname{Prox}_{r(\mathbf{A}^{\frac{1}{2}} \odot \mathbf{e})} (\mathbf{A}^{\frac{1}{2}} \odot \mathbf{B}), \ \mathbf{A}^{-1} \equiv \sum_{k} \hat{\mathbf{V}}_{k} \otimes \mathbf{\Sigma}_{k}, \ \mathbf{B} \equiv \sum_{k} \left(\boldsymbol{\mu}_{k} \hat{\mathbf{m}}_{k}^{\top} + \mathbf{\Xi}_{k} \odot \sqrt{\hat{\mathbf{Q}}_{k} \otimes \mathbf{\Sigma}_{k}} \right)$$

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

Main result (informal)

Theorem [B. Loureiro, G. Sicuro, CG, A. Pacco, F. Krzakala, L. Zdeborova '21]

Fixed-point of self-consistent equations

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

where
$$\mathbf{G} = \mathbf{A}^{\frac{1}{2}} \odot \operatorname{Prox}_{r(\mathbf{A}^{\frac{1}{2}} \odot \mathbf{e})} (\mathbf{A}^{\frac{1}{2}} \odot \mathbf{B}), \ \mathbf{A}^{-1} \equiv \sum_{k} \hat{\mathbf{V}}_{k} \otimes \mathbf{\Sigma}_{k}, \ \mathbf{B} \equiv \sum_{k} \left(\boldsymbol{\mu}_{k} \hat{\mathbf{m}}_{k}^{\top} + \mathbf{\Xi}_{k} \odot \sqrt{\hat{\mathbf{Q}}_{k} \otimes \mathbf{\Sigma}_{k}} \right)$$

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

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

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

Main result : the intuition (forget dimensions)

For any convex GLM (separable, covariance, matrix variable, etc ...)

Same global form of the result

For any convex GLM (separable, covariance, matrix variable, etc ...)

Same global form of the result

Problem defined by :

- estimator $\hat{\boldsymbol{w}}$
- "ground-truth" \pmb{w}_0/μ , random data with covariance Σ
- cost function L + r.

For any convex GLM (separable, covariance, matrix variable, etc ...)

Same global form of the result

Problem defined by :

- estimator $\hat{\boldsymbol{w}}$
- "ground-truth" \pmb{w}_0/μ , random data with covariance Σ
- cost function L + r.

Statement

 $\hat{\boldsymbol{w}} \sim \text{nonlinearity}_{L,r}(\alpha * \text{ground-truth} + \text{gaussian})$

With low-dim. closed form parameters.

Main result : a closer look

For any "well-behaved" observable ϕ_1 :

$$\phi_1(\boldsymbol{W}^{\star}) \xrightarrow{P} \mathbb{E}_{\Xi} [\phi_1(\boldsymbol{G})]$$

Main result : a closer look

For any "well-behaved" observable ϕ_1 :

$$\phi_1(\boldsymbol{W}^{\star}) \xrightarrow{P} \mathbb{E}_{\Xi} [\phi_1(\boldsymbol{G})]$$

$$\begin{split} \boldsymbol{G} &= \boldsymbol{\mathsf{A}}^{\frac{1}{2}} \odot \operatorname{Prox}_{r(\boldsymbol{\mathsf{A}}^{\frac{1}{2}} \odot \boldsymbol{\bullet})}(\boldsymbol{\mathsf{A}}^{\frac{1}{2}} \odot \boldsymbol{B}), \\ \boldsymbol{\mathsf{A}}^{-1} &\equiv \sum_{k} \hat{\boldsymbol{V}}_{k} \otimes \boldsymbol{\Sigma}_{k}, \qquad \boldsymbol{B} \equiv \sum_{k} \left(\boldsymbol{\mu}_{k} \hat{\boldsymbol{m}}_{k}^{\top} + \boldsymbol{\Xi}_{k} \odot \sqrt{\hat{\boldsymbol{Q}}_{k} \otimes \boldsymbol{\Sigma}_{k}} \right). \end{split}$$

Main result : a closer look

For any "well-behaved" observable ϕ_1 :

$$\phi_1(\boldsymbol{W}^{\star}) \xrightarrow{P} \mathbb{E}_{\Xi} [\phi_1(\boldsymbol{G})]$$

$$\begin{aligned} \boldsymbol{G} &= \boldsymbol{\mathsf{A}}^{\frac{1}{2}} \odot \operatorname{Prox}_{r(\boldsymbol{\mathsf{A}}^{\frac{1}{2}} \odot \boldsymbol{\bullet})}(\boldsymbol{\mathsf{A}}^{\frac{1}{2}} \odot \boldsymbol{B}), \\ \boldsymbol{\mathsf{A}}^{-1} &\equiv \sum_{k} \hat{\boldsymbol{\mathsf{V}}}_{k} \otimes \boldsymbol{\Sigma}_{k}, \qquad \boldsymbol{B} \equiv \sum_{k} \left(\boldsymbol{\mu}_{k} \hat{\boldsymbol{m}}_{k}^{\top} + \boldsymbol{\Xi}_{k} \odot \sqrt{\hat{\boldsymbol{\mathsf{Q}}}_{k} \otimes \boldsymbol{\Sigma}_{k}} \right). \end{aligned}$$

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

where

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

- very generic statement
- greatly simplifies with assumptions on covariances, separability of functions, ...
- in most cases reduces to low dimensional statement

Examples : synthetic random design problems



Figure 5: 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] [EJ. Candès, P. Sur '20] [F. Mignacco, F. Krzakala, Y. Lu, P. Urbani, L. Zdeborova '20][C. Thrampoulidis, S. Oymak, M. Soltanolkotabi '20]

Examples : real data



Figure 6: Binary classification on Mnist/Fashion-Mnist, odd vs even, Gaussian approximation and real data

Examples : real data





Figure 7: Adding more clusters to the Gaussian approximation

Figure 8: Idealized view

- correlated Gaussian designs are meaningful
- study the phenomenology : sample complexity, regularization, etc...
- more models : ensembling, boosting
- dataset/feature map distribution AND predictor geometry
- both important for Gaussian equivalence

Ongoing research : what feature maps and predictors have Gaussian equivalence properties ? Recall the problem : find a decoupled/product measure

- Guerra interpolation
- Gordon minmax
- cavity method
- approximate message passing iterations

Focus on AMP iterations

Recall the problem : find a decoupled/product measure

- Guerra interpolation interpolate with a known replica solution
- Gordon minmax Gaussian comparison inequalities
- cavity method effect of a single variable among d
- approximate message passing iterations more detail shortly

Focus on AMP iterations

Initially a relaxation of belief propagation equations

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

$$\mathbf{x}^{t+1} = \mathbf{A}\mathbf{m}^t - b_t \mathbf{m}^{t-1} \tag{10}$$

$$\mathbf{m}^t = f_t(\mathbf{x}^t) \tag{11}$$

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

$$b_t = \operatorname{div}\left[f_t(\mathbf{x}^t)\right] \tag{12}$$

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

$$\mathbf{x}^{t+1} = \mathbf{A}\mathbf{m}^t - b_t \mathbf{m}^{t-1} \tag{10}$$

$$\mathbf{m}^t = f_t(\mathbf{x}^t) \tag{11}$$

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

$$b_t = \operatorname{div}\left[f_t(\mathbf{x}^t)\right] \tag{12}$$

Theorem : for any $t, \mathbf{x}^t \xrightarrow[N \to \infty]{P_{Lk}} \mathbf{Z}^t \sim \mathcal{N}(\mathbf{0}, \kappa_{t,t} \mathbf{I}_N)$

where
$$\kappa_{t,t} = \lim_{N \to \infty} \frac{1}{N} \left[f^{t-1} (\mathbf{Z}^{t-1})^{\top} f^{t-1} (\mathbf{Z}^{t-1}) \right]$$
 (13)

First proof due to E.Bolthausen '09, '14, M. Bayati & A. Montanari '11

Sketch of proof : Bolthausen conditioning

$$\mathbf{x}^{t+1} = \mathbf{A}\mathbf{m}^t - b_t \mathbf{m}^{t-1}$$
(14)
$$\mathbf{m}^t = f_t(\mathbf{x}^t)$$
(15)

Define the σ -algebra $\mathfrak{S}_t = \sigma(\mathbf{x}^1, \mathbf{x}^2, ..., \mathbf{x}^t)$. 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

$$\mathbf{x}^{t+1} = \mathbf{A}\mathbf{m}^t - b_t \mathbf{m}^{t-1}$$
(14)
$$\mathbf{m}^t = f_t(\mathbf{x}^t)$$
(15)

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

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

Gaussian conditioning lemma

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

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

A bit of algebra leads to

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

A bit of algebra leads to

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

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

Advantages:

- very generic and adaptable (non-convex, multilayer, committee,...)
- \bullet algorithm + fixed points
- control of trajectory understood in the convex case

Advantages:

- very generic and adaptable (non-convex, multilayer, committee,...)
- \bullet algorithm + fixed points
- control of trajectory understood in the convex case

Disadvantage:

- state evolution proofs are tedious, case by case basis
- control of trajectory in non-convex/RS case ? (SK above AT line)

Advantages:

- very generic and adaptable (non-convex, multilayer, committee,...)
- algorithm + fixed points
- control of trajectory understood in the convex case

Disadvantage:

- state evolution proofs are tedious, case by case basis
- control of trajectory in non-convex/RS case ? (SK above AT line)

Solution

- unifying framework for AMP iterations
- generic, modular proof of SE equations

([C. Gerbelot, R. Berthier '21])

The oriented graph:



Arbitrary composition of this structure

The graph-based AMP iteration

$$\mathbf{x}_{\overrightarrow{e}}^{t+1} = \mathbf{A}_{\overrightarrow{e}} \mathbf{m}_{\overrightarrow{e}}^{t} - b_{\overrightarrow{e}}^{t} \mathbf{m}_{\overleftarrow{e}}^{t-1}, \qquad (16)$$

$$\mathbf{m}_{\vec{e}}^{t} = f_{\vec{e}}^{t} \left(\left(\mathbf{x}_{\vec{e}'}^{t} \right)_{\vec{e}': \vec{e}' \to \vec{e}} \right) \,, \tag{17}$$

where $b_{\overrightarrow{e}}^t$ is the so-called *Onsager term*

$$b_{\overrightarrow{e}}^{t} = \frac{1}{N} \operatorname{Tr} \frac{\partial f_{\overrightarrow{e}}^{t}}{\partial \mathbf{x}_{\overleftarrow{e}}} \left(\left(\mathbf{x}_{\overrightarrow{e}'}^{t} \right)_{\overrightarrow{e}': \overrightarrow{e}' \to \overrightarrow{e}} \right) \qquad \in \mathbb{R} \,. \tag{18}$$

The graph-based AMP iteration

$$\mathbf{x}_{\overrightarrow{e}}^{t+1} = \mathbf{A}_{\overrightarrow{e}} \mathbf{m}_{\overrightarrow{e}}^{t} - b_{\overrightarrow{e}}^{t} \mathbf{m}_{\overleftarrow{e}}^{t-1}, \qquad (16)$$

$$\mathbf{m}_{\vec{e}}^{t} = f_{\vec{e}}^{t} \left(\left(\mathbf{x}_{\vec{e}'}^{t} \right)_{\vec{e}': \vec{e}' \to \vec{e}} \right) \,, \tag{17}$$

where $b_{\overrightarrow{e}}^t$ is the so-called *Onsager term*

$$b_{\overrightarrow{e}}^{t} = \frac{1}{N} \operatorname{Tr} \frac{\partial f_{\overrightarrow{e}}^{t}}{\partial \mathbf{x}_{\overleftarrow{e}}} \left(\left(\mathbf{x}_{\overrightarrow{e}'}^{t} \right)_{\overrightarrow{e}': \overrightarrow{e}' \to \overrightarrow{e}} \right) \qquad \in \mathbb{R} \,. \tag{18}$$

Theorem (informal):

any Graph-based AMP iterations admits rigorous SE equations

Graph-based AMP: idea of the proof

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

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



and prove its SE equations.

- SE+converging trajectory \rightarrow properties of estimators
- design of the iteration is key (not discussed here)
- sample from non-convex priors/posteriors, multilayer, etc ...
- study the dynamics of learning algorithms, notably SGD [M.Celentano, C. Cheng, A. Montanari '21]

Thank you

Collaborators : Bruno Loureiro, Gabriele Sicuro, Raphaël Berthier, Lenka Zdeborova and Florent Krzakala