

Exact asymptotics and universality for gradient flows and other first order algorithms

Andrea Montanari

Stanford University

January 25, 2022



Michael Celentano



Chen Cheng



Yuchen Wu

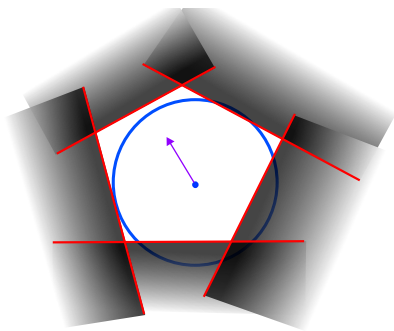
Toy problem

Find $\boldsymbol{\theta} \in \mathbb{R}^d$, such that $\|\boldsymbol{\theta}\|_2 = 1$, and

$$\begin{cases} \langle \boldsymbol{x}_1, \boldsymbol{\theta} \rangle \geq \kappa, \\ \vdots \\ \langle \boldsymbol{x}_n, \boldsymbol{\theta} \rangle \geq \kappa. \end{cases}$$

$$(\kappa < 0)$$

Toy problem ('Negative perceptron')



$$(\mathbf{x}_i)_{i \leq n} \sim \mathbf{N}(0, \mathbf{I}_d),$$

$$\mathcal{E}_{n,d}(\kappa) := \{\boldsymbol{\theta} \in \mathbb{S}^{d-1} : \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle \geq \kappa \quad \forall i \leq n\}.$$

Franz, Parisi 2016 (Physics); El Alaoui, Sellke, 2020; Baldi, Vershynin, 2021 (Math)

Toy problem ('Negative perceptron')

$$(\mathbf{x}_i)_{i \leq n} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d),$$

$$\mathcal{E}_{n,d}(\kappa) := \{\boldsymbol{\theta} \in \mathbb{S}^{d-1} : \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle \geq \kappa \forall i \leq n\}.$$

Questions

- ▶ is $\mathcal{E}_{n,d}(\kappa)$ non-empty?
- ▶ What is its geometry?
- ▶ Can we find $\boldsymbol{\theta} \in \mathcal{E}_{n,d}(\kappa)$ in polynomial time?

Gradient flow

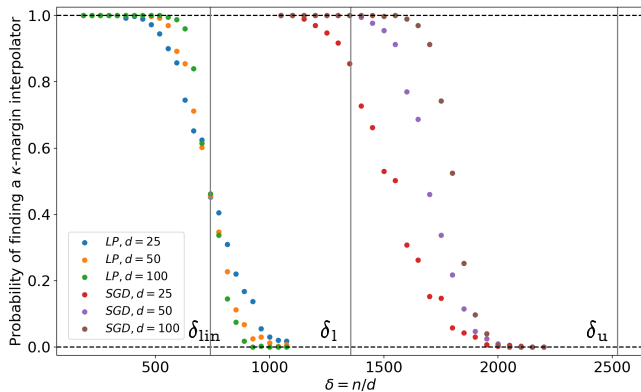
$$\begin{cases} \text{find} & \theta \in \mathbb{S}^{d-1}, \\ \text{such that} & \langle x_i, \theta \rangle \geq \kappa \quad \forall i \leq n \end{cases}$$

$$\dot{\theta}^t = \sum_{i=1}^n \rho(\langle x_i, \theta^t \rangle - \kappa \|\theta^t\|_2) (x_i - \kappa \hat{\theta}^t),$$

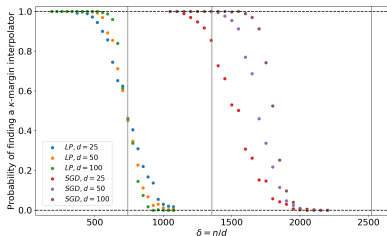
$$\hat{\theta}^t = \frac{\theta^t}{\|\theta^t\|_2}, \quad \rho(x) := \frac{1}{1 + e^x}.$$

If $\theta^t \rightarrow \infty$, then $\hat{\theta}^t \rightarrow \mathcal{E}_{n,d}(\kappa)$.

An experiment ($\kappa = -3$)



... and a theorem



Theorem (Zhong, Zhou, M, 2021)

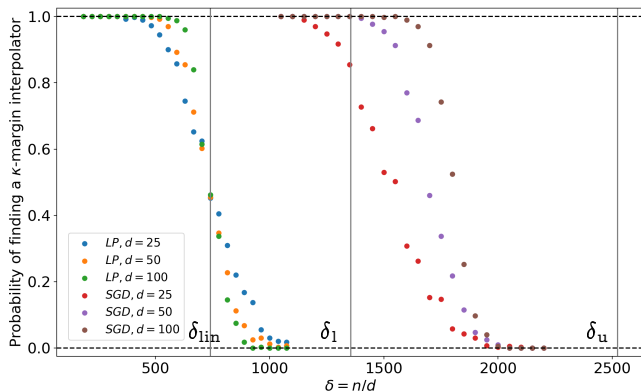
Assume $n, d \rightarrow \infty$ with $n, d \rightarrow \infty$ and define $\delta_s(\kappa)$ via

$$\delta < \delta_s(\kappa) \Leftrightarrow \liminf_{n \rightarrow \infty} \mathbb{P}_{n, n/\delta}(\mathcal{E}_{n, d}(\kappa) \neq \emptyset) > 0$$

Then

- ▶ $\delta_l(\kappa) \leq \delta_s(\kappa) \leq \delta_u(\kappa)$, with $\delta_l(\kappa), \delta_u(\kappa) = \Phi(\kappa)^{-1} \log |\kappa| (1 + o_\kappa(1))$.
- ▶ Linear prog. succeeds whp if $\delta < \delta_{\text{lin}}(\kappa) = \Phi(\kappa)^{-1} (1 + o_\kappa(1))$.

Questions



- ▶ Does gradient flow have a sharp threshold $\delta_{\text{GF}}(\kappa)$?
- ▶ Is $\delta_{\text{GF}}(\kappa) < \delta_s(\kappa)$?
- ▶ ...

More ambitious (2-layer neural nets)

$\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d)$, $y_i \sim \text{Unif}(\{+1, -1\})$, $a_\ell \in \{+1, -1\}$:

$$\begin{cases} \text{find} & \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k \in \mathbb{S}^{d-1}, \\ \text{such that} & y_i \sum_{\ell=1}^k a_\ell \sigma(\langle \boldsymbol{\theta}_\ell, \mathbf{x}_i \rangle) \geq \kappa \quad \forall i \leq n \end{cases}$$

- ▶ What is the interpolation threshold $\delta_s(k, \kappa)$?
- ▶ Does gradient flow have a sharp threshold $\delta_{\text{GF}}(k, \kappa)$?
- ▶ Is $\delta_{\text{GF}}(k, \kappa) < \delta_s(k, \kappa)$?
- ▶ ...

General first order flows

General First Order Flows

Data

$$\mathbf{X} = \begin{bmatrix} \cdots & \mathbf{x}_1^\top & \cdots \\ \cdots & \mathbf{x}_2^\top & \cdots \\ & \vdots & \\ \cdots & \mathbf{x}_n^\top & \cdots \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}.$$

Variable: $\theta \in \mathbb{R}^{d \times k}$

Flow

$$\frac{d\theta^t}{dt} = -\theta^t \Lambda^t - \mathbf{X}^\top \ell_t(\mathbf{X} \theta^t; \mathbf{y}),$$

General First Order Flows

Data

$$\mathbf{X} = \begin{bmatrix} \cdots & \mathbf{x}_1^\top & \cdots \\ \cdots & \mathbf{x}_2^\top & \cdots \\ & \vdots & \\ \cdots & \mathbf{x}_n^\top & \cdots \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}.$$

Variable: $\boldsymbol{\theta} \in \mathbb{R}^{d \times k}$

Flow

$$\frac{d\boldsymbol{\theta}^t}{dt} = -\boldsymbol{\theta}^t \Lambda^t - \mathbf{X}^\top \ell_t(\mathbf{X} \boldsymbol{\theta}^t; \mathbf{y}),$$

General First Order Flows

Data

$$\mathbf{X} = \begin{bmatrix} \cdots & \mathbf{x}_1^\top & \cdots \\ \cdots & \mathbf{x}_2^\top & \cdots \\ & \vdots & \\ \cdots & \mathbf{x}_n^\top & \cdots \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}.$$

Variable: $\boldsymbol{\theta} \in \mathbb{R}^{d \times k}$

Flow

$$\frac{d\boldsymbol{\theta}^t}{dt} = -\boldsymbol{\theta}^t \Lambda^t - \mathbf{X}^\top \ell_t(\mathbf{X} \boldsymbol{\theta}^t; \mathbf{y}),$$

General First Order Flows

$$\frac{d\boldsymbol{\theta}^t}{dt} = -\boldsymbol{\theta}^t \Lambda^t - \mathbf{X}^\top \boldsymbol{\ell}_t(\mathbf{X} \boldsymbol{\theta}^t; \mathbf{y}),$$

$$\Lambda : \mathbb{R} \rightarrow \mathbb{R}^{k \times k},$$

$$t \mapsto \Lambda_t$$

$$\boldsymbol{\ell}_t : \mathbb{R}^k \times \mathbb{R} \mapsto \mathbb{R}^k,$$

$$(\mathbf{u}, y) \mapsto \boldsymbol{\ell}_t(\mathbf{u}, y) \quad (\text{applies componentwise.})$$

- ▶ More general than gradient flow!
- ▶ Could add white noise term.

General First Order Flows

$$\frac{d\boldsymbol{\theta}^t}{dt} = -\boldsymbol{\theta}^t \Lambda^t - \mathbf{X}^\top \boldsymbol{\ell}_t(\mathbf{X} \boldsymbol{\theta}^t; \mathbf{y}),$$

$$\Lambda : \mathbb{R} \rightarrow \mathbb{R}^{k \times k},$$

$$t \mapsto \Lambda_t$$

$$\boldsymbol{\ell}_t : \mathbb{R}^k \times \mathbb{R} \mapsto \mathbb{R}^k,$$

$$(\mathbf{u}, y) \mapsto \boldsymbol{\ell}_t(\mathbf{u}, y) \quad (\text{applies componentwise.})$$

- ▶ More general than gradient flow!
- ▶ Could add white noise term.

General First Order Flows

$$\frac{d\boldsymbol{\theta}^t}{dt} = -\boldsymbol{\theta}^t \Lambda^t - \mathbf{X}^\top \ell_t(\mathbf{X}\boldsymbol{\theta}^t; \mathbf{y}).$$

Example #1: Negative perceptron

$$\dot{\boldsymbol{\theta}}^t = \sum_{i=1}^n \left(\mathbf{x}_i - \kappa \hat{\boldsymbol{\theta}}^t \right) \rho(\langle \mathbf{x}_i, \boldsymbol{\theta}^t \rangle - \kappa \|\boldsymbol{\theta}^t\|_2), \quad \hat{\boldsymbol{\theta}}^t := \frac{\boldsymbol{\theta}^t}{\|\boldsymbol{\theta}^t\|_2},$$

$$\dot{\boldsymbol{\theta}}^t = \mathbf{X}^\top \rho(\mathbf{X}\boldsymbol{\theta}^t - u_t \mathbf{1}) - \lambda_t \boldsymbol{\theta}^t,$$

$u_t, \lambda_t \in \mathbb{R}$ concentrate.

General First Order Flows

$$\frac{d\boldsymbol{\theta}^t}{dt} = -\boldsymbol{\theta}^t \Lambda^t - \mathbf{X}^\top \ell_t(\mathbf{X}\boldsymbol{\theta}^t; \mathbf{y}).$$

Example #1: Negative perceptron

$$\dot{\boldsymbol{\theta}}^t = \sum_{i=1}^n \left(\mathbf{x}_i - \kappa \hat{\boldsymbol{\theta}}^t \right) \rho(\langle \mathbf{x}_i, \boldsymbol{\theta}^t \rangle - \kappa \|\boldsymbol{\theta}^t\|_2), \quad \hat{\boldsymbol{\theta}}^t := \frac{\boldsymbol{\theta}^t}{\|\boldsymbol{\theta}^t\|_2},$$

$$\dot{\boldsymbol{\theta}}^t = \mathbf{X}^\top \rho(\mathbf{X}\boldsymbol{\theta}^t - u_t \mathbf{1}) - \lambda_t \boldsymbol{\theta}^t,$$

$u_t, \lambda_t \in \mathbb{R}$ concentrate.

General First Order Flows

$$\frac{d\boldsymbol{\theta}^t}{dt} = -\boldsymbol{\theta}^t \Lambda^t - \mathbf{X}^\top \ell_t(\mathbf{X}\boldsymbol{\theta}^t; \mathbf{y}).$$

Example #2: 2-layer neural network

$$\begin{aligned}\dot{\boldsymbol{\theta}}_j^t &= -\lambda_t \boldsymbol{\theta}_j^t - \nabla_{\boldsymbol{\theta}_j} \widehat{R}_n(\boldsymbol{\theta}^t) \quad j \leq k \\ \widehat{R}_n(\boldsymbol{\theta}^t) &= \frac{1}{n} \sum_{i=1}^n \left(y_i - \sum_{j \leq k} a_j \sigma(\langle \boldsymbol{\theta}_j^t, \mathbf{x}_i \rangle) \right)^2.\end{aligned}$$

General First Order Flows: Assumptions

$$\frac{d\boldsymbol{\theta}^t}{dt} = -\boldsymbol{\theta}^t \Lambda^t - \mathbf{X}^\top \ell_t(\mathbf{X} \boldsymbol{\theta}^t; \mathbf{y}).$$

- ▶ $n, d \rightarrow \infty$, $n/d \rightarrow \delta$, k fixed, $t = O(1)$
- ▶ $\{X_{ij}\}$ iid $\mathbb{E}\{X_{ij}\} = 0$, $\mathbb{E}\{X_{ij}^2\} = 1/d$, subgaussian.
- ▶ ℓ differentiable, with bounded Lipschitz Jacobian.
- ▶ $y_i = f(\boldsymbol{\theta}_*^\top \mathbf{x}_i; z_i)$, $\boldsymbol{\theta}_* \in \mathbb{R}^{d \times k_0}$, z_i independent of \mathbf{x}_i .
- ▶ Empirical distribution of rows of $\boldsymbol{\theta}^0$, $\boldsymbol{\theta}_*$ converges.

Our approach

- 1 Discretize time (first order methods)
- 2 Reduction to Approximate Message Passing¹ (AMP)
- 3 Apply existing asymptotic characterizations of AMP²
- 4 Take the continuous time limit.

¹Celentano, M, Wu, 2021

²Bayati, M, 2011; Chen, Lam, 2021

Our approach

- 1 Discretize time (first order methods)
- 2 Reduction to Approximate Message Passing³ (AMP)
- 3 Apply existing asymptotic characterizations of AMP⁴
- 4 Take the continuous time limit.

³Celentano, M, Wu, 2021

⁴Bayati, M, 2011; Chen, Lam, 2021

Time discretization

$$\frac{d\boldsymbol{\theta}^t}{dt} = -\boldsymbol{\theta}^t \Lambda^t - \mathbf{X}^\top \ell_t(\mathbf{X} \boldsymbol{\theta}^t; \mathbf{y}).$$

Euler scheme

$$\boldsymbol{\theta}_\varepsilon^{t+\varepsilon} = \boldsymbol{\theta}_\varepsilon^t - \varepsilon \boldsymbol{\theta}_\varepsilon^t \Lambda^t - \varepsilon \mathbf{X}^\top \ell_t(\mathbf{X} \boldsymbol{\theta}_\varepsilon^t; \mathbf{y}).$$

Claim

$$\lim_{\varepsilon \rightarrow 0} \limsup_{n, d \rightarrow \infty} \sup_{t \leq T} \frac{1}{n} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\varepsilon^t\|_2^2 = 0$$

Time discretization

$$\frac{d\boldsymbol{\theta}^t}{dt} = -\boldsymbol{\theta}^t \Lambda^t - \mathbf{X}^\top \ell_t(\mathbf{X} \boldsymbol{\theta}^t; \mathbf{y}).$$

Euler scheme

$$\boldsymbol{\theta}_\varepsilon^{t+\varepsilon} = \boldsymbol{\theta}_\varepsilon^t - \varepsilon \boldsymbol{\theta}_\varepsilon^t \Lambda^t - \varepsilon \mathbf{X}^\top \ell_t(\mathbf{X} \boldsymbol{\theta}_\varepsilon^t; \mathbf{y}).$$

Claim

$$\lim_{\varepsilon \rightarrow 0} \limsup_{n, d \rightarrow \infty} \sup_{t \leq T} \frac{1}{n} \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_\varepsilon^t\|_2^2 = 0$$

Time discretization

$$\boldsymbol{\theta}_\varepsilon^{t+\varepsilon} = \boldsymbol{\theta}_\varepsilon^t - \varepsilon \boldsymbol{\theta}_\varepsilon^t \Lambda^t - \varepsilon \mathbf{X}^\top \ell_t(\mathbf{X} \boldsymbol{\theta}_\varepsilon^t; \mathbf{y}).$$

Will drop \mathbf{y} hereafter

General First Order Methods

$$\begin{aligned}\mathbf{u}^{t+1} &= \mathbf{X}f_t(\mathbf{v}^1, \dots, \mathbf{v}^t), \\ \mathbf{v}^t &= \mathbf{X}^\top g_t(\mathbf{u}^1, \dots, \mathbf{u}^t), \\ \hat{\boldsymbol{\theta}}^t &= h_t(\mathbf{u}^1, \dots, \mathbf{u}^t).\end{aligned}$$

- ▶ $\mathbf{u}^t \in \mathbb{R}^{d \times k}$, $\mathbf{v}^t \in \mathbb{R}^{n \times k}$.
- ▶ Each iteration multiplication by \mathbf{X} or \mathbf{X}^\top .
- ▶ f_t, g_t, h_t Lipschitz, act entry-wise⁵.

Claim:

Euler scheme is a GFOM.

⁵More general assumptions, see M, Wu, 2022

General First Order Methods

$$\begin{aligned}\mathbf{u}^{t+1} &= \mathbf{X}f_t(\mathbf{v}^1, \dots, \mathbf{v}^t), \\ \mathbf{v}^t &= \mathbf{X}^\top g_t(\mathbf{u}^1, \dots, \mathbf{u}^t), \\ \hat{\boldsymbol{\theta}}^t &= h_t(\mathbf{u}^1, \dots, \mathbf{u}^t).\end{aligned}$$

- ▶ $\mathbf{u}^t \in \mathbb{R}^{d \times k}$, $\mathbf{v}^t \in \mathbb{R}^{n \times k}$.
- ▶ Each iteration multiplication by \mathbf{X} or \mathbf{X}^\top .
- ▶ f_t, g_t, h_t Lipschitz, act entry-wise⁵.

Claim:

Euler scheme is a GFOM.

⁵More general assumptions, see M, Wu, 2022

Approximate Message Passing (AMP)

$$\begin{aligned} \mathbf{u}^{t+1} &= \mathbf{X} \psi_t(\mathbf{v}^1, \dots, \mathbf{v}^t) - \sum_{s=1}^t \phi_s(\mathbf{u}^1, \dots, \mathbf{u}^s) \cdot \mathbf{a}_{t,s}, \\ \mathbf{v}^t &= \mathbf{X}^\top \phi_t(\mathbf{u}^1, \dots, \mathbf{u}^t) - \sum_{s=1}^{t-1} \psi_s(\mathbf{v}^1, \dots, \mathbf{v}^s) \cdot \mathbf{b}_{t,s}, \\ \hat{\boldsymbol{\theta}}^t &= \boldsymbol{\vartheta}_t(\mathbf{u}^1, \dots, \mathbf{u}^t). \end{aligned}$$

- ▶ $\psi_t, \phi_t, \boldsymbol{\vartheta}_t$ Lipschitz, act entry-wise⁶.
- ▶ $\mathbf{a}_{t,s}, \mathbf{b}_{t,s} \in \mathbb{R}^{k \times k}$ deterministic, with explicit formulas.
- ▶ Key property (state evolution)

$$\frac{1}{d} \sum_{i=1}^d \delta_{\mathbf{u}_i^1, \dots, \mathbf{u}_i^t} \Rightarrow \mathbf{N}(0, \boldsymbol{\Sigma}_{\leq t}), \quad +\text{recursion for } \boldsymbol{\Sigma}.$$

⁶More general assumptions, see Berthier, M, Nguyen 2018; M, Wu, 2022

Approximate Message Passing (AMP)

$$\begin{aligned}\mathbf{u}^{t+1} &= \mathbf{X}\psi_t(\mathbf{v}^1, \dots, \mathbf{v}^t) - \sum_{s=1}^t \phi_s(\mathbf{u}^1, \dots, \mathbf{u}^s) \cdot \mathbf{a}_{t,s}, \\ \mathbf{v}^t &= \mathbf{X}^\top \phi_t(\mathbf{u}^1, \dots, \mathbf{u}^t) - \sum_{s=1}^{t-1} \psi_s(\mathbf{v}^1, \dots, \mathbf{v}^s) \cdot \mathbf{b}_{t,s}, \\ \hat{\boldsymbol{\theta}}^t &= \boldsymbol{\vartheta}_t(\mathbf{u}^1, \dots, \mathbf{u}^t).\end{aligned}$$

Remarks

- ▶ Any AMP algorithm is a GFOM
- ▶ Any GFOM is equivalent to an AMP algorithm⁷
(With suitable post-processing $\boldsymbol{\vartheta}_t$)

⁷Celentano, M, Wu, 2021

Approximate Message Passing (AMP)

$$\begin{aligned}\mathbf{u}^{t+1} &= \mathbf{X}\psi_t(\mathbf{v}^1, \dots, \mathbf{v}^t) - \sum_{s=1}^t \phi_s(\mathbf{u}^1, \dots, \mathbf{u}^s) \cdot \mathbf{a}_{t,s}, \\ \mathbf{v}^t &= \mathbf{X}^\top \phi_t(\mathbf{u}^1, \dots, \mathbf{u}^t) - \sum_{s=1}^{t-1} \psi_s(\mathbf{v}^1, \dots, \mathbf{v}^s) \cdot \mathbf{b}_{t,s}, \\ \hat{\boldsymbol{\theta}}^t &= \boldsymbol{\vartheta}_t(\mathbf{u}^1, \dots, \mathbf{u}^t).\end{aligned}$$

Remarks

- ▶ Any AMP algorithm is a GFOM
- ▶ Any GFOM is equivalent to an AMP algorithm⁷
(With suitable post-processing $\boldsymbol{\vartheta}_t$)

⁷Celentano, M, Wu, 2021

DMFT asymptotics

Dynamical Mean Field Theory

$$\frac{d\boldsymbol{\theta}^t}{dt} = -\boldsymbol{\theta}^t \Lambda^t - \mathbf{X}^\top \ell_t(\mathbf{X} \boldsymbol{\theta}^t; \mathbf{y}).$$

General form: As $n, d \rightarrow \infty$,

$$(\boldsymbol{\theta}_i^t)_{t \leq T} \Rightarrow (\boldsymbol{\theta}^t)_{t \leq T}$$

where $(\boldsymbol{\theta}^t)_{t \leq T}$ is a diffusion process with memory.

Dynamical Mean Field Theory: History

Physics

- ▶ Sompolinskiy, Zippelius, 1981 (SK model)
- ▶ Crisanti, Horner, Sommers, 1993; Cugliandolo, Kurchan, 1993 (spherical spin glass)
- ▶ ...

Mathematics

- ▶ Ben Arous and Guionnet, 1995 (SK model)
- ▶ Ben Arous, Dembo, Guionnet 2006 (spherical spin glass)
- ▶ ... (Gaussian disorder)

Dynamical Mean Field Theory: Recent work

(Heuristic) Applications in high-dimensional statistics

- ▶ Agoritsas et al., 2018 (Perceptrons)
- ▶ Mannelli et al. 2020 (Tensor PCA)
- ▶ Mignacco et al. 2020 (Gaussian mixtures)
- ▶ ...

Universality

- ▶ Dembo, Lubetzky, Zeitouni, 2019 (Asymmetric interactions)
- ▶ Dembo, Gheissari, 2021 (Interacting diffusions)

Present work yields universality

Dynamical Mean Field Theory: Recent work

(Heuristic) Applications in high-dimensional statistics

- ▶ Agoritsas et al., 2018 (Perceptrons)
- ▶ Mannelli et al. 2020 (Tensor PCA)
- ▶ Mignacco et al. 2020 (Gaussian mixtures)
- ▶ ...

Universality

- ▶ Dembo, Lubetzky, Zeitouni, 2019 (Asymmetric interactions)
- ▶ Dembo, Gheissari, 2021 (Interacting diffusions)

Present work yields universality

Dynamical Mean Field Theory: Recent work

(Heuristic) Applications in high-dimensional statistics

- ▶ Agoritsas et al., 2018 (Perceptrons)
- ▶ Mannelli et al. 2020 (Tensor PCA)
- ▶ Mignacco et al. 2020 (Gaussian mixtures)
- ▶ ...

Universality

- ▶ Dembo, Lubetzky, Zeitouni, 2019 (Asymmetric interactions)
- ▶ Dembo, Gheissari, 2021 (Interacting diffusions)

Present work yields universality

DMFT process

$$\frac{d\theta^t}{dt} = -\theta^t \Lambda^t - \frac{1}{\delta} \mathbf{X}^\top \ell_t(\mathbf{X} \theta^t), \quad \theta^t \in \mathbb{R}^{d \times k}.$$

DMFT process ($\theta^t, r^t \in \mathbb{R}^k$)

$$\begin{aligned} \frac{d}{dt} \theta^t &= -(\Lambda^t + \Gamma^t) \theta^t - \int_0^t R_\ell(t, s) \theta^s ds + u^t, & u &\sim \text{GP}(0, C_\ell / \delta), \\ r^t &= -\frac{1}{\delta} \int_0^t R_\theta(t, s) \ell_s(r^s) ds + w^t, & w &\sim \text{GP}(0, C_\theta). \end{aligned}$$

DMFT: Fixed point conditions

DMFT process ($\theta^t, r^t \in \mathbb{R}^k$)

$$\frac{d}{dt}\theta^t = -(\Lambda^t + \Gamma^t)\theta^t - \int_0^t R_\ell(t, s)\theta^s ds + u^t, \quad u \sim \text{GP}(0, C_\ell/\delta),$$

$$r^t = -\frac{1}{\delta} \int_0^t R_\theta(t, s)l_s(r^s)ds + w^t, \quad w \sim \text{GP}(0, C_\theta).$$

Fixed point equations

$$C_\theta(t, s) = \mathbb{E} \left[\theta^t \theta^{s\top} \right], \quad R_\theta(t, s) = \mathbb{E} \left[\frac{\partial \theta^t}{\partial u^s} \right],$$

$$C_\ell(t, s) = \mathbb{E} \left[l_t(r^t) l_s(r^t)^\top \right], \quad R_\ell(t, s) = \mathbb{E} \left[\frac{\partial l_t(r^t)}{\partial w^s} \right],$$

$$\Gamma^t = \mathbb{E} \left[\nabla_r l_t(r^t; z) \right].$$

DMFT: Fixed point conditions

DMFT process ($\theta^t, r^t \in \mathbb{R}^k$)

$$\frac{d}{dt}\theta^t = -(\Lambda^t + \Gamma^t)\theta^t - \int_0^t R_\ell(t, s)\theta^s ds + u^t, \quad u \sim \text{GP}(0, C_\ell/\delta),$$

$$r^t = -\frac{1}{\delta} \int_0^t R_\theta(t, s)l_s(r^s)ds + w^t, \quad w \sim \text{GP}(0, C_\theta).$$

Fixed point equations

$$C_\theta(t, s) = \mathbb{E} \left[\theta^t \theta^{s\top} \right], \quad R_\theta(t, s) = \mathbb{E} \left[\frac{\partial \theta^t}{\partial u^s} \right],$$

$$C_\ell(t, s) = \mathbb{E} \left[l_t(r^t) l_s(r^t)^\top \right], \quad R_\ell(t, s) = \mathbb{E} \left[\frac{\partial l_t(r^t)}{\partial w^s} \right],$$

$$\Gamma^t = \mathbb{E} \left[\nabla_r l_t(r^t; z) \right].$$

DMFT: Fixed point conditions

Fixed point equations

$$\begin{aligned}C_\theta(t, s) &= \mathbb{E} \left[\theta^t \theta^{s\top} \right], & R_\theta(t, s) &= \mathbb{E} \left[\frac{\partial \theta^t}{\partial u^s} \right], \\C_\ell(t, s) &= \mathbb{E} \left[\ell_t(r^t) \ell_s(r^t)^\top \right], & R_\ell(t, s) &= \mathbb{E} \left[\frac{\partial \ell_t(r^t)}{\partial w^s} \right], \\ \Gamma^t &= \mathbb{E} \left[\nabla_r \ell_t(r^t; z) \right].\end{aligned}$$

$$(C_\theta, R_\theta) = \mathcal{T}(C_\ell, R_\ell, \Gamma), \quad (C_\ell, R_\ell, \Gamma) = \hat{\mathcal{T}}(C_\theta, R_\theta)$$

1st theorem: Existence and uniqueness

Theorem (Celentano, Chen, M, 2021)

- ▶ *The solution $(C_\theta, R_\theta, C_\ell, R_\ell, \Gamma)$ of the fixed point conditions exists.*
- ▶ *It is unique among all tuples such that (C_θ, R_θ) are bounded in all compact sets in $\mathbb{R}_{\geq 0}^2$.*
- ▶ *The DMFT processes $(\theta^t)_{t \geq 0}$, $(r^t)_{t \geq 0}$ are well defined with continuous sample paths.*
- ▶ *The map $\mathcal{T} \circ \hat{\mathcal{T}}$ is a contraction.*

2nd theorem: Convergence as $n, d \rightarrow \infty$

Theorem (Celentano, Chen, M, 2021)

Let d_W metrize weak convergence in $C([0, T], \mathbb{R}^k)$, and define $r^t := \mathbf{X}\theta^t \in \mathbb{R}^{n \times k}$.

Let $P_{\theta_0^T}, P_{r_0^T}$ be the laws of the DMFT processes.
Then, we have

$$\text{p-lim}_{n, d \rightarrow \infty} d_W \left(\frac{1}{d} \sum_{i=1}^d \delta_{(\theta_i)_0^T}, P_{\theta_0^T} \right) = 0,$$

$$\text{p-lim}_{n, d \rightarrow \infty} d_W \left(\frac{1}{n} \sum_{i=1}^n \delta_{(r_i)_0^T}, P_{r_0^T} \right) = 0.$$

Convergence as $n, d \rightarrow \infty$

Corollary

For any m , any $\psi : \mathbb{R}^m \rightarrow \mathbb{R}$ bounded continuous, and any $0 \leq t_1 \leq t_2 \leq \dots \leq t_m$, we have

$$\text{p-lim}_{n, d \rightarrow \infty} \frac{1}{d} \sum_{i=1}^d \psi(\theta_i^{t_1}, \dots, \theta_i^{t_m}) = \mathbb{E}\{\psi(\theta^{t_1}, \dots, \theta^{t_m})\}.$$

$$\underbrace{\frac{1}{d} \sum_{i=1}^d \psi(\theta_i^{t_1}, \dots, \theta_i^{t_m})}_{\text{avg over coordinates}} \xrightarrow{P} \underbrace{\mathbb{E}\{\psi(\theta^{t_1}, \dots, \theta^{t_m})\}}_{\text{expectation wrt DMFT}} .$$

Epilogue: Statistically optimal methods

Reduction

General First Order Methods \longrightarrow Approximate Message Passing

Q: Can we determine the optimal GFOM?

Statistical optimality

► Data

$$y_i = \varphi(\langle x_i, \theta_* \rangle, z_i), \quad x_i \sim N(0, I_d/d) \perp z_i \sim \text{Unif}([0, 1]).$$

► General GFOM

$$\begin{aligned} u^{t+1} &= X f_t(v^1, \dots, v^t), \\ v^t &= X^\top g_t(u^1, \dots, u^t), \\ \hat{\theta}^t &= h_t(u^1, \dots, u^t). \end{aligned}$$

$$\lim_{n, d \rightarrow \infty} \frac{1}{d} \inf_{\{f_s, g_s, h_s\} \in \mathcal{F}} \mathbb{E}\{\|\hat{\theta}^t - \theta_*\|_2^2\} = ???$$

Statistical optimality

► Data

$$y_i = \varphi(\langle \mathbf{x}_i, \boldsymbol{\theta}_* \rangle, z_i), \quad \mathbf{x}_i \sim \mathbf{N}(0, \mathbf{I}_d/d) \perp z_i \sim \text{Unif}([0, 1]).$$

► General GFOM

$$\begin{aligned} \mathbf{u}^{t+1} &= \mathbf{X} f_t(\mathbf{v}^1, \dots, \mathbf{v}^t), \\ \mathbf{v}^t &= \mathbf{X}^\top g_t(\mathbf{u}^1, \dots, \mathbf{u}^t), \\ \hat{\boldsymbol{\theta}}^t &= h_t(\mathbf{u}^1, \dots, \mathbf{u}^t). \end{aligned}$$

$$\lim_{n, d \rightarrow \infty} \frac{1}{d} \inf_{\{f_s, g_s, h_s\} \in \mathcal{F}} \mathbb{E}\{\|\hat{\boldsymbol{\theta}}^t - \boldsymbol{\theta}_*\|_2^2\} = ???$$

Statistical optimality

► Data

$$y_i = \varphi(\langle \mathbf{x}_i, \boldsymbol{\theta}_* \rangle, z_i), \quad \mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d/d) \perp z_i \sim \text{Unif}([0, 1]).$$

► General GFOM

$$\begin{aligned} \mathbf{u}^{t+1} &= \mathbf{X} f_t(\mathbf{v}^1, \dots, \mathbf{v}^t), \\ \mathbf{v}^t &= \mathbf{X}^\top g_t(\mathbf{u}^1, \dots, \mathbf{u}^t), \\ \hat{\boldsymbol{\theta}}^t &= h_t(\mathbf{u}^1, \dots, \mathbf{u}^t). \end{aligned}$$

$$\lim_{n, d \rightarrow \infty} \frac{1}{d} \inf_{\{f_s, g_s, h_s\} \in \mathcal{F}} \mathbb{E} \{ \|\hat{\boldsymbol{\theta}}^t - \boldsymbol{\theta}_*\|_2^2 \} = ???$$

Statistical optimality

► Data

$$y_i = \varphi(\langle \mathbf{x}_i, \boldsymbol{\theta}_* \rangle; z_i), \quad \mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d/d) \perp z_i \sim \text{Unif}([0, 1]).$$

► General GFOM

$$\begin{aligned} \mathbf{u}^{t+1} &= \mathbf{X} f_t(\mathbf{v}^1, \dots, \mathbf{v}^t), \\ \mathbf{v}^t &= \mathbf{X}^\top g_t(\mathbf{u}^1, \dots, \mathbf{u}^t), \\ \hat{\boldsymbol{\theta}}^t &= h_t(\mathbf{u}^1, \dots, \mathbf{u}^t). \end{aligned}$$

$$\lim_{n, d \rightarrow \infty} \frac{1}{d} \inf_{\{f_s, g_s, h_s\} \in \mathcal{F}} \mathbb{E} \{ \|\hat{\boldsymbol{\theta}}^t - \boldsymbol{\theta}_*\|_2^2 \} = ???$$

$\mathcal{F} :=$ separable Lipschitz functions

Example: Noiseless phase retrieval

$$y_i = \langle \mathbf{x}_i, \boldsymbol{\theta}_* \rangle^2$$

Spectral initialization⁸

$$D_n := \sum_{i=1}^n \mathcal{T}(y_i) \mathbf{x}_i \mathbf{x}_i^\top, \quad \mathcal{T}(y) = \frac{y-1}{y+\varepsilon},$$
$$\boldsymbol{\theta}^0 = c_{d,n} \cdot \mathbf{v}_1(D_n).$$

⁸Mondelli, Montanari, 2018

First order methods in phase retrieval

- ▶ Schniter, Rangan 2014
- ▶ Candés, Li, Soltanolkotabi, 2015
- ▶ Cai, Li, Ma, 2016
- ▶ Wang, Giannakis, Eldar, 2017
- ▶ Chen, Candés 2018
- ▶ Waldspurger, 2018
- ▶ Duchi, Ruan, 2019
- ▶ Maillard, Loureiro, Krzakala, Zdeborová, 2020
- ▶ Mondelli Venkatramanan, 2021
- ▶ **Review:** Fannjiang and Strohmer, 2020
- ▶ ...

Experiment with a real image: $d = 7560$



Original image.

Bayes AMP



$t = 2$



$t = 4$



$t = 8$

1 step prox-linear



$t = 2$

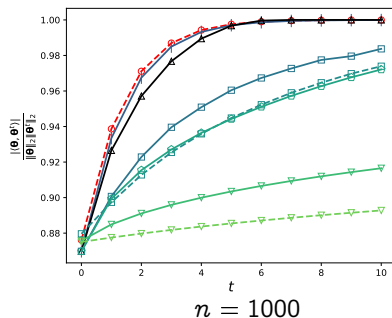
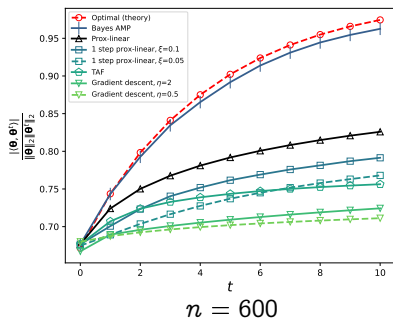


$t = 4$



$t = 8$

A simulation: Reconstruction accuracy $d = 400$



- ▶ Truncated Amplitude Flow
- ▶ Prox-linear⁹
- ▶ Gradient descent.
- ▶ Bayes-optimal AMP.

(Wang, Giannakis, Eldar, 2017)

(Duchi, Ruan, 2019)

⁹Not a GFOM.

Minimal error

- ▶ $y = \varphi(X\theta_*, z)$
- ▶ Side information v , $d^{-1} \sum_{i \leq d} \delta_{\theta_{*,i}, v_i} \Rightarrow P_{\Theta, V}$
- ▶ $n, d \rightarrow \infty$, $n/d \rightarrow \delta$ $\{X_{ij}\}$ bounded 4-th moment.

Define

$$\text{mmse}_{\Theta, V}(\alpha) := \mathbb{E}[\Theta^2] - \mathbb{E}\{\mathbb{E}[\Theta \mid \alpha\Theta + Z, V]^2\}.$$

and recursively:

$$\beta_s^2 = \frac{1}{\sigma_s^2} \mathbb{E}[\mathbb{E}[Z_0 \mid \varphi(\sigma_s Z_0 + \tilde{\sigma}_s Z_1, W), U, Z_1]^2], \quad \beta_s \geq 0,$$

$$\sigma_{s+1}^2 = \frac{1}{\delta} \text{mmse}_{\Theta, V}(\beta_s), \quad \tilde{\sigma}_{s+1}^2 = \frac{1}{\delta} (\mathbb{E}[\Theta^2] - \text{mmse}_{\Theta, V}(\beta_s)).$$

Minimal error

Theorem (Celentano, M, Wu, 2020; M, Wu, 2022)

For $t \in \mathbb{N}_{>0}$, let $\hat{\boldsymbol{\theta}}^t \in \mathbb{R}^d$ be the output of any GFOM after t iterations, the following holds:

$$\text{p-lim}_{n, d \rightarrow \infty} \frac{1}{d} \|\hat{\boldsymbol{\theta}}^t - \boldsymbol{\theta}_*\|_2^2 \geq \text{mmse}_{\Theta, V}(\beta_t).$$

Further, there exists a GFOM (Bayes AMP) which satisfies the above bound with equality.

Conclusion

Conclusion

First order methods comprise the most used algorithms in ML

- ▶ GFOMs as a useful abstraction/class.
- ▶ Reduction GFOM \rightarrow AMP
- ▶ Need to understand their behavior in high-dim/low SNR.
- ▶ New rigorous tool: DMFT (other available!)

Thanks!