

Dynamical phenomena in nonlinear learning

Andrea Montanari

Stanford University

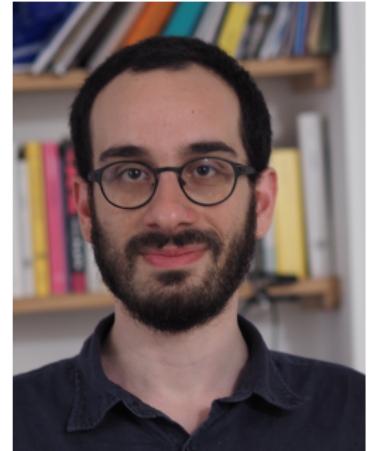
August 13, 2025



Michael Celentano



Chen Cheng



Pierfrancesco Urbani

Where we left yesterday

Two-layer neural networks

($m \equiv N$)

$$f(x; \theta) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\langle w_j, x \rangle), \quad \theta = ((a_i)_{i \leq m}, (w_i)_{i \leq m}) \in \mathbb{R}^m \times (\mathbb{S}^{d-1})^m.$$

- ▶ Square loss train error:

$$\hat{\mathcal{R}}_n(\theta) := \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i; \theta))^2.$$

- ▶ Square loss test error:

$$\mathcal{R}(\theta) := \frac{1}{2n} \mathbb{E}\{(y - f(x; \theta))^2\}.$$

Two-layer neural networks

($m \equiv N$)

$$f(x; \theta) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\langle w_j, x \rangle), \quad \theta = ((a_i)_{i \leq m}, (w_i)_{i \leq m}) \in \mathbb{R}^m \times (\mathbb{S}^{d-1})^m.$$

- ▶ Square loss train error:

$$\hat{\mathcal{R}}_n(\theta) := \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i; \theta))^2.$$

- ▶ Square loss test error:

$$\mathcal{R}(\theta) := \frac{1}{2n} \mathbb{E}\{(y - f(x; \theta))^2\}.$$

Single index model

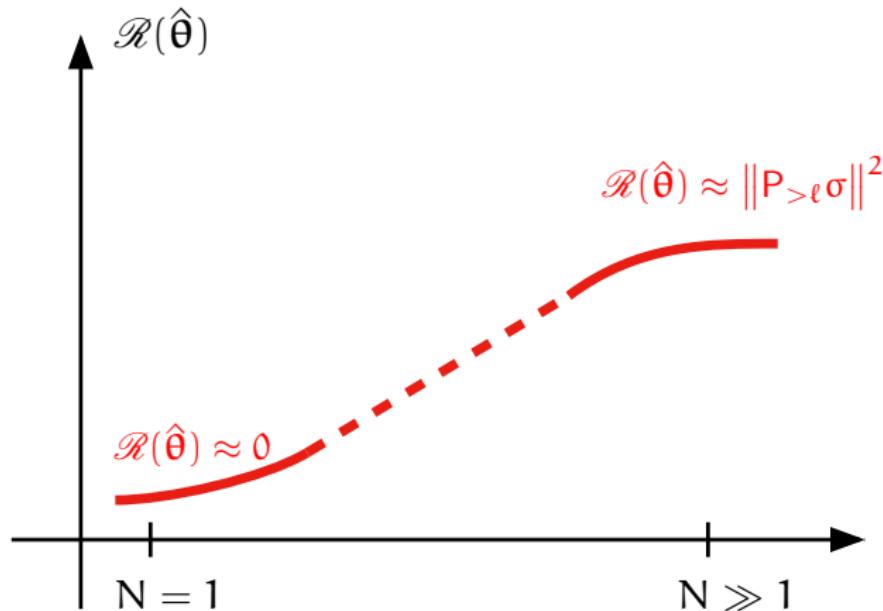
Data $\{(x_i, y_i) : i \leq n\}$ iid, $\varepsilon_i \sim N(0, \tau^2)$

$$x_i \sim N(0, I_d), \quad y_i = \varphi(\langle w_*, x_i \rangle) + \varepsilon_i$$

Two theories:

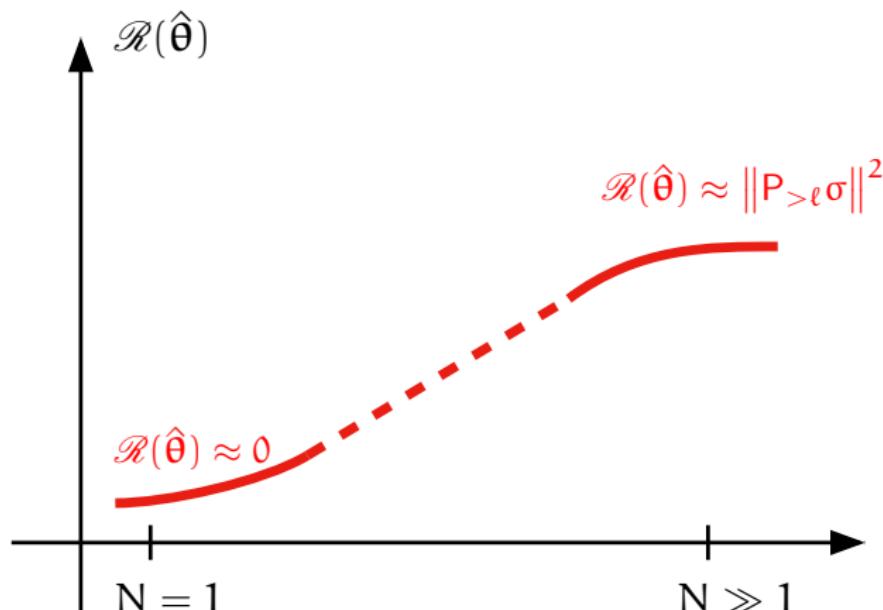
- ▶ Feature learning:
 - ▶ Weights w_i align quickly with w_*
 - ▶ No overfitting
- ▶ Neural Tangent:
 - ▶ Weights change minimally
 - ▶ Model is linear in the y_i (kernel method)
 - ▶ Overfitting/interpolation

A cartoon



$$f(x; \theta) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\langle w_j, x \rangle), \quad a_i^0 = \pm \gamma \sqrt{m}$$

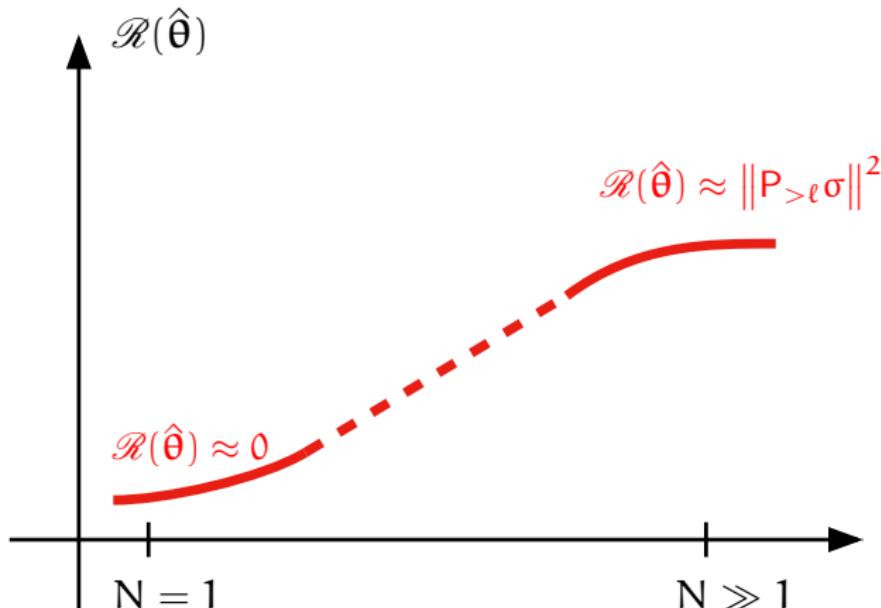
A cartoon



Are large networks necessarily bad?

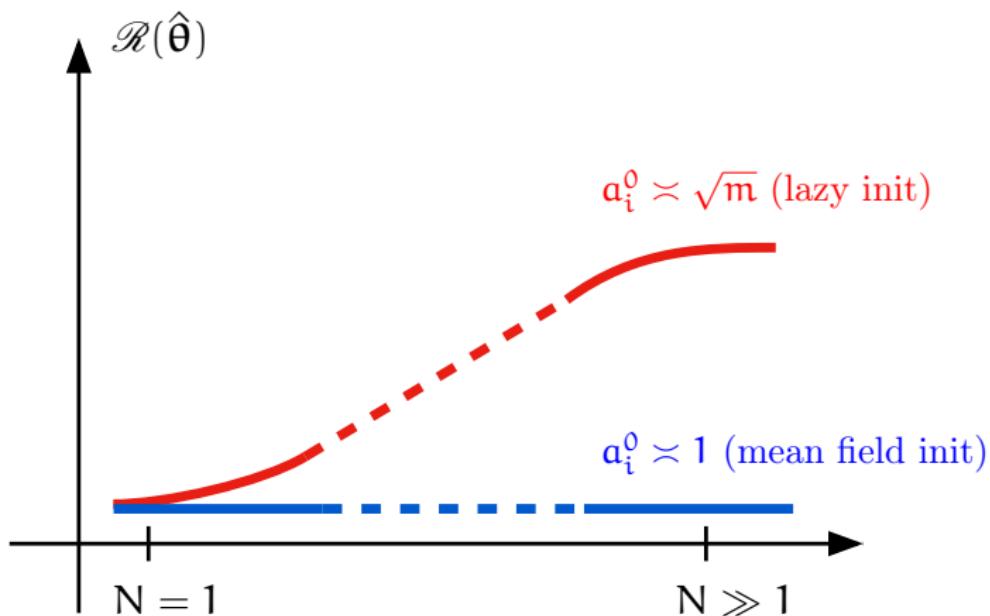
Spoiler

Spoiler



$$f(x; \theta) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\langle w_j, x \rangle), \quad a_i^0 = \pm \gamma \sqrt{m}$$

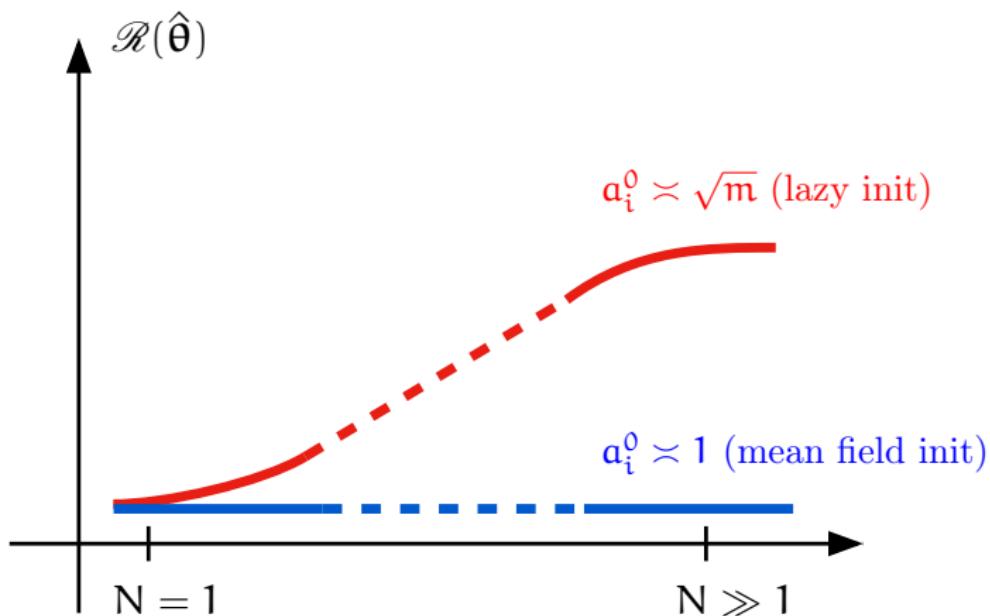
Spoiler



Good generalization: mean field initialization + early stopping

Lazy initialization is popular among practitioners.

Spoiler



Good generalization: mean field initialization + early stopping

Lazy initialization is popular among practitioners.

Key control parameters (large n, m, d)

t Training time,

$\alpha = \frac{n}{md}$ Overparametrization ratio,

a^0 Scale of 2nd layer weights at $t = 0$.

Questions

- Q1. For which region of α, a^0 does grad flow converge $\hat{R}_n \approx 0$?
- Q2. Does the selected model provide good generalization?
- Q3. When feature-learning/no-overfitting vs lazy-training/overfitting?
- Q4. How does the generalization error depend on network size and number of iterations?

Approach:

- ▶ Dynamical mean field theory (DMFT): $n, d \rightarrow \infty, n/d \rightarrow \bar{\alpha}$.
- ▶ Take limit $m, \bar{\alpha} \rightarrow \infty$, with $\bar{\alpha}/m \rightarrow \alpha$.

Questions

- Q1. For which region of α , a^0 does grad flow converge $\hat{R}_n \approx 0$?
- Q2. Does the selected model provide good generalization?
- Q3. When feature-learning/no-overfitting vs lazy-training/overfitting?
- Q4. How does the generalization error depend on network size and number of iterations?

Approach:

- ▶ Dynamical mean field theory (DMFT): $n, d \rightarrow \infty$, $n/d \rightarrow \bar{\alpha}$.
- ▶ Take limit $m, \bar{\alpha} \rightarrow \infty$, with $\bar{\alpha}/m \rightarrow \alpha$.

Questions

- Q1. For which region of α , a^0 does grad flow converge $\hat{R}_n \approx 0$?
- Q2. Does the selected model provide good generalization?
- Q3. When feature-learning/no-overfitting vs lazy-training/overfitting?
- Q4. How does the generalization error depend on network size and number of iterations?

Approach:

- ▶ Dynamical mean field theory (DMFT): $n, d \rightarrow \infty, n/d \rightarrow \bar{\alpha}$.
- ▶ Take limit $m, \bar{\alpha} \rightarrow \infty$, with $\bar{\alpha}/m \rightarrow \alpha$.

Questions

- Q1. For which region of α , a^0 does grad flow converge $\hat{R}_n \approx 0$?
- Q2. Does the selected model provide good generalization?
- Q3. When feature-learning/no-overfitting vs lazy-training/overfitting?
- Q4. How does the generalization error depend on network size and number of iterations?

Approach:

- ▶ Dynamical mean field theory (DMFT): $n, d \rightarrow \infty, n/d \rightarrow \bar{\alpha}$.
- ▶ Take limit $m, \bar{\alpha} \rightarrow \infty$, with $\bar{\alpha}/m \rightarrow \alpha$.

Questions

- Q1. For which region of α , a^0 does grad flow converge $\hat{R}_n \approx 0$?
- Q2. Does the selected model provide good generalization?
- Q3. When feature-learning/no-overfitting vs lazy-training/overfitting?
- Q4. How does the generalization error depend on network size and number of iterations?

Approach:

- ▶ Dynamical mean field theory (DMFT): $n, d \rightarrow \infty$, $n/d \rightarrow \bar{\alpha}$.
- ▶ Take limit $m, \bar{\alpha} \rightarrow \infty$, with $\bar{\alpha}/m \rightarrow \alpha$.

Outline

- 1 Rigorous DMFT
- 2 Gaussian approximation
- 3 Dynamics in pure noise
- 4 Dynamics under single-index model
- 5 Conclusion

Rigorous DMFT

A slightly more general setting

$\theta, \theta_* \in \mathbb{R}^{d \times m}$, regularization $t \mapsto \Lambda_t \in \mathbb{R}^{m \times m}$

$$\begin{aligned}\widehat{\mathcal{R}}_n(\theta) &= \frac{1}{n} \sum_{i=1}^n L(\theta^\top x_i; \theta_*^\top x_i, \varepsilon_i), \\ \dot{\theta}_t &= -\theta_t \Lambda_t^\top - \nabla \widehat{\mathcal{R}}_n(\theta_t).\end{aligned}$$

The case of 2-layer neural nets, k–index data:

$$\theta = [w_1 | w_2 | \cdots | w_m], \quad \theta_* = [w_{*,1} | \cdots | w_k | 0 | \cdots | 0],$$

$$L(r, u, \varepsilon) = \frac{1}{2} \left(\varphi(u) + \varepsilon - \frac{1}{m} \sum_{i=1}^m a_i \sigma(r_i) \right)^2.$$

A slightly more general setting

$\theta, \theta_* \in \mathbb{R}^{d \times m}$, regularization $t \mapsto \Lambda_t \in \mathbb{R}^{m \times m}$

$$\begin{aligned}\widehat{\mathcal{R}}_n(\theta) &= \frac{1}{n} \sum_{i=1}^n L(\theta^\top x_i; \theta_*^\top x_i, \varepsilon_i), \\ \dot{\theta}_t &= -\theta_t \Lambda_t^\top - \nabla \widehat{\mathcal{R}}_n(\theta_t).\end{aligned}$$

The case of 2-layer neural nets, k–index data:

$$\theta = [w_1 | w_2 | \cdots | w_m], \quad \theta_* = [w_{*,1} | \cdots | w_k | 0 | \cdots | 0],$$

$$L(r, u, \varepsilon) = \frac{1}{2} \left(\varphi(u) + \varepsilon - \frac{1}{m} \sum_{i=1}^m a_i \sigma(r_i) \right)^2.$$

A slightly more general setting

$\theta, \theta_* \in \mathbb{R}^{d \times m}$, regularization $t \mapsto \Lambda_t \in \mathbb{R}^{m \times m}$

$$\begin{aligned}\widehat{\mathcal{R}}_n(\theta) &= \frac{1}{n} \sum_{i=1}^n L(\theta^T x_i; \theta_*^T x_i, \varepsilon_i), \\ \dot{\theta}_t &= -\theta_t \Lambda_t^T - \nabla \widehat{\mathcal{R}}_n(\theta_t).\end{aligned}$$

DMFT process: $\vartheta_t, r_t \in \mathbb{R}^m$,

$$\begin{aligned}\dot{\vartheta}^t &= -(\Lambda^t + \Gamma^t)\vartheta^t - \int_0^t R_\ell(t, s)\vartheta^s ds - R_\ell(t, *)\vartheta^* + u^t, \\ r^t &= -\frac{1}{\alpha} \int_0^t R_\theta(t, s) \nabla_r L(r^s, w^*; \varepsilon) ds + w^t, \\ u^t &\sim \text{GP}(0, C_\ell/\delta), \quad w^t \sim \text{GP}(0, C_\theta).\end{aligned}$$

where $\Gamma_t, C_\ell(t, s), C_\theta(t, s), R_\ell(t, s), R_\ell(t, *), R_\theta(t, s) \in \mathbb{R}^{m \times m}$ are unique solution of ...

Assumptions

- ▶ \mathbf{X} has standardized sub-Gaussian independent entries
- ▶ ℓ is C^3 with bounded 2-nd, 3-rd derivatives.
- ▶ $n, d \rightarrow \infty$, $n/d \rightarrow \bar{\alpha}$
- ▶ $\theta_{0,i}, \theta_{*,i}$, i -th rows of θ_0, θ_* :

$$\frac{1}{d} \sum_{i=1}^d \delta_{\sqrt{d}\theta_{0,i}, \sqrt{d}\theta_{*,i}} \xrightarrow{W_2} \text{Law}(\vartheta_0, \vartheta_*) , \quad \frac{1}{n} \sum_{i=1}^n \delta_{\varepsilon_i} \xrightarrow{W_2} \text{Law}(\varepsilon) .$$

DMFT characterization

Theorem (Celentano, Cheng, M, 2021)

Under the above assumptions, for any T and any continuous bounded $\psi : \mathbb{R}^m \times C([0, T], \mathbb{R}^m) \rightarrow \mathbb{R}$ we have

$$\text{p-lim}_{n,d \rightarrow \infty} \frac{1}{d} \sum_{i=1}^d \psi(\sqrt{d}\theta_i^*, \sqrt{d}(\theta_i)_0^T) = \mathbb{E}\{\psi(\vartheta_{*,i}, (\vartheta_i)_0^T)\}.$$

- ▶ **Informally:** The DMFT process ϑ_t characterizes the distribution of rows of Θ_t .
- ▶ **Application to 2-layer nets:** $w_i(t)$: First layer weights in a 2-layer network

$$\text{p-lim}_{t,s} \langle w_i(t), w_j(s) \rangle = C_{ij}(t, s),$$

where the C_{ij} solve DMFT equations.

DMFT characterization

Theorem (Celentano, Cheng, M, 2021)

Under the above assumptions, for any T and any continuous bounded $\psi : \mathbb{R}^m \times C([0, T], \mathbb{R}^m) \rightarrow \mathbb{R}$ we have

$$\underset{n,d \rightarrow \infty}{\text{p-lim}} \frac{1}{d} \sum_{i=1}^d \psi(\sqrt{d}\theta_i^*, \sqrt{d}(\theta_i)_0^T) = \mathbb{E}\{\psi(\vartheta_{*,i}, (\vartheta_i)_0^T)\}.$$

- ▶ **Informally:** The DMFT process ϑ_t characterizes the distribution of rows of Θ_t .
- ▶ **Application to 2-layer nets:** $w_i(t)$: First layer weights in a 2-layer network

$$\underset{t,s}{\text{p-lim}} \langle w_i(t), w_j(s) \rangle = C_{ij}(t, s),$$

where the C_{ij} solve DMFT equations.

Application to 2-layer nets

$$C_{ij}^n(t, s) := \langle w_i(t), w_j(s) \rangle, \quad v_i^n(t) := \langle w_i(t), w_* \rangle, \quad a_i^n(t).$$

Theorem (Celentano, Cheng, M, 2021)

As $n, d \rightarrow \infty$, $n/d \rightarrow \bar{\alpha}$, uniformly over compacts,

$$C_{ij}^n \xrightarrow{p} C_{ij}, \quad v_i^n \xrightarrow{p} v_i, \quad a_i^n \xrightarrow{p} a_i.$$

Further $\mathbf{C} = (C_{ij} : i, j \leq m)$, $\mathbf{v} = (v_i : i \leq m)$, $\mathbf{a} = (a_i : i \leq m)$ are deterministic and the unique solution of DMFT equations.

Challenge: Solving the DMFT equations

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

$$r^t = -\frac{1}{\delta} \int_0^t R_\vartheta(t,s) \nabla L(r^s, w^*; z) ds + w^t, \quad w^t \sim \text{GP}(0, C_\vartheta),$$

$$R_\vartheta(t,s) = \mathbb{E} \left[\frac{\partial \vartheta^t}{\partial u^s} \right], \quad 0 \leq s \leq t < \infty,$$

$$R_\ell(t,s) = \mathbb{E} \left[\frac{\partial \nabla L(r^t, w^*; z)}{\partial w^s} \right], \quad 0 \leq s < t < \infty,$$

$$R_\ell(t,*) = \mathbb{E} \left[\frac{\partial \nabla L(r^t, w^*; z)}{\partial w^*} \right],$$

$$\Gamma^t = \mathbb{E} \left[\nabla^2 L(r^t, w^*; z) \right],$$

$$C_\vartheta(t,s) = \mathbb{E} \left[\vartheta^t \vartheta^s \top \right], \quad 0 \leq s \leq t < \infty \text{ or } s = *,$$

$$C_\ell(t,s) = \mathbb{E} \left[\nabla L(r^t, w^*; z) \nabla L(r^s, w^*; z) \top \right], \quad 0 \leq s \leq t < \infty.$$

Evaluating numerically the equations requires to compute expectation wrt the distribution of the processes ϑ , r .

Challenge: Solving the DMFT equations

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

$$r^t = -\frac{1}{\delta} \int_0^t R_\vartheta(t,s) \nabla L(r^s, w^*; z) ds + w^t, \quad w^t \sim \text{GP}(0, C_\vartheta),$$

$$R_\vartheta(t,s) = \mathbb{E} \left[\frac{\partial \vartheta^t}{\partial u^s} \right], \quad 0 \leq s \leq t < \infty,$$

$$R_\ell(t,s) = \mathbb{E} \left[\frac{\partial \nabla L(r^t, w^*; z)}{\partial w^s} \right], \quad 0 \leq s < t < \infty,$$

$$R_\ell(t,*) = \mathbb{E} \left[\frac{\partial \nabla L(r^t, w^*; z)}{\partial w^*} \right],$$

$$\Gamma^t = \mathbb{E} \left[\nabla^2 L(r^t, w^*; z) \right],$$

$$C_\vartheta(t,s) = \mathbb{E} \left[\vartheta^t \vartheta^s \top \right], \quad 0 \leq s \leq t < \infty \text{ or } s = *,$$

$$C_\ell(t,s) = \mathbb{E} \left[\nabla L(r^t, w^*; z) \nabla L(r^s, w^*; z) \top \right], \quad 0 \leq s \leq t < \infty.$$

Evaluating numerically the equations requires to compute expectation wrt the distribution of the processes ϑ , r .

Gaussian approximation

Take notice:

- ▶ Not Gaussian process regression.
- ▶ Not approximating ϑ_t , r_t by Gaussian processes.

Idea: Matching Gaussian model

$$\widehat{\mathcal{R}}_n(\theta) := \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i; \theta))^2 = \frac{1}{2n} \|F(\theta)\|^2.$$

- ▶ Replace $F(\theta)$ by Gaussian proc. $F^g(\theta)$ with matching mean and covariance.
- ▶ Study DMFT for this model.
- ▶ A piece of the covariance:

$$\mathbb{E}\{f(x; \theta_1)f(x; \theta_2)\} = \frac{1}{m^2} \sum_{i,j=1}^m h(w_i^\top w_j) a_i a_j.$$

Idea: Matching Gaussian model

$$\widehat{\mathcal{R}}_n(\theta) := \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i; \theta))^2 = \frac{1}{2n} \|F(\theta)\|^2.$$

- ▶ Replace $F(\theta)$ by Gaussian proc. $F^g(\theta)$ with matching mean and covariance.
- ▶ Study DMFT for this model.
- ▶ A piece of the covariance:

$$\mathbb{E}\{f(x; \theta_1)f(x; \theta_2)\} = \frac{1}{m^2} \sum_{i,j=1}^m h(w_i^\top w_j) a_i a_j.$$

DMFT eqs for Gaussian model are explicit

$$\begin{aligned} \frac{da}{dt}(t) &= \frac{\bar{\alpha}}{m} \hat{\varphi}(\mathbf{v}(t)) \int_0^t R_A(t, s) ds \\ &\quad - \frac{\bar{\alpha}}{m} \int_0^t R_A(t, s) a(s) \left[\frac{1}{m} h(C_d(t, s)) + \frac{m-1}{m} h(C_o(t, s)) \right] ds \\ &\quad - \frac{\bar{\alpha}}{m} \int_0^t C_A(t, s) a(s) \left[\frac{1}{m} h'(C_d(t, s)) R_d(t, s) + \frac{m-1}{m} h'(C_o(t, s)) R_o(t, s) \right] ds, \end{aligned} \quad (2.32)$$

$$\begin{aligned} \frac{d\mathbf{v}}{dt}(t) &= -\nu(t)\mathbf{v}(t) + \frac{\bar{\alpha}}{m} \nabla \hat{\varphi}(\mathbf{v}(t)) a(t) \int_0^t R_A(t, s) ds \\ &\quad - \frac{1}{m} \int_0^t [M_R^{(d)}(t, s) + (m-1)M_R^{(o)}(t, s)] \mathbf{v}(s) ds, \end{aligned} \quad (2.33)$$

$$\begin{aligned} \partial_t C_d(t, t') &= -\nu(t)C_d(t, t') + \frac{\bar{\alpha}}{m} \langle \nabla \hat{\varphi}'(\mathbf{v}(t)), \mathbf{v}(t') \rangle a(t) \int_0^t R_A(t, s) ds \\ &\quad - \frac{1}{m} \int_0^t [M_R^{(d)}(t, s) C_d(t', s) + (m-1)M_R^{(o)}(t, s) C_o(t', s)] ds \\ &\quad - \frac{1}{m} \int_0^{t'} [M_C^{(d)}(t, s) R_d(t', s) + (m-1)M_C^{(o)}(t, s) R_o(t', s)] ds, \end{aligned} \quad (2.34)$$

$$\begin{aligned} \partial_t C_o(t, t') &= -\nu(t)C_o(t, t') + \frac{\bar{\alpha}}{m} \langle \nabla \hat{\varphi}(\mathbf{v}(t)), \mathbf{v}(t') \rangle a(t) \int_0^t R_A(t, s) ds \\ &\quad - \frac{1}{m} \int_0^t [M_R^{(d)}(t, s) C_o(t', s) + M_R^{(o)}(t, s) C_d(t', s) + (m-2)M_R^{(o)}(t, s) C_o(t', s)] ds \\ &\quad - \frac{1}{m} \int_0^{t'} [M_C^{(d)}(t, s) R_o(t', s) + M_C^{(o)}(t, s) R_d(t', s) + (m-2)M_C^{(o)}(t, s) R_o(t', s)] ds, \end{aligned} \quad (2.35)$$

$$\begin{aligned} \partial_t R_d(t, t') &= -\nu(t)R_d(t, t') + \delta(t - t') \\ &\quad - \frac{1}{m} \int_{t'}^t [M_R^{(d)}(t, s) R_d(s, t') + (m-1)M_R^{(o)}(t, s) R_o(s, t')] ds, \end{aligned} \quad (2.36)$$

$$\begin{aligned} \partial_t R_o(t, t') &= -\nu(t)R_o(t, t') - \frac{1}{m} \int_{t'}^t [M_R^{(d)}(t, s) R_o(s, t') + M_R^{(o)}(t, s) R_d(s, t')] \\ &\quad + (m-2)M_R^{(o)}(t, s) R_o(s, t')] ds. \end{aligned} \quad (2.37)$$

(2) Equations for auxiliary functions. The memory kernels $M_R^{(x)}(t, s)$, $M_R^{(o)}(t, s)$ and $M_C^{(x)}(t, s)$, $M_C^{(o)}(t, s)$ are given by:

$$M_R^{(d)}(t, s) = \frac{\bar{\alpha}}{m} a(t)a(s) [R_A(t, s) h'(C_d(t, s)) + C_A(t, s) h''(C_d(t, s)) R_d(t, s)], \quad (2.38)$$

$$M_R^{(o)}(t, s) = \frac{\bar{\alpha}}{m} a(t)a(s) [R_A(t, s) h'(C_o(t, s)) + C_A(t, s) h''(C_o(t, s)) R_o(t, s)], \quad (2.39)$$

$$M_C^{(d)}(t, s) = \frac{\bar{\alpha}}{m} a(t)a(s) C_A(t, s) h'(C_d(t, s)), \quad (2.40)$$

$$M_C^{(o)}(t, s) = \frac{\bar{\alpha}}{m} a(t)a(s) C_A(t, s) h'(C_o(t, s)). \quad (2.41)$$

Further, $C_A(t, s)$, $R_A(t, s)$ are given by the same equations (2.19), where Σ_C , Σ_R are simplified as follows:

$$\begin{aligned} \Sigma_C(t, s) &= \tau^2 + \|\varphi\|^2 - a(t)\hat{\varphi}(\mathbf{v}(t)) - a(s)\hat{\varphi}(\mathbf{v}(s)) + \frac{a(t)a(s)}{m} h(C_d(t, s)) \\ &\quad + \frac{m-1}{m} a(t)a(s) h(C_o(t, s)) \\ \Sigma_R(t, s) &= \frac{a(t)a(s)}{m} h'(C_d(t, s)) R_d(t, s) + \frac{m-1}{m} a(t)a(s) h'(C_o(t, s)) R_o(t, s) \end{aligned} \quad (2.42)$$

Highly nonlinear!

DMFT eqs for Gaussian model are explicit

$$\begin{aligned} \frac{da}{dt}(t) &= \frac{\bar{\alpha}}{m} \hat{\varphi}(\mathbf{v}(t)) \int_0^t R_A(t, s) ds \\ &\quad - \frac{\bar{\alpha}}{m} \int_0^t R_A(t, s) a(s) \left[\frac{1}{m} h(C_d(t, s)) + \frac{m-1}{m} h(C_o(t, s)) \right] ds \\ &\quad - \frac{\bar{\alpha}}{m} \int_0^t C_A(t, s) a(s) \left[\frac{1}{m} h'(C_d(t, s)) R_d(t, s) + \frac{m-1}{m} h'(C_o(t, s)) R_o(t, s) \right] ds, \end{aligned} \quad (2.32)$$

$$\begin{aligned} \frac{d\mathbf{v}}{dt}(t) &= -\nu(t)\mathbf{v}(t) + \frac{\bar{\alpha}}{m} \nabla \hat{\varphi}(\mathbf{v}(t)) a(t) \int_0^t R_A(t, s) ds \\ &\quad - \frac{1}{m} \int_0^t [M_R^{(d)}(t, s) + (m-1)M_R^{(o)}(t, s)] \mathbf{v}(s) ds, \end{aligned} \quad (2.33)$$

$$\begin{aligned} \partial_t C_d(t, t') &= -\nu(t)C_d(t, t') + \frac{\bar{\alpha}}{m} \langle \nabla \hat{\varphi}'(\mathbf{v}(t)), \mathbf{v}(t') \rangle a(t) \int_0^t R_A(t, s) ds \\ &\quad - \frac{1}{m} \int_0^t [M_R^{(d)}(t, s) C_d(t', s) + (m-1)M_R^{(o)}(t, s) C_o(t', s)] ds \\ &\quad - \frac{1}{m} \int_0^{t'} [M_C^{(d)}(t, s) R_d(t', s) + (m-1)M_C^{(o)}(t, s) R_o(t', s)] ds, \end{aligned} \quad (2.34)$$

$$\begin{aligned} \partial_t C_o(t, t') &= -\nu(t)C_o(t, t') + \frac{\bar{\alpha}}{m} \langle \nabla \hat{\varphi}(\mathbf{v}(t)), \mathbf{v}(t') \rangle a(t) \int_0^t R_A(t, s) ds \\ &\quad - \frac{1}{m} \int_0^t [M_R^{(d)}(t, s) C_o(t', s) + M_R^{(o)}(t, s) C_d(t', s) + (m-2)M_R^{(o)}(t, s) C_o(t', s)] ds \\ &\quad - \frac{1}{m} \int_0^{t'} [M_C^{(d)}(t, s) R_o(t', s) + M_C^{(o)}(t, s) R_d(t', s) + (m-2)M_C^{(o)}(t, s) R_o(t', s)] ds, \end{aligned} \quad (2.35)$$

$$\begin{aligned} \partial_t R_d(t, t') &= -\nu(t)R_d(t, t') + \delta(t - t') \\ &\quad - \frac{1}{m} \int_{t'}^t [M_R^{(d)}(t, s) R_d(s, t') + (m-1)M_R^{(o)}(t, s) R_o(s, t')] ds, \end{aligned} \quad (2.36)$$

$$\begin{aligned} \partial_t R_o(t, t') &= -\nu(t)R_o(t, t') - \frac{1}{m} \int_{t'}^t [M_R^{(d)}(t, s) R_o(s, t') + M_R^{(o)}(t, s) R_d(s, t')] \\ &\quad + (m-2)M_R^{(o)}(t, s) R_o(s, t')] ds. \end{aligned} \quad (2.37)$$

(2) Equations for auxiliary functions. The memory kernels $M_R^{(x)}(t, s)$, $M_R^{(o)}(t, s)$ and $M_C^{(x)}(t, s)$, $M_C^{(o)}(t, s)$ are given by:

$$M_R^{(d)}(t, s) = \frac{\bar{\alpha}}{m} a(t)a(s) [R_A(t, s) h'(C_d(t, s)) + C_A(t, s) h''(C_d(t, s)) R_d(t, s)], \quad (2.38)$$

$$M_R^{(o)}(t, s) = \frac{\bar{\alpha}}{m} a(t)a(s) [R_A(t, s) h'(C_o(t, s)) + C_A(t, s) h''(C_o(t, s)) R_o(t, s)], \quad (2.39)$$

$$M_C^{(d)}(t, s) = \frac{\bar{\alpha}}{m} a(t)a(s) C_A(t, s) h'(C_d(t, s)), \quad (2.40)$$

$$M_C^{(o)}(t, s) = \frac{\bar{\alpha}}{m} a(t)a(s) C_A(t, s) h'(C_o(t, s)). \quad (2.41)$$

Further, $C_A(t, s)$, $R_A(t, s)$ are given by the same equations (2.19), where Σ_C , Σ_R are simplified as follows:

$$\begin{aligned} \Sigma_C(t, s) &= \tau^2 + \|\varphi\|^2 - a(t)\hat{\varphi}(\mathbf{v}(t)) - a(s)\hat{\varphi}(\mathbf{v}(s)) + \frac{a(t)a(s)}{m} h(C_d(t, s)) \\ &\quad + \frac{m-1}{m} a(t)a(s) h(C_o(t, s)) \\ \Sigma_R(t, s) &= \frac{a(t)a(s)}{m} h'(C_d(t, s)) R_d(t, s) + \frac{m-1}{m} a(t)a(s) h'(C_o(t, s)) R_o(t, s) \end{aligned} \quad (2.42)$$

Highly nonlinear!

A related (simpler) model

Random Gaussian Equations

- ▶ Physics: Fyodorov 2019; Urbani 2023; Kamali, Urnbani, 2023
- ▶ Mathematics: M, Subag 2023, 2024

DMFT eqs for Gaussian model are explicit

$$\begin{aligned}\frac{da}{dt}(t) &= \frac{\bar{\alpha}}{m} \dot{\varphi}(\mathbf{v}(t)) \int_0^t R_A(t, s) ds \\ &\quad - \frac{\bar{\alpha}}{m} \int_0^t R_A(t, s) a(s) \left[\frac{1}{m} h(C_d(t, s)) + \frac{m-1}{m} h(C_o(t, s)) \right] ds \\ &\quad - \frac{\bar{\alpha}}{m} \int_0^t C_A(t, s) a(s) \left[\frac{1}{m} h'(C_d(t, s)) R_d(t, s) + \frac{m-1}{m} h'(C_o(t, s)) R_o(t, s) \right] ds,\end{aligned}\tag{2.32}$$

$$\begin{aligned}\frac{dv}{dt}(t) &= -\nu(t) \mathbf{v}(t) + \frac{\bar{\alpha}}{m} \nabla \dot{\varphi}(\mathbf{v}(t)) a(t) \int_0^t R_A(t, s) ds \\ &\quad - \frac{1}{m} \int_0^t \left[M_R^{(d)}(t, s) + (m-1) M_R^{(o)}(t, s) \right] \mathbf{v}(s) ds,\end{aligned}\tag{2.33}$$

$$\begin{aligned}\partial_t C_d(t, t') &= -\nu(t) C_d(t, t') + \frac{\bar{\alpha}}{m} \langle \nabla \dot{\varphi}'(\mathbf{v}(t)), \mathbf{v}(t') \rangle a(t) \int_0^t R_A(t, s) ds \\ &\quad - \frac{1}{m} \int_0^t \left[M_R^{(d)}(t, s) C_d(t', s) + (m-1) M_R^{(o)}(t, s) C_o(t', s) \right] ds \\ &\quad - \frac{1}{m} \int_0^{t'} \left[M_C^{(d)}(t, s) R_d(t', s) + (m-1) M_C^{(o)}(t, s) R_o(t', s) \right] ds,\end{aligned}\tag{2.34}$$

$$\begin{aligned}\partial_t C_o(t, t') &= -\nu(t) C_o(t, t') + \frac{\bar{\alpha}}{m} \langle \nabla \dot{\varphi}(\mathbf{v}(t)), \mathbf{v}(t') \rangle a(t) \int_0^t R_A(t, s) ds \\ &\quad - \frac{1}{m} \int_0^t \left[M_R^{(d)}(t, s) C_o(t', s) + M_R^{(o)}(t, s) C_d(t', s) + (m-2) M_R^{(o)}(t, s) C_o(t', s) \right] ds \\ &\quad - \frac{1}{m} \int_0^{t'} \left[M_C^{(d)}(t, s) R_o(t', s) + M_C^{(o)}(t, s) R_d(t', s) + (m-2) M_C^{(o)}(t, s) R_o(t', s) \right] ds,\end{aligned}\tag{2.35}$$

$$\begin{aligned}\partial_t R_d(t, t') &= -\nu(t) R_d(t, t') + \delta(t - t') \\ &\quad - \frac{1}{m} \int_{t'}^t \left[M_R^{(d)}(t, s) R_d(s, t') + (m-1) M_R^{(o)}(t, s) R_o(s, t') \right] ds,\end{aligned}\tag{2.36}$$

$$\begin{aligned}\partial_t R_o(t, t') &= -\nu(t) R_o(t, t') - \frac{1}{m} \int_{t'}^t \left[M_R^{(d)}(t, s) R_o(s, t') + M_R^{(o)}(t, s) R_d(s, t') \right. \\ &\quad \left. + (m-2) M_R^{(o)}(t, s) R_o(s, t') \right] ds.\end{aligned}\tag{2.37}$$

(2) Equations for auxiliary functions. The memory kernels $M_R^{(s)}(t, s)$, $M_R^{(o)}(t, s)$ and $M_C^{(s)}(t, s)$, $M_C^{(o)}(t, s)$ are given by:

$$M_R^{(d)}(t, s) = \frac{\bar{\alpha}}{m} a(t) a(s) [R_A(t, s) h'(C_d(t, s)) + C_A(t, s) h''(C_d(t, s)) R_d(t, s)],\tag{2.38}$$

$$M_R^{(o)}(t, s) = \frac{\bar{\alpha}}{m} a(t) a(s) [R_A(t, s) h'(C_o(t, s)) + C_A(t, s) h''(C_o(t, s)) R_o(t, s)],\tag{2.39}$$

$$M_C^{(d)}(t, s) = \frac{\bar{\alpha}}{m} a(t) a(s) C_A(t, s) h'(C_d(t, s)),\tag{2.40}$$

$$M_C^{(o)}(t, s) = \frac{\bar{\alpha}}{m} a(t) a(s) C_A(t, s) h'(C_o(t, s)).\tag{2.41}$$

Further, $C_A(t, s)$, $R_A(t, s)$ are given by the same equations (2.19), where Σ_C , Σ_R are simplified as follows:

$$\begin{aligned}\Sigma_C(t, s) &= \tau^2 + \|\varphi\|^2 - a(t) \dot{\varphi}(\mathbf{v}(t)) - a(s) \dot{\varphi}(\mathbf{v}(s)) + \frac{a(t)a(s)}{m} h(C_d(t, s)) \\ &\quad + \frac{m-1}{m} a(t)a(s) h(C_o(t, s))\end{aligned}\tag{2.42}$$

$$\Sigma_R(t, s) = \frac{a(t)a(s)}{m} h'(C_d(t, s)) R_d(t, s) + \frac{m-1}{m} a(t)a(s) h'(C_o(t, s)) R_o(t, s)$$

- ▶ Symmetric initialization: $\mathbf{a}_i(0) = \mathbf{a}_0$, $\mathbf{w}_i(0) \sim \text{Unif}(\mathbb{S}^{d-1})$
- ▶ Numerical solution
- ▶ Singular asymptotics: $m \rightarrow \infty$, $t \rightarrow \infty$.

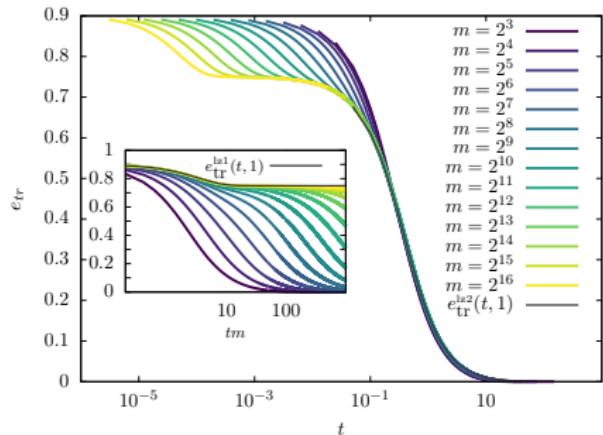
Dynamics in pure noise

Pure noise

$$(x_i, y_i) \sim N(0, I_d) \otimes N(0, \tau^2)$$

No learning, just optimization!

Fixed 2nd layer weights

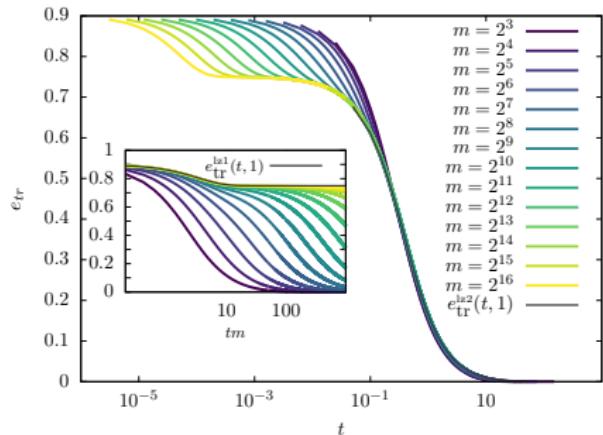


$$f(x; \theta) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\langle w_j, x \rangle), \quad a_i(t) = \gamma \sqrt{m}, \quad \text{fixed.}$$

$$\lim_{n,d \rightarrow \infty, n/d = \bar{\alpha}} \widehat{\mathcal{R}}_n(\theta(t)) = e_{tr}(t; m, \bar{\alpha}),$$

$$\lim_{m,\alpha \rightarrow \infty, \bar{\alpha}/m = \alpha} e_{tr}(t; m, \bar{\alpha}) = e_{tr}^{(2)}(t; \alpha).$$

Fixed 2nd layer weights

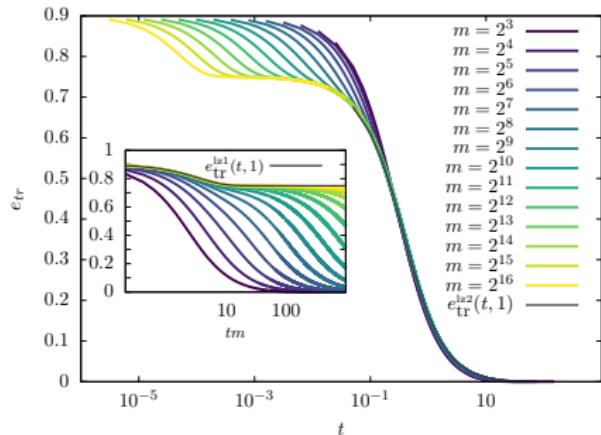


$$f(x; \theta) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\langle w_j, x \rangle), \quad a_i(t) = \gamma \sqrt{m}, \quad \text{fixed.}$$

$$\lim_{n,d \rightarrow \infty, n/d = \bar{\alpha}} \hat{\mathcal{R}}_n(\theta(t)) = e_{tr}(t; m, \bar{\alpha}),$$

$$\lim_{m,\alpha \rightarrow \infty, \bar{\alpha}/m = \alpha} e_{tr}(t; m, \bar{\alpha}) = e_{tr}^{(2)}(t; \alpha).$$

Fixed 2nd layer weights

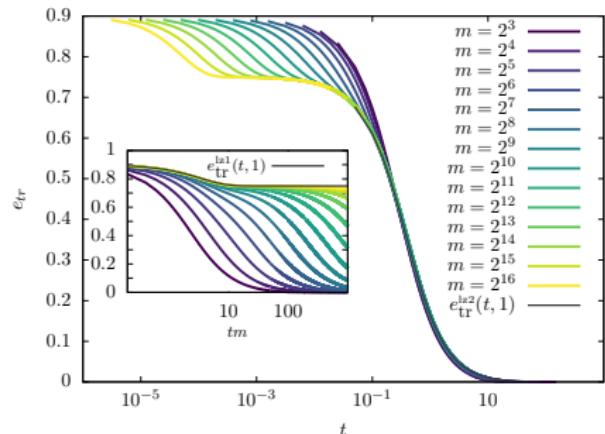


$$f(x; \theta) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\langle w_j, x \rangle), \quad a_i(t) = \gamma \sqrt{m}, \quad \text{fixed.}$$

$$\lim_{n,d \rightarrow \infty, n/d = \bar{\alpha}} \hat{\mathcal{R}}_n(\theta(t)) = e_{tr}(t; m, \bar{\alpha}),$$

$$\lim_{m,\alpha \rightarrow \infty, \bar{\alpha}/m = \alpha} e_{tr}(t; m, \bar{\alpha}) = e_{tr}^{(2)}(t; \alpha).$$

Fixed 2nd layer weights

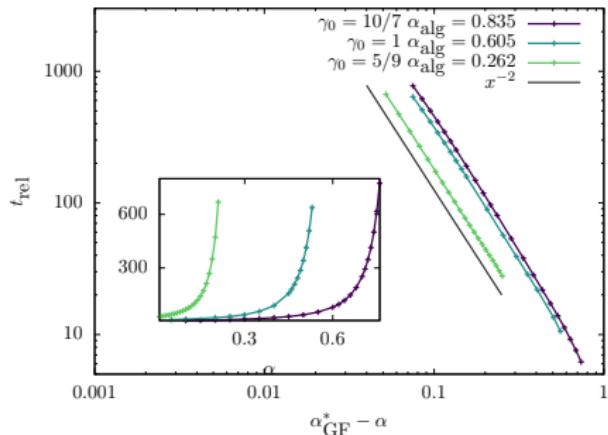


Algorithmic interpolation phase transition

$$\gamma < \gamma_{\text{GF}}(\alpha) \Rightarrow \lim_{t \rightarrow \infty} e_{\text{tr}}^{(2)}(t; \alpha) > 0,$$

$$\gamma > \gamma_{\text{GF}}(\alpha) \Rightarrow \lim_{t \rightarrow \infty} e_{\text{tr}}^{(2)}(t; \alpha) = 0.$$

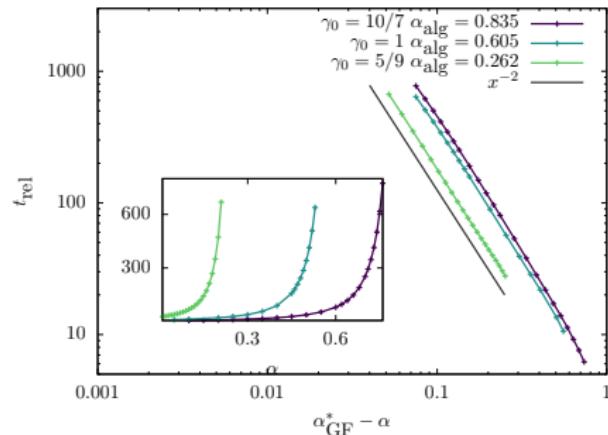
Fixed 2nd layer weights: Phase transition



$$\begin{aligned} t_{\text{rel}}(\alpha, \gamma) &\asymp (\gamma_{\text{GF}}(\alpha) - \gamma)^{-\nu} \\ &\asymp (\alpha_{\text{GF}}(\gamma) - \alpha)^{-\nu}. \end{aligned}$$

What happens if second-layer weights evolve?

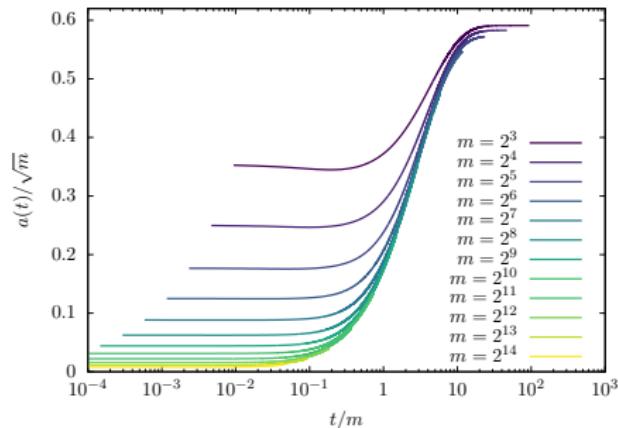
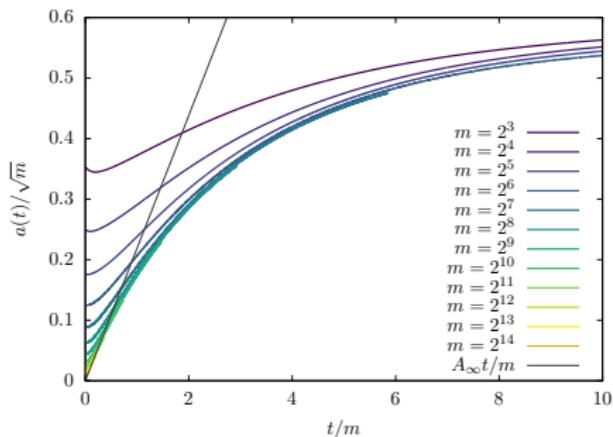
Fixed 2nd layer weights: Phase transition



$$\begin{aligned}t_{\text{rel}}(\alpha, \gamma) &\asymp (\gamma_{\text{GF}}(\alpha) - \gamma)^{-\nu} \\&\asymp (\alpha_{\text{GF}}(\gamma) - \alpha)^{-\nu}.\end{aligned}$$

What happens if second-layer weights evolve?

Evolving 2nd layer weights: $a_i(0) = a_0$

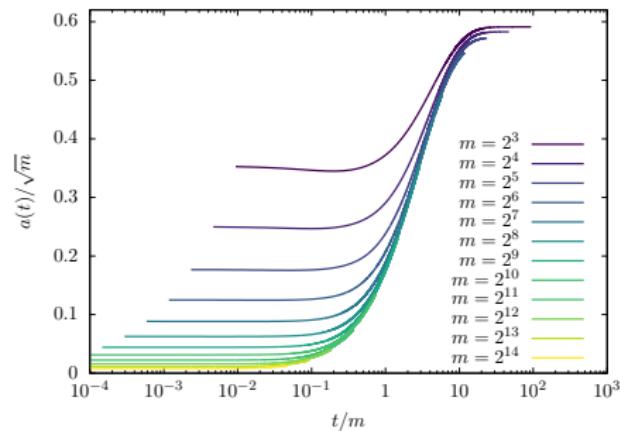
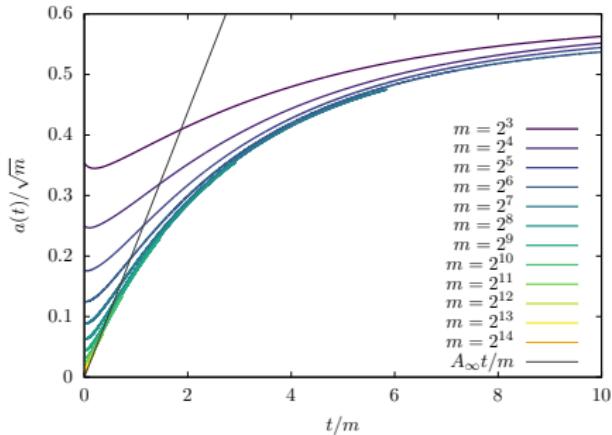


- ▶ Slowest of 3 dynamical regimes: $t/m = s = O(1)$.
- ▶ $a(t) = \sqrt{m}\gamma(t/m) + o(\sqrt{m})$
- ▶ Asymptotic behaviors:

$$s \rightarrow 0 : \quad \gamma(s) = \gamma_* s + o(s),$$

$$s \rightarrow \infty : \quad \gamma(s) \rightarrow \gamma_{GF}.$$

Evolving 2nd layer weights: $a_i(0) = a_0$

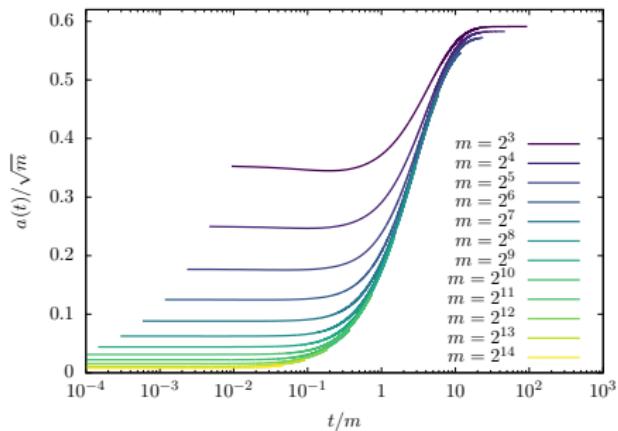
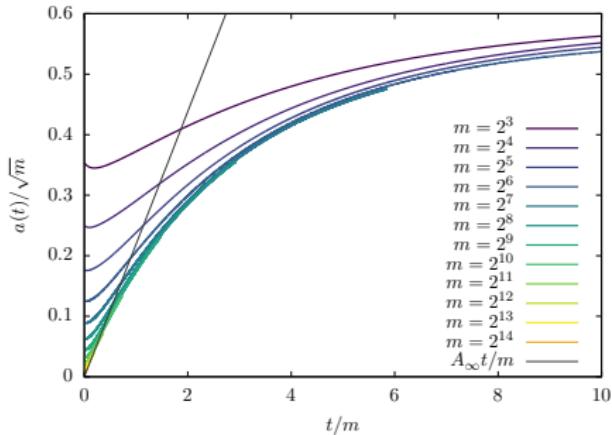


- ▶ Slowest of 3 dynamical regimes: $t/m = s = O(1)$.
- ▶ $a(t) = \sqrt{m}\gamma(t/m) + o(\sqrt{m})$
- ▶ Asymptotic behaviors:

$$s \rightarrow 0 : \quad \gamma(s) = \gamma_* s + o(s).$$

γ_* is the threshold energy of a p-spin model!

Evolving 2nd layer weights: $a_i(0) = a_0$

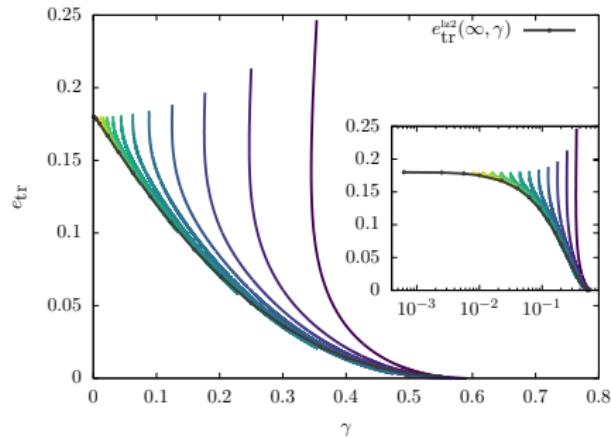


- ▶ Slowest of 3 dynamical regimes: $t/m = s = O(1)$.
- ▶ $a(t) = \sqrt{m}\gamma(t/m) + o(\sqrt{m})$
- ▶ Asymptotic behaviors:

$$s \rightarrow \infty : \quad \gamma(s) \rightarrow \gamma_{\text{GF}} .$$

Complexity $\gamma(t)$ grows slowly until interpolation threshold

Adiabatic evolution of $\gamma(t) = a(t)/\sqrt{m}$



- ▶ **Black:** Minimum empirical risk achieved at γ fixed.
- ▶ **Colored:** $a(t)$ evolving

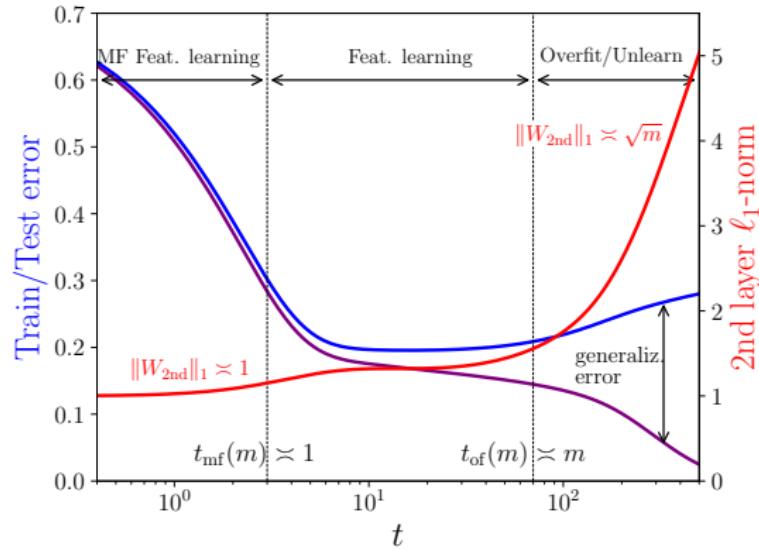
Dynamics under single-index model

Single-index model

Data $\{(x_i, y_i) : i \leq n\}$ iid, $\varepsilon_i \sim N(0, \tau^2)$

$$x_i \sim N(0, I_d), \quad y_i = \varphi(\langle w_*, x_i \rangle) + \varepsilon_i$$

Dynamical decoupling: Learning → Overfitting



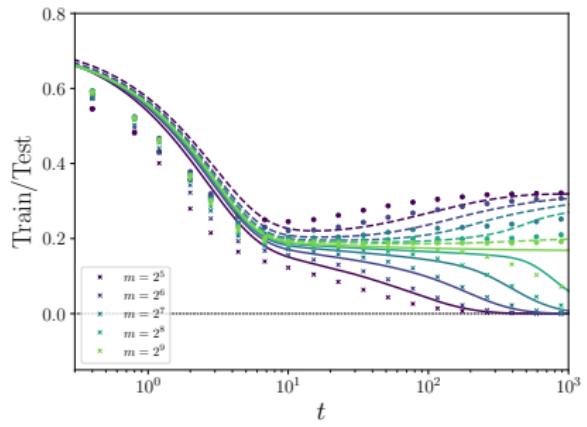
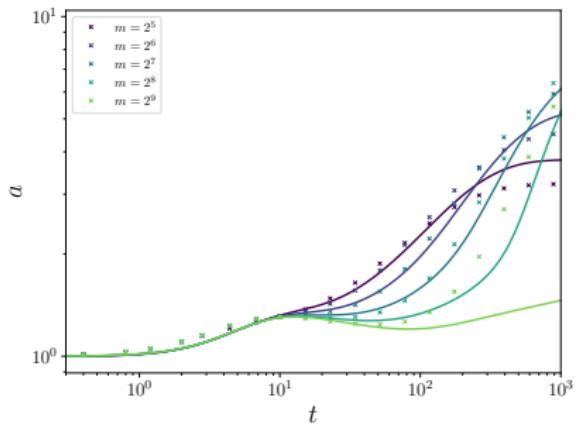
$t \asymp 1$:

- ▶ Test error \approx Train error
- ▶ Learns latent direction w_*
- ▶ Mean field asymptotically correct (Mei, M, Nguyen, 2018; Chizat, Bach, 2018; Rotskoff, Vanden Eijnden, 2018)

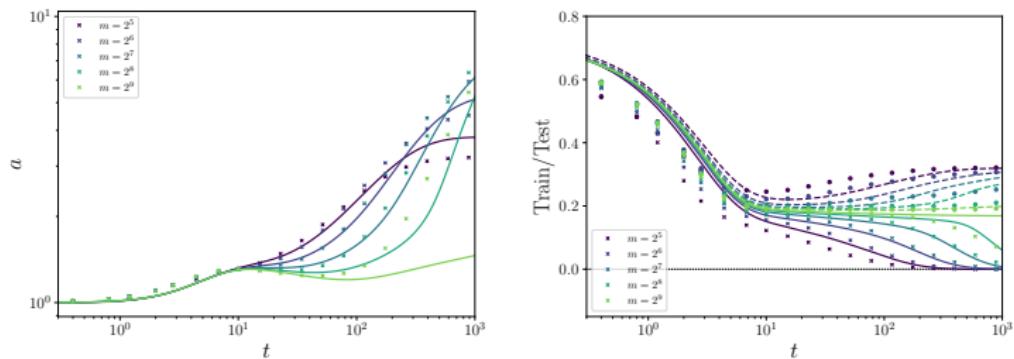
$t \asymp m$:

- ▶ Test error \gg Train error; Train error $\rightarrow 0$
- ▶ *Unlearns* latent direction w_*
- ▶ Neural tangent kernel learning

Comparing with SGD experiments



Overfitting mechanism: Adiabatic increase in complexity



► Bartlett 1996:

$$|\hat{\mathcal{R}}_n(\theta) - \mathcal{R}(\theta)| \lesssim \underbrace{\frac{1}{m} \|\mathbf{a}\|_1}_{\text{Rademacher complexity}} \cdot \frac{1}{\sqrt{\alpha m}},$$

Rademacher complexity

►

$$t = \Theta(1) \Rightarrow \frac{1}{m} \|\mathbf{a}(t)\|_1 = \mathbf{a}(t) = O(1) \Rightarrow \text{No overfitting}$$

Lower bounding the overfitting timescale

Theorem (M, Urbani, 2025)

Assume $\|\sigma\|_{\text{Lip}}, \|\sigma\|_\infty \leq L$, $|\varphi(0)|, \|\varphi\|_{\text{Lip}} \leq L$, $\|\alpha(0)\|_\infty \leq a_0$, $n \geq d \vee m$. Then, with prob $\geq 1 - 2 \exp(-cn)$ ($\hat{t} = t\alpha$)

$$\|\alpha(\hat{t})\|_\infty \leq a_0 + a_1 \hat{t}, \quad a_1 := C_0 L (\tau + a_0 L),$$

$$\mathcal{R}(\alpha(\hat{t}), \mathbf{W}(\hat{t})) - \widehat{\mathcal{R}}_n(\alpha(\hat{t}), \mathbf{W}(\hat{t})) \leq C_1 L^3 (a_0 + a_1 \hat{t})^2 \cdot \sqrt{\frac{d}{n}}.$$

Dynamics on the feature learning timescale

Mean field asymptotics: $(\mathbf{Q}_v := \mathbf{I}_k - vv^\top, \hat{\varphi}, h = \dots)$

$$\frac{d}{dt} v_i^{mf}(\hat{t}) = a_i^{mf}(\hat{t}) \mathbf{Q}_{v_i^{mf}} \left(\nabla \hat{\varphi}(v_i^{mf}(\hat{t})) - \frac{1}{m} \sum_{j=1}^m a_j^{mf}(\hat{t}) h'(\langle v_i^{mf}(\hat{t}), v_j^{mf}(\hat{t}) \rangle) v_j^{mf}(\hat{t}) \right),$$

$$\frac{d}{dt} a_i^{mf}(\hat{t}) = \hat{\varphi}(v_i^{mf}(\hat{t})) - \frac{1}{m} \sum_{j=1}^m a_j^{mf}(\hat{t}) h(\langle v_i^{mf}(\hat{t}), v_j^{mf}(\hat{t}) \rangle).$$

- ▶ Gradient flow wrt population risk.
- ▶ High-dimensional asymptotics: $\langle w_i(t), w_j(t) \rangle \approx \langle v_i(t), v_j(t) \rangle$.
- ▶ 2-dimensional under symmetric initialization

Dynamics on the feature learning timescale

Mean field asymptotics:

$$(\mathbf{Q}_v := \mathbf{I}_k - vv^\top, \hat{\varphi}, h = \dots)$$

$$\frac{d}{dt} v_i^{mf}(\hat{t}) = a_i^{mf}(\hat{t}) \mathbf{Q}_{v_i^{mf}} \left(\nabla \hat{\varphi}(v_i^{mf}(\hat{t})) - \frac{1}{m} \sum_{j=1}^m a_j^{mf}(\hat{t}) h'(\langle v_i^{mf}(\hat{t}), v_j^{mf}(\hat{t}) \rangle) v_j^{mf}(\hat{t}) \right),$$

$$\frac{d}{dt} a_i^{mf}(\hat{t}) = \hat{\varphi}(v_i^{mf}(\hat{t})) - \frac{1}{m} \sum_{j=1}^m a_j^{mf}(\hat{t}) h(\langle v_i^{mf}(\hat{t}), v_j^{mf}(\hat{t}) \rangle).$$

- ▶ Gradient flow wrt population risk.
- ▶ High-dimensional asymptotics: $\langle w_i(t), w_j(t) \rangle \approx \langle v_i(t), v_j(t) \rangle$.
- ▶ 2-dimensional under symmetric initialization

Dynamics on the feature learning timescale

Mean field asymptotics:

$$(\mathbf{Q}_v := \mathbf{I}_k - vv^\top, \hat{\varphi}, h = \dots)$$

$$\frac{d}{dt} v_i^{mf}(\hat{t}) = a_i^{mf}(\hat{t}) \mathbf{Q}_{v_i^{mf}} \left(\nabla \hat{\varphi}(v_i^{mf}(\hat{t})) - \frac{1}{m} \sum_{j=1}^m a_j^{mf}(\hat{t}) h'(\langle v_i^{mf}(\hat{t}), v_j^{mf}(\hat{t}) \rangle) v_j^{mf}(\hat{t}) \right),$$

$$\frac{d}{dt} a_i^{mf}(\hat{t}) = \hat{\varphi}(v_i^{mf}(\hat{t})) - \frac{1}{m} \sum_{j=1}^m a_j^{mf}(\hat{t}) h(\langle v_i^{mf}(\hat{t}), v_j^{mf}(\hat{t}) \rangle).$$

- ▶ Gradient flow wrt population risk.
- ▶ High-dimensional asymptotics: $\langle w_i(t), w_j(t) \rangle \approx \langle v_i(t), v_j(t) \rangle$.
- ▶ 2-dimensional under symmetric initialization

Dynamics on the feature learning timescale

$$\frac{d}{dt} \mathbf{v}_i^{mf}(\hat{t}) = a_i^{mf}(\hat{t}) \mathbf{Q}_{\mathbf{v}_i^{mf}} \left(\nabla \hat{\phi}(\mathbf{v}_i^{mf}(\hat{t})) - \frac{1}{m} \sum_{j=1}^m a_j^{mf}(\hat{t}) h'(\langle \mathbf{v}_i^{mf}(\hat{t}), \mathbf{v}_j^{mf}(\hat{t}) \rangle) \mathbf{v}_j^{mf}(\hat{t}) \right),$$

$$\frac{d}{dt} a_i^{mf}(\hat{t}) = \hat{\phi}(\mathbf{v}_i^{mf}(\hat{t})) - \frac{1}{m} \sum_{j=1}^m a_j^{mf}(\hat{t}) h(\langle \mathbf{v}_i^{mf}(\hat{t}), \mathbf{v}_j^{mf}(\hat{t}) \rangle).$$

Theorem (M, Urbani, 2025)

Under same assumptions, let $T_{lb} = c_0 (\log m)^{1/3} \wedge (\log n/d)^{1/3}$. Then with prob $\geq 1 - 2 \exp(-c_1 d)$,

$$\sup_{\hat{t} \leq T_{lb}} \frac{1}{m} \sum_{i=1}^m \left(|a_i(\hat{t}) - a_i^{mf}(\hat{t})| + \|\mathbf{v}_i(\hat{t}) - \mathbf{v}_i^{mf}(\hat{t})\| \right) \leq C \left(\frac{1}{m} \vee \frac{1}{d} \vee \frac{d}{n} \right)^{1/2-\delta}$$

Dynamics on the feature learning timescale

$$\frac{d}{dt} \mathbf{v}_i^{mf}(\hat{t}) = a_i^{mf}(\hat{t}) \mathbf{Q}_{\mathbf{v}_i^{mf}} \left(\nabla \hat{\phi}(\mathbf{v}_i^{mf}(\hat{t})) - \frac{1}{m} \sum_{j=1}^m a_j^{mf}(\hat{t}) h'(\langle \mathbf{v}_i^{mf}(\hat{t}), \mathbf{v}_j^{mf}(\hat{t}) \rangle) \mathbf{v}_j^{mf}(\hat{t}) \right),$$

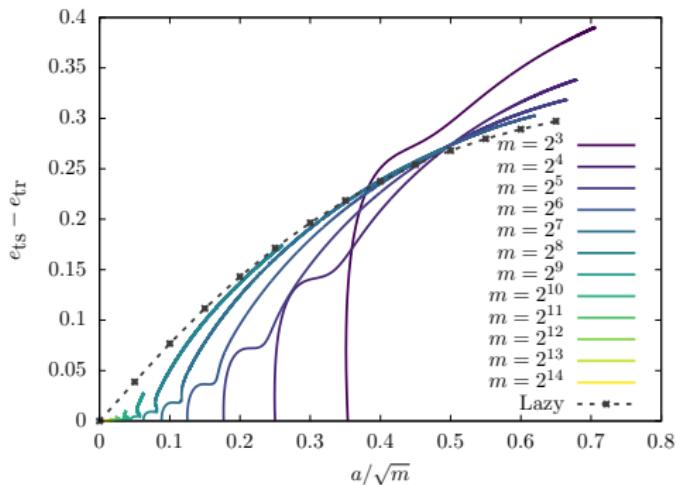
$$\frac{d}{dt} a_i^{mf}(\hat{t}) = \hat{\phi}(\mathbf{v}_i^{mf}(\hat{t})) - \frac{1}{m} \sum_{j=1}^m a_j^{mf}(\hat{t}) h(\langle \mathbf{v}_i^{mf}(\hat{t}), \mathbf{v}_j^{mf}(\hat{t}) \rangle).$$

Theorem (M, Urbani, 2025)

Under same assumptions, let $T_{lb} = c_0 (\log m)^{1/3} \wedge (\log n/d)^{1/3}$. Then with prob $\geq 1 - 2 \exp(-c_1 d)$,

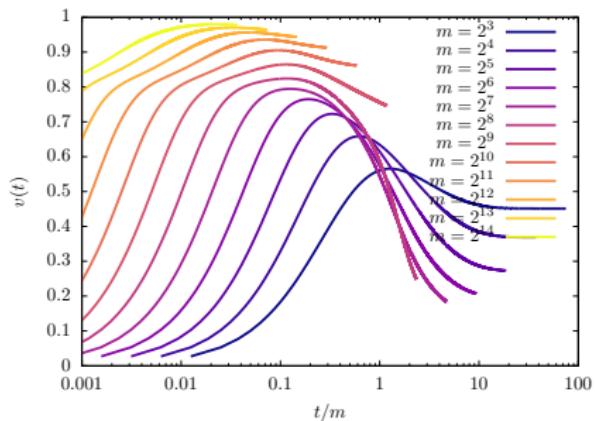
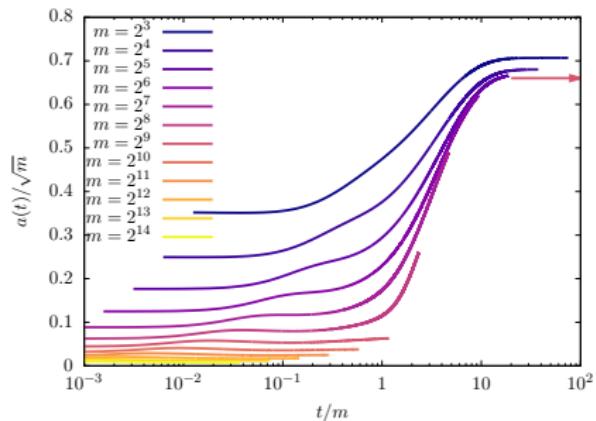
$$\sup_{\hat{t} \leq T_{lb}} \frac{1}{m} \sum_{i=1}^m \left(|a_i(\hat{t}) - a_i^{mf}(\hat{t})| + \|\mathbf{v}_i(\hat{t}) - \mathbf{v}_i^{mf}(\hat{t})\| \right) \leq C \left(\frac{1}{m} \vee \frac{1}{d} \vee \frac{d}{n} \right)^{1/2-\delta}$$

Mechanism: Adiabatic increase in complexity



- ▶ Slow timescale: Divergence of 2nd layer weights
- ▶ ⇔ Growth of Gaussian/Rademacher complexity

Unlearning: $t = \Theta(m)$



- ▶ Projection onto the latent direction vanishes.

Conclusion

Conclusion

- ▶ How can modern ML models generalize well?
- ▶ Requires to understand dynamics
- ▶ ...

Thank you!

Conclusion

- ▶ How can modern ML models generalize well?
- ▶ Requires to understand dynamics
- ▶ ...

Thank you!