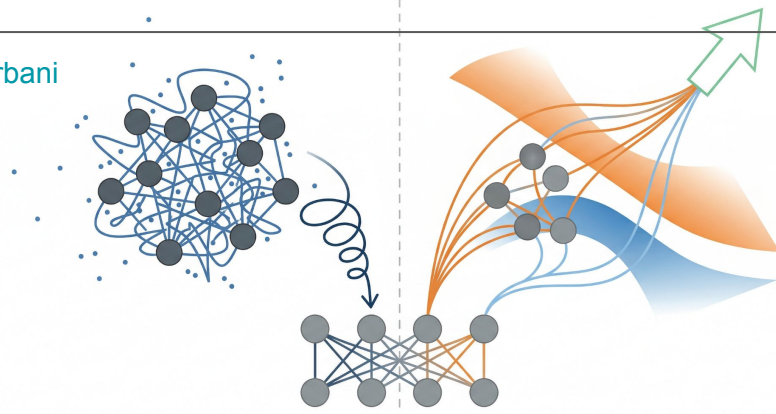


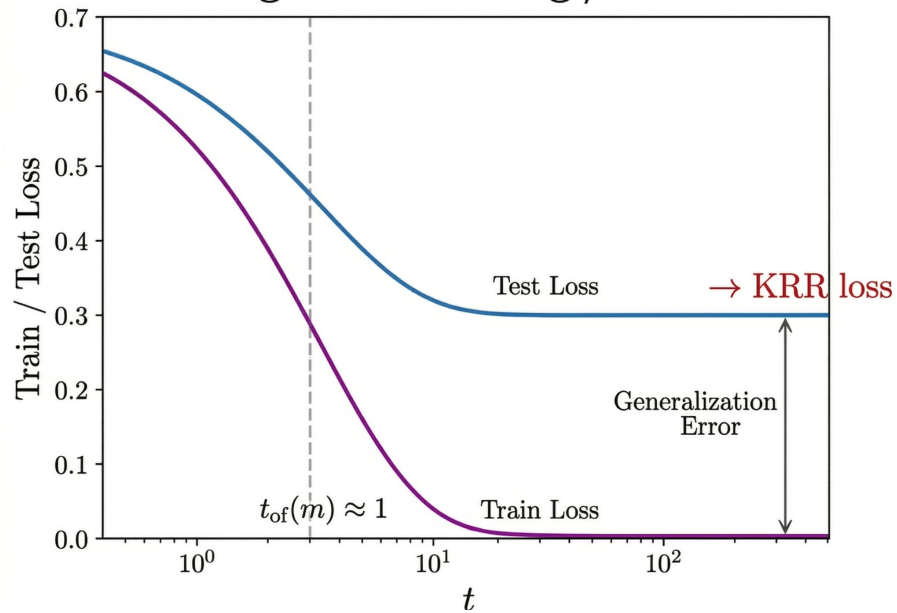
# Dynamical Decoupling of Generalization and Overfitting in Large Two-Layer Networks

Andrea Montanari, Pierfrancesco Urbani



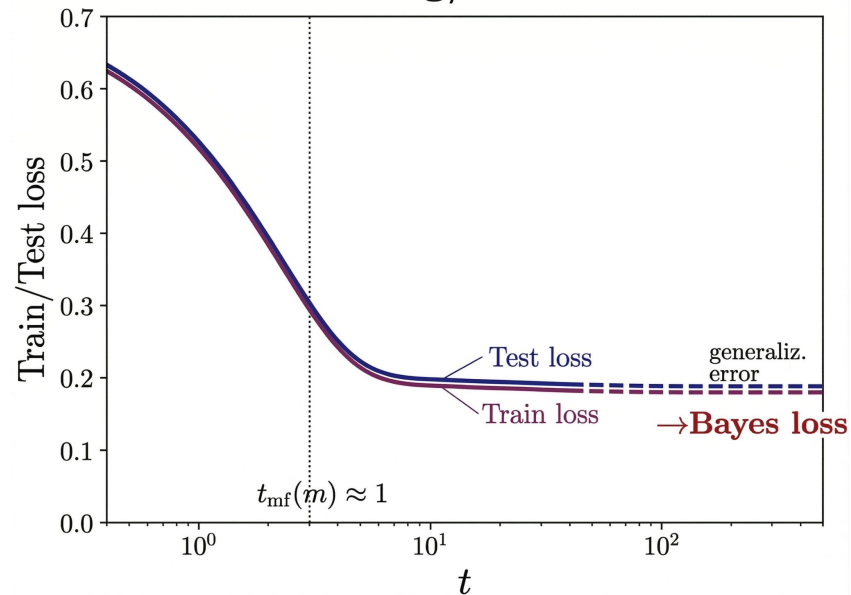
Mahdi Zamani

## Benign Overfitting / NTK



No feature learning

## Feature learning / Mean field



No overfitting?

# Setup

Model

$$f(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \frac{1}{m} \sum_{i=1}^m a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) \quad \|\mathbf{w}_i\|^2 = 1 \quad \boldsymbol{\theta} = (\mathbf{a}, \mathbf{W})$$

Data Generation

$$y_i = \varphi(\mathbf{U}^\top \mathbf{x}_i) + \varepsilon_i \quad \varepsilon_i \sim \mathbf{N}(0, \tau^2) \quad \mathbf{x}_i \sim \mathbf{N}(0, \mathbf{I}_d)$$

Empirical Risk

$$\widehat{\mathcal{R}}_n(\boldsymbol{\theta}) := \frac{1}{2n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2$$

Gradient Flow (GF)

$$\dot{\boldsymbol{\theta}}(t) = -\frac{n}{d} \mathbf{P}_{\boldsymbol{\theta}} \nabla \widehat{\mathcal{R}}_n(\boldsymbol{\theta}(t))$$

Overparameterization ratio

$$\alpha = \frac{n}{md}$$

# Pre-Techniques

## Gaussian Approximation

$$\widehat{\mathcal{R}}_n^g(\mathbf{a}, \mathbf{W}) = \frac{1}{2n} \sum_{i=1}^n (f_i^g(\mathbf{a}, \mathbf{W}) - \varphi_i^g - \varepsilon_i)^2$$

## Covariance Structure

$$\mathbb{E}\{f(\mathbf{x}; \mathbf{a}_1, \mathbf{W}_1) f(\mathbf{x}; \mathbf{a}_2, \mathbf{W}_2)\} = \frac{1}{m^2} \sum_{i,j=1}^m a_{1,i} a_{2,j} h(\langle \mathbf{w}_{1,i}, \mathbf{w}_{2,j} \rangle)$$

$$h(q) = \mathbb{E}\{\sigma(G_1)\sigma(G_2)\} \quad \mathbb{E}\{G_1 G_2\} = q$$

$$\mathbb{E}\{f(\mathbf{x}; \mathbf{a}, \mathbf{W}) \varphi(\mathbf{U}^\top \mathbf{x})\} = \frac{1}{m} \sum_{i=1}^m a_i \widehat{\varphi}(\mathbf{U}^\top \mathbf{w}_i)$$

$$\widehat{\varphi}(\mathbf{v}) := \mathbb{E}\left\{\sigma(\langle \mathbf{v}, \mathbf{G} \rangle + \sqrt{1 - \|\mathbf{v}\|^2} G_0) \varphi(\mathbf{G})\right\}$$

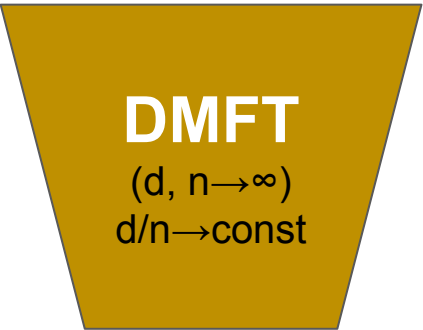
$\mathbf{G} \sim \mathbf{N}(0, \mathbf{I}_k)$  independent of  $G_0 \sim \mathbf{N}(0, 1)$

# DMFT

$$\begin{aligned} \dot{\mathbf{a}}(t) &= -\frac{n}{d} \nabla_{\mathbf{a}} \widehat{\mathcal{R}}_n(\mathbf{a}(t), \mathbf{W}(t)), \\ \dot{\mathbf{w}}_i(t) &= -\frac{n}{d} \nabla_{\mathbf{w}_i} \widehat{\mathcal{R}}_n(\mathbf{a}(t), \mathbf{W}(t)) - \nu_i(t) \mathbf{w}_i(t) \quad \forall i = 1, \dots, m. \end{aligned}$$

## Interacting ODEs

Dynamical Partition function  
 +  
 Grassmann Variables  
 +  
 Concentration of  
 Path Integral around its saddle point



## Average behavior of GF

$$\begin{aligned} \frac{da_i(t)}{dt} &= -\frac{\bar{\alpha}}{m} \int_0^t R_A(t, s) \left[ \frac{1}{m} \sum_{l=1}^m a_l(s) h(C_{li}(s, t)) - \hat{\varphi}(\mathbf{v}_i(t)) \right] ds \\ &\quad - \frac{\bar{\alpha}}{m} \int_0^t C_A(t, s) \frac{1}{m} \sum_{l=1}^m a_l(s) h'(C_{li}(s, t)) R_{il}(t, s) ds, \end{aligned} \tag{C.7}$$

# of equations O(m^2)

$$\frac{d\mathbf{v}_i(t)}{dt} = -\nu_i(t) \mathbf{v}_i(t) + \frac{\bar{\alpha}}{m} a_i(t) \nabla \hat{\varphi}(\mathbf{v}_i(t)) \int_0^t R_A(t, s) ds - \frac{1}{m} \sum_{j=1}^m \int_0^t M_{ij}^R(t, s) \mathbf{v}_j(s) ds, \tag{C.8}$$

Effectively one ODE  
 (particle) + noise + memory

$$\begin{aligned} \frac{\partial C_{ij}(t_a, t_b)}{\partial t_a} &= -\nu_i(t_a) C_{ij}(t_a, t_b) + \frac{\bar{\alpha}}{m} a_i(t_a) \langle \nabla \hat{\varphi}(\mathbf{v}_i(t_a)), \mathbf{v}_j(t_b) \rangle \int_0^{t_a} R_A(t_a, s) ds \\ &\quad - \frac{1}{m} \sum_{l=1}^m \int_0^{t_a} M_{il}^R(t_a, s) C_{lj}(s, t_b) ds - \frac{1}{m} \sum_{l=1}^m \int_0^{t_b} M_{il}^C(t_a, s) R_{jl}(t_b, s) ds, \end{aligned} \tag{C.9}$$

$$\frac{\partial R_{ij}(t_a, t_b)}{\partial t_a} = -\nu_i(t_a) R_{ij}(t_a, t_b) + \delta_{ij} \delta(t_a - t_b) - \frac{1}{m} \sum_{l=1}^m \int_{t_b}^{t_a} M_{il}^R(t_a, s) R_{lj}(s, t_b) ds. \tag{C.10}$$

# SymmDMFT

Can be significantly simplified by considering a symmetric initialization

$$\mathbf{w}_i^n(0) \sim \text{Unif}(\mathbb{S}^{d-1}) \quad a_i^n(0) = a_0$$

This symmetric initialization is preserved by GF and we end up with  $O(3 + 1)$  integro-differential equations and can derive Test and Train error's (limiting) evolution:

$$e_{\text{tr}}(t) = -\frac{1}{2}C_A(t, t),$$
$$e_{\text{ts}}(t) = \frac{1}{2} \left[ \tau^2 + \|\varphi\|^2 - 2a(t)\hat{\varphi}(\mathbf{v}(t)) + \frac{1}{m}a^2(t)h(1) + \frac{m-1}{m}a^2(t)h(C_o(t, t)) \right]$$

# Singular Perturbation Theory

However, it's still hard to derive the limiting distribution for the evolution of weights. i.e  $w(t)$  and  $a(t)$

So we also let  $m$  go to infinity!  $\sqrt{\epsilon}$

- (i) Hypothesize a certain asymptotic behavior of the DMFT solution in a specific time-scale
- (ii) Check consistency with the DMFT equations
- (iii) Check that this behavior is observed in the numerical solution of the DMFT equations.

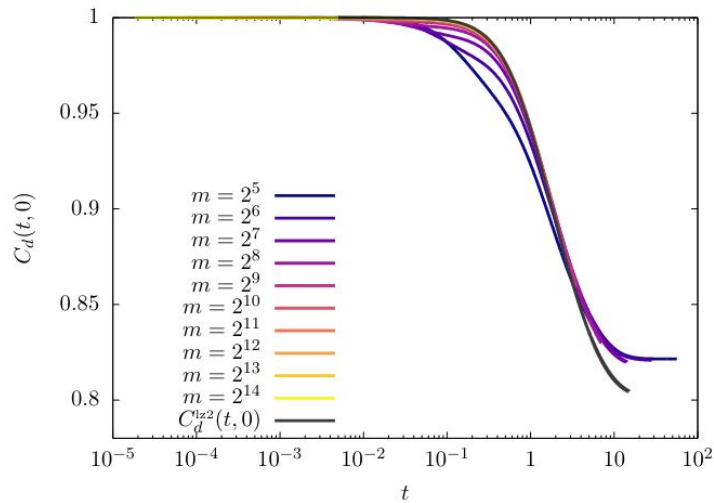
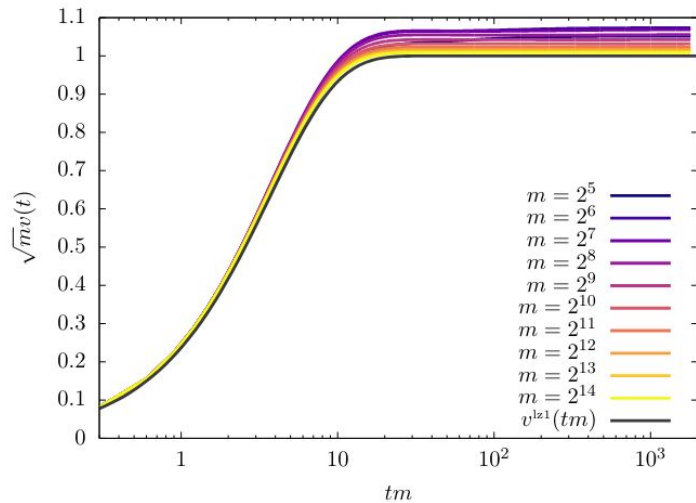
And this depends heavily on initialization of the second layer weights.

Lazy initialization

$$\gamma(t) = a(t)/\sqrt{m} \text{ (in particular, } \gamma(0) = \gamma_0)$$

# First dynamical regime: $t = O(1/m)$

On this timescale, the SymmDMFT equations are solved, up to higher order terms, by an ansatz



$$\|\mathbf{w}_i(t) - \mathbf{w}_i(0)\| = \Theta(1/\sqrt{m})$$

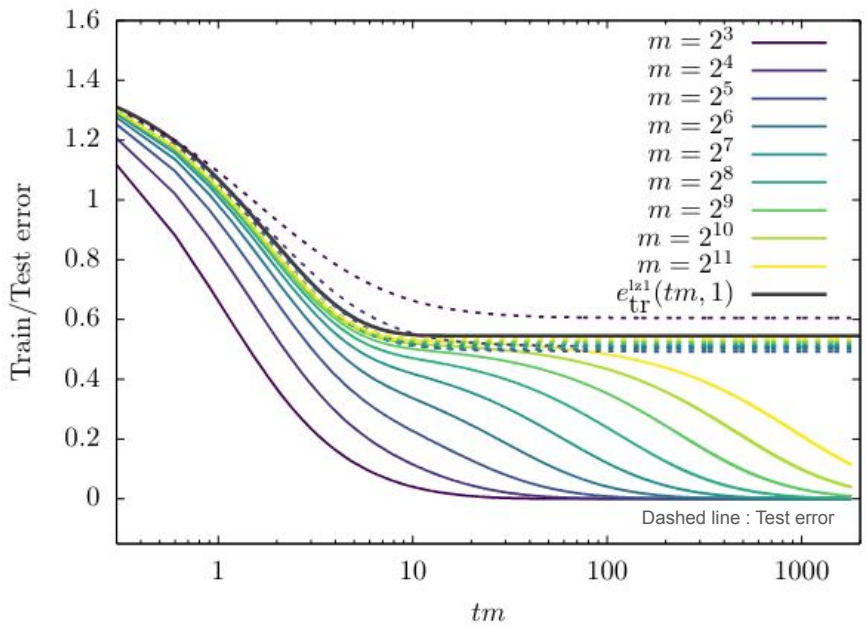
$$\|\mathbf{a}(t) - \mathbf{a}(0)\| = o_m(1)$$

$$\mathbf{v}_i(t) := \text{p-lim}_{n,d \rightarrow \infty} \mathbf{U}^\top \mathbf{w}_i^n(t)$$

$$C_{ij}(t_1, t_2) := \text{p-lim}_{n,d \rightarrow \infty} \langle \mathbf{w}_i^n(t_1), \mathbf{w}_j^n(t_2) \rangle$$

# Second dynamical regime: $t = \Theta(1)$

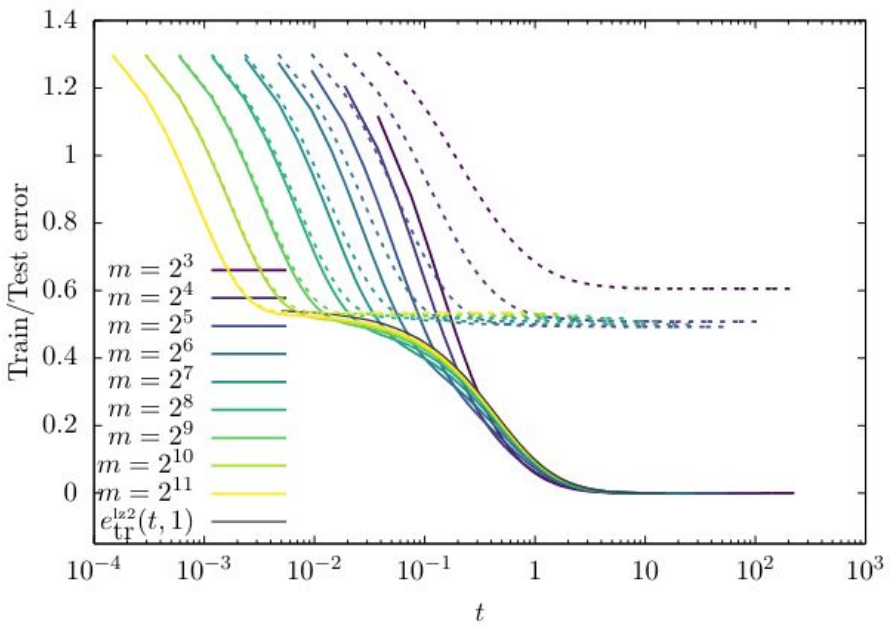
$t = \Theta(1/m)$



$train \approx test$

$train\ error\ change \approx O(1)$

$t = \Theta(1)$

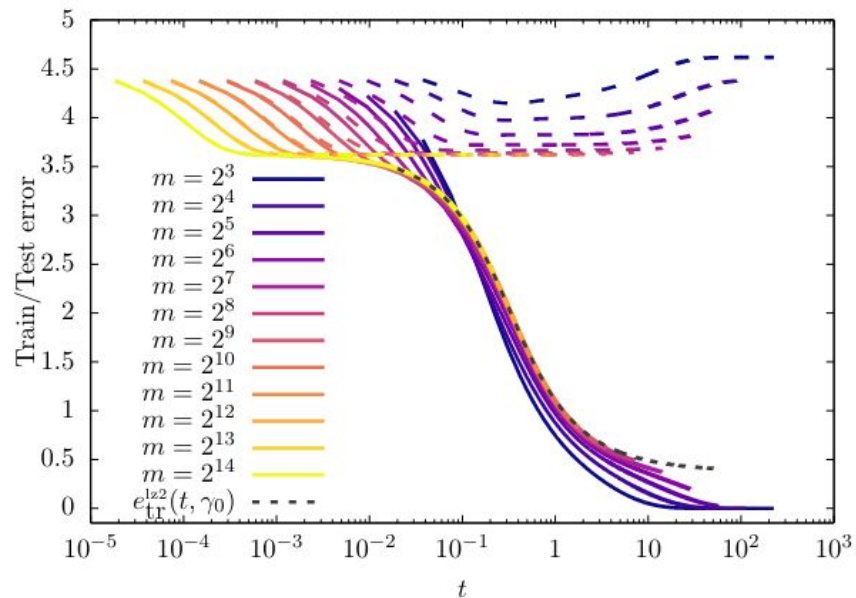


$$\|w_i(t) - w_i(0)\| = \Theta(1)$$

$$\|a(t) - a(0)\| = o_m(1)$$

No additional learning / memorization / overfitting

Third dynamical regime:  $t = \Theta(m)$



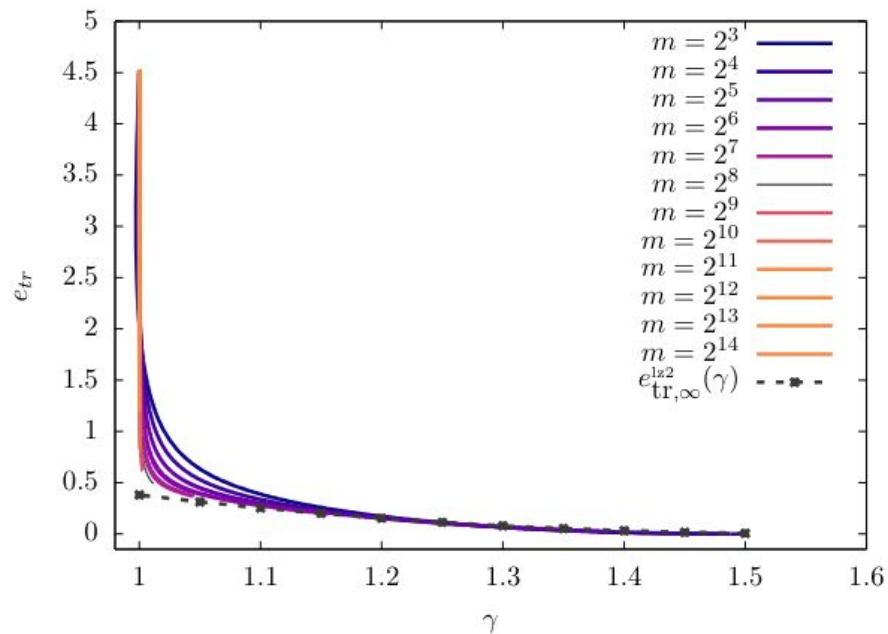
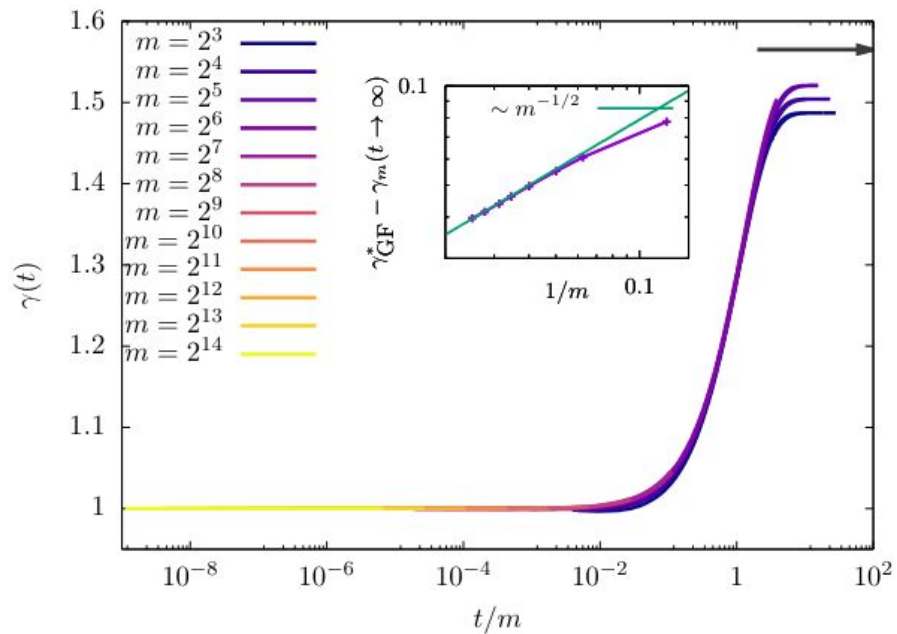
$$\|\mathbf{w}_i(t) - \mathbf{w}_i(0)\| = \Theta(1)$$

$$\|\mathbf{a}(t) - \mathbf{a}(0)\| = \Theta(1)$$

train  $\downarrow 0$

test  $\uparrow$

Third dynamical regime:  $t = \Theta(m)$



$$\gamma_0 < \gamma_{\text{GF}}^*(\alpha, \varphi, \tau) \quad a(t) = \gamma(t)\sqrt{m}$$

# Mean-Field initialization

$$a(0) = 1$$

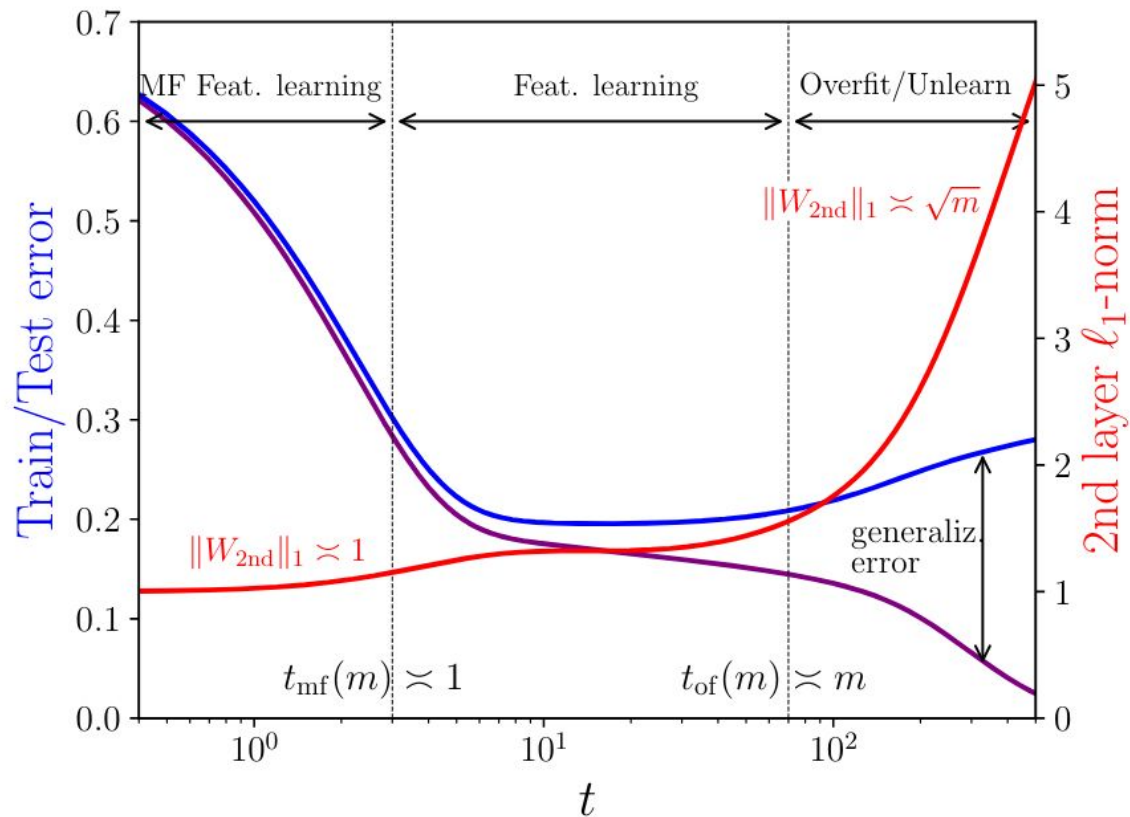
$t$	1st Layer Weights	2nd Layer Weights	Train Error	Test Error
$\Theta(1)$	$O(1)$	$O(1)$	$O(1)$ Reaches plateau	Test=Train=Bayes risk
$\Theta(m)$	$O(1)$	$O(\sqrt{m})$	Decays to 0	Increases

$$\varphi = \sigma \quad h'(0) > 0$$

$$v^{\text{mf1}}(t) \rightarrow 1$$

Fixed point

# Main Result



A lot more to learn from the paper

- Had some cool connections with spin glasses
- Lower bounding the overfitting timescale for finite  $n, d$  (MuP)
- Addressing implicit bias by the second layer weights
- Mathematical Insights

Thanks

# Questions