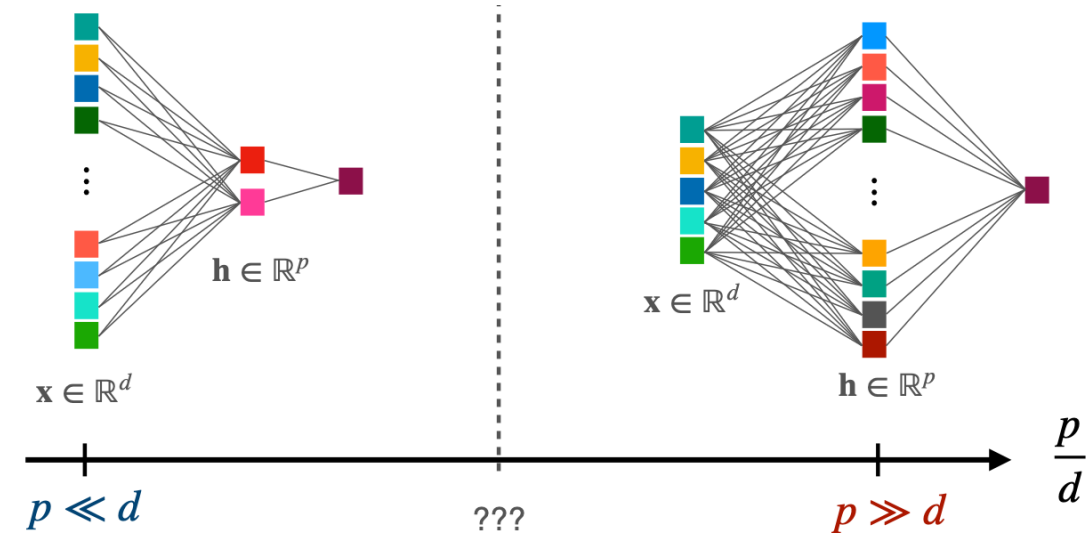




### Motivation

- Understanding the performance of SGD in neural networks is a major endeavor in machine learning, and significant progress was achieved in the context of large two-layer neural net
- High dimensional limit has been investigated first in the seminal work of [3], developing some ODEs for the dynamics.
- The optimization over wide two-layer neural networks can be rigorously studied using a well-defined PDE [4, 5]



**Aim:** drawing a precise connection between two limits

### Teacher-student model with online SGD

We introduce a teacher-student two-layer neural network model for studying the dynamics of the training with SGD:

- Input data is generated from independent Gaussian distributions:

$$\mathbf{x}^\nu \sim \mathcal{N}\left(\mathbf{0}_d, \frac{1}{d} \mathbf{I}_d\right)$$

Labels are generated by a **teacher network**

$$y^\nu = \frac{1}{k} \sum_{r=1}^k a_r^* \sigma^*(\mathbf{w}_r^{\top} \mathbf{x}^\nu) + \sqrt{\Delta} z^\nu, \quad z^\nu \sim \mathcal{N}(0, 1)$$

where  $\Delta$  is the artificial noise.

- The student network to be learned is

$$f_\Theta(\mathbf{x}) = \frac{1}{p} \sum_{i=1}^p a_i \sigma(\mathbf{w}_i^{\top} \mathbf{x})$$

- We are using the **square loss function**. The population risk is given by:

$$\mathcal{R}(\Theta) := \mathbb{E}_{(\mathbf{x}, y) \sim \rho} \left[ \frac{1}{2} (f_\Theta(\mathbf{x}) - y)^2 \right] + \frac{\Delta}{2}$$

- We are using the online stochastic gradient descent:

$$\Theta^{\nu+1} = \Theta^\nu - \gamma \nabla_{\Theta} \ell(f_{\Theta^\nu}(\mathbf{x}^\nu), y^\nu), \quad \nu \leq n$$

### Assumptions

#### Simplifying assumptions

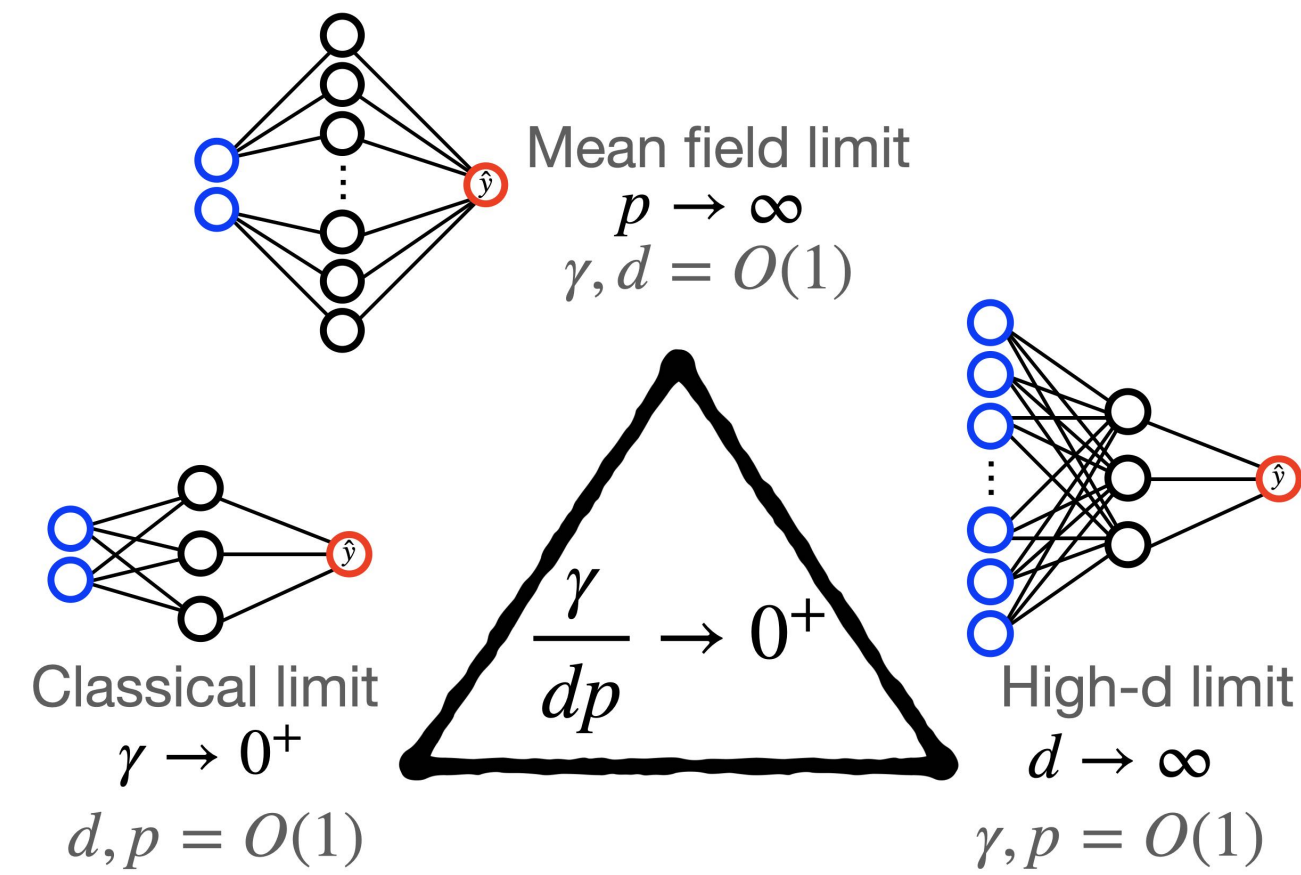
- $a_r^* = 1$  and  $a_i^\nu = 1$ ;
- $W^*$  is full-rank;
- $p$  is divisible by  $k$ ;

#### Technical assumptions

- With high probability:
- $\forall i \in [p], \|\mathbf{w}_i\| \leq K$ .
- $\|\sigma^{(i)}\|_\infty \leq K$  for  $i = 0, 1, 2$ .

### Main goal

We aim for a **common description for 3 different limits**



In particular a low-dimensional analysis for joint high-dimensional mean field limit!

### Low-dimensional sufficient statistics

Overlaps are **sufficient statistics**:

$$\Omega^\nu := \begin{pmatrix} Q^\nu & M^\nu \\ M^{\nu\top} & P \end{pmatrix} = \frac{1}{d} \begin{pmatrix} W^\nu W^{\nu\top} & W^\nu W^{*\top} \\ W^{*\top} W^\nu & W^* W^{*\top} \end{pmatrix} \in \mathbb{R}^{(p+k) \times (p+k)}$$

We can derive a closed set of stochastic processes

$$\begin{aligned} M_{ir}^{\nu+1} - M_{ir}^\nu &= \frac{\gamma}{pd} \sigma'(\lambda_i^\nu) \lambda_r^\nu \mathcal{E}^\nu \\ Q_{ij}^{\nu+1} - Q_{ij}^\nu &= \frac{\gamma}{pd} \left( \sigma'(\lambda_i^\nu) \lambda_j^\nu + \sigma'(\lambda_j^\nu) \lambda_i^\nu \right) \mathcal{E}^\nu \\ &\quad + \frac{\gamma^2 \|\mathbf{x}^\nu\|_2^2}{p^2 d} \sigma'(\lambda_i^\nu) \sigma'(\lambda_j^\nu) \mathcal{E}^{\nu 2} \end{aligned} \quad (\text{OV-SP})$$

where the *local fields* are jointly Gaussian vectors

$$(\boldsymbol{\lambda}^\nu, \boldsymbol{\lambda}^{\nu\nu}) \sim \mathcal{N}(\mathbf{0}_{p+k}, \Omega^\nu).$$

Informally, when  $\frac{\gamma}{pd}$  there is ODEs approximation

$$\begin{aligned} \frac{dM}{dt} &= \Psi^{(M)}(\Omega), \\ \frac{dQ}{dt} &= \Psi^{(GF)}(\Omega) + \frac{\gamma}{p} \Psi^{(Var)}(\Omega), \end{aligned} \quad (\text{SS-ODE})$$

$$\begin{aligned} \Psi_{ir}^{(M)}(\Omega) &= \mathbb{E}_{(\boldsymbol{\lambda}, \boldsymbol{\lambda}^*) \sim \mathcal{N}(\mathbf{0}_{p+k}, \Omega)} \left[ \sigma'(\lambda_i) \lambda_r^* \mathcal{E} \right] \\ \Psi_{ij}^{(GF)}(\Omega) &= \mathbb{E}_{(\boldsymbol{\lambda}, \boldsymbol{\lambda}^*) \sim \mathcal{N}(\mathbf{0}_{p+k}, \Omega)} \left[ \left( \sigma'(\lambda_i) \lambda_j + \sigma'(\lambda_j) \lambda_i \right) \mathcal{E} \right] \\ \Psi_{ij}^{(Var)}(\Omega) &= \mathbb{E}_{(\boldsymbol{\lambda}, \boldsymbol{\lambda}^*) \sim \mathcal{N}(\mathbf{0}_{p+k}, \Omega)} \left[ \sigma'(\lambda_i) \sigma'(\lambda_j) \mathcal{E}^2 \right] \end{aligned}$$

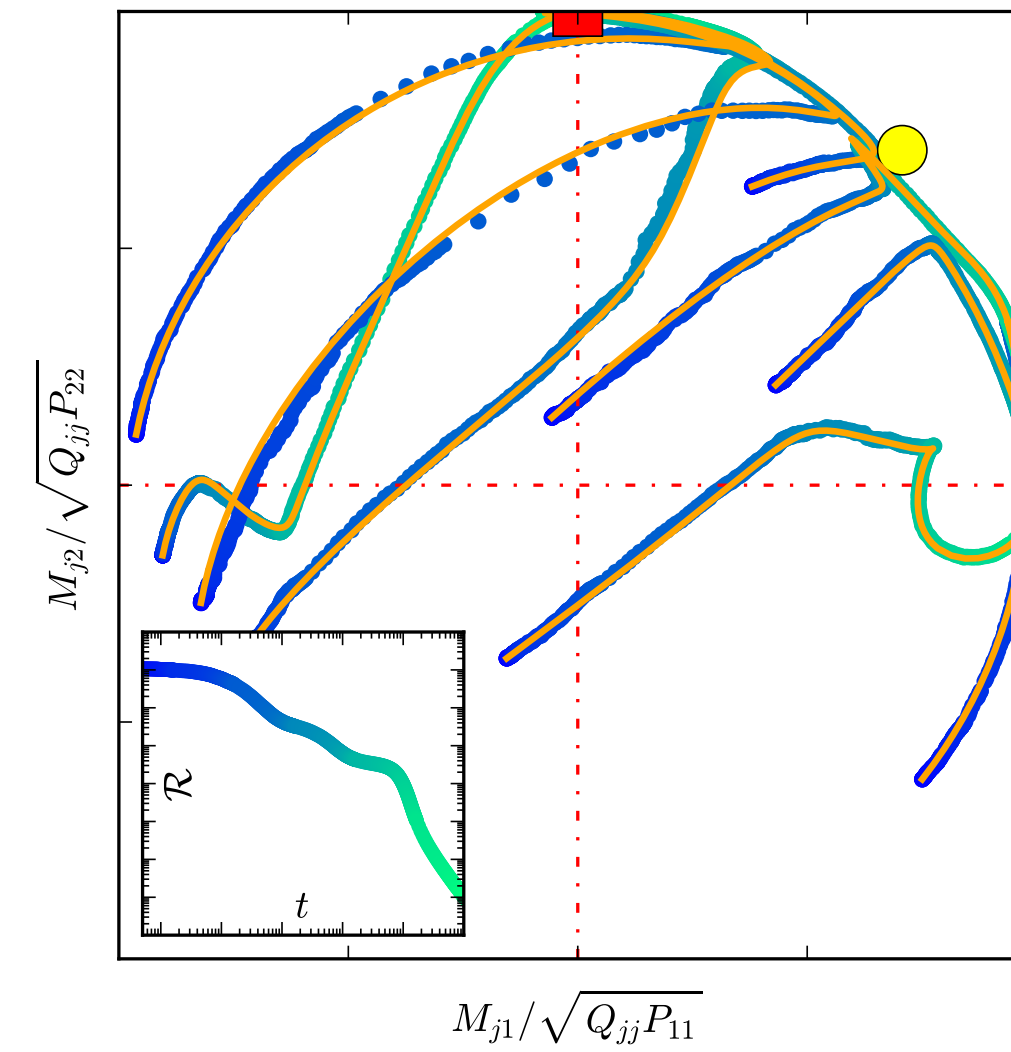
### Theorem (Veiga et al. [2])

Let  $\Omega^\nu$  be the random process of Eq. (OV-SP), and  $\Omega(t)$  the solution to the ODE (SS-ODE) with starting point  $\Omega(0) = \Omega^0$ . Define the stepsize  $\delta t = \frac{\gamma}{pd}$ , and assume that  $\gamma/p = O(1)$ . Then there exists a constant  $C > 0$  such that for any  $\nu \geq 0$ ,

$$\|\Omega^\nu - \Omega(\nu \delta t)\|_\infty \leq e^{C\nu \delta t} \sqrt{\frac{\gamma}{pd}}$$

### Key observation

(SS-ODE) hold whenever  $\frac{\gamma}{pd} \rightarrow 0$ , not only when  $d \rightarrow +\infty$ .



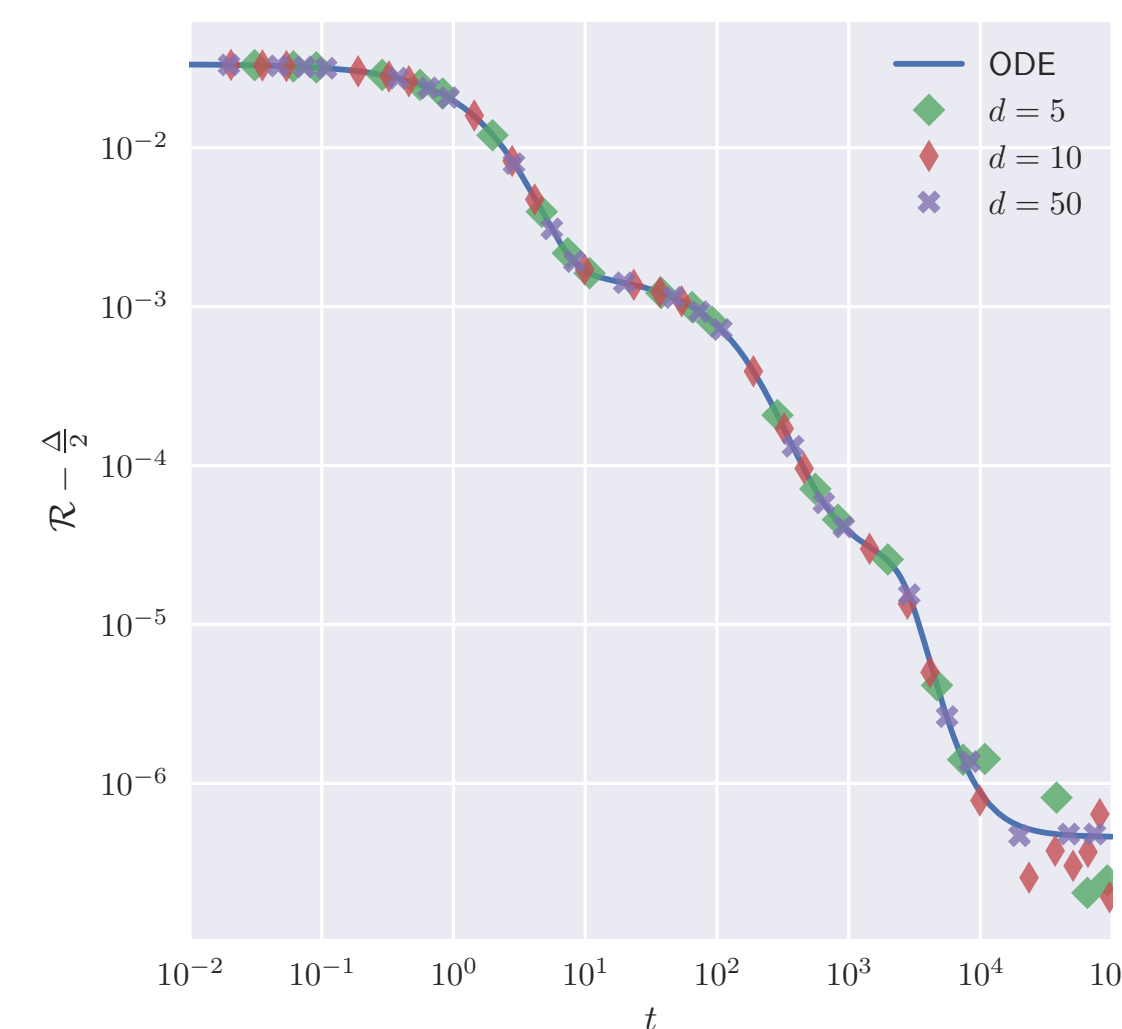
### Classical limit ( $\gamma \rightarrow +0$ )

The equations are equivalent to the gradient flow equations

$$\frac{d\mathbf{w}_i}{dt} = \nabla_{\mathbf{w}_i} \mathcal{R}(\Theta)$$

### Dimension Independence

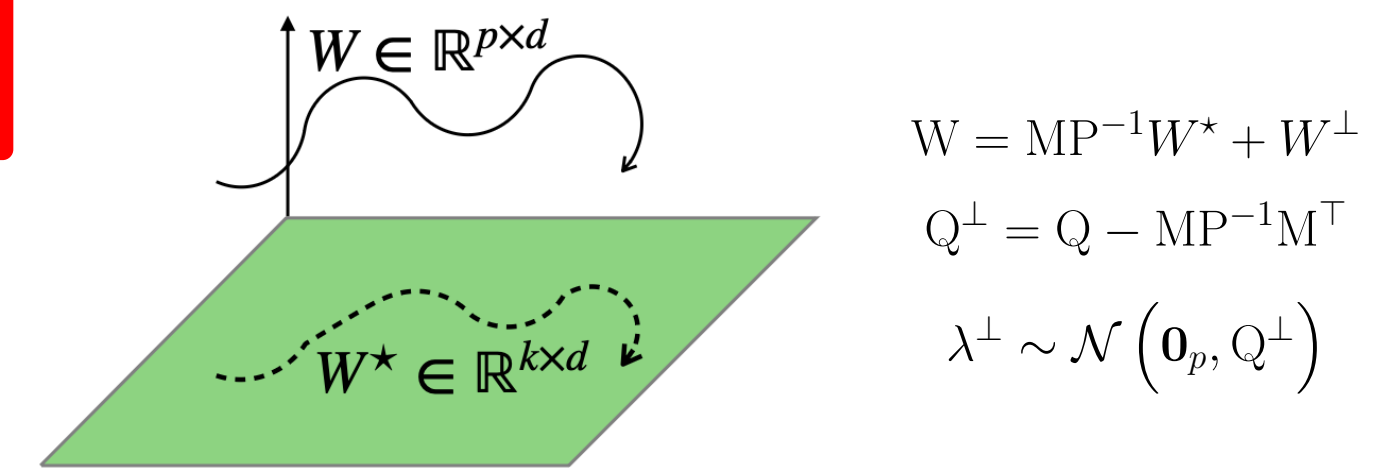
(SS-ODE) is independent of the data dimension  $d$ .



The trajectories are **exactly in the same** whether  $d$  is large or small, when starting from the same initial condition.

### Overparametrized regime ( $p \rightarrow +\infty$ )

Decompose the student between the teacher's space and its orthogonal



$$\begin{aligned} W &= MP^{-1}W^* + W^\perp \\ Q^\perp &= Q - MP^{-1}M^\top \\ \lambda^\perp &\sim \mathcal{N}(\mathbf{0}_p, Q^\perp) \end{aligned}$$

$Q^\perp$  can replace  $Q$  and the overlaps are still sufficient, with ODE:

$$\frac{dQ_{ij}^\perp}{dt} = \mathbb{E}_{(\boldsymbol{\lambda}^\perp, \boldsymbol{\lambda})} \left[ \left( \sigma'(\lambda_i) \lambda_j^\perp + \sigma'(\lambda_j) \lambda_i^\perp \right) \mathcal{E} \right] := \Psi_{ij}^\perp(\Omega).$$

### Low dimensional mean-field

**Assumption** ( $\mathbf{w}_1, \dots, \mathbf{w}_p$ ) are drawn i.i.d from an orthogonally invariant and  $\frac{K^2}{d}$ -subgaussian distribution.

**Ansatz** Just in the dynamic equation, the weights

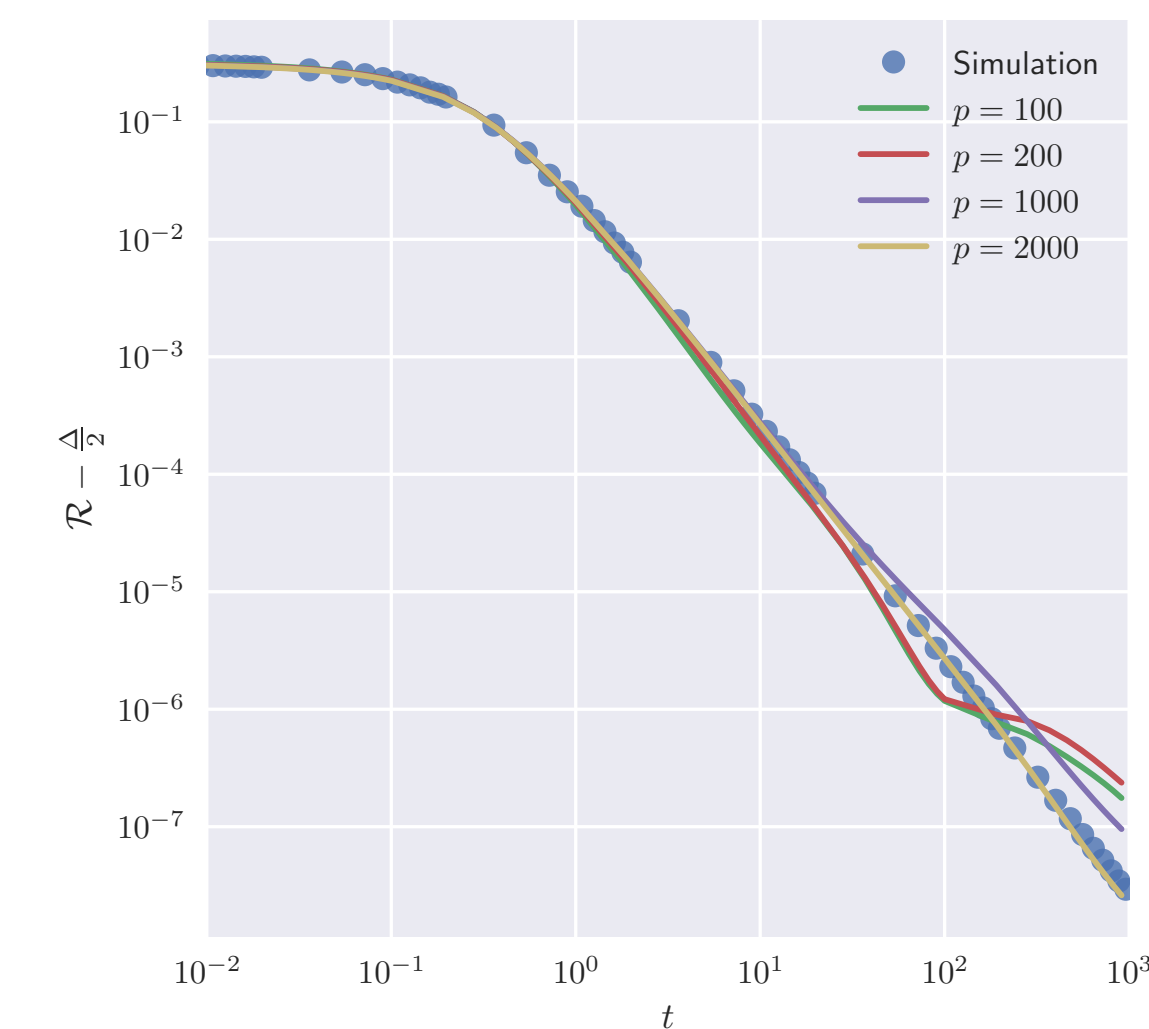
$$\mathbf{w}_i^\perp \approx \sqrt{Q_{ii}^\perp} \cdot \mathbf{g}_i, \quad \tilde{Q} = MP^{-1}M^\top + D \sqrt{Q^\perp} \Xi D \sqrt{Q^\perp}$$

where  $D \sqrt{Q^\perp} = \sqrt{\text{diag}(Q^\perp)}$  and

$$\Xi_{ii} = 1, \quad \Xi_{ij} = \langle \mathbf{g}_i, \mathbf{g}_j \rangle \quad \text{with } \mathbf{g}_i, \mathbf{g}_j \sim \text{Unif}(\mathbb{S}^{d-k-1})$$

**Reduced ODEs and risk**

$$\begin{aligned} \frac{dM}{dt} &= \mathbb{E}_\Xi \left[ \Psi^{(M)}(\tilde{\Omega}) \right] & \frac{dQ_{ii}^\perp}{dt} &= \mathbb{E}_\Xi \left[ \Psi_{ii}^\perp(\tilde{\Omega}) \right] & (\text{MF-ODE}) \\ \mathcal{R}(\tilde{\Theta}) &= \mathbb{E}_\Xi \left[ \mathcal{R}(\tilde{\Omega}) \right] \end{aligned}$$



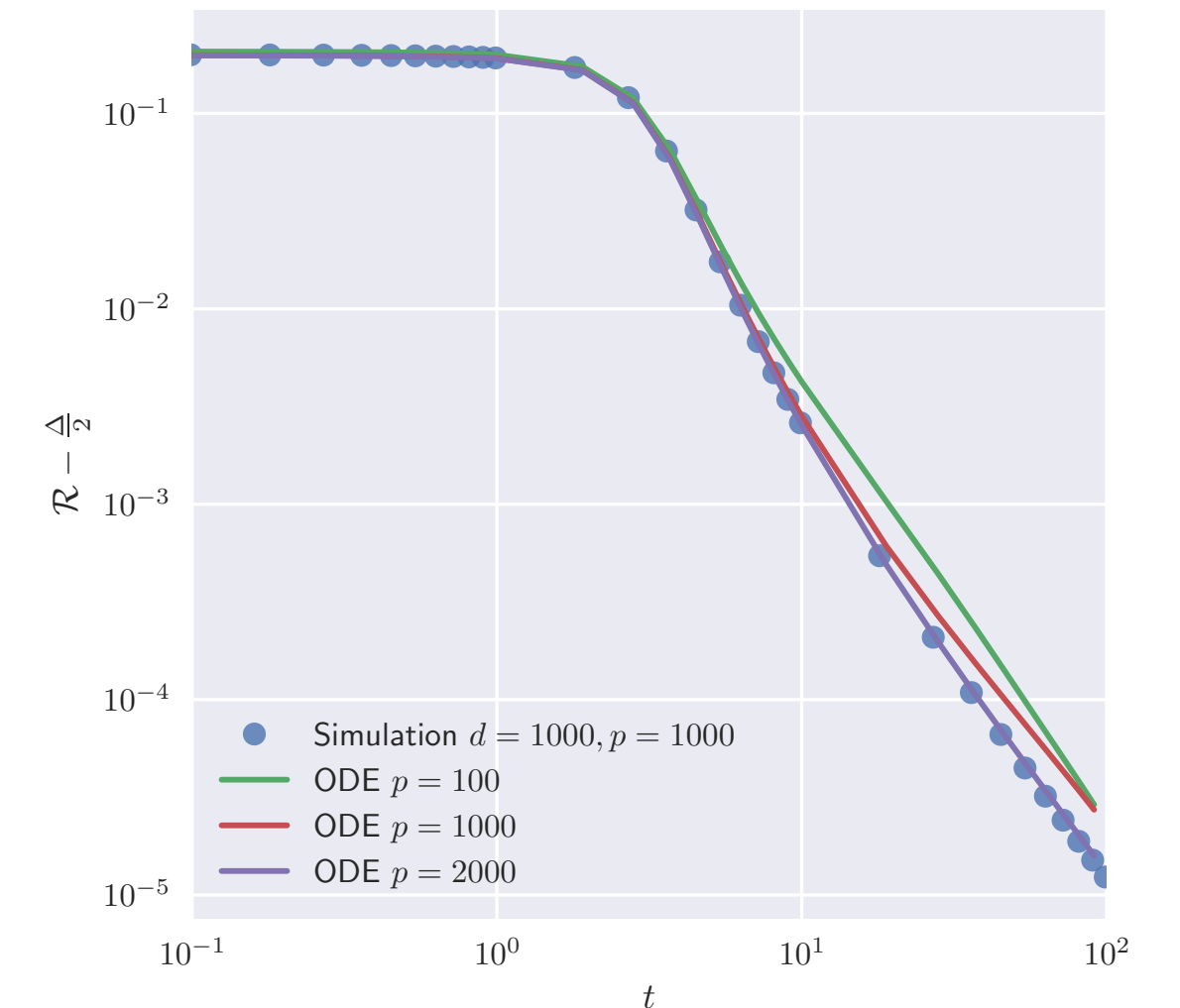
**Theorem.** Let  $\Omega(t)$  and  $\tilde{\Theta}(t)$  denote the solutions of (SS-ODE) and (MF-ODE), respectively. Then with probability at least  $1 - e^{-z^2}$  on the initialization:

$$\sup_{t \in [0, T]} \left| \mathcal{R}(\Omega(t)) - \mathcal{R}(\tilde{\Theta}(t)) \right| \leq C e^{CT} \left( \sqrt{\log(pT)} + z \right) / \sqrt{p}.$$

### High-dimensional mean-field

When  $d \rightarrow +\infty$  the random matrix  $\Xi$  is the identity. (MF-ODE) becomes

$$\frac{dM}{dt} = \Psi^{(M)}(\tilde{\Omega}) \quad \frac{dQ_{ii}^\perp}{dt} = \Psi_{ii}^\perp(\tilde{\Omega}). \quad (\text{HDMF-ODE})$$



**Theorem.** Let  $\Omega(t)$  and  $\tilde{\Theta}(t)$  denote the solutions of (SS-ODE) and (HDMF-ODE), respectively. Then with probability at least  $1 - e^{-z^2}$  on the initialization:

$$\sup_{t \in [0, T]} \left| \mathcal{R}(\Omega(t)) - \mathcal{R}(\tilde{\Theta}(t)) \right| \leq C e^{CT} \left( \frac{\sqrt{\log(pT)} + z}{\sqrt{p}} + \frac{1}{\sqrt{d}} \right)$$

(HDMF-ODE) are the **particle dynamics** of a measure over the sufficient statistics. Introducing a measure  $\mu_{(m,q)}$  over  $\mathbb{R}^{k+1}$  as

$$\mu_t := h_{\#} \mu_t \quad \text{where } h(\mathbf{w}) = \left( \frac{W^* \mathbf{w}}{d}, \|\mathbf{w}^\perp\|^2 \right)$$

the evolution can be written as Wasserstein GD

$$\partial_t \mu_{(m,q)} = \nabla_{(m,q)} \cdot \left( \mu_{(m,q)} \varphi(\cdot, \mu_{(m,q)}) \right)$$

### References

- [1] *From high-dimensional & mean-field dynamics to dimensionless ODEs: A unifying approach to SGD in two-layers networks*, Luca Arnaboldi, Ludovic Stephan, Florent Krzakala, Bruno Loureiro arXiv: 2302.05882 [stat.ML]
- [2] *Phase diagram of Stochastic Gradient Descent in high-dimensional two-layer neural networks*, Rodrigo Veiga et al. Advances in Neural Information Processing Systems, 2022
- [3] *On-line learning in soft committee machines*, David Saad, Sara A. Solla. Phys. Rev. E, 52:4225–4243, Oct 1995.
- [4] *On the global convergence of gradient descent for over-parameterized models using optimal transport*, L enaic Chizat and Francis Bach. Advances in Neural Information Processing Systems, 2018
- [5] *A mean field view of the landscape of two-layer neural networks*, Song Mei, Andrea Montanari, and Phan-Minh Nguyen Proceedings of the National Academy of Sciences, 115(33):E7665–E7671, 2018.