# Wonders of high-dimensions:
## the maths and physics of ML

**Bruno Loureiro**
Département d'Informatique
École Normale Supérieure & CNRS

brloureiro@gmail.com

ACDL 2023, 10-14.06.2023

# Part III

## Two layer neural networks in the rich regime

SGD dynamics of two-layer NNs

# Recall from last lecture...

Assuming that $a_{0,i} = O(1)$ and introducing a scaling:

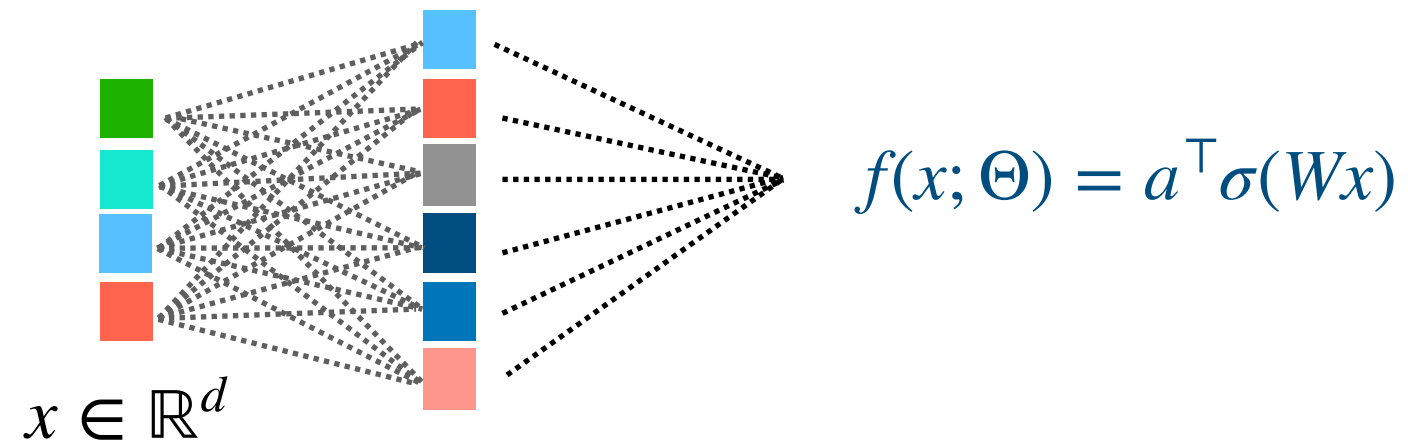$$f(x; \Theta) = \alpha(p) \sum_{i=1}^{p} a_i \sigma(w_i^\top x)$$

It can be shown that for $p \gg 1$:  [Chizat, Oyallon & Bach '19]

$$\mathbb{E}[\kappa(\Theta_0)] \lesssim \frac{1}{\sqrt{p}} + \frac{1}{p\alpha(p)}$$

Which means $f(x; \Theta) \approx \bar{f}_{\text{lin}}(x; \Theta_0)$ if $p\alpha(p) \to \infty$ as $p \to \infty$

a.k.a. "lazy" regime
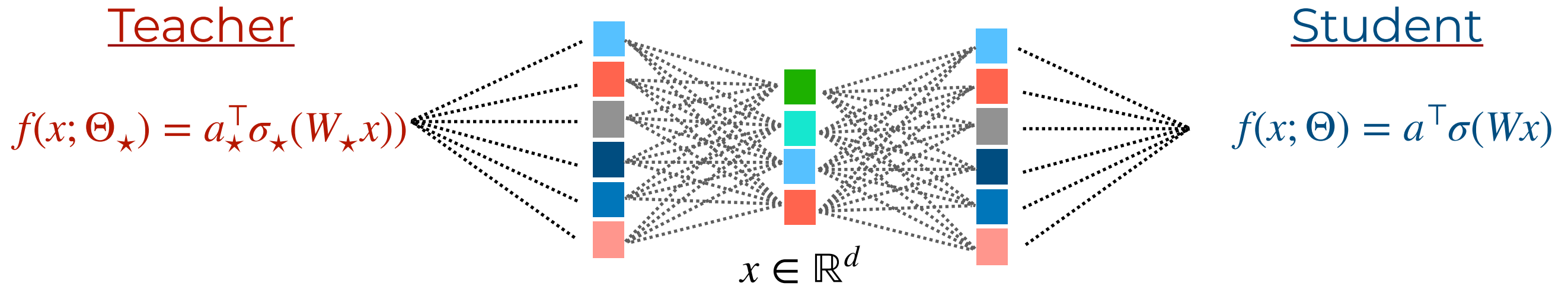
# Teacher-student setup



$$f(x; \Theta) = a^\top \sigma(Wx)$$

$x \in \mathbb{R}^d$

Hypothesis: $\quad f(x; \Theta) = \dfrac{1}{p} \displaystyle\sum_{i=1}^{p} a_i \sigma(w_i^\top x)$

# Teacher-student setup

Teacher

$$f(x; \Theta_\star) = a_\star^\top \sigma_\star(W_\star x))$$



$x \in \mathbb{R}^d$

Student

$$f(x; \Theta) = a^\top \sigma(Wx)$$

Hypothesis: $\quad f(x; \Theta) = \dfrac{1}{p} \sum_{i=1}^{p} a_i \sigma(w_i^\top x)$

Data: $(x^\nu, y^\nu)_{\nu \in [n]} \in \mathbb{R}^d \times \mathcal{Y}$ generated as:

$$y^\nu = \frac{1}{k} \sum_{r=1}^{k} a_{\star,r} \sigma(w_\star^\top x^\nu) + \sqrt{\Delta} z^\nu$$

$$x^\nu \sim \mathcal{N}(0, I_d)$$
$$z^\nu \sim \mathcal{N}(0, 1)$$

# Algorithm: SGD

Algorithm:    Let $b_k \subset [n]$ be mini-batch.

$$\Theta^{k+1} = \Theta^k - \gamma_k \nabla_{\Theta^k} \hat{\mathscr{R}}_{b_k} \left( \Theta^k \right)$$

$$\hat{\mathscr{R}}_b(\Theta) = \frac{1}{2|b|} \sum_{\nu \in b} \left( y^\nu - f(x^\nu; \Theta) \right)^2$$

mini-batch

# Algorithm: SGD

Algorithm:    Let $b_k \subset [n]$ be mini-batch.

$$\Theta^{k+1} = \Theta^k - \gamma_k \nabla_{\Theta^k} \hat{\mathscr{R}}_{b_k}\left(\Theta^k\right)$$

$$\hat{\mathscr{R}}_b(\Theta) = \frac{1}{2|b|} \sum_{\nu \in b} \left(y^\nu - f(x^\nu; \Theta)\right)^2$$

mini-batch

Gradient descent (GD)

$$b_k = [n], \quad \forall k$$

# Algorithm: SGD

Algorithm:    Let $b_k \subset [n]$ be mini-batch.

$$\Theta^{k+1} = \Theta^k - \gamma_k \nabla_{\Theta^k} \hat{\mathcal{R}}_{b_k} \left( \Theta^k \right)$$

$$\hat{\mathcal{R}}_b(\Theta) = \frac{1}{2|b|} \sum_{\nu \in b} \left( y^\nu - f(x^\nu; \Theta) \right)^2$$

mini-batch

Gradient descent (GD)

$$b_k = [n], \quad \forall k$$

$\gamma \to 0^+$ at fixed $d, p$:

$$\dot{\Theta}(t) = - \nabla_\Theta \hat{\mathcal{R}}_n \left( \Theta(t) \right)$$

# Algorithm: SGD

Algorithm:    Let $b_k \subset [n]$ be mini-batch.

$$\Theta^{k+1} = \Theta^k - \gamma_k \nabla_{\Theta^k} \hat{\mathscr{R}}_{b_k}\left(\Theta^k\right)$$

$$\hat{\mathscr{R}}_b(\Theta) = \frac{1}{2|b|} \sum_{\nu \in b} \left(y^\nu - f(x^\nu; \Theta)\right)^2$$

mini-batch

### Gradient descent (GD)

$$b_k = [n], \quad \forall k$$

$\gamma \to 0^+$ at fixed $d, p$:

$$\dot{\Theta}(t) = -\nabla_\Theta \hat{\mathscr{R}}_n\left(\Theta(t)\right)$$

### One-pass SGD

$b_k$ independent

# Algorithm: SGD

Algorithm:    Let $b_k \subset [n]$ be mini-batch.

$$\Theta^{k+1} = \Theta^k - \gamma_k \nabla_{\Theta^k} \hat{\mathcal{R}}_{b_k} \left( \Theta^k \right)$$

$$\hat{\mathcal{R}}_b(\Theta) = \frac{1}{2|b|} \sum_{\nu \in b} \left( y^\nu - f(x^\nu; \Theta) \right)^2$$

mini-batch

| Gradient descent (GD) | One-pass SGD |
|---|---|
| $b_k = [n], \quad \forall k$ | $b_k$ independent |
| $\gamma \to 0^+$ at fixed $d, p$: | $\gamma \to 0^+$ at fixed $d, p$: |
| $\dot{\Theta}(t) = -\nabla_\Theta \hat{\mathcal{R}}_n \left( \Theta(t) \right)$ | $\dot{\Theta}(t) = -\nabla_\Theta \mathcal{R} \left( \Theta(t) \right)$ |

# Algorithm: SGD

Algorithm:   Let $b_k \subset [n]$ be mini-batch.

$$\Theta^{k+1} = \Theta^k - \gamma_k \nabla_{\Theta^k} \hat{\mathcal{R}}_{b_k}\left(\Theta^k\right)$$

$$\hat{\mathcal{R}}_b(\Theta) = \frac{1}{2|b|} \sum_{\nu \in b} \left(y^\nu - f(x^\nu; \Theta)\right)^2$$

mini-batch

### Gradient descent (GD)

$$b_k = [n], \quad \forall k$$

$\gamma \to 0^+$ at fixed $d, p$:

$$\dot{\Theta}(t) = - \nabla_\Theta \hat{\mathcal{R}}_n\left(\Theta(t)\right)$$

### One-pass SGD

$b_k$ independent

$\gamma \to 0^+$ at fixed $d, p$:

$$\dot{\Theta}(t) = - \nabla_\Theta \mathcal{R}\left(\Theta(t)\right)$$

# Another look at SGD

Rewrite SGD:

$$\Theta^{k+1} = \Theta^k - \gamma_k \nabla_{\Theta^k} \mathscr{R}\left(\Theta^k\right) + \gamma_k \varepsilon^k$$

GD on population        Effective Noise

Where:

$$\varepsilon^k = \nabla_{\Theta^k} \left[ \mathscr{R}\left(\Theta^k\right) - \hat{\mathscr{R}}_{B_k}\left(\Theta^k\right) \right]$$

# Another look at SGD

Rewrite SGD:

$$\Theta^{k+1} = \underbrace{\Theta^k - \gamma_k \nabla_{\Theta^k} \mathscr{R}\left(\Theta^k\right)}_{\text{GD on population}} + \underbrace{\gamma_k \varepsilon^k}_{\substack{\text{Effective}\\\text{Noise}}}$$

Where:

$$\varepsilon^k = \nabla_{\Theta^k}\left[\mathscr{R}\left(\Theta^k\right) - \hat{\mathscr{R}}_{B_k}\left(\Theta^k\right)\right]$$

🤔 Question:   How to characterise this?

# Summary of setting

One-pass SGD for two-layer neural networks in the teacher-student setting.

Architecture:
$$f(x; \Theta) = \frac{1}{p} \sum_{i=1}^{p} a_i \sigma(w_i \cdot x)$$

Data model:
$$y^\nu = \frac{1}{k} \sum_{r=1}^{k} a_r^\star \sigma(w_r^\star \cdot x^\nu) + \sqrt{\Delta} z^\nu \qquad \begin{array}{l} x^\nu \sim \mathcal{N}(0, I_d) \\ z^\nu \sim \mathcal{N}(0, 1) \end{array}$$

Algorithm:
$$\Theta^{\nu+1} = \Theta^\nu - \gamma_\nu \nabla_{\Theta^\nu} (y^\nu - f(x^\nu; \Theta^\nu))^2$$

# Sufficient statistics

Goal: track population error exactly throughout the dynamics

$$\mathcal{R}(\Theta^\nu) = \frac{1}{2}\mathbb{E}_{x \sim \mathcal{N}(0, I_d)}\left[\left(\frac{1}{k}\sum_{r=1}^{k} a_{\star,r}\sigma(w_r^{*\top}x) - \frac{1}{p}\sum_{i=1}^{p} a_i^\nu \sigma(w_i^{\nu\top}x)\right)^2\right] + \frac{\Delta}{2}$$

# Sufficient statistics

Goal: track population error exactly throughout the dynamics

$$\mathscr{R}(\Theta^\nu) = \frac{1}{2} \mathbb{E}_{(\lambda_\star^\nu, \lambda^\nu) \sim \mathscr{N}(0,\Omega^\nu)} \left[ \left( \frac{1}{k} \sum_{r=1}^{k} a_{\star,r} \sigma(\lambda_{\star,r}^\nu) - \frac{1}{p} \sum_{i=1}^{p} a_i^\nu \sigma(\lambda_i^\nu) \right)^2 \right] + \frac{\Delta}{2}$$

Where:

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

# Sufficient statistics

Goal: track population error exactly throughout the dynamics

$$\mathscr{R}(\Theta^\nu) = \frac{1}{2}\mathbb{E}_{(\lambda_\star^\nu, \lambda^\nu)\sim\mathscr{N}(0,\Omega^\nu)}\left[\left(\frac{1}{k}\sum_{r=1}^{k}a_{\star,r}\sigma(\lambda_{\star,r}^\nu) - \frac{1}{p}\sum_{i=1}^{p}a_i^\nu\sigma(\lambda_i^\nu)\right)^2\right] + \frac{\Delta}{2}$$

Where:

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

💡 Key idea:

$$\text{One-pass SGD} \quad \longrightarrow \quad \Omega^{\nu+1} = \psi(\Omega^\nu)$$

# Tracking overlaps

Starting point: one-pass SGD

$$\Theta^{\nu+1} = \Theta^{\nu} - \frac{\gamma_{\nu}}{2} \nabla_{\Theta^{\nu}} (y^{\nu} - f(x^{\nu}; \Theta^{\nu}))^2$$

# Tracking overlaps

Starting point: one-pass SGD

$$\Theta^{\nu+1} = \Theta^{\nu} - \frac{\gamma_{\nu}}{2} \nabla_{\Theta^{\nu}} (y^{\nu} - f(x^{\nu}; \Theta^{\nu}))^2$$

$$= \Theta^{\nu} + \gamma_{\nu} (y^{\nu} - f(x^{\nu}; \Theta^{\nu})) \nabla_{\Theta^{\nu}} f(x^{\nu}; \Theta^{\nu})$$

# Tracking overlaps

Starting point: one-pass SGD

$$\Theta^{\nu+1} = \Theta^\nu - \frac{\gamma_\nu}{2} \nabla_{\Theta^\nu} (y^\nu - f(x^\nu; \Theta^\nu))^2$$

$$= \Theta^\nu + \gamma_\nu \textcolor{red}{(y^\nu - f(x^\nu; \Theta^\nu))} \nabla_{\Theta^\nu} f(x^\nu; \Theta^\nu)$$

$$= \Theta^\nu + \gamma_\nu \textcolor{red}{\mathscr{E}^\nu} \nabla_{\Theta^\nu} f(x^\nu; \Theta^\nu)$$

# Tracking overlaps

Starting point: one-pass SGD

$$\Theta^{\nu+1} = \Theta^{\nu} - \frac{\gamma_{\nu}}{2} \nabla_{\Theta^{\nu}} (y^{\nu} - f(x^{\nu}; \Theta^{\nu}))^2$$

$$= \Theta^{\nu} + \gamma_{\nu} \color{red}{(y^{\nu} - f(x^{\nu}; \Theta^{\nu}))} \color{black}{\nabla_{\Theta^{\nu}} f(x^{\nu}; \Theta^{\nu})}$$

$$= \Theta^{\nu} + \gamma_{\nu} \color{red}{\mathscr{E}^{\nu}} \color{black}{\nabla_{\Theta^{\nu}} f(x^{\nu}; \Theta^{\nu})}$$

But:

$$\nabla_{a_i} f(x; \Theta) = \frac{1}{p} \sigma(w_i^{\top} x)$$

$$\nabla_{w_i} f(x; \Theta) = \frac{1}{p} a_i \sigma'(w_i^{\top} x) x$$

# Tracking overlaps

Starting point: one-pass SGD

$$a_i^{\nu+1} = a_i^{\nu} + \frac{\gamma_{\nu}}{p} \mathscr{E}^{\nu} \sigma(w_i^{\nu\top} x^{\nu})$$

$$w_i^{\nu+1} = w_i^{\nu} + \frac{\gamma_{\nu}}{p} \mathscr{E}^{\nu} a_i \sigma'(w_i^{\nu\top} x^{\nu}) x^{\nu}$$

# Tracking overlaps

Starting point: one-pass SGD

$$a_i^{\nu+1} = a_i^{\nu} + \frac{\gamma_\nu}{p}\mathscr{E}^\nu\sigma(w_i^{\nu\top}x^\nu)$$

$$w_i^{\nu+1} = w_i^{\nu} + \frac{\gamma_\nu}{p}\mathscr{E}^\nu a_i\sigma'(w_i^{\nu\top}x^\nu)x^\nu$$

<u>Goal</u>: Go from this to equation for $\Omega^\nu$

# Tracking overlaps

Starting point: one-pass SGD

$$a_i^{\nu+1} = a_i^{\nu} + \frac{\gamma_\nu}{p} \mathscr{E}^\nu \sigma(w_i^{\nu\top} x^\nu)$$

$$w_i^{\nu+1} = w_i^{\nu} + \frac{\gamma_\nu}{p} \mathscr{E}^\nu a_i \sigma'(w_i^{\nu\top} x^\nu) x^\nu$$

<u>Goal</u>: Go from this to equation for $\Omega^\nu$

$$\mathscr{E}^\nu = (y^\nu - f(x^\nu; \Theta^\nu))$$

$$= \left( \frac{1}{k} \sum_{r=1}^{k} a_{\star,r} \sigma_\star(w_{\star,r}^\top x^\nu) - \sqrt{\Delta} z^\nu - \frac{1}{p} \sum_{i=1}^{p} a_i^\nu \sigma(w_i^{\nu\top} x^\nu) \right)$$

# Tracking overlaps

Starting point: one-pass SGD

$$a_i^{\nu+1} = a_i^\nu + \frac{\gamma_\nu}{p} \mathscr{E}^\nu \sigma(\lambda_i^\nu)$$

$$w_i^{\nu+1} = w_i^\nu + \frac{\gamma_\nu}{p} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) x^\nu$$

<u>Goal</u>: Go from this to equation for $\Omega^\nu$

$$\mathscr{E}^\nu = (y^\nu - f(x^\nu; \Theta^\nu))$$

$$= \left( \frac{1}{k} \sum_{r=1}^{k} a_{\star,r} \sigma_\star(\lambda_{\star,r}) - \sqrt{\Delta} z^\nu - \frac{1}{p} \sum_{i=1}^{p} a_i^\nu \sigma(\lambda_i^\nu) \right)$$

# Tracking overlaps

$$w_i^{\nu+1} = w_i^\nu + \frac{\gamma_\nu}{p} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) x^\nu$$

Equation for $M^\nu = \frac{1}{d} W_\star W^{\nu\top}$:

# Tracking overlaps

$$w_i^{\nu+1} = w_i^{\nu} + \frac{\gamma_\nu}{p}\mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu)x^\nu$$

Equation for $M^\nu = \frac{1}{d}W_\star W^{\nu\top}$:

$$w_{\star,r}^\top w_i^{\nu+1} = w_{\star,r}^\top w_i^{\nu} + \frac{\gamma_\nu}{p}\mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu)w_{\star,r}^\top x^\nu$$

# Tracking overlaps

$$w_i^{\nu+1} = w_i^{\nu} + \frac{\gamma_\nu}{p}\mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu)x^\nu$$

Equation for $M^\nu = \frac{1}{d}W_\star W^{\nu\top}$:

$$\frac{1}{d}w_{\star,r}^\top w_i^{\nu+1} = \frac{1}{d}w_{\star,r}^\top w_i^\nu + \frac{\gamma_\nu}{dp}\mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu)w_{\star,r}^\top x^\nu$$

# Tracking overlaps

$$w_i^{\nu+1} = w_i^{\nu} + \frac{\gamma_\nu}{p}\mathscr{E}^{\nu}a_i\sigma'(\lambda_i^{\nu})x^{\nu}$$

Equation for $M^{\nu} = \frac{1}{d}W_{\star}W^{\nu\top}$:

$$M_{ri}^{\nu+1} = M_{ri}^{\nu} + \frac{\gamma_\nu}{dp}\mathscr{E}^{\nu}a_i\sigma'(\lambda_i^{\nu})\lambda_{\star,r}^{\nu}$$

# Tracking overlaps

$$w_i^{\nu+1} = w_i^\nu + \frac{\gamma_\nu}{p} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) x^\nu$$

Equation for $M^\nu = \frac{1}{d} W_\star W^{\nu\top}$:

$$M_{ri}^{\nu+1} = M_{ri}^\nu + \frac{\gamma_\nu}{dp} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) \lambda_{\star,r}^\nu$$

Equation for $Q^\nu = \frac{1}{d} W^\nu W^{\nu\top}$:

$${\color{red}w_j^{\nu+1\top}} w_i^{\nu+1} = {\color{red}w_j^{\nu+1\top}} \left( w_i^\nu + \frac{\gamma_\nu}{p} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) x^\nu \right)$$

# Tracking overlaps

$$w_i^{\nu+1} = w_i^\nu + \frac{\gamma_\nu}{p} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) x^\nu$$

Equation for $M^\nu = \frac{1}{d} W_\star W^{\nu\top}$:

$$M_{ri}^{\nu+1} = M_{ri}^\nu + \frac{\gamma_\nu}{dp} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) \lambda_{\star,r}^\nu$$

Equation for $Q^\nu = \frac{1}{d} W^\nu W^{\nu\top}$:

$$w_j^{\nu+1\top} w_i^{\nu+1} = w_j^{\nu+1\top} \left( w_i^\nu + \frac{\gamma_\nu}{p} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) x^\nu \right)$$

$$= \left( w_i^\nu + \frac{\gamma_\nu}{p} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) x^\nu \right) \left( w_i^\nu + \frac{\gamma_\nu}{p} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) x^\nu \right)$$

# Tracking overlaps

$$w_i^{\nu+1} = w_i^\nu + \frac{\gamma_\nu}{p} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) x^\nu$$

Equation for $M^\nu = \frac{1}{d} W_\star W^{\nu\top}$:

$$M_{ri}^{\nu+1} = M_{ri}^\nu + \frac{\gamma_\nu}{dp} \mathscr{E}^\nu a_i \sigma'(\lambda_i^\nu) \lambda_{\star,r}^\nu$$

Equation for $Q^\nu = \frac{1}{d} W^\nu W^{\nu\top}$:

$$Q_{ji}^{\nu+1} = Q_{ji}^\nu + \frac{\gamma_\nu}{dp} \left( \mathscr{E}^\nu \sigma'(\lambda_j^\nu) \lambda_i^\nu + \mathscr{E}^\nu \sigma'(\lambda_i^\nu) \lambda_j^\nu \right)$$

$$+ \frac{\gamma_\nu^2}{dp^2} (\mathscr{E}^\nu)^2 \sigma'(\lambda_i^\nu) \sigma'(\lambda_j^\nu) ||x^\nu||_2^2$$

# Tracking overlaps: summary

$$w_i^{\nu+1} = w_i^\nu + \frac{\gamma_\nu}{p} \mathcal{E}^\nu a_i \sigma'(\lambda_i^\nu) x^\nu$$

Stochastic process
in $\mathbb{R}^{p \times d}$

$$M_{ri}^{\nu+1} - M_{ri}^\nu = \frac{\gamma_\nu}{dp} \Psi_M^{(\mathrm{GF})}(M, Q)$$

$$Q_{ji}^{\nu+1} - Q_{ji}^\nu = \frac{\gamma_\nu}{dp} \Psi_Q^{(\mathrm{GF})}(M, Q) + \frac{\gamma_\nu^2}{dp^2} \Psi_Q^{(\mathrm{var})}(M, Q)$$

Stochastic process
in $\mathbb{R}^{p(p+k)}$

# Concentration result

Define step-size $\delta t = \dfrac{\gamma}{dp}$ and $M(t), Q(t)$ such that:

$$M(\nu \delta t) = M^{\nu} \qquad\qquad Q(\nu \delta t) = Q^{\nu}$$

# Concentration result

Define step-size $\delta t = \dfrac{\gamma}{dp}$ and $M(t), Q(t)$ such that:

$$M(\nu \delta t) = M^\nu \qquad\qquad Q(\nu \delta t) = Q^\nu$$

Theorem [Veiga, Stephan, Loureiro, Krzakala, Zdeborová '22]

Then $\forall 0 \leq \nu \leq \left\lfloor \dfrac{n}{\delta t} \right\rfloor$ :

$$\mathbb{E} \left| \left| \Omega^\nu - \bar{\Omega}(\nu \delta t) \right| \right|_\infty \leq e^{C\nu\delta t} \sqrt{\frac{\gamma}{dp}}$$

Where $\bar{\Omega}(t) = \mathbb{E}[\Omega(t)]$ is the solution of an ODE:

$$\frac{d\bar{\Omega}(t)}{dt} = \mathbb{E}\left[ \psi(\bar{\Omega}(t)) \right]$$

# The different limiting regimes



Mean field limit
$$p \to \infty$$
$$\gamma, d = O(1)$$

Classical limit
$$\gamma \to 0^+$$
$$d, p = O(1)$$

$$\frac{\gamma}{dp} \to 0^+$$

High-d limit
$$d \to \infty$$
$$\gamma, p = O(1)$$

# The different limiting regimes



Mean field limit
$$p \to \infty$$
$$\gamma, d = O(1)$$

Classical limit
$$\gamma \to 0^+$$
$$d, p = O(1)$$

$$\frac{\gamma}{dp} \to 0^+$$

High-d limit
$$d \to \infty$$
$$\gamma, p = O(1)$$

# Classical regime



Classical limit
$\gamma \to 0^+$
$d, p = O(1)$

$$\Theta^{\nu+1} = \Theta^\nu - \frac{\gamma_\nu}{2} \nabla_{\Theta^\nu} (y^\nu - f(x^\nu; \Theta^\nu))^2$$
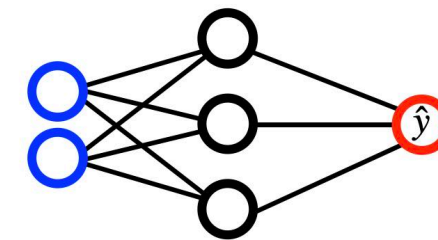
$\downarrow \quad \gamma \to 0^+$

[Robins & Monro '51]

$$\dot{\Theta}(t) = -\frac{1}{2} \nabla_\Theta \mathbb{E}_{(x,y)} \left[ (y - f(x; \Theta(t)))^2 \right]$$

# Classical regime



Classical limit
$\gamma \to 0^+$
$d, p = O(1)$

$$\Theta^{\nu+1} = \Theta^{\nu} - \frac{\gamma_{\nu}}{2} \nabla_{\Theta^{\nu}} (y^{\nu} - f(x^{\nu}; \Theta^{\nu}))^2$$

$\gamma \to 0^+$

[Robins & Monro '51]

$$\dot{\Theta}(t) = -\frac{1}{2} \nabla_{\Theta} \mathbb{E}_{(x,y)} \left[ (y - f(x; \Theta(t)))^2 \right]$$

Subleading in $\gamma$

$$\Theta^{k+1} = \Theta^k - \gamma_k \nabla_{\Theta^k} \mathscr{R} \left( \Theta^k \right) + \gamma_k \varepsilon^k$$

GD on population

Effective
Noise

# Classical regime

$$\Theta^{\nu+1} = \Theta^{\nu} - \frac{\gamma_{\nu}}{2} \nabla_{\Theta^{\nu}} (y^{\nu} - f(x^{\nu}; \Theta^{\nu}))^2$$

$\downarrow$ $\gamma \to 0^+$

[Robins & Monro '51]

$$\dot{\Theta}(t) = -\frac{1}{2} \nabla_{\Theta} \mathbb{E}_{(x,y)} \left[ (y - f(x; \Theta(t)))^2 \right]$$

$$\dot{M}_{ri}(t) = \mathbb{E}[\Psi_M^{(\mathrm{GF})}(\bar{M}, \bar{Q})]$$

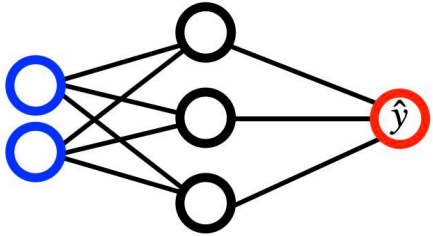$$\dot{\bar{Q}}_{ji}(t) = \mathbb{E}[\Psi_Q^{(\mathrm{GF})}(\bar{M}, \bar{Q})]$$

Classical limit
$\gamma \to 0^+$
$d, p = O(1)$

# Classical regime



Classical limit
$$\gamma \to 0^+$$
$$d, p = O(1)$$

# Classical regime



Classical limit
$$\gamma \to 0^+$$
$$d, p = O(1)$$

$t = 0.00$

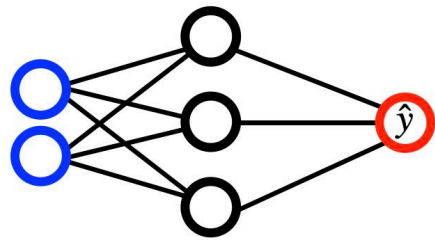$M_{j2}/\sqrt{Q_{jj}P_{22}}$

$\mathcal{R} - \frac{\Delta}{2}$

$t$

$M_{j1}/\sqrt{Q_{jj}P_{11}}$

# The different limiting regimes



Mean field limit
$$p \to \infty$$
$$\gamma, d = O(1)$$

$$\frac{\gamma}{dp} \to 0^+$$
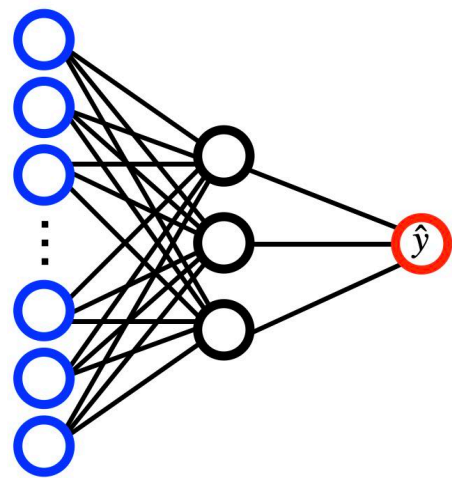
Classical limit
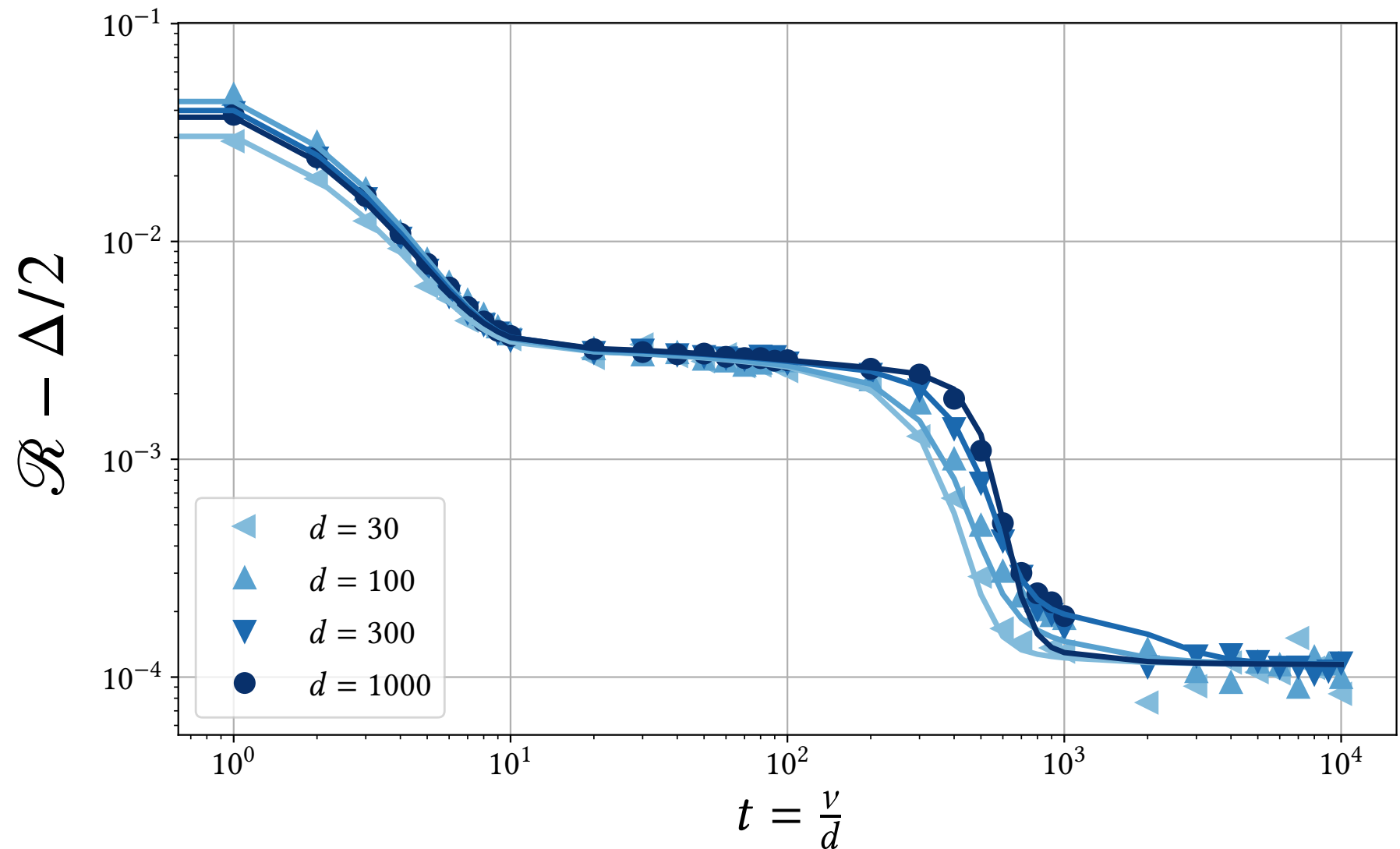$$\gamma \to 0^+$$
$$d, p = O(1)$$

High-d limit
$$d \to \infty$$
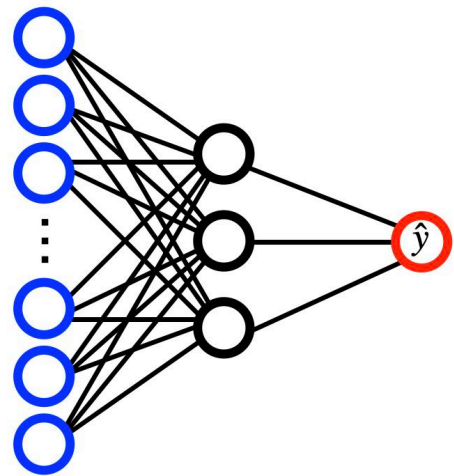$$\gamma, p = O(1)$$

# High-dimensional regime [Saad & Solla '95]



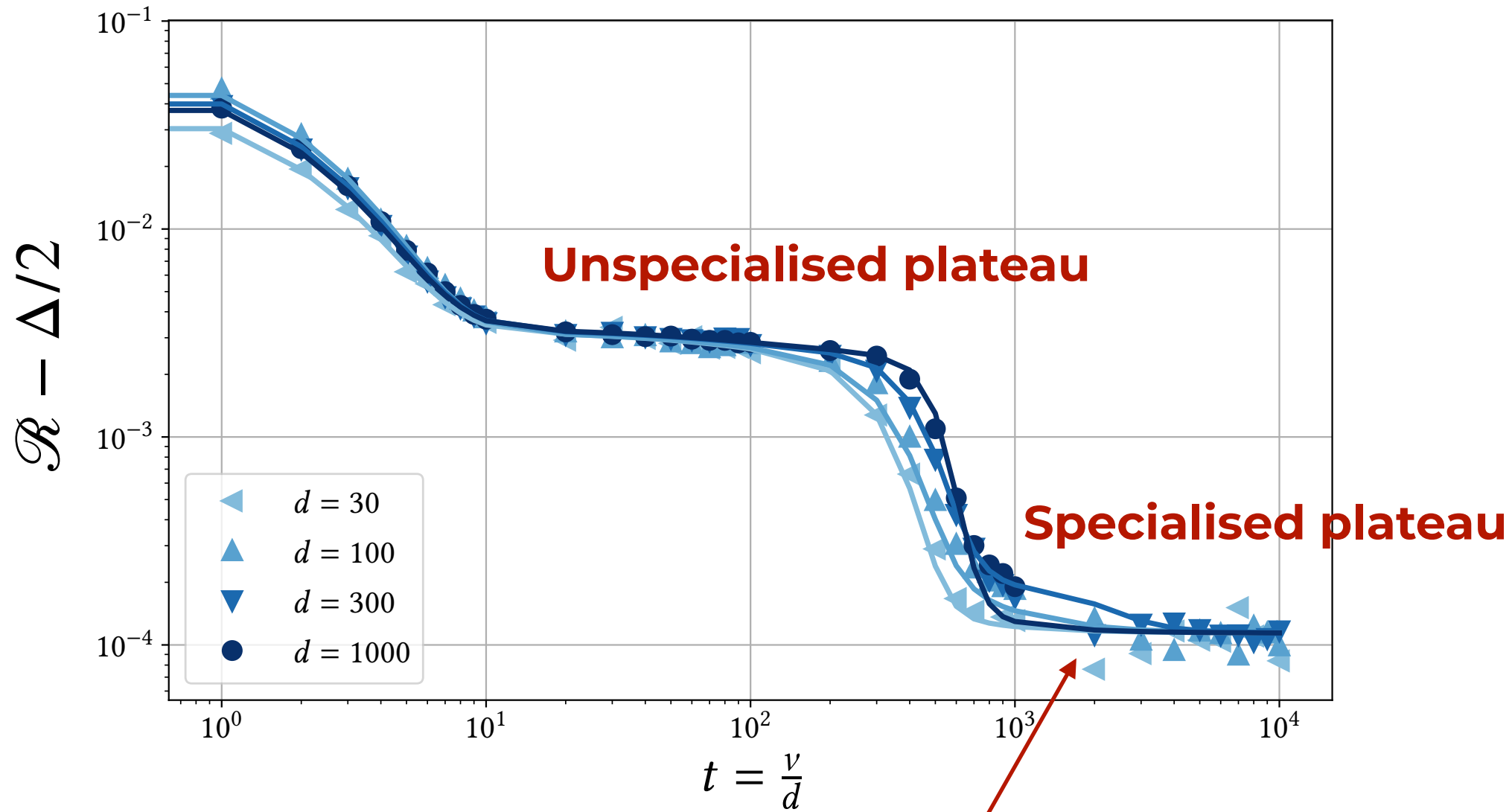High-d limit

$$d \rightarrow \infty$$

$$\gamma, p = O(1)$$

legend:
- $d = 30$
- $d = 100$
- $d = 300$
- $d = 1000$

y-axis: $\mathscr{R} - \Delta/2$

x-axis: $t = \frac{\nu}{d}$

$$\dot{\bar{M}}_{ri}(t) = \mathbb{E}[\Psi_M^{(\mathrm{GF})}(\bar{M}, \bar{Q})]$$

$$\dot{\bar{Q}}_{ji}(t) = \mathbb{E}[\Psi_Q^{(\mathrm{GF})}(\bar{M}, \bar{Q})] + \frac{\gamma}{p}\mathbb{E}\left[\Psi_Q^{\mathrm{var}}(\bar{M}, \bar{Q})\right]$$
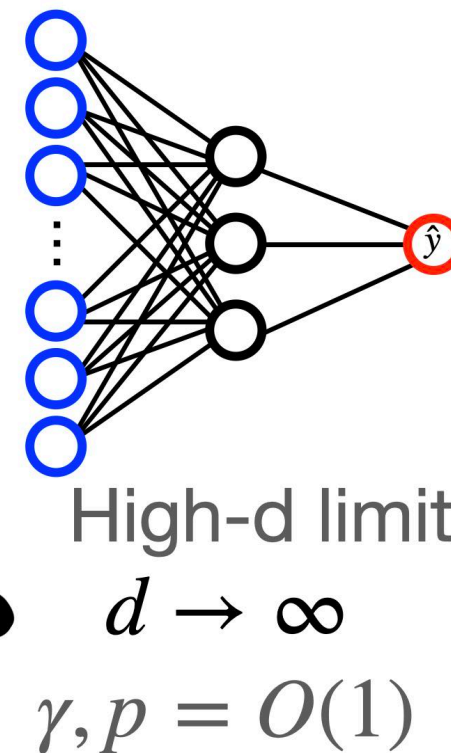
# High-dimensional regime [Saad & Solla '95]
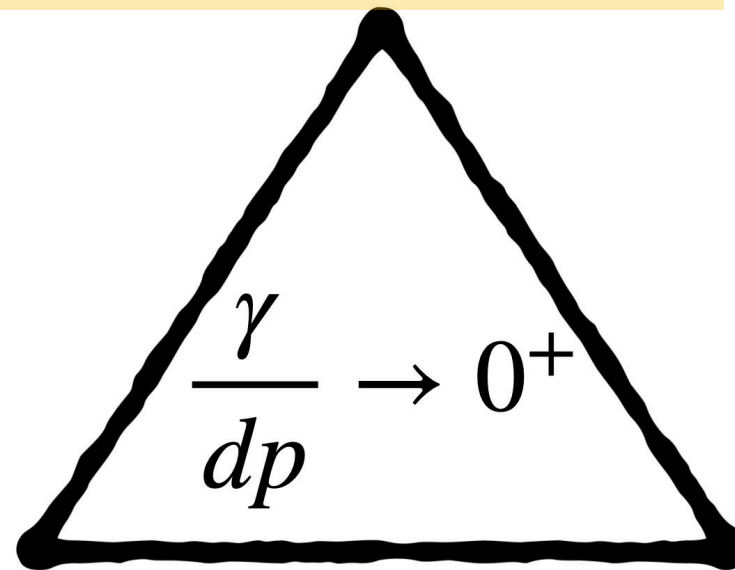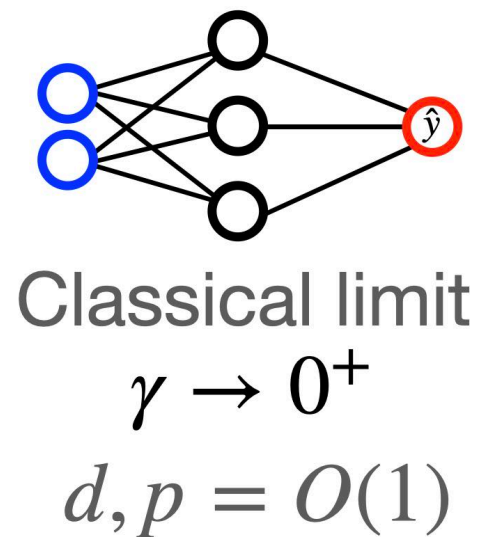


High-d limit

$$d \to \infty$$

$$\gamma, p = O(1)$$

Unspecialised plateau

Specialised plateau

$d = 30$
$d = 100$
$d = 300$
$d = 1000$

$$\mathscr{R} - \Delta/2$$

$$t = \frac{v}{d}$$

$$\mathscr{R}_\infty - \Delta/2 \propto \gamma\Delta$$

# The different limiting regimes



Mean field limit
$$p \to \infty$$
$$\gamma, d = O(1)$$

$$\frac{\gamma}{dp} \to 0^+$$

Classical limit
$$\gamma \to 0^+$$
$$d, p = O(1)$$
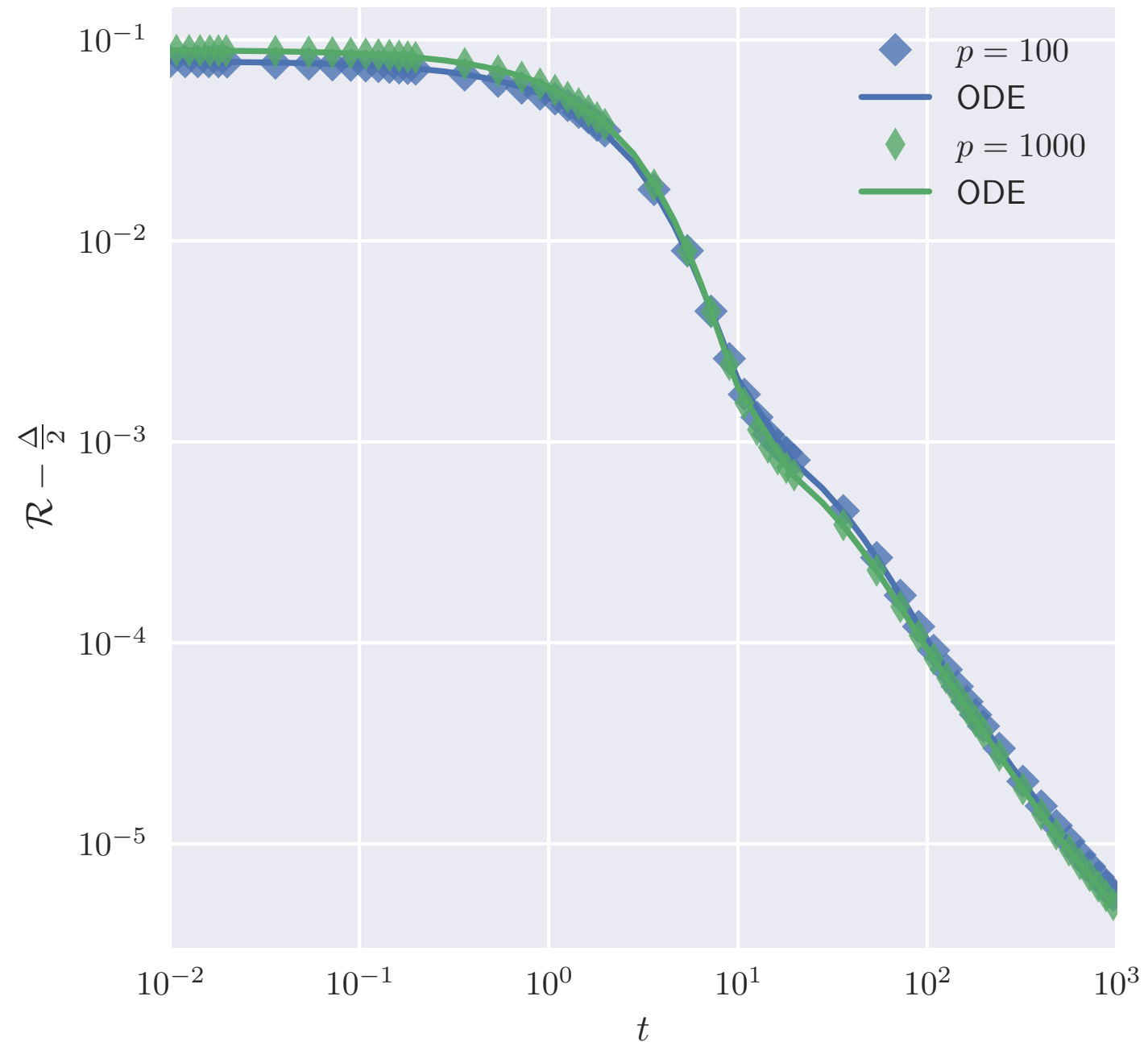
High-d limit
$$d \to \infty$$
$$\gamma, p = O(1)$$

# Mean-field regime



Mean field limit
$$p \to \infty$$
$$\gamma, d = O(1)$$

$$\dot{M}_{ri}(t) = \mathbb{E}[\Psi_M^{\mathrm{(GF)}}(\bar{M}, \bar{Q})] \qquad \dot{\bar{Q}}_{ji}(t) = \mathbb{E}[\Psi_Q^{\mathrm{(GF)}}(\bar{M}, \bar{Q})]$$

# Mean-field limit

## On the Global Convergence of Gradient Descent for Over-parameterized Models using Optimal Transport

**Lénaïc Chizat**
INRIA, ENS, PSL Research University
Paris, France
lenaic.chizat@inria.fr

**Francis Bach**
INRIA, ENS, PSL Research University
Paris, France
francis.bach@inria.fr

## TRAINABILITY AND ACCURACY OF NEURAL NETWORKS: AN INTERACTING PARTICLE SYSTEM APPROACH

GRANT M. ROTSKOFF AND ERIC VANDEN-EIJNDEN

## Mean field analysis of neural networks: A central limit theorem

Justin Sirignano[a,*], Konstantinos Spiliopoulos[b,1]

## A mean field view of the landscape of two-layer neural networks

Song Mei[a], Andrea Montanari[b,c,1], and Phan-Minh Nguyen[b]

# Mean-field limit

Idea: Define empirical density of weights:

$$\rho_p^\nu(a, w) = \frac{1}{p} \sum_{i=1}^{p} \delta(a_i - a_i^\nu)\delta(w_i - w_i^\nu)$$

# Mean-field limit

💡 <u>Idea:</u>  Define empirical density of weights:

$$\rho_p^\nu(a, w) = \frac{1}{p} \sum_{i=1}^{p} \delta(a_i - a_i^\nu)\delta(w_i - w_i^\nu)$$

The risk is linear in $\hat{\rho}_p$!

$$\mathscr{R}(\Theta) = \mathbb{E}\left(y - \int \hat{\rho}_p(\mathrm{d}a, \mathrm{d}w)a\sigma(w \cdot x)\right)^2$$

# Mean-field limit

💡 Idea: Define empirical density of weights:

$$\rho_p^\nu(a, w) = \frac{1}{p} \sum_{i=1}^{p} \delta(a_i - a_i^\nu)\delta(w_i - w_i^\nu)$$

The risk is linear in $\hat\rho_p$!

$$\mathscr{R}(\Theta) = \mathbb{E}\left(y - \int \hat\rho_p(\mathrm{d}a, \mathrm{d}w)a\sigma(w \cdot x)\right)^2$$

Show that, at fixed $d$ and $\gamma_k \ll 1/d$:

$$\text{One-pass SGD} \quad \xrightarrow{\ p \to \infty\ } \quad \boxed{\partial_t \rho_t = \gamma \nabla_\theta\big(\rho_t \nabla_\theta \ell(\theta; \rho_t)\big)}$$

"Mean-field" limit

[Mei, Montanari, Nguyen 18'; Chizat, Bach 18'; Rotskoff, Vanden-Eijnden 18'; Sirignano, Spiliopoulos 18']

# Global convergence

**Theorem 2 (Informal)** *If the support of the initial distribution includes all directions in $\mathbb{R}^{d+1}$, and if the function $\Psi$ is positively 2-homogeneous then if the Wasserstein gradient flow weakly converges to a distribution, it can only be to a global optimum of $F$.*

**From qualitative to quantitative results?** Our result states that for infinitely many particles, we can only converge to a global optimum (note that we cannot show that the flow always converges). However, it is only a qualitative result in comparison with what is known for convex optimization problems in Section 2.2:
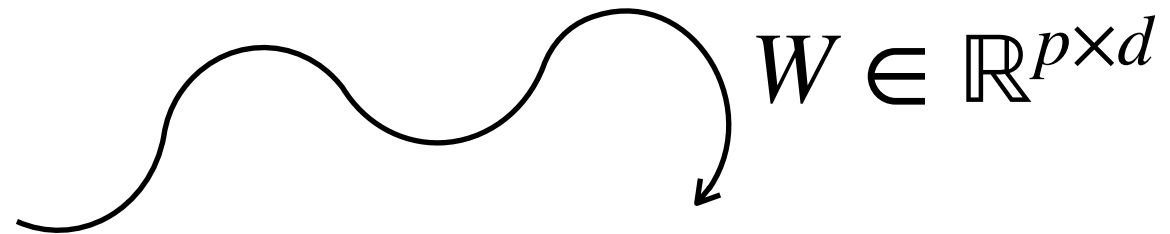
- This is only for $m = +\infty$, and we cannot provide an estimation of the number of particles needed to approximate the mean field regime that is not exponential in $t$ (see such results e.g. in [28]).

- We cannot provide an estimation of the performance as the function of time, that would provide an upper bound on the running time complexity.

[Mei, Montanari, Nguyen 18'; Chizat, Bach 18'; Rotskoff, Vanden-Eijnden 18'; Sirignano, Spiliopoulos 18']
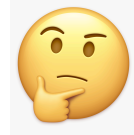
# Mean-field regime

🤔 But $Q \in \mathbb{R}^{p \times p}$ !!!
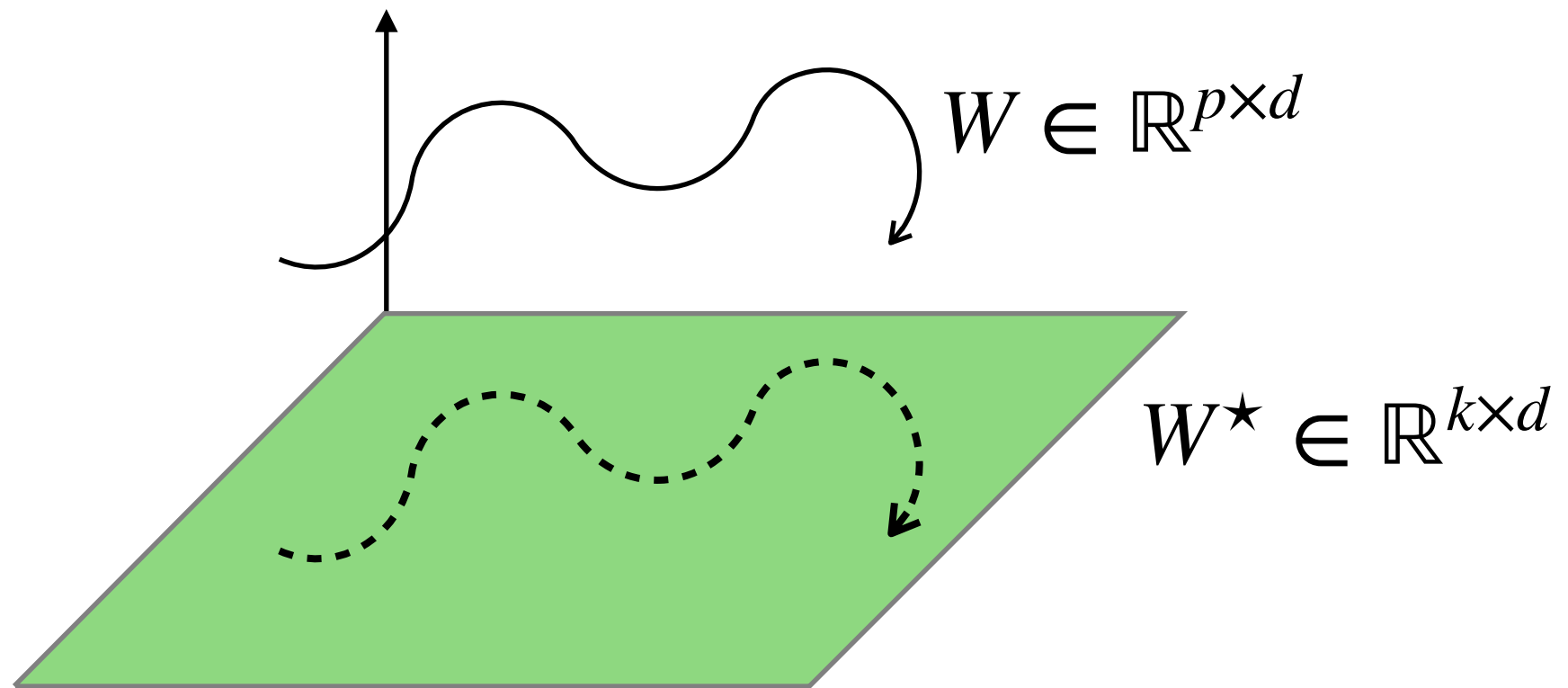
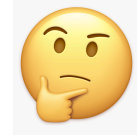$W \in \mathbb{R}^{p \times d}$

# Mean-field regime

🤔 But $Q \in \mathbb{R}^{p \times p}$ !!!

$$W = \boxed{MP^{-1}W^{\star}} + \boxed{W^{\perp}}$$

Teacher subspace



$W \in \mathbb{R}^{p \times d}$

$W^{\star} \in \mathbb{R}^{k \times d}$

# Mean-field regime

🤔 But $Q \in \mathbb{R}^{p \times p}$ !!!

$$W = \boxed{MP^{-1}W^{\star}} + \boxed{W^{\perp}} \longrightarrow Q \approx \boxed{MPM^{\top}} + \boxed{D_{\sqrt{q^{\perp}}} \Xi D_{\sqrt{q^{\perp}}}}$$

Teacher subspace

$\wr$

$\mathbb{S}^{d-k-1}$



$W \in \mathbb{R}^{p \times d}$

$W_{\star} \in \mathbb{R}^{k \times d}$

# Mean-field + high-d

Theorem [Arnaboldi, Stephan, Loureiro, Krzakala '23]

$$\mathbb{E}||Q(t) - MP^{-1}M^{\top} + \mathrm{diag}(Q^{\perp})||_{\infty} \leq e^{Ct}\left(p^{-1/2} + d^{-1/2}\right)$$

Suppressed in $1/\sqrt{d}$

$W \in \mathbb{R}^{p \times d}$

$W_{\star} \in \mathbb{R}^{k \times d}$

# Mean-field + high-d

$$\mathbb{E}||Q(t) - MP^{-1}M^{\top} + \mathrm{diag}(Q^{\perp})||_{\infty} \leq e^{Ct}\left(p^{-1/2} + d^{-1/2}\right)$$

This implies MF-like PDE for the sufficient statistics:

$$\hat{\mu}_t(m, q) = \frac{1}{p}\sum_{i=1}^{p}\delta(m - m_i(t))\delta(q - Q_{ii}^{\perp}(t))$$

$$\partial_t\hat{\mu}_p(m, q) = \nabla_{(m,q)} \cdot (\hat{\mu}_t\varphi(\,\cdot\,, \hat{\mu}_t))$$
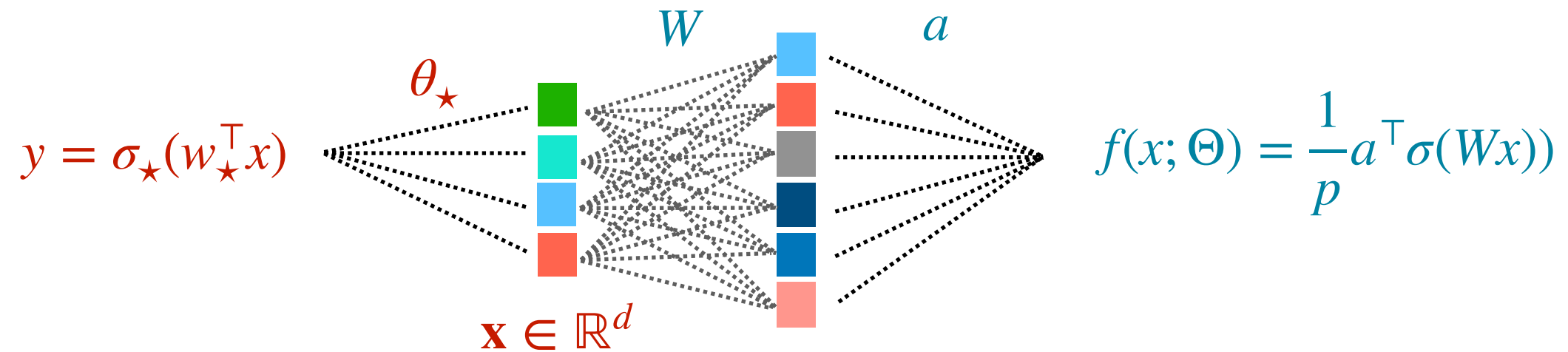
# Mean-field and high-d

# What can I do with that?

Consider simple case: $k = 1$ and $p \to \infty$



$\theta_\star$

$y = \sigma_\star(w_\star^\top x)$

$\mathbf{x} \in \mathbb{R}^d$

$W$

$a$

$f(x; \Theta) = \dfrac{1}{p} a^\top \sigma(Wx))$

# What can I do with that?

Consider simple case: $k = 1$ and $p \to \infty$



$\theta_\star$    $W$    $a$

$y = \sigma_\star(w_\star^\top x)$

$\mathbf{x} \in \mathbb{R}^d$

$f(x; \Theta) = \dfrac{1}{p} a^\top \sigma(Wx))$



$\mathscr{R}$

$O(\varepsilon^{1/2})$    $O(\varepsilon^{1/3})$    $O(\varepsilon^{1/4})$

$\mathscr{R}_{\text{init}}$

$\frac{1}{2}\varphi_1^2$

$\frac{1}{2}\varphi_2^2$

$\frac{1}{2}\varphi_3^2$

$0$   $O(\varepsilon)$   $\frac{1}{4|\sigma_1\varphi_1|}\varepsilon^{1/2}\log\frac{1}{\varepsilon}$   $c_2\varepsilon^{1/4}$   $c_3\varepsilon^{1/6}$   $t$

<u>Recall:</u> For $n \propto d$ and $x \sim \mathcal{N}(0, I_d)$

<u>RF / Kernels:</u> can only learn linear part

<u>NNs:</u> can learn non-linear components
if $(\sigma_\star, \sigma)$ "standard"
(full Hermite expansion).

# Follow-ups and challenges

➕ Characterisation of the stochastic dynamics close to fixed points as a coloured diffusion

[Ben Arous, Gheissari, Jagannath NeurIPS '22]

# Follow-ups and challenges

+ Characterisation of the stochastic dynamics
  close to fixed points as a coloured diffusion

[Ben Arous, Gheissari,
Jagannath NeurIPS '22]

+ Uniform control over the variance
  and convergence rates

# Follow-ups and challenges

➕ Characterisation of the stochastic dynamics close to fixed points as a coloured diffusion

[Ben Arous, Gheissari, Jagannath NeurIPS '22]

➕ Uniform control over the variance and convergence rates

➕ Phenomenology of the dynamics:
Functions of increasing complexity?
Distributions of increasing complexity?

[Abbe, Adsera, Misiakiewicz, COLT '22; Berthier, Montanari '23]

[Refinetti, Ingrosso, Goldt, arXiv '22]

# Follow-ups and challenges

➕ Characterisation of the stochastic dynamics
close to fixed points as a coloured diffusion

[Ben Arous, Gheissari,
Jagannath NeurIPS '22]

➕ Uniform control over the variance
and convergence rates

➕ Phenomenology of the dynamics:
Functions of increasing complexity?
Distributions of increasing complexity?

[Abbe, Adsera, Misiakiewicz,
COLT '22; Berthier, Montanari '23]

[Refinetti, Ingrosso, Goldt,
arXiv '22]

➕ Role of initialisation

# Lecture III: Summary

✓ Deterministic analysis of two-layer neural nets in the "rich" regime

# Lecture III: Summary

✓   Deterministic analysis of two-layer neural nets in the "rich" regime

✓   Phenomenology in the classical and high-dimensional regime

# Lecture III: Summary

✓ Deterministic analysis of two-layer neural nets
in the "rich" regime

✓ Phenomenology in the classical
and high-dimensional regime

✓ Interplay between effective SGD noise
and overparametrisation

# Lecture III: Summary

✓ Deterministic analysis of two-layer neural nets
in the "rich" regime

✓ Phenomenology in the classical
and high-dimensional regime

✓ Interplay between effective SGD noise
and overparametrisation

✓ Dimension free limits in the mean-field regime

# But this is only the tip of an iceberg…



[brloureiro@gmail.com](mailto:brloureiro@gmail.com)