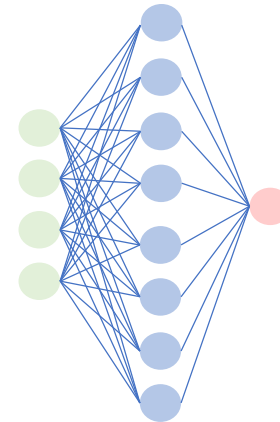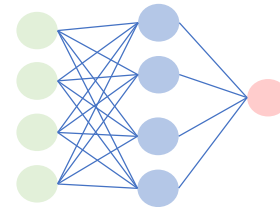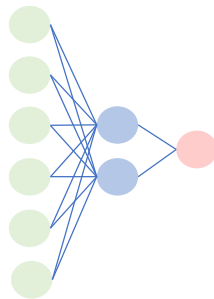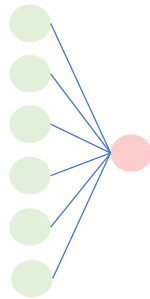# Learning (with) deep random networks

Hugo Cui

*SPOC lab, EPFL, Switzerland*

Cargèse 2023

*Learning curves of generic features maps for realistic datasets with a teacher-student model*, Loureiro, Gerbelot, **HC**, Goldt, Mézard, Krzakala, Zdeborová, NeurIPS 2021

*Bayes-optimal learning of deep random networks of extensive width*, **HC**, Krzakala, Zdeborová, ICML 2023

*Deterministic equivalent and gaussian universality of deep random features learning*, Schröder, **HC**, Dmitriev, Loureiro, ICML 2023

**Bruno Loureiro**
ENS

**Cédric Gerbelot**
NYU

**Sebastian Goldt**
SISSA

**Dominik Schröder**
ETH

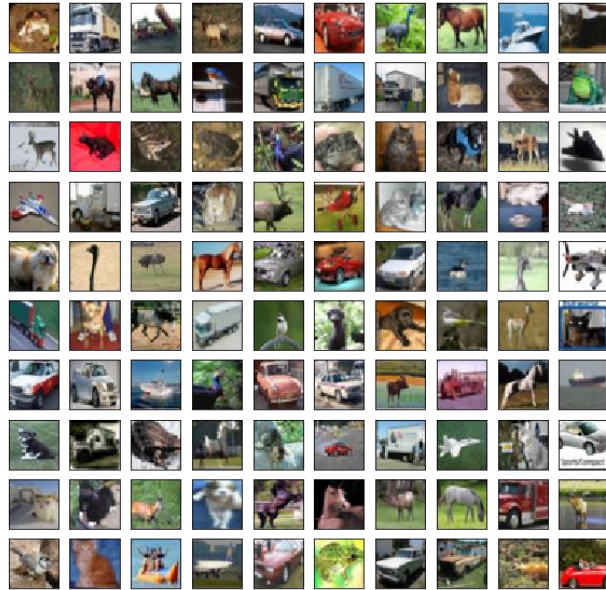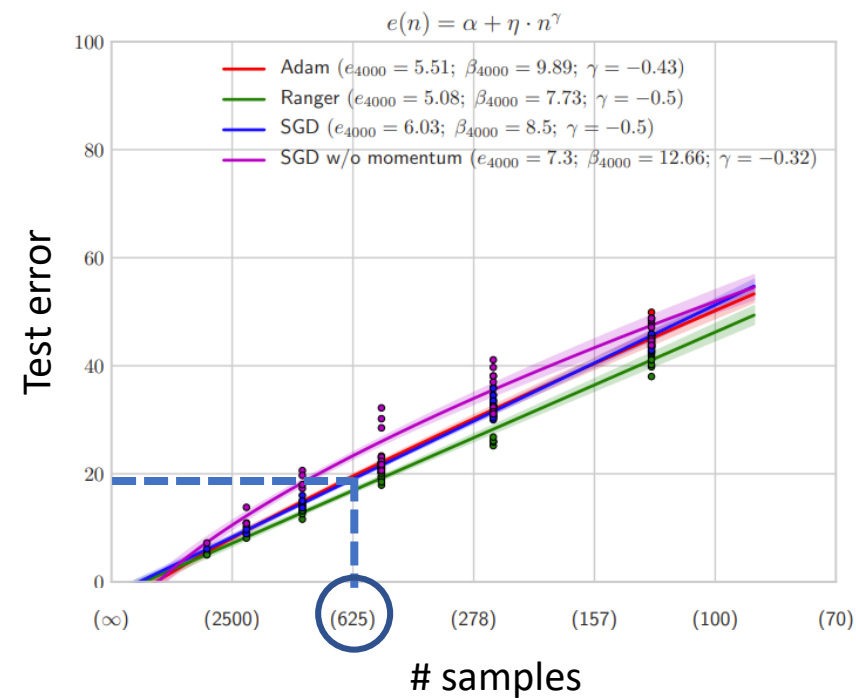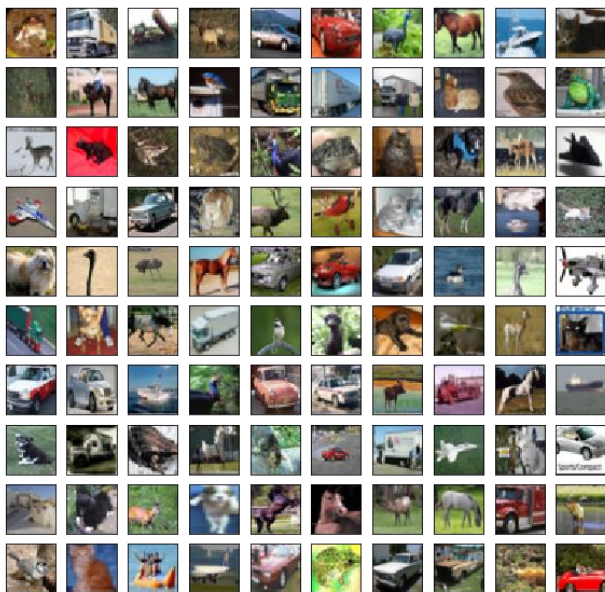**Florent Krzakala**
EPFL

**Lenka Zdeborová**
EPFL

**Marc Mézard**
Bocconi

**Daniil Dmitriev**
ETH

**Question**: **What is the best accuracy** one can achieve from 600 training samples?

$$e(n) = \alpha + \eta \cdot n^{\gamma}$$

Adam ($e_{4000} = 5.51$; $\beta_{4000} = 9.89$; $\gamma = -0.43$)
Ranger ($e_{4000} = 5.08$; $\beta_{4000} = 7.73$; $\gamma = -0.5$)
SGD ($e_{4000} = 6.03$; $\beta_{4000} = 8.5$; $\gamma = -0.5$)
SGD w/o momentum ($e_{4000} = 7.3$; $\beta_{4000} = 12.66$; $\gamma = -0.32$)

Test error

# samples

**Question**: **What is the best accuracy** one can achieve from 600 training samples?

**(Empirical) Answer**: Probably ≈ 82%, using good networks.

Hoeim et al., *Learning curves for Analysis of Deep Networks,* ICML 2020

For a train set $\mathcal{D} = \{x^\mu, y^\star(x^\mu)\}_{\mu=1}^n$ of given size $n$, what is ***the lowest achievable test error*** $\epsilon_g$ one can hope to achieve?

For a train set $\mathcal{D} = \{x^{\mu}, y^{\star}(x^{\mu})\}_{\mu=1}^{n}$ of given size $n$, what is *the lowest achievable test error* $\epsilon_g$ one can hope to achieve?
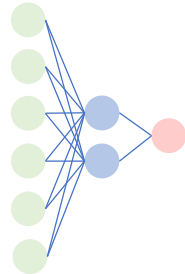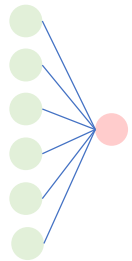
When the target function is *parametric*, the lowest (**Bayes-optimal**) test error is given by Bayesian inference.

# Theoretical testbeds: random neural networks



Barbier et al, *Optimal errors and phase transitions in high-dimensional generalized linear models,* PNAS 2017

8

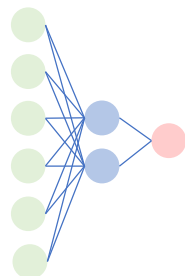# Theoretical testbeds: random neural networks

$$width \ll dimension$$

Barbier et al, *Optimal errors and phase transitions in high-dimensional generalized linear models,* PNAS 2017

Aubin et al, *The committee machine: Computational to statistical gaps,* NeurIPS 2019

# Theoretical testbeds: random neural networks

$width \ll dimension$

$width \gg dimension$

Barbier et al, *Optimal errors and phase transitions in high-dimensional generalized linear models,* PNAS 2017

Aubin et al, *The committee machine: Computational to statistical gaps,* NeurIPS 2019

Neal, *Priors for infinite nets*, Uni. Toronto 1996
Williams, *Computing with infinite networks, NeurIPS 1996*
Lee et. al., *Deep Neural Networks as GPs,* ICLR 2018

# Theoretical testbeds: random neural networks



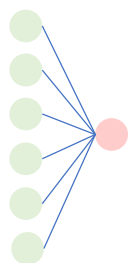*width* $\ll$ *dimension*

*width* $\sim$ *dimension*

*width* $\gg$ *dimension*

Barbier et al, *Optimal errors and phase transitions in high-dimensional generalized linear models,* PNAS 2017

Aubin et al, *The committee machine: Computational to statistical gaps,* NeurIPS 2019
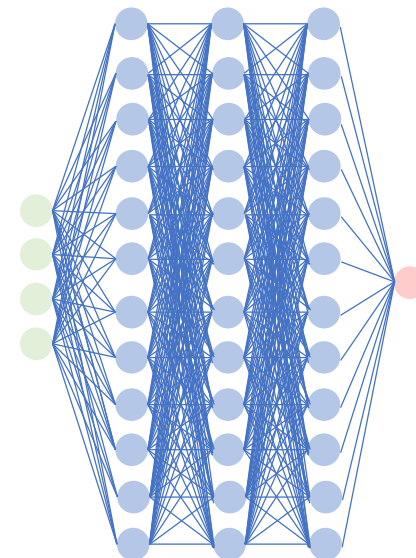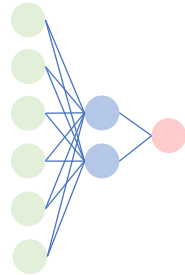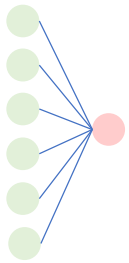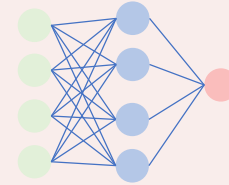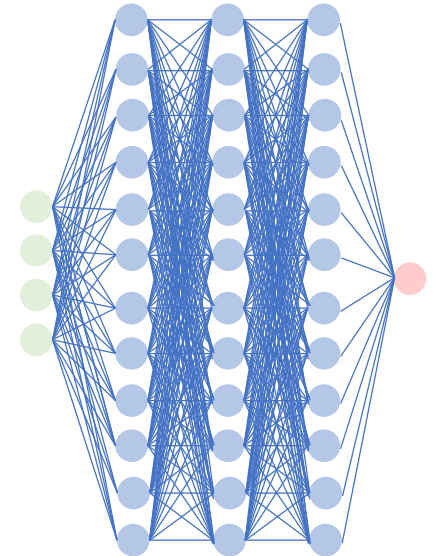
Neal, *Priors for infinite nets*, Uni. Toronto 1996
Williams, *Computing with infinite networks, NeurIPS 1996*
Lee et. al., *Deep Neural Networks as GPs,* ICLR 2018

_Some related works:_

**High-dimensional formulae for sign/ReLU Bayes regression**

Li and Sompolinsky, _Statistical mechanics of deep linear neural networks: The backpropagating kernel renormalization,_ PRX, 2021.

Ariosto et al., _Statistical mechanics of deep learning beyond the infinite-width limit._ ArXiv, abs/2209.04882, 2022.

**_(Non)-asymptotics for linear networks_**

Zavatone-Veth, Tong and Pehlevan, _Contrasting random and learned features in deep bayesian linear regression_, PRE 2022

Hanin and Zlokapa, _Bayesian interpolation with deep linear networks._ ArXiv, abs/2212.14457, 2022

**_Recent advances in neighbouring regimes_**

Camilli, Tieplova, Barbier, _Fundamental limits of overparametrized shallow networks for supervised learning,_ ArXiv, abs/2307.05635

$width \sim dimension$

*(Data)*    Gaussian data: $x \sim \mathcal{N}(0, \Sigma)$

*(Data)*   Gaussian data: $x \sim \mathcal{N}(0, \Sigma)$

*(Target)*   $y^\star(x) = f^\star \left( \dfrac{a^\top}{\sqrt{k_L}} \varphi_L \circ \cdots \circ \varphi_1(x) + \sqrt{\Delta}\xi \right)$

with layers   $\varphi_\ell(h) = \sigma_\ell \left( \dfrac{W_\ell}{\sqrt{k_{\ell-1}}} h \right)$

Odd activations $\sigma_\ell$

$(W_\ell)_{ij} \sim \mathcal{N}(0, \Delta_\ell), \quad a_i \sim \mathcal{N}(0, \Delta_a), \quad \xi \sim \mathcal{N}(0,1)$

(Data)    Gaussian data: $x \sim \mathcal{N}(0, \Sigma)$

(Target)

$$y^\star(x) = f^\star\left(\frac{a^\top}{\sqrt{k_L}} \varphi_L \circ \cdots \circ \varphi_1(x) + \sqrt{\Delta}\xi\right)$$

with layers    $\varphi_\ell(h) = \sigma_\ell\left(\frac{W_\ell}{\sqrt{k_{\ell-1}}} h\right)$

Odd activations $\sigma_\ell$

$(W_\ell)_{ij} \sim \mathcal{N}(0, \Delta_\ell), \ a_i \sim \mathcal{N}(0, \Delta_a), \ \xi \sim \mathcal{N}(0,1)$

(Train set)    Supervised learning with $n$ i.i.d samples $\mathcal{D} = \{x^\mu, y^\star(x^\mu)\}_{\mu=1}^n$
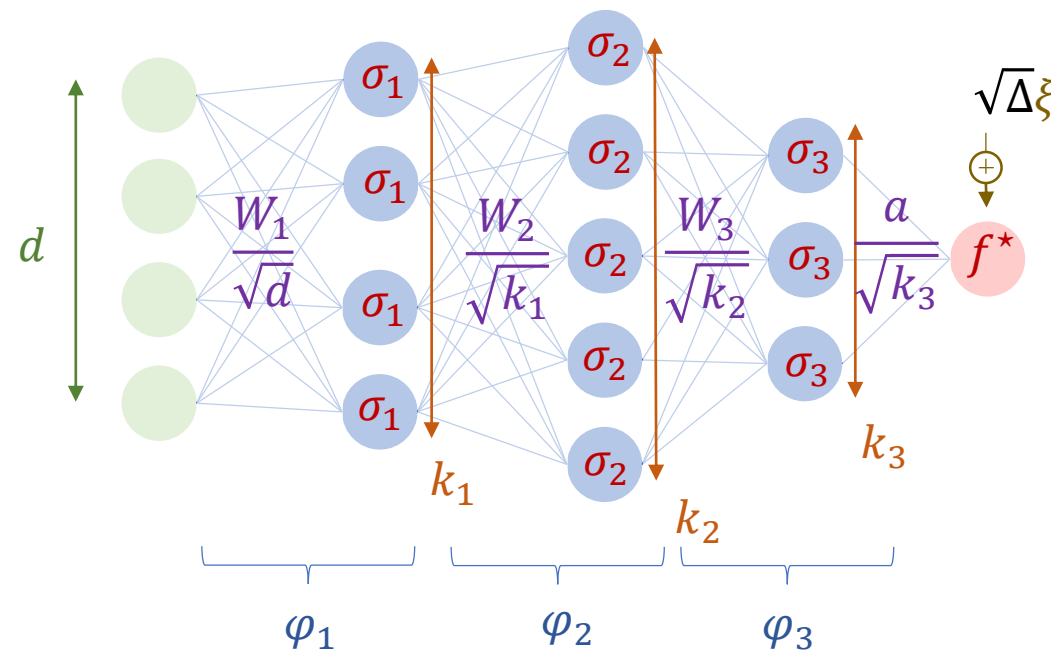
(Data)   Gaussian data: $x \sim \mathcal{N}(0, \Sigma)$

(Target)
$$y^\star(x) = f^\star\left(\frac{a^\top}{\sqrt{k_L}} \varphi_L \circ \cdots \circ \varphi_1(x) + \sqrt{\Delta}\xi\right)$$

with layers
$$\varphi_\ell(h) = \sigma_\ell\left(\frac{W_\ell}{\sqrt{k_{\ell-1}}} h\right)$$

Odd activations $\sigma_\ell$

$(W_\ell)_{ij} \sim \mathcal{N}(0, \Delta_\ell), \ a_i \sim \mathcal{N}(0, \Delta_a), \ \xi \sim \mathcal{N}(0,1)$

(Train set)   Supervised learning with $n$ i.i.d samples $\mathcal{D} = \{x^\mu, y^\star(x^\mu)\}_{\mu=1}^n$

**Proportional extensive-width limit**

$$n, d, k_1, \ldots, k_L \longrightarrow \infty \qquad \text{with} \qquad \alpha = \frac{n}{d}, \gamma_\ell = \frac{k_\ell}{d} = \mathcal{O}(1)$$

16

Suppose the architecture, priors, activations are known.
The best test error is then given by *Bayesian inference*:

*Bayes posterior*

$$\mathbb{P}\left(a, \{W_\ell\}_{\ell=1}^L \,\middle|\, \mathcal{D}\right) \propto e^{-\frac{||a||^2}{2\Delta_a} - \sum_{\ell=1}^L \frac{||W_\ell||_F^2}{2\Delta_\ell}} \times \prod_{\ell=1}^L \int \frac{d\xi\, e^{-\frac{\xi^2}{2}}}{\sqrt{2\pi}}\, \delta\left(y^\star(x^\mu) - f^\star\left(\frac{a^\top}{\sqrt{k_L}}\, \varphi_L \circ \cdots \circ \varphi_1(x) + \sqrt{\Delta}\xi\right)\right)$$

Suppose the architecture, priors, activations are known.
The best test error is then given by *Bayesian inference*:



*Bayes posterior*

$$\mathbb{P}\left(a, \{W_\ell\}_{\ell=1}^L \,\middle|\, \mathcal{D}\right) \propto e^{-\frac{||a||^2}{2\Delta_a} - \Sigma_{\ell=1}^L \frac{||W_\ell||_F^2}{2\Delta_\ell}} \times \prod_{\ell=1}^L \int \frac{d\xi\, e^{-\frac{\xi^2}{2}}}{\sqrt{2\pi}}\, \delta\left(y^\star(x^\mu) - f^\star\left(\frac{a^\top}{\sqrt{k_L}} \varphi_L \circ \cdots \circ \varphi_1(x) + \sqrt{\Delta}\xi\right)\right)$$

Regression ($f^\star = id$)

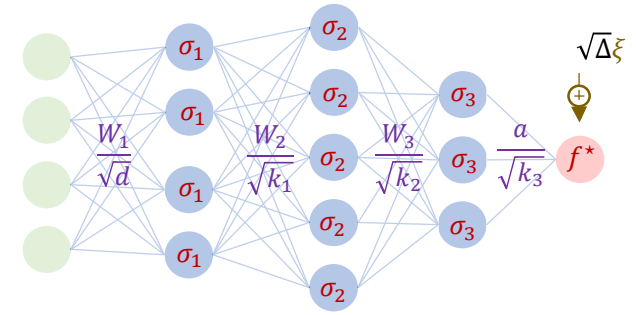$$\epsilon_{g,\mathrm{reg}}^{\mathrm{BO}} = \mathbb{E}_{\mathcal{D}, \{W_\ell^\star\}_{\ell=1}^L, \mathbf{a}_\star} \mathbb{E}_{\mathbf{x},y}\left[\left(y - \langle \hat{y}(\mathbf{x})\rangle_{\mathbf{a}, \{W_\ell\}_{\ell=1}^L \sim \mathbb{P}}\right)^2\right]$$

Classification ($f^\star = \mathrm{sign}$)

$$\epsilon_{g,\mathrm{class}}^{\mathrm{BO}} = \mathbb{E}_{\mathcal{D}, \{W_\ell^\star\}_{\ell=1}^L, \mathbf{a}_\star} \mathbb{P}_{\mathbf{x},y}\left[y \neq \mathrm{sign}\left(\langle \mathrm{sign}(\hat{y}(\mathbf{x}))\rangle_{\mathbf{a}, \{W_\ell\}_{\ell=1}^L \sim \mathbb{P}}\right)\right].$$

**Q1**. Can one provide a sharp asymptotic characterization of the Bayes-optimal error?

**Q2**. How do the test errors achieved by ERM algorithms in practice compare?

*Preliminaries*: Second-order statistics of random(-ish) neural nets

**A1** Bayes-optimal test errors

**A2** ERM test errors

*Preliminaries*: Second-order statistics of random(-ish) neural nets

*Why second order statistics?*

1. Appear naturally in the replica computation.

2. **Gaussian universality :** in a number of simple ERM settings, the test error only depends on the second order statistics of the data (*more later*)

Song Mei and Andrea Montanari. *Generalization Error of Random Features Regression: Precise Asymptotics and the Double Descent Curve. Commun*. Pure Appl. Math.,, 2022

Hong Hu and Yue M. Lu. *Universality Laws for High-Dimensional Learning with Random Features*. IEE Trans. Inf. Theory

Federica Gerace, Bruno Loureiro, Florent Krzakala, Marc Mezard, and Lenka Zdeborova. *Generalisation error in learning with random features and the hidden manifold model*. ICML 2020

*The shallow L=1 hidden layer case*

$x \sim \mathcal{N}(0, \Sigma)$

$\frac{F}{\sqrt{d}}$

$\sigma_1$

$\sigma_1$

$\sigma_1$

$\sigma_1$

*For fixed $F$, what is the covariance  $\Omega = \langle x_1 x_1^\top \rangle_x$ of the last layer post-activation wrt the Gaussian input randomness?*

$$x_1$$

$$x \sim \mathcal{N}(0, \Sigma)$$

$$\frac{F}{\sqrt{d}}$$

$$\sigma_1$$
$$\sigma_1$$
$$\sigma_1$$
$$\sigma_1$$

*For fixed $F$, what is the covariance $\Omega = \langle x_1 x_1^\top \rangle_x$ of the last layer post-activation wrt the Gaussian input randomness?*

*(Gaussian Equivalence Property)*

Defining

$$\kappa_1 = \mathbb{E}_{z \sim \mathcal{N}(0,1)}[\sigma_1(z)z]$$

$$\kappa_* = \sqrt{\mathbb{E}_{z \sim \mathcal{N}(0,1)}[\sigma_1(z)^2] - \kappa_1^2}$$

then simply

$$\Omega = \kappa_1^2 \frac{F\Sigma F^\top}{d} + \kappa_*^2 \mathbb{I}_k$$

Goldt, Mézard, Krzakala, and Zdeborová, *Modelling the influence of data structure on learning in neural networks.* PRX 2020

Draw two networks $W_1^a, \ldots, W_L^a, a^a$ and $W_1^b, \ldots, W_L^b, a^b$ i.i.d from the Bayes posterior.



What is the covariance $\Omega_L^{ab} = \langle x_L^a x_L^{b\top} \rangle_x$ ?

Draw two networks $W_1^a, \ldots, W_L^a, a^a$ and $W_1^b, \ldots, W_L^b, a^b$ i.i.d from the Bayes posterior.



What is the covariance $\Omega_L^{ab} = \left\langle x_L^a x_L^{b\top} \right\rangle_x$ ?

*(Deep Bayes conjecture)*

Defining

$$r_{\ell+1} = \Delta_{\ell+1} \mathbb{E}_{z \sim \mathcal{N}(0, r_\ell)} \left[ \sigma_\ell(z)^2 \right],$$

$$\kappa_1^{(\ell)} = \frac{1}{r_\ell} \mathbb{E}_{z \sim \mathcal{N}(0, r_\ell)} \left[ z \sigma_\ell(z) \right],$$

$$\kappa_*^{(\ell)} = \sqrt{\mathbb{E}_{z \sim \mathcal{N}(0, r_\ell)} \left[ \sigma_\ell(z)^2 \right] - r_\ell \left( \kappa_1^{(\ell)} \right)^2},$$

$\Omega_L^{ab}$ is given by the $L$th term of the recursion

$$\Omega_\ell^{ab} = \left( \kappa_1^{(\ell)} \right)^2 \frac{W_\ell^a \Omega_{\ell-1}^{ab} W_\ell^{b\top}}{k_{\ell-1}} + \delta_{ab} \left( \kappa_*^{(\ell)} \right)^2 \mathbb{I}_{k_\ell}$$

**HC**, Krzakala and Zdeborová, *Optimal learning of random networks of extensive width,* arXiv:2302.00375 (2023).

In terms of second-order activation statistics,



**Non-linear** deep network

In terms of second-order activation statistics,



**Non-linear** deep network

$\approx$

**Noisy, linear** deep network

**Q1**. Can one conjecture a sharp asymptotic characterization of the Bayes-optimal error?

**HC**, Krzakala and Zdeborová, *Optimal learning of random networks of extensive width,* arXiv:2302.00375 (2023).

$$y^{\star}(x) = f^{\star}\left(\frac{a^{\top}}{\sqrt{k_L}}\varphi_L \circ \cdots \circ \varphi_1(x) + \sqrt{\Delta}\mathcal{N}(0,1)\right)$$
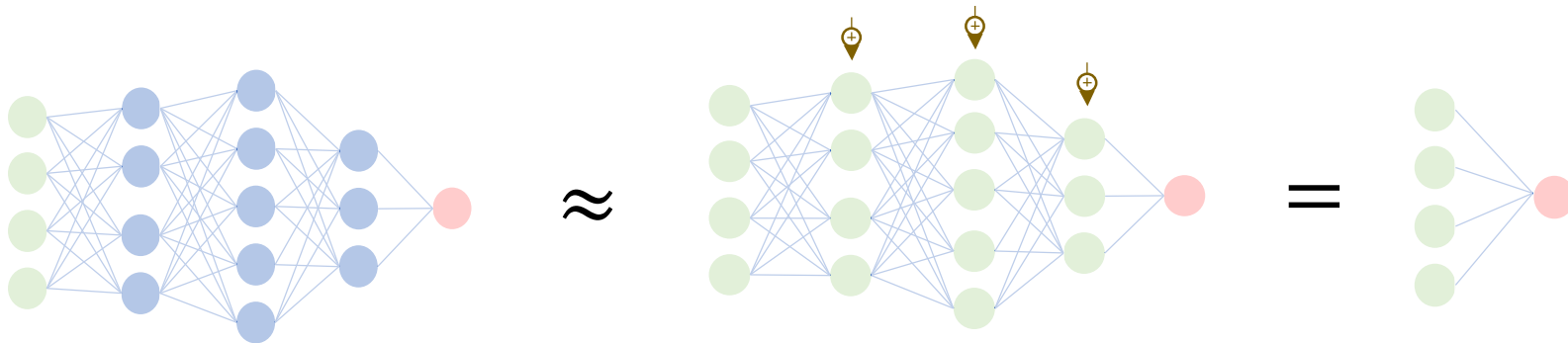
With layers $\varphi_\ell(h) = \sigma_\ell\left(\frac{W_\ell}{\sqrt{k_{\ell-1}}}\, h\right)$

$(W_\ell)_{ij} \sim \mathcal{N}(0,\Delta_\ell), \quad a_i \sim \mathcal{N}(0,\Delta_a)$

$$y^{\text{eq}}(x) = f^{\star}\left(\rho\frac{\theta^{\top}x}{\sqrt{d}} + \epsilon_r\mathcal{N}(0,1)\right)$$

With

$$\epsilon_r \equiv \sum_{\ell_0=1}^{L-1}(\kappa_*^{(\ell_0)})^2\Delta_a\prod_{\ell=\ell_0+1}^{L}(\kappa_1^{(\ell)})^2\Delta_\ell + (\kappa_*^{(L)})^2\Delta_a + \Delta$$

$$\rho \equiv \Delta_a\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2\Delta_\ell$$

$$\theta_i \sim \mathcal{N}(0,1)$$

***Conjecture*** : these two networks are characterized by the ***same Bayes optimal errors***

**Regression**

$$\epsilon_{g,\text{reg}}^{\text{BO}}=\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2\left(\Delta_a\left(\int z\mathrm{d}\mu(z)\right)\prod_{\ell=1}^{L}\Delta_\ell - q\right)+\epsilon_r$$
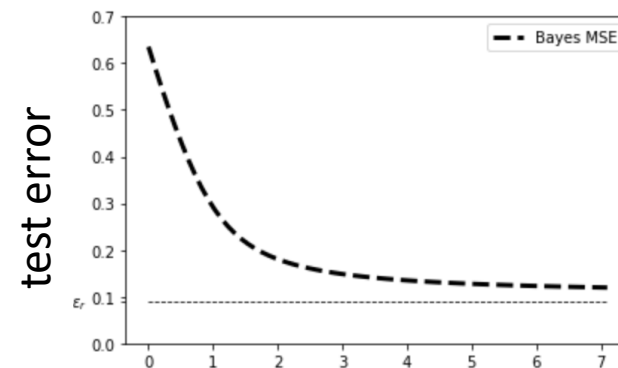
$$q = \frac{1}{2}\int \frac{\alpha \prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2 z^2\Delta_a^2 \prod_{\ell=1}^{L}\Delta_\ell^2}{\epsilon_{g,\text{reg}}^{\text{BO}} + \alpha \prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2 z\Delta_a \prod_{\ell=1}^{L}\Delta_\ell}\mathrm{d}\mu(z).$$
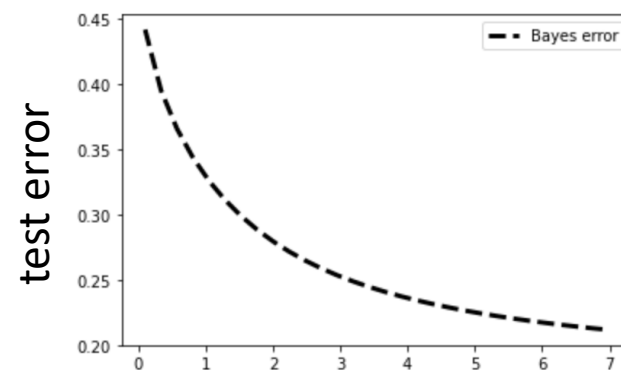
**Classification**

$$\epsilon_{g,\text{class}}^{\text{BO}}=\frac{1}{\pi}\arccos\left[\frac{\sqrt{\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2 q}}{\sqrt{\Delta_a\int z\mathrm{d}\mu(z)\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2\Delta_\ell + \epsilon_r}}\right]$$

$$\begin{cases} q = \int \frac{\hat{q}\Delta_a^2 \prod_{\ell=1}^{L}\Delta_\ell^2 z^2}{\hat{q}z\Delta_a \prod_{\ell=1}^{L}\Delta_\ell+1}\mathrm{d}\mu(z) \\[2mm] \hat{q} = \frac{2\alpha \prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2}{\Delta_a\int z\mathrm{d}\mu(z)\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2\Delta_\ell+\epsilon_r-\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2 q} \\[2mm] \int \frac{\mathrm{d}\xi}{(2\pi)^{\frac{3}{2}}}\frac{2e^{-\frac{1}{2}\frac{\Delta_a\int z\mathrm{d}\mu(z)\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2\Delta_\ell+\epsilon_r+\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2 q}{\Delta_a\int z\mathrm{d}\mu(z)\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2\Delta_\ell+\epsilon_r-\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2 q}\xi^2}}{1-\text{erf}\left(\frac{\prod_{\ell=1}^{L}\kappa_1^{(\ell)}\sqrt{q}\xi}{\sqrt{2\left(\Delta_a\int z\mathrm{d}\mu(z)\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2\Delta_\ell+\epsilon_r-\prod_{\ell=1}^{L}\left(\kappa_1^{(\ell)}\right)^2 q\right)}}\right)} \end{cases}$$

depth $= 3, \sigma = tanh$



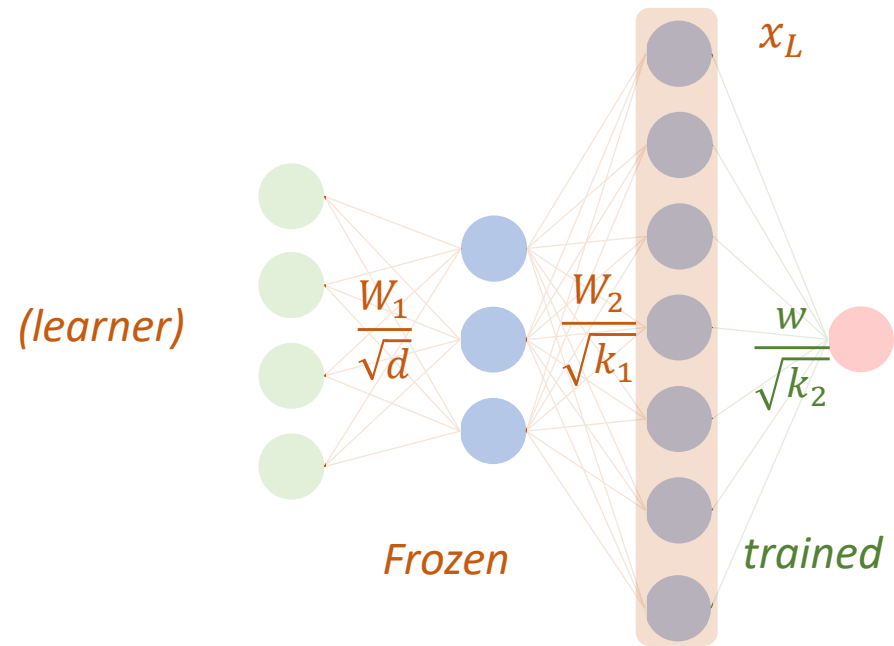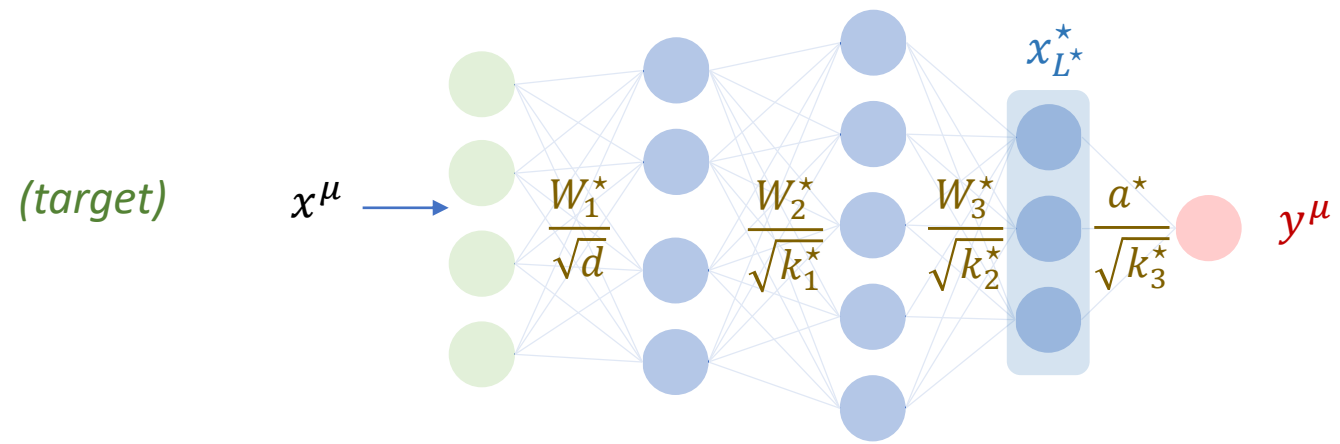test error vs # samples /dimension



test error vs # samples /dimension

$$\text{depth} = 2, \sigma = \text{ReLU} - \frac{1}{\sqrt{2\pi}}$$



Piccioli, Troiani and Zdeborová, *Sampling the posterior of neural networks,* arXiv:2306.02729

✓ **Q1**. Can one provide a sharp asymptotic characterization of the Bayes-optimal error?

**Q2**. How do the test errors achieved by ERM algorithms in practice compare?

*(target)*

$x^\mu \longrightarrow$

$\dfrac{W_1^\star}{\sqrt{d}}$ $\dfrac{W_2^\star}{\sqrt{k_1^\star}}$ $\dfrac{W_3^\star}{\sqrt{k_2^\star}}$ $x_{L^\star}^\star$ $\dfrac{a^\star}{\sqrt{k_3^\star}}$ $y^\mu$

*(learner)*

$\dfrac{W_1}{\sqrt{d}}$ $\dfrac{W_2}{\sqrt{k_1}}$ $x_L$ $\dfrac{w}{\sqrt{k_2}}$

*Frozen*  *trained*

*(target)*

$x^\mu \longrightarrow$

$\dfrac{W_1^\star}{\sqrt{d}}$   $\dfrac{W_2^\star}{\sqrt{k_1^\star}}$   $\dfrac{W_3^\star}{\sqrt{k_2^\star}}$   $x_{L^\star}^\star$   $\dfrac{a^\star}{\sqrt{k_3^\star}}$   $y^\mu$

$$\widehat{w} = \operatorname*{argmin}_{w} \left( \sum_{\mu=1}^{n} g\left( y^\mu, \frac{w^\top x_L^\mu}{\sqrt{k_L}} \right) + r(w) \right)$$

*(learner)*

$\dfrac{W_1}{\sqrt{d}}$   $\dfrac{W_2}{\sqrt{k_1}}$   $x_L$   $\dfrac{w}{\sqrt{k_2}}$

*Frozen*   *trained*

*(target)* $\quad x^\mu \longrightarrow$

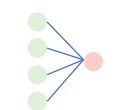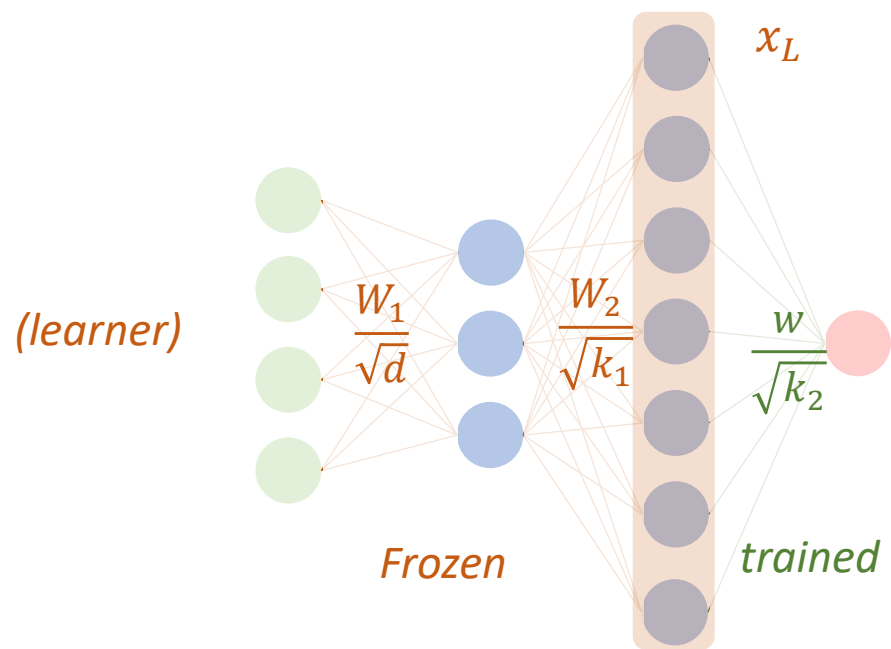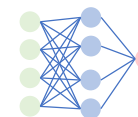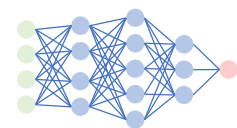$$\widehat{w} = \underset{w}{\operatorname{argmin}} \left( \sum_{\mu=1}^n g\left( y^\mu, \frac{w^\top x_L^\mu}{\sqrt{k_L}} \right) + r(w) \right)$$

- Ridge, LASSO, elastic net…
- Logistic / hinge/ ridge classification

- Random Features

- Deep Random Features

- Kernel regression/classification

*(learner)*

*Frozen*     *trained*

Introduce the **Gaussian clones** $u, v$ of $x_L, x_{L^\star}^\star$

$$u, v \sim \mathcal{N}\left(0, \begin{bmatrix} \langle x_L x_L^\top \rangle & \langle x_L x_{L^\star}^{\star\top} \rangle \\ \langle x_{L^\star}^\star x_L^\top \rangle & \langle x_{L^\star}^\star x_{L^\star}^{\star\top} \rangle \end{bmatrix}\right)$$

Schröder, **HC,** Dmitriev and Loureiro, *Deterministic equivalent and error universality of deep random features learning,* arXiv:2302.00401 (2023).

Introduce the **Gaussian clones** $u, v$ of $x_L, x_{L^\star}^\star$

$$u, v \sim \mathcal{N}\left(0, \begin{bmatrix} \langle x_L x_L^\top \rangle & \langle x_L x_{L^\star}^{\star\,\top} \rangle \\ \langle x_{L^\star}^\star x_L^\top \rangle & \langle x_{L^\star}^\star x_{L^\star}^{\star\,\top} \rangle \end{bmatrix}\right)$$

(ERM) $\quad \mathcal{D} = \left\{ x^\mu, y^\mu = f^\star\left(\dfrac{a_\star^\top x_{L^\star}^{\star\,\mu}}{\sqrt{k_{L^\star}^\star}}\right)\right\} \quad \widehat{w} = \underset{w}{\mathrm{argmin}}\left(\displaystyle\sum_{\mu=1}^n g\left(y^\mu, \dfrac{w^\top x_L^\mu}{\sqrt{k_L}}\right) + r(w)\right)$
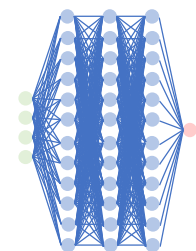
(ERMg) $\quad \mathcal{D}^G = \left\{ u^\mu, y^\mu = f^\star\left(\dfrac{a_\star^\top v^\mu}{\sqrt{k_{L^\star}^\star}}\right)\right\} \quad \widehat{w} = \underset{w}{\mathrm{argmin}}\left(\displaystyle\sum_{\mu=1}^n g\left(y^\mu, \dfrac{w^\top u^\mu}{\sqrt{k_L}}\right) + r(w)\right)$

Schröder, **HC,** Dmitriev and Loureiro, *Deterministic equivalent and error universality of deep random features learning,* arXiv:2302.00401 (2023).
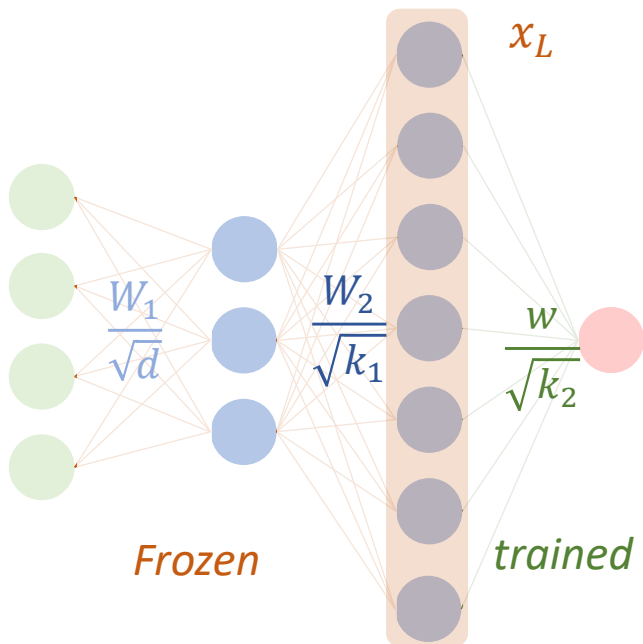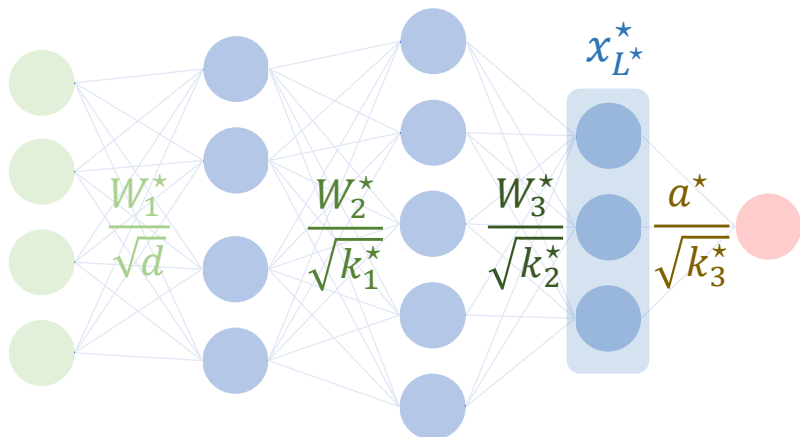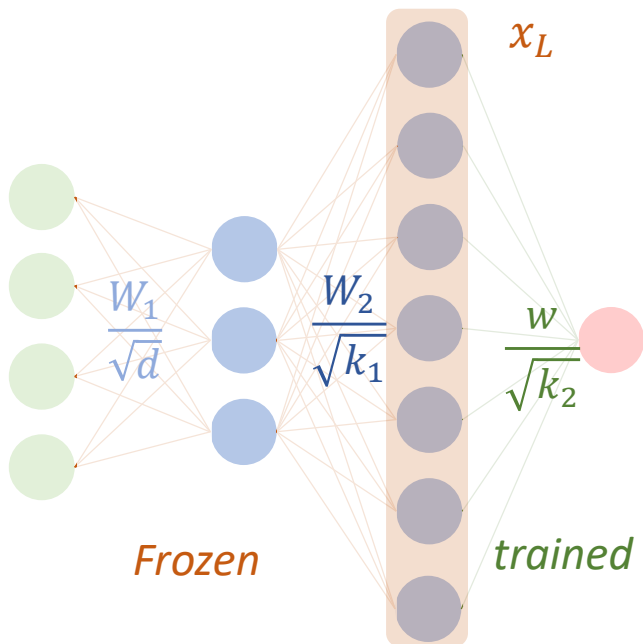
Introduce the **Gaussian clones** $u, v$ of $x_L, x_{L^\star}^\star$

$$u, v \sim \mathcal{N}\left(0, \begin{bmatrix} \langle x_L x_L^\top \rangle & \langle x_L x_{L^\star}^{\star\top} \rangle \\ \langle x_{L^\star}^\star x_L^\top \rangle & \langle x_{L^\star}^\star x_{L^\star}^{\star\top} \rangle \end{bmatrix}\right)$$

(ERM) $\quad \mathcal{D} = \left\{ x^\mu, y^\mu = f^\star\left(\frac{a_\star^\top x_{L^\star}^{\star\mu}}{\sqrt{k_{L^\star}^\star}}\right)\right\} \quad \widehat{w} = \underset{w}{\mathrm{argmin}}\left(\sum_{\mu=1}^n g\left(y^\mu, \frac{w^\top x_L^\mu}{\sqrt{k_L}}\right) + r(w)\right)$

(ERMg) $\quad \mathcal{D}^G = \left\{ u^\mu, y^\mu = f^\star\left(\frac{a_\star^\top v^\mu}{\sqrt{k_{L^\star}^\star}}\right)\right\} \quad \widehat{w} = \underset{w}{\mathrm{argmin}}\left(\sum_{\mu=1}^n g\left(y^\mu, \frac{w^\top u^\mu}{\sqrt{k_L}}\right) + r(w)\right)$

***Conjecture: (part 1) (Gaussian universality)*** The learning problems (ERM) and (ERMg) lead to the same test error and training loss.

Schröder, **HC,** Dmitriev and Loureiro, *Deterministic equivalent and error universality of deep random features learning,* arXiv:2302.00401 (2023).

$$u, v \sim \mathcal{N} \left( 0, \begin{bmatrix} \langle x_L x_L^\top \rangle & \langle x_L x_{L^\star}^\star {}^\top \rangle \\ \langle x_{L^\star}^\star x_L^\top \rangle & \langle x_{L^\star}^\star x_{L^\star}^\star {}^\top \rangle \end{bmatrix} \right)$$

**Conjecture: (part 2)** Furthermore, the covariances $\langle x_L x_L^\top \rangle$, $\langle x_{L^\star}^\star x_{L^\star}^\star {}^\top \rangle$ and $\langle x_{L^\star}^\star x_L^\top \rangle$ can be computed simply with the noisy equivalent model.

$$u, v \sim \mathcal{N}\left(0, \begin{bmatrix} \langle x_L x_L^\top \rangle & \langle x_L x_{L^\star}^{\star\,\top} \rangle \\ \langle x_{L^\star}^\star x_L^\top \rangle & \langle x_{L^\star}^\star x_{L^\star}^{\star\,\top} \rangle \end{bmatrix}\right)$$
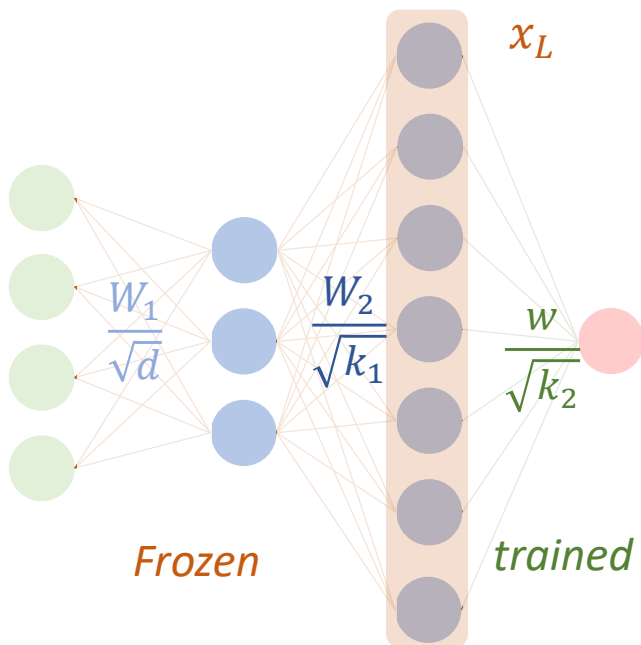
**Conjecture: (part 2)** Furthermore, the covariances $\langle x_L x_L^\top \rangle$, $\langle x_{L^\star}^\star x_{L^\star}^{\star\,\top} \rangle$ and $\langle x_{L^\star}^\star x_L^\top \rangle$ can be computed simply with the noisy equivalent model.

*Here for instance*

$$\langle x_L x_L^\top \rangle = \kappa_1^{(1)^2} \kappa_1^{(2)^2} \frac{W_2 W_1 \Sigma \, W_1^\top W_2^\top}{d k_1} + \kappa_*^{(1)^2} \kappa_1^{(2)^2} \frac{W_2 W_2^\top}{k_1} + \kappa_*^{(2)^2} \mathbb{I}_{k_1}$$

$$\langle x_{L^\star}^\star x_{L^\star}^{\star\,\top} \rangle = \kappa_1^{\star(1)^2} \kappa_1^{\star(2)^2} \kappa_1^{\star(3)^2} \frac{W_3^\star W_2^\star W_1^\star \Sigma \; W_1^{\star\top} W_2^{\star\top} W_3^{\star\top}}{d k_1^\star k_2^\star} +$$
$$\kappa_*^{\star(1)^2} \kappa_1^{\star(2)^2} \kappa_1^{\star(3)^2} \frac{W_3^\star W_2^\star W_2^{\star\top} W_3^{\star\top}}{k_1^\star k_2^\star} + \kappa_*^{\star(2)^2} \kappa_1^{\star(3)^2} \frac{W_3^\star W_3^{\star\top}}{k_2^\star} + \kappa_*^{\star(2)^2} \mathbb{I}_{k_2^\star}$$

$$\langle x_{L^\star}^\star x_L^\top \rangle = \kappa_1^{(1)} \kappa_1^{(2)} \, \kappa_1^{\star(1)} \, \kappa_1^{\star(2)} \, \kappa_1^{\star(3)} \frac{W_3^\star W_2^\star W_1^\star \Sigma \, W_1^\top W_2^\top}{d\sqrt{k_1 k_1^\star k_2^\star}}$$

So one just needs to solve the proxy ERM

$$u, v \sim \mathcal{N}\left(0, \begin{bmatrix} \langle x_L x_L^\top \rangle & \langle x_L x_{L^\star}^{\star\,\top} \rangle \\ \langle x_{L^\star}^\star x_L^\top \rangle & \langle x_{L^\star}^\star x_{L^\star}^{\star\,\top} \rangle \end{bmatrix}\right)$$
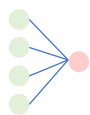
(ERMg)

$$\mathcal{D}^G = \left\{ u^\mu, y^\mu = f^\star\left(\frac{a_\star^\top v^\mu}{\sqrt{k_{L^\star}^\star}}\right) \right\}$$

$$\widehat{w} = \underset{w}{\mathrm{argmin}} \left( \sum_{\mu=1}^n g\left(y^\mu, \frac{w^\top u^\mu}{\sqrt{k_L}}\right) + r(w) \right)$$

***Theorem (informal)*** : The test error of the problem (ERMg) can be characterized in terms of three order parameters $q, m, V$ given as the solution of a system of self-consistent equations.

$$\begin{cases} V = \mathbb{E}_{(\omega,\bar\theta)\sim\mu}\left[\frac{\omega}{\lambda+\hat V\omega}\right] \\ m = \frac{\hat m}{\sqrt{\gamma}}\mathbb{E}_{(\omega,\bar\theta)\sim\mu}\left[\frac{\bar\theta^2}{\lambda+\hat V\omega}\right] \\ q = \mathbb{E}_{(\omega,\bar\theta)\sim\mu}\left[\frac{\hat m^2\bar\theta^2\omega+\hat q\omega^2}{(\lambda+\hat V\omega)^2}\right] \end{cases}, \quad \begin{cases} \hat V = \frac{\alpha}{V}(1 - \mathbb{E}_{s,h\sim\mathcal{N}(0,1)}[f_g'(V,m,q)]) \\ \hat m = \frac{1}{\sqrt{\rho\gamma}}\frac{\alpha}{V}\mathbb{E}_{s,h\sim\mathcal{N}(0,1)}\left[sf_g(V,m,q)-\frac{m}{\sqrt{\rho}}f_g'(V,m,q)\right] \\ \hat q = \frac{\alpha}{V^2}\mathbb{E}_{s,h\sim\mathcal{N}(0,1)}\left[\left(\frac{m}{\sqrt{\rho}}s+\sqrt{q-\frac{m^2}{\rho}}h-f_g(V,m,q)\right)^2\right] \end{cases}$$

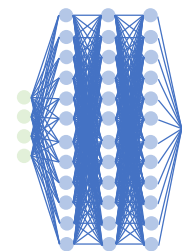Loureiro, Gerbelot, **HC**, Goldt, Krzakala, Mézard and Zdeborová, *Learning curves of generic feature maps for realistic datasets with a teacher-student model,* NeurIPS 2021

$$\epsilon_g = \rho \int z\,\mathrm{d}\mu(z) + q - 2\prod_{\ell=1}^{L}\kappa_1^{(\ell)}m + \epsilon_r$$

$$\begin{cases}\hat{V} = \frac{\alpha}{1+V}\\[4pt]\hat{q} = \alpha\frac{\epsilon_g}{(1+V)^2}\\[4pt]\hat{m} = \frac{\prod_{\ell=1}^{L}\kappa_1^{(\ell)}\alpha}{1+V}\end{cases}\qquad\begin{cases}V = \int\frac{z}{\lambda+\hat{V}z}\mathrm{d}\mu(z)\\[6pt]q = \int\frac{\Delta_a\prod_{\ell=1}^{L}\Delta_\ell\hat{m}^2 z^3+\hat{q}z^2}{(\lambda+\hat{V}z)^2}\mathrm{d}\mu(z)\\[6pt]m = \Delta_a\prod_{\ell=1}^{L}\Delta_\ell\hat{m}\int\frac{z^2}{\lambda+\hat{V}z}\mathrm{d}\mu(z)\end{cases}$$

$$\epsilon_g = \rho \int z\,\mathrm{d}\mu(z) + q - 2\prod_{\ell=1}^{L}\kappa_1^{(\ell)}m + \epsilon_r$$

$$\begin{cases}\hat{V} = \frac{\frac{\alpha}{\gamma}}{1+V}\\[4pt]\hat{q} = \frac{\alpha}{\gamma}\frac{\epsilon_g}{(1+V)^2}\\[4pt]\hat{m} = \sqrt{\Delta_a\prod_{\ell=1}^{L}\Delta_\ell}\sqrt{\gamma}\frac{\prod_{\ell=1}^{L}\kappa_1^{(\ell)}\frac{\alpha}{\gamma}}{1+V}\end{cases}$$

$$\begin{cases}V = \frac{1}{\hat{V}} - \frac{\lambda}{\hat{V}^2\kappa_1^2}g\left(-\frac{\lambda+\hat{V}\kappa_*^2}{\hat{V}\kappa_1^2}\right)\\[6pt]q = \frac{\hat{m}^2+\hat{q}}{\hat{V}^2} - \frac{1}{\kappa_1^2\hat{V}^2}\left(\frac{2\lambda(\hat{m}^2+\hat{q})}{\hat{V}} + \hat{m}^2\kappa_*^2\right)g\left(-\frac{\lambda+\hat{V}\kappa_*^2}{\hat{V}\kappa_1^2}\right)\\[6pt]\qquad + \frac{\lambda}{\kappa_1^4\hat{V}^3}\left(\frac{\lambda(\hat{m}^2+\hat{q})}{\hat{V}} + \hat{m}^2\kappa_*^2\right)g'\left(-\frac{\lambda+\hat{V}\kappa_*^2}{\hat{V}\kappa_1^2}\right)\\[6pt]m = \sqrt{\gamma}\frac{\hat{m}}{\hat{V}}\left[1 - \frac{1}{\kappa_1^2}\left(\frac{\lambda}{\hat{V}} + \kappa_*^2\right)g\left(-\frac{\lambda+\hat{V}\kappa_*^2}{\hat{V}\kappa_1^2}\right)\right].\end{cases}$$

$$\epsilon_g = \rho \int z\,\mathrm{d}\mu(z) + q - 2\prod_{\ell=1}^{L}\kappa_1^{(\ell)}m + \epsilon_r$$

$$\begin{cases}\hat{V} = \frac{\alpha}{1+V}\\[4pt]\hat{q} = \alpha\frac{\epsilon_g}{(1+V)^2}\\[4pt]\hat{m} = \alpha\frac{\prod_{\ell=1}^{L}\kappa_1^{(\ell)}}{1+V}\end{cases}\qquad\begin{cases}V = \frac{\kappa_*^2}{\lambda} + \frac{\kappa_1^2}{\lambda+\hat{V}\kappa_1^2}\\[6pt]q = \frac{\Delta_a\prod_{\ell=1}^{L}\Delta_\ell\hat{m}^2\kappa_1^4+\hat{q}\kappa_1^4}{(\lambda+\hat{V}\kappa_1^2)^2}\\[6pt]m = \Delta_a\prod_{\ell=1}^{L^\star}\Delta_\ell\hat{m}\frac{\kappa_1^2}{\lambda+\hat{V}\kappa_1^2}\end{cases}$$
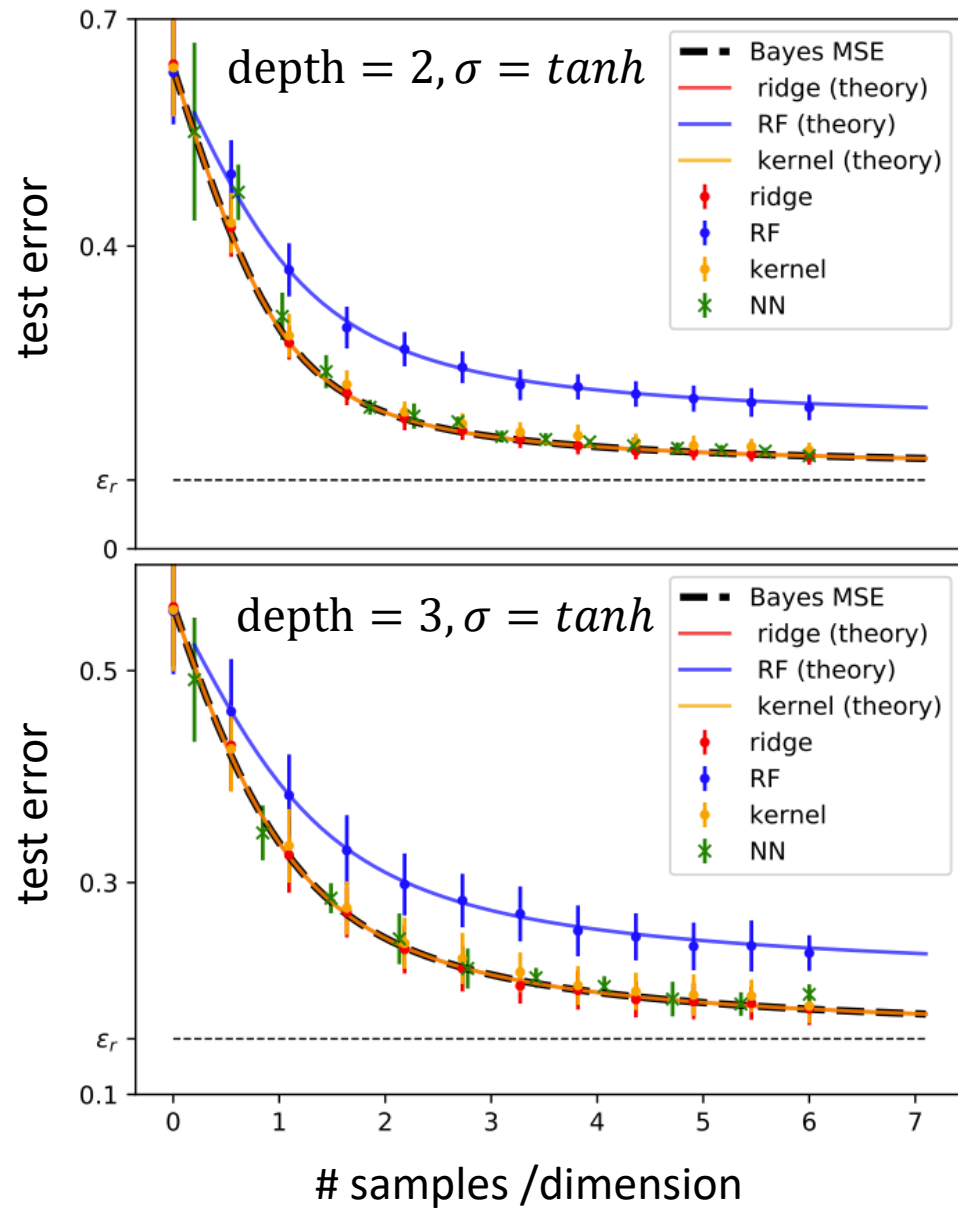
*See Daniil's poster!*

*Summary*:

✓ **Q1** We have sharp asymptotics for the Bayes optimal error of a deep, random network
  = *lowest information theoretically achievable error*

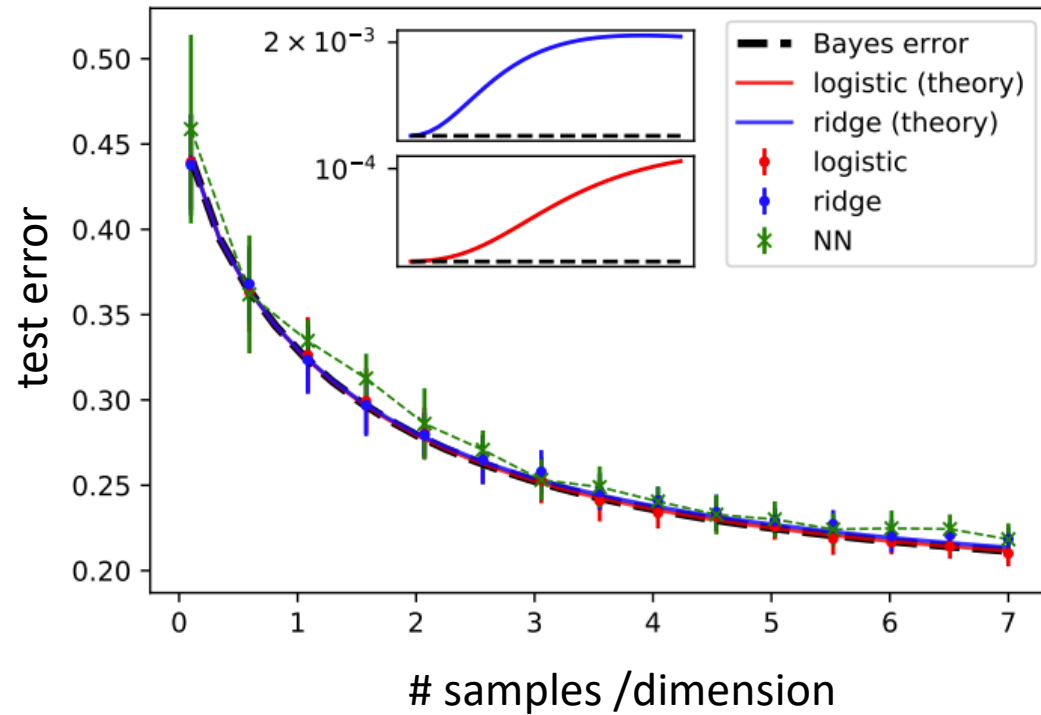✓ **Q2a** We have sharp asymptotics for test error of a large class of ERM algorithms on the same target.
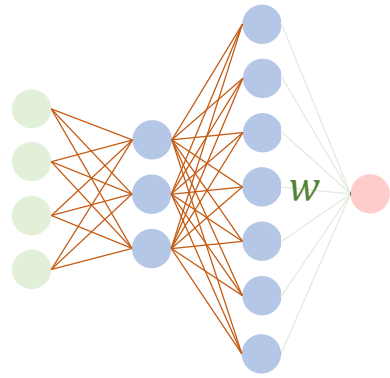
**Q2b** How do they compare?

Optimally regularized ridge regression and kernel regression **are Bayes optimal**.

$$\text{depth} = 3, \sigma = tanh$$



Optimally regularized logistic and ridge classification *are close to Bayes optimal*.
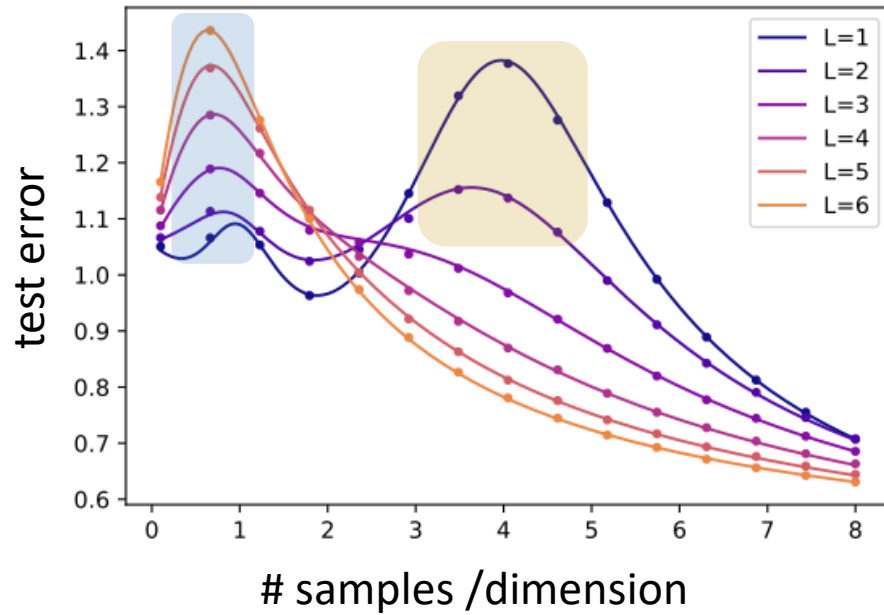
(learner)

(GEP) ≈ $w^\top A x + \xi$

signal
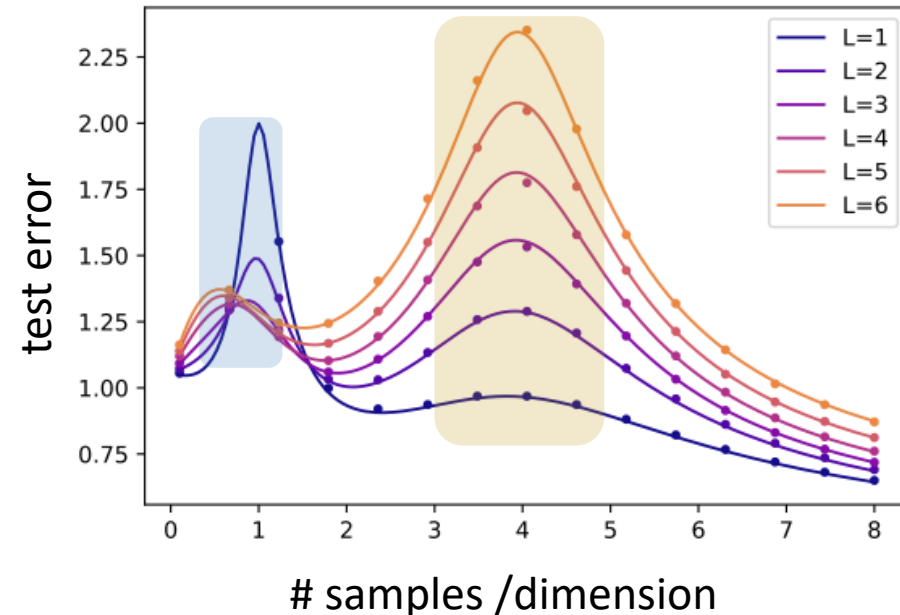
noise

When the signal is used to interpolate, the noise behaves as an depth-induced *implicit regularization*.

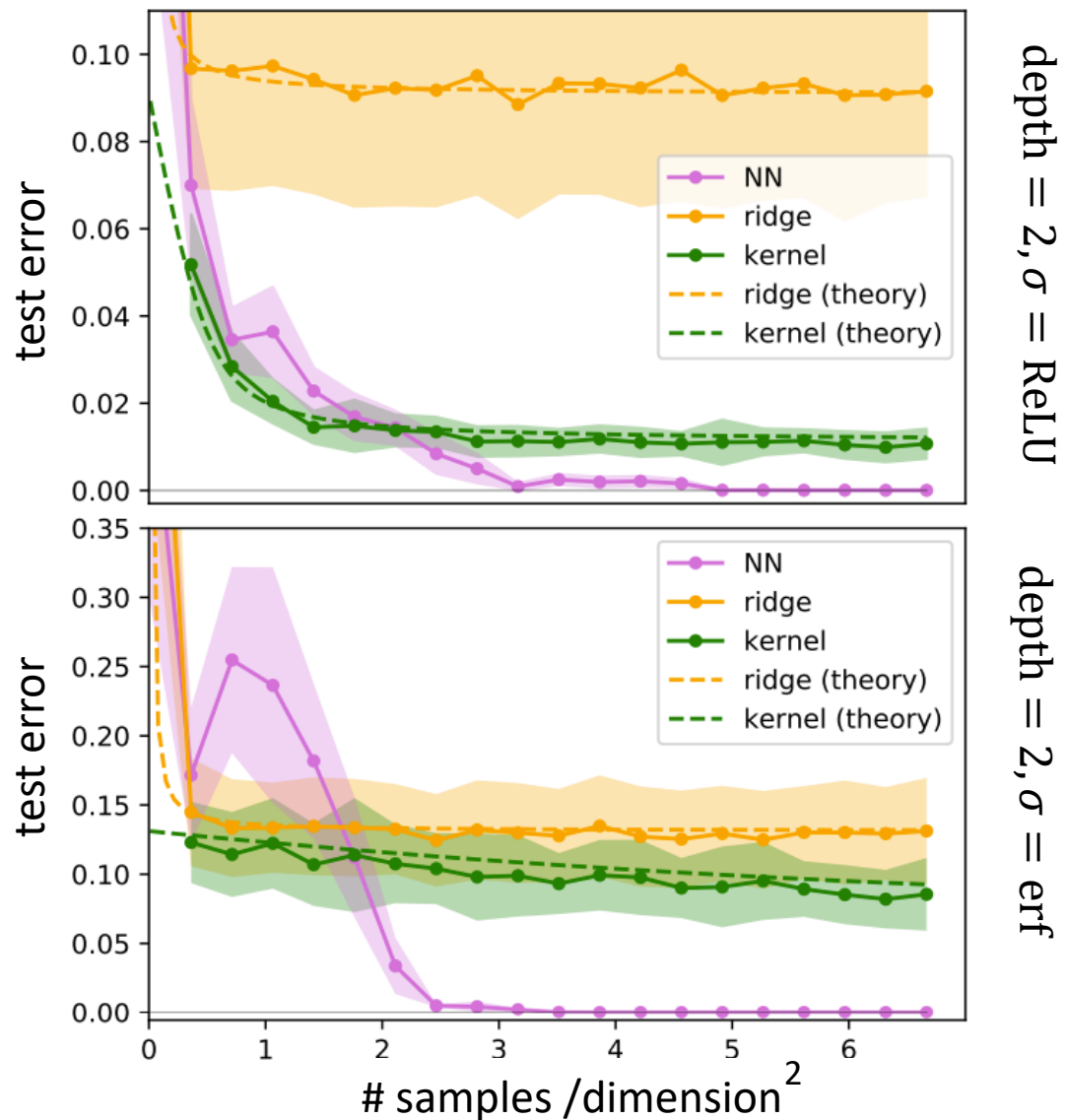A second peak appears when the noise is used to interpolate the train set.



D'Ascoli, Sagun and Biroli. *Triple descent and the two kinds of overfitting* J. Stat. Mech. 2021

**Q2** Can ERM methods achieve the Bayes error?

**A2** Yes, because in the $n \sim d$ regime ***only second-order statistics*** seem to be learnt, and in terms of those the target is equivalent to a single-layer network.

depth = 2, $\sigma$ = ReLU

depth = 2, $\sigma$ = erf

When $n \sim d^2$ , *higher-order statistics are learnt*, the Gaussian equivalences break down.

Misiakiewicz, *Sharp asymptotics of kernel ridge regression beyond the linear regime*, 2022

Hu and Lu. *Sharp asymptotics of kernel ridge regression beyond the linear regime*, 2022

Bordelon, Canatar, Pehlevan. *Spectrum dependent learning curves in kernel regression and wide neural networks*, 2020

*Takeaways*:

- In terms of *second order statistics* wrt a Gaussian input, a deep non-linear network is equivalent to a noisy linear network.

- Hence, In the $n \sim d$ regime, they are **characterized by the same Bayes / ERM errors.**

- Thus, single-layer ERM learners are Bayes optimal.

*Challenge /Future work:*

There is a need for a theory of finite-width architectures in *super linear regimes.*

- In terms of *second order statistics* wrt a Gaussian input, a deep non-linear network is equivalent to a noisy linear network.

- Hence, In the $n \sim d$ regime, they are **characterized by the same Bayes / ERM errors.**

- Thus, single-layer ERM learners are Bayes optimal.

*Challenge /Future work:*

There is a need for a theory of finite-width architectures in *super linear regimes.*

*Thank you for your attention !*