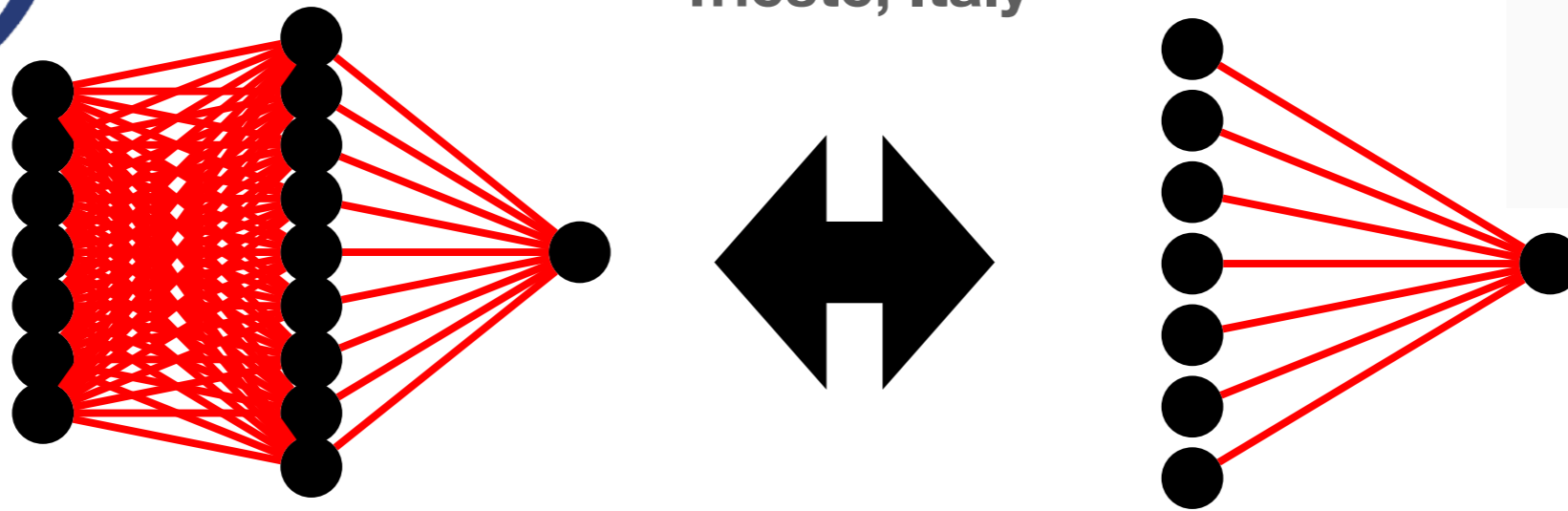


Fundamental limits of shallow neural networks in supervised learning



Jean Barbier

International Center for Theoretical Physics
Trieste, Italy



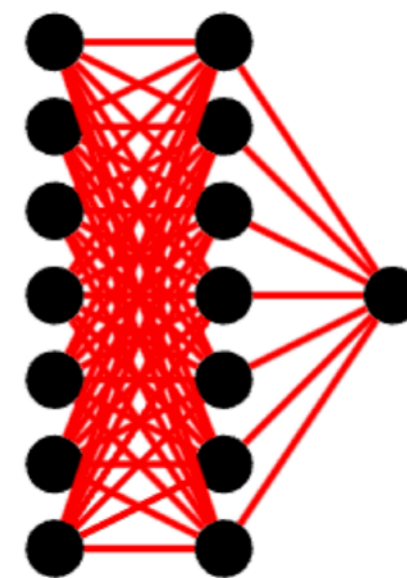
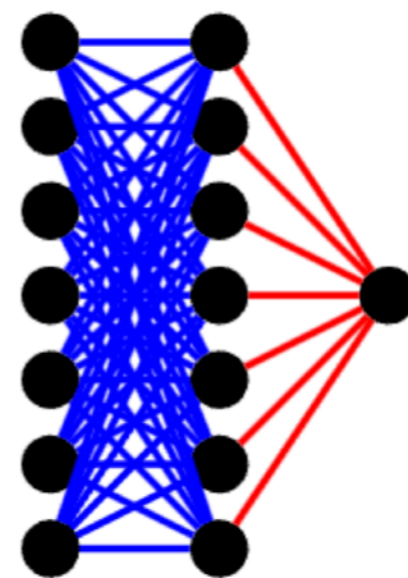
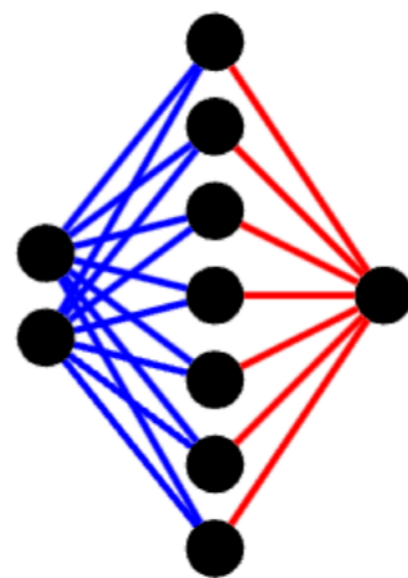
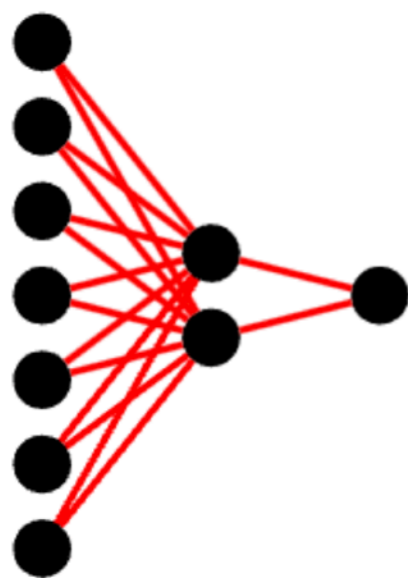
Joint with
Francesco Camilli
and Daria Tieplova
(ICTP)



Neural Networks Zoology

Committee machine

Random features models



(1)

(2)

(3)

(4)

(5)

(6)

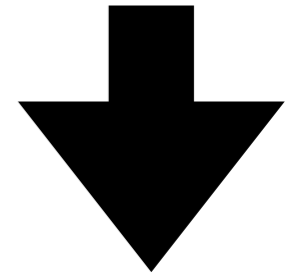
GLM/perceptron

Mean field regime

Fully learnable 2LNN
Proportional regime

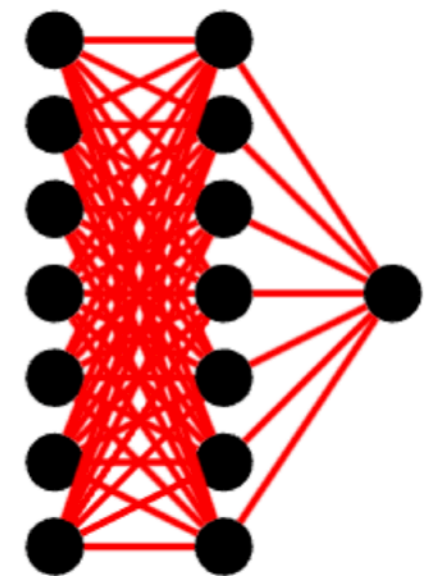
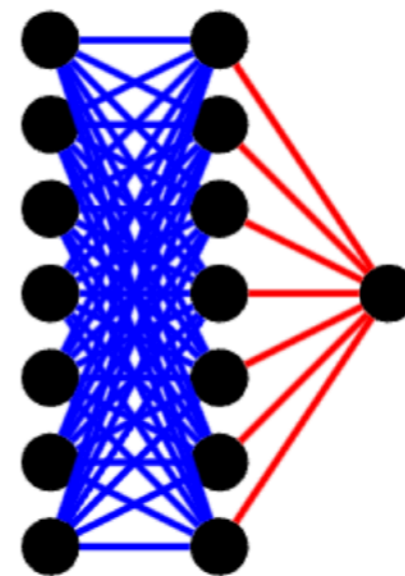
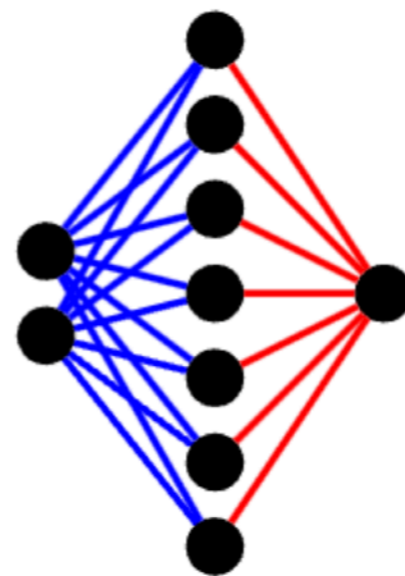
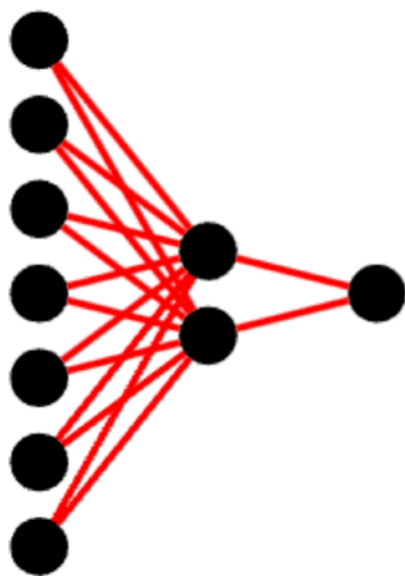
Neural Networks Zoology

But not « too much data »



Committee machine

Random features models



(1)

(2)

(3)

(4)

(5)

(6)

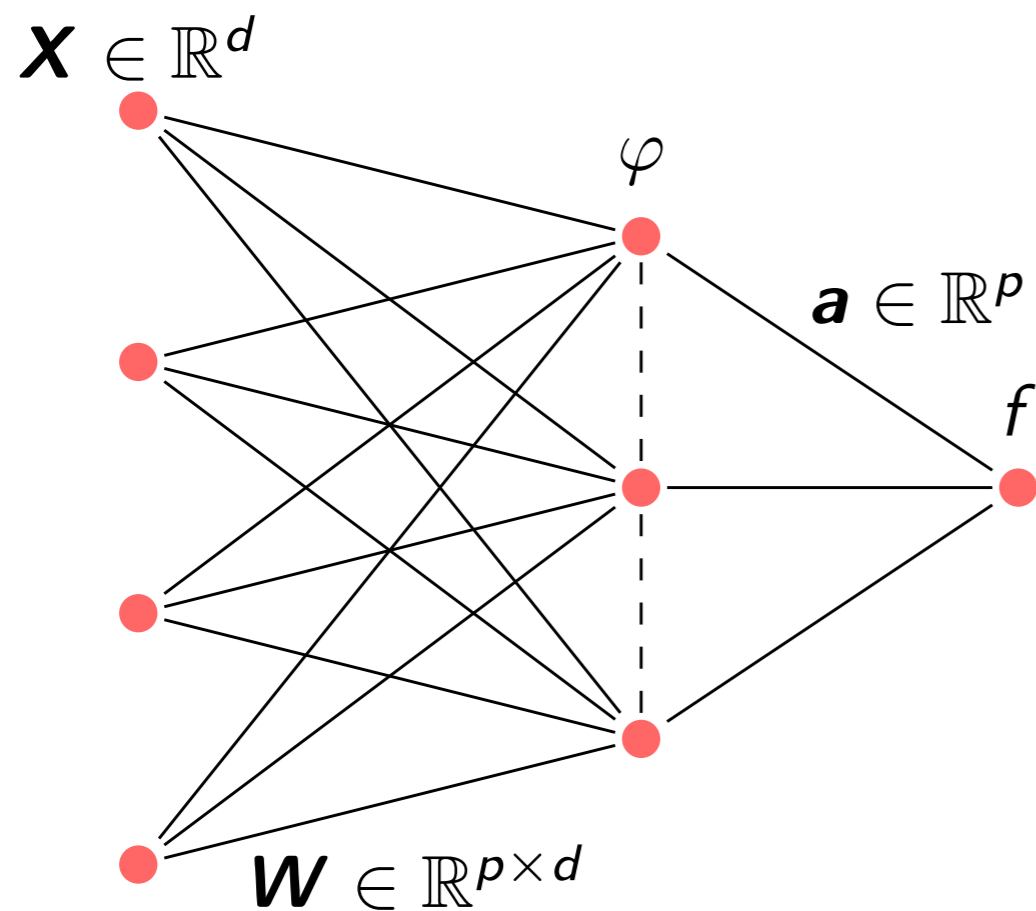
GLM/perceptron

Mean field regime

Fully learnable 2LNN
Proportional regime

Set-up

A two-layer NN



$$Y = f_{\mathbf{A}} \left(\frac{\mathbf{a}^T}{\sqrt{p}} \varphi \left(\frac{\mathbf{W}\mathbf{X}}{\sqrt{d}} \right) \right).$$

Supervised learning: Starting from a training set $\mathcal{D}_n = \{(\mathbf{X}_\mu, Y_\mu)_{\mu=1}^n\}$, we adjust the weights \mathbf{a} , \mathbf{W} s.t.

$$Y_\mu \approx f_{\mathbf{A}_\mu} \left(\frac{\mathbf{a}^T}{\sqrt{p}} \varphi \left(\frac{\mathbf{W}\mathbf{X}_\mu}{\sqrt{d}} \right) \right), \quad \forall \mu.$$

Main Goal

Produce the smallest possible generalization error:

$$\mathcal{E} = \left(Y_{\text{new}} - f_{\mathbf{A}_{\text{new}}} \left(\frac{\mathbf{a}^T}{\sqrt{p}} \varphi \left(\frac{\mathbf{W}\mathbf{X}_{\text{new}}}{\sqrt{d}} \right) \right) \right)^2$$

for a new couple $(\mathbf{X}_{\text{new}}, Y_{\text{new}})$.

What affects \mathcal{E} ?

Among the many factors that can affect it, the most relevant ones for us are:

- The size of the training set n : the more data points the more accurate we expect to be;
- The size of the network itself, parameterized by p ;
- The dimensionality of the input d ;
- The method used to train the network (e.g. ERM, SGD, **Bayes**)
- **The nature of the dataset**, i.e. what is the true underlying function $Y = f(\mathbf{X})$ the NN aims at approximating
- many more...

Central questions

What is the least possible \mathcal{E} ? When is it achieved?

Teacher-student setup

We assume that the training set is generated itself by a 2-layer **teacher network** with matching architecture:

$$Y_\mu = f_{\mathbf{A}_\mu} \left(\frac{\mathbf{a}^{*\top}}{\sqrt{p}} \varphi \left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}} \right) \right) + \underbrace{\sqrt{\Delta} Z_\mu}_{\text{label noise}}, \quad \forall \mu \leq n,$$

for some $\Delta > 0$, $Z_\mu \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$. **Prior on the weights:** $a_i^*, W_{ij}^* \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$.
Same as

$$Y_\mu \sim P_{\text{out}} \left(\cdot \mid \frac{\mathbf{a}^{*\top}}{\sqrt{p}} \varphi \left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}} \right) \right)$$

with

$$P_{\text{out}}(y \mid x) = \int \frac{dP_{\mathbf{A}}(\mathbf{A})}{\sqrt{2\pi\Delta}} \exp \left(-\frac{1}{2\Delta} (f_{\mathbf{A}}(x) - y)^2 \right)$$

Main theoretical restriction

We consider $\mathbf{X}_\mu \stackrel{\text{iid}}{\sim} \mathcal{N}(0, I_d)$, *structureless* input data.

Bayes-optimal student

Definition (informal)

A student network is Bayes-optimal if it is completely aware of the generative model

$$Y_\mu = f_{\mathbf{A}_\mu} \left(\frac{\mathbf{a}^{*\top}}{\sqrt{p}} \varphi \left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}} \right) \right) + \sqrt{\Delta} Z_\mu, \quad \forall \mu \leq n,$$

and it matches the teacher's architecture. In other words: **a part from the true weights, it knows everything there is know.**

A Bayes-optimal student has access to the **Bayes-posterior**:

$$dP(\boldsymbol{\theta} \mid \mathcal{D}_n) = \frac{1}{\mathcal{Z}(\mathcal{D}_n)} \prod_{\mu=1}^n P_{\text{out}} \left(Y_\mu \mid \frac{\mathbf{a}^\top}{\sqrt{p}} \varphi \left(\frac{\mathbf{W} \mathbf{X}_\mu}{\sqrt{d}} \right) \right) D\boldsymbol{\theta}$$

where $D\boldsymbol{\theta} = D\mathbf{a}D\mathbf{W}$ is the Gaussian prior on the weights.

Proposition (informal)

A Bayes-optimal student NN achieves the lowest expected generalization error

$$\mathbb{E}\mathcal{E} := \mathbb{E}\left(Y_{\text{new}} - \hat{Y}(\mathcal{D}_n, \mathbf{X}_{\text{new}})\right)^2$$

that is yielded by the BO predictor

$$\begin{aligned}\hat{Y}_{\text{Bayes}}(\mathcal{D}_n, \mathbf{X}_{\text{new}}) &= \mathbb{E}[Y_{\text{new}} \mid \mathcal{D}_n, \mathbf{X}_{\text{new}}] \\ &= \int dY Y P_{\text{out}}\left(Y \mid \frac{\mathbf{a}^\top}{\sqrt{p}} \varphi\left(\frac{\mathbf{W}\mathbf{X}_{\text{new}}}{\sqrt{d}}\right)\right) dP(\boldsymbol{\theta} \mid \mathcal{D}_n).\end{aligned}$$

Main information theoretic quantities

Recall

$$dP(\mathbf{a}, \mathbf{W} \mid \mathcal{D}_n) = \frac{1}{\mathcal{Z}(\mathcal{D}_n)} \prod_{\mu=1}^n P_{\text{out}} \left(Y_{\mu} \mid \frac{\mathbf{a}^{\top}}{\sqrt{p}} \varphi \left(\frac{\mathbf{W} \mathbf{X}_{\mu}}{\sqrt{d}} \right) \right) D\mathbf{a} D\mathbf{W} \rightarrow \langle \cdot \rangle.$$

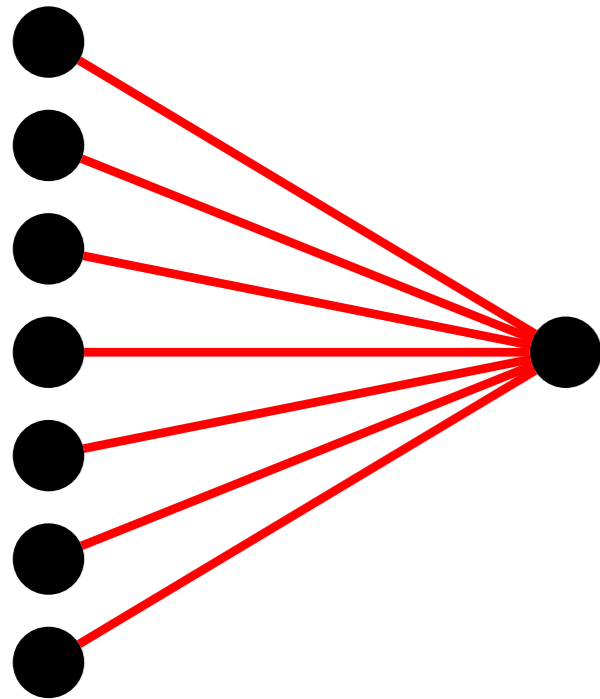
- **Partition function or evidence:**

$$\mathcal{Z}(\mathcal{D}_n) = \int \prod_{\mu=1}^n P_{\text{out}} \left(Y_{\mu} \mid \frac{\mathbf{a}^{\top}}{\sqrt{p}} \varphi \left(\frac{\mathbf{W} \mathbf{X}_{\mu}}{\sqrt{d}} \right) \right) D\mathbf{a} D\mathbf{W}$$

- **free entropy:** $\bar{f}_n = \frac{1}{n} \mathbb{E} \log \mathcal{Z}(\mathcal{D}_n)$
- **Mutual Information** per data point:

$$\begin{aligned} \frac{I_n(\mathbf{a}^*, \mathbf{W}^*; \mathcal{D}_n)}{n} &= \frac{H(\mathcal{D}_n)}{n} - \frac{H(\mathcal{D}_n \mid \mathbf{a}^*, \mathbf{W}^*)}{n} \\ &= -\bar{f}_n + \mathbb{E} \log P_{\text{out}} \left(Y_1 \mid \frac{\mathbf{a}^{*\top}}{\sqrt{p}} \varphi \left(\frac{\mathbf{W}^* \mathbf{X}_1}{\sqrt{d}} \right) \right) \end{aligned}$$

A simpler ancestor: the GLM



The teacher Generalized Linear Model is given by:

$$Y_{\mu}^{\circ} = f_{\mathbf{A}_{\mu}} \left(\rho \frac{\mathbf{v}^{*\top} \mathbf{X}_{\mu}}{\sqrt{d}} + \sqrt{\epsilon} \xi_{\mu}^{*} \right) + \sqrt{\Delta} Z_{\mu},$$

$$\text{or } Y_{\mu}^{\circ} \sim P_{\text{out}} \left(\cdot \mid \rho \frac{\mathbf{v}^{*\top} \mathbf{X}_{\mu}}{\sqrt{d}} + \sqrt{\epsilon} \xi_{\mu}^{*} \right)$$

with $v_i^*, \xi_{\mu}^* \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$, $\rho, \epsilon \geq 0$.

Free entropy:

$$\bar{f}_n^{\circ} = \frac{1}{n} \mathbb{E} \log \int \prod_{\mu=1}^n P_{\text{out}} \left(Y_{\mu} \mid \rho \frac{\mathbf{v}^{\top} \mathbf{X}_{\mu}}{\sqrt{d}} + \sqrt{\epsilon} \xi_{\mu} \right) D\mathbf{v} D\xi$$

Mutual information:

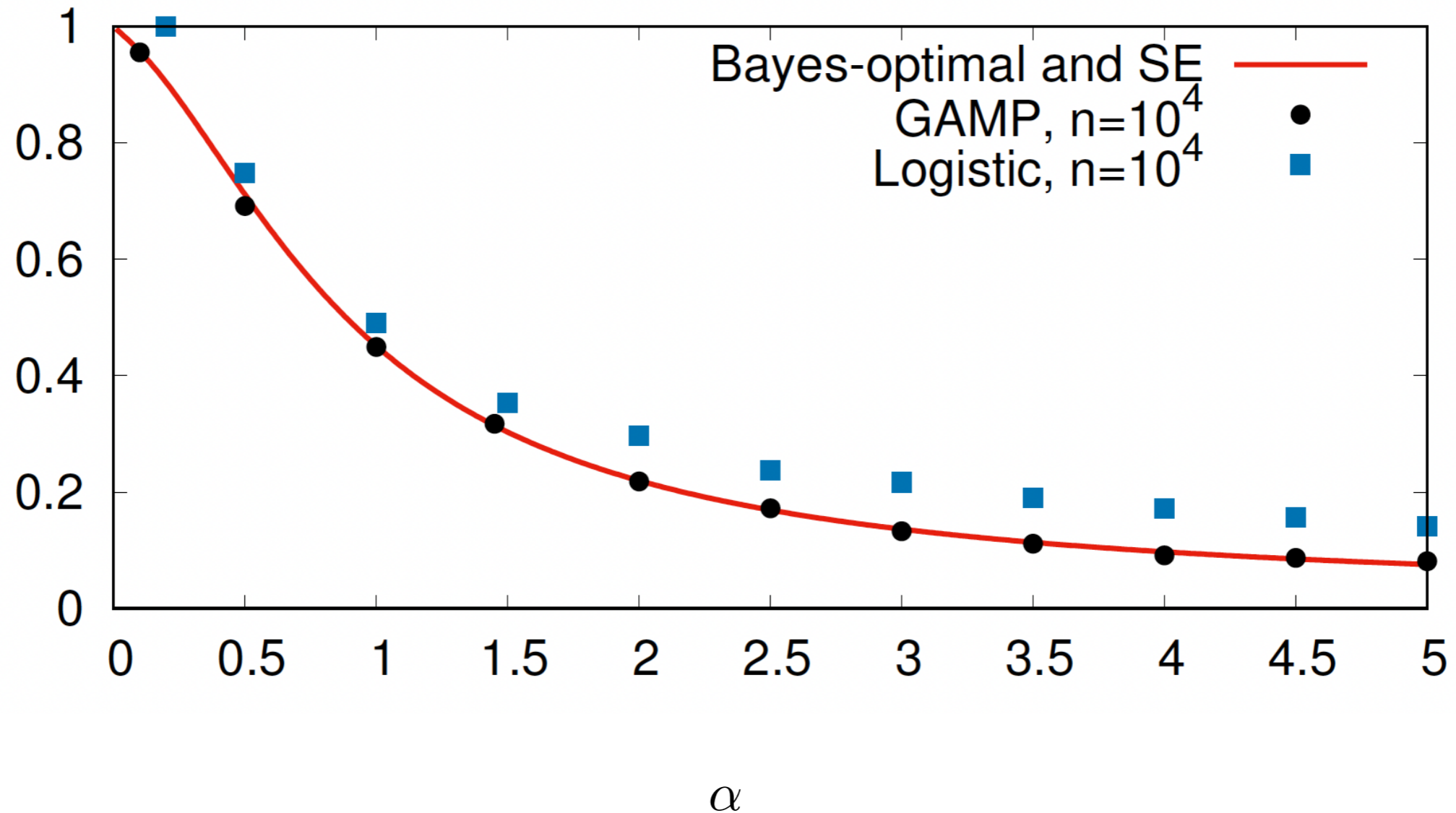
$$\frac{1}{n} I_n^{\circ}(\mathbf{v}^*, \boldsymbol{\xi}^*; \mathcal{D}_n^{\circ}) = -\bar{f}_n^{\circ} + \mathbb{E} \log P_{\text{out}} \left(Y_1 \mid \rho \frac{\mathbf{v}^{*\top} \mathbf{X}_1}{\sqrt{d}} + \sqrt{\epsilon} \xi_1^* \right).$$

Barbier, Miolane, Macris, Krzakala, Zdeborova PNAS 19'

Relevant scalings in the GLM

In the GLM the relevant scaling is $\alpha = \frac{n}{d} = O(1)$, and the free entropy (and MI) are given by an RS formula.

The $\mathbb{E}\mathcal{E}$ is related through a derivative to the MI.



Results

Free entropy equivalence

Recall

$$Y_\mu \sim P_{\text{out}}\left(\cdot \mid \frac{\mathbf{a}^{*\top}}{\sqrt{\rho}} \varphi\left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}\right)\right), \quad Y_\mu^\circ \sim P_{\text{out}}\left(\cdot \mid \rho \frac{\mathbf{v}^{*\top} \mathbf{X}_\mu}{\sqrt{d}} + \sqrt{\epsilon} \xi_\mu^*\right)$$

Theorem (Barbier, Camilli, Tiepova 23')

Let $\rho := \mathbb{E}_{\mathcal{N}(0,1)} \varphi'$ and $\epsilon^2 := \mathbb{E}_{\mathcal{N}(0,1)} \varphi^2 - \rho^2$. Assume

A1) $\varphi \in C^3(\mathbb{R}, \mathbb{R})$ is odd, $|\varphi'|, |\varphi''|, |\varphi'''| \leq \bar{K}$

A2) $|f|, |f'|, |f''| \leq \bar{K}$ P_A -almost surely.

Then:

$$|\bar{f}_n - \bar{f}_n^\circ| = O\left(\sqrt{\left(1 + \frac{n}{d}\right) \left(\frac{n}{\rho} + \frac{n}{d^{3/2}} + \frac{1}{\sqrt{d}}\right)}\right).$$

Mutual information

Corollary

Under the same hypothesis, the following holds:

$$\left| \frac{1}{n} I_n(\mathbf{a}^*, \mathbf{W}^*; \mathcal{D}_n) - \frac{1}{n} I_n^\circ(\mathbf{v}^*, \boldsymbol{\xi}^*; \mathcal{D}_n^\circ) \right| = O\left(\sqrt{\left(1 + \frac{n}{d}\right) \left(\frac{n}{p} + \frac{n}{d^{3/2}} + \frac{1}{\sqrt{d}}\right)}\right).$$

We can thus identify a scaling in which we expect the two-layer NN to collapse into the equivalent GLM:

$$\widetilde{\lim} \equiv \lim_{n,p,d \rightarrow \infty} \text{ s.t. } \left(1 + \frac{n}{d}\right) \left(\frac{n}{p} + \frac{n}{d^{3/2}} + \frac{1}{\sqrt{d}}\right) \rightarrow 0.$$

Mutual information

Corollary

Under the same hypothesis, the following holds:

$$\left| \frac{1}{n} I_n(\mathbf{a}^*, \mathbf{W}^*; \mathcal{D}_n) - \frac{1}{n} I_n^\circ(\mathbf{v}^*, \boldsymbol{\xi}^*; \mathcal{D}_n^\circ) \right| = O\left(\sqrt{\left(1 + \frac{n}{d}\right) \left(\frac{n}{p} + \frac{n}{d^{3/2}} + \frac{1}{\sqrt{d}}\right)}\right).$$

We can thus identify a scaling in which we expect the two-layer NN to collapse into the equivalent GLM:

$$\widetilde{\lim} \equiv \lim_{n,p,d \rightarrow \infty} \text{ s.t. } \left(1 + \frac{n}{d}\right) \left(\frac{n}{p} + \frac{n}{d^{3/2}} + \frac{1}{\sqrt{d}}\right) \rightarrow 0.$$

$p \sim d$ allowed! as long as $p \gg n$

Generalization Error

$$\widetilde{\lim} \equiv \lim_{n,p,d \rightarrow \infty} \text{ s.t. } \left(1 + \frac{n}{d}\right) \left(\frac{n}{p} + \frac{n}{d^{3/2}} + \frac{1}{\sqrt{d}}\right) \rightarrow 0.$$

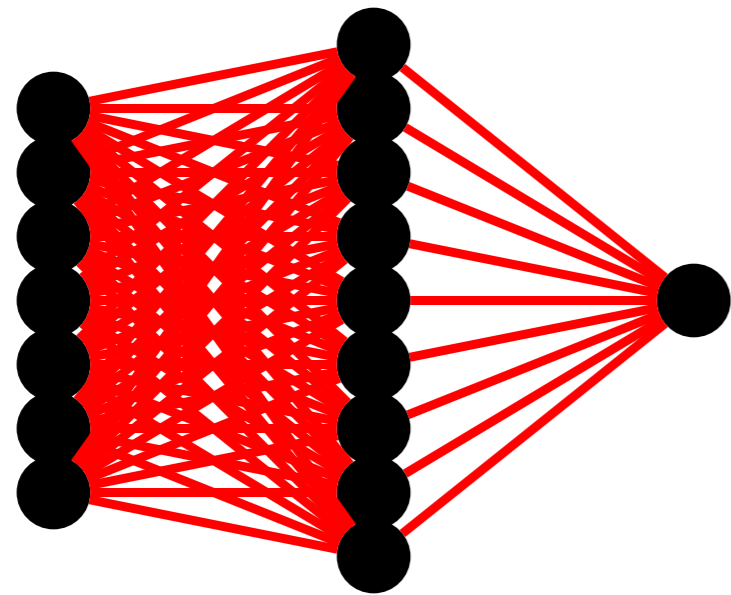
Corollary

Under the same hypothesis, the following holds:

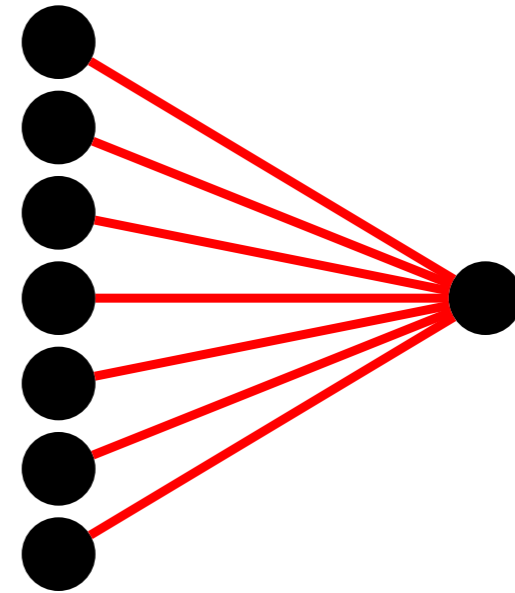
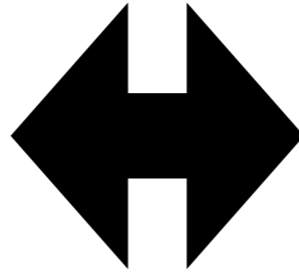
$$\widetilde{\lim} |\mathbb{E}\mathcal{E} - \mathbb{E}\mathcal{E}^\circ| = 0,$$

i.e. the 2-layer NN and the GLM share the same generalization error in the $\widetilde{\lim}$.

$p \sim d$ allowed! as long as $p \gg n$



Info. Theoretic



$$Y_{\mu} \sim P_{\text{out}} \left(\cdot \mid \frac{\mathbf{a}^{*\top}}{\sqrt{\rho}} \varphi \left(\frac{\mathbf{W}^* \mathbf{X}_{\mu}}{\sqrt{d}} \right) \right)$$

$$Y_{\mu}^{\circ} \sim P_{\text{out}} \left(\cdot \mid \rho \frac{\mathbf{v}^{*\top} \mathbf{X}_{\mu}}{\sqrt{d}} + \sqrt{\epsilon} \xi_{\mu}^* \right)$$

$$\mathcal{D}_n = \{(\mathbf{X}_{\mu}, Y_{\mu})_{\mu=1}^n\}$$

$$\mathcal{D}_n^{\circ} = \{(\mathbf{X}_{\mu}, Y_{\mu}^{\circ})_{\mu=1}^n\}$$

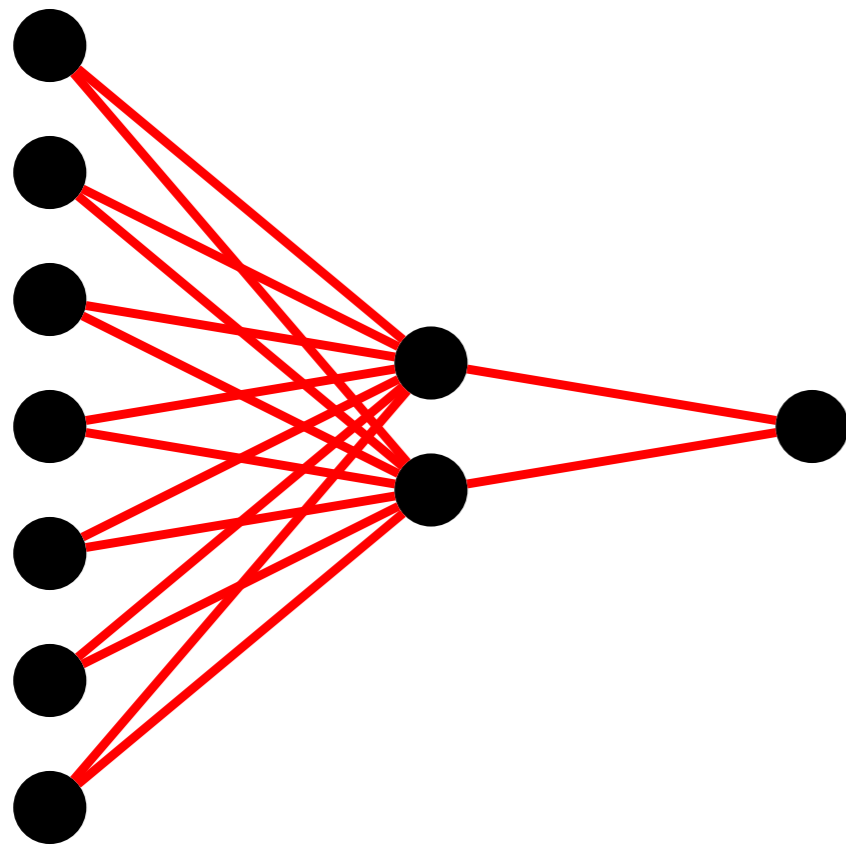
Remark

Our results **do not** imply that training a two-layer NN on \mathcal{D}_n° yields BO generalization error. (or vice versa)

Various scalings

Committee Machines

$$\widetilde{\text{lim}} \equiv \lim_{n,p,d \rightarrow \infty} \text{s.t.} \quad \left(1 + \frac{n}{d}\right) \left(\frac{n}{p} + \frac{n}{d^{3/2}} + \frac{1}{\sqrt{d}}\right) \rightarrow 0.$$



$$p = O(1), \quad \frac{n}{d} = O(1)$$

Warning!

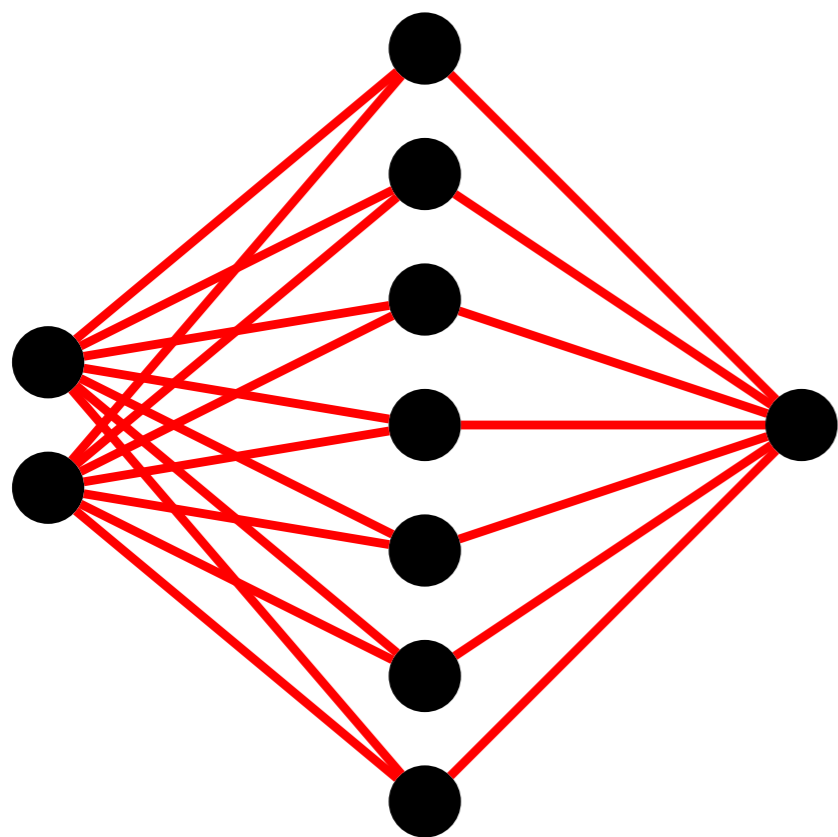
Only $n, d \rightarrow \infty$. Different from the
our scaling in $\widetilde{\text{lim}}$

If the data are really high-dimensional, the middle layer compresses them too much.

Aubin, Maillard, Barbier, Macris, Krzakala, Zdeborova NeurIPS 18'

Mean field regime

$$\widetilde{\lim} \equiv \lim_{n,p,d \rightarrow \infty} \text{ s.t. } \left(1 + \frac{n}{d}\right) \left(\frac{n}{p} + \frac{n}{d^{3/2}} + \frac{1}{\sqrt{d}}\right) \rightarrow 0.$$



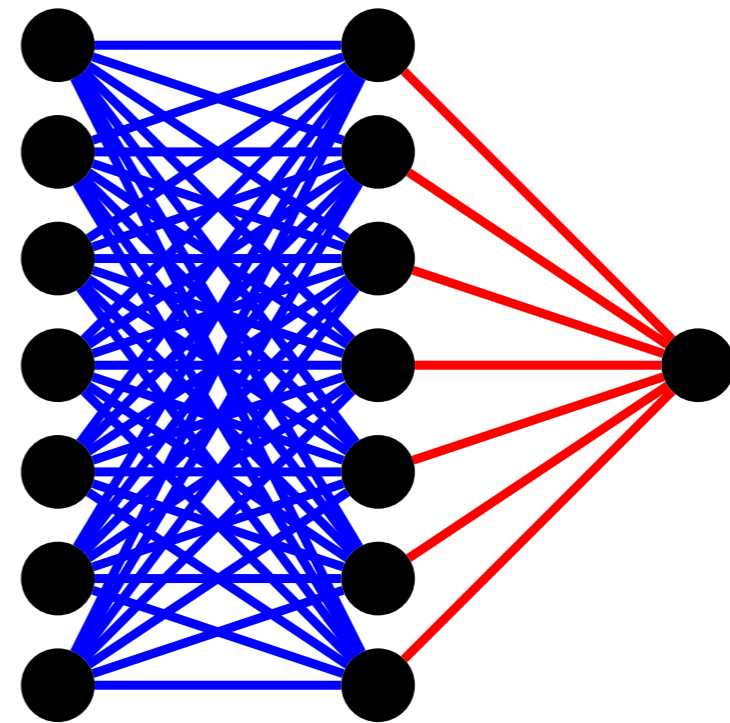
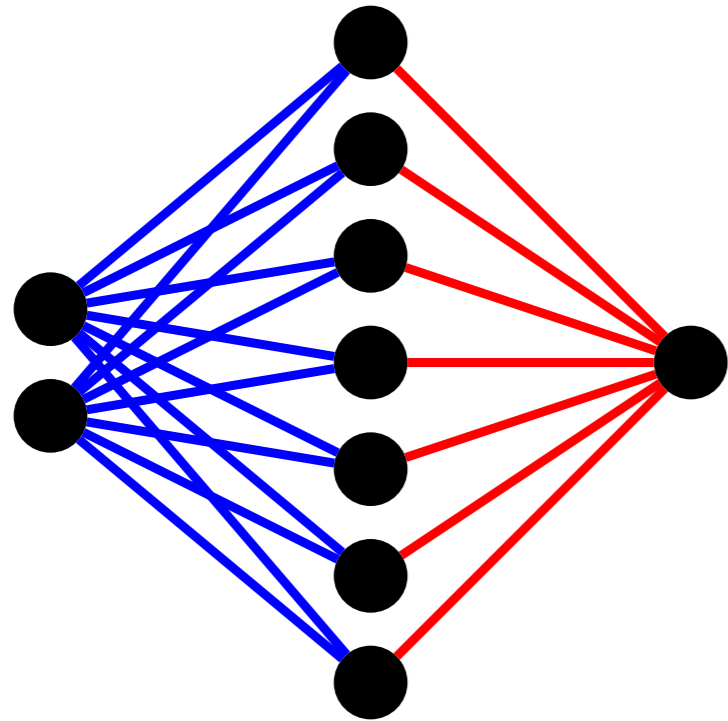
In the hypothesis $p \gg d$, one can track the empirical distribution of the weights of the network by means of a distributional equation **[Mei-Montanari-Nguyen-2008]**, regardless of n .
Chizat, Bach; ...

Remark

In our proof and in **[Mei-Montanari-Nguyen-2008]** p seems to play a crucial role!

Frozen hidden weights

Blue weights are **quenched**, they vary negligibly w.r.t. red weights, that are learnable or **annealed**.



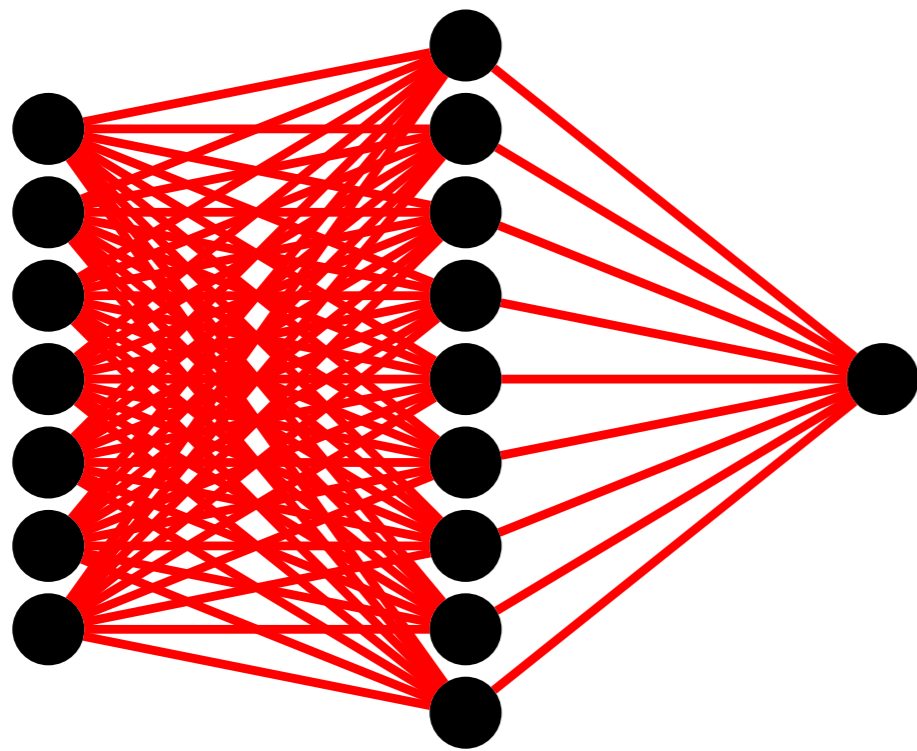
- SGD in MF at initial stage
- NTK Jacot, Hongler

- Random features models Rahmini et al
- NNs as Gaussian processes Neal

Parameter counting

Only p parameters to be learnt, w.r.t. $dp + p$ in our case.

Recent conjectures



Even for deeper networks.

Linear scalings

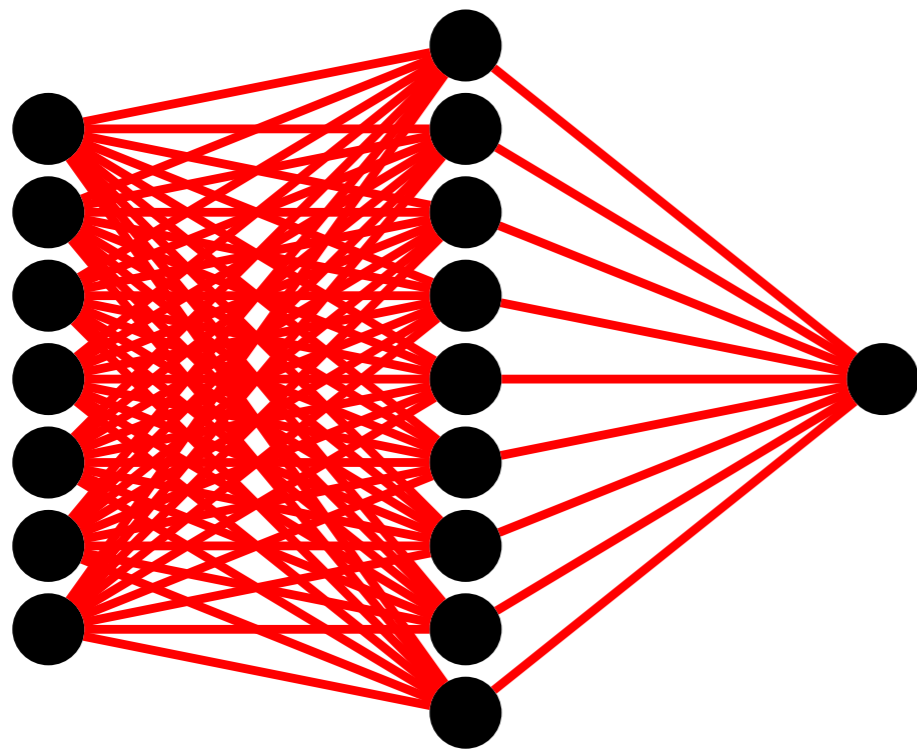
$$n, d, p \rightarrow \infty, \quad \frac{n}{p} = O(1) = \frac{d}{p}$$

$\tilde{\text{lim}}$ does **not** catch these scalings.

We need $p \gg n$

- [Li-Sampolinsky-2021] studied full training for **linear** networks;
- [Ariosto-Pacelli-Gherardi-Rotondo-2022] conjectures a formula for the ERM - generalization error;
- [Cui-Krzakala-Zdeborová-2023] builds on a **Gaussian Equivalence Principle** to compute the Bayes-optimal limits as we do

Recent conjectures



Even for deeper networks.

Linear scalings

$$n, d, p \rightarrow \infty, \quad \frac{n}{p} = O(1) = \frac{d}{p}$$

$\tilde{\text{lim}}$ does **not** catch these scalings.

We need $p \gg n$

- [Li-Sampolinsky-2021] studied full training for **linear** networks;
- [Ariosto-Pacelli-Gherardi-Rotondo-2022] conjectures a formula for the ERM - generalization error;
- [Cui-Krzakala-Zdeborová-2023] builds on a **Gaussian Equivalence Principle** to compute the Bayes-optimal limits as we do

Q: is our proof sub-optimal or is the replica prediction of Cui et al. in the linear scaling not exactly correct?

Proof idea

Universality: Gaussian Equivalence Principles

The most annoying part in the analysis is for sure the non-linearity in the middle layer. To get rid of it, some authors have noticed that

GEP (informal)

It amounts to the following replacement:

$$\varphi\left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}\right) \approx \rho \frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}} + \sqrt{\epsilon} \xi_\mu^*$$

with ξ_μ^* an independent standard Gaussian noise and

$$\rho = \mathbb{E}_{\mathcal{N}(0,1)} \varphi', \quad \epsilon = \mathbb{E}_{\mathcal{N}(0,1)} \varphi^2 - (\mathbb{E}_{\mathcal{N}(0,1)} \varphi')^2$$

In our setting, it is not clear to what extent this is applicable!

Interpolation

The interpolation has to keep all the ingredients together:

$$S_{t\mu} := \sqrt{1-t} \frac{\mathbf{a}^{*\top}}{\sqrt{p}} \varphi\left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}\right) + \sqrt{t} \rho \frac{\mathbf{v}^{*\top} \mathbf{X}_\mu}{\sqrt{d}} + \sqrt{t} \epsilon \xi_\mu^*,$$
$$s_{t\mu} := \sqrt{1-t} \frac{\mathbf{a}^\top}{\sqrt{p}} \varphi\left(\frac{\mathbf{W} \mathbf{X}_\mu}{\sqrt{d}}\right) + \sqrt{t} \rho \frac{\mathbf{v}^\top \mathbf{X}_\mu}{\sqrt{d}} + \sqrt{t} \epsilon \xi_\mu.$$

Interpolating dataset:

$$\mathcal{D}_{n,t} = \{(\mathbf{X}_\mu, Y_{t\mu})_{\mu=1}^n\}, \quad Y_{t\mu} \sim P_{\text{out}}(\cdot | S_{t\mu})$$

Interpolating free entropy:

$$\bar{f}_n(t) = \frac{1}{n} \mathbb{E}_{\mathcal{D}_{n,t}} \log \int D\mathbf{v} D\xi D\mathbf{a} D\mathbf{W} \prod_{\mu=1}^n P_{\text{out}}(Y_{t\mu} | s_{t\mu})$$

Interpolation

Interpolating free entropy:

$$\bar{f}_n(t) = \frac{1}{n} \mathbb{E}_{\mathcal{D}_{n,t}} \log \underbrace{\int D\mathbf{v} D\xi D\mathbf{a} D\mathbf{W} \prod_{\mu=1}^n P_{\text{out}}(Y_{t\mu} \mid s_{t\mu}(\mathbf{a}, \mathbf{W}, \mathbf{v}, \xi_{\mu}))}_{Z_t}$$

Main Goal

Prove that

$$\frac{d}{dt} \bar{f}_n(t) = O\left(\sqrt{\left(1 + \frac{n}{d}\right) \left(\frac{n}{p} + \frac{n}{d^{3/2}} + \frac{1}{\sqrt{d}}\right)}\right),$$

uniformly in $t \in [0, 1]$.

Cross terms

$$\frac{d}{dt} \bar{f}_n(t) = -A_1 + A_2 + A_3 + B$$

Let $u_y(x) = \log P_{\text{out}}(y | x)$, then

$$A_1 := \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu=1}^n u'_{Y_{t\mu}}(S_{t\mu}) \frac{\mathbf{a}^{*\top}}{\sqrt{(1-t)\rho}} \varphi\left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}\right),$$

$$A_2 := \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu=1}^n u'_{Y_{t\mu}}(S_{t\mu}) \rho \frac{\mathbf{v}^{*\top} \mathbf{X}_\mu}{\sqrt{td}},$$

$$A_3 := \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu=1}^n u'_{Y_{t\mu}}(S_{t\mu}) \sqrt{\frac{\epsilon}{t}} \xi_\mu^*,$$

$$B := \frac{1}{n} \mathbb{E}_{(t)} \left\langle \sum_{\mu=1}^n u'_{Y_{t\mu}}(s_{t\mu}) \frac{ds_{t\mu}}{dt} \right\rangle.$$

Let $u_y(x) = \log P_{\text{out}}(y | x)$, then

$$\begin{aligned} & \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu=1}^n u'_{Y_{t\mu}}(S_{t\mu}) \left[\frac{\mathbf{a}^{*\top}}{\sqrt{(1-t)p}} \varphi\left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}\right) - \rho \frac{\mathbf{a}^{*\top} \mathbf{W}^* \mathbf{X}_\mu}{\sqrt{pd}} \right] + \\ & + \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu=1}^n u'_{Y_{t\mu}}(S_{t\mu}) \rho \frac{\mathbf{a}^{*\top} \mathbf{W}^* \mathbf{X}_\mu}{\sqrt{pd}} \end{aligned}$$

Let $u_y(x) = \log P_{\text{out}}(y | x)$, then

$$\begin{aligned} & \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu=1}^n u'_{Y_{t\mu}}(S_{t\mu}) \left[\frac{\mathbf{a}^{*\top}}{\sqrt{(1-t)p}} \varphi\left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}\right) - \rho \frac{\mathbf{a}^{*\top} \mathbf{W}^* \mathbf{X}_\mu}{\sqrt{pd}} \right] + \\ & + \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu=1}^n u'_{Y_{t\mu}}(S_{t\mu}) \rho \frac{\mathbf{a}^{*\top} \mathbf{W}^* \mathbf{X}_\mu}{\sqrt{pd}} \end{aligned}$$

Gaussian integration by part (here w.r.t. \mathbf{a}^*)

$$= \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu, \nu=1}^n U_{\mu\nu} \left[\frac{\varphi(\boldsymbol{\alpha}_\mu)^\top \varphi(\boldsymbol{\alpha}_\nu) - \rho \boldsymbol{\alpha}_\mu^\top \varphi(\boldsymbol{\alpha}_\nu)}{p} \right]$$

$$\boldsymbol{\alpha}_\mu := \frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}$$

Let $u_y(x) = \log P_{\text{out}}(y | x)$, then

$$\begin{aligned} & \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu=1}^n u'_{Y_{t\mu}}(S_{t\mu}) \left[\frac{\mathbf{a}^{*\top}}{\sqrt{(1-t)p}} \varphi\left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}\right) - \rho \frac{\mathbf{a}^{*\top} \mathbf{W}^* \mathbf{X}_\mu}{\sqrt{pd}} \right] + \\ & + \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu=1}^n u'_{Y_{t\mu}}(S_{t\mu}) \rho \frac{\mathbf{a}^{*\top} \mathbf{W}^* \mathbf{X}_\mu}{\sqrt{pd}} \end{aligned}$$

Gaussian integration by part (here w.r.t. \mathbf{a}^*)

$$= \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu, \nu=1}^n U_{\mu\nu} \left[\frac{\varphi(\boldsymbol{\alpha}_\mu)^\top \varphi(\boldsymbol{\alpha}_\nu) - \rho \boldsymbol{\alpha}_\mu^\top \varphi(\boldsymbol{\alpha}_\nu)}{p} \right] \quad \boldsymbol{\alpha}_\mu := \frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}$$

Getting rid of non-linearities by approximation

$$\mathbb{E}_{\mathbf{W}^*} \varphi'(\boldsymbol{\alpha}_{\mu i}) = \rho + O\left(\frac{\|\mathbf{X}_\mu\|^2}{d} - 1\right),$$

$$\mathbb{E}_{\mathbf{W}^*} \varphi^2(\boldsymbol{\alpha}_{\mu i}) = \mathbb{E}_{\mathcal{N}(0,1)} \varphi^2 + O\left(\frac{\|\mathbf{X}_\mu\|^2}{d} - 1\right),$$

$$\mathbb{E}_{\mathbf{W}^*} \varphi(\boldsymbol{\alpha}_{\mu i}) \varphi(\boldsymbol{\alpha}_{\nu i}) = \mathbb{E}_{\mathcal{N}(0,1)} \varphi' \frac{\mathbf{X}_\mu^\top \mathbf{X}_\nu}{d} + O\left(\frac{\mathbf{X}_\mu^\top \mathbf{X}_\nu}{d} \left(\frac{\|\mathbf{X}_\mu\|^2}{d} - 1\right)\right) + O\left(\left(\frac{\mathbf{X}_\mu^\top \mathbf{X}_\nu}{\|\mathbf{X}_\nu\|^2}\right)^2\right) + O\left(\frac{(\mathbf{X}_\mu^\top \mathbf{X}_\nu)^2}{\|\mathbf{X}_\nu\|^2 d}\right),$$

Let $u_y(x) = \log P_{\text{out}}(y | x)$, then

$$\begin{aligned} & \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu=1}^n u'_{Y_{t\mu}}(S_{t\mu}) \left[\frac{\mathbf{a}^{*\top}}{\sqrt{(1-t)p}} \varphi\left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}\right) - \rho \frac{\mathbf{a}^{*\top} \mathbf{W}^* \mathbf{X}_\mu}{\sqrt{pd}} \right] + \\ & + \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu=1}^n u'_{Y_{t\mu}}(S_{t\mu}) \rho \frac{\mathbf{a}^{*\top} \mathbf{W}^* \mathbf{X}_\mu}{\sqrt{pd}} \end{aligned}$$

Gaussian integration by part (here w.r.t. \mathbf{a}^*)

$$= \frac{1}{2n} \mathbb{E}_{(t)} \log \mathcal{Z}_t \sum_{\mu, \nu=1}^n U_{\mu\nu} \left[\frac{\varphi(\boldsymbol{\alpha}_\mu)^\top \varphi(\boldsymbol{\alpha}_\nu) - \rho \boldsymbol{\alpha}_\mu^\top \varphi(\boldsymbol{\alpha}_\nu)}{p} \right] \quad \boldsymbol{\alpha}_\mu := \frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}$$

Getting rid of non-linearities by approximation

$$\mathbb{E}_{\mathbf{W}^*} \varphi'(\alpha_{\mu i}) = \rho + O\left(\frac{\|\mathbf{X}_\mu\|^2}{d} - 1\right),$$

$$\mathbb{E}_{\mathbf{W}^*} \varphi^2(\alpha_{\mu i}) = \mathbb{E}_{\mathcal{N}(0,1)} \varphi^2 + O\left(\frac{\|\mathbf{X}_\mu\|^2}{d} - 1\right),$$

$$\mathbb{E}_{\mathbf{W}^*} \varphi(\alpha_{\mu i}) \varphi(\alpha_{\nu i}) = \mathbb{E}_{\mathcal{N}(0,1)} \varphi' \frac{\mathbf{X}_\mu^\top \mathbf{X}_\nu}{d} + O\left(\frac{\mathbf{X}_\mu^\top \mathbf{X}_\nu}{d} \left(\frac{\|\mathbf{X}_\mu\|^2}{d} - 1\right)\right) + O\left(\left(\frac{\mathbf{X}_\mu^\top \mathbf{X}_\nu}{\|\mathbf{X}_\nu\|^2}\right)^2\right) + O\left(\frac{(\mathbf{X}_\mu^\top \mathbf{X}_\nu)^2}{\|\mathbf{X}_\nu\|^2 d}\right),$$

It is possible that the replica prediction is a good approximation in the proportional regimes, but only exact when inputs are orthogonal (while gaussians have a weak overlap)

$$\frac{1}{2} \mathbb{E}_{(t)} \frac{\overbrace{(\log \mathcal{Z}_t - \mathbb{E}_{(t)} \log \mathcal{Z}_t)}^{o\left(\sqrt{\frac{1}{n} + \frac{1}{d}}\right)}}{n} \sum_{\mu, \nu=1}^n U_{\mu\nu} \left[\frac{1}{p} \varphi\left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}\right)^\top \varphi\left(\frac{\mathbf{W}^* \mathbf{X}_\nu}{\sqrt{d}}\right) - \frac{\rho}{p} \varphi\left(\frac{\mathbf{W}^* \mathbf{X}_\mu}{\sqrt{d}}\right)^\top \frac{\mathbf{W}^* \mathbf{X}_\nu}{\sqrt{d}} \right]$$

Finally...

$$\sum_{\mu, \nu} U_{\mu\nu} [\text{RED} - \text{BLUE}]_{\mu\nu} \sim O\left(n \sqrt{\frac{1}{p} + \frac{1}{d^{3/2}}}\right)$$

Conclusion

Conclusion and perspectives (Work in progress)

- Well... what about $n/p = O(1)$? Are we actually able to reach physicists' conjectured regimes?
- Say we have partial information on \mathbf{W}^* through a Gaussian channel:
$$\tilde{\mathbf{W}} = \sqrt{\sigma} \mathbf{W}^* + \mathbf{Z}.$$
 - When $\sigma \rightarrow \infty \Rightarrow$ RF model in BO setting
 - When $\sigma = 0 \Rightarrow$ our setting

We could interpolate between the two!

- Why not more than two layers?
 - We only miss concentration of the free entropy, then an inductive argument and concentration by parts should yield the result in a similar scaling.
- What happens if we add structure to the data? $\mathbf{X} \sim \mathcal{N}(0, \Sigma)$ or mixtures.

(Some) References

- **Song Mei, Andrea Montanari and Phan-Minh Nguyen**, “*A mean field view of the landscape of two-layer neural networks*”, Proceedings of the National Academy of Sciences, v. 115, 2008;
- **Qianyi Li and Haim Sompolinsky**, “*Statistical Mechanics of Deep Linear Neural Networks: The Backpropagating Kernel Renormalization*”, Phys. Rev. X, v. 11, 2021;
- **S. Ariosto, R. Pacelli, M. Gherardi and P. Rotondo**, “*Statistical mechanics of deep learning beyond the infinite-width limit*”, arXiv:2209.04882, 2022;
- **Hugo Cui, Florent Krzakala and Lenka Zdeborová**, “*Optimal Learning of Deep Random Networks of Extensive-width*”, arXiv:2302.00375, 2023;
- **Hong Hu and Yue M Lu**, “*Universality laws for high-dimensional learning with random features*”, IEEE Transactions on Information Theory, 2022;
- **Sebastian Goldt, Bruno Loureiro, Galen Reeves, Florent Krzakala, Marc Mézard and Lenka Zdeborová**, “*The Gaussian equivalence of generative models for learning with shallow neural networks*”, Mathematical and Scientific Machine Learning PMLR, 2022.