

# Dense Hebbian neural networks: a replica symmetric picture of unsupervised learning

---

**Elena Agliari<sup>1, a</sup>, Linda Albanese<sup>b, f, g</sup>, Francesco Alemanno<sup>b, f</sup>, Andrea Alessandrelli<sup>c</sup>, Adriano Barra<sup>b, f</sup>, Fosca Giannotti<sup>d, e</sup>, Daniele Lotito<sup>c, f</sup>, Dino Pedreschi<sup>c</sup>.**

<sup>a</sup>*Dipartimento di Matematica, Sapienza Università di Roma, Piazzale Aldo Moro, 5, 00185, Roma, Italy*

<sup>b</sup>*Dipartimento di Matematica e Fisica, Università del Salento, Via per Arnesano, 73100, Lecce, Italy*

<sup>c</sup>*Dipartimento di Informatica, Università di Pisa, Lungarno Antonio Pacinotti, 43, 56126, Pisa, Italy*

<sup>d</sup>*Scuola Normale Superiore, Piazza dei Cavalieri 7, 56126, Pisa, Italy*

<sup>e</sup>*Istituto di Scienza e Tecnologie dell' Informazione, Via Giuseppe Moruzzi, 1, 56124 Pisa, Italy*

<sup>f</sup>*Istituto Nazionale di Fisica Nucleare, Campus Ecotekne, Via Monteroni, 73100, Lecce, Italy*

<sup>g</sup>*Scuola Superiore ISUFI, Campus Ecotekne, Via Monteroni, 73100, Lecce, Italy*

**ABSTRACT:** We consider dense, associative neural-networks trained with no supervision and we investigate their computational capabilities analytically, via statistical-mechanics tools, and numerically, via Monte Carlo simulations. In particular, we obtain a phase diagram summarizing their performance as a function of the control parameters (e.g. quality and quantity of the training dataset, network storage, noise) that is valid in the limit of large network size and structureless datasets. Moreover, we establish a bridge between macroscopic observables standardly used in statistical mechanics and loss functions typically used in the machine learning.

As technical remarks, from the analytical side, we extend Guerra's interpolation to tackle the non-Gaussian distributions involved in the post-synaptic potentials while, from the computational counterpart, we insert Plefka's approximation in the Monte Carlo scheme, to speed up the evaluation of the synaptic tensor, overall obtaining a novel and broad approach to investigate unsupervised learning in neural networks, beyond the shallow limit.

---

<sup>1</sup>Given her role as Editor of this journal, EA had no involvement in the peer-review of articles for which she was an author and had no access to information regarding their peer-review. Full responsibility for the peer-review process for this article was delegated to another Editor.

---

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Prelude: from Hebbian storing to Hebbian learning</b>	<b>3</b>
<b>3</b>	<b>Dense Hebbian neural networks in the unsupervised setting</b>	<b>5</b>
<b>4</b>	<b>Cost and Loss functions</b>	<b>10</b>
<b>5</b>	<b>Analytical findings</b>	<b>12</b>
5.1	High-load regime	14
5.2	Low-load regime	15
5.3	Additive noise in low-load regime	16
5.4	Low-entropy datasets in the high-load regime	17
<b>6</b>	<b>Numerical findings</b>	<b>18</b>
6.1	Stability analysis and Monte Carlo simulations	18
6.2	Critical load and bounds for the dataset size	22
<b>7</b>	<b>Conclusion and outlooks</b>	<b>24</b>
<b>A</b>	<b>Proof of Proposition 1</b>	<b>25</b>
<b>B</b>	<b>Plefka’s expansion of the effective Gibbs potential</b>	<b>29</b>
B.1	Plefka’s effective dynamics for Hebbian storing	30
B.2	Plefka’s effective dynamics for unsupervised Hebbian learning	32
<b>C</b>	<b>Evaluation of the momenta of the effective post-synaptic potential</b>	<b>33</b>
<b>D</b>	<b>List of symbols (in alphabetical order)</b>	<b>35</b>

---

## 1 Introduction

Since the seminal work on “Hebbian Learning” by John Hopfield in 1982 [35] and the paradigmatic discoveries on the landscape of spin-glasses by Giorgio Parisi around the 1980 [37], statistical mechanics has run as a theoretical tool to explain the emergent properties of large assemblies of neurons. The subsequent investigations by Amit, Gutfreund and Sompolinsky [17] have definitively established statistical-mechanics as a leading discipline for theoretical studying the collective properties of systems of neurons. In fact, in the last decades, the scientific literature has hosted a large number of contributions on *statistical mechanics of neural networks* (see e.g. [8, 26, 46]), including countless variations on the original Hopfield model (see e.g. [2, 5, 7, 13, 14, 19, 20, 30, 31]). Among these, Hebbian networks with higher-order interactions, also referred to as *dense neural networks* or *P-spin Hopfield models*,

were promptly introduced already in the 80’s by theoretical physicists (see e.g. [21, 32]) as well as by computer scientists (see e.g. [40]).

In the last few years, a renewed interest has raised for dense networks, in fact, on the one hand they have been shown to display intriguing properties as for applications (e.g., they turned out to be robust against adversarial attacks [36], they can perform pattern recognition at prohibitive signal-to-noise level [4]) and, on the other hand, many technical issues concerning their analytical and numerical investigation still deserve attention [18, 43, 44]. Further, recent advances in the usage of pair-wise Hebbian-like networks for learning tasks [6, 15] could be fruitfully extended to the higher-order case. In this work we aim to address these points.

More precisely, as for the analytical investigation, we apply interpolating techniques (see e.g., [5, 34]) and extend their range of applicability to include the challenging case of non-Gaussian local fields as it happens for dense networks. As for the numerical investigation, we propose a strategy based on Plefka expansion [38, 39] to overcome the strong difficulties implied by the update of a giant synaptic tensor in Monte Carlo (MC) simulations with a remarkable speed up.

Concerning the learning tasks, it should be recalled that, despite its name, “Hebbian learning” (in its standard meaning in Statistical Mechanics, that is provided by the Amit-Gutfreund-Sompolinsky theory of the Hopfield model [16]) is actually a *storing* rule and there is no training dataset or inference activity underlying. Yet, recently, the Hebbian rule has been shown to be recastable into a genuine learning rule allowing for both supervised and unsupervised modes and the resulting pair-wise neural network is feasible for a full statistical-mechanics investigation [6, 15]. In this work, we extend this framework to dense neural networks focusing on the unsupervised setting and referring to [1] for the supervised one; the analytical investigations are led under the replica-symmetry (RS) approximation and for structureless datasets. There are several results stemming from this study.

First, from an analytical perspective, we are able to summarize the behavior of the network into a phase diagram, namely to highlight in the space of the control parameters of the network (e.g., noise, storage, training-set size, training-set quality) the existence of different regions corresponding to different computational skills shown by the network. As we will deepen, this is a major reward of the statistical mechanics approach and yields pivotal information towards a *sustainable* artificial intelligence since its knowledge allows *a priori* setting the machine parameters in the best configuration for a given task. For instance, in this context we can assess the minimal size of the dataset, as function of the dataset quality and the amount of patterns to store, necessary for a successful training of the machine and we can also estimate the largest amount of information that the machine can safely handle.

Second, from a computational perspective, we inserted effective Plefka dynamics in Monte Carlo simulation to obtain a significant speed up for the update of the synaptic tensor (whose handling is otherwise cumbersome). This allowed us to test analytical prediction from the random theory confirming that these networks -if suitably trained- enjoy pattern recognition capabilities remarkably robust w.r.t. vast amount of noise (as compared to their shallow counterpart) as well as a supra-linear storage of patterns. However, nothing comes for free and the price to pay to enjoy these enhanced information processing capabilities lies in the huge amount of examples that the network has to experience before a correct learning can take place. Thresholds for learning and maximal storage capabilities also have all been confirmed numerically.

Finally, more of conceptual interest, we link quantifiers for machine retrieval (e.g., cost function and Mattis magnetization) to quantifiers for machine learning (e.g., loss function), making these two capabilities of neural networks -learning and retrieval- two aspects of a unified cognitive process and yielding a cross-fertilization between the two related fields (i.e., statistical mechanics and machine

learning).

The paper is structured as follows: in Sec. 2 we briefly review the state-of-the-art on Hebbian learning and in Sec. 3 we introduce the unsupervised dense Hebbian networks and define the related macroscopic observables necessary for its investigation. Next, in Sec. 4, we discuss the connection between performance quantifiers stemming from, respectively, statistical-mechanics and machine learning. Then, in Sec. 5, we study the model exploiting the statistical-mechanics framework, while the numerical investigation is addressed in Sec. 6, by relying on Monte Carlo simulations. Finally, in Sec. 7 we summarize and discuss our results. Technical details are collected in the Appendices.

Finally we stress again that the present manuscript is dedicated to the study of unsupervised learning, while the supervised protocol for dense networks is addressed in a twin work [1].

## 2 Prelude: from Hebbian storing to Hebbian learning

The standard Hopfield model is built upon  $N$  binary neurons, denoted as  $\boldsymbol{\sigma} = (\sigma_1, \sigma_2, \dots, \sigma_N) \in \{-1, +1\}^N$ , that are employed to reconstruct the information encoded in  $K$  binary vectors of length  $N$ , also called *patterns* and denoted as  $\boldsymbol{\xi}^\mu \in \{-1, +1\}^N$  with  $\mu = 1, \dots, K$ , whose entries are Rademacher random variables, namely, for the generic  $(i, \mu)$  entry

$$\mathbb{P}(\xi_i^\mu) = \frac{1}{2} \left[ \delta_{\xi_i^\mu, -1} + \delta_{\xi_i^\mu, +1} \right]. \quad (2.1)$$

This information is allocated in the synaptic matrix, namely in the couplings  $\mathbf{J} = \{J_{ij}\}_{i,j=1,\dots,N}$  among neurons, defined according to Hebb's rule

$$J_{ij} = \frac{1}{N} \sum_{\mu=1}^K \xi_i^\mu \xi_j^\mu, \quad (2.2)$$

ensuring that, under suitable conditions, the system can play as an associative memory (*vide infra*). The network has a Hamiltonian cost-function representation

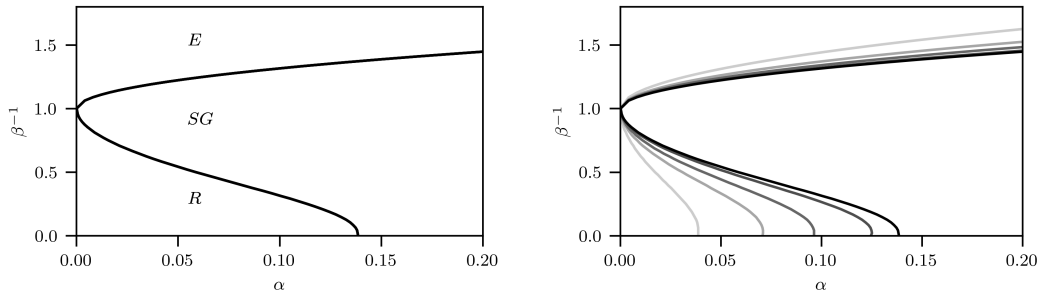
$$\mathcal{H}_{N,K}^{(H)}(\boldsymbol{\sigma}|\boldsymbol{\xi}) = - \sum_{i<j}^{N,N} J_{ij} \sigma_i \sigma_j = - \frac{N}{2} \sum_{\mu=1}^K m_\mu^2 + \frac{K}{2},$$

where in the last equivalence we used (2.2) and we introduced the Mattis magnetization

$$m_\mu := \frac{1}{N} \sum_{i=1}^N \xi_i^\mu \sigma_i. \quad (2.3)$$

The occurrence of a neural configurations  $\boldsymbol{\sigma}$  is ruled by the Boltzmann-Gibbs probability  $\propto e^{-\beta \mathcal{H}_{N,K}^{(H)}(\boldsymbol{\sigma}|\boldsymbol{\xi})}$  where  $\beta := 1/T$  tunes the degree of stochasticity and, in a physical jargon, represents the inverse of the temperature.

The relaxation to a state where  $m_\mu$  is close to 1 is interpreted as the *retrieval* of the pattern  $\boldsymbol{\xi}^\mu$ . As anticipated, the statistical-mechanical analysis allows summarizing the performance of a network into a phase diagram, which highlights the existence of qualitatively different behaviors of the system as its control parameters are tuned. Here the control parameters are the above-mentioned temperature  $T$ , also referred to as “fast noise”, and the load  $\alpha$  defined as  $\alpha := \lim_{N \rightarrow \infty} K/N$ , also referred to as “slow noise”. The phase diagram for the Hopfield model (see e.g., [17, 28]) is reported in Fig. 1 (left



**Figure 1:** Left: Phase diagram of the Hopfield model. Three regions corresponding to qualitatively different behaviors of the system are highlighted: ergodic (E, where fast noise prevails), spin-glass (SG, where slow noise prevails) and retrieval [R, where (a close neighbourhood of) each pattern  $\xi^\mu$  plays as an attractor for the neural configuration and consequently the system can perform pattern recognition as an associative memory]. In particular, in the retrieval region, if a noisy example of a pattern, say  $\eta^\mu$ , is inputted to the network (namely the neuron configuration is initialized as  $\sigma = \eta^\mu$ ), the latter reconstructs the original pattern (namely  $\sigma$  spontaneously relaxes (close) to  $\xi^\mu$ ). Right: Phase diagram of the unsupervised Hopfield model. The darker the lines, the better the quality  $r$  of the supplied dataset ( $r = 0.4, 0.5, 0.6, 1$ ). As the quality of the dataset improves, the retrieval region expands, and when the dataset quality saturates to 1 (black line) the phase diagram recovers the one in the left panel. Here  $M = 40$  and analogous results are found by retaining  $r$  constant and increasing  $M$ .

panel): one can notice that this machine is able to work as an associative memory solely in the retrieval region, corresponding to loads  $\alpha < \alpha_c \approx 0.14$ . Thus, it is pointless trying to allocate a larger amount of patterns as the machine will not be able to retrieve them: the knowledge of the phase diagram allows us to use this information *a priori*, before any trial is performed, thus potentially saving energy and CPU time<sup>1</sup>.

Despite the expression (2.2) is often named Hebbian *learning*, the above model has little to share with machine learning as there are no real learning processes underlying. These could however be introduced with minimal modification of Eq. (2.2), as we are going to explain. Let us treat each pattern  $\xi^\mu$  as an *archetype* and use it to generate a training set of  $M$  *examples* for each archetype. Denoting with  $\eta^{\mu,a} \in \{-1, +1\}^N$  the  $a$ -th example of the  $\mu$ -th archetype, we can write two generalizations of the above Hebbian rule, namely

$$J_{ij}^{(unsup)} = \frac{1}{NM} \sum_{\mu=1}^K \sum_{a=1}^M \eta_i^{\mu,a} \eta_j^{\mu,a}, \quad (2.4)$$

$$J_{ij}^{(sup)} = \frac{1}{NM^2} \sum_{\mu=1}^K \left( \sum_{a=1}^M \eta_i^{\mu,a} \right) \left( \sum_{a=1}^M \eta_j^{\mu,a} \right), \quad (2.5)$$

where, in the first expression there is no teacher that knows the labels and can cluster the examples archetype-wise as it happens in the second scenario, this is why the two generalizations are associated to, respectively, unsupervised and supervised protocols [6, 15].

<sup>1</sup>Optimized protocols in AI are especially longed for as, at present, training AI on large scale can result in conflicts with green policies [42].

In order to build our dataset  $\{\boldsymbol{\eta}^{\mu,a}\}_{a=1,\dots,M}^{\mu=1,\dots,K}$ , we generate  $M$  randomly-perturbed copies of each archetype, interpreted as *examples* and whose entries  $(i, \mu, a)$  are described by

$$\mathbb{P}(\eta_i^{\mu,a}|\xi_i^\mu) = \frac{1-r}{2}\delta_{\eta_i^{\mu,a},-\xi_i^\mu} + \frac{1+r}{2}\delta_{\eta_i^{\mu,a},\xi_i^\mu}, \quad (2.6)$$

in such a way that  $r \in [0, 1]$  assesses the training-set *quality*, that is, as  $r \rightarrow 1$  the example matches perfectly the archetype, whereas for  $r \rightarrow 0$  an example is, in the average, orthogonal to the related archetype.

A natural question is thus wondering the existence of a threshold  $M_\otimes$  beyond which the network can certainly infer the archetypes that gave rise to a newly experienced example. A schematic representation to figure out how the learning process of the archetype works in this kind of network is provided in Fig. 2.

The unsupervised, pairwise Hopfield model supplied with this kind of dataset has been investigated in details in [6, 15], obtaining a full statistical-mechanics description summarized in the phase diagram reported in Fig. 1 (right panel). Interestingly, one can see that, as the dataset is impaired (because either  $r$  or  $M$  is reduced), the retrieval region shrinks.

A useful quantity to assess the overall information content of the dataset  $\{\boldsymbol{\eta}^{\mu,a}\}_{a=1,\dots,M}^{\mu=1,\dots,K}$  is given by  $\rho = \frac{1-r^2}{Mr^2}$ , which in the following shall be referred to as *dataset entropy*. Strictly speaking,  $\rho$  is not an entropy, yet here we allow ourselves for this slight abuse of language because, as discussed in [15], the conditional entropy, that quantifies the amount of information needed to describe the original message  $\xi^\mu$  given the set of related examples  $\{\boldsymbol{\eta}^{\mu,a}\}_{a=1,\dots,M}$ , is a monotonically increasing function of  $\rho$ .

### 3 Dense Hebbian neural networks in the unsupervised setting

We consider a network with  $N$  Ising neurons  $\sigma_i \in \{-1, +1\}$  with  $i = 1, \dots, N$ ,  $K$  Rademacher archetypes  $\xi^\mu \in \{-1, +1\}^N$  and  $M$  noisy examples  $\boldsymbol{\eta}^{\mu,a} \in \{-1, +1\}^N$  per archetype  $\xi^\mu$  with  $a = 1, \dots, M$  and  $\mu = 1, \dots, K$ , whose entries are drawn according with, respectively, (2.1) and (2.6).

In the network considered here interactions among neurons are  $P$ -wise and their magnitude is obtained by generalizing (2.4), as captured by the next

**Definition 1.** *The cost-function (or Hamiltonian) of the dense Hebbian neural network in the unsupervised regime is*

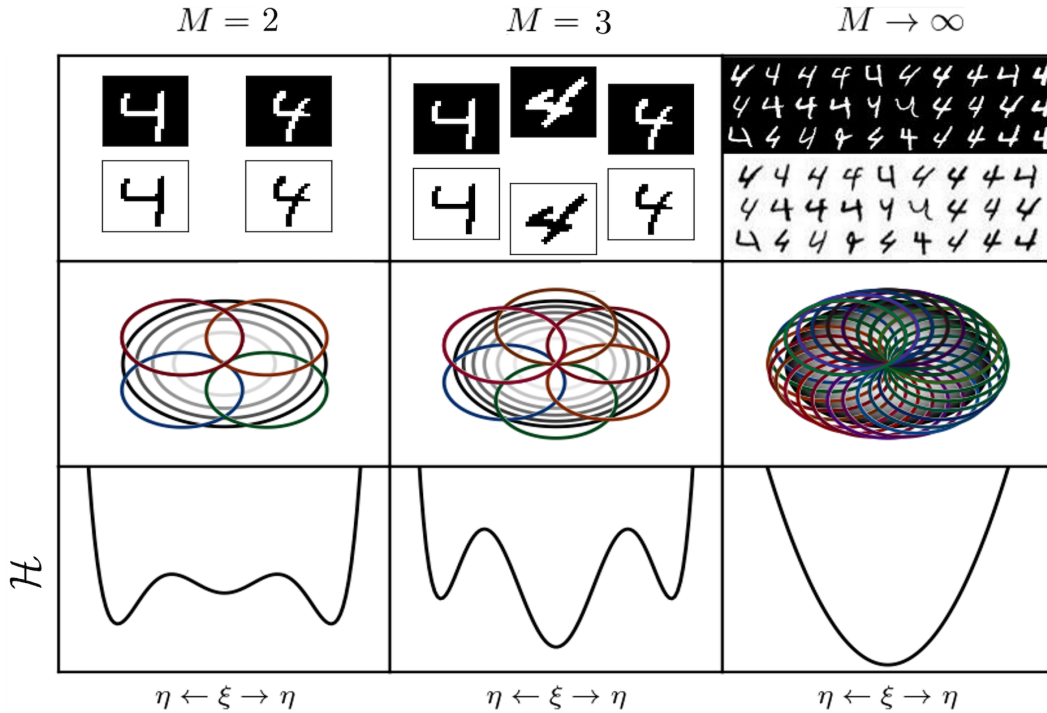
$$\mathcal{H}_{N,K,M,r}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\eta}) = -\frac{1}{\mathcal{R}^{P/2} M N^{P-1}} \sum_{\mu=1}^K \sum_{a=1}^M \left( \sum_{(i_1, \dots, i_P)}^{N, \dots, N} \eta_{i_1}^{\mu,a} \cdots \eta_{i_P}^{\mu,a} \sigma_{i_1} \cdots \sigma_{i_P} \right), \quad (3.1)$$

where  $P$  is the interaction order (assumed as even),  $\mathcal{R} = r^2 + \frac{1-r^2}{M}$  corresponds to  $\mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} [\sum_a \eta_i^{\mu,a} / (Mr)]^2$  (and plays as a normalization factor) and we also define  $\sum_{(i_1, \dots, i_P)}^{N, \dots, N} = \sum_{\substack{i_1, \dots, i_P \\ i_1 \neq \dots \neq i_P}}^{N, \dots, N}$  (namely the summation

in which only terms with all different "i" indices are taken into account). Further, the factor  $\frac{1}{N^{P-1}}$  in the right-hand side ensures the linear extensiveness of the Hamiltonian in the network size  $N$ .

The related partition function is defined as

$$\mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\eta}) = \sum_{\boldsymbol{\sigma}} \exp\left(-\beta \mathcal{H}_{N,K,M,r}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\eta})\right) =: \sum_{\boldsymbol{\sigma}} \mathcal{B}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\eta}), \quad (3.2)$$



**Figure 2:** Intuitive representation of the process of learning an archetype. In the upper row we show the neuron configurations corresponding to the supplied examples, in the middle row we schematically depict the attraction basins determined by these examples, and in the lower row we sketch a plausible cross-sectional view of the Hamiltonian function in a fictitious one dimensional representation of the configuration space. Going from left to right, as the number of examples  $M$  grows, the network learns to generalize from them by constructing a faithful representation of the generic archetype  $\xi$  (that has never been supplied to the network). The network tends to store at first each single example  $\eta^a$  without being able to retrieve the archetype (left column), thus the deepest minima of the Hamiltonian correspond to examples. Then, a minimum, close to  $\xi$ , appears and coexists with the other minima (middle column) and, finally, a unique stable minimum corresponding to the archetype emerges (right column). The variation of the energy landscape as  $M$  changes depends on the network architecture and on the dataset.

where  $\mathcal{B}_{N,K,M,r,\beta}^{(P)}(\sigma|\eta)$  is referred to as Boltzmann factor.

At finite network-size  $N$ , the quenched statistical pressure (or free energy<sup>2</sup> with no loss of generality.) of the model reads as

$$\mathcal{A}_{N,K,M,r,\beta}^{(P)} = \frac{1}{N} \mathbb{E} \log \mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\eta) \quad (3.3)$$

where  $\mathbb{E} = \mathbb{E}_{\xi} \mathbb{E}_{(\eta|\xi)}$  denotes the average over the realization of examples, namely over the distributions

<sup>2</sup>The free energy  $\mathcal{F}_{N,K,M,r,\beta}^{(P)}$  equals the statistical pressure, a factor  $-\beta$  apart, i.e.  $\mathcal{A}_{N,K,M,r,\beta}^{(P)} = -\beta \mathcal{F}_{N,K,M,r,\beta}^{(P)}$ . Thus, extremizing the former results in the same self-consistency equations for the macroscopic observables that we would obtain by extremizing the latter; in this paper we use mainly the statistical pressure with no loss of generality.

(2.1) and (2.6). By combining the quenched average  $\mathbb{E}[\cdot]$  and the Boltzmann average

$$\omega[\langle \cdot \rangle] := \frac{1}{\mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\eta})} \sum_{\boldsymbol{\sigma}} \langle \cdot \rangle \mathcal{B}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\eta}), \quad (3.4)$$

possibly replicated over two or more replicas<sup>3</sup>, that is,  $\Omega := \omega \times \omega \dots \times \omega$ , we get the expectation

$$\langle \cdot \rangle := \mathbb{E}\Omega(\cdot). \quad (3.5)$$

**Remark 1.** An integral representation of the partition function will be useful in the following numerical computations. Starting from Eq. (3.2), we apply the Hubbard-Stratonovich transformation to get

$$\mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\eta}) = \sum_{\boldsymbol{\sigma}} \int \prod_{\mu,a} d\tilde{\mu}(z_{\mu,a}) \exp \left[ \sqrt{\frac{\beta'}{\mathcal{R}^{P/2} M N^{P-1}}} \sum_{\mu>1}^K \sum_{a=1}^M \sum_{i_1, \dots, i_{P/2}}^{N, \dots, N} \eta_{i_1}^{\mu,a} \dots \eta_{i_{P/2}}^{\mu,a} \sigma_{i_1} \dots \sigma_{i_{P/2}} z_{\mu,a} \right] \quad (3.6)$$

where  $d\tilde{\mu}(z_{\mu,a}) = \frac{\exp(-z_{\mu,a}^2/2)}{\sqrt{2\pi}} dz_{\mu,a}$  is a Gaussian measure and we posed  $\beta' = 2\frac{\beta}{P!}$ . Moreover, we have exploited

$$\frac{P!}{2N^{P-1}} \sum_{(i_1, \dots, i_P)}^{N, \dots, N} (\Phi_{i_1}^\mu \dots \Phi_{i_P}^\mu) = \frac{1}{2N^{P-1}} \sum_{i_1, \dots, i_P}^{N, \dots, N} (\Phi_{i_1}^\mu \dots \Phi_{i_P}^\mu) + \mathcal{O}(N^{P/2-1}) \quad (3.7)$$

with  $\Phi_i^\mu$  is any finite random variable and we have neglected the subleading network-size terms.

We can think of the above transformation as a mapping between the original dense Hebbian network and a restricted Boltzmann machine (RBM) where  $K \times M$  hidden neurons  $z_{\mu,a}$  (equipped with a Gaussian prior) interact with the  $N$  visible neurons  $\boldsymbol{\sigma}$  grouped in sets each made of  $P/2$  neurons  $\sigma_{i_1} \dots \sigma_{i_{P/2}}$  with weight  $\eta_{i_1}^{\mu,a_1} \dots \eta_{i_{P/2}}^{\mu,a_{P/2}}$ . A schematic representation of the dense Hebbian network and its dual RBM are shown for a simple case in Fig. 3.

In our analytical investigation we leverage the asymptotic limit for the system size  $N$ , which shall be performed retaining the network load  $\alpha_b$  finite, as specified by the following

**Definition 2.** In the thermodynamic limit  $N \rightarrow \infty$ , the load is defined as

$$\lim_{N \rightarrow +\infty} \frac{K}{N^b} =: \alpha_b < \infty \quad (3.8)$$

with  $b \leq P-1$ <sup>4</sup>. We then distinguish between the so-called high-load regime, corresponding to  $b = P-1$ , namely to an amount of storable patterns that scales with the networks size as  $N^{P-1}$ , and a so-called low-load regime corresponding to  $b < P-1$ . As we will deepen, the resulting slower scaling for the amount of storable patterns allows for mitigating the effects of possible additive noise affecting synaptic strengths (see Sec. 5.3).

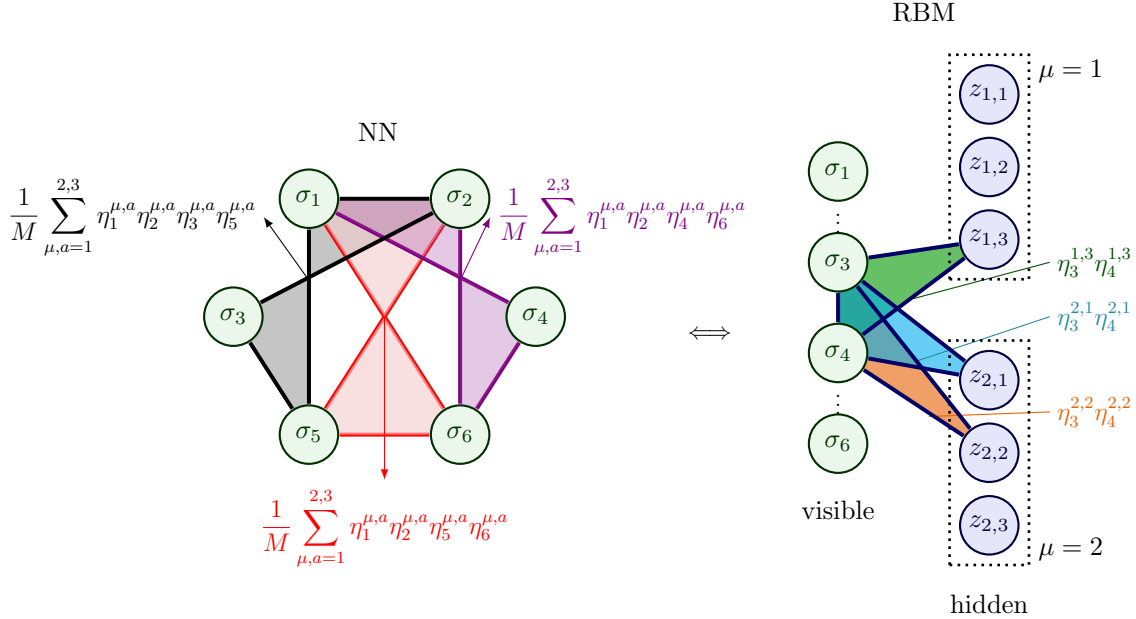
Further, the quenched statistical pressure in the thermodynamic limit is denoted as

$$\mathcal{A}_{\alpha_b, M, r, \beta}^{(P)} = \lim_{N \rightarrow \infty} \mathcal{A}_{N, K, M, r, \beta}^{(P)}. \quad (3.9)$$

<sup>3</sup>A replica is an independent copy of the system characterized by the same realization of disorder, namely by the same realization of the archetypes and examples. Thus, two replicas are sampled from the same distribution  $\mathcal{B}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\eta})$ . Comparing two copies allows us to determine whether slow noise prevails, that is, whether the interference between archetypes and examples prevents the system to retrieve, see Sec. 5

<sup>4</sup>The case  $b > P-1$  is known to lead to a black-out scenario [21, 25] not useful for computational purposes and shall be neglected here.





**Figure 3:** Representation of a dense unsupervised Hopfield network (NN, left) and its dual Restricted Boltzmann Machine (RBM, right), with  $N = 6$ ,  $K = 2$ ,  $M = 3$  and  $P = 4$ . Different groups of interacting neurons are depicted in different colors. As far as the RBM concerns, it is built with a visible layer made of  $N$  binary variables  $\{\sigma_i\}_{i=1,\dots,6}$  and a hidden layer made of  $K \times M$  Gaussian neurons  $\{z_{\mu,a}\}_{\mu=1,2}^{a=1,2,3}$ . In particular, any  $z_{\mu,a}$  can interact with sets of  $P/2$  (namely, 2 in this case) visible neurons  $\{\sigma_i, \sigma_j\}$  whose strength of interaction is  $\eta_i^{\mu,a} \eta_j^{\mu,a}$ . In the NN, the neurons interact 4-wise and the coupling strength for any set of variables  $\{\sigma_i, \sigma_j, \sigma_k, \sigma_l\}$  is  $\frac{1}{M} \sum_{\mu=1}^K \sum_{a=1}^M \eta_i^{\mu,a} \eta_j^{\mu,a} \eta_k^{\mu,a} \eta_l^{\mu,a}$ .

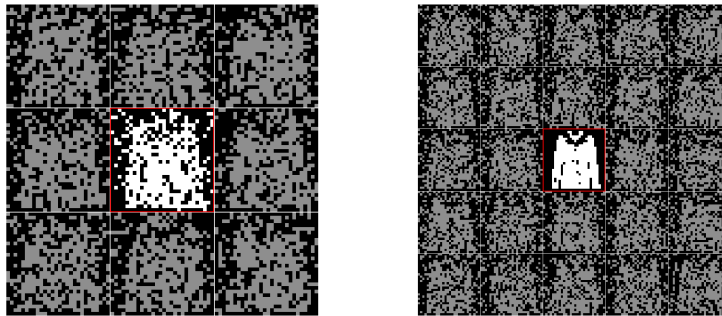
In order to further simplify the notation it is convenient to introduce a  $P$ -independent load denoted as  $\gamma$  and defined by

$$\alpha_{P-1} = \gamma \frac{2}{P!}. \quad (3.10)$$

We can notice that, as long as  $P$  is fixed, assuming that  $\alpha_{P-1} < \infty$  also means that  $\gamma < \infty$ .

We want to study the model defined in 1 looking at its learning and retrieval capabilities and, specifically, we aim to find out the thresholds for these capabilities to emerge. In other words, given a training dataset made of  $M \times K$  examples, each codified by a binary vector of size  $N$  and characterized by a quality  $r$ , and given a set of  $N$  binary neurons that interact  $P$ -wisely, we aim at answering the following questions:

- which is the minimum number of examples to be supplied to the network to ensure that it is able to infer the related archetypes and thus correctly generalize afterwards? (Note that we address this question while the network is handling simultaneously all the  $K \times M$  archetypes.)
- how many archetypes the network can learn and what happens if we load the network with a larger amount of information?



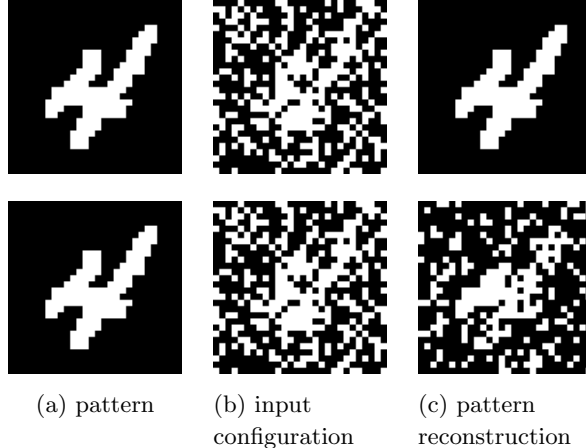
**Figure 4:** We consider a dense neural network with a degree of interactions  $P = 4$ , and we take just one picture of a coat from the Fashion-MNist dataset [45], thresholding the gray-scale values to obtain a binary representation of the image. Then, we generate other nine uncorrelated Rademacher archetypes, to reach a amount of patterns  $K = 10$ . Then, we produce  $M$  noisy examples for each archetype with quality level  $r = 0.5$  by flipping each pixel with probability as in (2.6). We focus on the cases where  $M = 8$  (left) and  $M = 24$  (right), and show the examples related to the Fashion-MNist coat, for these examples the white pixels are shown in gray. In both plots, they surround the final output of the network, that is the central image in the red box. The original picture of the coat is not shown but differs from the  $M = 24$  output by just two pixels. We notice that, among these two plots, solely in the one with  $M = 24$  the network is capable to correctly reconstruct the pattern starting from its noisy versions.

- can we account for training flaws in this system and, if so, how robust is the resulting pattern recognition capability of the network with respect to this kind of noise?

As we will see, the answers to the first two points highlight a conflicting role of the interaction order  $P$ : on the one hand, by increasing  $P$  the number of examples required for a sound training grows exponentially ( $\propto 1/(Pr^{2P})$ ), on the other hand, the number of storable archetypes also grows exponentially ( $\propto N^{P-1}$ ). Furthermore, to address the last question, we can introduce a supplementary, additive noise  $\omega$  to be applied to the couplings  $J_{i_1 i_2 \dots i_P}$  that mimics possible flaws occurring during the training and we show that there is an interplay between this kind of noise and the load: if we can afford a downgrade in terms of load (i.e.,  $b < P - 1$ ), then the network can work even in the presence of extensive noise (i.e.,  $\omega$  scaling with  $N$ ) on the weights. We anticipate that these features stem from, respectively, the vast available resources ( the  $K$  archetypes are allocated in a tensor made of  $N^P$  elements) and from the redundancy generated when employing over-sized resources.

These concepts are in part visualized in Figs. 4-5. Indeed, in Fig. 4 we show that, without a sufficient number  $M$  of examples, the network is incapable of generalizing an archetype starting from noisy versions of it. Moreover, it can happen that dense networks are able, from a noisy initial configuration, to recall the reference pattern better than their non-dense counterparts, as evidenced in Fig. 5, even if the number of examples given to the network and all the other parameters are equal.

To make the above statements quantitative, we need a set of observables to assess the retrieval ability of the network, therefore we state the following



**Figure 5:** Comparison between the retrieval capabilities exhibited by a dense network ( $P = 4$ , upper line) and by the pair-wise Hopfield model ( $P = 2$ , lower line). We built both the networks with  $N = 784$  Ising neurons and we chose a picture from the MNist dataset (i.e., the number 4). Then, we generated other 35 independent Rademacher archetypes, in order reach  $K = 36$ ; for each archetype the network unsupervisedly experienced  $M = 80$  examples at a level characterized by a quality  $r = 0.375$ . In the panels in the first column (a) we report the archetype, in the middle column (b) we report a noisy example inputted to the network and in the last column (c) we show the ultimate network’s reconstruction where it shines that, while the Hopfield model fails, the dense network correctly performs pattern recognition.

**Definition 3.** *The order parameters of the dense unsupervised Hebbian neural network introduced in Def. 1 are*

$$m_\mu := \frac{1}{N} \sum_{i=1}^N \xi_i^\mu \sigma_i \quad (3.11)$$

$$n_{\mu,a} := \frac{r}{\mathcal{R}} \frac{1}{N} \sum_{i=1}^N \eta_i^{\mu,a} \sigma_i, \quad (3.12)$$

$$q_{lm} := \frac{1}{N} \sum_{i=1}^N \sigma_i^{(l)} \sigma_i^{(m)}, \quad (3.13)$$

for  $\mu = 1, \dots, K$  and  $a = 1, \dots, M$ .

Note that the Mattis magnetization  $m_\mu$ , already defined in (2.3), quantifies the alignment of the network configuration  $\sigma$  with the archetype  $\xi^\mu$ ,  $n_{\mu,a}$  quantifies the alignment of the network configuration with the example  $\eta^{\mu,a}$ , and  $q_{lm}$  is the standard two-replica overlap between the replicas  $\sigma^{(l)}$  and  $\sigma^{(m)}$ .

## 4 Cost and Loss functions

Before proceeding with the investigation of the model, it is worth examining whether the order parameters introduced in the statistical-mechanics context to measure the ability of the system to learn and

retrieve archetypes from examples display any connection with the quantities usually employed in the machine-learning field. There, one typically introduces a *loss function*  $\mathcal{L}$ , namely a positive-definite function that maps any weight setting onto a real number representing some “cost” associated with that setting; during training the weights are tuned in such a way that  $\mathcal{L}$  is lowered and it reaches zero if and only if the network has learnt. Therefore, the goal of the learning stage is to vary the weights with some algorithm (e.g., contrastive divergence, back-propagation) in order to minimize  $\mathcal{L}$  as the various examples are provided to the network.

In the present case we want the system to learn to reconstruct archetypes from examples and the weights where the learnt information is allocated are the Hebbian couplings  $\mathbf{J}$ . Following an iterative procedure analogous to those adopted in a machine-learning context, we should prepare the system in some initial configuration  $(\boldsymbol{\sigma}^{(0)}, \mathbf{J}^{(0)})$ . Next, we should let the neurons (which are the fast degrees of freedom) relax to some  $\boldsymbol{\sigma}_{\mathbf{J}^{(0)}}^{(\text{eq})}$ , then evaluate the performance by some  $\mathcal{L}^{(0)} := \mathcal{L}(\boldsymbol{\sigma}_{\mathbf{J}^{(0)}}^{(\text{eq})})$  and modify couplings (e.g., via gradient descent) as  $\mathbf{J}^{(0)} \rightarrow \mathbf{J}^{(1)}$  in such a way that  $\mathcal{L}^{(1)} \leq \mathcal{L}^{(0)}$ , and so on so forth up to sufficiently small values of the loss function. For the present task, focusing on the  $\mu$ -th pattern, we envisage the following loss function

$$\mathcal{L}_\mu^{(n)} := \frac{1}{4N^2} \|\boldsymbol{\xi}^\mu + \boldsymbol{\sigma}_{\mathbf{J}^{(n)}}^{(\text{eq})}\|^2 \cdot \|\boldsymbol{\xi}^\mu - \boldsymbol{\sigma}_{\mathbf{J}^{(n)}}^{(\text{eq})}\|^2, \quad (4.1)$$

in such a way that  $\mathcal{L}_\mu^{(n)} \geq 0$ , and it reaches zero if and only if the system is retrieving the marked pattern (or its inverse, by gauge invariance). The loss function defined above can be recast in terms of the Mattis overlaps as

$$\mathcal{L}_\mu^{(n)} = [1 + m_\mu^{(n)}][1 - m_\mu^{(n)}], \quad (4.2)$$

where  $m_\mu^{(n)} := \frac{1}{N} \boldsymbol{\sigma}_{\mathbf{J}^{(n)}}^{(\text{eq})} \cdot \boldsymbol{\xi}^\mu$ , in such a way that the retrieval region in the phase diagram also highlights the set of values for the control parameters where  $\mathcal{L}_\mu^{(n)}$  is vanishing.

Further, we notice that the cost function of the the  $P$ -spin Hopfield model can be written as

$$\mathcal{H}_{N,K}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\xi}) = -\frac{N}{P!} \sum_{\mu=1}^K m_\mu^P = -\frac{N}{P!} \sum_{\mu=1}^K (1 - \mathcal{L}_\mu^*)^{\frac{P}{2}}, \quad (4.3)$$

where  $\mathcal{L}_\mu^*$  is the loss function evaluated for the generic configuration  $\boldsymbol{\sigma}$ . We should also keep in mind that we are considering a learning process where the available dataset is  $\{\boldsymbol{\eta}^{\mu,a}\}_{\mu=1,\dots,K}^{a=1,\dots,M}$  with  $\mu$  undisclosed. The most natural way to recover that framework is by replacing, in (4.3), archetypes with examples and then averaging over the latter; by doing so we recover the Hamiltonian (3.1).

Now, in the noiseless limit  $\beta \rightarrow \infty$ , the system spontaneously relaxes to configurations corresponding to the lowest energy which ensure that  $\mathcal{L}_\mu^* = 0$ . In order to see that only one of the losses that are summed up in the previous equation is minimised, and that the network will not attempt to minimise each of them at the same time, we can evaluate the average energy for the whole class of possible retrieval states. The most probable candidate states for retrieval are given by linear combinations of  $n$ -patterns:

$$\sigma_i = \text{sign} \left( \sum_{k=1}^n \xi_i^{\mu_k} \right), \quad (4.4)$$

their Mattis overlap is

$$m_{\mu_\ell} = \frac{1}{N} \sum_{i=1}^N \xi_i^{\mu_\ell} \text{sign} \left( \xi_i^{\mu_\ell} + \sum_{\substack{k=1 \\ k \neq \ell}}^n \xi_i^{\mu_k} \right). \quad (4.5)$$

Averaging over pattern realization, for any  $k \leq n$  we have

$$\mathbb{E}_{\xi}[m_{\mu_k}] = \int \frac{dz}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) \text{sign}(1 + \sqrt{n-1}z) = \text{erf}\left(\frac{1}{\sqrt{2(n-1)}}\right), \quad (4.6)$$

while  $\mathbb{E}_{\xi}[m_{\mu_k}] = 0$  for any  $k > n$ . We now estimate the expected energy for configurations like (4.4) obtaining:

$$\mathcal{H}_{N,K,M}^{(P)}(\sigma|\xi) = -\frac{nN}{P!} \left[ \text{erf}\left(\frac{1}{\sqrt{2(n-1)}}\right) \right]^P \quad (4.7)$$

Since this expression is a strictly increasing function of  $n$ , the only stable retrieval states are those with  $n = 1$ , thus their Mattis overlap will tend to 1, showing indeed that the network minimised only one of the  $\mathcal{L}_{\mu}^*$  losses at a time.

## 5 Analytical findings

In this section we solve the dense unsupervised Hopfield model introduced in Definition 1, specifically, we obtain an explicit expression for its quenched statistical pressure (i.e., the free energy) in terms of the order parameters of the theory. Then, by extremizing the free energy w.r.t. these order parameters we obtain a set of self-consistency equations for the latter whose inspection allows us to obtain the phase diagrams of the model. To this aim we exploit Guerra's interpolation technique [34] which allows us to get the free-energy explicitly. However, as we will see, some adjustments to the standard protocol are in order in the estimate of the noise distribution, which, unlike pairwise networks, can not be directly considered as a Gaussian random variable. The core problem is that the distributions of the post-synaptic potentials are not Gaussian here, hence standard universality of spin-glass noise [7, 27, 33] does not apply straightforwardly. To overcome this obstacle, we apply the Central Limit Theorem (CLT) in order to estimate it as a single Gaussian variable.

Before proceeding, we highlight that, as standard (see e.g., [28]) and with no loss of generality, in the following we will focus on the ability of the network to learn and retrieve the first archetype  $\xi^1$ . Thus, in the next expression, the contribution corresponding to  $\mu = 1$  shall be split from all the others and interpreted as the *signal* contribution, while the remaining ones make up the slow-noise contribution impairing both learning and retrieval of  $\xi^1$ , namely starting from Eq. (3.2), we apply the functional-generator technique<sup>5</sup> to get

$$\begin{aligned} \mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\eta) &= \lim_{J \rightarrow 0} \mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\xi^1, \eta; J) \\ &= \lim_{J \rightarrow 0} \sum_{\sigma} \exp \left[ J \sum_{i=1}^N \xi_i^1 \sigma_i + \frac{\beta'}{2 \mathcal{R}^{P/2} M N^{P-1}} \sum_{a=1}^M \left( \sum_{i=1}^N \eta_i^{a,1} \sigma_i \right)^P \right. \\ &\quad \left. + \frac{\beta' P!}{2 \mathcal{R}^{P/2} N^{P-1}} \sum_{\mu > 1}^K \sum_{(i_1, \dots, i_P)}^{N, \dots, N} \left( \frac{1}{M} \sum_{a=1}^M \eta_{i_1}^{\mu,a} \dots \eta_{i_P}^{\mu,a} \right) \sigma_{i_1} \dots \sigma_{i_P} \right]. \quad (5.1) \end{aligned}$$

Focusing only on the noise terms in round brackets in (5.1) we can apply the CLT and approximate it with a Gaussian variable with suitable first and second momenta. Therefore we can recast this term

<sup>5</sup>We add the term  $J \sum_i \xi_i^1 \sigma_i$  to generate the expectation of the Mattis magnetization  $m_1$ : the latter emerges by evaluating the derivative w.r.t.  $J$  of the quenched statistical pressure at  $J = 0$ .

as follows

$$\left( \frac{1}{M} \sum_{a=1}^M \eta_{i_1}^{\mu,a} \dots \eta_{i_P}^{\mu,a} \right) \sim r^P \sqrt{1 + \rho_P} \lambda_{i_1, \dots, i_P}^{\mu} \quad \text{with} \quad \lambda_{i_1, \dots, i_P}^{\mu} \sim \mathcal{N}(0, 1) \quad (5.2)$$

where  $\rho_P = \frac{1 - r^{2P}}{Mr^{2P}}$ . Remarkably, this reasoning shows that these dense networks exhibit the universality of the quenched noise [7, 27, 33], namely we can approximate the overall field experienced by a neuron (i.e., the post-synaptic potential) as a random Gaussian field.

Now, plugging (5.2) into (5.1) we reach a useful expression for the partition function for the unsupervised dense Hebbian neural network as

$$\begin{aligned} \mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\xi}^1, \boldsymbol{\eta}; J) &= \sum_{\boldsymbol{\sigma}} \exp \left[ J \sum_{i=1}^N \xi_i^1 \sigma_i + \frac{\beta'}{2\mathcal{R}^{P/2} M N^{P-1}} \sum_{a=1}^M n_{1,a}^P \right. \\ &\quad \left. + \frac{\beta' P! \sqrt{1 + \rho_P}}{2(1 + \rho)^{P/2} N^{P-1}} \sum_{\mu > 1}^K \left( \sum_{(i_1, \dots, i_P)}^{N, \dots, N} \lambda_{i_1, \dots, i_P}^{\mu} \sigma_{i_1} \dots \sigma_{i_P} \right) \right] \end{aligned} \quad (5.3)$$

where we exploit the relation  $\mathcal{R} = r^2(1 + \rho)$ .

We now proceed by applying Guerra's interpolation. The underlying idea behind this technique is to introduce a generalized free-energy which interpolates between the original one (which is the target of our investigation but we are not able to address it directly) and a simple one (which we can solve exactly). The latter is typically a one-body model mimicking the original one: the fields acting on neurons are chosen to exhibit statistical properties that simulate those experienced by neurons in the original model and due to the effect of other neurons. We thus find the solution of the simple model and we propagate the obtained solution back to the original model by the fundamental theorem of calculus. In this last passage we assume RS (*vide infra*), namely, we assume that the order-parameter fluctuations are negligible in the thermodynamic limit: this property makes the integral in the fundamental theorem of calculus analytical. Let us proceed by steps and give the next definitions

**Definition 4.** *Given the interpolating parameter  $t \in [0, 1]$ , the constants  $A, \psi \in \mathbb{R}$  to be set a posteriori, and the i.i.d. standard Gaussian variables  $Y_i \sim \mathcal{N}(0, 1)$  for  $i = 1, \dots, N$ , the interpolating partition function is given as*

$$\mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\xi}^1, \boldsymbol{\eta}; J, t) := \sum_{\boldsymbol{\sigma}} \mathcal{B}_{N,K,M,\beta}^{(P)}(\boldsymbol{\sigma} | \boldsymbol{\xi}^1, \boldsymbol{\eta}; J, t). \quad (5.4)$$

where  $\mathcal{B}_{N,K,M,r,\beta}^{(P)}$  is the related Boltzmann factor reads as

$$\begin{aligned} \mathcal{B}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\sigma} | \boldsymbol{\xi}^1, \boldsymbol{\eta}; J, t) &:= \exp \left[ J \sum_{i=1}^N \xi_i^1 \sigma_i + \frac{t\beta' N}{2M} (1 + \rho)^{P/2} \sum_{a=1}^M n_{1,a}^P + \psi(1 - t) N \sum_{a=1}^M n_{1,a} \right. \\ &\quad \left. + \sqrt{t} \frac{\beta' P! \sqrt{1 + \rho_P}}{2(1 + \rho)^{P/2} N^{P-1}} \sum_{\mu > 1}^K \sum_{(i_1, \dots, i_P)}^{N, \dots, N} \lambda_{i_1, \dots, i_P}^{\mu} \sigma_{i_1} \dots \sigma_{i_P} + \sqrt{1 - t} A \sum_{i=1}^N Y_i \sigma_i \right]; \end{aligned} \quad (5.5)$$

A generalized average follows from this generalized measure as

$$\omega_t[(\cdot)] := \frac{1}{\mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\xi}^1, \boldsymbol{\eta}; J, t)} \sum_{\boldsymbol{\sigma}} (\cdot) \mathcal{B}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\sigma} | \boldsymbol{\xi}^1, \boldsymbol{\eta}; J, t) \quad (5.6)$$

and

$$\langle(\cdot)\rangle_t := \mathbb{E}\{\omega_t[(\cdot)]\}, \quad (5.7)$$

where the expectation  $\mathbb{E}$  is now meant over any  $\lambda_{i_1, \dots, i_P}^\mu$  and  $Y_i$  too.

The interpolating quenched statistical pressure related to the partition function (5.4) is introduced as

$$\mathcal{A}_{N,K,M,r,\beta}^{(P)}(J, t) := \frac{1}{N} \mathbb{E} \left[ \ln \mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\xi}, \boldsymbol{\eta}; J, t) \right], \quad (5.8)$$

and, in the thermodynamic limit,

$$\mathcal{A}_{\alpha_b, M, r, \beta}^{(P)}(J, t) := \lim_{N \rightarrow \infty} \mathcal{A}_{N, K, M, r, \beta}^{(P)}(J, t). \quad (5.9)$$

Of course, by setting  $t = 1$  we recover the original model: the interpolating pressure recovers the original one (3.3), that is  $\mathcal{A}_{N, K, M, r, \beta}^{(P)}(J) = \mathcal{A}_{N, K, M, r, \beta}^{(P)}(J, t = 1)$ , and analogously for the partition function, the standard Boltzmann measure and the related averages.

As anticipated, the following analytical results are obtained under the RS hypothesis, namely assuming that, in the thermodynamic limit, the distribution of the generic order parameter  $X$  is centered at its expectation value w.r.t. the interpolating measure  $\bar{X} := \langle X \rangle_t$  with vanishing fluctuations for all  $t$ , that is,

$$\lim_{N \rightarrow \infty} \langle (X - \bar{X}) \rangle_t = 0. \quad (5.10)$$

Although this assumption is not fulfilled by this kind of systems (at least not everywhere in the space of control parameters), it is usually adopted as it yields only small quantitative corrections and, further, a full replica-symmetry-breaking theory for these systems is still under construction (see e.g., [2, 3, 11, 29, 41]).

We now proceed to determine the self-consistency equations for the order parameters by extremizing the quenched statistical pressure; to this aim it is mathematically convenient to take the thermodynamic limit and split the discussion in two cases: the high-storage regime  $b = P - 1$  (corresponding to the highest load allowed [9, 21, 25]) and the low-storage regime  $b < P - 1$ ; as stressed above, in both cases, we shall consider only even values of  $P$  and, specifically,  $P \geq 4$ <sup>6</sup>.

## 5.1 High-load regime

In this subsection we present the main analytical result obtained in the case where  $K/N^{P-1}$  remains finite in the thermodynamic limit, that is  $\alpha_{P-1}$  is finite and non-vanishing, see (3.8).

**Proposition 1.** *In the thermodynamic limit ( $N \rightarrow \infty$ ), under the RS assumption (5.10), the quenched statistical pressure for the unsupervised, dense neural-network described by (5.3) set in the high-load regime  $b = P - 1$  reads as*

$$\mathcal{A}_{\gamma, M, r, \beta}^{(P)}(J) = \mathbb{E} \left\{ \ln 2 \cosh \left[ J \xi^1 + \beta' \frac{P}{2} \bar{n}^{P-1} (1 + \rho)^{P/2-1} \hat{\eta} + Y \sqrt{\gamma \frac{\beta'^2 (1 + \rho_P) P}{(1 + \rho)^P} \frac{P}{2} \bar{q}^{P-1}} \right] \right\} \quad (5.11)$$

$$- \frac{\beta'}{2} (P - 1) (1 + \rho)^{P/2} \bar{n}^P + \gamma \frac{\beta'^2 (1 + \rho_P)}{4(1 + \rho)^P} (1 - P \bar{q}^{P-1} + (P - 1) \bar{q}^P).$$

<sup>6</sup>The case  $P = 2$  corresponds to the unsupervised Hopfield model treated in [6]; the assumption  $P \geq 4$  is used in the proof of Theorem 1 in Appendix A.

with  $\mathbb{E} = \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} \mathbb{E}_Y$ ,  $\hat{\eta} := \frac{1}{rM} \sum_{a=1}^M \eta^{1,a}$  and  $\bar{n}$  and  $\bar{q}$  fulfill the following self-consistency equations

$$\begin{aligned} \bar{n} &= \frac{1}{1+\rho} \mathbb{E} \left\{ \tanh \left[ \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} \hat{\eta} + Y \sqrt{\gamma \frac{\beta'^2 (1+\rho_P) P}{(1+\rho)^P} \frac{P}{2} \bar{q}^{P-1}} \right] \hat{\eta} \right\}, \\ \bar{q} &= \mathbb{E} \left\{ \tanh^2 \left[ \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} \hat{\eta} + Y \sqrt{\gamma \frac{\beta'^2 (1+\rho_P) P}{(1+\rho)^P} \frac{P}{2} \bar{q}^{P-1}} \right] \right\}. \end{aligned} \quad (5.12)$$

Furthermore, considering the auxiliary field  $J$  linked to  $\bar{m}$  as  $\bar{m} = \nabla_J \mathcal{A}_{\gamma, M, r, \beta}^{(P)}(J)|_{J=0}$ , we have

$$\bar{m} = \mathbb{E} \left\{ \tanh \left[ \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} \hat{\eta} + Y \sqrt{\gamma \frac{\beta'^2 (1+\rho_P) P}{(1+\rho)^P} \frac{P}{2} \bar{q}^{P-1}} \right] \xi^1 \right\}. \quad (5.13)$$

For the proof we refer to Appendix A.

## 5.2 Low-load regime

In this subsection we present the main analytical finding obtained by setting  $b < P - 1$  in (3.8).

**Proposition 2.** *In the thermodynamic limit ( $N \rightarrow \infty$ ), under the RS assumption (5.10), the quenched statistical pressure for the unsupervised, dense neural-network described by (5.3) set in the low-load regime  $b < P - 1$  reads as*

$$\mathcal{A}_{0, M, \beta, r}^{(P)}(J) = \mathbb{E} \left\{ \ln 2 \cosh \left[ J \xi^1 + \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} \hat{\eta} \right] \right\} - \frac{\beta'}{2} (P-1) (1+\rho)^{P/2} \bar{n}^P. \quad (5.14)$$

with  $\mathbb{E} = \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)}$  and  $\bar{n}$  fulfills the following self-consistency equation

$$\bar{n} = \frac{1}{1+\rho} \mathbb{E} \left\{ \tanh \left[ \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} \hat{\eta} \right] \hat{\eta} \right\}. \quad (5.15)$$

Furthermore, considering the auxiliary field  $J$  linked to  $\bar{m}$  as  $\bar{m} = \nabla_J \mathcal{A}_{0, M, \beta, r}^{(P)}(J)|_{J=0}$  we have

$$\bar{m} = \mathbb{E} \left\{ \tanh \left[ \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} \hat{\eta} \right] \xi^1 \right\}. \quad (5.16)$$

The proof is analogous to the case  $b = P - 1$  (see Appendix A), therefore we will omit it.

Note that, in this low-load regime, we are left with one single order parameter that measures the degree of order in the system, much as like in the Curie-Weiss model [22, 23]. In fact, here  $\bar{q} = 0$  and this indicates a loss of slow-noise or of “glassiness” in the network, analogously to what happens for the pairwise Hopfield model in the low-storage regime  $\lim_{N \rightarrow \infty} K/N = 0$ . However, in this dense generalization, this regime deserves a particular attention because, as we will see in Sec. 5.3, a relatively small load can release some resources to handle possible supplementary noise due, for instance, to flaws underlying storing [10]. We anticipate that, in that case, we ultimately recover self-consistence equations similar to those obtained in high-load ( $b = P - 1$ ), but the noise term, instead of stemming from the load, will be related to this additional disturbance.



### 5.3 Additive noise in low-load regime

As mentioned in the previous section and shown numerically in the next one, by increasing the interaction order  $P$  among neurons, the storage capacity increases arbitrarily ( $K \sim N^{P-1}$  – *high-load regime*). However, our model assumes that the coupling tensor  $\mathbf{J}$  is devoid of flaws, whereas, in general, the communication among neurons can be disturbed hence affecting the synaptic processes (e.g., see [4, 24]). It is then natural to question if unsupervised dense neural networks described by (5.3) are robust versus this kind of noise too.

Recalling the Hamiltonian (3.1) we can write

$$\begin{aligned} \mathcal{H}_{N,K,M}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\eta}) &= - \sum_{(i_1, \dots, i_P)}^{N, \dots, N} J_{i_1 \dots i_P} \sigma_{i_1} \dots \sigma_{i_P} \\ &= - \frac{1}{\mathcal{R}^{P/2} M N^{P-1}} \sum_{\mu=1}^K \sum_{a=1}^M \sum_{(i_1, \dots, i_P)}^{N, \dots, N} \eta_{i_1}^{\mu, a} \dots \eta_{i_P}^{\mu, a} \sigma_{i_1} \dots \sigma_{i_P}, \end{aligned} \quad (5.17)$$

where we outlined the entry of the coupling tensor  $\mathbf{J}$ . Then, following [10], we model the supplementary noise, by introducing an additional, random contribution as

$$J_{i_1 \dots i_P} \rightarrow \tilde{J}_{i_1 \dots i_P} = \eta_{i_1}^{\mu, a} \dots \eta_{i_P}^{\mu, a} + w \tilde{\eta}_{i_1 \dots i_P}^{\mu, a}, \quad (5.18)$$

with  $w \in \mathbb{R}$  and  $\tilde{\eta}_{i_1 \dots i_P}^{\mu, a} \sim_{\text{iid}} \mathcal{N}(0, 1)$ .

We investigate the effects of such a noise on the retrieval capabilities of the system and the existence of upper bounds for the amount of noise that the system can tolerate without losing its ability to play as an associative memory.

As shown in [10], if we can afford a downgrade in terms of load (i.e.  $b < P - 1$ ), we can consider the presence of extensive synaptic noise that grows algebraically with the network size  $N$ , namely

$$w = \tau N^\delta, \quad \text{with } \tau \in \mathbb{R} \quad \text{and } \delta \in \mathbb{R}^+. \quad (5.19)$$

In fact, following the same path presented in the first part of this section, including the noise defined in (5.18) and (5.19) yields self-consistency equations for the order parameters that display the same expression found in the high-storage regime (5.15)-(5.16), as long as we replace  $\beta'$  with  $\tau\beta'$  in the noise part, namely

$$\begin{aligned} \bar{n} &= \frac{1}{1 + \rho} \mathbb{E} \left\{ \tanh \left[ \beta' \frac{P}{2} \bar{n}^{P-1} (1 + \rho)^{P/2-1} \hat{\eta} + Y \beta' \tau \sqrt{\gamma \frac{(1 + \rho_P) P}{(1 + \rho)^P} \frac{P}{2} \bar{q}^{P-1}} \right] \hat{\eta} \right\}, \\ \bar{q} &= \mathbb{E} \left\{ \tanh^2 \left[ \beta' \frac{P}{2} \bar{n}^{P-1} (1 + \rho)^{P/2-1} \hat{\eta} + Y \beta' \tau \sqrt{\gamma \frac{(1 + \rho_P) P}{(1 + \rho)^P} \frac{P}{2} \bar{q}^{P-1}} \right] \right\}, \\ \bar{m} &= \mathbb{E} \left\{ \tanh \left[ \beta' \frac{P}{2} \bar{n}^{P-1} (1 + \rho)^{P/2-1} \hat{\eta} + Y \beta' \tau \sqrt{\gamma \frac{(1 + \rho_P) P}{(1 + \rho)^P} \frac{P}{2} \bar{q}^{P-1}} \right] \xi^1 \right\}. \end{aligned} \quad (5.20)$$

Then, one can show that, if we have an extensive noise (5.19) with  $\delta < \frac{P-1-b}{2}$ , the noise contribution in the hyperbolic tangent in the previous equations is vanishing and we recover the low-load scenario. In other words, there is an interplay between the load (ruled by  $b$ ), the interaction order

(ruled by  $P$ ) and the supplementary noise (ruled by  $\delta$ ). Thus, when one of these is enhanced, the others must be overall suitably downsized if we want to preserve the retrieval capability of the system.

Before concluding we stress that the results obtained in this subsection are not influenced by the dataset parameters  $M$  and  $r$ ; this implies that we can not leverage either the quality or the quantity of the dataset to mitigate the effects of this supplementary noise.

#### 5.4 Low-entropy datasets in the high-load regime

As explained in Sec. 2, the parameter  $\rho = (1 - r^2)/(Mr^2)$  quantifies the amount of information needed to describe the original message  $\xi^\mu$  given the set of related examples  $\{\eta^{\mu,a}\}_{a=1,\dots,M}$ . In this section we focus on the case  $\rho \ll 1$  that corresponds to a low-entropy dataset or, otherwise stated, to a high-informative dataset. The advantage of this analysis is that, under this condition, we obtain a relation between  $\bar{n}$  (a natural order parameter of the model) and  $\bar{m}$  (a practical order parameter of the model)<sup>7</sup>, thus, the self-consistency equation for  $\bar{n}$  can be recast into a self-consistency equation for  $\bar{m}$  and its numerical solution versus the control parameters allows us to get the phase diagram for the system more straightforwardly.

As explained in Appendix A, we start from the self-consistency equations found in the high-storage regime (5.12)-(5.13) and we exploit the CLT to write  $\hat{\eta} \sim 1 + \lambda\sqrt{\rho}$ . In this way we reach the simpler expressions

$$(1 + \rho)\bar{n} = \bar{m} + \beta' \frac{P}{2} \rho (1 + \rho)^{P/2-1} (1 - \bar{q})\bar{n}^{P-1}, \quad (5.21)$$

$$\bar{q} = \mathbb{E}_Z [\tanh^2 g(\beta, Z, \bar{n})], \quad (5.22)$$

$$\bar{m} = \mathbb{E}_Z [\tanh g(\beta, Z, \bar{n})], \quad (5.23)$$

where

$$g(\beta, Z, \bar{n}) = \beta' \frac{P}{2} \bar{n}^{P-1} (1 + \rho)^{P/2-1} + \beta' Z \sqrt{\rho \frac{P^2}{4} \bar{n}^{2P-2} (1 + \rho)^{P-2} + \gamma \frac{(1 + \rho_P) P}{(1 + \rho)^P} \frac{P}{2} \bar{q}^{P-1}} \quad (5.24)$$

and  $Z \sim \mathcal{N}(0, 1)$  is a standard Gaussian variable.

Focusing on the argument of the hyperbolic tangent (5.24), we can split it into three parts: the first one represents the amplification of the signal; the second one reflects the use of perturbed version of the retrieved pattern and not the pattern themselves; the third one is the noise linked to the presence of the other patterns.

Further, in the retrieval region, where  $1 - \bar{q}$  is vanishing, as long as  $\rho \ll 1$ , we can truncate the right-hand-side of (5.21) into  $\bar{n}(1 + \rho) \sim \bar{m}$ . This leads to significant advantages in the computation time required to get a numerical solution of the self-consistency equations. In fact, by using  $\bar{n}(1 + \rho) = \bar{m}$ , in the argument of hyperbolic tangent, we get

$$g(\beta, Z, \bar{m}) = \tilde{\beta} \frac{P}{2} \bar{m}^{P-1} + \tilde{\beta} Z \sqrt{\rho \left( \frac{P}{2} \bar{m}^{P-1} \right)^2 + \gamma (1 + \rho_P) \frac{P}{2} \bar{q}^{P-1}}, \quad (5.25)$$

<sup>7</sup>It is worth recalling that the model is supplied only with examples – upon which  $\{n^{\mu,a}\}$  are defined – while it is not aware of archetypes – upon which  $\{m^\mu\}$  are defined. The former constitute natural order parameters and, in fact, the Hamiltonian  $\mathcal{H}_{N,K,M,r}^{(P)}$  in (3.1) can be written in terms of the example overlaps. The latter are practical order parameters through which we can assess the capabilities of the network.

where we put

$$\tilde{\beta} = \frac{\beta'}{(1+\rho)^{\frac{P}{2}}} = \frac{2\beta}{P!} \frac{1}{(1+\rho)^{\frac{P}{2}}}. \quad (5.26)$$

The consequent, remarkable reward of this truncation consists in retaining only two of the three self-consistency equations, namely only the ones for  $\bar{q}$  and  $\bar{m}$ , while the resulting error by this truncation is numerically small, as checked in Fig. 6 where we plot  $\bar{n}$  versus  $r$  for different values of the parameters and compare the outcomes obtained with and without the truncation.

**Remark 2.** We stress that here we are focusing only on the low entropy dataset limit because in the high entropy scenario, i.e.  $\rho \gg 1$ , the (5.21) will reduce to

$$\bar{n} = \tilde{\beta} \frac{P}{2} \rho (1+\rho)^{P-2} (1-\bar{q}) \bar{n}^{P-1}, \quad (5.27)$$

whose solutions are  $\bar{n} = 0$  or

$$\bar{n} = \frac{1}{(1+\rho)} \left( \tilde{\beta} \frac{P}{2} \rho (1-\bar{q}) \right)^{-\frac{1}{P-2}}. \quad (5.28)$$

Replacing this expression in (5.24) we get a signal term that reads as

$$\tilde{\beta} \frac{P}{2} (1+\rho)^{P-1} \bar{n}^{P-1} \underset{\rho \gg 1}{\sim} \rho^{-\frac{P-1}{P-2}}. \quad (5.29)$$

Therefore, the signal term in (5.24) will be strongly suppressed and the retrieval process will no longer be possible.

We now further handle the Eqs. (5.21) by computing their zero-temperature limit. As detailed in the Appendix A, by taking the limit  $\beta \rightarrow \infty$  in Eqs. (5.22) and (5.23) we get

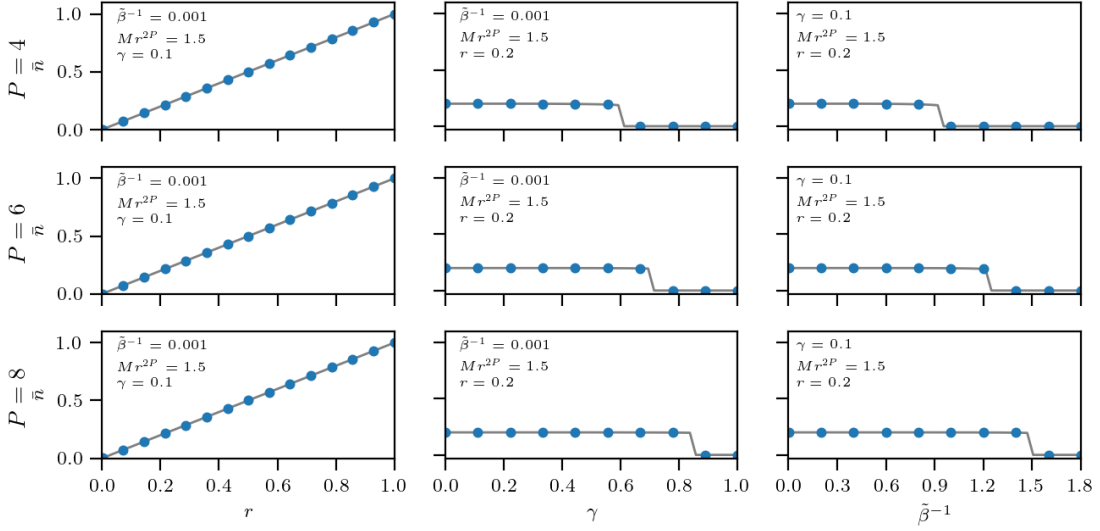
$$\begin{aligned} \bar{m} &= \operatorname{erf} \left( \frac{P \bar{m}^{P-1}}{2G} \right), \quad \bar{q} = 1, \\ G &= \sqrt{2 \left[ \rho \left( \frac{P \bar{m}^{P-1}}{2} \right)^2 + \gamma (1+\rho_P) \frac{P}{2} \right]}. \end{aligned} \quad (5.30)$$

## 6 Numerical findings

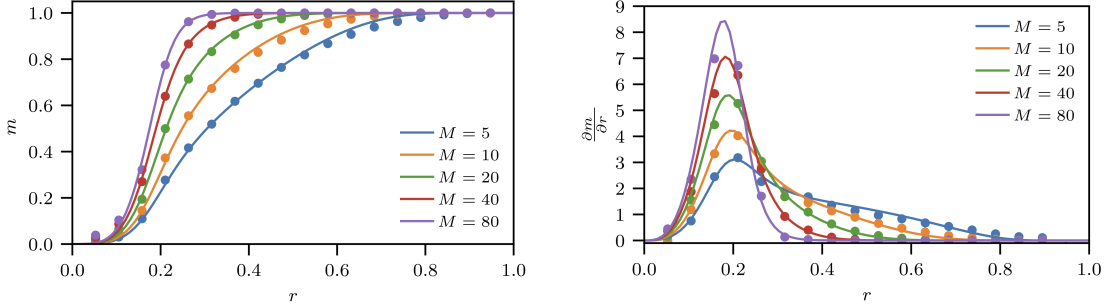
In this section we present some results useful to check the effective performance of the network. First, we use the stability analysis to find an explicit expression for  $\bar{m}$  in the noiseless limit also via this path. Next, we estimate the minimum number of examples, as a function of  $P, r$ , and  $\gamma$ , that we need for a successful retrieval. Finally, we show some outcomes obtained by MC simulations, concerning the reconstruction of the archetypes and critical load.

### 6.1 Stability analysis and Monte Carlo simulations

In this section we carry on a stability analysis in the noiseless limit: we suppose that the network is in a retrieval configuration, say  $\sigma = \xi^1$  without loss of generality, we evaluate the local field  $h_i(\xi^1)$  acting on the generic neuron  $\sigma_i$ , and check that  $h_i(\xi^1)\sigma_i > 0$  is satisfied for any  $i = 1, \dots, N$ ; this condition ensures the stability of the retrieval configuration.



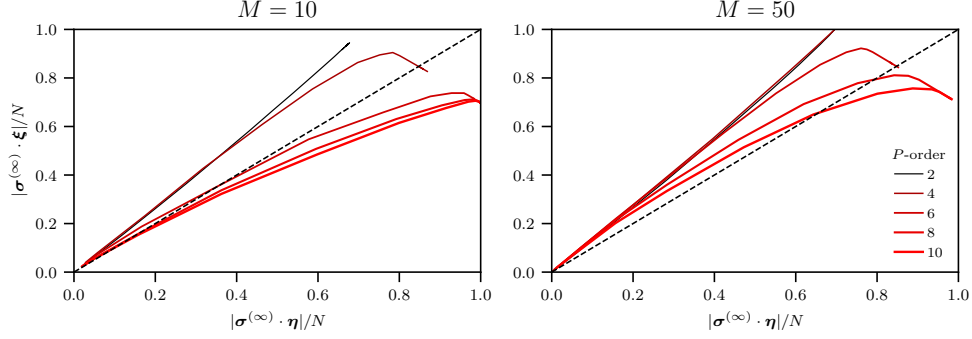
**Figure 6:** We compute  $\bar{n}$  solving numerically (5.21) for different values of  $P$ ,  $\tilde{\beta}$ ,  $\gamma$  and  $M$ , as reported in each panel, and we plot it versus  $r$ . We compare the results obtained using the exact expression of  $\bar{n}$  (blue dots) and those obtained using the approximated one (namely  $\bar{n}(1 + \rho) = \bar{m}$ , solid grey line). We notice that there is completely agreement between the exact expression and the approximated one, whatever the value of  $P$  considered.



**Figure 7:** Comparisons between observables evaluated by MC simulations equipped with Plefka’s dynamic (lines) and stability analysis (dots, see (6.12)). The number of examples  $M$  varies as specified by the legend, while the number of neurons and patterns and are kept fixed at  $N = 6000$ ,  $K = 100$ . As a consequence, the load is fixed below the critical value,  $\gamma < \gamma_c$ . In particular, we report the archetype magnetization  $m$  and its susceptibility  $\partial_r m$  at various training-set sizes  $M$  by making the noise  $r$  in the training set vary from 0, where all the example are pure random noise, to 1, where there is no difference among examples and archetype. We note that in the small noise limit  $r \rightarrow 1$  the network always perfectly retrieves the archetype as expected, whereas, for  $r \rightarrow 0$ , no retrieval is possible.

We start by rearranging the cost function (3.1) exploiting the mean-field nature of the model, namely

$$-\beta \mathcal{H}_{N,K,M,r}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\eta}) = \sum_{i=1}^N h_i(\boldsymbol{\sigma})\sigma_i \quad (6.1)$$



**Figure 8:** Analysis of the capacity of the network to reconstruct an archetype generalizing from corrupted versions of it. The dataset is generated by  $K = 10$  Rademacher archetypes, each of size  $N = 784$ , whence  $M = 10$  (left panel) and  $M = 50$  (right panel) examples are built for each archetype by setting  $r = 0.7$ . Then, we let the system relax from an initial configuration  $\sigma^{(0)}$  – chosen as a corrupted version of one of the examples, say  $\eta^a$  – to the thermalized configuration  $\sigma^{(\infty)}$  by Plefka’s dynamics. We determine the overlap between  $\sigma^{(\infty)}$  and  $\eta^a$ , as well as the overlap between  $\sigma^{(\infty)}$  and the archetype  $\xi$ . These quantities are then averaged over different initializations. Notice that these overlaps, i.e., the normalized scalar products, provide a measure of resemblance between the involved vectors; for instance, the Hamming distance between  $\xi$  and  $\sigma^{(0)}$  is nothing but  $\frac{1}{2}(N - \xi \cdot \sigma^{(0)})$ . The dashed line represents the identity and plays a reference: above this line the system has escaped from the attraction basin of the example  $\eta^a$  and has moved closer to the archetype. We notice that, as the interaction degree  $P$  increases (see the legend), the attractivity of the archetype is impaired. If we want that the curve remains above the threshold as  $P$  increases, the number of examples has to be increased accordingly.

where the local field  $h_i(\sigma)$  acting on the  $i$ -th spin is

$$h_i(\sigma) = \frac{1}{\mathcal{R}^{P/2} M N^{P-1}} \sum_{\mu=1}^K \sum_{a=1}^M \sum_{(i_2, \dots, i_P) \neq i}^{N, \dots, N} \eta_{i_1}^{\mu, a} \dots \eta_{i_P}^{\mu, a} \sigma_{i_2} \dots \sigma_{i_P}. \quad (6.2)$$

Calling  $O^{(n)}$  the  $n$ -th iteration of the MC Markov chain scheme regarding the generic observable  $O$  and starting by a Cauchy condition where the neurons are aligned with the first pattern, i.e.  $\sigma^{(0)} = \xi^1$ , we update the neural configuration as

$$\sigma_i^{(n+1)} = \sigma_i^{(n)} \text{sign} \left[ \tanh \left( \sigma_i^{(n)} h_i^{(n)}(\sigma^{(n)}) \right) + \Gamma_i \right] \quad \text{with } \Gamma_i \sim \mathcal{U}[-1; +1] \quad (6.3)$$

and, performing the zero fast-noise limit  $\beta \rightarrow \infty$ , we have

$$\sigma_i^{(n+1)} = \sigma_i^{(n)} \text{sign} \left[ \sigma_i^{(n)} h_i^{(n)}(\sigma^{(n)}) \right]. \quad (6.4)$$

The one-step MC approximation for the magnetization is then

$$m_1^{(2)} := \frac{1}{N} \sum_{i=1}^N \xi_i^1 \sigma_i^{(2)} = \frac{1}{N} \sum_{i=1}^N \text{sign} \left( \xi_i^1 h_i^{(1)}(\xi^1) \right), \quad (6.5)$$

and, in the thermodynamic limit ( $N \rightarrow \infty$ ) the argument of the sign function in the r.h.s. of Eq. (6.5) can be approximated, by the CLT<sup>8</sup>, as  $\xi_i^1 h_i^{(1)} \sim \mu_1 + z_i \sqrt{\mu_2 - \mu_1^2}$ , where  $z_i \sim \mathcal{N}(0, 1)$ , and

$$\mu_1 := \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} \left[ \xi_i^1 h_i^{(1)}(\boldsymbol{\xi}^1) \right] \quad (6.6)$$

$$\mu_2 := \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} \left\{ \left[ h_i^{(1)}(\boldsymbol{\xi}^1) \right]^2 \right\}. \quad (6.7)$$

Then, recalling

$$\int_{-\infty}^{+\infty} \frac{dz}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} \text{sign} \left( \mu_1 + z \sqrt{\mu_2 - \mu_1^2} \right) = \text{erf} \left( \frac{\mu_1}{\sqrt{2(\mu_2 - \mu_1^2)}} \right),$$

for  $N \gg 1$ , we get

$$m_1^{(2)} \sim \text{erf} \left( \frac{\mu_1}{\sqrt{2(\mu_2 - \mu_1^2)}} \right). \quad (6.8)$$

As reported in Appendix C, first and second momenta of  $\xi_i^1 h_i^{(1)}(\boldsymbol{\xi}^1)$  read as

$$\mu_1 = \frac{1}{(1 + \rho)^{P/2}} \quad (6.9)$$

$$\mu_2 = \left( \frac{1}{(1 + \rho)^{P/2}} \right)^2 \left[ \alpha_{P-1}(P-1)!(1 + \rho_P) + 1 + \rho \right] \quad (6.10)$$

where, we introduced  $\rho_P := \frac{1-r^{2P}}{Mr^{2P}}$  as a generalization of the dataset entropy  $\rho = \frac{1-r^2}{Mr^2}$ . Using the  $P$ -independent load (see Eq. (3.10)), we can write

$$\mu_2 - \mu_1^2 = \left( \frac{1}{(1 + \rho)^{P/2}} \right)^2 \left[ \frac{2}{P} \gamma (1 + \rho_P) + \rho \right], \quad (6.11)$$

and we obtain the following explicit expression for the one-step MC magnetization

$$m_1^{(2)} \sim \text{erf} \left\{ \left[ \frac{4\gamma}{P} (1 + \rho_P) + 2\rho \right]^{-\frac{1}{2}} \right\}. \quad (6.12)$$

Thus, we can recast the condition determining if the network successfully retrieves one of the archetypes by requiring that this one-step MC magnetization is larger than  $\text{erf}(\Theta)$  where  $\Theta \in \mathbb{R}^+$  is a tolerance level, thus we obtain

$$\frac{1}{\sqrt{2 \left[ \gamma \frac{2}{P} (1 + \rho_P) + \rho \right]}} > \Theta. \quad (6.13)$$

Otherwise stated, in order to retrieve (under a confidence level  $\Theta$ ) a given archetype starting from a perturbed versions of it, one has to fulfil the following condition

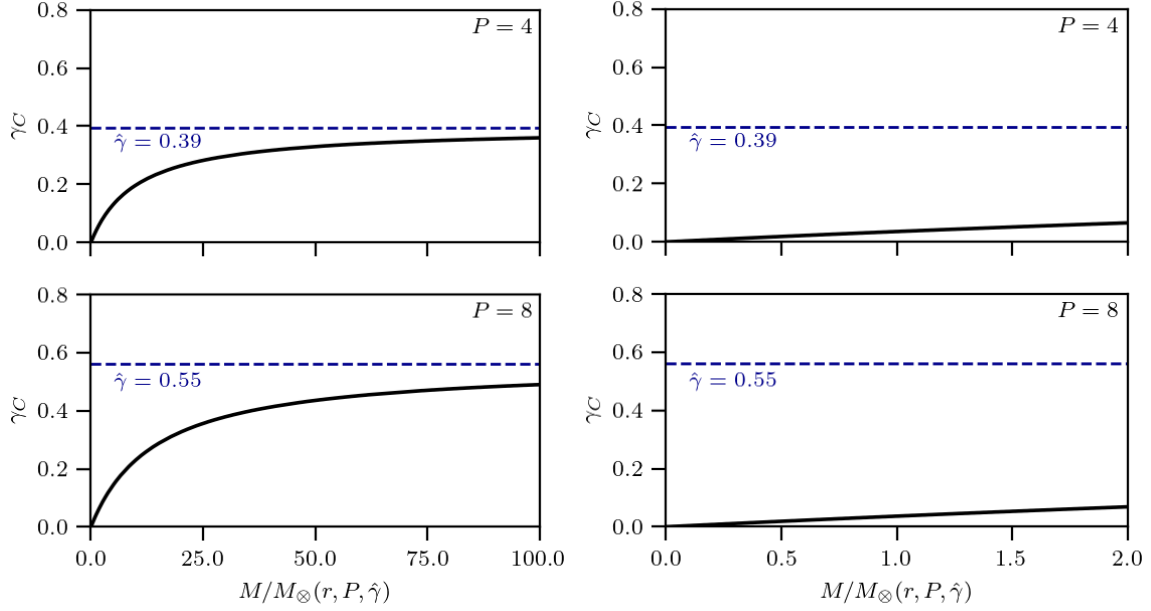
$$1 > 2\Theta^2 \left[ \rho + \gamma \frac{2}{P} (1 + \rho_P) \right]. \quad (6.14)$$

Setting the confidence level  $\Theta = 1/\sqrt{2}$ , which corresponds to the condition

$$\mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} [\xi_i^1 h_i^{(1)}(\boldsymbol{\xi}^1)] > \sqrt{\text{Var}[\xi_i^1 h_i^{(1)}(\boldsymbol{\xi}^1)]} \quad (6.15)$$

(namely we have a non-null magnetization in (6.8)), the previous relation determines a lower bound for  $M$  that is denoted as  $M_\otimes(r, P, \gamma)$ .

<sup>8</sup>Again, we have a sum of variables that are not Gaussian, but whose momenta are vanishing fast enough with  $N$  to make the CLT applicable. However, unlike the case discussed in Sec. 5, the overall sum is mathematically more treatable (because here the variables  $z_{\mu,a}$  are missing) and we can estimate directly first and second moments.



**Figure 9:** We numerically solve the self equations in the  $T \rightarrow 0$  limit, namely (5.30), for  $r = 0.2$ . In these four panels we plot the critical load  $\gamma_c$  vs the order of magnitude number of examples  $M$  w.r.t.  $M_\infty(r, P, \hat{\gamma})$ , where  $\hat{\gamma}$  is the critical load of the standard dense Hebbian network with the same interaction order, for different degrees of interaction  $P = 4, 8$ , at work with the simpler storing protocol. We notice that  $\gamma_c$  increases with  $M$  and, for  $M \gg M_\infty(r, P, \gamma = \hat{\gamma})$ , it saturates to  $\hat{\gamma}$  [6], that is represented by the horizontal dashed line.

## 6.2 Critical load and bounds for the dataset size

We now discuss a few special cases for  $M_\infty(r, P, \gamma)$  under the assumption  $r \ll 1$ . In the low-load regime  $\gamma = 0$ , the expression in (6.14) becomes

$$M > \left(\frac{1}{\sqrt{2}}\right)^2 \frac{(1-r^2)}{r^2} \sim \frac{1}{2} \frac{1}{r^2} \implies M_\infty(r, P, 0) = \frac{1}{2r^2} \quad (6.16)$$

where the last equality holds for  $r \ll 1$ .

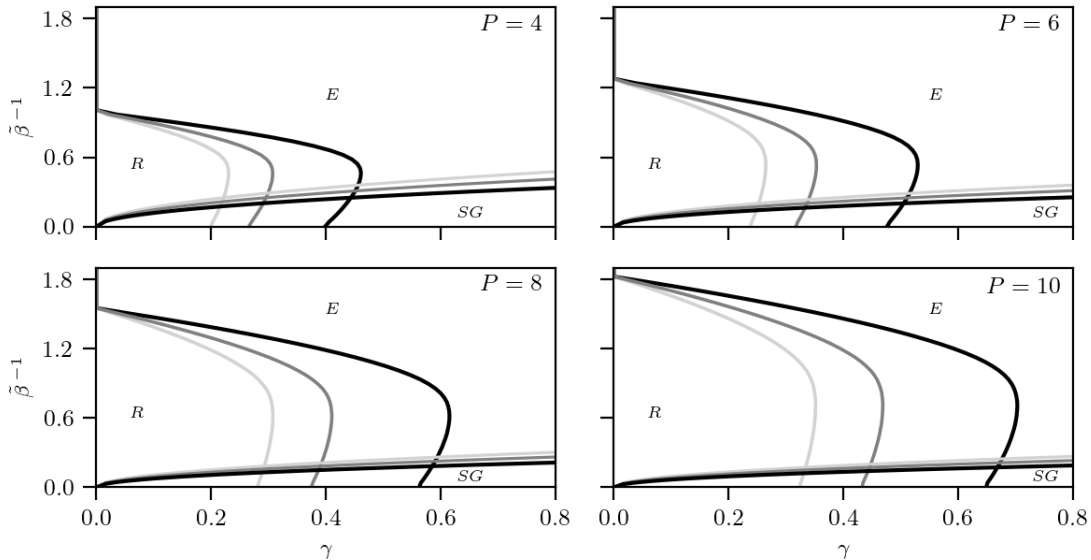
If  $\gamma \neq 0$  and  $P = 2$ , i.e. Hopfield classic case, as shown in [10]

$$M > \Theta^2 \left( \frac{(1-r^2)}{r^2} + \gamma \frac{1}{r^4} \right) \sim \frac{1}{2} \gamma \frac{1}{r^4} \implies M_\infty(r, 2, \gamma) = \gamma \frac{1}{2r^4}. \quad (6.17)$$

Finally, if  $\gamma \neq 0$  and  $P > 2$ , we have

$$M > \left(\frac{1}{\sqrt{2}}\right)^2 \left( \frac{(1-r^2)}{r^2} + \gamma \frac{2}{P} \frac{1}{r^{2P}} \right) \sim \frac{1}{2} \gamma \frac{2}{P} \frac{1}{r^{2P}} \implies M_\infty(r, P, \gamma) = \gamma \frac{1}{P} \frac{1}{r^{2P}}. \quad (6.18)$$

This result implies that we need a larger number of examples if we use dense networks w.r.t. Hopfield pairwise networks, further, we stress that -whatever the case- we always end up with power laws thresholds for learning relating the critical amount of examples to the dataset noise.



**Figure 10:** Each panel represents the phase diagram of the dense Hebbian networks trained -with no supervision- at different values of  $P$ , as specified. In any case we set  $r = 0.2$  and we compare outcomes for  $M = M_{\otimes}(r, P, \gamma)$  (light gray) and  $M = 2M_{\otimes}(r, P, \gamma)$  (dark gray), where  $M_{\otimes}(r, P, \gamma)$  is given by the expression in equation (6.18). We notice that, as  $P$  increases, the transition lines approach those of the dense Hebbian storage limit (black solid line). Moreover, the instability region, caused by the overlap between retrieval and spin-glass regions, decreases as  $P$  increases. Note that the recess of the maximal storage as the temperature goes to zero is a signature of replica symmetry breaking, that is not addressed here (see [11]).

We notice that if in (5.24) we set either  $r \rightarrow 1$  and  $M \rightarrow 1$  (i.e., we give the original patterns and not the examples, see Eq. (2.6), and, so, we have  $n_{1,a} = m_1$  in Def. (3)), or  $M \gg M_{\otimes}(r, P, \gamma)$  (i.e., we give the network a very large number of examples), we recover dense Hebbian neural network at work with the simpler storing protocol.

The results obtained analytically in this section are corroborated by numerical simulations and the related outputs are collected in Figs. 7-10. We stress that, to avoid the computational-expensive updating of the synaptic tensor in (3.1), we implemented a Plefka’s dynamics in our MC scheme. This is an effective dynamics that allows us to keep track of the evolution of the network’s order parameters at the level of their mean values: we refer to Appendix B for more details. Let us now comment these numerical results.

A corroboration of the goodness of Plefka’s dynamics is given in Fig. 7, where we show a comparison between MC simulation with Plefka’s dynamics and stability analysis.

Figure 8 shows evidence that, for a given choice of  $N, K, M$ , and  $r$ , when the degree of interaction  $P$  increases, the network is no longer able to generalise the archetype from the examples. This suggests that, if we want to retain the reconstruction capabilities, the number of examples  $M$  should scale with  $P$ , as showed in (6.18).

Tying in with this speech, in Fig. 9 we plot the number of examples  $M$  w.r.t. the critical load  $\gamma_c$ , for different values of  $P$ . We recall that the critical load  $\gamma_c$  is the load beyond which a black-out scenario



emerges, namely  $\lim_{\gamma \rightarrow \gamma_c^-} \bar{m} \neq 0$  and  $\lim_{\gamma \rightarrow \gamma_c^+} \bar{m} = 0$ . The black dotted line represents the critical load of the Hopfield dense neural network as a reference. We can see that, when  $M$  is chosen following the prescription in (6.18), we can reach the performances of Hopfield dense network.

Finally, in Fig. 10 we show the phase diagrams in the space  $(\tilde{\beta}, \gamma)$  for different values of  $P$  (each corresponding to a different panel). Interestingly, as  $M$  increases, the retrieval zone gets wider.

## 7 Conclusion and outlooks

In this paper we investigated the information processing capabilities of dense Hebbian networks endowed with couplings stemmed by an unsupervised protocol for learning. In order to have a mathematically tractable theory, this is developed for random structureless datasets. The network is made of  $N$  neurons that interact in groups of  $P$  units with a strength encoded by the synaptic tensor  $\mathbf{J}^{(unsup)}$ , whose generic entry (retaining only the leading order) reads as

$$J_{i_1 i_2 \dots i_P}^{(unsup)} \sim \frac{1}{MN^{P-1}} \sum_{\mu=1}^K \sum_{a=1}^M \eta_{i_1}^{\mu,a} \eta_{i_2}^{\mu,a} \dots \eta_{i_P}^{\mu,a}, \quad (7.1)$$

where  $\{\boldsymbol{\eta}^\mu\}_{\mu=1, \dots, K}$  is the dataset available, made of  $K$  subsets of examples (labeled by  $a$ ) referred to  $K$  unknown archetypes. The quality of the dataset, namely how “far” these examples are, in the average, from the related archetype, is ruled by  $r$  (such that by setting  $M = 1$  and  $r = 1$  we recover the standard dense Hebbian network under the simpler storage prescription [11, 21, 25]).

Hereafter we summarize the main outcomes of our work. As far as general neural network’s theory is concerned

1. The dense Hebbian network under the simpler storage prescription is well-known to be able to store a number of patterns that grows as  $K \sim N^{P-1}$ . This *high-load regime*, is preserved when the standard Hebbian coupling is replaced by the unsupervised Hebbian coupling. One can still introduce a load  $\alpha_{P-1} = \lim_{N \rightarrow \infty} \frac{K}{N^{P-1}}$  and determine a critical value beyond which a black-out scenario emerges: notably, this value does not depend on the dataset properties and it is solely a network’s characteristic.
2. For a correct learning -and subsequent retrieval- of the archetype, there exists a threshold value  $M_\otimes$  to overcome and this scales as  $M_\otimes \propto 1/(P r^{2P})$ . Thus, when using these dense machines in the unsupervised regime, one can actually reconstruct up to  $K \sim N^{P-1}$  archetypes only under the condition of a suitably large number of required examples available. In other words, increasing the number of retrievable patterns by a factor  $N$  (which means increasing the interaction order of one unit) requires an amplification in the number of examples per archetype of about  $1/r^2$ . Increasing the dataset size beyond  $M_\otimes$  leads to a wider retrieval region, namely to a larger critical load and to a larger critical temperature. The large cost in terms of available data is a peculiarity of the unsupervised regime, in fact, in supervised dense networks,  $M_\otimes$  does not scale with  $P$ , see [1].
3. There is another intriguing feature displayed by dense networks that is preserved in the unsupervised regime. Indeed, the reconstruction is feasible also in the presence of an extensive noise affecting its coupling and yielding to a signal-to-noise ratio increasing algebraically with  $N$ . Again, to mitigate the effects of this noise one has to move to a low-load regime in order to generate redundancy in the information allocated in the coupling tensor.

As far as the computational and mathematical technicalities are concerned

1. in this dense scenario, the post-synaptic potential does not have a Gaussian shape and standard techniques (e.g. replica trick, interpolation approaches) do not work straightforwardly, however, it is possible to adapt them by applying the CLT and restoring an effective Gaussian framework whose validity is corroborated by numerical simulations.
2. as  $P$  grows, the synaptic tensors become very expensive to evaluate (and prohibitive to be updated during learning dynamics): to overcome this problem, we adapted the Plefka's approximation to the case, resulting in a remarkable speed up of the simulations, yet preserving an extremely good accuracy in the results.

Overall these technical extensions are of broad generality and can be applied to several other neural networks.

## A Proof of Proposition 1

In order to prove Proposition 1, we need the following

**Lemma 1.** *The  $t$  derivative of the interpolating quenched pressure (5.8) is given by*

$$\begin{aligned} \frac{d\mathcal{A}_{N,K,M,r,\beta}^{(P)}}{dt} &:= \frac{\beta'}{2M}(1+\rho)^{P/2} \sum_{a=1}^M \left( \langle n_{1,a}^P \rangle_t - \frac{2M\psi}{\beta'(1+\rho)^{P/2}} \langle n_{1,a} \rangle_t \right) \\ &\quad - \frac{A^2}{2} \left( 1 - \langle q_{12} \rangle_t \right) + \frac{\beta'^2(1+\rho_P)}{4(1+\rho)^P} \frac{KP!}{2N^{P-1}} \left( 1 - \langle q_{12} \rangle_t \right). \end{aligned} \quad (\text{A.1})$$

where we use  $\rho_P = \frac{1-r^{2P}}{Mr^{2P}}$ .

*Proof.* Deriving Eq. (5.8) with respect to  $t$ , we get

$$\begin{aligned} \frac{d\mathcal{A}_{N,K,M,r,\beta}^{(P)}}{dt} &= \frac{1}{N} \mathbb{E} \frac{1}{\mathcal{Z}_{N,K,M,r,\beta}^{(P)}} \sum_{\sigma} \mathcal{B}_{N,K,M,r,\beta}^{(P)} \left[ \frac{\beta' N}{2M} (1+\rho)^{P/2} \sum_{a=1}^M n_{1,a}^P - \psi N \sum_{a=1}^M n_{1,a} \right. \\ &\quad \left. - \frac{1}{2\sqrt{1-t}} A \sum_{i=1}^N Y_i \sigma_i + \frac{1}{2} \frac{\beta' P! \sqrt{(1+\rho_P)}}{2t(1+\rho)^{P/2} N^P} \sum_{\mu \geq 2}^K \sum_{(i_1, \dots, i_P)}^{N, \dots, N} \lambda_{i_1 \dots i_P}^\mu \sigma_{i_1} \dots \sigma_{i_P} \right] \\ &= \frac{\beta'}{2M} (1+\rho)^{P/2} \sum_{a=1}^M \langle n_{1,a}^P \rangle_t - \psi \sum_{a=1}^M \langle n_{1,a} \rangle_t + D_1 + D_2. \end{aligned} \quad (\text{A.2})$$

Now, using Stein's lemma<sup>9</sup> on the random variables  $Y_i$  and  $\lambda_{i_1 \dots i_P}^\mu$ , we may rewrite the last two terms of (A.2) as

<sup>9</sup>This lemma, also known as Wick's theorem, applies to standard Gaussian variables, say  $J \sim \mathcal{N}(0,1)$ , and states that, for a generic function  $f(J)$  for which the two expectations  $\mathbb{E}(Jf(J))$  and  $\mathbb{E}(\partial_J f(J))$  both exist, then

$$\mathbb{E}(Jf(J)) = \mathbb{E} \left( \frac{\partial f(J)}{\partial J} \right). \quad (\text{A.3})$$

$$D_1 = -\frac{1}{2N\sqrt{1-t}}B \sum_{i=1}^N \mathbb{E} \partial_{Y_i} \left[ \frac{1}{\mathcal{Z}_{N,K,M,r,\beta}^{(P)}} \sum_{\sigma} \mathcal{B}_{N,K,M,r,\beta}^{(P)} \sigma_i \right] = -\frac{A^2}{2} (1 - \langle q_{12} \rangle_t), \quad (\text{A.4})$$

$$\begin{aligned} D_2 &= \frac{\beta' \sqrt{(1+\rho_P)}}{4t(1+\rho)^{P/2} 2N^{P+1}} \sum_{\mu \geq 2}^K \sum_{(i_1, \dots, i_P)}^{N, \dots, N} \mathbb{E} \partial_{\lambda_{i_1, \dots, i_P}^\mu} \left[ \frac{1}{\mathcal{Z}_{N,K,M,r,\beta}^{(P)}} \sum_{\sigma} \mathcal{B}_{N,K,M,r,\beta}^{(P)} \sigma_{i_1} \cdots \sigma_{i_P} \right] \\ &= \frac{\beta'^2 (1+\rho_P)}{4(1+\rho)^P} \frac{KP!}{2N^{P-1}} (1 - \langle q_{12} \rangle_t). \end{aligned} \quad (\text{A.5})$$

Rearranging together (A.4) and (A.5) we obtain the thesis.  $\square$

**Assumption 1.** *As a consequence of the RS assumption, for the generic order parameter  $X$ , being  $\Delta X := X - \bar{X}$ , the deviation w.r.t. the expectation value, then*

$$\langle (\Delta X)^2 \rangle_t \xrightarrow{N \rightarrow \infty} 0$$

and, clearly, the RS approximation also implies that, in the thermodynamic limit,  $\langle \Delta X \Delta Y \rangle_t \rightarrow 0$  for any generic pair of order parameters  $X, Y$ . Moreover in the thermodynamic limit, we have  $\langle (\Delta X)^k \rangle_t \rightarrow 0$  for  $k \geq 2$ .

Hereafter, in order to lighten the notation, we will drop the subscript  $t$ . In the following we can use the relation

$$\langle x^P \rangle - P \bar{x}^{P-1} \langle x \rangle = -(P-1) \bar{x}^P + \sum_{k=2}^P \binom{P}{k} \langle (x - \bar{x})^k \rangle \bar{x}^{P-k}, \quad (\text{A.6})$$

for any order parameter  $x$  with equilibrium value  $\bar{x}$ , which is computed straightforwardly by Newton's binomial [12].

Using these relations, if we fix the constants  $\psi, A$ , appearing in the interpolating partition function introduced in Definition 4, as

$$\psi = \beta' \frac{P}{2M} (1+\rho)^{P/2} \bar{n}^{P-1}, \quad A^2 = \frac{\beta'^2 (1+\rho_P)}{(1+\rho)^P} \frac{P}{2} \frac{KP!}{2N^{P-1}} \bar{q}^{P-1}, \quad (\text{A.7})$$

we can rewrite the derivative of the interpolating pressure w.r.t.  $t$  as

$$\begin{aligned} \frac{d\mathcal{A}_{N,K,M,r,\beta}^{(P)}}{dt} &:= -\frac{\beta'}{2} (P-1) (1+\rho)^{P/2} \bar{n}^P + \frac{\beta'^2 (1+\rho_P)}{4(1+\rho)^P} \frac{KP!}{2N^{P-1}} (1 - P\bar{q}^{P-1} + (P-1)\bar{q}^P) \\ &+ \frac{\beta'}{2M} (1+\rho)^{P/2} \sum_{a=1}^M \sum_{k=2}^P \binom{P}{k} \langle (n_{1,a} - \bar{n})^k \rangle \bar{n}^{P-k} \\ &+ \frac{\beta'^2 (1+\rho_P)}{4(1+\rho)^P} \frac{KP!}{2N^{P-1}} \sum_{k=2}^P \binom{P}{k} \langle (q_{12} - \bar{q})^k \rangle \bar{q}^{P-k}. \end{aligned} \quad (\text{A.8})$$

*Proof.* (of Proposition 1) Exploiting the fundamental theorem of calculus, we can relate  $\mathcal{A}_{N,K,M,\beta,r}^{(P)}(t = 1)$  and  $\mathcal{A}_{N,K,M,\beta,r}^{(P)}(t = 0)$  as

$$\mathcal{A}_{N,K,M,\beta,r}^{(P)} = \mathcal{A}_{N,K,M,r,\beta}^{(P)}(t=1) = \mathcal{A}_{N,K,M,r,\beta}^{(P)}(t=0) + \int_0^1 \partial_s \mathcal{A}_{N,K,M,r,\beta}^{(P)}(s) \Big|_{s=t} dt. \quad (\text{A.9})$$

We have just computed the derivative w.r.t.  $t$ , which is (A.8); all we need is to recover the one-body term:

$$\begin{aligned} \mathcal{A}_{N,K,M,r,\beta}^{(P)}(t=0) = \mathbb{E} \left\{ \ln 2 \cosh \left[ J\xi^1 + \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} \hat{\eta} \right. \right. \\ \left. \left. + Y \sqrt{\frac{\beta'^2 (1+\rho_P) P}{(1+\rho)^P} \frac{KP!}{2N^{P-1}} \bar{q}^{P-1}} \right] \right\}. \end{aligned} \quad (\text{A.10})$$

Putting (A.8) and (A.10) in (A.9), we find

$$\begin{aligned} \mathcal{A}_{N,K,M,r,\beta}^{(P)} = \mathbb{E} \left\{ \ln 2 \cosh \left[ J\xi^1 + \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} \hat{\eta} + Y \sqrt{\frac{\beta'^2 (1+\rho_P) P}{(1+\rho)^P} \frac{KP!}{2N^{P-1}} \bar{q}^{P-1}} \right] \right\} \\ - \frac{\beta'}{2} (P-1) (1+\rho)^{P/2} \bar{n}^P + \frac{\beta'^2 (1+\rho_P) KP!}{4(1+\rho)^P 2N^{P-1}} (1 - P\bar{q}^{P-1} + (P-1)\bar{q}^P) + \int_0^1 V_{N,M}(t) dt. \end{aligned} \quad (\text{A.11})$$

where  $\mathbb{E} = \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} \mathbb{E}_Y$ ,  $\hat{\eta} = \frac{1}{rM} \sum_{a=1}^M \eta^{1,a}$  and the potential  $V_{N,M}(t)$  is

$$V_{N,M}(t) = \frac{\beta' (1+\rho)^{P/2}}{2M} \sum_{a,k=1}^{M,P} \binom{P}{k} \langle (n_{1,a} - \bar{n})^k \rangle \bar{n}^{P-k} + \frac{\beta'^2 (1+\rho_P) KP!}{8(1+\rho)^P N^{P-1}} \sum_{k=2}^P \binom{P}{k} \langle (q_{12} - \bar{q})^k \rangle \bar{q}^{P-k} \quad (\text{A.12})$$

Now, we know that for  $b = P-1$  we have  $K = \alpha_{P-1} N^{P-1} + O(N^{P-1-\epsilon})$ ,  $\epsilon > 0$ ; thus, neglecting the lower terms, we have

$$\begin{aligned} \mathcal{A}_{\alpha_{P-1},M,r,\beta}^{(P)} = \mathbb{E} \left\{ \ln 2 \cosh \left[ J\xi^1 + \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} \hat{\eta} + Y \sqrt{\alpha_{P-1} \frac{P! \beta'^2 (1+\rho_P) P}{2(1+\rho)^P} \bar{q}^{P-1}} \right] \right\} \\ - \frac{\beta'}{2} (P-1) (1+\rho)^{P/2} \bar{n}^P + \alpha_{P-1} \frac{P! \beta'^2 (1+\rho_P)}{2 \cdot 4(1+\rho)^P} (1 - P\bar{q}^{P-1} + (P-1)\bar{q}^P). \end{aligned} \quad (\text{A.13})$$

Finally, we maximise the statistical pressure in (A.13) w.r.t. the order parameters and we find

$$\begin{aligned} \bar{n} = \frac{1}{1+\rho} \mathbb{E} \left\{ \tanh \left[ \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} \hat{\eta} + Y \sqrt{\alpha_{P-1} \frac{P! \beta'^2 (1+\rho_P) P}{2(1+\rho)^P} \bar{q}^{P-1}} \right] \hat{\eta} \right\}, \\ \bar{q} = \mathbb{E} \left\{ \tanh^2 \left[ \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} \hat{\eta} + Y \sqrt{\alpha_{P-1} \frac{P! \beta'^2 (1+\rho_P) P}{2(1+\rho)^P} \bar{q}^{P-1}} \right] \right\}. \end{aligned} \quad (\text{A.14})$$

Putting the definition of P-independent load from (3.10) in (A.14) we reach the thesis.  $\square$

**Corollary 1.** *In the large dataset limit, in the high-storage regime, the RS self-consistency equations can be expressed as*

$$\begin{aligned}\bar{n} &= \frac{\bar{m}}{1+\rho} + \beta' \frac{P}{2} \rho (1+\rho)^{P/2-2} (1-\bar{q}) \bar{n}^{P-1}. \\ \bar{q} &= \mathbb{E}_Z [\tanh^2 g(\beta, Z, \bar{n})]. \\ \bar{m} &= \mathbb{E}_Z [\tanh g(\beta, Z, \bar{n})].\end{aligned}\tag{A.15}$$

where

$$g(\beta, Z, \bar{n}) = \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} + \beta' Z \sqrt{\rho \frac{P^2}{4} \bar{n}^{2P-2} (1+\rho)^{P-2} + \alpha_{P-1} \frac{P!}{2} \frac{(1+\rho_P)}{(1+\rho)^P} \frac{P}{2} \bar{q}^{P-1}}.\tag{A.16}$$

*Proof.* In large dataset limit we can use the CLT so that we have

$$\mathbb{E}_{(\eta|\xi)}[\hat{\eta}] = \xi^1, \quad \mathbb{E}_{(\eta|\xi)}[(\hat{\eta})^2] - \left(\mathbb{E}_{(\eta|\xi)}[\hat{\eta}]\right)^2 = \rho(\xi^1)^2\tag{A.17}$$

thus we get

$$\hat{\eta} \sim \xi^1 (1 + \lambda \sqrt{\rho}).\tag{A.18}$$

where  $\lambda$  is a standard Gaussian variable  $\lambda \sim \mathcal{N}(0, 1)$ . Now, replacing (A.18) in the self-consistency equation for  $\bar{n}$  in (5.12), applying Stein's lemma and exploiting the self-consistency equations for  $\bar{m}$  and  $\bar{q}$  in (5.12) we get (A.15). Replacing this new expression of  $\bar{n}$  in the argument of the hyperbolic tangent of (5.12) and exploiting the parity of the hyperbolic tangent, we can explicitly compute the mean over  $\xi$

$$\begin{aligned}\bar{n} &= \frac{1}{1+\rho} \mathbb{E}_{\xi, \lambda, Y} \left\{ \tanh \left[ \beta' \left( \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} (1 + \lambda \sqrt{\rho}) \xi^1 + Y \sqrt{\frac{\alpha_{P-1} P! (1+\rho_P)}{2(1+\rho)^P} \frac{P}{2} \bar{q}^{P-1}} \right) \right] (1 + \lambda \sqrt{\rho}) \xi^1 \right\} \\ &= \frac{1}{1+\rho} \mathbb{E}_{\lambda, Y} \left\{ \tanh \left[ \beta' \left( \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1} (1 + \lambda \sqrt{\rho}) + Y \sqrt{\alpha_{P-1} \frac{P!}{2} \frac{(1+\rho_P)}{(1+\rho)^P} \frac{P}{2} \bar{q}^{P-1}} \right) \right] (1 + \lambda \sqrt{\rho}) \right\}\end{aligned}\tag{A.19}$$

Now we use the relation

$$\mathbb{E}_{\lambda, Y} [F(a_1 + \lambda a_2 + Y a_3)] = \mathbb{E}_Z \left[ F \left( a_1 + Z \sqrt{a_2^2 + a_3^2} \right) \right],\tag{A.20}$$

with  $\lambda, Y$  and  $Z$  i.i.d. Gaussian random variables,  $\lambda, Y, Z \sim \mathcal{N}(0, 1)$  and we have put  $F(a_1 + \lambda a_2 + Y a_3) = \tanh(a_1 + \lambda a_2 + Y a_3)$ ,  $a_1 = \beta' \frac{P}{2} \bar{n}^{P-1} (1+\rho)^{P/2-1}$ ,  $a_2 = a_1 \sqrt{\rho}$ ,  $a_3 = \beta' \sqrt{\alpha_{P-1} \frac{P!}{2} \frac{\beta'^2 (1+\rho_P)}{(1+\rho)^P} \frac{P}{2} \bar{q}^{P-1}}$

In this way we can reduce the number of Gaussian averages to a single one and reach the thesis.  $\square$

**Corollary 2.** *The self-consistency equations of the dense Hebbian neural networks in the unsupervised setting in the high-storage regime and in null-temperature limit  $\beta \rightarrow \infty$  are*

$$\begin{aligned}\bar{m} &= \operatorname{erf} \left[ \frac{P}{2} \frac{\bar{m}^{P-1}}{G} \right], \quad \bar{q} = 1, \\ G &= \sqrt{2 \left[ \rho \left( \frac{P}{2} \bar{m}^{P-1} \right)^2 + \alpha_{P-1} \frac{P!}{2} (1+\rho_P) \frac{P}{2} \right]}.\end{aligned}\tag{A.21}$$

*Proof.* For this proof it is convenient to introduce an additional term  $\tilde{\beta}x$  in the argument of the hyperbolic tangent ( $g(\tilde{\beta}, Z, \tilde{m})$ ) in (A.16)

$$\begin{aligned}\bar{q} &= \mathbb{E}_Z \left[ \tanh^2 \left( \tilde{\beta} \frac{P}{2} \tilde{m}^{P-1} + \tilde{\beta} Z \sqrt{\rho \left( \frac{P}{2} \tilde{m}^{P-1} \right)^2 + \alpha_{P-1} \frac{P!}{2} (1 + \rho_P) \frac{P}{2} \bar{q}^{P-1} + \tilde{\beta} x} \right) \right]. \\ \bar{m} &= \mathbb{E}_Z \left[ \tanh \left( \tilde{\beta} \frac{P}{2} \tilde{m}^{P-1} + \tilde{\beta} Z \sqrt{\rho \left( \frac{P}{2} \tilde{m}^{P-1} \right)^2 + \alpha_{P-1} \frac{P!}{2} (1 + \rho_P) \frac{P}{2} \bar{q}^{P-1} + \tilde{\beta} x} \right) \right].\end{aligned}\tag{A.22}$$

Also, we notice that, as  $\tilde{\beta} \rightarrow \infty$ , in the previous equations  $\bar{q} \rightarrow 1$ , thus in order to correctly perform the limit we introduce the reparametrization

$$\bar{q} = 1 - \frac{\delta\bar{q}}{\tilde{\beta}} \quad \text{as } \tilde{\beta} \rightarrow \infty.\tag{A.23}$$

Using (A.23) in (A.22) we obtain

$$\begin{aligned}\bar{m} &= \mathbb{E}_Z \left[ \tanh \left( \tilde{\beta} \frac{P}{2} \bar{m}^{P-1} + \tilde{\beta} Z \sqrt{\rho \left( \frac{P}{2} \bar{m}^{P-1} \right)^2 + \alpha_b \alpha_{P-1} \frac{P!}{2} (1 + \rho_P) \frac{P}{2} \left( 1 - \frac{\delta\bar{q}}{\tilde{\beta}} \right)^{P-1} + \tilde{\beta} x} \right) \right]. \\ 1 - \frac{\delta\bar{q}}{\tilde{\beta}} &= \mathbb{E}_Z \left[ \tanh^2 \left( \tilde{\beta} \frac{P}{2} \bar{m}^{P-1} + \tilde{\beta} Z \sqrt{\rho \left( \frac{P}{2} \bar{m}^{P-1} \right)^2 + \alpha_{P-1} \frac{P!}{2} (1 + \rho_P) \frac{P}{2} \left( 1 - \frac{\delta\bar{q}}{\tilde{\beta}} \right)^{P-1} + \tilde{\beta} x} \right) \right].\end{aligned}\tag{A.24}$$

Taking advantage of the new parameter  $x$ , we can recast the last equation in  $\delta\bar{q}$  as a derivative of the magnetization  $\bar{m}$  :

$$\frac{\partial \bar{m}}{\partial x} = \tilde{\beta} \left[ 1 - \left( 1 - \frac{\delta\bar{q}}{\tilde{\beta}} \right) \right] = \delta\bar{q}.\tag{A.25}$$

Thanks to this correspondence between the self equation for  $\bar{m}$  and the one for  $\delta\bar{q}$ , we can focus only on the self equation for  $\bar{m}$  and we can proceed with the  $\tilde{\beta} \rightarrow \infty$ . Thus, as  $\tilde{\beta} \rightarrow \infty$ , we have

$$\bar{m} = \mathbb{E}_Z \left[ \text{sign} \left( \frac{P}{2} \bar{m}^{P-1} + Z \sqrt{\rho \left( \frac{P}{2} \bar{m}^{P-1} \right)^2 + \alpha_{P-1} \frac{P!}{2} (1 + \rho_P) \frac{P}{2}} \right) \right]\tag{A.26}$$

$$\bar{q} \rightarrow 1.$$

Where we have restored  $x$  to zero. Finally, we can rearrange (A.26) using the relation

$$\mathbb{E}_z \text{sign}[b_1 + z b_2] = \text{erf} \left[ \frac{b_1}{\sqrt{2} b_2} \right],\tag{A.27}$$

where  $b_1 = \frac{P}{2} \bar{m}^{P-1}$  and  $b_2 = \sqrt{\rho \left( \frac{P}{2} \bar{m}^{P-1} \right)^2 + \alpha_{P-1} \frac{P!}{2} (1 + \rho_P) \frac{P}{2}}$ , hence reaching the thesis.  $\square$

## B Plefka's expansion of the effective Gibbs potential

When dealing with dense networks, MC simulations become prohibitively slow due to the computationally expensive update of their synaptic tensor (see the cost function, Eq. 3.1) and, clearly, the

higher the interaction order  $P$  the slower their convergence. To speed up simulations an alternative route where these computations can be avoided should be walked. Implementing Plefka's dynamics within a MC scheme can be a way out as the latter is an effective dynamics that allows us to keep track of the evolution of the network order-parameters at the level of their mean values.

The purpose of this section is thus to follow the path developed in [38, 39], namely we switch from the free energy to the Gibbs potential via Legendre transform, then we compute the expansion of the Gibbs potential that allows us to write effective (coarse grained) dynamics for the mean of the Mattis magnetization and we use such an expression within the update rule of the MC iterations.

We split the following analysis in two parts: first, we show the computation in classic dense networks, then we generalize the approach for unsupervised learning.

### B.1 Plefka's effective dynamics for Hebbian storing

The system is described by its Hamiltonian

$$\mathcal{H}_{N,K}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\xi}; \mathbf{z}) = -\frac{1}{\sqrt{N^{P-1}}} \sum_{\mu} \sum_{i_1 \dots i_{P/2}} \xi_{i_1}^{\mu} \dots \xi_{i_{P/2}}^{\mu} \sigma_{i_1} \dots \sigma_{i_{P/2}} z_{\mu} \quad (\text{B.1})$$

which is the interaction part of the integral representation of the Hamiltonian of the dense Hebbian network in [11]. Moreover,  $\{z_{\mu}\}$  is an additional set of real Gaussian variable we have added to describe the model.

In order to use Plefka's dynamic, we introduce a control parameter  $\varphi$  and we defined a new Hamiltonian as

$$\mathcal{H}_{N,K}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\xi}; \mathbf{z}, \mathbf{h}, \tilde{\mathbf{h}}, \varphi) = \varphi \mathcal{H}_{N,K}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\xi}; \mathbf{z}) - \frac{1}{\sqrt{N^{P-1}}} \sum_i h_i \sigma_i - \sum_{\mu} \tilde{h}_{\mu} z_{\mu} \quad (\text{B.2})$$

in such a way that if  $\varphi = 0$  we have an Hamiltonian representing non-interacting units, whereas if  $\varphi = 1$  we have an Hamiltonian representing full-interacting units; for this reason  $\varphi$  is referred to as interaction strength. Moreover,  $\{h_i\}$  and  $\{\tilde{h}_i\}$  are external fields which act, respectively, on  $\boldsymbol{\sigma}$  and  $\mathbf{z}$ . Using the expression of the modified Hamiltonian (B.2) we write down the partition function as

$$\begin{aligned} \mathcal{Z}_{N,K,\beta}^{(P)}(\boldsymbol{\xi}; \mathbf{h}, \tilde{\mathbf{h}}, \varphi) &= \sum_{\boldsymbol{\sigma}} \int \prod_{\mu} dz_{\mu} \sqrt{\frac{\beta'}{2\pi}} \exp \left( -\varphi \frac{\beta'}{\sqrt{N^{P-1}}} \sum_{\mu} \sum_{i_1 \dots i_{P/2}} \xi_{i_1}^{\mu} \dots \xi_{i_{P/2}}^{\mu} \sigma_{i_1} \dots \sigma_{i_{P/2}} z_{\mu} \right. \\ &\quad \left. - \frac{\beta'}{2} \sum_{\mu} z_{\mu}^2 + \beta' \sum_{\mu} \tilde{h}_{\mu} z_{\mu} + \frac{\beta'}{\sqrt{N^{P-1}}} \sum_i h_i \sigma_i \right). \end{aligned} \quad (\text{B.3})$$

where  $\beta' = 2\beta/P!$ .

We now consider the Gibbs potential for this model, that is the Legendre transformation of the free energy constrained w.r.t. the magnetizations  $m_i = \langle \sigma_i \rangle$  and  $\langle z_{\mu} \rangle$  averaged w.r.t. Boltzmann distribution  $\propto \exp^{-\beta H(\varphi)}$ ; for our model the Gibbs potential is given by

$$\mathcal{G}_{N,K,\beta}^{(P)}(\boldsymbol{\xi}; \mathbf{z}, \mathbf{h}, \tilde{\mathbf{h}}, \varphi) = -\frac{1}{\beta'} \ln \mathcal{Z}(\varphi) + \sum_{\mu} \tilde{h}_{\mu} \langle z_{\mu} \rangle + \frac{1}{\sqrt{N^{P-1}}} \sum_i h_i m_i. \quad (\text{B.4})$$

Now, we shorten the notation as  $\mathcal{G}_{N,K,\beta}^{(P)}(\boldsymbol{\xi}; \mathbf{z}, \mathbf{h}, \tilde{\mathbf{h}}, \varphi) \equiv \mathcal{G}(\varphi)$ , and expand the last expression around  $\varphi = 0$  as

$$\mathcal{G}(\varphi) = \mathcal{G}(0) + \sum_{n=1}^{\infty} \frac{\varphi^n \mathcal{G}^{(n)}}{n!} \quad \mathcal{G}^{(n)} = \left. \frac{\partial^n \mathcal{G}(\varphi)}{\partial \varphi^n} \right|_{\varphi=0}. \quad (\text{B.5})$$

For our computations we stop the expansion to the first order and we start to find all the terms we need. The non-interacting Gibbs potential  $\mathcal{G}(0)$  reads as

$$\mathcal{G}(0) = -\frac{N}{\beta} \log 2 - \frac{1}{\beta'} \sum_i \log \cosh \left( \frac{\beta'}{\sqrt{N^{P-1}}} h_i \right) + \sum_{\mu} \tilde{h}_{\mu} \langle z_{\mu} \rangle + \frac{1}{\sqrt{N^{P-1}}} \sum_i h_i \langle \sigma_i \rangle - \frac{1}{2} \sum_{\mu} \tilde{h}_{\mu}^2. \quad (\text{B.6})$$

If we extremize  $\mathcal{G}(0)$  w.r.t. local fields, namely  $h_i$  and  $\tilde{h}_{\mu}$  for  $\mu = 1, \dots, K$ ,  $i = 1, \dots, N$ , we find expressions of them read as

$$h_i = \frac{\sqrt{N^{P-1}}}{\beta'} \tanh^{-1}(m_i) \implies h_i = \frac{\sqrt{N^{P-1}}}{2\beta'} \ln \left( \frac{1+m_i}{1-m_i} \right), \quad (\text{B.7})$$

$$\tilde{h}_{\mu} = \langle z_{\mu} \rangle; \quad (\text{B.8})$$

where we have use the relation

$$\tanh^{-1}(x) = \ln \frac{1+x}{1-x}.$$

Putting (B.7) and (B.8) in the non-interacting Gibbs potential (B.6) we get

$$\mathcal{G}(0) = -\frac{N}{\beta'} \log 2 + \frac{1}{2\beta'} \sum_i [(1-m_i) \log(1-m_i) + (1+m_i) \log(1+m_i)] + \frac{1}{2} \sum_{\mu} \langle z_{\mu} \rangle^2. \quad (\text{B.9})$$

Now we have found  $\mathcal{G}(0)$ , all we need is the first-order contribution which is

$$\left. \frac{\partial \mathcal{G}(\varphi)}{\partial \varphi} \right|_{\varphi=0} = -\frac{1}{\sqrt{N^{P-1}}} \sum_{i_1, \dots, i_{P/2}} \sum_{\mu} \xi_{i_1}^{\mu} \dots \xi_{i_{P/2}}^{\mu} \langle z_{\mu} \rangle m_{i_1} \dots m_{i_{P/2}}. \quad (\text{B.10})$$

Therefore the first-order expression of Gibbs potential for the full interacting system, namely for  $\varphi = 1$  is

$$\begin{aligned} \mathcal{G}(\varphi = 1) &= -\frac{N}{\beta'} \log 2 + \frac{1}{2\beta'} \sum_i [(1-m_i) \log(1-m_i) + (1+m_i) \log(1+m_i)] + \frac{1}{2} \sum_{\mu} \langle z_{\mu} \rangle^2 \\ &\quad - \frac{1}{\sqrt{N^{P-1}}} \sum_{i_1, \dots, i_{P/2}} \sum_{\mu} \xi_{i_1}^{\mu} \dots \xi_{i_{P/2}}^{\mu} \langle z_{\mu} \rangle m_{i_1} \dots m_{i_{P/2}}. \end{aligned} \quad (\text{B.11})$$

Extremizing (B.11) w.r.t.  $m_i$  and  $\langle z_{\mu} \rangle$ , we find the respective self-consistency equations:

$$\frac{\partial \mathcal{G}}{\partial m_i} = 0 \implies m_i = \tanh \left[ \beta' \frac{P}{2} \frac{1}{\sqrt{N^{P-1}}} \sum_{\mu} \xi_i^{\mu} \langle z_{\mu} \rangle \left( \sum_j \xi_j^{\mu} m_j \right)^{P/2-1} \right], \quad (\text{B.12})$$

$$\frac{\partial \mathcal{G}}{\partial \langle z_{\mu} \rangle} = 0 \implies \langle z_{\mu} \rangle = \frac{1}{\sqrt{N^{P-1}}} \left( \sum_i \xi_i^{\mu} m_i \right)^{P/2} \quad (\text{B.13})$$

These equations are then used “in tandem” to make the system evolve: starting from an initial configuration  $(\boldsymbol{\sigma}^{(0)}, \boldsymbol{z}^{(0)})$ , we evaluate the related  $m_i^{(0)}$  and  $z_{\mu}^{(0)}$  for any  $i$  and  $\mu$ , and we use them in



(B.12) to get  $m_i^{(1)}$ , the latter is then used in (B.13) to get  $z_\mu^{(1)}$ , and we proceed this way, bouncing from (B.12) to (B.13), up to thermalization. We stress that these equations allow to implement an effective dynamics that avoid the computation of spin configurations. In fact, this coarse grained dynamics only cares about the Boltzmann average of each spin direction, whose behaviour is given by (B.12). Even if at first glance (B.2) seems to require three set of auxiliary variables,  $\{z_\mu\}$ ,  $\{h_i\}$  and  $\{\tilde{h}_i\}$ , the extremization of the Gibbs potential at first order fixes the external fields  $\{h_i\}$  and  $\{\tilde{h}\}$ . The gaussian variables  $\{z_\mu\}$  act as latent dynamical variables that evolve according to (B.13). In such an iterative MC scheme these hidden degrees of freedom are suitably updated in order to effectively retrieve the pattern that constitutes the signal.

## B.2 Plefka's effective dynamics for unsupervised Hebbian learning

The system is described by the Hamiltonian

$$\mathcal{H}_{N,K,M,r}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\eta}; \mathbf{z}) = -\sqrt{\frac{1}{N^{P-1}\mathcal{R}^{P/2}M}} \sum_{\mu>1}^K \sum_{a=1}^M \left( \sum_{i_1, \dots, i_{P/2}}^{N, \dots, N} \eta_{i_1}^{\mu,a} \dots \eta_{i_{P/2}}^{\mu,a} \sigma_{i_1} \dots \sigma_{i_{P/2}} \right) z_{\mu,a} \quad (\text{B.14})$$

which is the Hamiltonian corresponding of the interacting term in (3.6). Moreover,  $\{z_{\mu,a}\}$  is an additional set of real variable computed by Gaussian distribution we have added to describe the model. Mirroring computations in Subsection B.1, in order to use Plefka's dynamic, we introduce a control parameter  $\varphi$  and we defined a new Hamiltonian as

$$\mathcal{H}_{N,K,M,r}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\eta}; \mathbf{z}, \mathbf{h}, \tilde{\mathbf{h}}, \varphi) = \varphi \mathcal{H}_{N,K,M}^{(P)}(\boldsymbol{\sigma}|\boldsymbol{\eta}; \mathbf{z}) - \sqrt{\frac{1}{N^{P-1}\mathcal{R}^{P/2}M}} \sum_i h_i \sigma_i - \sum_{\mu,a} \tilde{h}_{\mu,a} z_{\mu,a} \quad (\text{B.15})$$

where  $\varphi$  describes the interaction strength. We stress that if  $\varphi = 0$  we have the Hamiltonian of non-interacting terms. Moreover,  $\{h_i\}$  and  $\{\tilde{h}_{\mu,a}\}$  are external fields who act, respectively, on  $\boldsymbol{\sigma}$  and  $\mathbf{z}$ .

The expression of the modified Hamiltonian (B.15) can be used to write down the partition function as

$$\begin{aligned} \mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\eta}; \mathbf{h}, \tilde{\mathbf{h}}, \varphi) &= \sum_{\boldsymbol{\sigma}} \int \prod_{\mu,a} dz_{\mu,a} \sqrt{\frac{\tilde{\beta}}{2\pi}} \exp \left( -\frac{\tilde{\beta}}{2} \sum_{\mu,a} z_{\mu,a}^2 + \tilde{\beta} \sum_{\mu,a} \tilde{h}_{\mu,a} z_{\mu,a} \right. \\ &\left. + \frac{\tilde{\beta}}{\sqrt{N^{P-1}r^{2PM}}} \sum_i h_i \sigma_i - \varphi \frac{\tilde{\beta}}{\sqrt{N^{P-1}r^{2PM}}} \sum_{\mu,a} \sum_{i_1, \dots, i_{P/2}} \eta_{i_1}^{\mu,a} \dots \eta_{i_{P/2}}^{\mu,a} \sigma_{i_1} \dots \sigma_{i_{P/2}} z_{\mu,a} \right) \end{aligned} \quad (\text{B.16})$$

where we define  $\tilde{\beta}$  as in (5.26).

The Gibbs potential is defined as the Legendre transformation of the free energy constrained w.r.t. the magnetization  $m_i = \langle \sigma_i \rangle$  and  $\langle z_{\mu,a} \rangle$  averaged w.r.t. the Boltzmann distribution  $P(\boldsymbol{\sigma}) \sim \exp^{-\beta H(\varphi)}$  (B.4). We have

$$\mathcal{G}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\eta}; \mathbf{z}, \mathbf{h}, \tilde{\mathbf{h}}, \varphi) - \frac{1}{\tilde{\beta}} \ln \mathcal{Z}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\eta}; \mathbf{h}, \tilde{\mathbf{h}}, \varphi) + \sum_{\mu} \tilde{h}_{\mu} \langle z_{\mu} \rangle + \frac{1}{\sqrt{N^{P-1}r^{2PM}}} \sum_i h_i m_i \quad (\text{B.17})$$

and we write Plefka's expansion as in (B.5). Also in this case we stop the expansion at the first order. Thus, shortening the notation as  $\mathcal{G}_{N,K,M,r,\beta}^{(P)}(\boldsymbol{\eta}; \mathbf{z}, \mathbf{h}, \tilde{\mathbf{h}}, \varphi) \equiv G(\varphi)$ , we compute the non-interacting Gibbs potential  $\mathcal{G}(0)$  and the first order contribution  $\left. \frac{\partial \mathcal{G}(\varphi)}{\partial \varphi} \right|_{\varphi=0}$ . Let us start from  $\mathcal{G}(0)$ .

$$\begin{aligned} \mathcal{G}(0) &= -\frac{N}{\tilde{\beta}} \log 2 - \frac{1}{\tilde{\beta}} \sum_i \log \cosh \left( \frac{\tilde{\beta}}{\sqrt{N^{P-1} r^{2P} M}} h_i \right) + \sum_{\mu,a} \tilde{h}_{\mu,a} \langle z_{\mu,a} \rangle \\ &\quad + \frac{1}{\sqrt{N^{P-1} r^{2P} M}} \sum_i h_i \langle \sigma_i \rangle - \frac{1}{2} \sum_{\mu,a} \tilde{h}_{\mu,a}^2 \end{aligned} \quad (\text{B.18})$$

If we extremize  $\mathcal{G}(0)$  w.r.t. the local fields, namely  $h_i$  and  $\tilde{h}_{\mu,a}$  for  $i = 1, \dots, N$ ,  $a = 1, \dots, M$ , we can find their expressions, that read as

$$h_i = \frac{\sqrt{N^{P-1} r^{2P} M}}{2\tilde{\beta}} \log \left( \frac{1 + m_i}{1 - m_i} \right), \quad (\text{B.19})$$

$$\tilde{h}_{\mu,a} = \langle z_{\mu,a} \rangle; \quad (\text{B.20})$$

While the first-order derivative of Gibbs potential w.r.t.  $\varphi$  is

$$\left. \frac{\partial \mathcal{G}(\varphi)}{\partial \varphi} \right|_{\varphi=0} = -\frac{1}{\sqrt{N^{P-1} r^{2P} M}} \sum_{i_1, \dots, i_{P/2}} \sum_{\mu,a} \eta_{i_1}^{\mu,a} \dots \eta_{i_{P/2}}^{\mu,a} \langle z_{\mu,a} \rangle m_{i_1} \dots m_{i_{P/2}}. \quad (\text{B.21})$$

Therefore,  $\mathcal{G}(\varphi)$  is rewritten using Plefka's expansion as

$$\begin{aligned} \mathcal{G}(\varphi) &= -\frac{N}{\tilde{\beta}} \log 2 + \frac{1}{2\tilde{\beta}} \sum_i [(1 - m_i) \log(1 - m_i) + (1 + m_i) \log(1 + m_i)] \\ &\quad + \frac{1}{2} \sum_{\mu,a} \langle z_{\mu,a} \rangle^2 - \frac{\varphi}{\sqrt{N^{P-1} r^{2P} M}} \sum_{\mu,a} \left( \sum_i \eta_i^{\mu,a} m_i \right)^{P/2} \langle z_{\mu,a} \rangle. \end{aligned} \quad (\text{B.22})$$

To conclude, we can compute the self-consistence equations w.r.t.  $m_i$  and  $\langle z_{\mu,a} \rangle$  extremizing the first order expression of  $\mathcal{G}(\varphi = 1)$  read as

$$m_i = \tanh \left[ \tilde{\beta} \frac{P}{2} \frac{1}{\sqrt{N^{P-1} r^{2P} M}} \sum_{\mu,a} \eta_i^{\mu,a} \langle z_{\mu,a} \rangle \left( \sum_j \eta_j^{\mu,a} m_j \right)^{P/2-1} \right], \quad (\text{B.23})$$

$$\langle z_{\mu,a} \rangle = \frac{1}{\sqrt{N^{P-1} r^{2P} M}} \left( \sum_i \eta_i^{\mu,a} m_i \right)^{P/2} \quad (\text{B.24})$$

These equations are then used “in tandem” to make the system evolve, as explained in the previous subsection. This iterative MC updating scheme leaves the network free to arrange the hidden degrees of freedom  $\{z_{\mu,a}\}$  in such a way that the mean values of the neurons  $m_i$  converge to the correspondent element of the archetype vector, provided that the network is posed in the retrieval region of the phase diagram.

## C Evaluation of the momenta of the effective post-synaptic potential

The purpose of this section is to evaluate the first and second momenta of the expression  $\xi_i^1 h_i^{(1)}(\boldsymbol{\xi}^1)$ , that are referred to as, respectively,  $\mu_1$  and  $\mu_2$  and are used in Sec. 6.1; specifically,

$$\mu_1 := \mathbb{E}_{\xi} \mathbb{E}_{(\eta|\xi)} \left[ \xi_i^1 h_i^{(1)}(\boldsymbol{\xi}^1) \right] = \frac{1}{\mathcal{R}^{P/2} M N^{P-1}} \sum_{\mu,a=1}^{K,M} \sum_{(i_2, \dots, i_P)} \mathbb{E}_{\xi} \mathbb{E}_{(\eta|\xi)} \left[ (\xi_{i_1}^1 \dots \xi_{i_P}^1) \eta_{i_1}^{\mu,a} \dots \eta_{i_P}^{\mu,a} \right], \quad (\text{C.1})$$

$$\begin{aligned} \mu_2 := \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} \left[ \{h_i^{(1)}(\boldsymbol{\xi}^1)\}^2 \right] &= \frac{1}{N^{2P-2} \mathcal{R}^P M^2} \sum_{\mu, \nu=1}^K \sum_{a, b=1}^{M, M} \sum_{(i_2, \dots, i_P)} \sum_{(j_2, \dots, j_P)} \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} \left[ (\xi_{i_2}^1 \xi_{j_2}^1 \dots \xi_{i_P}^1 \xi_{j_P}^1) \right. \\ &\quad \left. \eta_{i_1}^{\mu, a} \eta_{i_1}^{\nu, b} \left( \eta_{i_2}^{\mu, a} \eta_{j_2}^{\nu, b} \dots \eta_{i_P}^{\mu, a} \eta_{j_P}^{\nu, b} \right) \right]. \end{aligned} \quad (\text{C.2})$$

As for  $\mu_1$ , using  $\mathbb{E}_{(\eta|\xi)}[\eta_i^{\mu, a}] = r \xi_i^\mu$ ,

$$\mu_1 = \frac{1}{\mathcal{R}^{P/2} M N^{P-1}} \sum_{\mu=1}^K \sum_{(i_2, \dots, i_P)} \mathbb{E}_\xi \left[ M r^P (\xi_{i_1}^1 \dots \xi_{i_P}^1) (\xi_{i_1}^\mu \dots \xi_{i_P}^\mu) \right], \quad (\text{C.3})$$

since  $\mathbb{E}_\xi[\xi_i^\mu] = 0$  the only non-zero terms are those with  $\mu = 1$  and the expression simplifies into

$$\mu_1 = \frac{r^P}{\mathcal{R}^{P/2} N^{P-1}} \sum_{(i_2, \dots, i_P)} \mathbb{E}_\xi \left[ (\xi_{i_1}^1 \dots \xi_{i_P}^1)^2 \right] = \frac{1}{(1 + \rho)^{P/2}}. \quad (\text{C.4})$$

As for  $\mu_2$ , since  $\mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)}[\eta_{i_1}^{\mu, a} \eta_{i_1}^{\nu, b}] = r^2 \delta^{\mu\nu}$ , the only non-zero terms are those where  $\mu = \nu$ , thus

$$\begin{aligned} \mu_2 &= \frac{1}{N^{2P-2} \mathcal{R}^P M^2} \sum_{\mu=1}^K \sum_{a, b=1}^{M, M} \sum_{(i_2, \dots, i_P)} \sum_{(j_2, \dots, j_P)} \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} \left[ (\xi_{i_2}^1 \xi_{j_2}^1 \dots \xi_{i_P}^1 \xi_{j_P}^1) \right. \\ &\quad \left. \eta_{i_1}^{\mu, a} \eta_{i_1}^{\mu, b} \left( \eta_{i_2}^{\mu, a} \eta_{j_2}^{\mu, b} \dots \eta_{i_P}^{\mu, a} \eta_{j_P}^{\mu, b} \right) \right] \\ &= A_{\mu=1} + B_{\mu>1}, \end{aligned} \quad (\text{C.5})$$

where, in the last line, we highlighted the contributions stemming from terms with, respectively,  $\mu = 1$  ( $A_{\mu=1}$ ) and  $\mu > 1$  ( $B_{\mu>1}$ ). These are evaluated hereafter:

$$\begin{aligned} A_{\mu=1} &= \frac{1}{N^{2P-2} \mathcal{R}^P M^2} \sum_{\mu=1}^K \sum_{(i_2, \dots, i_P)} \sum_{(j_2, \dots, j_P)} \left\{ \sum_{a=1}^M \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} \left[ (\xi_{i_2}^1 \xi_{j_2}^1 \dots \xi_{i_P}^1 \xi_{j_P}^1) \right. \right. \\ &\quad \left. \left. (\eta_{i_2}^{\mu, a} \eta_{j_2}^{1, a} \dots \eta_{i_P}^{1, a} \eta_{j_P}^{1, a}) \right] + \sum_{a \neq b}^M \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} \left[ (\xi_{i_2}^1 \xi_{j_2}^1 \dots \xi_{i_P}^1 \xi_{j_P}^1) \eta_{i_1}^{1, a} \eta_{i_1}^{1, b} \left( \eta_{i_2}^{\mu, a} \eta_{j_2}^{1, b} \dots \eta_{i_P}^{1, a} \eta_{j_P}^{1, b} \right) \right] \right\} \\ &= \frac{1}{N^{2P-2} \mathcal{R}^P M^2} \sum_{\mu=1}^K \sum_{(i_2, \dots, i_P)} \sum_{(j_2, \dots, j_P)} \mathbb{E}_\xi \left[ \sum_{a=1}^M r^{2P-2} (\xi_{i_2}^1 \xi_{j_2}^1 \dots \xi_{i_P}^1 \xi_{j_P}^1)^2 + \sum_{a \neq b}^M r^{2P} (\xi_{i_1}^1 \xi_{i_2}^1 \xi_{j_2}^1 \dots \xi_{i_P}^1 \xi_{j_P}^1)^2 \right] \\ &= \frac{r^{2P}}{\mathcal{R}^P M} [r^{-2} + (M-1)] = \frac{r^{2P}}{\mathcal{R}^P} \left[ 1 + \frac{1-r^2}{r^2 M} \right] = \left( \frac{1}{(1+\rho)^{P/2}} \right)^2 (1+\rho); \end{aligned} \quad (\text{C.6})$$

for  $B_{\mu>1}$ , splitting the case  $a = b$  and  $a \neq b$ , we get

$$\begin{aligned}
B_{\mu>1} &= \frac{1}{N^{2P-2}\mathcal{R}^P M^2} \sum_{\mu>1} \sum_{(i_2, \dots, i_P)} \sum_{(j_2, \dots, j_P)} \left\{ \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} \sum_{a=1}^M [(\xi_{i_2}^1 \xi_{j_2}^1 \dots \xi_{i_P}^1 \xi_{j_P}^1) \right. \\
&\quad \left. (\eta_{i_2}^{\mu,a} \eta_{j_2}^{\mu,a} \dots \eta_{i_P}^{\mu,a} \eta_{j_P}^{\mu,a})] + \mathbb{E}_\xi \mathbb{E}_{(\eta|\xi)} \sum_{a \neq b}^M (\xi_{i_2}^1 \xi_{j_2}^1 \dots \xi_{i_P}^1 \xi_{j_P}^1) \eta_{i_1}^{\mu,a} \eta_{i_1}^{\mu,b} (\eta_{i_2}^{\mu,a} \eta_{j_2}^{\mu,b} \dots \eta_{i_P}^{\mu,a} \eta_{j_P}^{\mu,b}) \right\} \\
&= \frac{1}{N^{2P-2}\mathcal{R}^P M^2} \sum_{\mu>1} \sum_{(i_2, \dots, i_P)} \sum_{(j_2, \dots, j_P)} \mathbb{E}_\xi \left\{ r^{2P-2} \sum_{a=1}^M (\xi_{i_2}^1 \xi_{j_2}^1 \dots \xi_{i_P}^1 \xi_{j_P}^1) (\xi_{i_2}^\mu \xi_{j_2}^\mu \dots \xi_{i_P}^\mu \xi_{j_P}^\mu) \right. \\
&\quad \left. + r^{2P} \sum_{a \neq b}^M (\xi_{i_2}^1 \xi_{j_2}^1 \dots \xi_{i_P}^1 \xi_{j_P}^1) (\xi_{i_2}^\mu \xi_{j_2}^\mu \dots \xi_{i_P}^\mu \xi_{j_P}^\mu) \right\}
\end{aligned} \tag{C.7}$$

as far as  $\mathbb{E}_\xi(\xi_i^\mu \xi_j^\mu) = \delta_{ij}$ , the only non-zero terms are the ones in which the sum over  $i$  and  $j$  will be equal in pairs:

$$\begin{aligned}
B_{\mu>1} &= \frac{(P-1)!}{N^{2P-2}\mathcal{R}^P M^2} \sum_{\mu>1} \sum_{(i_2, \dots, i_P)} \mathbb{E}_\xi \left\{ [r^{2P-2} M + r^{2P} M(M-1)] (\xi_{i_2}^1 \dots \xi_{i_P}^1)^2 (\xi_{i_2}^\mu \dots \xi_{i_P}^\mu)^2 \right\} \\
&= \frac{r^{2P}}{N^{P-1}\mathcal{R}^P} (P-1)! K \left( 1 + \frac{1 - r^{2P}}{Mr^{2P}} \right) = \frac{r^{2P}}{N^{P-1}\mathcal{R}^P} (P-1)! K (1 + \rho_P),
\end{aligned} \tag{C.8}$$

and, if we set  $K = \alpha_{P-1} N^{P-1}$  we have

$$B_{\mu>1} = \left( \frac{1}{(1+\rho)^{P/2}} \right)^2 (P-1)! \alpha_{P-1} (1 + \rho_P). \tag{C.9}$$

Putting together (C.6) and (C.9) we reach the explicit expression of  $\mu_2$ .

## D List of symbols (in alphabetical order)

- $\mathcal{A}$  is the statistical quenched pressure (i.e.  $\mathcal{A} = -\beta\mathcal{F}$ )
- $\alpha_b$  is the storage of the network defined as  $\alpha_b = \lim_{N \rightarrow +\infty} K/N^b$
- $\mathcal{B}(\boldsymbol{\sigma}; t)$  is the Boltzmann factor, defined as  $\mathcal{B}(\boldsymbol{\sigma}; t) = \exp[\beta H(\boldsymbol{\sigma}; t)]$
- $\beta \in \mathbb{R}^+$  is the (fast) noise in the network (such that for  $\beta \rightarrow 0$  the behavior of the network is a pure random walk while for  $\beta \rightarrow +\infty$  it is a steepest descent toward the minima)
- $\mathbb{E}$  denotes the average over all the (quenched) coupling variables
- $\boldsymbol{\eta} \in \{-1, +1\}^{K \times M}$  are the noisy examples, namely noisy versions of the archetypes  $\boldsymbol{\xi}^\mu$
- $\mathcal{F}$  is the free energy (i.e.  $\mathcal{F} = -\beta^{-1}\mathcal{A}$ )

- $\gamma$  is the  $P$  independent part of  $\alpha_b$ , namely  $\alpha_b = \gamma \frac{2}{P-1}$
- $\mathcal{H}$  is the cost function (or Hamiltonian) defining the model
- $K \in \mathbb{N}$  is the amount of archetypes  $\xi$  to learn and retrieve
- $N \in \mathbb{N}$  is the amount of neurons in the network, i.e. the network size
- $M \in \mathbb{N}$  is the amount of examples per archetype, i.e. the training set
- $m_\mu$  is the Mattis magnetization of the archetype  $\xi^\mu$  defined as  $\frac{1}{N} \sum_i \xi_i^\mu \sigma_i$
- $\bar{m}$  is the asymptotic value of the Mattis magnetization of the signal,  $m$ , in the thermodynamic limit, i.e.  $\lim_{N \rightarrow \infty} P(m_\mu) = \delta(m - \bar{m})$
- $n_{a,\mu}$  is the magnetization of the example  $\eta^{\mu,a}$  defined as  $\frac{1}{N} \sum_i \eta_i^{\mu,a} \sigma_i$
- $\bar{n}$  is the asymptotic value of the magnetization of each example related to the signal,  $n_{1,a}$ , in the thermodynamic limit, i.e.  $\lim_{N \rightarrow \infty} P(n_{1,a}) = \delta(n_{1,a} - \bar{n})$
- $\omega_t(O(\boldsymbol{\sigma}))$  is the generalized Boltzmann measure, namely  $\omega_t(O(\boldsymbol{\sigma})) = \frac{1}{\mathcal{Z}_N} \sum_{\boldsymbol{\sigma}} O(\boldsymbol{\sigma}) \mathcal{B}(\boldsymbol{\sigma}; t)$
- $w$  is the noise added to synaptic tensor  $\boldsymbol{\eta}$  defined as  $w = \tau N^\delta$  where  $\delta = \frac{(P-1)-b}{2}$
- $P$  is the degree of interaction among neurons in the network (e.g.  $P = 2$  is the standard pairwise scenario)
- $q_{lm}$  is the overlap among two replicas defined as  $\frac{1}{N} \sum_i \sigma_i^{(l)} \sigma_i^{(m)}$
- $\bar{q}$  is the asymptotic value of  $q_{lm}$  in the thermodynamic limit and under the replica symmetric assumption, i.e.  $\lim_{N \rightarrow \infty} P(q_{lm}) = \delta(q_{lm} - \bar{q})$
- $\mathcal{R}$  is defined as  $\mathcal{R} = r^2 + \frac{1-r^2}{M}$
- $r$  assesses the training set quality such that for  $r \rightarrow 1$  the example matches perfectly the archetype whereas for  $r = 0$  solely noise remains.
- $\rho$  quantifies the entropy of the training set, namely  $\rho = \frac{1-r^2}{r^2 M}$
- $\rho_P$  is a generalization of  $\rho$  defined as  $\rho_P = \frac{1-r^{2P}}{r^{2P} M}$
- $t \in (0, 1)$  is the parameter for Guerra's interpolation: when  $t = 1$  we recover the original model, whereas for  $t = 0$  we compute the one-body terms
- $\mathcal{Z}$  is the partition function
- $\langle O(\boldsymbol{\sigma}) \rangle$  is the generalize average defined as  $\langle O(\boldsymbol{\sigma}) \rangle = \mathbb{E} \omega_t(O(\boldsymbol{\sigma}))$

## Acknowledgments

E.A. acknowledges financial support from Sapienza University of Rome (RM120172B8066CB0) and from PNRR MUR project n. PE0000013-FAIR.

A.B. acknowledges financial support from Ministero degli Affari Esteri e della Cooperazione Internazionale Italy-Israel (F85F21006230001) and PRIN grant *Statistical Mechanics of Learning Machines: from algorithmic and information-theoretical limits to new biologically inspired paradigms* n. 20229T9EAT.

L.A. acknowledges financial support from INdAM –GNFM Project (CUP E53C22001930001) and from PRIN grant *Stochastic Methods for Complex Systems* n. 2017JFFHS

D.L. acknowledges INdAM and C.N.R. (National Research Council), and A.A. acknowledges UniSalento, both for financial support via PhD-AI.

E.A., L.A., F.A., A.A., A.B acknowledge the stimulating research environment provided by the Alan Turing Institute’s Theory and Methods Challenge Fortnights event “Physics-informed Machine Learning”.

## References

- [1] E. Agliari, L. Albanese, F. Alemanno, A. Alessandrelli, A. Barra, F. Giannotti, D. Lotito, and D. Pedreschi. Dense Hebbian neural networks: a replica symmetric picture of supervised learning. *preprint arXiv*, 2022.
- [2] E. Agliari, L. Albanese, F. Alemanno, and A. Fachechi. A transport equation approach for deep neural networks with quenched random weights. *Journal of Physics A: Mathematical and Theoretical*, 54, 2021.
- [3] E. Agliari, L. Albanese, A. Barra, and G. Ottaviani. Replica symmetry breaking in neural networks: A few steps toward rigorous results. *Journal of Physics A: Mathematical and Theoretical*, 53, 2020.
- [4] E. Agliari, F. Alemanno, A. Barra, M. Centonze, and A. Fachechi. Neural networks with a redundant representation: Detecting the undetectable. *Physical Review Letters*, 124:28301, 2020.
- [5] E. Agliari, F. Alemanno, A. Barra, and A. Fachechi. Generalized guerra’s interpolation schemes for dense associative neural networks. *Neural Networks*, 128:254–267, 2020.
- [6] E. Agliari, F. Alemanno, A. Barra, and G. D. Marzo. The emergence of a concept in shallow neural networks. *Neural Networks*, 148:232–253, 2022.
- [7] E. Agliari, A. Barra, C. Longo, and D. Tantari. Neural networks retrieving boolean patterns in a sea of Gaussian ones. *Journal of Statistical Physics*, 168:1085–1104, 2017.
- [8] E. Agliari, A. Barra, P. Sollich, and L. Zdeborova. Machine learning and statistical physics: theory, inspiration, application. *J. Phys. A: Math. and Theor.*, Special, 2020.
- [9] E. Agliari, A. Fachechi, and C. Marullo. Nonlinear PDEs approach to statistical mechanics of dense associative memories. *Journal of Mathematical Physics*, 63(10):103304, 2022.
- [10] E. Agliari and G. D. Marzo. Tolerance versus synaptic noise in dense associative memories. *European Physical Journal Plus*, 135, 2020.
- [11] L. Albanese, F. Alemanno, A. Alessandrelli, and A. Barra. Replica symmetry breaking in dense hebbian neural networks. *J. Stat. Phys.*, 189(2):1–41, 2022.
- [12] L. Albanese and A. Alessandrelli. On gaussian spin glass with p-wise interactions. *Journal of Mathematical Physics*, 63:43302, 2022.

- [13] D. Alberici, F. Camilli, P. Contucci, and E. Mingione. The solution of the deep boltzmann machine on the nishimori line. *Comm. Math. Phys.*, 387(2):1191–1214, 2021.
- [14] D. Alberici, P. Contucci, and E. Mingione. Deep boltzmann machines: rigorous results at arbitrary depth. *Ann. H. Poinc.*, 22(8):2619–2642, 2021.
- [15] F. Alemanno, M. Aquaro, I. Kanter, A. Barra, and E. Agliari. Supervised hebbian learning. *Europhysics Letters*, *in press*, 2022.
- [16] D. J. Amit. *Modeling brain function: The world of attractor neural networks*. Cambridge university press, 1989.
- [17] D. J. Amit, H. Gutfreund, and H. Sompolinsky. Storing infinite numbers of patterns in a spin-glass model of neural networks. *Physical Review Letters*, 55:1530–1533, 1985.
- [18] A. Auffinger, G. B. Arous, and J. Cerny. Random matrices and complexity of spin glasses. *Communications on Pure and Applied Mathematics*, 66(2):165–201, 2013.
- [19] A. Auffinger and W. Chen. Free energy and complexity of spherical bipartite models. *J. Stat. Phys.*, 157(1):40–49, 2014.
- [20] A. Auffinger and Y. Zhou. The spherical p+s spin glass at zero temperature. *arXiv preprint*, page 2209.03866, 2022.
- [21] P. Baldi and S. S. Venkatesh. Number of stable points for spin-glasses and neural networks of higher orders. *Physical Review Letters*, 58, 1987.
- [22] J. Barbier and N. Macris. The adaptive interpolation method for proving replica formulas. Applications to the Curie–Weiss and wigner spike models. *J. Phys. A: Math. & Theor.*, 52:294002, 2019.
- [23] A. Barra. The mean field Ising model trough interpolating techniques. *J. Stat. Phys.*, 132:787–809, 2008.
- [24] A. Battista and R. Monasson. Capacity-resolution trade-off in the optimal learning of multiple low-dimensional manifolds by attractor neural networks. *Physical Review Letters*, 124, 2020.
- [25] A. Bovier and B. Niederhauser. The spin-glass phase-transition in the Hopfield model with p-spin interactions. *Advances in Theoretical and Mathematical Physics*, 5:1001–1046, 8 2001.
- [26] G. Carleo and et al. Machine learning and the physical sciences. *Rev. Mod. Phys.*, 91:045002., 2019.
- [27] P. Carmona and Y. Hu. Universality in Sherrington-Kirkpatrick’s spin glass model. *Annales de l’institut Henri Poincare (B) Probability and Statistics*, 42, 2006.
- [28] A. C. C. Coolen, R. Kühn, and P. Sollich. *Theory of neural information processing systems*. OUP Oxford, 2005.
- [29] A. Crisanti, D. J. Amit, and H. Gutfreund. Saturation level of the Hopfield model for neural network. *Europhysics Letters*, 2:337–341, 8 1986.
- [30] A. Decelle, S. Hwang, J. Rocchi, and D. Tantari. Annealing and replica-symmetry in deep boltzmann machines. *Scientific Reports*, 11(1):1–13, 2021.
- [31] A. Decelle and F. Ricci-Tersenghi. Solving the inverse Ising problem by mean-field methods in a clustered phase space with many states. *Physical Review E*, 94(1):012112, 2016.
- [32] E. Gardner. Multiconnected neural network models. *Journal of Physics A: General Physics*, 20, 1987.
- [33] G. Genovese. Universality in bipartite mean field spin glasses. *Journal of Mathematical Physics*, 53, 2012.
- [34] F. Guerra. Broken replica symmetry bounds in the mean field spin glass model. *Communications in Mathematical Physics*, 233:1–12, 2003.

- [35] J. J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences of the United States of America*, 79:2554–2558, 1982.
- [36] D. Krotov and J. Hopfield. Dense associative memory is robust to adversarial inputs. *Neural Computation*, 30:3151–3167, 2018.
- [37] M. Mézard, G. Parisi, and M. A. Virasoro. *Spin glass theory and beyond: an introduction to the Replica Method and its applications*, volume 9. World Scientific Publishing Company, 1987.
- [38] T. Plefka. Convergence condition of the TAP equation for the infinite-ranged Ising spin glass model. *Journal of Physics A: Mathematical and General*, 15, 1982.
- [39] T. Plefka. Expansion of the gibbs potential for quantum many-body systems: General formalism with applications to the spin glass and the weakly nonideal bose gas. *Physical Review E - Statistical, Nonlinear and Soft Matter Physics*, 73, 2006.
- [40] T. J. Sejnowski. Higher-order boltzmann machines. In *AIP Conference Proceedings*, volume 151, pages 398–403. American Institute of Physics, 1986.
- [41] H. Steffan and R. Kühn. Replica symmetry breaking in attractor neural network models. *Zeitschrift für Physik B Condensed Matter*, 95, 1994.
- [42] E. Strubell, A. Ganesh, and A. McCallum. Energy and policy considerations for deep learning in NLP. *ACL 2019 - 57th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference*, 2020.
- [43] E. Subag. The complexity of spherical  $p$ -spin models – A second moment approach. *The Annals of Probability*, 45(5):3385–3450, 2017.
- [44] E. Subag and O. Zeitouni. The extremal process of critical points of the pure  $p$ -spin spherical spin glass model. *Probability theory and related fields*, 168(3):773–820, 2017.
- [45] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. *arXiv preprint arXiv:1708.07747*, pages 1–6, 2017.
- [46] L. Zdeborova. Understanding deep learning is also a job for physicists. *Nature Physics*, 16:602–604, 2020.