# Thermodynamics of Restricted Boltzmann Machines and Related Learning Dynamics

A. Decelle G. Fissore C. Furtlehner

#### **Abstract**

We investigate the thermodynamic properties of a Restricted Boltzmann Machine (RBM), a simple energy-based generative model used in the context of unsupervised learning. Assuming the information content of this model to be mainly reflected by the spectral properties of its weight matrix W, we try to make a realistic analysis by averaging over an appropriate statistical ensemble of RBMs.

First, a phase diagram is derived. Otherwise similar to that of the Sherrington-Kirkpatrick (SK) model with ferromagnetic couplings, the RBM's phase diagram presents a ferromagnetic phase which may or may not be of compositional type depending on the kurtosis of the distribution of the components of the singular vectors of W.

Subsequently, the learning dynamics of the RBM is studied in the thermodynamic limit. A "typical" learning trajectory is shown to solve an effective dynamical equation, based on the aforementioned ensemble average and explicitly involving order parameters obtained from the thermodynamic analysis. In particular, this let us show how the evolution of the dominant singular values of W, and thus of the unstable modes, is driven by the input data. At the beginning of the training, in which the RBM is found to operate in the linear regime, the unstable modes reflect the dominant covariance modes of the data. In the non-linear regime, instead, the selected modes interact and eventually impose a matching of the order parameters to their empirical counterparts estimated from the data.

Finally, we illustrate our considerations by performing experiments on both artificial and real data, showing in particular how the RBM operates in the ferromagnetic compositional phase.

#### 1 Introduction

The Restricted Boltzmann Machine (RBM) [1] is an important machine learning tool used in many applications, by virtue of its ability to model complex probability distributions. It is a neural network which serves as a generative model, in the sense that it is able to approximate the probability distribution corresponding to the empirical distribution of any set of high-dimensional data points living in a discrete or real space of dimension  $N \gg 1$ . From the theoretical point of view, the RBM is of high interest as it is one of the simplest neural network generative models and the probability distribution that it defines presents a simple analytic form. Moreover, there are clear

1 Introduction 2

connections between RBMs and well known disordered systems in statistical physics. As an example, when data are composed by vectors with binary components the discrete RBM takes the form of an heterogeneous Ising model composed of one layer of visible units (the observable variables) connected to one layer of hidden units (the latent or hidden variables building up the dependencies between the visible ones), in which couplings and fields are obtained from the training data through a learning procedure. In order to build more powerful models, RBMs can be stacked to form "deep" architectures. In such a case, they can form a multi-layer generative model known as a Deep Boltzmann Machine (DBM) [2] or they can be stacked and trained layerwise as a pre-training procedure for neural networks [3]. The standard learning algorithms in use are the contrastive divergence [4] (CD) and the refined Persistence CD [5] (PCD), which are based on a quick Monte Carlo estimation of the response function of the RBM and are efficient and well documented [6]. Nevertheless, despite some interesting interpretations of CD in terms of non-equilibrium statistical physics [7], the learning of RBMs remains a set of obscure recipes from the statistical physics point of view: hyperparameters (like the size of the hidden layer) are supposed to be set empirically without any theoretical guidelines.

Historically, statistical physics played a central role in studying the theoretical foundations of neural networks. In particular, during the 1980s many works on the Hopfield model [8, 9, 10, 11] managed to define its learning capacity and to compute the number of independent patterns that it could store. It is worth noticing that, as RBMs are ultimately defined as a Boltzmann distribution with pairwise interactions on a bipartite graph, they can be studied in a way similar to that used for the Hopfield model. The analogy is even stronger since connections between the Hopfield model and RBMs have been made explicit when using Gaussian hidden variables [12], here the number of patterns of the Hopfield model corresponding to the number of hidden units. Motivated by a renewed excitement for neural networks, recent works actually propose to exploit the statistical physics formulation of the RBM to understand what is its learning capacity and how mean-field methods can be exploited to improve the model. In [13, 14, 15], mean-field based learning methods using TAP equations are developed. TAP solutions are usually expected to define a decomposition of the measure in terms of pure thermodynamical states and are useful both as an algorithm to compute the marginals of the variables of the model and to identify the pure states when they are yet unknown. For instance, in a sparse explicit Boltzmann machine (i.e. without latent variables) this implicit clustering can be done by means of belief propagation fixed points <sup>1</sup> with simple empirical learning rules [16]. In [17, 18], an analysis of the static properties of RBMs is done assuming a given weight matrix W, in order to understand collective phenomena in the latent representation, i.e. the way latent variables organize themselves in a compositional phase [19, 20] to represent actual data. These analysis make use of the replica trick (or equivalent) making the common assumption that the components of the weight matrix W are i.i.d.; despite the fact that this approach may give some insights into the retrieval phase, this approximation is problematic since, as far as a realistic RBM is concerned (an RBM learned on data), the learning mechanism introduces correlations within the weights of W and then it seems rather crude to continue to assume the independence and hope to understand the realistic statistical properties of the model.

Concerning the learning procedure of neural networks, many recent statistical physics based analyses have been proposed, most of them within teacher-student set-

<sup>&</sup>lt;sup>1</sup> a somewhat different form of the TAP equations

1 Introduction 3

ting [21]. This imposes a rather strong assumption on the data in the sense that it is assumed that these are generated from a model belonging to the parametric family of interest, hiding as a consequence the role played by the data themselves in the procedure. From the analysis of related linear models [22, 23], it is already a well established fact that a selection of the most important modes of the singular values decomposition (SVD) of the data is performed in the linear case. In fact in the simpler context of linear feed-forward models the learning dynamics can be fully characterized by means of the SVD of the data matrix [24], showing in particular the emergence of each mode by order of importance with respect to the corresponding singular values.

First steps to follow this guideline have been done in [25], in the context of a general RBM and to address the shortcomings of previous analyses, in particular concerning the assumptions over the weights distribution. To this end it has been proposed to characterize both the learned RBM and the learning process itself by means of the SVD spectrum of the weight matrix in order to single out the information content of the RBM. It is assumed that the SVD spectrum is split in a continuous bulk of singular vectors corresponding to noise and a set of outliers that represent the information content. By doing this it is possible to go beyond the usual unrealistic assumption of i.i.d. weights made for analyzing RBMs. Proceeding along this direction, in the present work we first present a thermodynamic analysis of RBMs under the more realistic assumptions over the weight matrix that we propose. Then, on the same basis, the learning dynamics of RBMs is studied by direct analysis of the dynamics of the SVD modes, both in the linear and non-linear regimes.

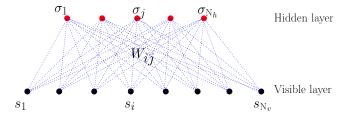


Fig. 1: bipartite structure of the RBM.

The paper is organized as follows: in Section 2 we introduce the RBM model and its associated learning procedures. Section 3 presents the static thermodynamical properties of the RBM with realistic hypothesis on its weights: a statistical ensemble of weight matrices is discussed in Section 3.1; mean-field equations in the replica-symmetric (RS) framework are given in Section 3.2 and the corresponding phase diagram is studied in Section 3.3 with a proper delimitation of the RS domain where the learning procedure is supposed to take place. The ferromagnetic phase is studied in great details in 3.4 by looking in particular at the conditions leading to a compositional phase. Section 4 is devoted to the learning dynamics. In Section 4.1, a deterministic learning equation is derived in the thermodynamic limit and a set of dynamical parameters is shown to emerge naturally from the SVD of the weight matrix. This equation is analyzed for linear RBMs in Section 4.2 in order to identify the unstable deformation modes of W that result in the first emerging patterns at the beginning of the learning process; the non-linear regime is described in Section 4.3, on the basis of the thermodynamic analysis, by numerically solving the effective learning equations in simple cases. Our analysis is finally illustrated and validated in Section 5 by actual tests on the MNIST dataset.

# 2 The RBM and its associated learning procedure

An RBM is a Markov random field with pairwise interactions defined on a bipartite graph formed by two layers of non-interacting variables: the visible nodes and the hidden nodes representing respectively data configurations and latent representations (see Figure 1). The former noted  $\mathbf{s} = \{s_i, i = 1 \dots N_v\}$  correspond to explicit representations of the data while the latter noted  $\mathbf{\sigma} = \{\sigma_j, j = 1 \dots N_h\}$  are there to build arbitrary dependencies among the visible units. They play the role of an interacting field among visible nodes. Usually the nodes are binary-valued (of Boolean type or Bernoulli distributed) but Gaussian distributions or more broadly arbitrary distributions on real-valued bounded support are also used [26], ultimately making RBMs adapted to more heterogeneous data sets. Here to simplify we assume that visible and hidden nodes will be taken as binary variables  $s_i, \sigma_j \in \{-1,1\}$  (using  $\pm 1$  values gives the advantage of working with symmetric equations hence avoiding to deal with the "hidden" biases on the variables that appear when considering binary  $\{0,1\}$  variables). Like in the Hopfield model [8], which can actually be cast into an RBM [12], an energy function is defined for a configuration of nodes

$$E(\boldsymbol{s}, \boldsymbol{\sigma}) = -\sum_{i,j} s_i W_{ij} \sigma_j + \sum_{i=1}^{N_v} \eta_i s_i + \sum_{j=1}^{N_h} \theta_j \sigma_j$$
 (1)

and this is exploited to define a joint distribution between visible and hidden units, namely the Boltzmann distribution

$$p(s,\sigma) = \frac{e^{-E(s,\sigma)}}{Z}$$
 (2)

where W is the weight matrix and  $\eta$  and  $\theta$  are biases, or external fields on the variables.  $Z = \sum_{s,\sigma} e^{-E(s,\sigma)}$  is the partition function of the system. The joint distribution between visible variables is then obtained by summing over hidden ones. In this context, learning the parameters of the RBM means that, given a dataset of M samples composed of  $N_v$  variables, we ought to infer values to W,  $\eta$  and  $\theta$  such that new generated data obtained by sampling this distribution should be similar to the input data. The general method to infer the parameters is to maximize the log likelihood of the model, where the pdf (2) has first been summed over the hidden variables

$$\mathcal{L} = \sum_{j} \langle \log(2\cosh(\sum_{i} W_{ij} s_i - \theta_j)) \rangle_{\text{Data}} - \sum_{i} \eta_i \langle s_i \rangle_{\text{Data}} - \log(Z).$$
 (3)

Different learning methods have been set up and proven to work efficiently, in particular the contrastive divergence (CD) algorithm from Hinton [4] and more recently TAP based learning [13]. They all correspond to expressing the gradient ascent on the likelihood as

$$\Delta W_{ij} = \gamma \left( \langle s_i \sigma_j p(\sigma_j | \mathbf{s}) \rangle_{\text{Data}} - \langle s_i \sigma_j \rangle_{p_{\text{BBM}}} \right) \tag{4}$$

$$\Delta \eta_i = \gamma \left( \langle s_i \rangle_{p_{\text{RBM}}} - \langle s_i \rangle_{\text{Data}} \right) \tag{5}$$

$$\Delta\theta_j = \gamma \left( \langle \sigma_j \rangle_{p_{\text{RBM}}} - \langle \sigma_j p(\sigma_j | \boldsymbol{s}) \rangle_{\text{Data}} \right)$$
 (6)

where  $\gamma$  is the learning rate. The main problem are the  $\langle \cdots \rangle_{P\text{RBM}}$  terms on the right hand side of (4-6). These are not tractable and the various methods basically differ in their way of estimating those terms (Monte-Carlo Markov chains, naive mean-field, TAP...). For an efficient learning the  $\langle \cdots \rangle_{\text{Data}}$  terms must also be approximated by making use of random mini-batches of data at each step.

# 3 Static thermodynamical properties of an RBM

#### 3.1 Statistical ensemble of RBMs

When analyzing the thermodynamical properties of RBMs, it is common to assume that the weights  $W_{ij}$  are i.i.d. random variables, like for example in [20, 17, 18]. This generally leads to a Marchenko-Pastur (MP) distribution [27] of the singular values of W, which is unrealistic.

In order to clarify our notation, let us recall the definition of the singular value decomposition (SVD). As a generalization of eigenmodes decomposition to rectangular matrices, the SVD for a RBM is given by

$$\mathbf{W} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \tag{7}$$

where  $\mathbf{U}$  is an orthogonal  $N_v \times N_h$  matrix whose columns are the left singular vectors  $\mathbf{u}^{\alpha}$ ,  $\mathbf{V}$  is an orthogonal  $N_h \times N_h$  matrix whose columns are the right singular vectors  $\mathbf{v}^{\alpha}$  and  $\mathbf{\Sigma}$  is a diagonal matrix whose elements are the singular values  $w_{\alpha}$ . The separation into left and right singular vectors is due to the rectangular nature of the decomposed matrix, and the similarity with eigenmodes decomposition is revealed by the following SVD equations

$$\mathbf{W}\mathbf{v}^{\alpha} = w_{\alpha}\mathbf{u}^{\alpha}$$
$$\mathbf{W}^{T}\mathbf{u}^{\alpha} = w_{\alpha}\mathbf{v}^{\alpha}$$

In [25] it is argued that the MP distribution of SVD modes actually corresponds to the noise of the weight matrix, while the information content of the RBM is better expressed by the presence of SVD modes outside of this bulk. This leads us to write the weight matrix as

$$W_{ij} = \sum_{\alpha=1}^{K} w_{\alpha} u_i^{\alpha} v_j^{\alpha} + r_{ij} \tag{8}$$

where the  $w_{\alpha} = O(1)$  are isolated singular values (describing a rank K matrix), the  $\boldsymbol{u}^{\alpha}$  and  $\boldsymbol{v}^{\alpha}$  are the dominant eigenvectors of the SVD decomposition and the  $r_{ij} = \mathcal{N}(0,\sigma^2/L)$  are i.i.d. terms corresponding to noise, with  $L = \sqrt{N_h N_v}$ . The  $\{u^{\alpha}\}$  and  $\{v^{\alpha}\}$  are two sets of respectively  $N_v$  and  $N_h$ -dimensional orthonormal vectors, which means that their components are respectively  $O(1/\sqrt{N_v})$  and  $O(1/\sqrt{N_h})$ , and  $K \leq N_v, N_h$ . We assume  $N_h < N_v$  to be the rank of W and  $w_{\alpha} > 0$  and O(1) for all  $\alpha$ . Note that in the limit  $N_v \to \infty$  and  $N_h \to \infty$  with  $\kappa \stackrel{\text{def}}{=} N_h/N_v$  fixed and  $K/L \to 0$ ,  $WW^T$  has a spectrum density  $\rho(\lambda)$  composed of a Marchenko-Pastur bulk of eigenvalues and of set of discrete modes:

$$\rho(\lambda) = \frac{L}{2\pi\sigma^2} \frac{\sqrt{(\lambda^+ - \lambda)(\lambda - \lambda^-)}}{\kappa\lambda} \mathbb{1}_{\{\lambda \in [\lambda^-, \lambda^+]\}} + \sum_{\alpha = 1}^K \delta(\lambda - w_\alpha^2),$$

with

$$\lambda^{\pm} \stackrel{\text{def}}{=} \sigma^2 \left( \kappa^{\frac{1}{4}} \pm \kappa^{-\frac{1}{4}} \right)^2.$$

The interpretation for the noise term  $r_{ij}$  is given by the presence of an extensive number of modes at the bottom of the spectrum, along which the variables won't be able to condense but that still contribute to the fluctuations. In the present form our model of RBM is similar to the Hopfield model and recent generalizations [28], the patterns being represented by the SVD modes outside of the bulk. The main difference, in addition to the bipartite structure of the graph, is the non-degeneracy of the singular values  $w_{\alpha}$ . The choice made here is to consider K finite, giving  $W_{ij} = O(1/N)$  which means that the thresholds  $\theta_j$  (having the meaning of feature detectors) should be O(1) because feature j is detected when an extensive number of spins  $S_i$  is aligned with  $W_{ij}$ . In addition, this allows us to assume simple distributions for the components of  $u^{\alpha}$  and  $v^{\alpha}$  (for instance, considering them i.i.d.). Altogether, this defines the statistical ensemble of RBM to which we restrict our analysis of the learning procedure.

Another approach would be to consider  $K = N_h$  extensive, thereby assuming that all modes can potentially condense even though they are associated to dominated singular values. In that case, the separation between the condensed modes and the rest should be made when order parameters are introduced and the noise would then correspond to uncondensed modes. If the number of condensed modes is assumed to be extensive, then we should instead consider an average over the orthogonal group which would lead to a slightly different mean-field theory [29, 30].

# 3.2 Replica symmetric Mean-field equation

Our analysis in the thermodynamic limit follows classical treatments using replicas, like [31, 9] for the Hopfield model or [17] for bipartite models. The starting point is to express the average over u, v and  $r_{ij}$  of the log partition function Z in (2) with the help of the replica trick:

$$\mathsf{E}_{u,v,r}[\log(Z)] = \lim_{p \to 0} \frac{d}{dp} \mathsf{E}_{u,v,r}[Z^p].$$

First the average over  $r_{ij}$  yields

$$\exp \left[\frac{\sigma^2}{2L} \Bigl(\sum_a s_i^a \sigma_j^a\Bigr)^2\right] = \exp \Bigl[\frac{\sigma^2}{2L} \Bigl(p + \sum_{a \neq b} s_i^a s_i^b \sigma_j^a \sigma_j^b\Bigr)\Bigr].$$

After this averaging, 4 sets of order parameters  $\{(m_{\alpha}^a, \bar{m}_{\alpha}^a), a=1, \ldots p, \alpha=1, \ldots K\}$  and  $\{(Q_{ab}, \bar{Q}_{ab}), a, b=1, \ldots p, a\neq b\}$  are introduced with the help of two distinct Hubbard-Stratonovich transformations. The first one corresponds to

$$\exp\left[\frac{\sigma^2}{2L}\left(\sum_{i,j,a\neq b} s_i^a s_i^b \sigma_j^a \sigma_j^b\right)\right] = \int \prod_{a\neq b} \frac{dQ_{ab} d\bar{Q}_{ab}}{2\pi}$$

$$\times \exp\left[-\frac{L\sigma^2}{2} \sum_{a\neq b} \left(Q_{ab} \bar{Q}_{ab} - \frac{Q_{ab}}{N_v} \sum_i s_i^a s_i^b - \frac{\bar{Q}_{ab}}{N_h} \sum_i \sigma_j^a \sigma_j^b\right)\right].$$

The second one is aimed at extracting magnetization's contributions correlated with the modes:

$$\exp\left(L\sum_{\alpha}w_{\alpha}s_{\alpha}^{a}\sigma_{\alpha}^{a}\right) \propto \int \prod_{\alpha} \frac{dm_{\alpha}^{a}d\bar{m}_{\alpha}^{a}}{2\pi} \times \exp\left(-L\sum_{\alpha}w_{\alpha}\left(m_{\alpha}^{a}\bar{m}_{\alpha}^{a} - m_{\alpha}^{a}s_{\alpha}^{a} - \bar{m}_{\alpha}^{a}\sigma_{\alpha}^{a}\right)\right),$$

with

$$s_{\alpha}^{a} \stackrel{\text{def}}{=} \frac{1}{\sqrt{L}} \sum_{i} s_{i} u_{i}^{\alpha} \quad \text{and} \quad \sigma_{\alpha}^{a} \stackrel{\text{def}}{=} \frac{1}{\sqrt{L}} \sum_{j} \sigma_{j}^{a} v_{j}^{\alpha},$$
 (9)

These variables represent the following quantities:

$$m_{\alpha}^{a} \sim E_{u,v,r}(\langle \sigma_{\alpha}^{a} \rangle)$$
  $\bar{m}_{\alpha}^{a} \sim E_{u,v,r}(\langle s_{\alpha}^{a} \rangle)$   $Q_{ab} \sim E_{u,v,r}(\langle \sigma_{i}^{a} \sigma_{i}^{b} \rangle)$   $\bar{Q}_{ab} \sim E_{u,v,r}(\langle s_{j}^{a} s_{j}^{b} \rangle),$ 

namely the correlations of the hidden [resp. visible] states with the left [resp. right] singular vectors and the Edward-Anderson (EA) order parameters measuring the correlation between replicas of hidden or visible states.  $\mathsf{E}_u$  and  $\mathsf{E}_v$  denote an average w.r.t. the rescaled components  $u \simeq \sqrt{N_v} u_i^\alpha$  and  $v \simeq \sqrt{N_h} v_j^\alpha$  of the SVD modes. The transformations involve pairs of complex integration variables because of the asymmetry introduced by the two-layers structure in contrast to fully connected models.

We obtain the following representation:

$$\begin{split} \mathsf{E}_{u,v,r}[Z^p] &= \int \prod_{a,\alpha} \frac{dm_{\alpha}^a d\bar{m}_{\alpha}^a}{2\pi} \prod_{a \neq b} \frac{dQ_{ab} d\bar{Q}_{ab}}{2\pi} \\ &\times \exp \left\{ -L \left( \sum_{a,\alpha} w_{\alpha} m_{\alpha} \bar{m}_{\alpha} + \frac{\sigma^2}{2} \sum_{a \neq b} Q_{ab} \bar{Q}_{ab} - \frac{1}{\sqrt{\kappa}} A[m,Q] - \sqrt{\kappa} B[\bar{m},\bar{Q}] \right) \right\} \end{split}$$

with  $\kappa = N_h/N_v$  and

$$A[m,Q] \stackrel{\text{def}}{=} \log \left[ \sum_{S^a \in \{-1,1\}} \mathsf{E}_u \left( e^{\frac{\sqrt{\kappa}\sigma^2}{2} \sum_{a \neq b} Q_{ab} S^a S^b + \kappa^{\frac{1}{4}} \sum_{a,\alpha} (w_\alpha m_\alpha^a - \eta_\alpha) u^\alpha S^a} \right) \right], \quad (10)$$

$$B[\bar{m}, \bar{Q}] \stackrel{\text{def}}{=} \log \left[ \sum_{S^a \in \{-1, 1\}} \mathsf{E}_v \left( e^{\frac{\sqrt{\kappa}\sigma^2}{2} \sum_{a \neq b} \bar{Q}_{ab} \sigma^a \sigma^b + \kappa^{-\frac{1}{4}} \sum_{a, \alpha} (w_\alpha \bar{m}_\alpha^a - \theta_\alpha) v^\alpha \sigma^a} \right) \right], \tag{11}$$

(12)

with

$$\theta_{\alpha} \stackrel{\text{def}}{=} \frac{1}{\sqrt{L}} \sum_{j} \theta_{j} v_{j}^{\alpha} = O(1).$$

Since  $\{v^{\alpha}\}$  is an incomplete basis we also need to take care of the potential residual

transverse parts  $\eta^{\perp}$  and  $\theta^{\perp}$ , such that the following decompositions hold:

$$\eta_i = \eta_i^{\perp} + \sqrt{L} \sum_{\alpha} \eta_{\alpha} u_i^{\alpha}, \tag{13}$$

$$\theta_j = \theta_j^{\perp} + \sqrt{L} \sum_{\alpha} \theta_{\alpha} v_j^{\alpha}. \tag{14}$$

To keep things tractable, both  $\eta^{\perp}$  and  $\theta^{\perp}$  will be considered negligible in the sequel. Taking into account these components would lead to the addition of a random field to the effective RS field of the variables and eventually to a richer set of saddle point solutions. Note that the order of magnitude of  $\eta_{\alpha}$  and  $\theta_{\alpha}$  is at this stage an assumption. If  $\eta_i$  and  $u_i^{\alpha}$  (or  $\theta_j$  and  $v_j^{\alpha}$ ) were uncorrelated they would scale as  $1/\sqrt{L}$ . Moreover, regarding the ensemble average, we will consider  $\eta_{\alpha}$  and  $\theta_{\alpha}$  fixed in the sequel.

The thermodynamic properties are obtained by first making a saddle point approximation possible by letting first  $L \to \infty$  and taking the limit  $p \to 0$  afterwards. We restrict here the discussion to RS saddle points [32]. The breakdown of RS can actually be determined by computing the so-called AT line [33] (see Appendix A). At this point we assume a non-broken replica symmetry. The set  $\{Q_{ab}, \bar{Q}_{ab}\}$  reduces then to a pair  $(q, \bar{q})$  of spin glass parameters, i.e.  $Q_{ab} = q$  and  $\bar{Q}_{ab} = \bar{q}$  for all  $a \neq b$ , while quenched magnetizations on the SVD directions are now represented by  $\{(m_{\alpha}, \bar{m}_{\alpha}), \alpha = 1, \ldots K\}$ .

Taking the limit  $p \to 0$  yields the following limit for the free energy:

$$f[m, \bar{m}, q, \bar{q}] = \sum_{\alpha} w_{\alpha} m_{\alpha} \bar{m}_{\alpha} - \frac{\sigma^{2}}{2} q \bar{q} + \frac{\sigma^{2}}{2} (q + \bar{q})$$
$$- \frac{1}{\sqrt{\kappa}} \mathsf{E}_{u,x} \Big[ \log 2 \cosh \big( h(x, u) \big) \Big] - \sqrt{\kappa} \mathsf{E}_{v,x} \Big[ \log 2 \cosh \big( \bar{h}(x, v) \big) \Big]. \tag{15}$$

Assuming a replica-symmetric phase, the saddle-point equations are given by

$$m_{\alpha} = \kappa^{\frac{1}{4}} \mathsf{E}_{v,x} \Big[ v^{\alpha} \tanh \big( \bar{h}(x,v) \big) \Big], \qquad q = \mathsf{E}_{v,x} \Big[ \tanh^{2} \big( \bar{h}(x,v) \big) \Big]$$
 (16)

$$\bar{m}_{\alpha} = \kappa^{-\frac{1}{4}} \mathsf{E}_{u,x} \Big[ u^{\alpha} \tanh \big( h(x,u) \big) \Big], \qquad \bar{q} = \mathsf{E}_{u,x} \Big[ \tanh^2 \big( h(x,u) \big) \Big]$$
 (17)

where

$$h(x,u) \stackrel{\text{def}}{=} \kappa^{\frac{1}{4}} \left( \sigma \sqrt{q} x + \sum_{\gamma} (w_{\gamma} m_{\gamma} - \eta_{\gamma}) u^{\gamma} \right)$$

$$\bar{h}(x,v) \stackrel{\text{def}}{=} \kappa^{-\frac{1}{4}} \left( \sigma \sqrt{\bar{q}} x + \sum_{\gamma} (w_{\gamma} \bar{m}_{\gamma} - \theta_{\gamma}) v^{\gamma} \right),$$

and  $\kappa = N_h/N_v$ , with  $\mathsf{E}_{u,x}$  and  $\mathsf{E}_{v,x}$  denoting an average over the Gaussian variable  $x = \mathcal{N}(0,1)$  and the rescaled components  $u \sim \sqrt{N_v} u_i^\alpha$  and  $v \sim \sqrt{N_h} v_j^\alpha$  of the SVD modes. We note that the equations are symmetric under the exchange  $\kappa \to \kappa^{-1}$ , simultaneously with  $m \leftrightarrow \bar{m}$ ,  $q \leftrightarrow \bar{q}$  and  $\eta \leftrightarrow \theta$ , given that u and v have the same distribution. In addition, for independently distributed  $u_i^\alpha$  and  $v_j^\alpha$  and vanishing fields  $(\eta = \theta = 0)$ , solutions corresponding to non-degenerate magnetizations have symmetric counterparts: each pair of non-vanishing magnetizations can be negated independently as  $(m_\alpha, \bar{m}_\alpha) \to (-m_\alpha, -\bar{m}_\alpha)$ , generating new solutions. So to one solution presenting n condensed modes, there correspond  $2^n$  distinct solutions.

## 3.3 Phase Diagram

The fixed point equations (16, 17) can be solved numerically to tell us how the variables condensate on the SVD modes within each equilibrium state of the distribution and whether a spin-glass or a ferromagnetic phase is present. The important point here is that with K finite and a non-degenerate spectrum the mode with highest singular value dominates the ferromagnetic phase.

In absence of bias  $(\eta = \theta = 0)$  and once  $1/\sigma$  is interpreted as temperature and  $w_{\alpha}/\sigma$  as ferromagnetic couplings, we get a phase diagram similar to that of the Sherrington-Kirkpatrick (SK) model with three distinct phases (see Figure 2)

- a paramagnetic phase  $(q = \bar{q} = m_{\alpha} = \bar{m}_{\alpha} = 0)$  (P),
- a ferromagnetic phase  $(q, \bar{q}, m_{\alpha}, \bar{m}_{\alpha} \neq 0)$  (F),
- a spin glass phase  $(q, \bar{q} \neq 0; m_{\alpha} = \bar{m}_{\alpha} = 0)$  (SG).

In general, the lines separating the different phases correspond to second order phase transitions and can be obtained by a stability analysis of the Hessian of the free energy. They are related to unstable modes of the linearized mean-field equations and correspond to an eigenvalue of the Hessian becoming negative.

The (SG-P) line is obtained by looking at the Hessian in the  $(q, \bar{q})$  sector:

$$H_{q\bar{q}} \underset{\substack{m=0\\ a=0}}{=} -\frac{1}{2} \begin{bmatrix} \sigma^2 & \frac{\sigma^4}{\sqrt{\kappa}} \\ \sqrt{\kappa}\sigma^4 & \sigma^2 \end{bmatrix}$$

from what results that the spin glass phase develops when  $\sigma \geq 1^2$ . This transition line is understood tacking directly into account the spectral properties of the weight matrix. Classically, this is done with the help of the linearized TAP equations and exploiting the Marchenko-Pastur distribution [32]. In our context, the linearized TAP equations read

$$\begin{bmatrix} \mu \\ \nu \end{bmatrix} = \begin{bmatrix} -\sqrt{\kappa}\sigma^2 & W^T \\ W & -\frac{\sigma^2}{\sqrt{\kappa}} \end{bmatrix} \begin{bmatrix} \mu \\ \nu \end{bmatrix}$$

given the variance  $\sigma^2/L$  of the weights in absence of dominant modes. Then we can show that the paramagnetic phase becomes unstable when the highest eigenvalue of the matrix on the rhs is equal to 1: if  $\lambda$  is a singular value of W, the corresponding eigenvalues  $\Lambda^{\pm}$  verify the relation

$$\left(\frac{\Lambda^{\pm}}{\sqrt{\kappa}} \pm \sigma^2\right) \left(\sqrt{\kappa} \Lambda^{\pm} \pm \sigma^2\right) = \lambda^2.$$

from which it is clear that the largest eigenvalue  $\Lambda_{max}$  corresponds to the largest singular value  $\lambda_{max}$ . Owing to the Marchenko-Pastur distribution  $\lambda_{max} = \sigma^2(\sqrt{\kappa} + 1)(1 + 1/\sqrt{\kappa})$  so  $\Lambda_{max}$  verifies

$$\left(\frac{\Lambda_{max}}{\sqrt{\kappa}} + \sigma^2\right) \left(\sqrt{\kappa} \Lambda_{max} + \sigma^2\right) = \sigma^2 (\sqrt{\kappa} + 1) \left(\frac{1}{\sqrt{\kappa}} + 1\right).$$

 $\Lambda_{max} = 1$  is readily obtained for  $\sigma^2 = 1$ .

<sup>&</sup>lt;sup>2</sup> Note that in [17] a dependence  $\sqrt{\kappa(1-\kappa)}$  ( $\sqrt{\alpha(1-\alpha)}$ in their notation) is found. This dependence is hidden in our definition of  $\sigma^2$  giving  $L=\sqrt{N_vN_h}$  times the variance of  $r_{ij}$  instead of  $N_v+N_h$  as in their case.

For the (F-SG) frontier we can look at the sector  $(m_{\alpha}, \bar{m}_{\alpha})$  corresponding to the emergence of a single mode  $\alpha$  (written in the spin-glass phase):

$$\begin{split} H_{\alpha\alpha} &= \begin{bmatrix} w_{\alpha} & w_{\alpha}^2 \mathsf{E}_{v,x} \Big[ (v^{\alpha})^2 \operatorname{sech}^2 \big( \bar{h}(x,v) \big) \big) \Big] \\ w_{\alpha}^2 \mathsf{E}_{u,x} \Big[ (u^{\alpha})^2 \operatorname{sech}^2 \big( h(x,u) \big) \Big] & w_{\alpha} \end{bmatrix} \\ &= \begin{bmatrix} w_{\alpha} & w_{\alpha}^2 (1-q) \\ w_{\alpha}^2 (1-\bar{q}) & w_{\alpha} \end{bmatrix} \end{split}$$

From this it is clear that the first mode to become unstable is the mode  $\alpha$  with highest singular value  $w_{\alpha}$  and this occurs when q and  $\bar{q}$ , solutions of (16,17), verify

$$(1-q)(1-\bar{q})w_{\alpha}^{2}=1.$$

As for the SK model, this line appears to be well below the de Almeida-Thouless (AT) line, which is the line above which the RS solution is stable (see Figure 2, and Appendix A for the computation of the AT line). This means that in principle a replica symmetry breaking treatment would be necessary to properly separate the two phases. However, we will leave aside this point as we are mainly interested in the practical aspects, namely the ability of the RBM to learn arbitrary data, and so we are mostly concerned with the ferromagnetic phase above the AT line.

For the (P-F) line we consider the same sector of the Hessian but now written in the paramagnetic phase, i.e. setting q=0 in the above equation, and this simply yields the emergence of the single mode  $\alpha$  for  $w_{\alpha}=1$ .

Note that all of this is independent on how the statistical average over u and v is performed. Instead, as we shall see later on, the way of averaging influences the nature of the ferromagnetic phase.

Regarding the stability of the RS solution, the computation of the AT line reported in Appendix A is similar to the classical one made for the SK model, though slightly more involved. In fact we were not able to fully characterize, in replica space, all the possible instabilities of the Hessian which would potentially lead to a breakdown of the replica symmetry. At least the one responsible for the ordinary SK model RS breakdown has a counterpart in the bipartite case that gives a necessary condition for the stability of the RS solution:

$$\frac{1}{\sigma^2} > \sqrt{\mathsf{E}_{x,u} \Big( \mathrm{sech}^4 \Big( h(x,u) \Big) \Big) \mathsf{E}_{x,v} \Big( \mathrm{sech}^4 \Big( \bar{h}(x,v) \Big) \Big)},$$

For  $\kappa=1$  the terms below the radical become identical and the condition reduces to the one of the SK model, except for the u averages which are not present in the SK model. In Figure 2, is shown the influence on the phase diagram of the value of  $\kappa$  and of the type of average made on u and v.

# 3.4 Nature of the Ferromagnetic phase

Some subtleties arise when considering various ways of averaging over the components of the singular vectors. In [19, 20] is emphasized the importance for networks to be able to reproduce compositional states structured by combination of hidden variables. In our representation, we don't have direct access to this property but, in some sense, to the dual one, which is given by states corresponding to combinations of modes. Their presence and their structure are rather sensitive to the way the average over u

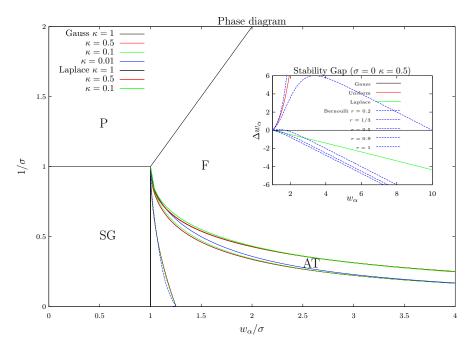


Fig. 2: Phase diagram in absence of bias and with a finite number of modes, with Gaussian and Laplace distributions for u and v. The dotted line separates the spin glass phase from the ferromagnetic phase under the RS hypothesis. The RS phase is unstable below the AT line. The influence of  $\kappa$  on the AT and SG-F lines is shown. In all cases, the hypothetical SG-F line lies well inside the broken RS phase. Inset: high temperature ( $\sigma = 0$ ) stability gap  $\Delta w_{\alpha}$  corresponding to a fixed point associated to a mode  $\beta$ , expressed as a function of  $w_{\alpha}$  and considering various distributions.

and v is performed. In this respect the case in which  $u^{\alpha}$  and  $v^{\alpha}$  have i.i.d. Gaussian components is very special: all fixed points associated to dominated modes can be shown to be unstable and fixed points associated to combinations of modes are not allowed. To see this, first notice that in such a case the magnetization's part of the saddle point equations (16,17) read

$$m_{\alpha} = (w_{\alpha}\bar{m}_{\alpha} - \theta_{\alpha})(1 - q) \tag{18}$$

$$\bar{m}_{\alpha} = (w_{\alpha} m_{\alpha} - \eta_{\alpha})(1 - \bar{q}). \tag{19}$$

Since the role of the bias is mainly to introduce some asymmetry between otherwise degenerated fixed points obtained by sign reversal of at least one pair  $(m_{\alpha}, \bar{m}_{\alpha})$ , let us analyze the situation without fields, i.e. by setting  $\eta = \theta = 0$ . We immediately see that as long as the singular values are non degenerate, only one single mode may condense at a time. Indeed if mode  $\alpha$  condenses we necessarily have

$$w_{\alpha}^{2}(1-q)(1-\bar{q})=1,$$

and this can be verified only by one mode at a time. Looking at the stability of the fixed points, we see that only the fixed point associated to the largest singular value is actually stable (details reported after the introduction of lemma 3.1).

For other distributions like uniform Bernoulli or Laplace, instead, stable fixed points associated to many different single modes or combinations of modes can exist and contribute to the thermodynamics. In order to analyze this question in more general terms we first rewrite the mean-field equations in a convenient way which require some preliminary remarks. We restrict the discussion to i.i.d. variables so that we can consider single variable distributions. Joint distributions will be distinguished from single variable distributions by the use of bold:  $\mathbf{u} = \{u^{\alpha}, \alpha = 1, \dots, K\}$ , K being the (finite) number of modes susceptible of condensing.

Given the distribution p and assuming it to be even, we define a related distribution  $p^*$  attached to mode  $\alpha$ :

$$p^{\star}(u) \stackrel{\text{def}}{=} - \int_{-\infty}^{u} x p(x) dx = \int_{|u|}^{\infty} x p(x) dx, \tag{20}$$

This distribution has some useful properties.

**Lemma 3.1.** Given that p is centered with unit variance and kurtosis  $\kappa_u$ ,  $p^*$  is a centered probability distribution with variance

$$\int_{-\infty}^{\infty} u^2 p^{\star}(u) du = \frac{\kappa_u}{3}.$$

**Proof.** Consider the moments of  $p^*$ . For n odd they vanish while for n even they read:

$$\int_{-\infty}^{+\infty} u^n p^*(u) du = 2 \int_0^{\infty} u^n p^*(u) du$$
$$= 2 \int_0^{\infty} du u^n \int_u^{\infty} x p(x) dx$$
$$= 2 \int_0^{\infty} x p(x) dx \int_0^x u^n du$$
$$= \frac{1}{n+1} \int_{-\infty}^{\infty} x^{n+2} p(x) dx,$$

i.e. the  $n_{th}$  even moments of  $p^*$  relate to moments of order n+2 of p. The lemma then follows from the fact that p has unit variance.

In this respect, the Gaussian averaging is special because we have  $\kappa_u = 3$  and  $p^* = p$ . Then the mean-field equations (16,17) corresponding to the magnetizations can be rewritten in a form similar to (18,19) by introducing the variables  $q_{\alpha}$  and  $\bar{q}_{\alpha}$ :

$$m_{\alpha} = (w_{\alpha}\bar{m}_{\alpha} - \theta_{\alpha})(1 - q_{\alpha}), \tag{21}$$

$$\bar{m}_{\alpha} = (w_{\alpha} m_{\alpha} - \eta_{\alpha})(1 - \bar{q}_{\alpha}), \tag{22}$$

with

$$q_{\alpha} = \int dx \frac{e^{-x^{2}/2}}{\sqrt{2\pi}} d\mathbf{v} p_{\alpha}(\mathbf{v}) \tanh^{2} \left( \kappa^{-\frac{1}{4}} \left( \sigma \sqrt{\bar{q}} x + \sum_{\gamma} (w_{\gamma} \bar{m}_{\gamma} - \theta_{\gamma}) v^{\gamma} \right) \right), \tag{23}$$

$$\bar{q}_{\alpha} = \int dx \frac{e^{-x^2/2}}{\sqrt{2\pi}} d\mathbf{u} p_{\alpha}(\mathbf{u}) \tanh^2 \left( \kappa^{\frac{1}{4}} \left( \sigma \sqrt{q} x + \sum_{\gamma} (w_{\gamma} m_{\gamma} - \eta_{\gamma}) u^{\gamma} \right) \right), \tag{24}$$

where

$$p_{\alpha}(\mathbf{u}) \stackrel{\text{def}}{=} p^{\star}(u^{\alpha}) \prod_{\beta \neq \alpha} p(u^{\beta}).$$

This rewriting will prove very useful also in the next section when analyzing the learning dynamics.

Let us now assume, in absence of bias, a non-degenerate fixed point associated to some given mode  $\beta$  with finite  $(m_{\beta}, \bar{m}_{\beta})$  and  $m_{\alpha} = \bar{m}_{\alpha} = 0, \forall \alpha \neq \beta$ . The fixed point equation imposes the relation

$$w_{\beta} = \frac{1}{\sqrt{(1 - q_{\beta})(1 - \bar{q}_{\beta})}} \stackrel{\text{def}}{=} w(q_{\beta}, \bar{q}_{\beta}). \tag{25}$$

The stability of such a fixed point with respect to any other mode  $\alpha$  is related to the positive definiteness of the following block of the Hessian

$$H_{\alpha\alpha} = \begin{bmatrix} w_{\alpha} & w_{\alpha}^{2} \mathsf{E}_{v,x} \left[ (v^{\alpha})^{2} \operatorname{sech}^{2} \left( \bar{h}(x,v) \right) \right] \\ w_{\alpha}^{2} \mathsf{E}_{u,x} \left[ (u^{\alpha})^{2} \operatorname{sech}^{2} \left( h(x,u) \right) \right] & w_{\alpha} \end{bmatrix}$$

with, in the present case

$$h(x,u) = \kappa^{\frac{1}{4}} \left( \sigma \sqrt{q} x + w_{\beta} \bar{m}_{\beta} u^{\beta} \right)$$
 and  $\bar{h}(x,v) = \kappa^{-\frac{1}{4}} \left( \sigma \sqrt{\bar{q}} x + w_{\beta} \bar{m}_{\beta} v^{\beta} \right)$ 

This reduces to

$$H_{\alpha\alpha} = \begin{bmatrix} w_{\alpha} & w_{\alpha}^{2}(1-q) \\ w_{\alpha}^{2}(1-\bar{q}) & w_{\alpha} \end{bmatrix}.$$

Therefore for the Gaussian averaging case, since  $q_{\beta} = q$ ,  $\bar{q}_{\beta} = \bar{q}$  and given (25), we necessarily have

$$1 - (1 - q)(1 - \bar{q})w_{\alpha}^2 = 1 - \frac{w_{\alpha}^2}{w_{\beta}^2} < 0$$
 for  $w_{\alpha} > w_{\beta}$ ,

i.e. the Hessian has negative eigenvalues. This means that if the mode  $\beta$  is dominated by another mode  $\alpha$ , the magnetization  $(m_{\alpha}, \bar{m}_{\alpha})$  will develop until  $(1-q)(1-\bar{q})w_{\alpha}^2=1$ , while  $m_{\beta}$  will vanish.

For the general case of i.i.d. variables, assuming  $u^{\alpha}$  and  $v^{\alpha}$  obey the same distribution p, let F and  $F_{\alpha}$  be the cumulative distributions associated respectively to p and  $p_{\alpha}$ 

$$F(u) \stackrel{\text{def}}{=} \int_{-\infty}^{u} p(x) dx$$

$$F_{\alpha}(u) \stackrel{\text{def}}{=} \int d\mathbf{u} \ \theta(u - u^{\alpha}) p_{\alpha}(\mathbf{u}) dx = -\int_{-\infty}^{u} du^{\alpha} \int_{-\infty}^{u^{\alpha}} x p(x) dx.$$

Given the values of  $(q, \bar{q})$  obtained from the fixed point associated to mode  $\beta$ , we have the following property:

#### Proposition 3.2. If

(i) 
$$F_{\beta}(u) < F(u)$$
,  $\forall u \in \mathbb{R}^+$  then  $q_{\beta} > q$  and  $\bar{q}_{\beta} > \bar{q}$ ,

(ii) 
$$F_{\beta}(u) > F(u)$$
,  $\forall u \in \mathbb{R}^+$  then  $q_{\beta} < q$  and  $\bar{q}_{\beta} < \bar{q}$ ,

which in turn implies

$$w(q, \bar{q}) < w_{\beta}$$
 (i) and  $w(q, \bar{q}) > w_{\beta}$  (ii)

with

$$w(q,\bar{q}) \stackrel{\text{def}}{=} \frac{1}{\sqrt{(1-q)(1-\bar{q})}}.$$

**Proof.** This is obtained by straightforward by parts integration respectively over u and v in equations (16,17), relative to magnetizations.

In other words if  $F_{\beta}$  dominates F on  $\mathbb{R}^+$  then there is a positive stability gap defined as

$$\Delta w_{\beta} \stackrel{\text{def}}{=} w(q, \bar{q}) - w_{\beta} \tag{26}$$

such that there is a non-empty range for higher values of  $w_{\alpha} \in [w_{\beta}, w(q, \bar{q})]$  for which the fixed point associated to mode  $\beta$  corresponds to a local minimum of the free energy. Note that property (i) [resp. (ii)] is analogous (in the sense that it implies it) to  $p_{\beta}$  having a larger [resp. smaller] variance than p, i.e.  $\kappa_u > 3$  [resp.  $\kappa_u < 3$ ]. Therefore distributions p with negative relative kurtosis ( $\kappa_u - 3$ ) will tend to favor the presence of metastable states, while the situation will tend to be more complex for probabilities with positive relative kurtosis. Indeed, in the latter case the fixed point associated to the highest mode  $\alpha_{max}$  might not correspond to a stable state if lower modes in the range  $[w(q,\bar{q}),w_{\alpha_{max}}]$  are present, and fixed points associated to combinations of modes have to be considered. Note that in contrary with the Gaussian case, this can happen because  $q_{\alpha}$  is different for each mode and therefore more flexibility is offered by equations (21,22) than from equations (18,19).

Let us give some examples. Denote by  $\gamma_u \stackrel{\text{def}}{=} \kappa_u - 3$  the relative kurtosis. As already said the Gaussian distribution is a special case with  $\gamma_u = 0$ . In addition, for instance for p corresponding to Bernoulli, Uniform or Laplace, we have the following properties illustrated in the inset of Figure 2:

• Bernoulli  $(\gamma_u = -2)$ :

$$p(u) = \frac{1}{2} (\delta(u+1) + \delta(u-1)), \qquad F(u) = \frac{1}{2} (\theta(u+1) + \theta(u-1))$$
$$p_{\alpha}(u) = \frac{1}{2} \theta(1 - u^{2}), \qquad F_{\alpha}(u) = \frac{1}{2} \theta(1 - u^{2})(u+1) + \theta(u-1)$$

then  $F_{\alpha}(u) > F(u)$  for u > 0, yielding a positive stability gap.

• Uniform  $(\gamma_u = -6/5)$ :

$$p(u) = \frac{1}{2\sqrt{3}}\theta(3-u^2), \qquad F(u) = \frac{1}{2\sqrt{3}}\theta(3-u^2)(u+\sqrt{3}) + \theta(u-\sqrt{3})$$
$$p_{\alpha}(u) = \frac{1}{4\sqrt{3}}\theta(3-u^2)(3-u^2), \qquad F_{\alpha}(u) = \frac{1}{4\sqrt{3}}\theta(3-u^2)(3u-\frac{u^3}{3}+2\sqrt{3}) + \theta(u-\sqrt{3}).$$

It can be verified that  $F_{\alpha}(u) > F(u)$  for u > 0, yielding again a positive stability gap.

• Laplace  $(\gamma_u = 3)$ :

$$\begin{split} p(u) &= \frac{1}{\sqrt{2}} e^{-\sqrt{2}|u|}, \qquad F(u) = \frac{1}{2} + \frac{u}{2|u|} \left( 1 - e^{-\sqrt{2}|u|} \right) \\ p_{\alpha}(u) &= \frac{1}{2} \left( |u| + \frac{1}{\sqrt{2}} \right) e^{-\sqrt{2}|u|}, \qquad F_{\alpha}(u) = F(u) - \frac{u}{2\sqrt{2}} e^{-\sqrt{2}|u|}. \end{split}$$

Here we have  $F_{\alpha}(u) < F(u)$  for u > 0, yielding a negative stability gap.

These three examples fall either in condition (i) or (ii), with a stability gap  $\Delta w_{\beta}$  that is either always positive or always negative, independently of  $w_{\beta}$ . We can also provide examples for which the stability condition may vary with  $w_{\beta}$ . Consider for instance a sparse Bernoulli distribution, with  $r \in [0, 1]$  a sparsity parameter:

$$p(u) = \frac{r}{2} \left( \delta(u + \frac{1}{\sqrt{r}}) + \delta(u - \frac{1}{\sqrt{r}}) \right) + (1 - r)\delta(u).$$

The relative kurtosis is in this case

$$\gamma_u(r) = \frac{1}{r} - 3.$$

Looking at F(u) and  $F_{\alpha}(u)$  it is seen that both conditions (i) and (ii) are not fulfilled, except for r=1 which corresponds to the plain Bernoulli case. As we see in the inset of Figure 2, for r<1/3 the stability gap is always negative, meaning that a unimodal ferromagnetic phase is not stable, and it is replaced by a compositional ferromagnetic phase at all temperatures. Instead, for r>1/3 and at sufficiently high temperature (low  $w_{\alpha}$ ) the single mode fixed point dominate the ferromagnetic phase.

**Laplace distribution:** let us look at the properties of the phase diagram in the case of singular vectors' components being Laplace i.i.d., case in which a negative stability gap is expected and it may lead to a compositional phase. For this we need the expression for a sum of Laplace variables to compute the averages involved in (16,17). For this purpose, we define the following distributions:

$$f(s) = \int \prod_{\gamma} du^{\gamma} \frac{\lambda_{\gamma}}{2} e^{-\lambda_{\gamma} |u^{\gamma}|} \, \delta(s - \sum_{\gamma} u^{\gamma}),$$

$$g_{\alpha}(s) = \int du^{\alpha} \frac{\lambda_{\alpha}}{4} (\lambda_{\alpha} |u^{\alpha}| + 1) e^{-\lambda_{\alpha} |u^{\alpha}|} \prod_{\alpha \neq \alpha} du^{\gamma} \frac{\lambda_{\gamma}}{2} e^{-\lambda_{\gamma} |u^{\gamma}|} \, \delta(s - \sum_{\alpha} u^{\gamma}).$$

Their Laplace transform upon decomposing into partial fractions reads:

$$\tilde{f}(\omega) = \prod_{\gamma} \frac{\lambda_{\gamma}^{2}}{\lambda_{\gamma}^{2} - \omega^{2}} = \sum_{\gamma} C_{\gamma} \frac{\lambda_{\gamma}^{2}}{\lambda_{\gamma}^{2} - \omega^{2}}$$

and

$$\begin{split} \tilde{g}_{\alpha}(\omega) &= \frac{\lambda_{\alpha}^{2}}{\lambda_{\alpha}^{2} - \omega^{2}} \prod_{\gamma} \frac{\lambda_{\gamma}^{2}}{\lambda_{\gamma}^{2} - \omega^{2}} \\ &= C_{\alpha} \frac{\lambda_{\alpha}^{4}}{(\lambda_{\alpha}^{2} - \omega^{2})^{2}} + \sum_{\gamma \neq \alpha} C_{\gamma} \frac{\lambda_{\gamma}^{2} \lambda_{\alpha}^{2}}{\lambda_{\alpha}^{2} - \lambda_{\gamma}^{2}} \Big( \frac{1}{\lambda_{\gamma}^{2} - \omega^{2}} - \frac{1}{\lambda_{\alpha}^{2} - \omega^{2}} \Big). \end{split}$$

where

$$C_{\gamma} \stackrel{\text{def}}{=} \prod_{\delta \neq \gamma} \frac{\lambda_{\delta}^2}{\lambda_{\delta}^2 - \lambda_{\gamma}^2}.$$

From these decompositions we immediately identify

$$f(s) = \frac{1}{2} \sum_{\gamma} C_{\gamma} \lambda_{\gamma} e^{-\lambda_{\gamma}|s|},$$

$$g_{\alpha}(s) = \frac{\lambda_{\alpha} C_{\alpha}}{4} (\lambda_{\alpha}|s| + 1)e^{-\lambda_{\alpha}|s|} + \frac{1}{2} \sum_{\gamma \neq \alpha} C_{\gamma} \frac{\lambda_{\gamma} \lambda_{\alpha}}{\lambda_{\alpha}^{2} - \lambda_{\gamma}^{2}} (\lambda_{\alpha} e^{-\lambda_{\gamma}|s|} - \lambda_{\gamma} e^{-\lambda_{\alpha}|s|}).$$

This results in the following decomposition of the EA parameters:

$$q = \int dx ds \frac{e^{-\sqrt{2}|s|-x^2/2}}{2\sqrt{\pi}} \sum_{\gamma} C_{\gamma}[\bar{m}] \tanh^2(\bar{h}_{\gamma}(x,s))$$
(27)

$$q_{\alpha} = \int dx ds \frac{e^{-\sqrt{2}|s|-x^{2}/2}}{2\sqrt{\pi}} \left[ \frac{1}{\sqrt{2}} (|s| + \frac{1}{\sqrt{2}}) C_{\alpha}[\bar{m}] \tanh^{2}(\bar{h}_{\alpha}(x,s)) \right]$$
(28)

$$+\sum_{\gamma\neq\alpha}C_{\gamma}[\bar{m}]\frac{(w_{\gamma}\bar{m}_{\gamma}-\theta_{\gamma})^{2}\tanh^{2}(\bar{h}_{\gamma}(x,s))-(w_{\alpha}\bar{m}_{\alpha}-\theta_{\alpha})^{2}\tanh^{2}(\bar{h}_{\alpha}(x,s))}{(w_{\gamma}\bar{m}_{\gamma}-\theta_{\gamma})^{2}-(w_{\alpha}\bar{m}_{\alpha}-\theta_{\alpha})^{2}}\Big]$$
(29)

with

$$\bar{h}_{\gamma}(x,s) \stackrel{\text{def}}{=} \kappa^{-\frac{1}{4}} \left( \sigma \sqrt{\bar{q}} x + (w_{\gamma} \bar{m}_{\gamma} - \theta_{\gamma}) s \right)$$

and

$$C_{\gamma}[\bar{m}] \stackrel{\text{def}}{=} \prod_{\delta \neq \gamma} \frac{(w_{\gamma}\bar{m}_{\gamma} - \theta_{\gamma})^2}{(w_{\gamma}\bar{m}_{\gamma} - \theta_{\gamma})^2 - (w_{\delta}\bar{m}_{\delta} - \theta_{\delta})^2}.$$

This allows for an efficient resolution of the mean-field equations (16,17,21,22), which let us observe the appearance of a purely compositional phase in the ferromagnetic domain when the modes at the top of the spectrum get close enough. In order to characterize this phase, we consider the stability gap  $\Delta^{(n)}(w_{\alpha})$  for which the range  $[w_a - \Delta^{(n)}(w_{\alpha}), w_a]$  lies below the highest mode  $w_a$ , such that the ferromagnetic states correspond to the condensation of n distinct modes present in this interval, including the highest.

In addition, this will prove useful when analyzing the learning dynamics described in the next section.

# 4 Learning dynamics of the RBM

## 4.1 Learning dynamics in the thermodynamic limit

A mean field analysis of the learning dynamics has been proposed in [25], in the form of phenomenological equations obtained after averaging over some parameters of the RBM, i.e. by choosing a well defined statistical ensemble of RBMs and using self-averaging properties in the thermodynamic limit. Here we rederive these equations, we add some details and then explore their properties in the light of the preceding section.

First we project the gradient ascent equations (4-6) onto the bases  $\{u_{\alpha}(t) \in \mathbb{R}^{N_v}\}$  and  $\{v_{\alpha}(t) \in \mathbb{R}^{N_h}\}$  defined by the SVD of W. Discarding stochastic fluctuations usually inherent to the learning procedure and letting the learning rate  $\gamma \to 0$ , the continuous version of (4-6) can be recast as follows:

$$\frac{1}{L} \left( \frac{dW}{dt} \right)_{\alpha\beta} = \langle s_{\alpha} \sigma_{\beta} \rangle_{\text{Data}} - \langle s_{\alpha} \sigma_{\beta} \rangle_{\text{RBM}}, \tag{30}$$

$$\frac{1}{\sqrt{L}} \left( \frac{d\eta}{dt} \right)_{\alpha} = \langle s_{\alpha} \rangle_{\text{RBM}} - \langle s_{\alpha} \rangle_{\text{Data}}, \tag{31}$$

$$\frac{1}{\sqrt{L}} \left( \frac{d\theta}{dt} \right)_{\alpha} = \langle \sigma_{\alpha} \rangle_{\text{RBM}} - \langle \sigma_{\alpha} \rangle_{\text{Data}}, \tag{32}$$

with  $s_{\alpha}$  and  $\sigma_{\alpha}$  given in (9). We also have

$$\left(\frac{dW}{dt}\right)_{\alpha\beta} = \delta_{\alpha,\beta} \frac{dw_{\alpha}}{dt} + (1 - \delta_{\alpha,\beta}) \left(w_{\beta}(t)\Omega_{\beta\alpha}^{v}(t) + w_{\alpha}(t)\Omega_{\alpha\beta}^{h}\right)$$

$$\frac{1}{\sqrt{L}} \left(\frac{d\eta}{dt}\right)_{\alpha} = \frac{d\eta_{\alpha}}{dt} - \sum_{\beta} \Omega_{\alpha\beta}^{v} \eta_{\beta}$$

$$\frac{1}{\sqrt{L}} \left(\frac{d\theta}{dt}\right)_{\alpha} = \frac{d\theta_{\alpha}}{dt} - \sum_{\beta} \Omega_{\alpha\beta}^{h} \theta_{\beta}$$

where

$$egin{align} \Omega^v_{lphaeta}(t) &= -\Omega^v_{etalpha} \stackrel{ ext{def}}{=} rac{doldsymbol{u}^{lpha,T}}{dt} oldsymbol{u}^eta \ & \Omega^h_{lphaeta}(t) &= -\Omega^h_{etalpha} \stackrel{ ext{def}}{=} rac{doldsymbol{v}^{lpha,T}}{dt} oldsymbol{v}^eta \end{aligned}$$

By eliminating  $\left(\frac{dw}{dt}\right)_{\alpha\beta}$ ,  $\left(\frac{d\eta}{dt}\right)_{\alpha}$  and  $\left(\frac{d\theta}{dt}\right)_{\alpha}$  we get the following set of dynamical equations:

$$\frac{1}{L}\frac{dw_{\alpha}}{dt} = \langle s_{\alpha}\sigma_{\alpha}\rangle_{\text{Data}} - \langle s_{\alpha}\sigma_{\alpha}\rangle_{\text{RBM}}$$
(33)

$$\frac{d\eta_{\alpha}}{dt} = \langle s_{\alpha} \rangle_{\text{RBM}} - \langle s_{\alpha} \rangle_{\text{Data}} + \sum_{\beta} \Omega_{\alpha\beta}^{\nu} \eta_{\beta}$$
 (34)

$$\frac{d\theta_{\alpha}}{dt} = \langle \sigma_{\alpha} \rangle_{\text{RBM}} - \langle \sigma_{\alpha} \rangle_{\text{Data}} + \sum_{\beta} \Omega_{\alpha\beta}^{h} \theta_{\beta}$$
 (35)

along with the infinitesimal rotation generators of the left and right singular vectors

$$\Omega_{\alpha\beta}^{v}(t) = -\frac{1}{w_{\alpha} + w_{\beta}} \left(\frac{dW}{dt}\right)_{\alpha\beta}^{A} + \frac{1}{w_{\alpha} - w_{\beta}} \left(\frac{dW}{dt}\right)_{\alpha\beta}^{S}$$
(36)

$$\Omega_{\alpha\beta}^{h}(t) = \frac{1}{w_{\alpha} + w_{\beta}} \left(\frac{dW}{dt}\right)_{\alpha\beta}^{A} + \frac{1}{w_{\alpha} - w_{\beta}} \left(\frac{dW}{dt}\right)_{\alpha\beta}^{S}$$
(37)

where

$$\left(\frac{dW}{dt}\right)_{\alpha\beta}^{A,S} \stackrel{\text{def}}{=} \frac{1}{2} \left( \langle s_{\alpha} \sigma_{\beta} \rangle_{\text{Data}} \pm \langle s_{\beta} \sigma_{\alpha} \rangle_{\text{Data}} \mp \langle s_{\beta} \sigma_{\alpha} \rangle_{\text{RBM}} - \langle s_{\alpha} \sigma_{\beta} \rangle_{\text{RBM}} \right).$$

The dynamics of learning is now expressed in the reference frame defined by the singular vectors of W. The skew-symmetric rotation generators  $\Omega^{v,h}_{\alpha\beta}(t)$  of the basis vectors (induced by the dynamics) tell us how data rotate relatively to this frame. Given the initial conditions, these help us keeping track of the representation of data in this frame. Note that these equations become singular when some degeneracy occurs in W because then the SVD is not uniquely defined. Except from the numerical point of view, where some regularizations might be needed, this does not constitute an issue. In fact only rotations among non-degenerate modes are meaningful, while the rest corresponds to gauge degrees of freedom.

At this point our set of dynamical equations (33-37) is written in a general form. Our goal is to find the typical trajectory of the RBM within a certain statistical ensemble. For this reason, we make the hypothesis that the learning dynamics is represented by a trajectory in the space  $\{w_{\alpha}(t), \eta_{\alpha}(t), \theta_{\alpha}(t), \Omega_{\alpha\beta}^{v,h}(t)\}$ , while the specific realization of  $u_i^{\alpha}$ ,  $v_i^{\alpha}$  and  $r_{ij}$  in (8) can be considered irrelevant and only the way they are distributed is important. We are then allowed to perform an average over  $u_i^{\alpha}$ ,  $v_i^{\alpha}$ and  $r_{ij}$  with respect to some simple distributions, as long as this average is correlated with the data. By this we mean that the components  $s_{\alpha}$  of any given sample are kept fixed while averaging. In the end, what really matters are the strength and the rotation of the SVD modes, respectively determined by  $w_{\alpha}(t)$  and  $\Omega_{\alpha\beta}^{v,h}(t)$ . As a simplification and also by lack of understanding of what intrinsically drives their evolution, the distributions of  $u_i^{\alpha}$  and  $v_j^{\alpha}$  will be considered stationary in the sequel. Concerning  $r_{ij}$ , we allow its variance  $\sigma^2/L$  to vary with time in order to give a minimal description of how the MP bulk evolves during the learning. The detailed dynamics of  $\sigma$  will be derived later in Section 4.3. Using the same notation of Section 3.4 and in particular using the rescaling  $v \sim \sqrt{N_h} v_i^{\alpha}$ , the empirical terms take the form:

$$\langle \sigma_{\alpha} \rangle_{\text{Data}} = \langle (s_{\alpha} w_{\alpha} - \theta_{\alpha}) (1 - q_{\alpha} [\mathbf{s}]) \rangle_{\text{Data}}$$
 (38)

$$\langle s_{\alpha}\sigma_{\beta}\rangle_{\text{Data}} = \langle s_{\alpha}(s_{\beta}w_{\beta} - \theta_{\beta})(1 - q_{\beta}[\mathbf{s}])\rangle_{\text{Data}}$$
(39)

where

$$q_{\alpha}[\mathbf{s}] \stackrel{\text{def}}{=} \int dx \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}} d\mathbf{v} p_{\alpha}(\mathbf{v}) \tanh^2\left(\kappa^{-\frac{1}{4}} \left(\sigma x + \sum_{\gamma} (w_{\gamma} s_{\gamma} - \theta_{\gamma}) v^{\gamma}\right)\right),$$

Note that the last equation actually depends on the activation function (hyperbolic tangent in this case), and the term  $\sigma x$  corresponds to  $\sum_k r_{kj}s_k$  and is obtained by central limit theorem from the independence of the  $r_{kj}$ .  $q_{\alpha}[\mathbf{s}]$  is the empirical counterpart of the EA parameters q and  $q_{\alpha}$  already encountered in Section 3.4, and for simple i.i.d. distributions like Gaussian or Laplace it can be estimated easily. The main point here is that the empirical terms (38,39) define operators whose decomposition over the SVD modes of W functionally depends only on  $w_{\alpha}$ ,  $\theta_{\alpha}$  and on the projection of the data over the SVD modes of W. These terms are driving the dynamics in a precise way. The adaptation of the RBM to this driving force is given by the  $\langle \dots \rangle_{\text{RBM}}$  terms in (33,34,35), which can be estimated in the thermodynamic limit (see Section 4.3) as a function of  $w_{\alpha}$ ,  $\theta_{\alpha}$  and  $\eta_{\alpha}$  alone, by means of the order parameters  $(m_{\alpha}, \bar{m}_{\alpha})$  given

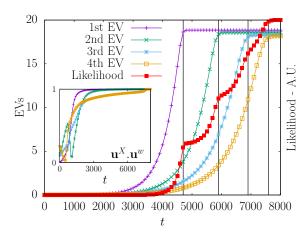


Fig. 3: Time evolution of the eigenvalues and of the likelihood in the linear model. We observe very clearly how the different modes emerge from the bulk and how the likelihood increases with each learned eigenvalue. In the inset, the scalar product of the vectors  $\boldsymbol{u}$  obtained from the SVD of the data and from the weights  $\boldsymbol{w}$ . The  $\boldsymbol{u}$ s of  $\boldsymbol{w}$  are aligned with the SVD of the data at the end of the learning.

in Section 3.2 and once the mean-field equations (16,17) have been solved. Of course, all of this is based on the hypothesis that the RBM stays in the RS domain during learning. Experimental evidence supports this hypothesis (see Section 5).

#### 4.2 Linear instabilities

At the beginning of the learning, the elements of the weight matrix W are usually small; therefore, we can analyze the linear behavior of the RBM in order to understand what happens. In particular, we will see that the dynamics of a non-linear RBM at the beginning of the learning can be understood by looking at the stability analysis of the learning process. The purpose of this analysis is to identify which "deformation modes" of the weight matrix are the most unstable, and how they are related to the input data. Additionally, a good feature of the linear case is that no averaging is needed, the dynamics being actually independent on the particular realization of the components  $u_i^{\alpha}$  and  $v_j^{\beta}$ . Also, always relative to the linear case, no distinction has to be made between dominant modes and other modes to be treated as the noise component of equation (8), we can simply put all of the modes on the same footing.

Let us analyze the linear regime for an RBM with binary units. The derivation is done by rescaling all the weights and fields by a common "inverse temperature"  $\beta$  and letting this go to zero in equation (4). In principle, the stability analysis would lead to assume both the weights and the magnetizations to be small. However, we can assume only the magnetizations to be small and consider a slightly more general case with no approximations. Such a case is analogous to a linear RBM whose magnetizations undergo Gaussian fluctuations, and it is derived by keeping up to quadratic terms of

the magnetizations in the mean field free energy:

$$F_{MF}(\mu,\nu) \simeq \frac{1}{2} \sum_{i=1}^{N} (1+\mu_i) \log(1+\mu_i) + (1-\mu_i) \log(1-\mu_i)$$

$$+ \frac{1}{2} \sum_{j=1}^{M} (1+\nu_j) \log(1+\nu_j) + (1-\nu_j) \log(1-\nu_j)$$

$$- \sum_{i,j} \left( W_{ij} \mu_i \nu_j - \frac{1}{2} W_{ij}^2 (\mu_i^2 + \nu_j^2) \right) + \sum_{i=1}^{N} \eta_i \mu_i + \sum_{j=1}^{M} \theta_j \nu_j$$

$$= \frac{1}{2\sigma_v^2} \sum_{i=1}^{N} \mu_i^2 + \frac{1}{2\sigma_h^2} \sum_{i=1}^{M} \nu_i^2 - \sum_{i,j} W_{ij} \mu_i \nu_j + \sum_{i=1}^{N} \eta_i \mu_i + \sum_{j=1}^{M} \theta_j \nu_j.$$

where the variances  $(\sigma_v^2, \sigma_h^2)$  of respectively visible and hidden variables read  $(N_h < N_v)$ :

$$\sigma_v^{-2} = 1 + \sum_i W_{ij}^2 \simeq 1 + \sum_{\alpha} w_{\alpha}^2$$
 (40)

$$\sigma_h^{-2} = 1 + \sum_i W_{ij}^2 = 1 + \sum_{\alpha} w_{\alpha}^2. \tag{41}$$

We omitted the quadratic term in  $W_{ij}$  coming from the TAP contribution to the free energy, which is optional for our stability analysis. In absence of this term the modes evolve strictly independently, while taking it into account leads to a correction to individual variances which couples the modes.

Magnetizations  $(\mu, \nu)$  of visible and hidden variables have now Gaussian fluctuations with covariance matrix

$$C(\mu_v, \mu_h) \stackrel{\text{def}}{=} \begin{bmatrix} \sigma_v^{-2} & -W \\ -W^T & \sigma_h^{-2} \end{bmatrix}^{-1}$$

We can discard the biases of the data and the related fields  $(\theta_{\alpha}, \eta_{\alpha})$  with a proper centering of the variables, and we consider equation (33) directly involving the covariance matrix of the data expressed in the frame defined by the SVD modes of W

$$\langle s_{\alpha} \sigma_{\beta} \rangle_{\text{Data}} = \sigma_h^2 w_{\beta} \langle s_{\alpha} s_{\beta} \rangle_{\text{Data}}.$$

From  $C(\mu_v, \mu_h)$  we get the other terms yielding the following equations:

$$\frac{dw_{\alpha}}{dt} = w_{\alpha}\sigma_{h}^{2} \left( \langle s_{\alpha}^{2} \rangle_{\text{Data}} - \frac{\sigma_{v}^{2}}{1 - \sigma_{v}^{2}\sigma_{h}^{2}w_{\alpha}^{2}} \right) 
\Omega_{\alpha\beta}^{v,h} = (1 - \delta_{\alpha\beta})\sigma_{h}^{2} \left( \frac{w_{\beta} - w_{\alpha}}{w_{\alpha} + w_{\beta}} \mp \frac{w_{\beta} + w_{\alpha}}{w_{\alpha} - w_{\beta}} \right) \langle s_{\alpha}s_{\beta} \rangle_{\text{Data}}$$

Note that these equations are exact for a linear RBM, since they can be derived without any reference to the coordinates of  $u_{\alpha}$  and  $v_{\alpha}$  over which we average in the non-linear regime. These equations tell us that the learning dynamics drives the rotation of the

vectors  $\boldsymbol{u}^{\alpha}$  (and  $\boldsymbol{v}^{\alpha}$ ) until they are aligned to the principal components of the data, i.e. until  $\langle s_{\alpha}s_{\beta}\rangle_{\text{Data}}$  becomes diagonal. Calling  $\hat{w}_{\alpha}^2$  the empirical variance of the data, the system reaches the following equilibrium values:

$$w_{\alpha}^2 = \begin{cases} \frac{\hat{w}_{\alpha}^2 - \sigma_v^2}{\sigma_v^2 \sigma_h^2 \hat{w}_{\alpha}^2} & \text{if} & \hat{w}_{\alpha}^2 > \sigma_v^2, \\ 0 & \text{if} & \hat{w}_{\alpha}^2 \leq \sigma_v^2. \end{cases}$$

assuming  $(\sigma_v, \sigma_h)$  fixed. From this we see that the RBM selects the strongest SVD modes of the data. The linear instabilities correspond to directions along which the variance of the data is above the threshold  $\sigma_v^2$ , and they determine the development of the unstable deformation modes of the weight matrix; during the learning process, these modes will eventually interact following the usual mechanism of non-linear pattern formation encountered for instance in reaction-diffusion processes [34]. Other possible deformations are damped to zero. The linear RBM will therefore learn all the principal components that passed the threshold (up to  $N_h$ ). Note that this selection mechanism is already known to occur for linear auto-encoders [23] or other similar linear Boltzmann machines [22]. On Fig. 3 we can see the eigenvalues being learned one by one in a linear RBM.

If we take into account the expressions (40,41) for  $(\sigma_v, \sigma_h)$ , we see that the system cannot reach a stable solution except for the case in which all the modes are below the threshold at the beginning. Otherwise the modes that are excited first will eventually grow like  $\sqrt{t}$  for a large time, and the excitation threshold will tend to zero for all modes.

In any case, by the definition of a multivariate Gaussian, this simple non-linear analysis describes a unimodal distribution. In order to properly understand the dynamics and the steady-state regime of a non-linear RBM, a well suited mean-field theory is required.

# 4.3 Non-linear regime

In the linear regime, some specific modes are selected and at some point they start to interact in a non-trivial manner. As seen explicitly in (39), the empirical terms in (4-6) involve higher order statistics of the data and then the Gaussian estimation with  $\sigma_v^2 = \sigma_h^2 = 1$  of the RBM response terms  $\langle s_\alpha \rangle_{\text{RBM}}$  and  $\langle s_\alpha \sigma_\beta \rangle_{\text{RBM}}$  is no longer valid when the interactions kick in. Schematically, the linear regime is valid as long as the RBM is found in the paramagnetic phase. But as soon as one mode passes the linear threshold, the system enters the ferromagnetic phase. Then the proper estimation of the response terms follows from the thermodynamic analysis performed in Section 3, and depends on the assumptions made on the statistical properties of the components of the singular vectors of the weight matrix. In the case of Gaussian i.i.d. components, given the analysis proposed in Section 3.4, we know that the mode with the highest singular value completely dominates the ferromagnetic phase: we expect one single ferromagnetic state characterized by magnetizations aligned to this mode only, while magnetizations correlated to other modes vanish. To be precise, this is the correct picture without fields ( $\eta = \theta = 0$ ) but we don't expect this picture to drastically change in the case of non-vanishing fields. In fact, solving the mean-field equations in presence of the fields show the appearance of meta-stable states correlated with single dominated modes; however, the free energy difference with respect to the ground state, i.e. the state correlated with the mode with the highest singular value,

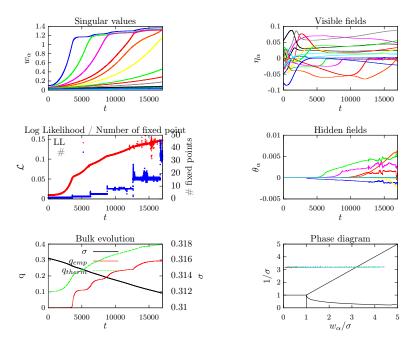


Fig. 4: Predicted mean evolution of an RBM of size  $(N_v, N_h) = (1000, 500)$  learned on a synthetic dataset of  $10^4$  samples of size  $N_v = 1000$  obtained from a multimodal distribution with 20 clusters randomly defined on a submanifold of dimension d=15. The dynamics follows the projected magnetizations in this reduced space with help of 15 modes. We observe a kind of pressure on top singular values from lower ones.

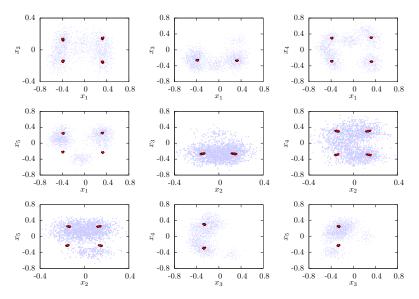


Fig. 5: Scatter plots of the mean-field magnetizations (in red) and the samples (in blue) in various plan projections defined by pairs of left eigenvectors of W. This case corresponds to an RBM of size  $(N_v, N_h) = (100, 50)$  learned on a synthetic dataset of  $10^4$  samples of size  $N_v = 100$  obtained from a multimodal distribution with 11 clusters randomly defined on a submanifold of dimension d = 5. The scatter plot is obtained at a point where 5 modes have already condensed and 16 saddle point solutions have been found.

is of order  $O(L(w_{\alpha} - w_{max}))$ , which means that the contribution of those meta-stable states become rapidly negligible with large system size.

To draw a realistic picture of the learning process we now consider Laplace i.i.d. components for the SVD modes that, as seen in Section 3.4, allow the ferromagnetic phase to be of compositional type. The reason for this is that the Laplace distribution leads to less interference among modes than the Gaussian distribution, so that the modes will weakly interact in the mean-field equations. Solving equations (21,22,27,29) in absence of fields yields the following picture: one fixed point solution will typically have non-vanishing magnetizations  $\{m_{\alpha}, \bar{m}_{\alpha}\}$  for all  $\alpha$  such that  $w_{\alpha} \in [w_{max} - \Delta w, w_{max}]$ , where  $\Delta w$  is approximately the gap  $\Delta w(q, \bar{q})$  defined in (26). This solution is a degenerate ground state, all other solutions being obtained by independently reversing the signs of the condensed magnetizations  $(m_{\alpha}, \bar{m}_{\alpha})$ . Hence for K condensed modes we get a degeneracy of  $2^K$ . When the fields are included, all these fixed points are displaced in the direction of the fields, and some of them may disappear. In the end we are left with a potentially large amount of nearly degenerate states able to cover the empirical distribution of the data, at least in some simple cases.

Coming back to the learning dynamics the terms corresponding to the response of

the RBM in (4,6) are estimated in the thermodynamic limit by means of the previously defined order parameters:

$$\langle s_{\alpha} \rangle_{\text{RBM}} = \frac{1}{Z_{\text{Therm}}} \sum_{\omega} e^{-Lf(m^{\omega}, \bar{m}^{\omega}, q^{\omega}, \bar{q}^{\omega})} \bar{m}_{\alpha}^{\omega} \stackrel{\text{def}}{=} \langle \bar{m}_{\alpha} \rangle_{\text{Therm}},$$
$$\langle s_{\alpha} \sigma_{\beta} \rangle_{\text{RBM}} = \frac{1}{Z_{\text{Therm}}} \sum_{\omega} e^{-Lf(m^{\omega}, \bar{m}^{\omega}, q^{\omega}, \bar{q}^{\omega})} \bar{m}_{\alpha}^{\omega} m_{\beta}^{\omega} \stackrel{\text{def}}{=} \langle \bar{m}_{\alpha} m_{\beta} \rangle_{\text{Therm}}.$$

Here  $\langle \dots \rangle_{\text{Therm}}$  denotes the thermodynamical average and the partition function is expressed, in the thermodynamic limit, as

$$Z_{\text{Therm}} \stackrel{\text{def}}{=} \sum_{\omega} e^{-Lf(m^{\omega}, \bar{m}^{\omega}, q^{\omega}, \bar{q}^{\omega})}$$

The index  $\omega$  runs over all the stable fixed point solutions of (16,17) weighted accordingly to the free energy given by (15). These are the dominant contributions as long as free energy differences are O(1), and the internal fluctuations given by each fixed point are comparatively of order O(1/L). In addition, the dynamics of the bulk can be characterized by empirically defining  $\sigma^2$ :

$$\sigma^2 = \frac{1}{L} \sum_{ij} r_{ij}^2,$$

whose evolution is:

$$\begin{split} \frac{d\sigma^2}{dt} &= \frac{1}{L} \sum_{ij} r_{ij} \frac{dW_{ij}}{dt}, \\ &= \frac{1}{L} \sum_{ij} r_{ij} \left[ \langle s_i \tanh \left( \sum_k r_{kj} s_k + \kappa^{-\frac{1}{4}} \sum_{\alpha} (w_{\alpha} s_{\alpha} - \theta_{\alpha}) v_j^{\alpha} \sqrt{L} \right) \rangle_{\text{Data}} - \langle s_i \sigma_j \rangle_{\text{RBM}} \right] \end{split}$$

given the independence of  $r_{i*}$  (resp.  $r_{*j}$ ) and  $u_i^{\alpha}$  (resp.  $v_i^{\alpha}$ ).

Exploiting the self-averaging properties of both the empirical and the response terms with respect to  $r_{ij}$ ,  $u_i^{\alpha}$  and  $v_j^{\alpha}$  yields

$$\frac{1}{L^2} \sum_{ij} r_{ij} \langle s_i \sigma_j \rangle_{\text{Data}} = \frac{\sigma^2}{L} \left( 1 - \langle q[\mathbf{s}] \rangle_{\text{Data}} \right)$$

$$\frac{1}{L^2} \sum_{ij} r_{ij} \langle s_i \sigma_j \rangle_{\text{RBM}} = \frac{\sigma^2}{L} \left( 1 - \langle q \rangle_{\text{Therm}} \right),$$

with

$$q[\mathbf{s}] \stackrel{\text{def}}{=} \int dx \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}} d\mathbf{v} p(\mathbf{v}) \tanh^2 \left(\kappa^{-\frac{1}{4}} \left(\sigma x + \sum_{\gamma} (w_{\gamma} s_{\gamma} - \theta_{\gamma}) v^{\gamma}\right)\right).$$

Summarizing, our equations take the suggestive form

$$\frac{1}{L}\frac{dw_{\alpha}}{dt} = \langle s_{\alpha}(w_{\alpha}s_{\alpha} - \theta_{\alpha})(1 - q_{\alpha}[\mathbf{s}])\rangle_{\text{Data}} - \langle \bar{m}_{\alpha}(w_{\alpha}\bar{m}_{\alpha} - \theta_{\alpha})(1 - q_{\alpha})\rangle_{\text{Therm}},$$
(42)

$$\frac{d\eta_{\alpha}}{dt} = \langle \bar{m}_{\alpha} \rangle_{\text{Therm}} - \langle s_{\alpha} \rangle_{\text{Data}} + \sum_{\beta} \Omega_{\alpha\beta}^{\nu} \eta_{\beta}, \tag{43}$$

$$\frac{d\theta_{\alpha}}{dt} = \langle (w_{\alpha}\bar{m}_{\alpha} - \theta_{\alpha})(1 - q_{\alpha}) \rangle_{\text{Therm}} - \langle (w_{\alpha}s_{\alpha} - \theta_{\alpha})(1 - q_{\alpha}[\mathbf{s}]) \rangle_{\text{Data}} + \sum_{\beta} \Omega_{\alpha\beta}^{h} \theta_{\beta},$$
(44)

$$\frac{d\sigma^2}{dt} = \sigma^2 \Big( \langle q \rangle_{\text{Therm}} - \langle q[\mathbf{s}] \rangle_{\text{Data}} \Big), \tag{45}$$

with  $\Omega^{v,h}$  taking the form of a difference between a data averaging  $\langle \dots \rangle_{\text{Data}}$  and a thermodynamical averaging  $\langle \dots \rangle_{\rm Therm}$  involving only order parameters. Note here that the  $w_{\alpha}$  variables, with respect to the other variables, evolve on a faster time scale. This is our final and main result, which might possibly help improving current learning algorithms of RBMs. From this, it is clear what the learning of an RBM is aimed at: the equations will converge once the dataset is clustered in such a way that each cluster is represented by a solution of the mean-field equations with magnetizations  $\bar{m}_{\alpha}$  and EA parameters  $q_{\alpha}$  corresponding respectively to their empirical counterparts  $\langle s_{\alpha} \rangle$  and  $\langle q_{\alpha}[\mathbf{s}] \rangle$  representing cluster magnetization and variance. In particular, these clusters can somehow be regarded as the attractors in the context of feed-forward networks, defining a partition of the data. This can be seen by starting from random configurations and letting the system evolve using the TAP equations or a MCMC method. At the end the system will end up in one of those clusters (characterized by a fixed point of the mean-field equations). Note that this is the reason why the RBM needs to reach a ferromagnetic phase with many states to be able to match the empirical term in (4) and reach convergence.

Additionally, the log likelihood (3) can be estimated in the thermodynamic limit (after normalization by L).

$$\mathcal{L} = \left\langle \sqrt{\kappa} \mathsf{E}_{x,v} \left[ \log \cosh \left( \kappa^{-\frac{1}{4}} \left( \sigma x + \sum_{\alpha} (w_{\alpha} s_{\alpha} - \theta_{\alpha}) v^{\alpha} \right) \right) \right] \right\rangle_{\mathrm{Data}} - \left\langle \sum_{\alpha} \eta_{\alpha} s_{\alpha} \right\rangle_{\mathrm{Data}} - \frac{1}{L} \log \left( Z_{\mathrm{Therm}} \right),$$

As an example, for a multimodal data distribution with a finite number of clusters embedded in a high dimensional configuration space, the SVD modes of W that will develop are the one pointing to the directions of the magnetizations defined by these clusters (which will be almost surely orthogonal, given the high dimensionality of the embedding space). In this simple case the RBM will evolve, as in the linear case, to a state in which the empirical term becomes diagonal, while the singular values will adjust to match the proper magnetization in each fixed point.

We have integrated equations (42,43,44,45,36,37) in simple cases by using the Laplace averaging of the components of the SVD modes and using for the EA parameters the expressions given in (27,29). Basically, the hidden distribution to be

modeled is defined by

$$P(s) = \sum_{c=1}^{C} p_c \prod_{i=1}^{N} \frac{e^{h_i^c s_i}}{2 \cosh(h_i^c)},$$
(46)

i.e. a multimodal distribution composed of C clusters of independent variables, where the magnetization of each variable i in cluster c is given by  $m_i^c = \tanh(h_i^c)$ . Each cluster is weighted by some probability  $p_c$ . In addition we assume these magnetization vectors  $m^c$  to be embedded in a low dimensional space of dimension  $d \ll N$ . d defines the rank of W. The initial conditions for W are such that the left singular vectors  $\{u_{\alpha}, \alpha = 1, \dots d\}$  span this low dimensional space. An example of the typical dynamics obtained in the case at hand is shown in Figure 4. In contrast to the linear problem where singular values evolve independently, here we distinctively witness the interaction between singular values: a kind of pressure is exerted by lower modes on higher ones resulting in successive bumps in the dynamics of the top modes. The number of states is roughly multiplied by two each time a mode condenses and get close enough to the top modes. Concerning the dynamics of the fields, we don't really observe convergence towards stable directions. Some (possibly numerical) instability is observed when many modes condense, with both the fields and the number of fixed point solutions becoming very noisy. It is also interesting to see how the magnetizations related to the states are distributed with respect to the dataset. On Figure 5 we see that the fixed points tend (as expected) to settle within dense regions of sample points. However, our coarse description shows some limitations for more complex situations, the number of adjustable parameters being too limited to be able to match arbitrary distributions of clusters. It is then appropriate to think about this behaviour in a mean sense; at least, it is able to reproduce a realistic learning dynamics of the singular values of the weight matrix.

#### 5 Numerical Experiments

Given the comprehensive theoretical analysis of the RBM model given in the previous sections, we are now able to provide a meaningful description of the learning dynamics for a RBM trained with k-steps contrastive divergence (CDk) [4]. The observations presented in this section will serve as a validation for the theoretical analysis. First, to provide a more direct comparison to section 4.3, we will look at the learning dynamics of an RBM trained on a set of simple synthetic data. Subsequently, we will test the model against real world data by training on the MNIST dataset.

#### 5.1 Synthetic dataset

As a simple case, we trained the RBM over the same dataset defined in fig. 4, derived from the simple multimodal distribution in eq. 46 (see Appendix B for details). Thus we set  $N_v = 1000$ ,  $N_h = 500$  and we trained using  $10^4$  samples with an effective dimension d=15 organized in 20 separate clusters. The weights are initialized from a Gaussian distribution with standard deviation  $\sigma = 10^{-3}$ , while the hidden bias is initialized to 0 and the visible bias is initialized with the empirical mean of the data

$$\eta_i = \frac{1}{2} \log \left( \frac{p_i}{1 - p_i} \right)$$

where  $p_i$  is the empirical probability of activation for the  $i_{th}$  hidden node.

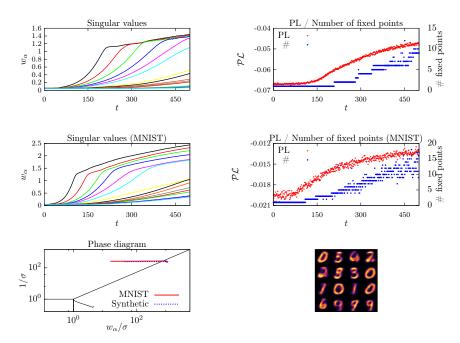


Fig. 6: Experimental evolution of an RBM during training for a synthetic dataset (top plots, to compare to Fig. 4) and for MNIST (central plots). The bottom left plot shows the learning trajectories in the phase diagram, while the bottom right image shows some examples of fixed point solutions for MNIST (we note the presence of some spurious fixed points).

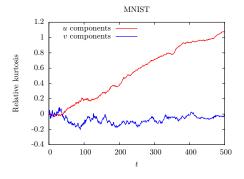


Fig. 7: Relative kurtosis of the components of the modes after training on MNIST.

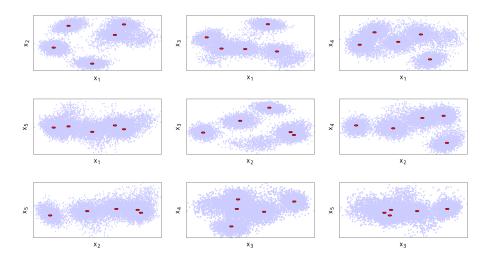


Fig. 8: Scatter plots of samples (blue) and fixed points (red) in various plan projections defined by pairs of left eigenvectors of W. The dataset is the same as in Fig. 5 and in this case 5 modes have condensed and 7 fixed point solutions have been found.

Finally, the training set is divided into batches of size 20, 5 Gibbs sampling steps are used (CD5) and the learning rate  $\gamma$  is kept low in order to reduce noise,  $\gamma = 5 \times 10^{-8}$ . The results of the analysis are shown in fig. 6. We see that the dynamics of the singular values obtained by direct integration of the mean-field equations (Fig. 4) are very well reproduced, the only difference being a slightly higher pressure on the strongest modes. The number of fixed point solutions also seems to follow the same trend but more noise is present, an indication of the fact that the RBM has a tendency to learn spurious fixed points during the training. The learning trajectory on the phase diagram is also of interest; we see that the RBM is initialized in the paramagnetic state as expected and the effect of the learning is to drive the model to the ferromagnetic phase. Once in the ferromagnetic phase, the trajectory slows down and the model is assessed near the critical line between paramagnetic and ferromagnetic states, where the estimate of the weights is most stable (according to [35]). Finally, in Fig. 8 we see how the RBM is able to generate a proper clustering of the data over the spectral modes. In particular, the TAP fixed points of the trained model are well distributed and able to cover the full data distribution, improving over the typical behaviour for Laplace distributed weights that emerged with our theoretical analysis (Fig. 5).

#### 5.2 MNIST dataset

The MNIST dataset is composed by 70000 handwritten digits (60000 for training, 10000 for testing) of size  $28 \times 28$  pixels. Being highly multimodal, we expect this dataset to push the limits of our spectral analysis. For the training, the initialization of the model is the same one used for the synthetic data, 10000 training samples are used (taken at random from the dataset) and the values of the other hyperparameters

6 Discussion 29

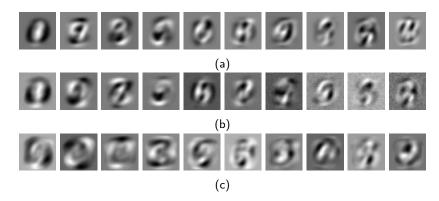


Fig. 9: (a) Principal components extracted from the training set (starting from the second, as the first one is encoded into the visible bias). (b) The first 10 modes of a RBM trained for 1 epoch (with  $\gamma \simeq 0.1$ ). (c) Same as (b) but after a 10 epochs training.

are as follows:  $N_v = 784$ ,  $N_h = 100$ , batch size = 20,  $\gamma = 5 \times 10^{-7}$ . With respect to the linear regime (described in section 4.2) we see in Fig. 9 how the RBM is able to learn the SVD of the dataset quite precisely at the beginning of the training, then the learning dynamics quickly enter the non-linear regime. Even in this highly multimodal scenario, our findings over simple synthetic data seem to be confirmed, as seen in Fig. 6. The high number of modes, however, determines an increase in the magnitude of the singular values of condensed modes and seems to destabilize a bit the learning, making the computation of fixed points less reliable. In fact, as a high number of modes are condensing, the model is not able to get rid of all the spurious fixed points. This problem can be mitigated by using an even smaller learning rate, at the cost of slowing down the training. Probably, using a variable learning rate could be a more practical solution (descreasing the learning rate from time to time to let the model eliminate unneeded fixed points). Concerning the (relative) kurtosis of the mode components distributions, we did not observe a very stable and systematic behavior. Either we see small fluctuations around zero, either some excursions occur and a finite value in the range [0,3] is building up either for the u or the v components, coherently to the compositional phase interpretation given previously. The latter is the case for MNIST, as shown in Fig. 7. Additionally the transverse part of the fields, meaning orthogonal to the condensed modes, is usually not completely negligible, in contrary to what we assume in (13,14). This clearly constitutes a limitation of our analysis. These transverse components offer more flexibility for generating and selecting fixed points and interfere in some non-trivial way with the kurtosis property, which possibly explains why we don't get a systematic behavior.

#### 6 Discussion

Before drawing some perspectives, let us summarize the main outcomes of the present work:

• (i) thermodynamic properties of realistic RBMs: our analysis focused on a non-i.i.d. ensemble of weight matrices, whose derivation has been inspired

6 Discussion 30

by empirical observations obtained by training RBMs on real data.

• (ii) RS equations and compositional phase: we found a way of writing the RS equations for the RBM (in particular with equations (21,22,23,24)) which leads to a simple characterization of the ferromagnetic phase where the RBM is assumed to operate. Schematically, a negative relative kurtosis for the distribution of the singular vectors' components favors the proliferation of metastable states, while a positive one tends to favor a compositional phase. In particular, we were able to precisely address a concrete case presenting the compositional phase by considering a Laplace distribution for the singular vectors' components.

- (iii) a set of equations representing a typical learning dynamics that defines a trajectory in  $\{w_{\alpha}(t), \eta_{\alpha}(t), \theta_{\alpha}(t), \Omega_{\alpha\beta}^{v,h}(t), \sigma^2(t)\}$ . The spectrum of the dominant singular values, represented by  $\{w_{\alpha}(t)\}$  and expressing the information content of the RBM, is playing the main role. The bulk of dominated modes corresponding to noise sees its dynamics summarized by the evolution of  $\sigma^2(t)$ . Rotations of dominant singular vectors during the learning process are given by  $\Omega^{v,h}$  while the projections of the biases along the main modes are given by  $\eta$  and  $\theta$ . These equations have been obtained by averaging over the components of left and right SVD vectors of the weight matrix, while keeping fixed the quantities considered to be relevant. This averaging actually corresponds to a standard self-averaging assumption in a RS phase.
- (iv) a clustering interpretation of the training process is obtained through equations (42,43,44,45) where it is explicitly shown the kind of matching that the RBM is trying to perform between the order parameters obtained from the fixed point solutions and their empirical counterparts in the non-linear regime. A natural clustering of the data can actually be defined by assigning to each sample the fixed point obtained after initializing the fixed point equations with a visible configuration corresponding to that same sample.

The main picture emerging from the present analysis is that of a set of clusters corresponding to the fixed points of the RBM, which try to uniformly cover the support of the dataset. A full understanding of the mechanism by which the RBM manages to properly cover the dataset is still lacking, even though the case of Laplace distributed singular vectors' components gives some insights. By comparison, real RBMs have more flexibility than the simple "mean Laplace RBM" considered in Section 3.4 and they can produce a good covering of the data manifold. We were not yet able to precisely pinpoint the main ingredients for that mechanism, even though we suspect the transverse biases (orthogonal to the modes) of the hidden units to be the missing ingredient in our analysis.

From the theoretical point of view we would like to see how these results can be adapted to more complex models like DBM or generative models based on convolutional networks. In particular we would like to understand whether adding more layers can facilitate the covering of the dataset by fixed points. From the practical point of view these results might help to orientate the choice of the hyper-parameters used for training an RBM and to refine the criteria for assessing the quality of a learned RBM. For instance, the choice of the number of hidden variables is dictated by two considerations: the effective rank of W, i.e. the number of relevant modes to be considered, and the level of interaction between these modes. Using less hidden variables gives more compact RBMs and reduces the rank of W to its needed value, but it also leads to modes with stronger interactions, which means less flexibility for generating a good covering of fixed points.

A AT line 31

## A AT line

The stability of the RS solution to the mean-field equations is studied along the lines of [33] by looking at the Hessian of the replicated version of the free energy and identifying eigenmodes from symmetry arguments. Before taking the limit  $p \to 0$  the free energy reads

$$f[m, \bar{m}, Q, \bar{Q}] = \sum_{a,\alpha} w_{\alpha} m_{\alpha}^{a} \bar{m}_{\alpha}^{a} + \frac{\sigma^{2}}{2} \sum_{a \neq b} Q_{ab} \bar{Q}_{ab} - \frac{1}{\sqrt{\kappa}} A_{p}[m, Q] - \sqrt{\kappa} B_{p}[\bar{m}, \bar{Q}],$$

with  $A_p$  and  $B_p$  given in (10,11). Assuming the small perturbations

$$m_{\alpha}^{a}=m_{\alpha}+\epsilon_{\alpha}^{a}$$
  $\bar{m}_{\alpha}^{a}=\bar{m}_{\alpha}+\bar{\epsilon}_{\alpha}^{a}$  
$$Q_{ab}=q+\eta_{ab}$$
  $\bar{Q}_{ab}=\bar{q}+\bar{\eta}_{ab},$ 

around the saddle point  $(m_{\alpha}, \bar{m}_{\alpha}, q, \bar{q})$ , the perturbed free energy reads

$$\Delta f = \sum_{a,\alpha} w_{\alpha} \bar{\epsilon}_{\alpha}^{a} \epsilon_{\alpha}^{a} + \frac{\sigma^{2}}{2} \sum_{a \neq b} \bar{\eta}_{ab} \eta_{ab} + \sum_{a,b,\alpha,\beta} \left[ \left( \delta_{ab} \bar{A}_{\alpha\beta} + \bar{\delta}_{ab} \bar{B}_{\alpha\beta} \right) \epsilon_{\alpha}^{a} \epsilon_{\beta}^{b} + CT \right]$$

$$+ \sum_{a \neq b,c,\alpha} \left[ \left( \left( \delta_{ab} + \delta_{ac} \right) \bar{C}_{\alpha} + \left( 1 - \delta_{ac} - \delta_{bc} \right) \bar{D}_{\alpha} \right) \epsilon_{\alpha}^{c} \eta_{ab} + CT \right]$$

$$+ \sum_{a \neq b,c,\alpha} \left[ \left( \delta_{(ab)(cd)} \bar{E}_{0} + \mathbb{1}_{\{a \in (cd) \oplus b \in (cd)\}} \bar{E}_{1} + \mathbb{1}_{\{(ab) \cap (cd) = \emptyset\}} \bar{E}_{2} \right) \eta_{ab} \eta_{cd} + CT \right],$$

where CT means "conjugate term" in the sense  $\epsilon \leftrightarrow \bar{\epsilon}$ ,  $A_{\alpha\beta} \leftrightarrow \bar{A}_{\alpha\beta}...$ , where  $\bar{b}_{ab} \stackrel{\text{def}}{=} 1 - \delta_{ab}$  and the operators are given by

$$A_{\alpha\beta} \stackrel{\text{def}}{=} (\delta_{\alpha\beta} - m_{\alpha} m_{\beta}) w_{\alpha} w_{\beta} \qquad B_{\alpha\beta} \stackrel{\text{def}}{=} \left( \mathsf{E}_{x,v} \left( v^{\alpha} v^{\beta} \tanh^{2}(\bar{h}(x,v)) \right) - m_{\alpha} m_{\beta} \right) w_{\alpha} w_{\beta}$$

$$C_{\alpha} \stackrel{\text{def}}{=} \frac{\kappa^{1/4} \sigma^{2}}{2} m_{\alpha} (1 - q) w_{\alpha} \qquad D_{\alpha} \stackrel{\text{def}}{=} \frac{\kappa^{1/4} \sigma^{2}}{2} \left( \mathsf{E}_{x,v} \left( v^{\alpha} \tanh^{3}(\bar{h}(x,v)) \right) - m_{\alpha} q \right) w_{\alpha}$$

$$E_{0} \stackrel{\text{def}}{=} \frac{\sqrt{\kappa} \sigma^{4}}{4} (1 - q^{2}) \qquad E_{1} \stackrel{\text{def}}{=} \frac{\sqrt{\kappa} \sigma^{4}}{4} q (1 - q) \qquad E_{2} \stackrel{\text{def}}{=} \frac{\sqrt{\kappa} \sigma^{4}}{4} \left( \mathsf{E}_{x,v} \left( \tanh^{4}(\bar{h}(x,v)) \right) - q^{2} \right)$$
with
$$h(x,u) \stackrel{\text{def}}{=} \kappa^{1/4} \left( \sqrt{q} \sigma x + \sum_{\alpha} (m_{\alpha} w_{\alpha} - \eta_{\alpha}) u^{\alpha} \right),$$

Conjugate quantities are obtained by replacing  $m_{\alpha}$  by  $\bar{m}_{\alpha}$ , q by  $\bar{q}$ ,  $u^{\alpha}$  by  $v^{\alpha}$ ,  $\eta_{\alpha}$  by  $\theta_{\alpha}$  and  $\kappa$  by  $1/\kappa$ . As for the SK model, the  $2Kp \times 2Kp$  Hessian thereby defined can be diagonalized with the help of three similar sets of eigenmodes corresponding to different permutation symmetries in replica space.

The first set corresponds to 2K+2 replica symmetric modes defined by  $\eta_{\alpha}^{a}=\eta_{\alpha}$ 

A AT line 32

and  $\eta_{ab} = \eta$  solving the linear system

$$\left(\frac{w_{\alpha}}{2} - \lambda\right)\bar{\epsilon}_{\alpha} - \frac{1}{2}\bar{A}_{\alpha\alpha}\epsilon_{\alpha} + \sum_{\beta} \left(\bar{A}_{\alpha\beta} + (p-1)\bar{B}_{\alpha\beta}\right)\epsilon_{\beta} + \left((p-1)\bar{C}_{\alpha} + \frac{(p-1)(p-2)}{2}\bar{D}_{\alpha}\right)\eta = 0$$

$$\left(\frac{w_{\alpha}}{2} - \lambda\right)\epsilon_{\alpha} - \frac{1}{2}A_{\alpha\alpha}\bar{\epsilon}_{\alpha} + \sum_{\beta} \left(A_{\alpha\beta} + (p-1)B_{\alpha\beta}\right)\bar{\epsilon}_{\beta} + \left((p-1)C_{\alpha} + \frac{(p-1)(p-2)}{2}D_{\alpha}\right)\bar{\eta} = 0$$

$$\left(\frac{\sigma^{2}}{2} - \lambda\right)\bar{\eta} + \sum_{\alpha} \left(\bar{C}_{\alpha} + \frac{p-2}{2}\bar{D}_{\alpha}\right)\epsilon_{\alpha} + 2\left(\bar{E}_{0} + 2(p-2)\bar{E}_{1} + \frac{(p-2)(p-3)}{2}\bar{E}_{2}\right)\eta = 0$$

$$\left(\frac{\sigma^{2}}{2} - \lambda\right)\eta + \sum_{\alpha} \left(C_{\alpha} + \frac{p-2}{2}D_{\alpha}\right)\bar{\epsilon}_{\alpha} + 2\left(E_{0} + 2(p-2)E_{1} + \frac{(p-2)(p-3)}{2}E_{2}\right)\bar{\eta} = 0$$

with eigenvalue  $\lambda$  solving a polynomial equation of degree 2K+2 corresponding to a vanishing determinant in the above system.

The second set corresponds to a broken replica symmetry where one replica  $a_0$  is different from the others

$$(\epsilon_{\alpha}^{a}, \bar{\epsilon}_{\alpha}^{a}) = \begin{cases} (\epsilon_{\alpha}, \bar{\epsilon}_{\alpha}) & \text{for } a \neq a_{0} \\ (1-p)(\epsilon_{\alpha}, \bar{\epsilon}_{\alpha}) & \text{for } a = a_{0} \end{cases} \qquad (\eta_{ab}, \bar{\eta}_{ab}) = \begin{cases} (\eta, \bar{\eta}) & \text{for } a, b \neq a_{0} \\ (1-\frac{p}{2})(\eta, \bar{\eta}) & \text{for } a = a_{0} \text{ or } b = a_{0} \end{cases}$$

This set has dimension (2K+2)(p-1). Its parameterization is obtained by imposing orthogonality with the previous one. The corresponding system reads

$$\left(\frac{w_{\alpha}}{2} - \lambda\right)\bar{\epsilon}_{\alpha} - \frac{1}{2}\bar{A}_{\alpha\alpha}\epsilon_{\alpha} + \sum_{\beta}(\bar{A}_{\alpha\beta} - \bar{B}_{\alpha\beta})\epsilon_{\beta} + \frac{p-2}{2}(\bar{C}_{\alpha} - \bar{D}_{\alpha})\eta = 0$$

$$\left(\frac{w_{\alpha}}{2} - \lambda\right)\epsilon_{\alpha} - \frac{1}{2}A_{\alpha\alpha}\bar{\epsilon}_{\alpha} + \sum_{\beta}(A_{\alpha\beta} - B_{\alpha\beta})\bar{\epsilon}_{\beta} + \frac{p-2}{2}(C_{\alpha} - D_{\alpha})\bar{\eta} = 0$$

$$\left(\frac{\sigma^{2}}{2} - \lambda\right)\bar{\eta} + \sum_{\alpha}(\bar{C}_{\alpha} - \bar{D}_{\alpha})\epsilon_{\alpha} + 2(\bar{E}_{0} + (p-4)\bar{E}_{1} - (p-3)\bar{E}_{2})\eta = 0$$

$$\left(\frac{\sigma^{2}}{2} - \lambda\right)\eta + \sum_{\alpha}(C_{\alpha} - D_{\alpha})\bar{\epsilon}_{\alpha} + 2(E_{0} + (p-4)E_{1} - (p-3)E_{2})\bar{\eta} = 0$$

Finally the eigenmodes of the Hessian are made complete by considering a broken symmetry where two replicas  $a_0$  and  $a_1$  are different from the others, with the following parameterization dictated again by orthogonality constraints with the previous sets:

$$(\epsilon_{\alpha}^{a}, \bar{\epsilon}_{\alpha}^{a}) = 0, \qquad (\eta_{ab}, \bar{\eta}_{ab}) = \begin{cases} (\eta, \bar{\eta}) & \text{for } a, b \neq a_{0} \\ \frac{3-p}{2}(\eta, \bar{\eta}) & \text{for } a \in a_{0}, a_{1} \text{ or } b \in a_{0}, a_{1} \\ \frac{(p-2)(p-3)}{2}(\eta, \bar{\eta}) & \text{for } (a, b) = (a_{0}, a_{1}). \end{cases}$$

The dimension of this set is now p(p-3), and it represents eigenvectors iff the following

system of equations is satisfied

$$\left(\frac{\sigma^2}{2} - \lambda\right)\bar{\eta} + 2(\bar{E}_0 - 2\bar{E}_1 + \bar{E}_2)\eta = 0$$

$$\left(\frac{\sigma^2}{2} - \lambda\right)\eta + 2(E_0 - 2E_1 + E_2)\bar{\eta} = 0$$

The corresponding eigenvalues read

$$\lambda = \frac{\sigma^2}{2} \pm 2\sqrt{(\bar{E}_0 - 2\bar{E}_1 + \bar{E}_2)(E_0 - 2E_1 + E_2)},$$

with degeneracy p(p-3)/2. Finally the RS stability condition reads

$$\frac{1}{\sigma^2} > \sqrt{\mathsf{E}_{x,u} \Big( \mathrm{sech}^4 \Big( h(x,u) \Big) \Big) \mathsf{E}_{x,v} \Big( \mathrm{sech}^4 \Big( \bar{h}(x,v) \Big) \Big)},$$

which reduces to the same form of the AT line for the SK model when  $\kappa=1$ , except for the u and v averages that are specific to our model. As seen in Figure 2 the influence of  $\kappa$  is very limited.

# B Synthetic dataset

The multimodal distribution modeling the N-dimensional synthetic data is

$$P(s) = \sum_{c=1}^{C} p_c \prod_{i=1}^{N} \frac{e^{h_i^c s_i}}{2 \cosh(h_i^c)},$$
(47)

where C is the number of clusters,  $p_c$  is a weight and  $\mathbf{h}^c$  is a hidden field for cluster c. The values for  $p_c$  are taken at random and normalized, while to compute  $h_i^c$  we take into account the magnetizations  $m_i^c = \tanh(h_i^c)$ . Expanding over the spectral modes, we can set an effective dimension d by constraining the sum to the range  $\alpha = 1, \ldots, d$ 

$$m_i^c = \sum_{\alpha=1}^d m_\alpha^c u_i^\alpha \tag{48}$$

Clusters' magnetizations  $m_{\alpha}^{c}$  are drawn at random between [-1,1] and normalized with the factor

$$Z = \sqrt{\frac{\sum_{\alpha} m_{\alpha}^{2}}{d \cdot r}}, \quad r = \tanh(\eta)$$
 (49)

where r is introduced to decrease the clusters' polarizations (in our simulations, we used  $\eta=0.3$ ). The spectral basis  $u_i^{\alpha}$  is obtained by drawing at random d N-dimensional vectors and applying the Gram-Schmidt process (which can be safely employed as N is supposedly big and thus the initial vectors are nearly orthogonal). The hidden fields are then obtained from the magnetizations

$$h_i^c = \tanh^{-1}(m_i^c) \tag{50}$$

and the samples are generated by choosing a cluster according to  $p_c$  and setting the visible variables to  $\pm 1$  according to

$$p(s_i = 1) = \frac{1}{1 + e^{-2h_i^c}} \tag{51}$$

#### References

- [1] P. Smolensky. In Parallel Distributed Processing: Volume 1 by D. Rumelhart and J. McLelland, chapter 6: Information Processing in Dynamical Systems: Foundations of Harmony Theory. 194-281. MIT Press, 1986.
- [2] R. Salakhutdinov and G. Hinton. Deep Boltzmann machines. In Artificial Intelligence and Statistics, pages 448–455, 2009.
- [3] G.E. Hinton and R.R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- [4] G. E. Hinton. Training products of experts by minimizing contrastive divergence. Neural computation, 14:1771–1800, 2002.
- [5] T. Tieleman. Training restricted Boltzmann machines using approximations to the likelihood gradient. In *Proceedings of the 25th International Conference on Machine Learning*, ICML '08, pages 1064–1071, New York, NY, USA, 2008. ACM.
- [6] G.E. Hinton. A Practical Guide to Training Restricted Boltzmann Machines, pages 599–619. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012.
- [7] D.S.P. Salazar. Nonequilibrium thermodynamics of restricted Boltzmann machines. *Phys. Rev. E*, 96:022131, 2017.
- [8] 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(8):2554–2558, 1982.
- [9] D. J. Amit, H. Gutfreund, and H. Sompolinsky. Statistical mechanics of neural networks near saturation. *Annals of Physics*, 173(1):30–67, 1987.
- [10] E. Gardner. Maximum storage capacity in neural networks. EPL (Europhysics Letters), 4(4):481, 1987.
- [11] E. Gardner and B. Derrida. Optimal storage properties of neural network models. Journal of Physics A: Mathematical and General, 21(1):271, 1988.
- [12] B. Barra, A. Bernacchia, E. Santucci, and P. Contucci. On the equivalence of Hopfield networks and Boltzmann machines. *Neural Networks*, 34:1–9, 2012.
- [13] G. Marylou, E.W. Tramel, and F. Krzakala. Training restricted Boltzmann machines via the Thouless-Anderson-Palmer free energy. In *Proceedings of the 28th International Conference on Neural Information Processing Systems*, NIPS'15, pages 640–648, 2015.
- [14] H. Huang and T. Toyoizumi. Advanced mean-field theory of the restricted Boltzmann machine. *Physical Review E*, 91(5):050101, 2015.
- [15] C. Takahashi and M. Yasuda. Mean-field inference in gaussian restricted Boltzmann machine. *Journal of the Physical Society of Japan*, 85(3):034001, 2016.
- [16] C. Furtlehner, J.-M. Lasgouttes, and A. Auger. Learning multiple belief propagation fixed points for real time inference. *Physica A: Statistical Mechanics and its Applications*, 389(1):149–163, 2010.
- [17] A. Barra, G. Genovese, P. Sollich, and D. Tantari. Phase diagram of restricted Boltzmann machines and generalized Hopfield networks with arbitrary priors. arXiv:1702.05882, 2017.

- [18] H. Huang. Statistical mechanics of unsupervised feature learning in a restricted Boltzmann machine with binary synapses. *Journal of Statistical Mechanics: Theory and Experiment*, 2017(5):053302, 2017.
- [19] E. Agliari, A. Barra, A. Galluzzi, F. Guerra, and F. Moauro. Multitasking associative networks. Phys. Rev. Lett., 109:268101, 2012.
- [20] R. Monasson and J. Tubiana. Emergence of compositional representations in restricted Boltzmann machines. Phys. Rev. Let., 118:138301, 2017.
- [21] L. Zdeborová and F. Krzakala. Statistical physics of inference: thresholds and algorithms. Advances in Physics, 65(5):453–552, 2016.
- [22] M. E. Tipping and C. M. Bishop. Mixtures of probabilistic principal component analyzers. *Neural Comput.*, 11(2):443–482, 1999.
- [23] H. Bourlard and Y. Kamp. Auto-association by multilayer perceptrons and singular value decomposition. *Biological Cybernetics*, 59(4):291–294, 1988.
- [24] A. M. Saxe, J. L. McClelland, and S. Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. arXiv:1312.6120, 2014.
- [25] A. Decelle, G. Fissore, and C. Furtlehner. Spectral dynamics of learning in restricted Boltzmann machines. EPL, 119(6):60001, 2017.
- [26] E.W. Tramel, M. Gabrié, A. Manoel, F. Caltagirone, and F. Krzakala. A Deterministic and Generalized Framework for Unsupervised Learning with Restricted Boltzmann Machines. arXiv:1702.03260, 2017.
- [27] V. A. Marčenko and L. A. Pastur. Distribution of eigenvalues for some sets of random matrices. *Mathematics of the USSR-Sbornik*, 1(4):457, 1967.
- [28] M. Mézard. Mean-field message-passing equations in the Hopfield model and its generalizations. *Phys. Rev. E*, 95:022117, 2017.
- [29] G. Parisi and M. Potters. Mean-field equations for spin models with orthogonal interaction matrices. *Journal of Physics A: Mathematical and General*, 28(18):5267, 1995.
- [30] M. Opper and O. Winther. Adaptive and self-averaging Thouless-Anderson-Palmer mean field theory for probabilistic modeling. *Physical Review E*, 64:056131, 2001.
- [31] D. J. Amit, H. Gutfreund, and H. Sompolinsky. Spin-glass models of neural networks. Phys. Rev. A, 32:1007–1018, 1985.
- [32] M. Mézard, G. Parisi, and M. A. Virasoro. Spin Glass Theory and Beyond. World Scientific, Singapore, 1987.
- [33] J. R. L. Almeida and D. J. Thouless. Stability of the Sherrington-Kirkpatrick solution of a spin glass model. J. Phys. A:Math. Gen., 11(5):983-990, 1978.
- [34] P. C. Hohenberg and M. C. Cross. An introduction to pattern formation in nonequilibrium systems, pages 55–92. Springer Berlin Heidelberg, Berlin, Heidelberg, 1987.
- [35] I. Mastromatteo and M. Marsili. On the criticality of inferred models. *Journal of Statistical Mechanics: Theory and Experiment*, 2011(10):P10012, 2011.