

The Decimation Scheme for Symmetric Matrix Factorization

Francesco Camilli* and Marc Mézard†

**The Abdus Salam International Center for Theoretical Physics, Trieste, Italy*

†*Bocconi University, Milan, Italy*

August 1, 2023

Abstract

Matrix factorization is an inference problem that has acquired importance due to its vast range of applications that go from dictionary learning to recommendation systems and machine learning with deep networks. The study of its fundamental statistical limits represents a true challenge, and despite a decade-long history of efforts in the community, there is still no closed formula able to describe its optimal performances in the case where the rank of the matrix scales linearly with its size. In the present paper, we study this extensive rank problem, extending the alternative ‘decimation’ procedure that we recently introduced, and carry out a thorough study of its performance. Decimation aims at recovering one column/line of the factors at a time, by mapping the problem into a sequence of neural network models of associative memory at a tunable temperature. Though being sub-optimal, decimation has the advantage of being theoretically analyzable. We extend its scope and analysis to two families of matrices. For a large class of compactly supported priors, we show that the replica symmetric free entropy of the neural network models takes a universal form in the low temperature limit. For sparse Ising prior, we show that the storage capacity of the neural network models diverges as sparsity in the patterns increases, and we introduce a simple algorithm based on a ground state search that implements decimation and performs matrix factorization, with no need of an informative initialization.

Contents

1	Introduction	2
2	Decimation	4
2.1	An assumption on retrieval accuracy	7
3	Decimation free entropies	7
3.1	Fixed point equations	11
3.2	Remarks	11
4	Low temperature limits	13
4.1	Sparse prior	13
4.2	Continuous priors	15
5	Phase diagrams for the first decimation step	17

*fcamilli@ictp.it, †marc.mezard@unibocconi.it

6	Numerical tests	19
6.1	Testing the saddle point equations with AMP	19
6.2	Expected decimation performance	21
6.3	A ground state oracle for sparse Ising priors	23
6.4	Reversed decimation	24
7	Related works	25
7.1	Unlearning and dreaming	25
7.2	Sub-linear rank	26
7.3	Channel universality properties	26
8	Conclusion and outlooks	27

1 Introduction

The factorization of a matrix into two, or more, factors represents a building block for many machine learning and inference problems. A well-known instance of it is *dictionary learning* [1–4], which aims at representing a matrix as a product of two factor matrices, where the first, called *dictionary*, is very sparse, and the second, called *feature matrix*, has columns that form an over-complete basis of a euclidean space. As a result, each vector stored in the initial matrix is represented as a linear combination of few elements of the feature matrix. Matrix factorization is also at the basis of recommendation systems [5], and in general proves to be very effective whenever we want to reconstruct missing elements in a matrix of data, be it an image, a correlation matrix, or a matrix of preferences [6–8]. Other applications of matrix factorization include, but are not limited to, sparse principal component analysis [9], blind source separation [10], matrix completion [11, 12], robust principal component analysis [13]

In more specific terms, matrix factorization is the problem of reconstructing the two factors \mathbf{A} , \mathbf{B} of a matrix \mathbf{AB} from a potentially noisy observation of the latter, say \mathbf{Y} . One would like to answer two main questions: (i) in what regimes of sizes of \mathbf{A} , \mathbf{B} and noise is it possible to reconstruct the two factors (up to a permutation of the lines of \mathbf{A} and the columns of \mathbf{B})? (ii) Do there exist efficient algorithms that achieve a good performance?

In the present paper we focus on symmetric matrix factorization in which the two factors to retrieve are identical. Consider an $N \times P$ matrix $(\xi_i^\mu)_{\substack{\mu \leq P \\ i \leq N}} = \boldsymbol{\xi} \in \mathbb{R}^{N \times P}$ whose elements are independently and identically distributed according to a given prior probability P_ξ , that we suppose to be symmetric, with unit variance and compact support: $\mathbb{E}\xi = 0$, $\mathbb{E}\xi^2 = 1$, $|\xi| \leq C$ for some $C > 0$. Secondly, let $(Z_{ij})_{i,j \leq N} = (Z_{ji})_{i,j \leq N} = \mathbf{Z}$ be a Wigner matrix, that is $Z_{ij} = Z_{ji} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1 + \delta_{ij})$. Symmetric matrix factorization can thus be formulated as an inference problem: a Statistician needs to recover $\boldsymbol{\xi}$ given the noisy observations

$$\mathbf{Y} = \frac{\boldsymbol{\xi}\boldsymbol{\xi}^\top}{\sqrt{N}} + \sqrt{\Delta}\mathbf{Z}. \tag{1}$$

The strength of the noise \mathbf{Z} w.r.t. that of the signal is tuned by $\Delta \geq 0$. In the following we will need to single out the P column vectors inside $\boldsymbol{\xi}$, denoted by $\boldsymbol{\xi}^\mu$, and we shall refer to them as *patterns*. Despite the model is presented here in a stylized way, i.e. with the two factors being identical and with completely factorized prior, we believe this setting represents a fundamental first step in the understanding of the general problem. Concerning in particular the assumption of a factorized prior, this is often used also in concrete situations. Indeed, for instance, the L^2 norm regulators appearing in the empirical risk used to train neural networks are inherited from a zero temperature limit of a Statistical Mechanics problem that has the empirical risk as a Hamiltonian with factorized prior on the weights of the network, as clarified by [14].

A very popular setting to tackle an inference problem is the Bayes-optimal one, in which the Statistician tasked with the reconstruction of $\boldsymbol{\xi}$ knows the generating process of the observations \mathbf{Y} , namely they know that \mathbf{Z} is Gaussian, they know N, P, Δ and the probability distribution of factors P_ξ . This Bayes-optimal setting is of utmost relevance as it provides the information-theoretic optimal performance. Indeed, the posterior mean

estimator $\mathbb{E}[\mathbf{X}\mathbf{X}^\top|\mathbf{Y}]$, where

$$dP(\boldsymbol{\xi} = \mathbf{X} | \mathbf{Y}) = \frac{1}{\mathcal{Z}(\mathbf{Y})} \prod_{i \leq N, \mu \leq P} dP_\xi(X_i^\mu) \exp \left[\frac{1}{2\sqrt{N}\Delta} \text{Tr} \mathbf{Y} \mathbf{X} \mathbf{X}^\top - \frac{1}{4\Delta N} \text{Tr}(\mathbf{X} \mathbf{X}^\top)^2 \right], \quad (2)$$

is the one that minimizes the mean square error loss on the reconstruction of $\boldsymbol{\xi}\boldsymbol{\xi}^\top$. The normalization of the distribution $\mathcal{Z}(\mathbf{Y})$ is called *partition function* and the associated *free entropy* is defined as

$$\Phi_{N,P} = \frac{1}{NP} \mathbb{E} \log \mathcal{Z}(\mathbf{Y}). \quad (3)$$

The free entropy has a central role. In fact, from the thermodynamic point of view, it can be used to identify what macrostates dominate probability and are thus selected at thermodynamic equilibrium. These macrostates are usually identified by the values of some global order parameters, such as $\text{Tr} \mathbf{X} \mathbf{X}^\top \boldsymbol{\xi} \boldsymbol{\xi}^\top / N^2$, which measures the average alignment of a sample from the posterior and the ground truth $\boldsymbol{\xi}$ we want to estimate. On the other hand, the free entropy is in close relationship with the *mutual information* $I(\boldsymbol{\xi}; \mathbf{Y})$ between the data and the ground truth. This information theoretic quantity quantifies the amount of residual information about the ground truth that is still available in the data after they have been corrupted by the noise.

If the rank P is finite, the model (1) is typically referred to as *spiked Wigner model*, first introduced as model for Principal Component Analysis (PCA) [15]. The spectral properties of low rank perturbations of high-rank matrices (such as the Wigner matrix \mathbf{Z}) are by now largely understood in random matrix theory, and they can give rise to the celebrated BBP transition [16], further studied and extended in [17–24]. Thanks to the effort of a wide interdisciplinary community, we also have a control on the asymptotic behaviour of the posterior measure (2) and an exact formula for the free entropy associated to the low-rank problem [25–32] (recently extended to rotational invariant noise [33]), which yields the Bayes-optimal limit of the noise allowing the reconstruction of the low-rank spike. Finally, a particular class of algorithms, known as *Approximate Message Passing* (AMP) [34–38], is able to perform factorization up to this Bayes-optimal limit.

Here we are interested in the extensive rank regime where $P, N \rightarrow \infty$ with fixed ratio $P/N = \alpha$. In the hypothesis of a rotationally invariant noise \mathbf{Z} , the spectral properties of \mathbf{Y} are governed by the free-convolution [39] of the spectral densities of \mathbf{Z} and $\boldsymbol{\xi}\boldsymbol{\xi}^\top$. On the information theoretic side instead, there still is no accepted closed formula that expresses $\Phi_{N,P}$. Hence, the information theoretic limits are currently out of reach, and the Minimum Mean Square Error (MMSE) for this estimation problem is not known. Among the past attempts, we must mention the line of works [40–44], whose proposed solution, as pointed out in [45, 46], provides only an approximation of the correct limit. In fact, the authors of [46] build a perturbative approach that highlights the presence of relevant correlations neglected in the previous works. A further attempt to produce a closed replica formula was put forward in [47], but, as [40], it involves uncontrolled approximations.

The main obstacle in the computation of the asymptotics of (3) is the fact that it is a matrix model, and, in particular, the term $\text{Tr}(\mathbf{X}\mathbf{X}^\top)^2$ couples both the “rank, or patterns indices” μ , and the “dimension, or particle site indices” i . We will use here a different approach that we introduced and studied recently [48] in the simplest case where the factors’ elements ξ_i^μ are independent binary variables. Instead of the Bayes-optimal setting we use a simpler procedure, that we call *decimation*. At the cost of giving up on Bayes-optimality, decimation solves this problem and allows us to identify an iterative scheme to estimate pattern by pattern, giving an estimate of $\boldsymbol{\xi}$ through a sequential estimation of its columns, and, more importantly, whose asymptotic performance turns out to be completely analyzable. In the case of binary patterns we could thus show that matrix factorization is possible in a part of the phase diagram where α and Δ are small enough. Here we generalize this approach to arbitrary distributions of the patterns’ elements.

Organization of the paper and main contributions In Section 2 we define the decimation scheme, laying the ground for the replica computation of Section 3. In Section 4, we compute the low temperature limits for two classes of priors: sparse Ising and a generic absolutely continuous, symmetric and bounded support prior.

Surprisingly, the free entropies of the neural network models arising from decimation evaluated at the equilibrium value of the order parameters have a universal form, but in general not the same numerical value.

As we shall argue in the following, the starting point of the decimation procedure, i.e. the initial value of the parameters α and Δ , is of crucial importance for its success. Therefore, in Section 5 we analyze the phase diagrams for the initial step of decimation. For the sparse Ising prior, we show that as sparsity increases, the storage capacity of the sequential neural network models of decimation diverges. For the class of continuous priors we highlight the presence of a thermodynamic transition, where there is a non-trivial overlap between a sample from the Gibbs measure and the sought pattern, and a performance transition, where Gibbs sampling can outperform the null-estimator.

In Section 6 we provide numerical evidence in support of the replica theory. We introduce the Decimated AMP algorithm (DAMP), in order to verify the predictions of the replica theory, and we relate the replica symmetric order parameters to the mean square error on the reconstruction of the patterns, as well as to the matrix mean square error for matrix denoising, showing that decimation can outperform Rotational Invariant Estimators (RIEs) [49–51] in this task. Furthermore, this Section contains the pseudo-code of a ground state oracle, an algorithm that is indeed able to find all the patterns one by one, with no need of informative initialization, contrary to DAMP.

Section 7 contains a comparison with recent relevant works that are related to the present one. Finally, Section 8 gathers the conclusions and future perspectives.

2 Decimation

Let us give a closer look at the probability distribution (2). For the purpose of the theoretical analysis we can replace Y_{ij} with the r.h.s. of (1), getting

$$dP(\boldsymbol{\xi} = \mathbf{X} \mid \mathbf{Y}) = \frac{1}{\mathcal{Z}(\mathbf{Y})} \prod_{i \leq N, \mu \leq P} [dP_{\xi}(X_i^{\mu})] e^{-\beta[\sum_{\mu}(E_1(\mathbf{X}^{\mu}) + E_2(\mathbf{X}^{\mu}) + E_3(\mathbf{X}^{\mu})) + \sum_{\mu < \nu} E_4(\mathbf{X}^{\mu}, \mathbf{X}^{\nu})]} \quad (4)$$

where $\beta = \frac{1}{\Delta}$, $\mathbf{X}^{\mu} = (X_i^{\mu})_{i \leq N}$ and

$$E_1(\mathbf{x}) = - \sum_{i,j=1}^N J_{ij} x_i x_j \quad ; \quad J_{ij} = \frac{1}{N} \sum_{\nu} \xi_i^{\nu} \xi_j^{\nu} \quad (5)$$

$$E_2(\mathbf{x}) = - \sum_{i,j=1}^N \frac{\sqrt{\Delta}}{2\sqrt{N}} Z_{ij} x_i x_j \quad (6)$$

$$E_3(\mathbf{x}) = \frac{1}{4N} \left[\sum_i x_i^2 \right]^2 \quad (7)$$

$$E_4(\mathbf{x}, \mathbf{x}') = \frac{1}{2N} \left[\sum_i x_i x'_i \right]^2. \quad (8)$$

Here one should be careful not to confuse ξ_i^{μ} which is the 'ground-truth' matrix from which the signal \mathbf{Y} was generated, and X_i^{μ} which is a random variable distributed according to the measure $dP(\boldsymbol{\xi} = \mathbf{X} \mid \mathbf{Y})$, so that the expectation value of X_i^{μ} gives the best possible approximation to ξ_i^{μ} .

Looking at the above decomposition, we notice that, if we could drop the term $E_4(\mathbf{X}^{\mu}, \mathbf{X}^{\nu})$, we would have a system of P decoupled problems, one for each value of μ , described by an energy $E_1(\mathbf{X}^{\mu}) + E_2(\mathbf{X}^{\mu}) + E_3(\mathbf{X}^{\mu})$. The energy E_1 is that of a spin glass with N variables x_i , each with an a-priori measure $P_{\xi}(x_i)$, interacting by pairs through a matrix of couplings J_{ij} which has a Hebbian form determined by the ground-truth patterns $\boldsymbol{\xi}$. The energy E_2 is a random spin glass term created by measurement noise. The energy E_3 is a global penalty that ensures that the norm of \mathbf{X} does not get too large; one can also incorporate it into the local measure using a Lagrange multiplier. Altogether, the system described by $E_1 + E_2 + E_3$ is a spin glass Hamiltonian with an

interaction which is a noisy version of a Hebbian interaction. This is typical of problems that have been studied as neural networks for associative memory, following the seminal work by Hopfield [52]. The present one is a generalization of the Hopfield model, where the stored patterns components ξ_i^μ are no longer binary but have a more general distribution which can be continuous. Based on our knowledge of associative memories, one can expect that, when the noise strength Δ and the number of patterns per variable $\alpha = P/N$ are small enough, there can exist a 'retrieval' phase, in which the configurations \mathbf{x} that minimize $E_1(\mathbf{x}) + E_2(\mathbf{x}) + E_3(\mathbf{x})$ are close to the stored patterns ξ_i^μ . This is certainly the case for binary patterns as shown in [48]. Assuming that such a retrieval phase exists, one can understand the use of the fourth energy term, E_4 . In fact one can interpret (2) as follows: we start from P replicas of an associative memory each with energy $E_1(\mathbf{X}^\mu) + E_2(\mathbf{X}^\mu) + E_3(\mathbf{X}^\mu)$. These copies interact by pairs through the term $E_4(\mathbf{X}^\mu, \mathbf{X}^\nu)$ which is a repulsive term. If one works in the retrieval phase of the associative memory, then at low temperature the ground state will be found when each replica \mathbf{X}^μ is close to one of the patterns $\xi^{\pi(\mu)}$. As there are P retrieval states and P replicas, all the $\pi(\mu)$ must be distinct from one another, and therefore π is a permutation. In such a scenario, one would have found a phase where the factors can be reconstructed.

Decimation is based precisely on this idea. It works as a sequence of P estimations, each one studying a probability distribution which is that of a neural network model of associative memory. More precisely, one looks for one column ξ^μ of ξ at a time.

To fix ideas, let us start by discussing the search of a first pattern, using a Gibbs measure in the form

$$dP(\mathbf{x} | \mathbf{Y}) = \frac{dP_\xi(\mathbf{x})}{\mathcal{Z}_0(\mathbf{Y})} \exp \left(\beta \left[\frac{1}{2N} \sum_{\mu=1}^P \left(\sum_{i=1}^N \xi_i^\mu x_i \right)^2 + \frac{\sqrt{\Delta}}{2\sqrt{N}} \sum_{i,j=1}^N Z_{ij} x_i x_j - \frac{\|\mathbf{x}\|^4}{4N} \right] \right). \quad (9)$$

Here we have introduced a factor β that plays the role of an inverse absolute temperature for this Boltzmann-Gibbs measure. We could use $\beta = 1/\Delta$ as in the Bayes-optimal approach, but as we shall see taking the large β limit can also be a good choice.

When using this approach with variables x_i that are not constrained on the hypercube $\{-1, 1\}^N$ or in general on a sphere, it is also useful to introduce another term in the exponential that favours \mathbf{x} -configurations with square norm equal to N , as we know that the original signal is centered and with unit variance. Hence, the Boltzmann-Gibbs measure that we use to find a first pattern is actually $dP_\xi(\mathbf{x})e^{-\beta E(\mathbf{x}|\mathbf{Y})}/\mathcal{Z}_0$ with an energy function

$$-E(\mathbf{x}|\mathbf{Y}) = \frac{\sqrt{\Delta}}{2\sqrt{N}} \sum_{i,j=1}^N Z_{ij} x_i x_j + \frac{N}{2} \sum_{\mu=1}^P (m^\mu(\mathbf{x}))^2 - \frac{\|\mathbf{x}\|^4}{4N} - \frac{\lambda}{4N} (\|\mathbf{x}\|^2 - N)^2 \quad (10)$$

where we have introduced the *Mattis magnetization*

$$m^\mu(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \xi_i^\mu x_i. \quad (11)$$

λ is a parameter penalizing (if positive) configurations with $\|\mathbf{x}\|^2 \neq N$, as mentioned before. If $\lambda \rightarrow +\infty$ then the spins are constrained on a sphere. Let us now assume that we are able to sample a configuration $\boldsymbol{\eta}^P$ from the Boltzmann-Gibbs measure with energy (10) that, without loss of generality (we shall relabel the patterns in such a way that the permutation π is the identity), we take as an estimate of ξ^P . How do we find the estimate of the other ξ^μ , $\mu < P$?

If $\boldsymbol{\eta}^P$ is a good estimate of ξ^P , the corresponding rank one contribution $\boldsymbol{\eta}^P \boldsymbol{\eta}^{P\top}$ should be close (in Frobenius norm) to $\xi^P \xi^{P\top}$. Then, if we subtract it from the Hebbian coupling $E_1(X)$, we can hope that the ground state of the new associative memory problem will now have only $P - 1$ ground states, each close to one of the patterns ξ^μ , $\mu = 1, \dots, P - 1$. This new associative memory problem therefore has $P - 1$ stored patterns instead of P so that the well known phenomenon of *pattern interference* [53, 54], which limits the storage capacity, will be reduced.

Based on this intuition, we define the decimation procedure as follows: after having found the first estimate of a pattern, we modify the coupling matrix as

$$\mathbf{Y}_1 = \mathbf{Y} - \frac{\boldsymbol{\eta}^P \boldsymbol{\eta}^{P\top}}{\sqrt{N}}, \quad (12)$$

which gives a modified energy function

$$-E(\mathbf{x}|\mathbf{Y}_1) = \frac{\sqrt{\Delta}}{2\sqrt{N}} \sum_{i,j=1}^N Z_{ij} x_i x_j + \frac{N}{2} \sum_{\mu=1}^P (m^\mu(\mathbf{x}))^2 - \frac{N}{2} (p^P(\mathbf{x}))^2 - \frac{\|\mathbf{x}\|^4}{4N} - \frac{\lambda}{4N} (\|\mathbf{x}\|^2 - N)^2 \quad (13)$$

where, here and in the following

$$p^\mu(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \eta_i^\mu x_i. \quad (14)$$

The same reasoning as above applies to this second step.

In general, if the first R ($= 0, 1, 2, \dots, P-1$) patterns have already been estimated, the decimation assumes to produce the estimate of the $R+1$ -th pattern sampling from the Boltzmann Gibbs measure

$$d\mu_R(\mathbf{x}) = \frac{dP_\xi(\mathbf{x})}{\mathcal{Z}_R} \exp(-\beta E(\mathbf{x}|\mathbf{Y}_R)) \quad (15)$$

where

$$\mathbf{Y}_R = \mathbf{Y} - \sum_{\mu=P-R+1}^P \frac{\boldsymbol{\eta}^\mu \boldsymbol{\eta}^{\mu\top}}{\sqrt{N}} \quad (16)$$

and

$$-E(\mathbf{x}|\mathbf{Y}_R) = \frac{\sqrt{\Delta}}{2\sqrt{N}} \sum_{i,j=1}^N Z_{ij} x_i x_j + \frac{N}{2} \sum_{\mu=1}^P (m^\mu(\mathbf{x}))^2 - \frac{N}{2} \sum_{\mu=P-R+1}^P (p^\mu(\mathbf{x}))^2 - \frac{\|\mathbf{x}\|^4}{4N} - \frac{\lambda}{4N} (\|\mathbf{x}\|^2 - N)^2. \quad (17)$$

The energy function above has some desirable features. First, the summation of the squared Mattis' magnetizations attracts mass of the distribution towards those configurations that are most aligned with one of the columns of $\boldsymbol{\xi}$, which are our goal. Secondly, if the R estimates $\boldsymbol{\eta}^\mu$, with $\mu = P-R+1, \dots, P$ are reliable, in a sense we shall specify later, the summation containing the squared $(p^\mu(\mathbf{x}))^2$ repels the mass of the probability distribution from those configurations that are similar to previously estimated patterns, preventing the sampling from finding a pattern more than once.

We notice at this point that there are three noise sources in this procedure:

- (a) the original Wigner matrix \mathbf{Z} ;
- (b) pattern interference whose strength, as discussed above, is increasing with the ratio $\alpha = P/N$;
- (c) the imperfect retrieval of patterns in the previous steps of decimation.

(c) is maybe the least obvious one. At each step, we subtract a rank one contribution $\boldsymbol{\eta}^\mu \boldsymbol{\eta}^{\mu\top} / \sqrt{N}$ that is not exactly $\boldsymbol{\xi}^\mu \boldsymbol{\xi}^{\mu\top} / \sqrt{N}$. This introduces an additional form of noise that depends on the quality of the previous reconstructions.

In order to monitor the strength of this third noise, we introduce the *retrieval accuracy* of a pattern $\boldsymbol{\xi}^\mu$:

$$m^\mu = \frac{\boldsymbol{\xi}^\mu \cdot \boldsymbol{\eta}^\mu}{N}, \quad \mu = P-R+1, \dots, P. \quad (18)$$

These quantities turn out to be order parameters of the previous decimation steps. Indeed, they are nothing but Mattis' magnetizations of typical samples from (15) with a pattern. Hence, each decimation step has its own free entropy and we will determine the new retrieval accuracy via consistency equations arising from the maximization of it, namely we look for those macrostates that dominate probability in the $N \rightarrow \infty$ limit. In addition to m^μ we will have other order parameters appearing. In particular, there will be one, denoted by r , tuning the amplitude of the overall noise, that, according to the considerations above, must comprise the three contributions coming from sources (a), (b) and (c).

2.1 An assumption on retrieval accuracy

In order to carry out the computations we need some information on the statistics of the retrieved configurations $\boldsymbol{\eta}^\mu$. We assume that an ‘‘oracle’’ algorithm will produce $\boldsymbol{\eta}^\mu$ with an asymptotic measure given by

$$\eta_i^\mu \sim \langle \cdot \rangle_{\xi_i^\mu, Z} = \frac{\int dP_\xi(x) e^{(Z\sqrt{r} + \beta m^\mu \xi_i^\mu)x - \frac{r+u}{2}x^2}(\cdot)}{\int dP_\xi(x) e^{(Z\sqrt{r} + \beta m^\mu \xi_i^\mu)x - \frac{r+u}{2}x^2}}, \quad \xi_i^\mu \sim P_\xi, Z \sim \mathcal{N}(0, 1) \text{ independent of other noises,} \quad (19)$$

where m^μ , *i.e.* the retrieval accuracy for $\boldsymbol{\eta}^\mu$, and r, u must be determined self-consistently. (19) amounts to requiring that, asymptotically, the sites are decoupled and they feel an effective external random magnetic field, that is Gaussian with a mean shifted by the ground truth ξ_i^μ . Define for later convenience the quantities

$$\mathbb{E}_{\boldsymbol{\eta}|\boldsymbol{\xi}}[\eta_i^\mu] = m_i^\mu, \quad \mathbb{E}_{\boldsymbol{\eta}|\boldsymbol{\xi}}[(\eta_i^\mu)^2] = v_i^\mu. \quad (20)$$

Then (19) has the following implications:

$$\mathbb{E}_\xi[\eta_i^\mu] = \mathbb{E}_\xi \mathbb{E}_{\boldsymbol{\eta}|\boldsymbol{\xi}}[\eta_i^\mu] = 0, \quad \mathbb{E}_\xi[\xi_i^\mu m_i^\nu] = m^\mu \delta_{\mu,\nu}, \quad \mathbb{E}_\xi[v_i^\mu] = v^\mu \quad (21)$$

that will be self-consistent with the fixed point equations for each decimation step. We shall see from the replica computation that this assumption holds inductively: if it is true at the R -th decimation step, then we are able to decouple the site indices also for the step $R + 1$, and the resulting spin-glass model has an effective random magnetic field of the same form.

3 Decimation free entropies

In this section we compute the large N limit of the free entropy

$$\Phi = \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E} \log \int dP_\xi(\mathbf{x}) \exp[-\beta E(\mathbf{x}|\mathbf{Y}_R)], \quad (22)$$

where \mathbb{E} is taken w.r.t. all the disorder: $\mathbf{Z}, \boldsymbol{\xi}, \boldsymbol{\eta}$, and recall that R is the number of patterns that were already estimated. This is done using the *replica method* [55]. We thus introduce

$$\mathbb{E} \mathcal{Z}_N^n := \mathbb{E}_{\mathbf{Z}} \mathbb{E}_{\boldsymbol{\xi}, \boldsymbol{\eta}} \int \prod_{a=1}^n dP_\xi(\mathbf{x}_a) \exp \left[-\beta \sum_{a=1}^n E(\mathbf{x}_a|\mathbf{Y}_R) \right]. \quad (23)$$

We decompose this computation and start with the first noise terms in (17), and the related $\mathbb{E}_{\mathbf{Z}}$ average

$$\begin{aligned} \mathbb{E}_{\mathbf{Z}} \exp \left(\frac{\beta\sqrt{\Delta}}{2\sqrt{N}} \sum_{i,j=1}^N Z_{ij} \sum_{a=1}^n x_{a,i} x_{a,j} \right) &= \exp \left(\frac{\beta^2 \Delta}{4N} \sum_{i,j=1}^N \sum_{a,b=1}^n x_{a,i} x_{a,j} x_{b,i} x_{b,j} \right) = \\ &= \exp \left(\frac{N\beta^2 \Delta}{4} \sum_{a \neq b}^n Q^2(\mathbf{x}_a, \mathbf{x}_b) + \beta^2 \Delta \frac{\|\mathbf{x}_a\|^4}{4N} \right). \end{aligned} \quad (24)$$

where $Q(\mathbf{x}, \mathbf{x}') = (1/N) \sum_i x_i x'_i$. For future convenience, we introduce the ‘‘decimation time’’ $t = R/P$, i.e. the fraction of patterns already estimated. Now we take care of the penalizing p -terms in (17). After replicating, their contribution to the partition function is

$$A := \prod_{\mu=P(1-t)+1}^P \prod_{a=1}^n e^{-\frac{N\beta}{2} (p^\mu(\mathbf{x}_a))^2} = \prod_{\mu=P(1-t)+1}^P \prod_{a=1}^n \int \frac{ds_a^\mu}{\sqrt{2\pi}} e^{-\frac{(s_a^\mu)^2}{2} + i\sqrt{\frac{\beta}{N}} s_a^\mu \sum_{j=1}^N \eta_j^\mu x_{a,j}}. \quad (25)$$

Notice that, thanks to the introduction of the auxiliary Gaussian variables $(s_a^\mu)_{a \leq n, P(1-t) < \mu \leq P}$, the exponential is now decoupled over the particle indices j . Consider then the expectation of A w.r.t. $\boldsymbol{\eta}$, given $\boldsymbol{\xi}$ with the assumptions (21):

$$\mathbb{E}_{\boldsymbol{\eta}|\boldsymbol{\xi}}[A] = \prod_{\mu=P(1-t)+1}^P \prod_{a=1}^n \int \frac{ds_a^\mu}{\sqrt{2\pi}} \exp\left(-\frac{(s_a^\mu)^2}{2} + \sum_{i=1}^N \log \mathbb{E}_{\eta_i^\mu|\xi_i^\mu} e^{i\sqrt{\frac{\beta}{N}} \eta_i^\mu \sum_{a=1}^n s_a^\mu x_{a,i}}\right). \quad (26)$$

Now we can expand the exponential inside the log up to second order, the remaining terms will be of sub-leading order and thus neglected in the following:

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\eta}|\boldsymbol{\xi}}[A] &= \prod_{\mu=P(1-t)+1}^P \prod_{a=1}^n \int \frac{ds_a^\mu}{\sqrt{2\pi}} \exp\left(-\frac{(s_a^\mu)^2}{2} + \sum_{a=1}^n i s_a^\mu \sqrt{\frac{\beta}{N}} \sum_{i=1}^N m_i^\mu x_{a,i} - \frac{\beta}{2} \sum_{a,b=1}^n s_a^\mu s_b^\mu \sum_{i=1}^N \frac{(v_i^\mu - (m_i^\mu)^2)}{N} x_{a,i} x_{b,i}\right) \\ &= \prod_{\mu=P(1-t)+1}^P \prod_{a=1}^n \int \frac{ds_a^\mu}{\sqrt{2\pi}} \exp\left[-\frac{1}{2} \sum_{a,b=1}^n s_a^\mu s_b^\mu \left(\delta_{ab} + \beta \sum_{i=1}^N \frac{(v_i^\mu - (m_i^\mu)^2)}{N} x_{a,i} x_{b,i}\right) + \sum_{a=1}^n i s_a^\mu \sqrt{\frac{\beta}{N}} \sum_{i=1}^N m_i^\mu x_{a,i}\right]. \end{aligned} \quad (27)$$

To continue, we assume condensation on a finite number of patterns, say the first k . We focus now on the remaining ones, namely for $\mu > k$:

$$B := \exp\left[\frac{\beta N}{2} \sum_{a=1}^n \sum_{\mu=k+1}^P (m^\mu(\mathbf{x}_a))^2\right] = \int \prod_{\mu=k+1}^P \prod_{a=1}^n \frac{dz_a^\mu}{\sqrt{2\pi}} \exp\left[-\sum_{a=1}^n \sum_{\mu=k+1}^P \frac{(z_a^\mu)^2}{2} + \sqrt{\frac{\beta}{N}} \sum_{a=1}^n \sum_{\mu=k+1}^P z_a^\mu \sum_{i=1}^N x_{a,i} \xi_i^\mu\right]. \quad (28)$$

Putting A and B together, their overall average over $(\boldsymbol{\xi}^\mu)_{\mu > k}$ takes the form

$$\begin{aligned} \mathbb{E}_{(\boldsymbol{\xi}^\mu)_{\mu > k}}[AB] &= \int \prod_{\mu=P(1-t)+1}^P \prod_{a=1}^n \frac{ds_a^\mu}{\sqrt{2\pi}} \int \prod_{\mu=k+1}^P \prod_{a=1}^n \frac{dz_a^\mu}{\sqrt{2\pi}} e^{-\sum_{a=1}^n \left(\sum_{\mu=P(1-t)+1}^P \frac{(s_a^\mu)^2}{2} + \sum_{\mu=k+1}^P \frac{(z_a^\mu)^2}{2}\right)} \\ &\quad \exp\left[\sum_{i=1}^N \sum_{\mu=k+1}^P \log \mathbb{E}_{\xi_i^\mu} e^{\sqrt{\frac{\beta}{N}} \sum_{a=1}^n x_{a,i} (\xi_i^\mu z_a^\mu + i\theta(\mu-P+R) m_i^\mu s_a^\mu) - \theta(\mu-P+R) \sum_{a,b=1}^n s_a^\mu s_b^\mu \frac{\beta(v_i^\mu - (m_i^\mu)^2) x_{a,i} x_{b,i}}{2N}}\right], \end{aligned} \quad (29)$$

where θ is Heaviside’s step function. If we call $\mathbb{E}_{\boldsymbol{\xi}} m_i^{\mu 2} =: \bar{M}^{\mu 2}$, a further expansion of the exponential yields:

$$\begin{aligned} \mathbb{E}_{(\boldsymbol{\xi}^\mu)_{\mu > k}}[AB] &= \int \prod_{\mu=P(1-t)+1}^P \prod_{a=1}^n \frac{ds_a^\mu}{\sqrt{2\pi}} \exp\left[-\frac{1}{2} \sum_{\mu=P(1-t)+1}^P \mathbf{s}^\mu \cdot (\mathbb{1} + \beta(v_{\tau^\mu} - \bar{M}^{\mu 2})Q) \mathbf{s}^\mu\right] \\ &\quad \int \prod_{\mu=k+1}^P \prod_{a=1}^n \frac{dz_a^\mu}{\sqrt{2\pi}} \exp\left\{-\sum_{\mu=k+1}^P \sum_{a=1}^n \frac{(z_a^\mu)^2}{2} + \frac{\beta}{2} \sum_{\mu=k+1}^P \sum_{a,b=1}^n z_a^\mu z_b^\mu Q(\mathbf{x}_a, \mathbf{x}_b) + \right. \\ &\quad \left. + i\beta \sum_{\mu=P(1-t)+1}^P \mathbb{E}_{\boldsymbol{\xi}}[\xi_1^\mu m_1^\mu] \sum_{a,b=1}^n z_a^\mu s_b^\mu Q(\mathbf{x}_a, \mathbf{x}_b) - \frac{\beta}{\Delta} \sum_{\mu=P(1-t)+1}^P \sum_{a,b=1}^n (\bar{M}^{\mu 2})^2 s_a^\mu s_b^\mu Q(\mathbf{x}_a, \mathbf{x}_b)\right\} \end{aligned} \quad (30)$$

We can now perform a Gaussian integration over the variables $\mathbf{z}^\mu = (z_a^\mu)_{a \leq n}$:

$$\begin{aligned} \mathbb{E}_{(\boldsymbol{\xi}^\mu)_{\mu > k}}[AB] &= \int \prod_{\mu=P(1-t)+1}^P \prod_{a=1}^n \frac{ds_a^\rho}{\sqrt{2\pi}} \exp \left[-\frac{1}{2} \sum_{\mu=P(1-t)+1}^P \mathbf{s}^\mu \cdot \left(\mathbb{1} + \beta v^\mu Q + \beta^2 Q \frac{\mathbb{E}_{\boldsymbol{\xi}}^2[\xi_1^\mu m_1^\mu]}{\mathbb{1} - \beta Q} Q \right) \mathbf{s}^\mu \right] \\ &\quad \times \exp \left[-\frac{\alpha N}{2} \log \det (\mathbb{1} - \beta Q) \right]. \end{aligned} \quad (31)$$

Finally, after an integration over the remaining Gaussian variables \mathbf{s}^μ , and using (21), we get

$$\mathbb{E}_{(\boldsymbol{\xi}^\mu)_{\mu > k}}[AB] = \exp \left[-\frac{\alpha(1-t)N}{2} \log \det (\mathbb{1} - \beta Q) - \frac{1}{2} \sum_{\mu=P(1-t)+1}^P \log \det (\mathbb{1} + \beta Q(v_{\tau^\mu} - 1) - (v_{\tau^\mu} - m_{\tau^\mu}^2)\beta^2 Q^2) \right], \quad (32)$$

where $\tau^\mu = (1 - (\mu - 1)/P)$, and $m_{\tau^\mu} = m^\mu$ are the previous retrieval accuracies. It remains to analyze the contribution given by $(\boldsymbol{\xi}^\mu)_{\mu \leq k}$:

$$C := \exp \left[\frac{\beta N}{2} \sum_{a=1}^n \sum_{\mu=1}^k (m^\mu(\mathbf{x}_a))^2 \right] = \int \prod_{a=1}^n \prod_{\mu=1}^k dm_a^\mu \sqrt{\frac{\beta N}{2\pi}} \exp \left[\sum_{a=1}^n \sum_{\mu=1}^k \left(-N\beta \frac{(m_a^\mu)^2}{2} + \beta m_a^\mu \sum_{i=1}^N \xi_i^\mu x_{a,i} \right) \right]. \quad (33)$$

Before plugging the contributions coming from A , B and C into $\mathbb{E}\mathcal{Z}_N^n$ we need to introduce a collection of Dirac deltas to fix the desired order parameters, that are organized in the overlap matrix $(Q(\mathbf{x}_a, \mathbf{x}_b))_{a,b=1}^n$:

$$1 = \int \prod_{a \leq b \leq n} dq_{ab} \delta(Q(\mathbf{x}_a, \mathbf{x}_b) - q_{ab}) = \int \prod_{a \leq b \leq n} \frac{N dr_{ab} dq_{ab}}{4\pi i} \exp \left[-\frac{1}{2} \sum_{a,b=1}^n r_{ab} (Nq_{ab} - \sum_i x_{a,i} x_{b,i}) \right]. \quad (34)$$

Hence, the averaged replicated partition function, at leading exponential order in N , takes the form

$$\begin{aligned} \mathbb{E}\mathcal{Z}_N^n &= \int \prod_{a \leq b \leq n} \frac{N dr_{ab} dq_{ab}}{4\pi i} \int \prod_{a=1}^n \prod_{\mu=1}^k dm_a^\mu \sqrt{\frac{N\beta}{2\pi}} \exp \left[-\frac{N}{2} \sum_{a,b} r_{ab} q_{ab} - \frac{\beta N}{2} \sum_{a=1}^n \sum_{\mu=1}^k (m_a^\mu)^2 \right] \\ &\quad \times \exp \left[-\frac{1}{2} \sum_{\mu=P(1-t)+1}^P \log \det (\mathbb{1} + \beta Q(v_{\tau^\mu} - 1) - (v_{\tau^\mu} - m_{\tau^\mu}^2)\beta^2 Q^2) \right] \\ &\quad \times \exp \left[-\frac{\alpha(1-t)N}{2} \log \det (\mathbb{1} - \beta Q) + N\beta^2 \Delta \sum_{a \neq b, 1}^n \frac{q_{ab}^2}{4} + N\beta \sum_{a=1}^n \left(-\frac{\lambda}{4} (1 - q_{aa})^2 + \frac{\beta \Delta - 1}{4} q_{aa}^2 \right) \right] \\ &\quad \times \left(\int \prod_{\mu=1}^k dP_\xi(\boldsymbol{\xi}^\mu) \prod_{a=1}^n dP_\xi(x_a) \exp \left[\frac{1}{2} \sum_{a,b=1}^n r_{ab} x_a x_b + \beta \sum_{\mu=1}^k \sum_{a=1}^n m_a^\mu \xi^\mu x_a \right] \right)^N, \end{aligned} \quad (35)$$

where we denote $Q = (q_{ab})_{a,b=1}^n$. We can finally express the replicated free entropy with a variational principle

coming from a saddle point argument applied to the formula above:

$$\begin{aligned}
\Phi_n &:= \lim_{N \rightarrow \infty} \Phi_{N,n} = \frac{1}{n} \text{Extr} \left\{ -\frac{1}{2} \sum_{a,b} r_{ab} q_{ab} - \frac{\beta}{2} \sum_{a=1}^n \sum_{\mu=1}^k (m_a^\mu)^2 - \frac{\alpha(1-t)N}{2} \log \det (\mathbb{1} - \beta Q) \right. \\
&+ \beta \sum_{a=1}^n \left(\frac{\beta \Delta - 1}{4} q_{aa}^2 - \frac{\lambda}{4} (1 - q_{aa})^2 \right) - \frac{\alpha t}{2R} \sum_{\mu=P(1-t)+1}^P \log \det [\mathbb{1} + \beta Q (v_{\tau\mu} - 1) - (v_{\tau\mu} - m_{\tau\mu}^2) \beta^2 Q^2] \\
&\left. + \beta^2 \Delta \sum_{a \neq b, 1}^n \frac{q_{ab}^2}{4} + \log \int \prod_{\mu=1}^k \mathbb{E}_{\xi^\mu} \int \prod_{a=1}^n dP_\xi(x_a) \exp \left[\frac{1}{2} \sum_{a,b=1}^n r_{ab} x_a x_b + \beta \sum_{\mu=1}^k \sum_{a=1}^n m_a^\mu \xi^\mu x_a \right] \right\}. \quad (36)
\end{aligned}$$

The normalized sum over $\mu = P(1-t) + 1, \dots, P$ on the second line can be turned into an integral $\int_0^t d\tau \dots$ in the large N limit. The extremization is taken w.r.t. the collection of parameters $(r_{ab}, q_{ab})_{a,b=1}^n, (m_a^\mu)_{a=1, \mu=1}^{n,k}$. Within the replica symmetric ansatz

$$\begin{cases} r_{ab} = r, & a \neq b \\ r_{aa} = -u \end{cases} \quad \begin{cases} q_{ab} = q, & a \neq b \\ q_{aa} = v \end{cases} \quad m_a^\mu = m^\mu, \quad Q = \begin{pmatrix} v & q & q & \dots & q \\ q & v & q & \dots & q \\ q & q & v & \dots & q \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ q & q & q & \dots & v \end{pmatrix} \in \mathbb{R}^{n \times n}. \quad (37)$$

The determinants of $\mathbb{1} - \beta Q$ and $\mathbb{1} + \beta Q (v_{\tau\mu} - 1) - (v_{\tau\mu} - m_{\tau\mu}^2) \beta^2 Q^2$ are easily computed:

$$\det(\mathbb{1} - \beta Q) = (1 - \beta(v - q))^n \left[1 - n \frac{\beta q}{1 - \beta(v - q)} \right] \quad (38)$$

$$\begin{aligned}
\det(\mathbb{1} + \beta Q (v_{\tau\mu} - 1) - (v_{\tau\mu} - m_{\tau\mu}^2) \beta^2 Q^2) &= [1 + \beta(v_{\tau\mu} - 1)(v - q) - (v_{\tau\mu} - m_{\tau\mu}^2) \beta^2 (v - q)^2]^{n-1} \\
&\times \left[1 + \beta(v_{\tau\mu} - 1)(v - q + nq) - (v_{\tau\mu} - m_{\tau\mu}^2) \beta^2 (v - q + nq)^2 \right]. \quad (39)
\end{aligned}$$

Further simplifications occur for the other terms in the replicated free entropy. In particular the remaining log integral is:

$$\begin{aligned}
&\int \prod_{\mu=1}^k \mathbb{E}_{\xi^\mu} \int \prod_{a=1}^n dP_\xi(x_a) \exp \left[\frac{r}{2} \sum_{a \neq b, 1}^n x_a x_b - \frac{u}{2} \sum_{a=1}^n x_a^2 + \beta \sum_{\mu=1}^k m^\mu \xi^\mu \sum_{a=1}^n x_a \right] = \\
&= \mathbb{E}_Z \int \prod_{\mu=1}^k \mathbb{E}_{\xi^\mu} \int \prod_{a=1}^n dP_\xi(x_a) \exp \left[\sqrt{r} Z x_a - \frac{u+r}{2} x_a^2 + \beta \sum_{\mu=1}^k m^\mu \xi^\mu x_a \right] = \\
&= \mathbb{E}_Z \mathbb{E}_\xi \left[\int dP_\xi(x) \exp \left((Z \sqrt{r} + \beta \mathbf{m} \cdot \boldsymbol{\xi}) x - \frac{u+r}{2} x^2 \right) \right]^n \quad (40)
\end{aligned}$$

where $Z \sim \mathcal{N}(0, 1)$, $\boldsymbol{\xi} = (\xi^1, \dots, \xi^k)$, $\mathbf{m} = (m^1, \dots, m^k)$. Finally, expanding at first order in n one has:

$$\begin{aligned}
\Phi &:= \text{Extr} \left\{ \frac{r q + u v}{2} - \beta \sum_{\mu=1}^k \frac{(m^\mu)^2}{2} - \frac{\beta^2 \Delta q^2}{4} - \frac{\alpha(1-t)}{2} \left[\log(1 - \beta(v - q)) - \frac{\beta q}{1 - \beta(v - q)} \right] \right. \\
&- \frac{\alpha t}{2} \int_0^t d\tau \left[\log(1 + \beta(v_\tau - 1)(v - q) - (v_\tau - m_\tau^2) \beta^2 (v - q)^2) + \frac{\beta q (v_\tau - 1) - 2\beta^2 q (v - q) (v_\tau - m_\tau^2)}{1 + \beta(v_\tau - 1)(v - q) - (v_\tau - m_\tau^2) \beta^2 (v - q)^2} \right] \\
&\left. + \beta \left(\frac{\beta \Delta - 1}{4} v^2 - \frac{\lambda}{4} (1 - v)^2 \right) + \mathbb{E}_{Z, \boldsymbol{\xi}} \log \int dP_\xi(x) \exp \left((Z \sqrt{r} + \beta \mathbf{m} \cdot \boldsymbol{\xi}) x - \frac{u+r}{2} x^2 \right) \right\}. \quad (41)
\end{aligned}$$

The correct stationary parameters v, m, q, u, r will be those that maximize the free entropy. Hence it is clear that if $\lambda \rightarrow \infty$ we recover the constraint $v = 1$.

3.1 Fixed point equations

Let us introduce the following notation:

$$\langle \cdot \rangle_{t, \xi} \equiv \langle \cdot \rangle_t := \frac{\int dP_\xi(x) \exp((Z\sqrt{r} + \beta \mathbf{m} \cdot \xi)x - \frac{r+u}{2}x^2)(\cdot)}{\int dP_\xi(y) \exp((Z\sqrt{r} + \beta \mathbf{m} \cdot \xi)y - \frac{r+u}{2}y^2)}, \quad (42)$$

where the subscript t emphasizes that we have already reconstructed $R = tP$ patterns. The stationarity conditions coming from (41) are

$$v = \mathbb{E}_\xi \langle X^2 \rangle_t \quad (43)$$

$$m^\mu = \mathbb{E}_\xi \xi^\mu \langle X \rangle_t, \quad \mu = 1, \dots, k \quad (44)$$

$$q = \mathbb{E}_\xi \langle X \rangle_t^2 \quad (45)$$

$$r = \frac{\alpha(1-t)\beta^2 q}{(1-\beta(v-q))^2} + \beta^2 \Delta q + \alpha t \int_0^t d\tau \left[\frac{2q\beta^2(v_\tau - m_\tau^2)}{1 + \beta(v_\tau - 1)(v - q) - (v_\tau - m_\tau^2)\beta^2(v - q)^2} + q \frac{\beta^2[v_\tau - 1 - 2\beta(v - q)(v_\tau - m_\tau^2)]^2}{[1 + \beta(v_\tau - 1)(v - q) - (v_\tau - m_\tau^2)\beta^2(v - q)^2]^2} \right] \quad (46)$$

$$u = \beta\lambda(v - 1) + \beta(1 - \beta\Delta)v - \alpha(1 - t)\beta \frac{1 - \beta(v - 2q)}{(1 - \beta(v - q))^2} - \alpha t \int_0^t d\tau \left[\frac{2v\beta^2(v_\tau - m_\tau^2) - \beta(v_\tau - 1)}{1 + \beta(v_\tau - 1)(v - q) - (v_\tau - m_\tau^2)\beta^2(v - q)^2} + q \frac{\beta^2[v_\tau - 1 - 2\beta(v - q)(v_\tau - m_\tau^2)]^2}{[1 + \beta(v_\tau - 1)(v - q) - (v_\tau - m_\tau^2)\beta^2(v - q)^2]^2} \right]. \quad (47)$$

Notice that the effect of decimation is visible only in the variables u and r that affect the local measure (19). With a close look to the expression of r we can recognize the three predicted independent noise contribution. The first term is due to pattern interference (noise (b)), and we see that it decreases as t approaches 1. The second term can be identified with the noise contribution (a), which is due to the original Gaussian noise \mathbf{Z} . The decimation noise contribution (noise (c)) is instead given by the third term, that is expressed in integral form, which correctly takes into account all the history of the process. As anticipated above, the success of decimation is determined by the interplay between noises (b) and (c). Since, as we shall see in Section 6, the retrieval accuracies remain close to one in the range of parameters α, Δ where the first step of decimation is feasible, the noise contribution (c) will be small. In addition, solving the previous equations for each decimation step shows that the benefit we gain due to the reduction of pattern interference is higher than the penalty we pay for introducing noise with decimation. As a consequence, decimation proves to be a viable strategy for matrix factorization.

For all practical purposes, we will make finite size simulations and use the discretized form present in (36) of the integral accounting for decimation contributions, starting from step 0, when no pattern has been retrieved yet. Finally, notice that mixed states solutions are possible, with the estimates aligning to more than 1 pattern, *i.e.* several m^μ 's in (44) are non-vanishing. This is not desirable in inference, since one wants to estimate one pattern at a time with the best possible performance.

3.2 Remarks

First of all, we clarify the relation between our formula and the low-rank formula for the spiked Wigner model. Therefore, let us set $\beta = 1/\Delta$, $P = 1$, which means $\alpha = 0$, and $\lambda = 0$. In this case the free entropy reads

$$\Phi := \text{Extr} \left\{ \frac{rq + uv}{2} - \frac{m^2}{2\Delta} - \frac{q^2}{4\Delta} + \mathbb{E}_{Z, \xi} \log \int dP_\xi(x) \exp \left(\left(Z\sqrt{r} + \frac{m}{\Delta} \xi \right) x - \frac{u+r}{2} x^2 \right) \right\} \quad (48)$$

Extremizing w.r.t. q and v we readily find:

$$r = \frac{q}{\Delta}, \quad u = 0. \quad (49)$$

Plugging this result inside the free entropy yields

$$\Phi := \text{Extr} \left\{ \frac{q^2}{4\Delta} - \frac{m^2}{2\Delta} + \mathbb{E}_{Z,\xi} \log \int dP_\xi(x) \exp \left(\left(Z \sqrt{\frac{q}{\Delta}} + \frac{m\xi}{\Delta} \right) x - \frac{q}{2\Delta} x^2 \right) \right\}. \quad (50)$$

Finally, extremization w.r.t. q and m yields two coupled equations

$$m = \mathbb{E}_\xi \xi \langle X \rangle_t |_{r=\frac{q}{\Delta}, u=0}, \quad q = \mathbb{E}_\xi \langle X \rangle_t^2 |_{r=\frac{q}{\Delta}, u=0} \quad (51)$$

that admit a self consistent solution satisfying a single equation

$$m = q = \mathbb{E}_\xi \xi \langle X \rangle_t |_{r=\frac{m}{\Delta}, u=0} \quad (52)$$

which is exactly the known fixed point equation for the overlap in the spiked Wigner model.

Secondly, we need to ensure a proper scaling w.r.t. β . In particular the limit $\lim_{\beta \rightarrow \infty} \frac{\Phi}{\beta}$ must be well defined at any decimation step. The only terms in the free entropy that could give rise to overscalings in β are

$$\frac{rq + uv}{2} - \frac{\beta^2 \Delta q}{4} + \frac{\beta^2 \Delta v}{4}, \quad \frac{r + u}{2}. \quad (53)$$

The latter in particular appears at the exponent in the gas free entropy in the last line of (41). Both the fixed point equations for u and r contain terms proportional to β^2 . This issue though is only apparent, and the fixed point remains well defined. To show this let us rewrite the first problematic term as follows:

$$\frac{rq + uv}{2} - \frac{\beta^2 \Delta q}{4} + \frac{\beta^2 \Delta v}{4} = \frac{-r(v - q) + (u + r)v}{2} + \frac{\beta^2 \Delta (v - q)}{4}. \quad (54)$$

In the limit $\beta \rightarrow \infty$ the term

$$-\frac{\beta q}{1 - \beta(v - q)} \quad (55)$$

arising from the square bracket in the first line of (41) forces $q \rightarrow v$ in such a way that $\beta(v - q) < 1$ remains of order $O(1)$. Hence $\frac{\beta^2 \Delta (v - q)}{4}$ and $r(v - q) = (r/\beta)\beta(v - q)$ are at most of order $O(\beta)$ as they should. It remains to verify that $u + r = O(\beta)$:

$$u + r = \beta \lambda (v - 1 + \beta v) - \beta^2 \Delta (v - q) - \frac{\alpha \beta}{1 - \beta(v - q)} - \alpha t \int_0^t d\tau \left[\frac{2\beta^2 (v - q)(v_\tau - m_\tau^2) - \beta(v_\tau - 1)}{1 + \beta(v_\tau - 1)(v - q) - (v_\tau - m_\tau^2)\beta^2 (v - q)^2} \right]. \quad (56)$$

Again, thanks to the fact that $\beta(v - q) < 1$, the correct scaling occurs.

Thirdly, we notice that for Gaussian prior, when patterns are generated from $P_\xi = \mathcal{N}(0, 1)$, retrieval is impossible if $\alpha > 0$. In fact, from the fixed point equation for m^μ , one can perform a Gaussian integration by parts on the ξ^μ obtaining:

$$m^\mu = m^\mu \beta (\mathbb{E} \langle X^2 \rangle_t - \mathbb{E} \langle X \rangle_R^2) = m^\mu \beta (v - q) \quad (57)$$

which entails $m^\mu = 0$ or $\beta(v - q) = 1$. The latter though is not possible because it would cause the free entropy to diverge to minus infinity. Hence, the only possibility is to have negligible alignment with all the patterns, $m^\mu = 0$. On the contrary if $\alpha = 0$, the diverging contribution disappears, and setting $\beta = 1/\Delta$ yields the usual PCA estimator overlap $m = q = 1 - \Delta$.

4 Low temperature limits

4.1 Sparse prior

Let us express the $\beta \rightarrow \infty$ limit of the free entropy with a prior of the form

$$P_\xi = (1 - \rho)\delta_0 + \frac{\rho}{2} [\delta_{-1/\sqrt{\rho}} + \delta_{1/\sqrt{\rho}}], \quad \rho \in (0, 1). \quad (58)$$

The case $\rho = 1$ shall be discussed separately in the end. For future convenience we introduce the notations

$$C := \beta(v - q) \in [0, 1), \quad \bar{r} := r/\beta^2, \quad U := \frac{u + r}{\beta} \quad (59)$$

where q is intended as the stationary value of the overlap solving the fixed point equations. Denote $\mathbf{m} = (m^\mu)_{\mu=1}^k$, where k is the maximum number of condensed patterns. In the low temperature limit the free entropy, re-scaled by β , and evaluated at the stationary values of the parameters involved has the form

$$\frac{1}{\beta}\Phi = -\frac{\lambda(v-1)^2}{4} - \frac{\bar{r}C}{2} + \frac{Uv}{2} + \frac{\alpha(1-t)v}{2(1-C)} - \frac{v^2}{4} - \frac{\mathbf{m}^2}{2} + \frac{\Delta Cv}{2} + \psi + \frac{\alpha tv}{2} \int_0^t d\tau \frac{2C(v_\tau - m_\tau^2) - (v_\tau - 1)}{1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2} \quad (60)$$

where

$$\psi = \frac{1}{\beta} \mathbb{E}_{\xi, Z} \log \left[1 - \rho + \rho \cosh \frac{\beta}{\sqrt{\rho}} \left(Z\sqrt{\bar{r}} + \mathbf{m} \cdot \xi \right) \exp \left(-\frac{\beta U}{2\rho} \right) \right]. \quad (61)$$

When $\beta \rightarrow \infty$ we have to distinguish two cases in the Z average:

$$\psi = O\left(\frac{1}{\beta}\right) + \frac{1}{\beta} \mathbb{E}_\xi \left(\int_{-\mathbf{m} \cdot \xi / \sqrt{\bar{r}} + U/2\sqrt{\bar{r}\rho}}^{\infty} + \int_{-\infty}^{-\mathbf{m} \cdot \xi / \sqrt{\bar{r}} - U/2\sqrt{\bar{r}\rho}} \right) \frac{dz e^{-\frac{z^2}{2}}}{\sqrt{2\pi}} \log \left[1 - \rho + \rho \cosh \frac{\beta}{\sqrt{\rho}} \left(z\sqrt{\bar{r}} + \mathbf{m} \cdot \xi \right) e^{-\frac{\beta U}{2\rho}} \right]. \quad (62)$$

The $O(\beta^{-1})$ instead comes from integration on the interval $[-\mathbf{m} \cdot \xi / \sqrt{\bar{r}} - U/2\sqrt{\bar{r}\rho}, -\mathbf{m} \cdot \xi / \sqrt{\bar{r}} + U/2\sqrt{\bar{r}\rho}]$ of the same integrand, that can be easily bounded.

Let us now focus on the first integral in (62). The hyperbolic cosine and the exponential in U dominate on the other terms in the log. Taking into account the exponential growth in the selected range of z -values the first integral can be approximated with:

$$\begin{aligned} \mathbb{E}_\xi \int_{-\mathbf{m} \cdot \xi / \sqrt{\bar{r}} + U/2\sqrt{\bar{r}\rho}}^{\infty} \frac{dz}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} \left(\frac{Z\sqrt{\bar{r}} + \mathbf{m} \cdot \xi}{\sqrt{\rho}} - \frac{U}{2\rho} \right) &= \sqrt{\frac{\bar{r}}{2\pi\rho}} \mathbb{E}_\xi e^{-\frac{1}{2\bar{r}} \left(\frac{U}{2\sqrt{\rho}} - \mathbf{m} \cdot \xi \right)^2} + \\ &+ \mathbb{E}_\xi \left(\frac{\mathbf{m} \cdot \xi}{\sqrt{\rho}} - \frac{U}{2\rho} \right) \int_{-\mathbf{m} \cdot \xi / \sqrt{\bar{r}} + U/2\sqrt{\bar{r}\rho}}^{\infty} \frac{dz}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}. \end{aligned} \quad (63)$$

The second integral in (62) can be treated similarly. Putting all the terms together one gets

$$\begin{aligned} \frac{1}{\beta}\Phi &= -\frac{\bar{r}C}{2} + \frac{\Delta Cv}{2} + \frac{Uv}{2} + \frac{\alpha(1-t)v}{2(1-C)} - \frac{v^2 + \lambda(v-1)^2}{4} - \frac{\mathbf{m}^2}{2} + \sqrt{\frac{2\bar{r}}{\pi\rho}} \mathbb{E}_\xi e^{-\frac{1}{2\bar{r}} \left(\frac{U}{2\sqrt{\rho}} - \mathbf{m} \cdot \xi \right)^2} \\ &+ \mathbb{E}_\xi \frac{\mathbf{m} \cdot \xi}{\sqrt{\rho}} \operatorname{erf} \left(\frac{\mathbf{m} \cdot \xi + \frac{U}{2\sqrt{\rho}}}{\sqrt{2\bar{r}}} \right) - \frac{U}{2\rho} \mathbb{E}_\xi \left[1 - \operatorname{erf} \left(\frac{\mathbf{m} \cdot \xi + \frac{U}{2\sqrt{\rho}}}{\sqrt{2\bar{r}}} \right) \right] + \frac{\alpha tv}{2} \int_0^t d\tau \frac{2C(v_\tau - m_\tau^2) - (v_\tau - 1)}{1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2}. \end{aligned} \quad (64)$$

Using the fact that all the parameters are evaluated at their stationary values, the previous formula can be further simplified by looking at the limiting version of the fixed point equations. In particular we have that

$$C = \sqrt{\frac{2}{\pi\rho\bar{r}}}\mathbb{E}_\xi \exp\left(-\left(\frac{U/2\sqrt{\rho}-\mathbf{m}\cdot\xi}{\sqrt{2\bar{r}}}\right)^2\right). \quad (65)$$

The value of \bar{r} can be found directly from (46) by multiplying it by β^{-2} :

$$\bar{r} = \frac{\alpha(1-t)v}{(1-C)^2} + \Delta v + \alpha tv \int_0^t d\tau \left[\frac{2(v_\tau - m_\tau^2)}{1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2} + \frac{[v_\tau - 1 - 2C(v_\tau - m_\tau^2)]^2}{[1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2]^2} \right]. \quad (66)$$

Deriving w.r.t. v we get the equation for $U = \frac{u+r}{\beta}$:

$$U = -\Delta C + v + \lambda(v-1) - \frac{\alpha(1-t)}{(1-C)} - \alpha t \int_0^t d\tau \frac{2C(v_\tau - m_\tau^2) - (v_\tau - 1)}{1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2}. \quad (67)$$

From a derivative w.r.t. U we get an equations for v :

$$v = \frac{1}{\rho}\mathbb{E}_\xi \left[1 - \operatorname{erf}\left(\frac{\mathbf{m}\cdot\xi + \frac{U}{2\sqrt{\rho}}}{\sqrt{2\bar{r}}}\right) \right]. \quad (68)$$

We can solve this equation in order to get U as a function of v , for instance by dichotomy. Finally, from (44) and (61)

$$\mathbf{m} = \mathbb{E}_\xi \langle X \rangle_{Z,\xi} = \frac{\partial\psi}{\partial\mathbf{m}} = \mathbb{E}_\xi \frac{\xi}{\sqrt{\rho}} \operatorname{erf}\left(\frac{\mathbf{m}\cdot\xi - U/2\sqrt{\rho}}{\sqrt{2\bar{r}}}\right). \quad (69)$$

If we insert these conditions in (64) we get

$$\frac{\Phi}{\beta} = \frac{\alpha(1-t)v}{2(1-C)^2} + \Delta Cv - \frac{v^2 + \lambda(v-1)^2}{4} + \frac{\mathbf{m}^2}{2} + \frac{\alpha tv}{2} \int_0^t d\tau \frac{4C(v_\tau - m_\tau^2) - (v_\tau - 1)[1 - (v_\tau - m_\tau^2)C^2]}{[1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2]^2}. \quad (70)$$

A numerical procedure to find a solution to the previous system of equations is to solve simultaneously (65) and (68) plugging into them the definitions of \bar{r} and U for a fixed m . Then one can iterate (69).

Notice that, when λ is finite, the problem is not continuous in $\rho = 1$, namely sending $\beta \rightarrow +\infty$ before or after setting $\rho = 1$ is different. This can be seen as a consequence of the non commutation of the two limits $\lim_{\beta \rightarrow \infty}$ and $\lim_{\rho \rightarrow 1}$ for the quantity $(1-\rho)^{1/\beta}$. In fact, for $\rho = 1$ the $O(\beta^{-1})$ contribution in ψ that was discarded before, is no longer negligible. Considering that contribution too would yield a free entropy of the form:

$$\begin{aligned} \frac{1}{\beta}\Phi = & -\frac{\bar{r}C}{2} + \frac{\Delta Cv}{2} + \frac{Uv}{2} + \frac{\alpha(1-t)v}{2(1-C)} - \frac{v^2 + \lambda(v-1)^2}{4} - \frac{\mathbf{m}^2}{2} + \sqrt{\frac{2\bar{r}}{\pi\rho}}\mathbb{E}_\xi e^{-\frac{1}{2\bar{r}}(\theta(1-\rho)\frac{U}{2\sqrt{\rho}}-\mathbf{m}\cdot\xi)^2} \\ & + \mathbb{E}_\xi \frac{\mathbf{m}\cdot\xi}{\sqrt{\rho}} \operatorname{erf}\left(\frac{\mathbf{m}\cdot\xi + \theta(1-\rho)\frac{U}{2\sqrt{\rho}}}{\sqrt{2\bar{r}}}\right) - \frac{U}{2\rho}\mathbb{E}_\xi \left[1 - \operatorname{erf}\left(\frac{\mathbf{m}\cdot\xi + \theta(1-\rho)\frac{U}{2\sqrt{\rho}}}{\sqrt{2\bar{r}}}\right) \right] \\ & + \frac{\alpha tv}{2} \int_0^t d\tau \frac{2C(v_\tau - m_\tau^2) - (v_\tau - 1)}{1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2}, \quad (71) \end{aligned}$$

where we set $\theta(0) = 0$. We see quickly that now, if $\rho = 1$, $v = 1$ is automatically enforced, whereas it was not so before. This discontinuous behaviour disappears if one sends $\lambda \rightarrow +\infty$ from the very beginning, as studied in [48].

4.2 Continuous priors

Consider the same definitions of \bar{r}, C, U as above. In this section we deal with priors that are symmetric and absolutely continuous over the Lebesgue measure, with density $p(x)$. We require the density to be finite at the boundaries of the support $[-a, a]$, or to go to zero with at most polynomial speed, and to be non-vanishing in the interior of the support. An example is the uniform distribution over $[-\sqrt{3}, \sqrt{3}]$. The prior dependent part in the free entropy is still

$$\psi := \frac{1}{\beta} \mathbb{E}_{Z, \xi} \log \int dP_{\xi}(x) e^{\beta(Z\sqrt{\bar{r}} + \mathbf{m} \cdot \xi)x - \frac{\beta U}{2} x^2}. \quad (72)$$

We separate the quenched Gaussian integral from the expectation w.r.t. ξ , and we perform the following changes of variables: $z \mapsto z/\sqrt{\bar{r}}$, $z \mapsto z - \mathbf{m} \cdot \xi$. This yields

$$\begin{aligned} \psi &= \frac{1}{\beta} \mathbb{E}_{\xi} \int \frac{dz}{\sqrt{2\pi\bar{r}}} e^{-\frac{(z-\mathbf{m} \cdot \xi)^2}{2\bar{r}}} \log \int_{-a}^a dx p(x) e^{-\frac{\beta U}{2} (x - \frac{z}{\bar{r}})^2 + \frac{\beta z^2}{2\bar{r}}} = \\ &= \frac{\bar{r} + \mathbf{m}^2}{2U} + \frac{1}{\beta} \mathbb{E}_{\xi} \int \frac{dz}{\sqrt{2\pi\bar{r}}} e^{-\frac{(z-\mathbf{m} \cdot \xi)^2}{2\bar{r}}} \log \int_{-a}^a dx p(x) e^{-\frac{\beta U}{2} (x - \frac{z}{\bar{r}})^2} =: \frac{\bar{r} + \mathbf{m}^2}{2U} + \bar{\psi}. \end{aligned} \quad (73)$$

The integral inside the logarithm in $\bar{\psi}$ can be computed by Laplace's approximation when β is large. However, the location of the maximum of the exponent depends on the value of z . In particular if $z \in [-Ua, Ua]$ then the maximum point falls inside the support of $p(x)$. Otherwise, given the quadratic nature of the exponent, the maximum in x will be attained at the boundaries of the support $-a$ and a . Hence the z -integral must be divided into three segments. Let us first consider:

$$\text{I} = \frac{1}{\beta} \mathbb{E}_{\xi} \int_{-Ua}^{Ua} \frac{dz}{\sqrt{2\pi\bar{r}}} e^{-\frac{(z-\mathbf{m} \cdot \xi)^2}{2\bar{r}}} \log \int_{-a}^a dx p(x) e^{-\frac{\beta U}{2} (x - \frac{z}{\bar{r}})^2} \xrightarrow{\beta \rightarrow \infty} 0 \quad (74)$$

because the exponent equals 0 at the maximum. Hence no exponential contribution in β is given, that is able to contrast the $1/\beta$ in front.

Let us turn to a second contribution:

$$\text{II} = \frac{1}{\beta} \mathbb{E}_{\xi} \int_{Ua}^{+\infty} \frac{dz}{\sqrt{2\pi\bar{r}}} e^{-\frac{(z-\mathbf{m} \cdot \xi)^2}{2\bar{r}}} \log \int_{-a}^a dx p(x) e^{-\frac{\beta U}{2} (x - \frac{z}{\bar{r}})^2} \xrightarrow{\beta \rightarrow \infty} -\frac{U}{2} \mathbb{E}_{\xi} \int_{Ua}^{+\infty} \frac{dz}{\sqrt{2\pi\bar{r}}} e^{-\frac{(z-\mathbf{m} \cdot \xi)^2}{2\bar{r}}} \left(a - \frac{z}{U}\right)^2 \quad (75)$$

From the square in the integrand we get three sub-contributions.

$$\text{IIA} = -\frac{Ua^2}{2} \mathbb{E}_{\xi} \int_{Ua}^{+\infty} \frac{dz}{\sqrt{2\pi\bar{r}}} e^{-\frac{(z-\mathbf{m} \cdot \xi)^2}{2\bar{r}}} = -\frac{Ua^2}{4} \text{erfc}\left(\frac{Ua - \mathbf{m} \cdot \xi}{\sqrt{2\bar{r}}}\right) \quad (76)$$

where the last step follows from a simple change of variables. The second one, with a shift in the integration variable, is

$$\text{IIB} = a \mathbb{E}_{\xi} \int_{Ua - \mathbf{m} \cdot \xi}^{+\infty} \frac{dz}{\sqrt{2\pi\bar{r}}} e^{-\frac{z^2}{2\bar{r}}} (z + \mathbf{m} \cdot \xi) = a \sqrt{\frac{\bar{r}}{2\pi}} \mathbb{E}_{\xi} e^{-\frac{(Ua - \mathbf{m} \cdot \xi)^2}{2\bar{r}}} + a \mathbb{E}_{\xi} \mathbf{m} \cdot \xi \text{erfc}\left(\frac{Ua - \mathbf{m} \cdot \xi}{\sqrt{2\bar{r}}}\right). \quad (77)$$

Finally, with the same shift in the integration variable, we get a third contribution:

$$\begin{aligned} \text{IIC} &= -\frac{1}{2U} \mathbb{E}_{\xi} \int_{Ua - \mathbf{m} \cdot \xi}^{+\infty} \frac{dz}{\sqrt{2\pi\bar{r}}} e^{-\frac{z^2}{2\bar{r}}} (z^2 + 2z\mathbf{m} \cdot \xi + (\mathbf{m} \cdot \xi)^2) = -\frac{1}{2U} \sqrt{\frac{\bar{r}}{2\pi}} \mathbb{E}_{\xi} (Ua + \mathbf{m} \cdot \xi) e^{-\frac{(Ua - \mathbf{m} \cdot \xi)^2}{2\bar{r}}} \\ &\quad - \frac{1}{4U} \mathbb{E}_{\xi} (\mathbf{m} \cdot \xi)^2 \text{erfc}\left(\frac{Ua - \mathbf{m} \cdot \xi}{\sqrt{2\bar{r}}}\right) - \frac{\bar{r}}{4U} \mathbb{E}_{\xi} \text{erfc}\left(\frac{Ua - \mathbf{m} \cdot \xi}{\sqrt{2\bar{r}}}\right). \end{aligned} \quad (78)$$

Now, it remains to compute the last gaussian integral:

$$\text{III} = \frac{1}{\beta} \mathbb{E}_{\boldsymbol{\xi}} \int_{-\infty}^{Ua} \frac{dz}{\sqrt{2\pi\bar{r}}} e^{-\frac{(z-\mathbf{m}\cdot\boldsymbol{\xi})^2}{2\bar{r}}} \log \int_{-a}^a dx p(x) e^{-\frac{\beta U}{2} (x-\frac{z}{U})^2}. \quad (79)$$

Thanks to the parity of $p(x)$, if we perform the changes of variables $z \mapsto -z$, $\boldsymbol{\xi} \mapsto -\boldsymbol{\xi}$, $x \mapsto -x$ we find that $\text{II}=\text{III}$. Hence we can finally recompose ψ :

$$\psi = \frac{\bar{r} + \mathbf{m}^2}{2U} + 2\text{II} = -\frac{Ua^2}{2} + \frac{1}{U} \sqrt{\frac{\bar{r}}{2\pi}} \mathbb{E}_{\boldsymbol{\xi}} (Ua - \mathbf{m} \cdot \boldsymbol{\xi}) e^{-\frac{(Ua - \mathbf{m} \cdot \boldsymbol{\xi})^2}{2\bar{r}}} + \mathbb{E}_{\boldsymbol{\xi}} \frac{\bar{r} + (Ua - \mathbf{m} \cdot \boldsymbol{\xi})^2}{2U} \text{erf}\left(\frac{Ua - \mathbf{m} \cdot \boldsymbol{\xi}}{\sqrt{2\bar{r}}}\right). \quad (80)$$

and the final form of the asymptotic free entropy is

$$\begin{aligned} \frac{\Phi}{\beta} \xrightarrow{\beta \rightarrow \infty} & -\frac{\bar{r}C}{2} + \frac{U(v-a^2)}{2} - \frac{\mathbf{m}^2}{2} + \frac{\alpha(1-t)v}{2(1-C)} + \frac{\Delta Cv}{2} - \frac{v^2 + \lambda(v-1)^2}{4} + \frac{1}{U} \sqrt{\frac{\bar{r}}{2\pi}} \mathbb{E}_{\boldsymbol{\xi}} (Ua - \mathbf{m} \cdot \boldsymbol{\xi}) e^{-\frac{(Ua - \mathbf{m} \cdot \boldsymbol{\xi})^2}{2\bar{r}}} \\ & + \mathbb{E}_{\boldsymbol{\xi}} \frac{\bar{r} + (Ua - \mathbf{m} \cdot \boldsymbol{\xi})^2}{2U} \text{erf}\left(\frac{Ua - \mathbf{m} \cdot \boldsymbol{\xi}}{\sqrt{2\bar{r}}}\right) + \frac{\alpha tv}{2} \int_0^t d\tau \frac{2C(v_\tau - m_\tau^2) - (v_\tau - 1)}{1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2}. \end{aligned} \quad (81)$$

The saddle point equations can be obtained by deriving the previous formula. The gradient w.r.t. \mathbf{m} yields:

$$\mathbf{m} = \mathbb{E}_{\boldsymbol{\xi}} \frac{\boldsymbol{\xi}}{U} \left[-\sqrt{\frac{2\bar{r}}{\pi}} e^{-\frac{(Ua - \mathbf{m} \cdot \boldsymbol{\xi})^2}{2\bar{r}}} + (Ua - \mathbf{m} \cdot \boldsymbol{\xi}) \text{erf}\left(\frac{\mathbf{m} \cdot \boldsymbol{\xi} - Ua}{\sqrt{2\bar{r}}}\right) \right]. \quad (82)$$

The derivative w.r.t. \bar{r} gives the equation for C :

$$C = \frac{1}{U} \mathbb{E}_{\boldsymbol{\xi}} \text{erf}\left(\frac{Ua - \mathbf{m} \cdot \boldsymbol{\xi}}{\sqrt{2\bar{r}}}\right). \quad (83)$$

Deriving w.r.t. U yields an equation for v :

$$\frac{a^2 - v}{2} = \frac{1}{U^2} \sqrt{\frac{\bar{r}}{2\pi}} \mathbb{E}_{\boldsymbol{\xi}} (Ua + \mathbf{m} \cdot \boldsymbol{\xi}) e^{-\frac{(Ua - \mathbf{m} \cdot \boldsymbol{\xi})^2}{2\bar{r}}} - \mathbb{E}_{\boldsymbol{\xi}} \left[\frac{\bar{r} + (Ua - \mathbf{m} \cdot \boldsymbol{\xi})^2}{2U^2} - \frac{a}{U} (Ua - \mathbf{m} \cdot \boldsymbol{\xi}) \right] \text{erf}\left(\frac{Ua - \mathbf{m} \cdot \boldsymbol{\xi}}{\sqrt{2\bar{r}}}\right). \quad (84)$$

In all the previous equations \bar{r} and U must be considered as the following functions:

$$\bar{r} = \frac{\alpha(1-t)v}{(1-C)^2} + \Delta v + \alpha tv \int_0^t d\tau \left[\frac{2(v_\tau - m_\tau^2)}{1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2} + \frac{[v_\tau - 1 - 2C(v_\tau - m_\tau^2)]^2}{[1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2]^2} \right] \quad (85)$$

$$U = -\Delta C + v + \lambda(v-1) - \frac{\alpha(1-t)}{(1-C)} - \alpha t \int_0^t d\tau \frac{2C(v_\tau - m_\tau^2) - (v_\tau - 1)}{1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2}. \quad (86)$$

Equations (83) and (84) shall be solved simultaneously at any iteration step for \mathbf{m} . This will yield a convergent algorithm to solve the system of equations.

To evaluate the free entropy at the solution of the previous system of saddle point equations we first enforce equation (84), obtaining:

$$\begin{aligned} \frac{\Phi}{\beta} \xrightarrow{\beta \rightarrow \infty} & -\frac{\bar{r}C}{2} + \frac{U(v-a^2)}{2} - \frac{\mathbf{m}^2}{2} + \frac{\alpha(1-t)v}{2(1-C)} + \frac{\Delta Cv}{2} - \frac{v^2 + \lambda(v-1)^2}{4} + \frac{1}{U} \sqrt{\frac{\bar{r}}{2\pi}} \mathbb{E}_{\boldsymbol{\xi}} (Ua - \mathbf{m} \cdot \boldsymbol{\xi}) e^{-\frac{(Ua - \mathbf{m} \cdot \boldsymbol{\xi})^2}{2\bar{r}}} \\ & + \mathbb{E}_{\boldsymbol{\xi}} \frac{\bar{r} + (Ua - \mathbf{m} \cdot \boldsymbol{\xi})^2}{2U} \text{erf}\left(\frac{Ua - \mathbf{m} \cdot \boldsymbol{\xi}}{\sqrt{2\bar{r}}}\right) + \frac{\alpha tv}{2} \int_0^t d\tau \frac{2C(v_\tau - m_\tau^2) - (v_\tau - 1)}{1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2}. \end{aligned} \quad (87)$$

Using the equation for C (83) we see that the first term in the first line and the first term in the second line can be summed together. After some algebra, imposing also (82) we get

$$\frac{\Phi}{\beta} \xrightarrow{\beta \rightarrow \infty} \frac{\bar{r}C}{2} + \frac{\mathbf{m}^2}{2} + \frac{\alpha(1-t)v}{2(1-C)} + \frac{\Delta Cv}{2} - \frac{v^2 + \lambda(v-1)^2}{4} + \frac{\alpha tv}{2} \int_0^t d\tau \frac{2C(v_\tau - m_\tau^2) - (v_\tau - 1)}{1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2}. \quad (88)$$

Finally, inserting also (85) we get

$$\frac{\Phi}{\beta} = \frac{\alpha(1-t)v}{2(1-C)^2} + \Delta Cv - \frac{v^2 + \lambda(v-1)^2}{4} + \frac{\mathbf{m}^2}{2} + \frac{\alpha tv}{2} \int_0^t d\tau \frac{4C(v_\tau - m_\tau^2) - (v_\tau - 1)[1 - (v_\tau - m_\tau^2)C^2]}{[1 + (v_\tau - 1)C - (v_\tau - m_\tau^2)C^2]^2}. \quad (89)$$

which surprisingly coincides with (70).

5 Phase diagrams for the first decimation step

The starting point of the decimation process is of crucial importance for its success. In fact, if we were to subtract an estimate $\boldsymbol{\eta}\boldsymbol{\eta}^\top/\sqrt{N}$ from the observations \mathbf{Y} where $\boldsymbol{\eta}$ had a negligible alignment with all the patterns, we would actually introducing further noise without decreasing the rank of the hidden matrix: decimation would be bound to fail.

At the 1-st step ($R = 0$ or $t = 0$) the replica symmetric decimation free entropy is simply that of a Hopfield model with Gaussian noise:

$$\Phi(t=0) := \text{Extr} \left\{ \frac{rq + uv}{2} - \beta \sum_{\mu=1}^k \frac{(m^\mu)^2}{2} - \frac{\beta^2 \Delta q^2}{4} - \frac{\alpha}{2} \left[\log(1 - \beta(v - q)) - \frac{\beta q}{1 - \beta(v - q)} \right] \right. \quad (90)$$

$$\left. + \beta \left(\frac{\beta \Delta - 1}{4} v^2 - \frac{\lambda}{4} (1 - v)^2 \right) + \mathbb{E}_{Z, \boldsymbol{\xi}} \log \int dP_\xi(x) \exp \left((Z\sqrt{r} + \beta \mathbf{m} \cdot \boldsymbol{\xi}) x - \frac{u+r}{2} x^2 \right) \right\}. \quad (91)$$

The set of fixed point equations then simplifies remarkably to

$$v = \mathbb{E}_{\boldsymbol{\xi}} \langle X^2 \rangle_t, \quad m^\mu = \mathbb{E}_{\boldsymbol{\xi}} \xi \langle X \rangle_t, \quad q = \mathbb{E}_{\boldsymbol{\xi}} \langle X \rangle_t^2 \quad (92)$$

$$r = \frac{\alpha \beta^2 q}{(1 - \beta(v - q))^2} + \beta^2 \Delta q, \quad u = \beta \lambda (v - 1) + \beta (1 - \beta \Delta) v - \alpha \beta \frac{1 - \beta(v - 2q)}{(1 - \beta(v - q))^2}. \quad (93)$$

where we have assumed condensation onto only one pattern.

Starting from these equations, one can specialize to the different 0 temperature limits that exhibit interesting features. For instance in the left panel of Figure 1, we see how the phase diagram at 0 temperature changes as sparsity increases when $\lambda \rightarrow \infty$ for the sparse Ising prior. It appears that sparsity increases the retrieval region and also the storage capacity. From the right panel we indeed see that the critical storage capacity in the noiseless limit $\Delta = 0$ diverges when $\rho \rightarrow 0$. This observation can be turned into an analytical statement as follows. To begin with, we notice that

$$C = \frac{2(1-\rho)}{\sqrt{2\pi\bar{r}\rho}} e^{-\frac{U^2}{8\bar{r}\rho}} + \frac{\rho}{\sqrt{2\pi\bar{r}\rho}} \left[e^{-\left(\frac{U/2+m}{\sqrt{2\bar{r}\rho}}\right)^2} + e^{-\left(\frac{U/2-m}{\sqrt{2\bar{r}\rho}}\right)^2} \right] \xrightarrow{\rho \rightarrow 0} 0, \quad (94)$$

exponentially fast, and

$$\bar{r} \xrightarrow{\rho \rightarrow 0} v(\alpha + \Delta). \quad (95)$$

As a consequence the equation (67) for U reduces to:

$$U = v + \lambda(v - 1) - \alpha \quad \Rightarrow \quad v = \frac{U + \alpha + \lambda}{\lambda + 1}. \quad (96)$$

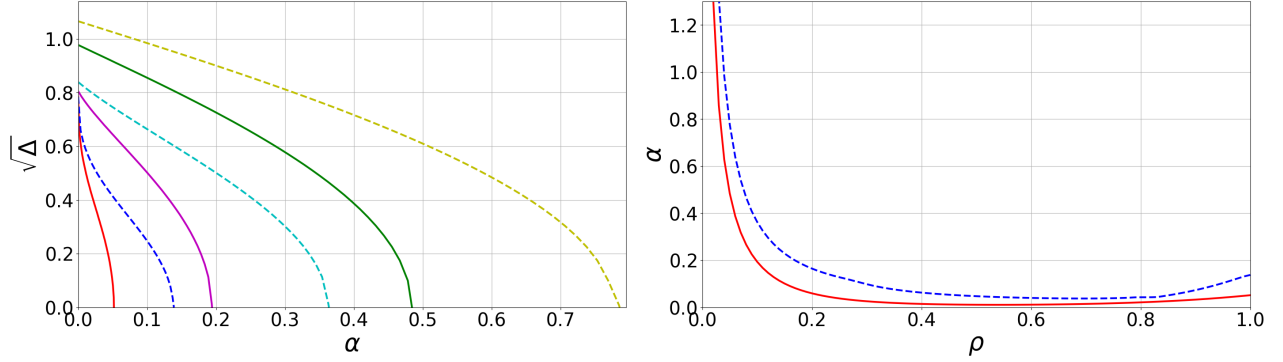


Figure 1: **Left panel:** Phase diagram for the first step of decimation in the case of sparse Ising prior. The lines show the zero temperature phase diagram for different values of the sparsity parameter ρ (using $\lambda \rightarrow \infty$). Dashed lines plot the storage capacity as a function of Δ . Solid lines signal the thermodynamic transition from the glassy phase to the retrieval phase, when configurations with non vanishing magnetizations with the patterns become thermodynamically stable. The blue and red lines are for $\rho = 1$; cyan and magenta for $\rho = 0.1$; green and yellow for $\rho = 0.05$. **Right panel:** zero temperature storage capacity α_c and critical thermodynamic storage α_F , in dashed blue and solid red lines respectively, versus sparsity ρ in the case $\Delta = 0$ (using $\lambda \rightarrow \infty$). This plot tracks the behaviour of the intersection of the dashed and solid lines with the x -axis in the left panel as ρ varies in $(0, 1]$.

We argue that U is always positive, as it serves as a norm regulator on the estimator, and we verified this statement numerically. This implies that v is always strictly positive. Equation (68) can thus be rewritten as an equation for U that reads as:

$$\frac{U + \alpha + \lambda}{\lambda + 1} = \frac{1}{\rho} - \frac{1 - \rho}{\rho} \operatorname{erf}\left(\frac{U}{2\sqrt{2\rho\bar{r}}}\right) - \frac{1}{2} \left[\operatorname{erf}\left(\frac{U/2 - m}{\sqrt{2\bar{r}\rho}}\right) + \operatorname{erf}\left(\frac{U/2 + m}{\sqrt{2\bar{r}\rho}}\right) \right]. \quad (97)$$

The error function saturates exponentially fast to 1 when $\rho \rightarrow 0$, and this entails

$$\frac{U + \alpha + \lambda}{\lambda + 1} = 1 - \frac{1}{2} \left[\operatorname{erf}\left(\frac{U/2 - m}{\sqrt{2\bar{r}\rho}}\right) + \operatorname{erf}\left(\frac{U/2 + m}{\sqrt{2\bar{r}\rho}}\right) \right] + O(e^{-K/\rho}) \quad (98)$$

for some positive constant K , and up to logarithmic corrections at the exponent in the remainder. The argument in the square brackets can go either to 0 or to 2 depending on the signs of the arguments in the error functions. However, the second possibility, that would correspond to $U/2 > |m|$, is not possible, since the l.h.s. cannot converge to 0 thanks to the positivity of U . Hence, the only alternative we have is that $U/2 < |m|$, which is also verified numerically. This implies that the limiting equation for $\rho \rightarrow 0$ appears as

$$\frac{U + \alpha + \lambda}{\lambda + 1} = 1 \quad \Rightarrow \quad \lim_{\rho \rightarrow 0} U = 1 - \alpha \quad \Rightarrow \quad \lim_{\rho \rightarrow 0} v = 1. \quad (99)$$

Finally, using the condition $U/2 < |m|$, the limit of the magnetization can be easily computed from (69):

$$m = \frac{1}{2} \left[\operatorname{erf}\left(\frac{m - U/2}{\sqrt{2\bar{r}\rho}}\right) + \operatorname{erf}\left(\frac{U/2 + m}{\sqrt{2\bar{r}\rho}}\right) \right] \xrightarrow{\rho \rightarrow 0} 1. \quad (100)$$

The behaviour depicted so far of the variables m, C, v, \bar{r} and U has been verified numerically for various values of λ, α and Δ .

In Figure 2 we plot the phase diagram for a continuous uniform prior supported on $[-\sqrt{3}, \sqrt{3}]$ with $\lambda = 0$. We verified that once that a magnetization $m \neq 0$ is a solution to the fixed point equations, then it is also thermodynamically stable, namely its free entropy is automatically bigger than that of the $m = 0$ solution, contrary to what happens for the discrete priors discussed above. The dashed line here does not signal a proper

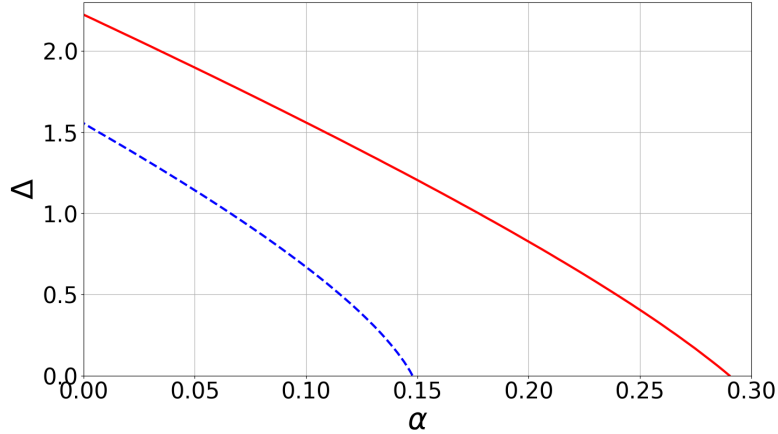


Figure 2: Zero temperature phase diagram for uniform prior supported on $[-\sqrt{3}, \sqrt{3}]$ and $\lambda = 0$. The solid line represents the thermodynamic phase transition. Below it, probability is dominated by those ‘retrieval’ states that have a non vanishing Mattis magnetization with one pattern. The dashed blue line represents a performance transition: below it the mean configuration of the Boltzmann-Gibbs measure has a better performance in reconstructing the pattern than the null estimator $\eta_{null} = 0$.

phase transition, but it is the location of the phase space where the mean square error in the reconstruction of the single pattern outperforms the null estimator $\eta_{null} = 0$, namely when:

$$\text{MSE}(\boldsymbol{\eta}; \boldsymbol{\xi}) = \frac{1}{N} \|\boldsymbol{\xi} - \langle \boldsymbol{\eta} \rangle\|^2 \simeq 1 + v - 2m < 1, \quad (101)$$

where the approximate equality holds true in the $N \rightarrow \infty$ and $\beta \rightarrow \infty$ limit. Notice that the performance of a Bayes-optimal estimator is always upper bounded by 1 thanks to the Nishimori identities, hence it is always at least as good as the null estimator.

6 Numerical tests

6.1 Testing the saddle point equations with AMP

In order to test our theoretical predictions, we need an algorithm that is able to sample from the Boltzmann-Gibbs measure, or at least that can estimate its marginals, namely the local magnetizations. Approximate message passing is an algorithm that serves the purpose. Furthermore, one needs to integrate the decimation scheme into it. The resulting algorithm was called *decimated AMP* (see Algorithm 1), which first appeared informally in [56], and then refined in [57].

It is possible to derive a suitable AMP from the set of belief propagation equations for the Boltzmann-Gibbs measure:

$$\hat{m}_{(ij) \rightarrow i}^t(x_i) \propto \int dx_j \hat{m}_{j \rightarrow (ij)}^t(x_j) \exp \left[\frac{\beta}{\sqrt{N}} Y_{ij} x_i x_j - \frac{\beta(1+\lambda)}{2N} x_i^2 x_j^2 \right] \quad (102)$$

$$m_{i \rightarrow (ij)}^{t+1}(x_i) \propto dP_{\xi}(x_i) \exp \left(\frac{\beta \lambda x_i^2}{2} \right) \prod_{k \neq i, j} \hat{m}_{(ki) \rightarrow i}^t(x_i), \quad (103)$$

by expanding in N and keeping the leading order. The resulting algorithm, which takes as input an appropriate

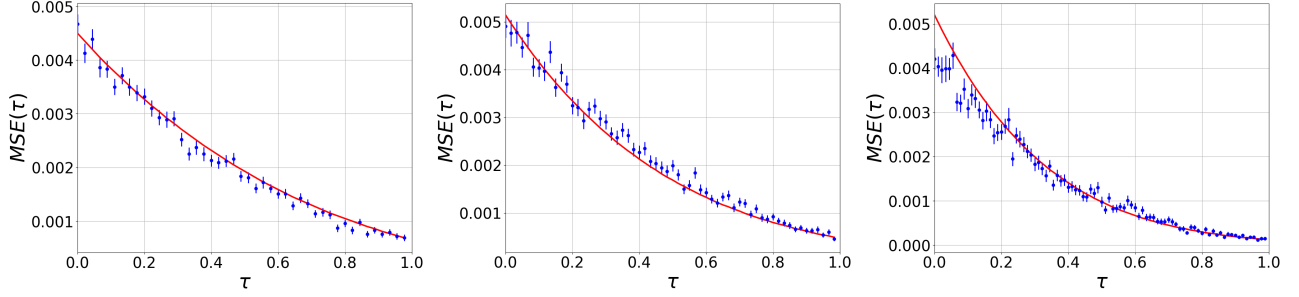


Figure 3: Mean Square Error of decimation in the case of sparse Ising priors: theory versus Decimated AMP algorithm. The red solid curves are the expected pattern MSE predicted by theory as a function of the decimation time (i.e. the number of decoded patterns). The blue data points and error bars are obtained by running DAMP over $n = 300$ independent instances. $N = 1500$, $\lambda = 0$ in all plots. **Left panel:** $\rho = 1$, $\alpha = 0.03$ namely $P = 45$, $\Delta = 0.08$ and $\beta = 10$. **Middle panel:** $\rho = 0.2$, $\alpha = 0.04$ namely $P = 60$, $\Delta = 0.09$ and $\beta = 8$. **Right panel:** $\rho = 0.15$, $\alpha = 0.06$ namely $P = 90$, $\Delta = 0.1$ and $\beta = 8$.

initialization and the data, reads:

$$\mathbf{x}^{t+1} = f(\mathbf{A}^t, \mathbf{B}^t), \quad \mathbf{v}^{t+1} = \partial_a f(\mathbf{A}^t, \mathbf{B}^t) \quad (104)$$

$$\mathbf{A}^t = \frac{\beta}{\sqrt{N}} \mathbf{Y} \mathbf{x}^t - \frac{\beta^2}{N} \mathbf{x}^{t-1} \circ (\mathbf{Y}^{\circ 2} \mathbf{v}^t) \quad (105)$$

$$\mathbf{B}^t = \frac{\beta}{N} ((1 - \mathbf{Y}^{\circ 2}) \mathbf{v} + \|\mathbf{x}^t\|^2) + \frac{\beta \lambda}{N} \sum_{i=1}^N (v_i^t + (x_i^t)^2 - 1) \quad (106)$$

where constants are summed element/component-wise, \circ is the Hadamard entry-wise product (or power), and as denoisers we have chosen the local means

$$f(a, b) = \frac{\int dP_\xi(x) x \exp(ax - \frac{bx^2}{2})}{\int dP_\xi(y) \exp(ay - \frac{by^2}{2})} \quad (107)$$

that are also applied component-wise to vectors. We denote this algorithm in a compact way by $\text{AMP}(\mathbf{Y}, \mathbf{x}^0, \mathbf{v}^0)$, and it is run until the marginals stabilize with a certain tolerance. The above AMP is used to estimate the first and second moment marginals of the Boltzmann-Gibbs measure: $x_i^\infty \simeq \langle x_i \rangle$, $v_i^\infty \simeq \langle x_i^2 \rangle - \langle x_i \rangle^2$. Of course the very same algorithm can be run on the set of modified observations \mathbf{Y}_R in (16), which is accessible to the statistician at every decimation step.

Algorithm 1 Decimated AMP (DAMP)

Require: N, P or $\alpha, \mathbf{Y}, \xi, \epsilon$
while $\mu \leq P$ **do**
 $\mathbf{g} \leftarrow \mathcal{N}(0, \mathbf{1}_N)$
 $\mathbf{x}^0 \leftarrow \sqrt{1 - \epsilon^2} \mathbf{g} + \epsilon \xi^\mu$
 $\mathbf{v}^0 \leftarrow 1 - 0.9(\mathbf{x}^0)^{\circ 2}$
 $\langle \boldsymbol{\eta}^\mu \rangle_{R=\mu-1}, \langle (\boldsymbol{\eta}^\mu)^{\circ 2} \rangle_{R=\mu-1} - \langle \boldsymbol{\eta}^\mu \rangle_{R=\mu-1}^{\circ 2} \leftarrow \text{AMP}(\mathbf{Y}_{R=\mu-1}, \mathbf{x}^0, \mathbf{v}^0)$
 $\mathbf{Y}_{R=\mu} = \mathbf{Y}_{R=\mu-1} - \frac{\langle \boldsymbol{\eta}^\mu \rangle_{R=\mu-1} \langle \boldsymbol{\eta}^\mu \rangle_{R=\mu-1}^\top}{\sqrt{N}}$
end while
Return $(\langle \boldsymbol{\eta}^\mu \rangle_{R=\mu-1}, \langle (\boldsymbol{\eta}^\mu)^{\circ 2} \rangle_{R=\mu-1})_{1 \leq \mu \leq P}$.

It is a known fact, that in the Hopfield model AMP needs to be initialized sufficiently close to the patterns to converge, and here we experience the same behavior starting from the first step of decimation until the end.

Hence DAMP is not suitable as an inference algorithm as it needs an informative initialization, whose correlation with the pattern sought is ϵ in Algorithm 1. Nevertheless, DAMP can be considered as a tool to verify that our replica computations are correct and that decimation is able to retrieve all the patterns, which means it does not corrupt itself too much.

In Figure 3 we plot the predicted theoretical curves of the expected MSE on the reconstruction on the single pattern

$$\mathbb{E}\text{MSE}(\boldsymbol{\xi}^\mu; \boldsymbol{\eta}^\mu) = \frac{1}{N} \|\boldsymbol{\xi}^\mu - \langle \boldsymbol{\eta}^\mu \rangle_{t|tP=\mu-1}\|^2 \simeq 1 + q_t - 2m_t \quad (108)$$

in red, where the subscript t indicates that we are at the decimation time t . The blue data points and error bars are obtained from an average of 300 instances of DAMP run on independently generated data. We considered different values of sparsity and the regularization parameter λ was always set to 0. In every case the theoretical curve seems to reproduce accurately the behaviour of the pattern MSE, yielding a good confirmation of our RS theory.

6.2 Expected decimation performance

In this section, we compare the expected denoising performance of decimation with the typical performance of a Rotation Invariant Estimator (RIE) introduced in [49]. A RIE is characterized by the fact that it provides an estimate of the original matrix $\boldsymbol{\xi}\boldsymbol{\xi}^\top$ which has the same eigenbasis as the one of the data matrix \mathbf{Y} . Once the eigenbasis is established, one only has to produce an estimate on the spectrum based on that of \mathbf{Y} . As such, the RIE is a purely spectral estimator and it does not exploit the prior knowledge on the signal components. Among the possible RIEs, the one that acts optimally on the spectrum of \mathbf{Y} is

$$\hat{\boldsymbol{\lambda}} = \boldsymbol{\lambda}_{\mathbf{Y}} - 2\Delta\mathcal{H}[\rho_{\mathbf{Y}}](\boldsymbol{\lambda}_{\mathbf{Y}}) \quad (109)$$

where $\hat{\boldsymbol{\lambda}}$ and $\boldsymbol{\lambda}_{\mathbf{Y}}$ are the vector of the eigenvalues of the estimate and of $\mathbf{Y}\sqrt{N}$ respectively, $\mathcal{H}[\rho_{\mathbf{Y}}]$ is the Hilbert transform of the spectral density of \mathbf{Y}/\sqrt{N} .

We shall measure the performance of an estimator \mathbf{S} , whose eigenvalues are of order 1 by convention, with the matrix MSE:

$$\text{mMSE}(\mathbf{S}; \boldsymbol{\xi}) = \frac{1}{N} \mathbb{E} \left\| \mathbf{S} - \frac{\boldsymbol{\xi}\boldsymbol{\xi}^\top}{\sqrt{NP}} \right\|_F^2, \quad (110)$$

and the matrix norm is the Frobenius' norm. The estimator produced by decimation would thus be

$$\mathbf{S}_{\text{dec}} := \sum_{\mu=1}^P \frac{\langle \boldsymbol{\eta}^\mu \rangle_{R=\mu-1} \langle \boldsymbol{\eta}^\mu \rangle_{R=\mu-1}^\top}{\sqrt{NP}} \quad (111)$$

In order to make the comparison we need to connect the mMSE predicted by the theory for the decimation estimator with the definition (110), namely to re-express the latter in terms of the order parameters of the decimation free entropies. This can be done as follows, leveraging the assumption (19). By expanding the square in the mMSE definition evaluated at \mathbf{S}_{dec} we recognize three main contributions:

$$\frac{1}{N^2P} \sum_{i,j=1}^N \sum_{\mu,\nu=1}^P \mathbb{E}[\xi_i^\mu \xi_j^\mu \xi_i^\nu \xi_j^\nu] = \frac{1+\alpha}{2} + o_N(1) \quad (112)$$

$$\frac{1}{N^2P} \sum_{i,j=1}^N \sum_{\mu,\nu=1}^P \mathbb{E}[\xi_i^\mu \langle \eta_j^\mu \rangle \xi_i^\nu \langle \eta_j^\nu \rangle] \quad (113)$$

$$\frac{1}{N^2P} \sum_{i,j=1}^N \sum_{\mu,\nu=1}^P \mathbb{E}[\langle \eta_i^\mu \rangle \langle \eta_j^\mu \rangle \langle \eta_i^\nu \rangle \langle \eta_j^\nu \rangle] \quad (114)$$

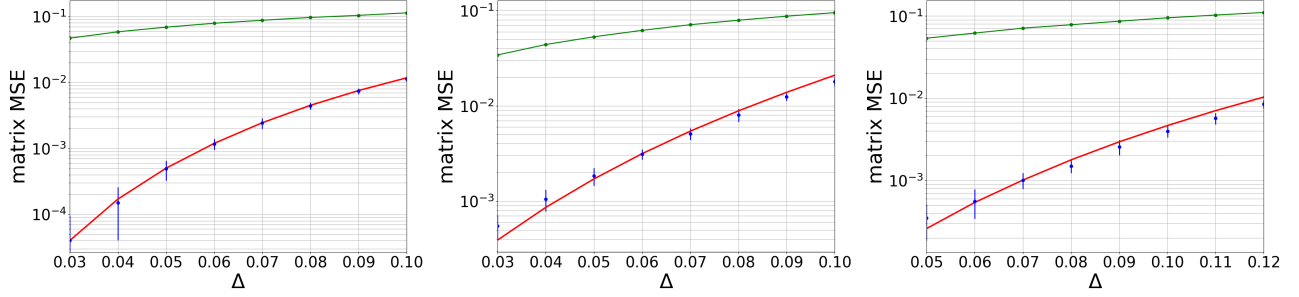


Figure 4: Matrix MSE as a function of Δ for sparse Ising priors with various sparsities. In green the denoising performance of a RIE, obtained by averaging over 30 independent samples. Error bars, corresponding to one standard deviation, are too small to be seen. In red, the performance predicted for an algorithm implementing decimation. The blue data points are obtained averaging over 30 DAMP's outputs, run on independently generated data. Error bars correspond to one standard deviation. In all cases $\lambda = 0$, $\beta = 8$ and $N = 1500$. **Left panel:** $\rho = 1$, $\alpha = 0.03$ namely $P = 45$ and $\Delta = 0.08$. **Middle panel:** $\rho = 0.2$, $\alpha = 0.07$ namely $P = 105$ and $\Delta = 0.09$. **Right panel:** $\rho = 0.15$, $\alpha = 0.07$ namely $P = 105$ and $\Delta = 0.1$.

where we dropped the subscripts in the Gibbs brackets for convenience. While the first one can be computed right away using the properties of the prior, the other two require some extra effort. Concerning (113) we have:

$$\begin{aligned} \frac{1}{N^2 P} \sum_{i,j=1}^N \sum_{\mu,\nu=1}^P \mathbb{E}[\xi_i^\mu \langle \eta_j^\mu \rangle \xi_i^\nu \langle \eta_j^\nu \rangle] &= \frac{1}{N^2 P} \sum_{i,j=1}^N \sum_{\mu,\nu=1}^P [\delta_{\mu\nu} \xi_i^\mu \langle \eta_j^\mu \rangle \xi_i^\nu \langle \eta_j^\nu \rangle + \delta_{ij} \mathbb{E}(\xi_i^\mu)^2 \langle \eta_i^\mu \rangle^2] = \\ &= \frac{1}{P} \sum_{\mu=1}^P (m^\mu)^2 + \frac{\alpha}{P} \sum_{\mu=1}^P q^\mu + o_N(1) \end{aligned} \quad (115)$$

where we have enforced (19) and q^μ and m^μ are the overlap and Mattis magnetization respectively coming from the μ -th decimation step. Let us now turn to (114). Using similar arguments one can argue that:

$$\frac{1}{N^2 P} \sum_{i,j=1}^N \sum_{\mu,\nu=1}^P \mathbb{E}[\langle \eta_i^\mu \rangle \langle \eta_j^\mu \rangle \langle \eta_i^\nu \rangle \langle \eta_j^\nu \rangle] = \frac{1}{P} \sum_{\mu=1}^P (q^\mu)^2 + \alpha \left(\frac{1}{P} \sum_{\mu=1}^P q^\mu \right)^2 + o_N(1) \quad (116)$$

Therefore, collecting all the contributions one gets the asymptotic prediction:

$$\text{mMSE}(\mathbf{S}_{\text{dec}}; \boldsymbol{\xi}) \simeq \frac{1}{P} \sum_{\mu=1}^P (1 + (q^\mu)^2 - 2(m^\mu)^2) + \alpha \left(1 - \frac{1}{P} \sum_{\mu=1}^P q^\mu \right)^2. \quad (117)$$

In Figure 4 we compare the performance of the RIE, in green, against the theoretical performance predicted for decimation in red, and the blue data points are obtained using the estimator produced by decimation (DAMP). As we can see there is a good agreement between DAMP and the theory, and both outperform the RIE as we expected. The RIE appears more robust to both noises (a) and (b), tuned by Δ and α respectively. On the contrary, the performance of decimation deteriorates quickly as soon as we get out of the retrieval region in the phase diagrams Figure 1-2, and the amount of noise it can bear is strongly affected by the nature of the signal (sparse Ising or continuous). However, one must bear in mind that RIEs are suitable only for matrix denoising, and no information is reconstructed on the signal factor $\boldsymbol{\xi}$. Moreover, we notice that the performance of the RIE does not change sensibly from the left to the right panel ($\rho = 1$ to $\rho = 0.15$), and this is coherent with its purely spectral nature. In fact, the empirical spectral distribution of $\boldsymbol{\xi} \boldsymbol{\xi}^T / \sqrt{NP}$ always converges to a Marchenko-Pastur law because of the completely factorized prior on the elements of $\boldsymbol{\xi}$. Hence, the small changes from the left to the right panel are mostly due to the slight increment in the noise level Δ and the aspect ratio (or load) α .

6.3 A ground state oracle for sparse Ising priors

Our ground state oracle is based on an iterated simulated annealing (SA) routine that can be found in Algorithm 2, which is a refinement of the one in [48].

Algorithm 2 Simulated annealing (SA)

Require: N , \mathbf{Y} , threshold, $\beta_{\max} \in \mathbb{R}$, niter ($\in \mathbb{N}$), maxr ($\in \mathbb{N}$), restarts ($\in \mathbb{N}$)

```

ityr  $\leftarrow$  0
found  $\leftarrow$  False
while ityr < 300 do and found == False
  stop  $\leftarrow$  0
   $\beta \leftarrow$  0
   $\mathbf{s} \leftarrow$  random sample from  $\prod_{i=1}^N P_{\xi}$ 
  ityr  $\leftarrow$  ityr + 1
  if ityr + restarts > maxr then
    return  $\mathbf{s}$ , ityr
  end if
  if ityr % 20 = 0 then
    threshold  $\leftarrow$  threshold  $\cdot$  0.9975
  end if
  while  $k <$  niter do
     $k \leftarrow k + 1$ 
     $\beta \leftarrow 1 + \frac{k}{\text{niter}} \cdot \beta_{\max}$ 
     $\mathbf{h} \leftarrow \frac{\mathbf{Y}}{\sqrt{N}} \mathbf{s}$ 
     $V \leftarrow \frac{\|\mathbf{s}\|^2}{N} + \frac{\lambda}{N} (\|\mathbf{s}\|^2 - 1)$ 
     $\mathbf{Z}_{\text{loc}} \leftarrow (1 - \rho) \mathbf{1} + \rho \cosh(\beta \mathbf{h}) e^{-\frac{\beta V}{2}}$  (Scalar functions are applied component-wise to vectors.)
    sample  $\mathbf{ss}$  from  $\exp(\beta \mathbf{h} \cdot (\cdot) - \frac{\beta V}{2}) / \mathbf{Z}_{\text{loc}}$ 
    if  $\|\mathbf{s} - \mathbf{ss}\| < 10^{-3}$  then
       $\mathbf{s} \leftarrow \mathbf{ss}$ 
      stop  $\leftarrow$  stop + 1 (Updates become negligible.)
      if stop > 5 then
        if  $-E(\mathbf{s} | \mathbf{Y}) >$  threshold then
          return  $\mathbf{s}$ , ityr
        else
          break (wrong energy, try again)
        end if
      end if
    end if
  end while
  else
    stop  $\leftarrow$  0
     $\mathbf{s} \leftarrow \mathbf{ss}$ 
  end if
end while
end while

```

The energy landscape at the various steps of decimation is very similar to that of the Hopfield model. Consequently, algorithms that search for minima get frequently stuck in metastable states, which have a low overlap with the patterns. SA is not immune to this phenomenon. Therefore, we equip our SA routine with an acceptance criterion of the configuration output by the algorithm, that is based on the computation of the energy:

$$-E(\mathbf{s} | \mathbf{Y}_R) = \frac{1}{2\sqrt{N}} \mathbf{s}^T \mathbf{Y}_R \mathbf{s} - \frac{\|\mathbf{s}\|^4}{4N} - \frac{\lambda}{4N} (\|\mathbf{s}\|^2 - 1)^2 \quad (118)$$

which is nothing the energy of our model at the R -th decimation step. Notice that this quantity is accessible by the Statistician and it is thus correct to use it as an input for a candidate algorithm. In Algorithm 2 niter is the maximum number of temperature updates we allow, maxr is instead the maximum number of restarts allowed, considering also the restarts coming from previous pattern searches. The reason why we introduced this additional control is that typically when a bad configuration is accepted as a pattern estimate by mistake,

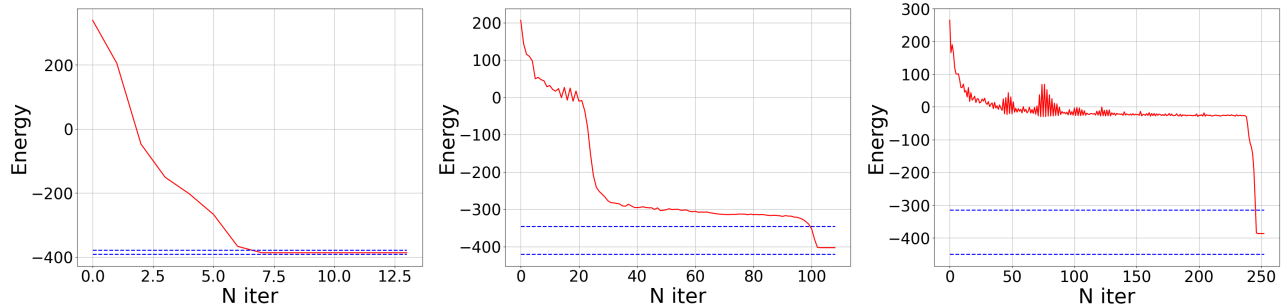


Figure 5: Energy landscape exploration of the Simulated Annealing applied to sparse Ising priors. On the vertical axis we have the energy value as a function of the number of iterations (temperature updates) of SA on the horizontal axis. For all the three plots $N = 1500$, $\alpha = 0.01$ (namely only 15 patterns to be found), $\Delta = 0.05$ and $\lambda = -0.08$. From the left to the right: $\rho = 1, 0.3, 0.15$. The patterns were reconstructed exactly in all three cases. SA finds immediately the patterns for low sparsities $\rho \sim 1$. As soon as sparsity increases, a lot of configurations start to exhibit an almost vanishing energy (recall that the noise shifts this value). The dashed blue lines mark the highest and the lowest pattern energy. As we can see the band they identify is narrow with low sparsity, and it becomes wider for higher values of sparsity due to more intense fluctuations.

the ensuing searches for other patterns require even more restarts. The above SA routine has to be combined with decimation, so once a configuration is accepted as a pattern the observations are modified $\mathbf{Y} \leftarrow \mathbf{Y} - \frac{\mathbf{s}\mathbf{s}^T}{\sqrt{N}}$ and the routine is restarted. In order to make sure we really find patterns, we thus run all the algorithm (SA plus decimation) multiple times, typically five, and then we accept the output that required the least number of restarts to be produced. This procedure is costly, and as noticed already in [48], it requires an exponential number of restarts.

Algorithm 2 suffers from the same issues as the one in [48]. For instance, the overall decimation procedure still requires an exponential (in N) number of restarts. However, the presence of sparsity introduces further non-trivial complications. In fact, the signal components are no longer constrained on the hypercube, and this allows for fluctuations in the norm of the outputs that reflect in fluctuations on the average energy of the patterns. Specifically, the more sparse the signal is, the wider the gap between the highest and the lowest energy of the patterns. These fluctuations can challenge the energy restarting criterion in our SA routine, that can thus confuse a metastable state for a pattern.

Furthermore, one observes that when too few patterns are stored or remain in \mathbf{Y} , it is harder for the SA routing to find them. If, for instance, we only have one pattern left, the Hebbian matrix $\xi\xi^T$, which is supposed to attract the \mathbf{x} -configurations towards the pattern, has only a fraction ρ^2 of non-zero components. This gives rise to a large number of configurations that have degenerate energy, close to 0. The energy landscape thus appears as a golf course, flat almost everywhere, except for a pit, corresponding to the pattern left. From our numerical experiments, this effect seems to hold also for more than one, but still few, patterns stored. See Figure 5.

6.4 Reversed decimation

In all the tests we have run, the performance of decimation in reconstructing the patterns improves along the procedure itself. The last patterns are always better estimated than the first ones, and this supports the idea that decimation effectively decreases the pattern interference. In particular, it is clear that the quality of reconstruction of one pattern depends on the previous “history” of the process.

Once the procedure exhausts the patterns, one can imagine to run it again backwards, keeping the last half of the patterns that were reconstructed with higher accuracy. As illustrated in Figure 6, this improves the reconstruction performance also for the first half of the patterns. One can then re-iterate the same procedure, keeping only the first 1/2 and the last 1/4 of the patterns, that are now the best reconstructed ones. This in

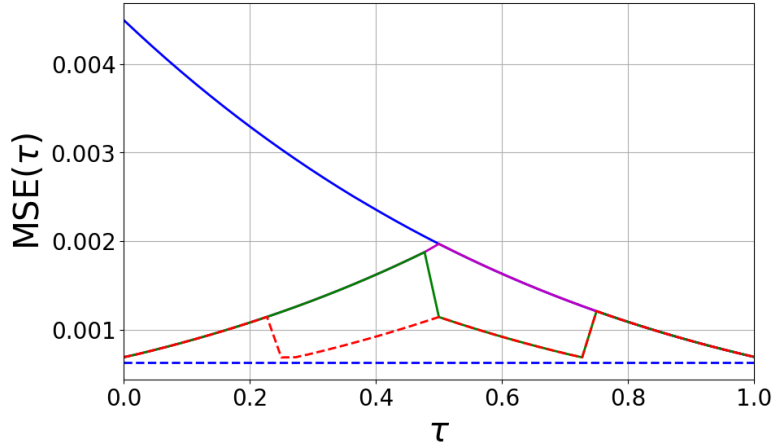


Figure 6: Improvement in performance obtained re-iterating decimation for Rademacher prior. In this example $\Delta = 0.08$, $\alpha = 0.03$, $\rho = 1$ and $\beta = 10$. The blue line is the first run, where the expected MSE on the reconstruction of the single patterns decreases along decimation. The magenta curve is instead obtained by fixing the last half of pattern MSEs, and running decimation backwards. Starting from the magenta line, we obtained the green solid line by fixing the first half and the last quarter of MSEs, and then running decimation for finding the third quarter of MSEs. Finally, the red dashed line was obtained from the green line running decimation again, with fixed first quarter and last half of MSEs. The blue dashed line is the expected MSE predicted by the rank one formula. Coherently, the last decimation steps approach the rank-one formula MSE from above, because the interference noise has been almost completely eliminated, except for noise of decimation itself, that is responsible for the final small gap.

turn leads to a further improvement in the reconstruction also for the middle patterns. This reasoning can be iterated ad libitum.

In Figure 6 we see how performance improves in the various rounds of decimation, and we compare it to the performance predicted by the rank-one formula, i.e. what we should have for any sub-linear rank ($\alpha = 0$, see Section 7). We see that, little by little, the performance approaches that of the rank-one formula.

7 Related works

7.1 Unlearning and dreaming

As evident from Figure 1, without having strong sparsity, the storage capacity of the model is not very large, and the network is far from being able to store an over-complete basis of \mathbb{R}^N . In an attempt to solve this issue one can pre-process the observation matrix with Hebbian unlearning [58, 59], with which decimation itself bears some similarity. Unlearning consists in iterating a zero temperature dynamics until convergence, which is likely to occur at a spurious state $\boldsymbol{\eta}$ that is then removed from the observations $\mathbf{Y} \leftarrow \mathbf{Y} - \varepsilon \boldsymbol{\eta} \boldsymbol{\eta}^\top / \sqrt{N}$, with a small ε . If run for an appropriate number of times, unlearning acts on the energy landscape penalizing spurious metastable states. This procedure has two fundamental parameters to be tuned: ε and the number of times D it is iterated [60]. If ε or D are too large one risks to remove also the wanted patterns.

Apart from numerical evidence, there is little theoretical understanding of the unlearning procedure as illustrated above. However, there are other convenient iterative ways of modifying the Hebbian matrix [61–64] that converge to the so called pseudo-inverse learning rule (or modifications of it) [65–67], which in turn is able to increase the storage capacity to $\alpha_c = 1$.

Despite the apparent similarities, the goal of decimation is very different from that of unlearning. Its aim is to find a pattern, and not a metastable state, and to remove it completely (or almost completely) from \mathbf{Y} , which amounts to set $\varepsilon = 1$ (or close to 1) above. Furthermore, it is worth stressing that, unlike classical unlearning,

we have a theoretical control on decimation, namely we can track its behaviour step by step.

7.2 Sub-linear rank

In a recent work [57] the authors discuss the denoising of large matrices in the same setting as ours, with a main focus on the case $P = N^\delta$, $\delta \in (0, 1)$, i.e. a sub-linear rank regime. In the mentioned paper, it is stated that, as long as the prior on the $N \times P$ matrix $\boldsymbol{\xi}$ is completely factorized over the matrix elements, the mutual information between $\boldsymbol{\xi}$ and the data is given by the rank-one replica formula for *any* sub-linear rank regime, in agreement with [68]. Though not explicitly stated in our previous work [48], our findings indeed suggest the same result, as it can be deduced from Section 3.2. In fact our free entropy, which is in close relation with the mutual information between observations and signal, takes the same form for any P such that $P/N \rightarrow 0$. Furthermore, for $\alpha = 0$ and $\beta = 1/\Delta$, the fixed point equations admit a self-consistent solution that satisfies the Nishimori identities, which suggests that Bayes-optimality is recovered. From the form of the free entropy (41), it is also evident that the effect of decimation is visible only for truly extensive rank. The reason is that, if we penalize a finite number of directions in a space of dimension growing to infinity, the system can easily find other favoured directions to thermalize in. In other words, the $p^\mu(\mathbf{x})$'s in (17) give a sub-extensive contribution that can be neglected in any sub-linear rank regime.

Another delicate point is the definition of DAMP. We stress that in (105) and (106) the presence of a high-rank spike inside \mathbf{Y} can induce non-trivial modifications both in \mathbf{A} and \mathbf{B} . More specifically, it is known that, for instance, the Onsager reaction in (105) containing \mathbf{Y}^{o2} has different asymptotically equivalent formulations. In the case of a Gaussian channel with a low-rank spike \mathbf{Y}^{o2} can be replaced by an all-ones matrix. This is due to the fact that the rank of the spike is not large enough to induce modifications in the spectrum of the noise matrix. In the high-rank regime, on the contrary, the extensive rank starts to play a role and gives rise to important contributions in the reaction term. Moreover, the reaction term changes also along the decimation procedure, in which one further perturbs the data matrix with the high rank matrix of the decimation estimates $\sum_{\mu=P-R+1}^P \frac{\boldsymbol{\eta}^\mu \boldsymbol{\eta}^{\mu\top}}{\sqrt{N}}$. Hence, the formulation in (105)-(106) turns out to be convenient. The low-rank regime is insensitive to the aforementioned changes.

Despite we were not able to prove it, Figure 6 suggests that re-iterating decimation in a proper way could lead to a performance similar to that predicted by the low rank replica symmetric formula. One may be led to think that reversed decimation yields Bayes-optimal performance. This is however not true. In fact, in the high rank case the spike induces a non-negligible perturbation of the spectrum of the noise matrix that can be used to perform inference (this deformation is captured by the RIE for instance) especially for large α 's, where decimation fails.

7.3 Channel universality properties

Low-rank spiked models are known to fulfill channel universality [69–71], namely for any well-behaved $P_{\text{out}}(y | x)$ and data generated with the rule

$$Y_{ij} \sim P_{\text{out}}\left(\cdot \mid \sum_{\mu=1}^P \frac{\xi_i^\mu \xi_j^\mu}{\sqrt{N}}\right) \quad (119)$$

the mutual information between the data \mathbf{Y} and $\boldsymbol{\xi}$ can be computed through an equivalent Gaussian channel as in (1) with a properly tuned noise intensity Δ . The proof of this equivalence requires two concomitant behaviours, *i*) universality in the likelihood, and *ii*) universality in the quenched disorder (i.e. the law of the data \mathbf{Y}), and holds as long as $P^3/\sqrt{N} \rightarrow 0$ [70]. Informally, the main idea is to expand $P_{\text{out}}\left(\cdot \mid \sum_{\mu=1}^P \frac{\xi_i^\mu \xi_j^\mu}{\sqrt{N}}\right)$ around 0 in its second entry up to second order, since for low-rank spikes $\sum_{\mu=1}^P \frac{\xi_i^\mu \xi_j^\mu}{\sqrt{N}}$ is small for any fixed couple of indices i, j . On the contrary, in the high-rank setting the higher moments of the spike start to matter, meaning that the previous expansion fails, and universality breaks down.

In our mismatched setting one can still count on the universality of the likelihood *for a single decimation step*. In fact, here the Statistician assumes to observe a low-rank spike, that is they consider

$$Y_{ij} \sim P_{\text{out}}\left(\cdot \mid \frac{x_i x_j}{\sqrt{N}}\right) \quad (120)$$

whereas the data are generated through (1). The free entropy of the related model reads as

$$\frac{1}{N} \mathbb{E}[\log \mathcal{Z}_R - \sum_{i,j} \log P_{\text{out}}(Y_{ij} \mid 0)] = \frac{1}{N} \mathbb{E} \log \int dP_{\xi}(\mathbf{x}) \exp \left[\sum_{i,j} \left(\log P_{\text{out}}\left(Y_{ij} \mid \frac{x_i x_j}{\sqrt{N}}\right) - \log P_{\text{out}}(Y_{ij} \mid 0) \right) \right] \quad (121)$$

where $\sum_{i,j} \log P_{\text{out}}(Y_{ij} \mid 0)$ has been subtracted to have a proper scaling. From the above equation one readily realizes that an expansion up to second order of P_{out} yields the desired equivalent quadratic model, for which our computations hold. However, we stress that exploiting this universality produces errors of $O(N^{-1/2})$. These errors accumulate along the $P = O(N)$ steps of decimation resulting in potentially non-negligible deviations from the original model towards the end of the procedure.

8 Conclusion and outlooks

Building on the results of [48], we have extended the analysis of the decimation procedure to a wide class of priors on the matrix elements of the factors ξ for symmetric matrix factorization. We provided exhaustive numerical evidence in support of our replica theory, via the introduction of DAMP, whose performance in pattern retrieval, and matrix denoising matches the one predicted by the theory. Our numerical experiments confirm that decimation is a viable strategy for matrix factorization. In particular, as long as the first step is feasible, i.e. the procedure is started at a point of the phase diagram where there is a non-vanishing Mattis magnetization with one of the patterns, decimation is able to find all of them, up to a permutation. We stress again that DAMP is not an appropriate algorithm for inference, since it needs a strongly informative initialization. Nevertheless, in the case of sparse Ising priors, we were able to find a ground state oracle that is able to find all the patterns in suitable regions of the phase space of the decimation neural network models. The latter still suffers from an exponential complexity: it needs an exponential number of restarts (in N) in order to find all the patterns and discard correctly the spurious states it may get stuck in.

The idea of reversed decimation and unlearning are insightful perspectives. In fact, in order to increase the storage capacity of the neural networks, or equivalently to widen the region of the phase space where we can perform matrix factorization, one could pre-process the Hebbian interaction matrix using a local updating rule, as the ones described in [63, 72]. In these works, besides the usual “forgetting” mechanism, the authors also consider a consolidation of the memories, which avoids the risk of corrupting the Hebbian interaction too much. This pre-processing could be combined with reversed decimation in order to obtain a better performing procedure that is also more robust to pattern interference.

Finally, in an upcoming work, we shall tackle the asymmetric problem, which is closer to practical applications. Here, the Statistician has to reconstruct two independent matrices $\mathbf{F} \in \mathbb{R}^{N \times P}$ and $\mathbf{X} \in \mathbb{R}^{P \times M}$ from the observations

$$\mathbf{Y} = \frac{1}{\sqrt{N}} \mathbf{F} \mathbf{X} + \sqrt{\Delta} \mathbf{Z} \in \mathbb{R}^{N \times M} \quad (122)$$

in the scaling limit $N, M, P \rightarrow \infty$ with $P/N = \alpha > 0$ and $P/M = \gamma > 0$.

Acknowledgments

We would like to thank Enzo Marinari and Federico Ricci-Tersenghi for their suggestions on the reversed decimation, Enzo Marinari and Marco Benedetti for discussions on unlearning, as well as Florent Krzakala, Lenka Zdeborová and Jean Barbier for many fruitful discussions on matrix factorization. MM acknowledges financial support by the PNRR-PE-AI FAIR project funded by the NextGeneration EU program.

References

- [1] B. A. Olshausen and D. J. Field, "Emergence of simple-cell receptive field properties by learning a sparse code for natural images," *Nature*, vol. 381, no. 6583, pp. 607–609, 1996.
- [2] B. A. Olshausen and D. J. Field, "Sparse coding with an overcomplete basis set: A strategy employed by v1?" *Vision Research*, vol. 37, no. 23, pp. 3311–3325, 1997, ISSN: 0042-6989. DOI: 10.1016/S0042-6989(97)00169-7.
- [3] K. Kreutz-Delgado, J. F. Murray, B. D. Rao, K. Engan, T.-W. Lee, and T. J. Sejnowski, "Dictionary Learning Algorithms for Sparse Representation," *Neural Computation*, vol. 15, no. 2, pp. 349–396, Feb. 2003, ISSN: 0899-7667. DOI: 10.1162/089976603762552951.
- [4] J. Mairal, F. Bach, J. Ponce, and G. Sapiro, "Online dictionary learning for sparse coding," in *Proceedings of the 26th Annual International Conference on Machine Learning*, ser. ICML '09, Montreal, Quebec, Canada: Association for Computing Machinery, 2009, pp. 689–696, ISBN: 9781605585161. DOI: 10.1145/1553374.1553463.
- [5] A. Mnih and R. R. Salakhutdinov, "Probabilistic matrix factorization," in *Advances in Neural Information Processing Systems*, J. Platt, D. Koller, Y. Singer, and S. Roweis, Eds., vol. 20, Curran Associates, Inc., 2007.
- [6] J. Mairal, M. Elad, and G. Sapiro, "Sparse representation for color image restoration," *IEEE Transactions on Image Processing*, vol. 17, no. 1, pp. 53–69, 2008. DOI: 10.1109/TIP.2007.911828.
- [7] J. Wright, Y. Ma, J. Mairal, G. Sapiro, T. S. Huang, and S. Yan, "Sparse representation for computer vision and pattern recognition," *Proceedings of the IEEE*, vol. 98, no. 6, pp. 1031–1044, 2010. DOI: 10.1109/JPROC.2010.2044470.
- [8] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *Trans. Img. Proc.*, vol. 15, no. 12, pp. 3736–3745, Dec. 2006, ISSN: 1057-7149. DOI: 10.1109/TIP.2006.881969.
- [9] H. Zou, T. Hastie, and R. Tibshirani, "Sparse principal component analysis," *Journal of Computational and Graphical Statistics*, vol. 15, no. 2, pp. 265–286, 2006. DOI: 10.1198/106186006X113430.
- [10] A. Belouchrani, K. Abed-Meraim, J.-F. Cardoso, and E. Moulines, "A blind source separation technique using second-order statistics," *IEEE Transactions on Signal Processing*, vol. 45, no. 2, pp. 434–444, 1997. DOI: 10.1109/78.554307.
- [11] E. Candès and B. Recht, "Exact matrix completion via convex optimization," *Commun. ACM*, vol. 55, no. 6, pp. 111–119, Jun. 2012, ISSN: 0001-0782. DOI: 10.1145/2184319.2184343.
- [12] E. J. Candès and T. Tao, "The power of convex relaxation: Near-optimal matrix completion," *IEEE Transactions on Information Theory*, vol. 56, no. 5, pp. 2053–2080, 2010. DOI: 10.1109/TIT.2010.2044061.
- [13] E. J. Candès, X. Li, Y. Ma, and J. Wright, "Robust principal component analysis?" *J. ACM*, vol. 58, no. 3, Jun. 2011, ISSN: 0004-5411. DOI: 10.1145/1970392.1970395.
- [14] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013. DOI: 10.1109/TPAMI.2013.50.
- [15] I. M. Johnstone, "On the distribution of the largest eigenvalue in principal components analysis," *The Annals of Statistics*, vol. 29, no. 2, pp. 295–327, 2001. DOI: 10.1214/aos/1009210544.
- [16] J. Baik, G. Ben-Arous, and S. Péché, "Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices," *The Annals of Probability*, vol. 33, no. 5, pp. 1643–1697, 2005. DOI: 10.1214/009117905000000233.
- [17] J. Baik and J. W. Silverstein, "Eigenvalues of large sample covariance matrices of spiked population models," *Journal of multivariate analysis*, vol. 97, no. 6, pp. 1382–1408, 2006.
- [18] S. Péché, "The largest eigenvalue of small rank perturbations of hermitian random matrices. probab. theory relat. fields 134, 127-173," *Probability Theory and Related Fields*, vol. 134, pp. 127–173, Jan. 2006. DOI: 10.1007/s00440-005-0466-z.
- [19] D. Féral and S. Péché, "The largest eigenvalue of rank one deformation of large Wigner matrices," *Communications in mathematical physics*, vol. 272, no. 1, pp. 185–228, 2007.
- [20] M. Capitaine, C. Donati-Martin, and D. Féral, "The largest eigenvalues of finite rank deformation of large Wigner matrices: Convergence and nonuniversality of the fluctuations," *The Annals of Probability*, vol. 37, no. 1, pp. 1–47, 2009.
- [21] R. R. Nadakuditi and J. W. Silverstein, "Fundamental limit of sample generalized eigenvalue based detection of signals in noise using relatively few signal-bearing and noise-only samples," *IEEE Journal of Selected Topics in Signal Processing*, vol. 4, no. 3, pp. 468–480, 2010. DOI: 10.1109/JSTSP.2009.2038310.
- [22] F. Benaych-Georges and R. R. Nadakuditi, "The eigenvalues and eigenvectors of finite, low rank perturbations of large random matrices," *Advances in Mathematics*, vol. 227, no. 1, pp. 494–521, 2011, ISSN: 0001-8708. DOI: <https://doi.org/10.1016/j.aim.2011.02.007>.
- [23] F. Benaych-Georges and R. R. Nadakuditi, "The singular values and vectors of low rank perturbations of large rectangular random matrices," *Journal of Multivariate Analysis*, vol. 111, pp. 120–135, 2012.
- [24] Z. Bai and J. Yao, "On sample eigenvalues in a generalized spiked population model," *Journal of Multivariate Analysis*, vol. 106, pp. 167–177, 2012.

- [25] T. Lesieur, F. Krzakala, and L. Zdeborová, “Mmse of probabilistic low-rank matrix estimation: Universality with respect to the output channel,” in *2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, 2015, pp. 680–687. DOI: 10.1109/ALLERTON.2015.7447070.
- [26] M. Lelarge and L. Miolane, “Fundamental limits of symmetric low-rank matrix estimation,” *Probability Theory and Related Fields*, vol. 173, pp. 859–929, 2017.
- [27] J. Barbier, N. Macris, and L. Miolane, “The layered structure of tensor estimation and its mutual information,” in *55th Annual Allerton Conference on Communication, Control, and Computing*, 2017.
- [28] J. Barbier and N. Macris, “The adaptive interpolation method for proving replica formulas. applications to the curie–weiss and wigner spike models,” *Journal of Physics A: Mathematical and Theoretical*, vol. 52, no. 29, p. 294 002, 2019. DOI: 10.1088/1751-8121/ab2735.
- [29] J. Barbier and N. Macris, “The adaptive interpolation method: a simple scheme to prove replica formulas in Bayesian inference,” *Probability Theory and Related Fields*, vol. 174, 2019.
- [30] J. Barbier, M. Dia, N. Macris, F. Krzakala, and L. Zdeborová, “Rank-one matrix estimation: analysis of algorithmic and information theoretic limits by the spatial coupling method,” *arXiv e-prints*, arXiv:1812.02537, 2018.
- [31] A. E. Alaoui, F. Krzakala, and M. Jordan, “Fundamental limits of detection in the spiked Wigner model,” *The Annals of Statistics*, vol. 48, no. 2, pp. 863–885, 2020. DOI: 10.1214/19-AOS1826.
- [32] F. Camilli, P. Contucci, and E. Mingione, “Central limit theorem for the overlaps on the nishimori line,” *arXiv preprint arXiv:2305.19943*, 2023.
- [33] J. Barbier, F. Camilli, M. Mondelli, and M. Sáenz, “Fundamental limits in structured principal component analysis and how to reach them,” *Proceedings of the National Academy of Sciences*, vol. 120, no. 30, e2302028120, 2023. DOI: 10.1073/pnas.2302028120.
- [34] M. Mezard, “The space of interactions in neural networks: Gardner’s computation with the cavity method,” *Journal of Physics A: Mathematical and General*, vol. 22, no. 12, p. 2181, Jun. 1989. DOI: 10.1088/0305-4470/22/12/018.
- [35] Y. Kabashima, “A CDMA multiuser detection algorithm on the basis of belief propagation,” *Journal of Physics A: Mathematical and General*, vol. 36, no. 43, pp. 11 111–11 121, Oct. 2003. DOI: 10.1088/0305-4470/36/43/030.
- [36] D. L. Donoho, A. Maleki, and A. Montanari, “Message-passing algorithms for compressed sensing,” *Proceedings of the National Academy of Sciences*, vol. 106, no. 45, pp. 18 914–18 919, 2009. DOI: 10.1073/pnas.0909892106.
- [37] A. Fletcher and S. Rangan, “Iterative reconstruction of rank-one matrices in noise,” *Information and Inference: A Journal of the IMA*, vol. 7, pp. 531–562, Sep. 2018. DOI: 10.1093/imaiai/iax014.
- [38] S. Rangan and A. K. Fletcher, “Iterative estimation of constrained rank-one matrices in noise,” in *2012 IEEE International Symposium on Information Theory Proceedings*, 2012, pp. 1246–1250. DOI: 10.1109/ISIT.2012.6283056.
- [39] D. Voiculescu, “Addition of certain non-commuting random variables,” *Journal of Functional Analysis*, vol. 66, no. 3, pp. 323–346, 1986, ISSN: 0022-1236. DOI: [https://doi.org/10.1016/0022-1236\(86\)90062-5](https://doi.org/10.1016/0022-1236(86)90062-5).
- [40] Y. Kabashima, F. Krzakala, M. Mézard, A. Sakata, and L. Zdeborová, “Phase transitions and sample complexity in bayes-optimal matrix factorization,” *IEEE Transactions on Information Theory*, vol. 62, no. 7, pp. 4228–4265, 2016. DOI: 10.1109/TIT.2016.2556702.
- [41] J. T. Parker, P. Schniter, and V. Cevher, “Bilinear generalized approximate message passing—part i: Derivation,” *IEEE Transactions on Signal Processing*, vol. 62, no. 22, pp. 5839–5853, 2014. DOI: 10.1109/TSP.2014.2357776.
- [42] J. T. Parker, P. Schniter, and V. Cevher, “Bilinear generalized approximate message passing—part ii: Applications,” *IEEE Transactions on Signal Processing*, vol. 62, no. 22, pp. 5854–5867, 2014. DOI: 10.1109/TSP.2014.2357773.
- [43] Q. Zou, H. Zhang, and H. Yang, “Multi-layer bilinear generalized approximate message passing,” *IEEE Transactions on Signal Processing*, vol. 69, pp. 4529–4543, 2021. DOI: 10.1109/TSP.2021.3100305.
- [44] C. Lucibello, F. Pittorino, G. Perugini, and R. Zecchina, “Deep learning via message passing algorithms based on belief propagation,” *Machine Learning: Science and Technology*, vol. 3, no. 3, p. 035 005, Jul. 2022. DOI: 10.1088/2632-2153/ac7d3b.
- [45] H. C. Schmidt, “Statistical Physics of Sparse and Dense Models in Optimization and Inference,” Ph.D. dissertation, Université Paris Saclay (COMUE), Oct. 2018. [Online]. Available: <https://theses.hal.science/tel-03227132>.
- [46] A. Maillard, F. Krzakala, M. Mézard, and L. Zdeborová, “Perturbative construction of mean-field equations in extensive-rank matrix factorization and denoising,” *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2022, no. 8, p. 083 301, Aug. 2022. DOI: 10.1088/1742-5468/ac7e4c.
- [47] J. Barbier and N. Macris, “Statistical limits of dictionary learning: Random matrix theory and the spectral replica method,” *Phys. Rev. E*, vol. 106, p. 024 136, 2 Aug. 2022. DOI: 10.1103/PhysRevE.106.024136.
- [48] F. Camilli and M. Mézard, “Matrix factorization with neural networks,” *Phys. Rev. E*, vol. 107, p. 064 308, 6 Jun. 2023. DOI: 10.1103/PhysRevE.107.064308.

- [49] J. Bun, R. Allez, J.-P. Bouchaud, and M. Potters, “Rotational invariant estimator for general noisy matrices,” *IEEE Transactions on Information Theory*, vol. 62, no. 12, pp. 7475–7490, 2016. DOI: 10.1109/TIT.2016.2616132.
- [50] E. Troiani, V. Erba, F. Krzakala, A. Maillard, and L. Zdeborov’a, “Optimal denoising of rotationally invariant rectangular matrices,” *ArXiv*, vol. abs/2203.07752, 2022.
- [51] F. Pourkamali and N. Macris, “Rectangular rotational invariant estimator for general additive noise matrices,” *arXiv preprint arXiv:2304.12264*, 2023.
- [52] J. J. Hopfield, “Neural networks and physical systems with emergent collective computational abilities.,” *Proceedings of the National Academy of Sciences*, vol. 79, no. 8, pp. 2554–2558, 1982. DOI: 10.1073/pnas.79.8.2554. eprint: <https://www.pnas.org/doi/pdf/10.1073/pnas.79.8.2554>.
- [53] D. J. Amit, H. Gutfreund, and H. Sompolinsky, “Spin-glass models of neural networks,” *Phys. Rev. A*, vol. 32, pp. 1007–1018, 2 Aug. 1985. DOI: 10.1103/PhysRevA.32.1007.
- [54] D. J. Amit, H. Gutfreund, and H. Sompolinsky, “Storing infinite numbers of patterns in a spin-glass model of neural networks,” *Phys. Rev. Lett.*, vol. 55, pp. 1530–1533, 14 Sep. 1985. DOI: 10.1103/PhysRevLett.55.1530.
- [55] M. Mézard, G. Parisi, and M. Virasoro, *Spin Glass Theory and Beyond*. WORLD SCIENTIFIC, 1986. DOI: 10.1142/0271.
- [56] F. Camilli, “New perspectives in statistical mechanics and high-dimensional inference,” Ph.D. dissertation, Alma Mater Studiorum - Università di Bologna and Paris Sciences et Lettres, Mar. 2023. DOI: 10.48676/unibo/amsdottorato/10592.
- [57] F. Pourkamali, J. Barbier, and N. Macris, “Matrix Inference in Growing Rank Regimes,” *arXiv e-prints*, Jun. 2023. arXiv: 2306.01412 [cs.IT].
- [58] J. J. Hopfield, D. I. Feinstein, and R. G. Palmer, “‘unlearning’ has a stabilizing effect in collective memories,” *Nature*, vol. 304, pp. 158–159, 1983. DOI: 10.1038/304158a0.
- [59] J. Van Hemmen, L. Ioffe, R. Kühn, and M. Vaas, “Increasing the efficiency of a neural network through unlearning,” *Physica A: Statistical Mechanics and its Applications*, vol. 163, no. 1, pp. 386–392, 1990, ISSN: 0378-4371. DOI: [https://doi.org/10.1016/0378-4371\(90\)90345-S](https://doi.org/10.1016/0378-4371(90)90345-S).
- [60] M. Benedetti, E. Ventura, E. Marinari, G. Ruocco, and F. Zamponi, “Supervised perceptron learning vs unsupervised hebbian unlearning: Approaching optimal memory retrieval in hopfield-like networks,” *The Journal of Chemical Physics*, vol. 156, no. 10, p. 104107, 2022. DOI: 10.1063/5.0084219.
- [61] V. S. Dotsenko, N. D. Yarunin, and E. A. Dorotheyev, vol. 24, no. 10, p. 2419, May 1991. DOI: 10.1088/0305-4470/24/10/026.
- [62] A. Plakhov and S. Semenov, “The modified unlearning procedure for enhancing storage capacity in hopfield network,” in *[Proceedings] 1992 RNNS/IEEE Symposium on Neuroinformatics and Neurocomputers*, 1992, 242–251 vol.1. DOI: 10.1109/RNNS.1992.268563.
- [63] E. Agliari, F. Alemanno, A. Barra, and A. Fachechi, “Dreaming neural networks: Rigorous results,” *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2019, no. 8, p. 083503, Aug. 2019. DOI: 10.1088/1742-5468/ab371d.
- [64] A. Fachechi, A. Barra, E. Agliari, and F. Alemanno, “Outperforming rbm feature-extraction capabilities by “dreaming” mechanism,” *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–10, 2022. DOI: 10.1109/TNNLS.2022.3182882.
- [65] T. Kohonen, “Self-organization and associative memory,” Springer Berlin, Heidelberg, 1984. DOI: 10.1007/978-3-642-88163-3.
- [66] L. Personnaz, I. Guyon, and G. Dreyfus, “Information storage and retrieval in spin-glass like neural networks,” *Journal de Physique Lettres*, vol. 46, Jan. 1985. DOI: 10.1051/jphyslet:01985004608035900.
- [67] I. Kanter and H. Sompolinsky, “Associative recall of memory without errors,” *Phys. Rev. A*, vol. 35, pp. 380–392, 1 Jan. 1987. DOI: 10.1103/PhysRevA.35.380.
- [68] J. Husson and J. Ko, “Spherical Integrals of Sublinear Rank,” *arXiv e-prints*, Aug. 2022. arXiv: 2208.03642 [math.PR].
- [69] F. Krzakala, J. Xu, and L. Zdeborová, “Mutual information in rank-one matrix estimation,” in *2016 IEEE Information Theory Workshop (ITW)*, 2016, pp. 71–75. DOI: 10.1109/ITW.2016.7606798.
- [70] A. Guionnet, J. Ko, F. Krzakala, and L. Zdeborová, “Low-rank matrix estimation with inhomogeneous noise,” *arXiv preprint arXiv:2208.05918*, 2022.
- [71] A. Guionnet, J. Ko, F. Krzakala, and L. Zdeborová, “Estimating rank-one matrices with mismatched prior and noise: universality and large deviations,” *arXiv e-prints*, Jun. 2023. arXiv: 2306.09283 [math.PR].
- [72] A. Fachechi, E. Agliari, and A. Barra, “Dreaming neural networks: forgetting spurious memories and reinforcing pure ones,” *arXiv e-prints*, Oct. 2018. arXiv: 1810.12217 [cs.NE].