

# Learning probability distributions from smooth observables and the maximum entropy principle: some remarks

Tomoyuki Obuchi<sup>1,†</sup>, Rémi Monasson<sup>2</sup>

<sup>1</sup>Department of Computational Intelligence and Systems Science, Tokyo Institute of Technology, Yokohama, Japan

<sup>2</sup>Laboratoire de Physique Théorique de l'Ecole Normale Supérieure, associé au CNRS et l'Université Pierre et Marie Curie, 75005 Paris, France

E-mail: <sup>†</sup>obuchi@sp.dis.titech.ac.jp

**Abstract.** The maximum entropy principle (MEP) is a very useful working hypothesis in a wide variety of inference problems, ranging from biological to engineering tasks. To better understand the reasons of the success of MEP, we propose a statistical-mechanical formulation to treat the space of probability distributions constrained by the measures of (experimental) observables. In this paper we first review the results of a detailed analysis of the simplest case of randomly chosen observables. In addition, we investigate by numerical and analytical means the case of smooth observables, which is of practical relevance. Our preliminary results are presented and discussed with respect to the efficiency of the MEP.

## 1. Introduction

The maximum entropy (ME) principle, which was proposed by Jaynes in 1957 to rederive statistical mechanics from an information-theoretic viewpoint [1], recently stands out as a useful working hypothesis in a wide variety of inference problems [2, 3, 4, 5, 6, 7, 8, 9, 10]. However, a better comprehension of why and when ME applies is still needed. In this respect, we have recently introduced a statistical-mechanical formulation, in which we treated the space of all probability distributions satisfying a set of moment-matching conditions [11]. Moment-matching conditions were defined through the measured observables, and provided information about the target distribution of interest. It was assumed in that work that the observables were identically and independently distributed random variables, which enabled us to perform analytical investigations and to get a precise quantitative characterization of the volume of the space of probability distributions and of the distance between certain probability distributions and the target one as the number of moment-matching conditions varies. An interesting, but disappointed outcome of the analysis was that the ME distribution was not ‘closer’ to the target distribution than most of the other distributions satisfying the conditions. It was however clear that the failure of the ME principle was due to the unrealistic nature of the observable statistics. In practical applications of the ME, indeed, the observables are not chosen randomly, but are designed to represent some underlying features of the target distribution.

In the present paper, we try to relax the above randomization hypothesis. This strategy leads us to introduce a bias in choosing observables. In particular, we consider the case of ‘smooth’



observables, that is, taking values that cannot differ much between two nearby microscopic configurations of the system. Though we have not been able to extend the detailed analysis of [11] to this (much) more involved case, we report here some analytical and numerical results we have recently obtained.

The organization of this paper is as follows. In the Sec. 2, we present our statistical-physics framework, and discuss different distributions of observables (random and smooth); we also give an overview of our previous results on the random case. In Sec. 3, we report numerical calculations based on a Monte Carlo method in several different situations. In Sec. 4 we report some attempts at solving the problem analytically. Last of all, we discuss our findings and their possible implications for the ME. in the last Section.

## 2. Formalism and overview of results in the random observable case

### 2.1. Basic formulation

Let us consider a system described by configurations of  $N$  binary spins,  $\mathbf{s} = \{s_i = \pm 1\}_{i=1}^N$ , which we call target system. The probability distribution of the target system, the target distribution, is denoted by  $\hat{p}_{\mathbf{s}}$ . Trial distributions, which we will fit based on the measurements of observables (moments). will be generally written as  $p_{\mathbf{s}}$ .

We assume that a measurement gives a value of a linear combination of some moments of the target distribution. Thus we can write a measurement of an observable  $\{v_{\mathbf{s}}\}_{\mathbf{s}}$  as

$$\sum_{\mathbf{s}} v_{\mathbf{s}} \hat{p}_{\mathbf{s}} = \mathbf{v} \cdot \hat{\mathbf{p}}. \quad (1)$$

where the summation runs over all the spin configuration and we introduce a vector notation of the  $2^N$  dimension corresponding to all the spin configurations. The moment-matching (MM) conditions for  $M$  observables  $\{\mathbf{v}^{\mu}\}_{\mu=1}^M$  are then written as

$$\mathbf{v}^{\mu} \cdot \mathbf{p} = \mathbf{v}^{\mu} \cdot \hat{\mathbf{p}}, \quad \forall \mu. \quad (2)$$

Those equations impose that the trial distribution  $\mathbf{p}$  reproduces the measurements of the average values observable  $\mathbf{v}^{\mu}$  over the target distribution.

We are interested in the phase space of trial distributions constrained by the MM conditions (2). To investigate this, we introduce the following probability distribution on trial distributions

$$\rho(\mathbf{p} | \Gamma, E, \{\mathbf{v}^{\mu}\}_{\mu=1}^M, \hat{\mathbf{p}}) = \frac{1}{V} \prod_{\mathbf{s}} \theta(p_{\mathbf{s}}) \delta\left(\sum_{\mathbf{s}} p_{\mathbf{s}} - 1\right) \exp\left\{-\frac{1}{2E} \sum_{\mu=1}^M (\mathbf{v}^{\mu} \cdot (\mathbf{p} - \hat{\mathbf{p}}))^2 + \Gamma S(\mathbf{p})\right\}, \quad (3)$$

where  $S(\mathbf{p})$  is the Shannon entropy whose mathematical expression is

$$S(\mathbf{p}) = - \sum_{\mathbf{s}} p_{\mathbf{s}} \log p_{\mathbf{s}}. \quad (4)$$

The denominator  $V$  is the normalization factor, or the partition function,

$$V(\Gamma, E, \{\mathbf{v}^{\mu}\}_{\mu=1}^M, \hat{\mathbf{p}}) = \int_0^\infty \prod_{\mathbf{s}} dp_{\mathbf{s}} \delta\left(\sum_{\mathbf{s}} p_{\mathbf{s}} - 1\right) \exp\left\{-\frac{1}{2E} \sum_{\mu=1}^M (\mathbf{v}^{\mu} \cdot (\mathbf{p} - \hat{\mathbf{p}}))^2 + \Gamma S(\mathbf{p})\right\}. \quad (5)$$

The factor  $\exp\left(-\frac{1}{2E} \sum_{\mu=1}^M (\mathbf{v}^{\mu} \cdot (\mathbf{p} - \hat{\mathbf{p}}))^2\right)$  corresponds to the MM conditions and  $E$  is the tolerance against the error on the MM conditions. If  $\mathbf{p}$  satisfies the MM conditions, the

corresponding probability value  $\rho(\mathbf{p})$  is large, and small otherwise. Especially in the  $E \rightarrow 0$  limit, trial distributions exactly satisfying all the MM conditions are the only ones to ‘survive’ and to contribute to  $V$ . We call those distributions ‘good’ distributions.

All good distributions have an equal weight value of  $\rho$  at  $\Gamma = 0$ , but for large  $\Gamma \gg 0$  good distributions with larger entropies are assigned larger weights. Hence  $\Gamma$  is a bias towards distributions with large entropies. Studying the two extreme cases  $\Gamma = 0$  and  $\Gamma \rightarrow \infty$ , for which only the ME distribution survives, we can compare typical good distributions and the ME distribution. Actually, we observe some drastic changes for specific intermediate values of  $\Gamma$ , which correspond to phase transitions; hence, the whole range of  $\Gamma > 0$  is of interest.

## 2.2. Distribution of observables

We now need to specify how the observables  $\{\mathbf{v}^\mu\}_\mu$  are chosen. In [11], we assumed the observables were identically independently distributed (i.i.d.) Gaussian random variable

$$P_{\text{rand}}(\mathbf{v}) = \prod_s \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}v_s^2}. \quad (6)$$

The analysis of this random case, which is briefly reported in Section sec. 2.3, yields some important notions and order parameters. However, it is unrealistic as the components of observables are uncorrelated with each other. In particular, two configurations  $\mathbf{s}$  and  $\mathbf{s}'$  differing by a single spin, say,  $s_i$ , have totally uncorrelated observable values,  $v_{\mathbf{s}}$  and  $v_{\mathbf{s}'}$ . This is clearly not true for most systems: changing one among  $N$  spins generically does not affect dramatically the system properties.

An inspiring idea comes from the Fourier transform in the spin-configuration space. Let us consider an arbitrary function of spin configurations  $f_{\mathbf{s}}$  and the Fourier transform characterized by a wavenumber vector of binary components  $\mathbf{k} = \{k_i = 0, 1\}_{i=1}^N$ . The explicit expression of the Fourier transform is

$$f_{\mathbf{k}} = \sum_{\mathbf{s}} \left( \frac{1}{\sqrt{2^N}} \prod_{i=1}^N s_i^{k_i} \right) f_{\mathbf{s}} = \mathbf{w}^{\mathbf{k}} \cdot \mathbf{f}, \quad (7)$$

where we defined the  $\mathbf{s}$ -component of  $\mathbf{w}^{\mathbf{k}}$  to be equal to  $(\prod_{i=1}^N s_i^{k_i})/\sqrt{2^N}$ . The modes  $\{\mathbf{w}^{\mathbf{k}}\}_{\mathbf{k}}$  constitute a complete orthogonal set, and thus any vector in the phase space can be expanded on this set. Each mode has a physical meaning: a mode with a wavenumber  $k = |\mathbf{k}| = \sum_i k_i$  corresponds to a  $k$ -spin multiplet. In particular, if  $f_{\mathbf{s}}$  is a probability distribution,  $f_{\mathbf{k}}$  is the average value of the  $k$ -spin multiplet (divided by  $1/\sqrt{2^N}$ ).

If a random vector drawn from eq. (6) is expanded in the Fourier space, the corresponding weights are statistically uniform on all the modes, which implies that the corresponding MM condition includes correlations of very high-order multiplets of spins with relevant weights. This is not natural because low-order multiplets are often considered to be more important than high-order ones in a general situation. Actually in most of recent applications ranging from biological to engineering problems [2, 3, 4, 5, 6, 7, 8, 9, 10], the ME models are constructed based on observations of singlets and doublets,  $|\mathbf{k}| = 1$  and 2, of spins. Hence, a naive extension from the random case to a more realistic situation is to bias the weights on some (low- $|\mathbf{k}|$ )’s Fourier modes in generating observables.

A simple choice along the above idea is using the following distribution

$$P(\mathbf{v}) = C \exp \left\{ -\frac{1}{2} \sum_{\mathbf{s}, \mathbf{t}} M_{\mathbf{s}\mathbf{t}} v_{\mathbf{s}} v_{\mathbf{t}} \right\}, \quad (8)$$

instead of eq. (6). In the distribution above, the matrix  $M_{st}$  is chosen such to make sure that components of the observable  $\mathbf{v}$  corresponding to nearby configurations have similar values. In other words, the presence of  $M$  should make the observable *smooth* over the configuration space. Alternatively, and perhaps more simply, we may require that the Fourier expression of eq. (8) becomes diagonal and the variance only depends on the absolute value of  $k = |\mathbf{k}|$

$$P(\mathbf{v}) = C \exp \left\{ -\frac{1}{2} \sum_{\mathbf{k}} \frac{v_{\mathbf{k}}^2}{V_{\mathbf{k}}} \right\}. \quad (9)$$

Expression eq. (9) implies that the matrix entry  $M_{st}$  depends on  $\mathbf{s}$  and  $\mathbf{t}$  through the overlap  $\sum_i s_i t_i$  only. If we choose  $V_{\mathbf{k}}$  to be a rapidly-decreasing function of the modulus  $k$ , modes with large  $k$  will be suppressed. There are many different ways to introduce a  $k$ -dependence in  $V_{\mathbf{k}}$ . We try several different expressions for  $V_{\mathbf{k}}$  in the numerical studies reported in Section sec. 3.

### 2.3. Overview of results for random observables

We now briefly recall the results of [11] corresponding to the random case defined in eq. (6). The average value of logarithm of  $V$  in the large-system-size limit can be computed with the replica method (within the replica symmetric (RS) ansatz),

$$F = [\log V]_{\text{rand}} = \lim_{n \rightarrow 0} \frac{1}{n} \log [V^n]_{\text{rand}}, \quad (10)$$

where the square brackets  $[\cdots]_{\text{rand}}$  denote the average over the observables (6). Our analysis showed that  $F$  depends on the target distribution only through the entropy curve

$$\sigma(\omega) = \frac{1}{N} \sum_{\mathbf{s}} \delta(\omega - \omega_{\mathbf{s}}), \quad (11)$$

where we assume the target distribution obeys an exponential scale for  $N \gg 1$ , *i.e.* that

$$\omega_{\mathbf{s}} \equiv -\frac{1}{N} \log \hat{p}_{\mathbf{s}}, \quad (12)$$

has a well-defined limit for  $N \rightarrow \infty$ . This assumption is true for common physical systems. We give a typical shape of  $\sigma(\omega)$  as the left panel of Fig. 1, which corresponds to the independent spin model (ISM) with  $H = 0.5$  as the target distribution

$$\hat{p}_{\mathbf{s}}^{\text{ISM}} = \frac{e^{H \sum_i s_i}}{(2 \cosh H)^N}. \quad (13)$$

The characteristic values of  $\omega$ ,  $\omega_k$  with  $k = 0, 1, 2$ , are defined by

$$\left. \frac{d\sigma(\omega)}{d\omega} \right|_{\omega=\omega_k} = k, \quad (14)$$

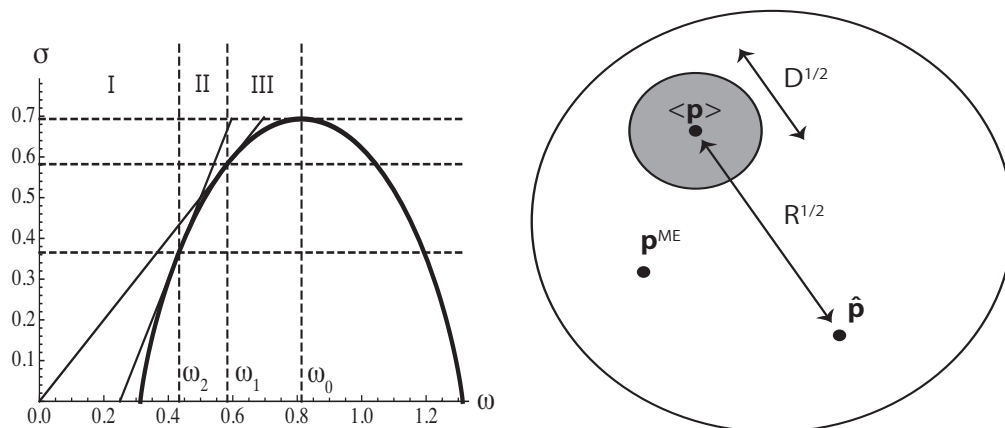
and they are connected to certain phase transitions (see [11] for details of the transitions). These characteristic values of  $\omega$  separate the entropy curve in three regions (the rightmost region is not relevant), each of which is called I, II, and III as shown in Fig. 1.

The analysis in the random case revealed that the following three order parameters are useful

$$Q(\{\mathbf{v}^\mu\}_\mu, \hat{\mathbf{p}}) \equiv \sum_{\mathbf{s}} \left\langle (p_{\mathbf{s}} - \hat{p}_{\mathbf{s}})^2 \right\rangle, \quad (15)$$

$$R(\{\mathbf{v}^\mu\}_\mu, \hat{\mathbf{p}}) \equiv \sum_{\mathbf{s}} (\langle p_{\mathbf{s}} \rangle - \hat{p}_{\mathbf{s}})^2, \quad (16)$$

$$D(\{\mathbf{v}^\mu\}_\mu, \hat{\mathbf{p}}) \equiv Q - R = \sum_{\mathbf{s}} \left( \langle p_{\mathbf{s}}^2 \rangle - \langle p_{\mathbf{s}} \rangle^2 \right), \quad (17)$$



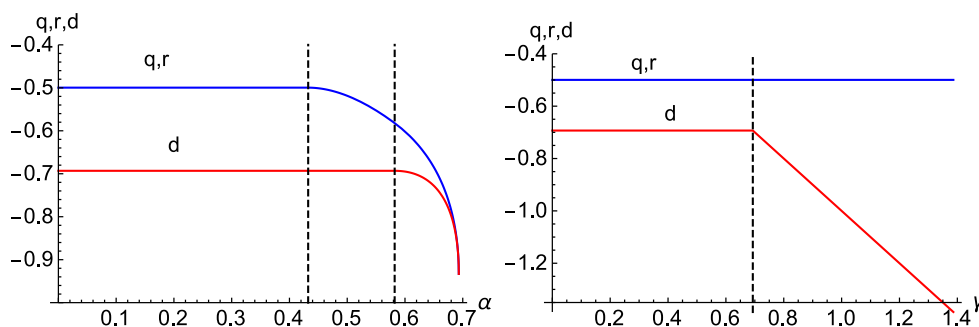
**Figure 1.** (Left) The entropy curve of the ISM with  $H = 0.5$ . There are three characteristic values of  $\omega_k$  with  $k = 0, 1, 2$ , each of which is defined as a tangential point of the slope  $k$  (tangential lines are denoted by solid straight lines). (Right) A schematic diagram of the phase space. The large circle represents the set of all good probability vectors. The shaded area represents the typical fluctuating region of  $\mathbf{p}$ . The target distribution exists apart from  $\langle \mathbf{p} \rangle$  by  $\sqrt{R}$ . The ME distribution  $\mathbf{p}^{\text{ME}}$  also exists somewhere inside the large circle.

where the angular brackets  $\langle \cdots \rangle$  denote the average over  $\rho(\mathbf{p})$ . The physical significance of each order parameter is clear:  $Q$  measures the averaged square distance between  $\mathbf{p}$  and  $\hat{\mathbf{p}}$ .  $R$  also quantifies the distance between the target distribution and the averaged  $\mathbf{p}$  and is similar to  $Q$ , but  $R$ 's meaning is clearer since it does not include the fluctuation of  $\mathbf{p}$ . On the other hand,  $D$  purely measures the fluctuation of  $\mathbf{p}$  around the averaged value. We depict the phase space of good probability vectors with the order parameters in the right panel of Fig. 1.

The control and order parameters scale exponentially with  $N$  in the large  $N$  limit

$$M = e^{N\alpha}, \quad \Gamma = e^{N\gamma}, \quad Q = e^{Nq}, \quad R = e^{Nr}, \quad D = e^{Nd}. \quad (18)$$

Here we plot  $q, r$  and  $d$  of the ISM against  $\alpha$  and  $\gamma$  in Fig. 2. Clearly  $q = r$  holds, which is

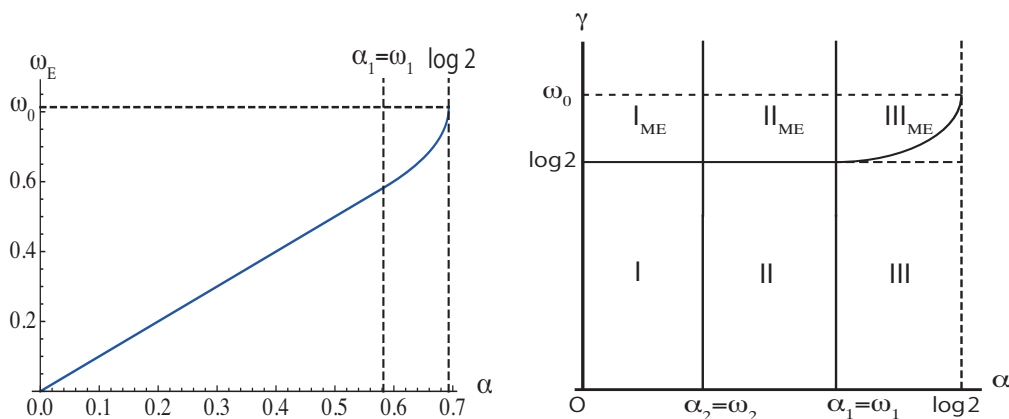


**Figure 2.** Plots of  $q, r$  and  $d$  against  $\alpha$  at  $\Gamma = 0$  (left) and against  $\gamma$  at a small value of  $\alpha$  (right). Dotted vertical lines denote the transition points.

always true in the random case. Our crucial observation about the ME principle is that  $q$  (or  $r$ ) does not depend on  $\gamma$ , that is, on the strength of the entropy bias at all as seen in Fig. 2, which is the case for all values of  $\alpha$ . This means that the distances from the target distribution

to the ME one and to a typical (randomly chosen) ‘good’ distribution are equal. Thus the ME does not help to infer the target distribution. This is the consequence of the random nature of the observables.

An interesting finding in the analysis of the random case was the existence of the so-called ‘learning edge’. The learning edge  $\omega_E$  is the scale separating accurately inferred values of  $\omega_s$  (see eq. (12)) from the ones which cannot be inferred: for target probability values with  $\omega_s \leq \omega_E$  we can infer those values, and for  $\omega_s > \omega_E$  we cannot. The learning edge is a strictly monotonically increasing function of  $\alpha$  (see Fig. 3 as an example for the ISM), whose behaviour reveals the existence of phase transitions. When the learning edge takes values of  $\omega_k$  with  $k = 0, 1, 2$ , the



**Figure 3.** (Left) Plot of the learning edge  $\omega_E$  against  $\alpha$  at  $\Gamma = 0$  for the ISM with  $H = 0.5$ . The dotted lines represent the transition points of the learning edge. At  $\alpha = \log 2$ , all the values of target probabilities are learned, implying the learning edge is not well defined for  $\alpha > \log 2$ . (Right) A schematic phase diagram of the random case. Each phase is named after the location of the learning edge.

phase transitions occur. We call each phase by the location of the learning edge: if  $\omega_E < \omega_2$  then the phase is I, and phases II and III are defined as well (see the left panel of Fig. 1). The resultant phase diagram is given in the right panel of Fig. 3. In the phases subscripted by ME, the order parameter  $d$  is dominated by the ME bias  $\Gamma$  in contrast to the phases without the subscript.

An important finding of the calculation in the random case is that the marginal density of a small number (compared to  $2^N$ ) of probabilities of configurations factorizes a product of single-configuration probability densities  $\rho(\mathbf{p}) \approx \prod_s \rho_s(p_s)$ . Unfortunately, this nice property, at the origin of the existence of the learning edge, does not hold for the case of smooth observable, which makes the analysis much more complicated.

### 3. Numerical studies

#### 3.1. Algorithm for sampling the probability distribution space

We briefly summarize our Monte Carlo (MC) algorithm. We randomly change the probability vector from  $\mathbf{p}$  to  $\mathbf{p}' = \mathbf{p} + \Delta\mathbf{p}$  in the phase space to perform sampling of the space of all distributions. The change of the vector,  $\Delta\mathbf{p}$ , must satisfy the following conditions

**Orthogonality** ( $E = 0$ ):  $\Delta\mathbf{p} \cdot \mathbf{v}^\mu = 0, \forall \mu$ .

**Normalization:**  $\sum_s \Delta p_s = \Delta\mathbf{p} \cdot \mathbf{1} = 0$  where  $\mathbf{1} = (1, 1, \dots, 1)^t$  and the symbol  $^t$  denotes the transpose.

**Positivity:**  $p_s + \Delta p_s \geq 0, \forall s$ .

Combining the orthogonality with the initial condition  $\mathbf{p}(0) = \hat{\mathbf{p}}$ , we see that the vector chosen in this algorithm always satisfies the constraints  $\mathbf{v}^\mu \cdot \mathbf{p} = \mathbf{v}^\mu \cdot \hat{\mathbf{p}}, \forall \mu$ , which is needed since we are basically interested in  $E = 0$ . These orthogonality and normalization conditions restrict the possible move directions. On the other hand, the positivity condition is maintained by appropriately choosing the move width  $x$ . With a chosen move direction  $\mathbf{w}$ , the minimum and maximum values of  $x$  are written as

$$x_{\max} = \max_s \min \left( \frac{-p_s}{\omega_s}, \frac{1-p_s}{w_s} \right), \quad x_{\min} = \min_s \max \left( \frac{-p_s}{\omega_s}, \frac{1-p_s}{w_s} \right). \quad (19)$$

The intermediate values between these two bounds are randomly chosen with equal weight, and thus the trial move becomes  $\Delta \mathbf{p} = x\mathbf{w}$ . Finally, we calculate the entropy difference  $\Delta S = S(\mathbf{p}') - S(\mathbf{p})$  and determines whether the move  $\mathbf{p} \rightarrow \mathbf{p}'$  is accepted or not by the probability

$$p_{\text{accept}} = \min(1, e^{\Gamma \Delta S}). \quad (20)$$

*3.1.1. Observables drawn from the Fourier basis* In the analytical treatment, we put  $\mathbf{v}$  as a random vector drawn from eq. (8), each component of which can take a continuous value. However for numerical implementation, it is more convenient to choose observables from a set of vectors constituting a complete orthogonal set, because the above orthogonal and normalization conditions can be easily satisfied in that way.

We choose the Fourier basis as such a set. The random case corresponds to the situation where all the modes (except for  $\mathbf{w}^0$ ) are drawn uniformly at random. In the studies reported below, we bias the weights of the Fourier modes based on the values of corresponding wavenumber  $k = |\mathbf{k}|$ , see eq. (9).

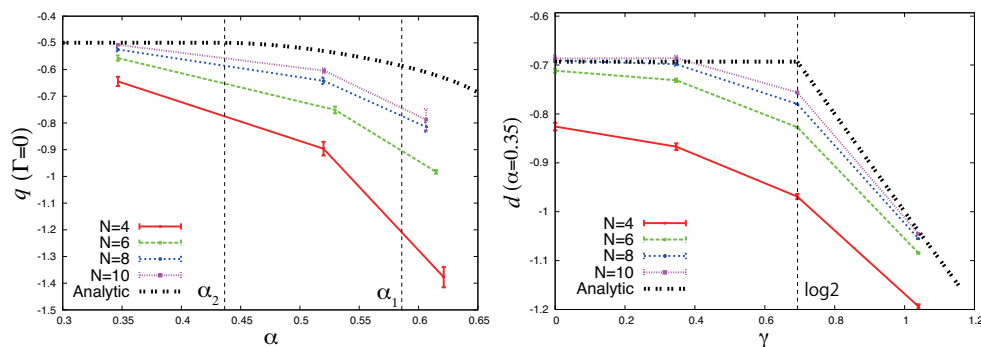
*3.1.2. Lessons from the random case: equilibration and finite-size effects* We define one MC step by one trial of changing the probability vector  $\mathbf{p}$ . As reported in [11], the total MC steps for sufficient equilibration and sampling grow fast as the system size increases. For example for  $N = 10$ , we need  $O(10^8)$  MC steps typically, which requires a couple of days. Due to this expensiveness, we mainly run our simulations for sizes less than or equal to  $N = 10$ .

This system size might sound small, but our numerical simulation on the random case showed that even for  $N = 10$  numerical result could be quantitatively compared to the analytical one. As a demonstration, we display the order parameters of the random case calculated by the MC simulation for finite  $N$  in Fig. 4, which exhibits a clear indication of the phase transitions even for size  $N = 10$ . This result in the random case provides a guide for simulating more general cases. Our simulations given below are based on the equilibration criterion invented in the random case, and we expect that the treated sizes  $N = 8$  or  $10$  capture the behavior in the large-size-system limit.

### 3.2. Exponentially-decaying weight and the smoothness parameter

Here we assume that each mode of the Fourier basis is chosen with a probability weight  $e^{-b|\mathbf{k}|}$ , and we call  $b$  smoothness parameter. The pseudo code to construct the set of observables is as follows:

- (i) Prepare a set of Fourier modes consisting of all the modes except for  $|\mathbf{k}| = 0$ . We call this set “pool”. Similarly one prepares an empty set, called “observables”.
- (ii) Calculate the normalization constant as  $Z = \sum_{\mathbf{k} \in \text{pool}} e^{-b|\mathbf{k}|}$ .
- (iii) Choose one mode of the wavenumber vector  $\mathbf{k}$  from the pool with probability  $p_{\mathbf{k}} = e^{-b|\mathbf{k}|}/Z$  and put the chosen mode into the observables set. Remove the mode from the pool.



**Figure 4.** Numerically-evaluated order parameters  $q$  (left) and  $d$  (right) of the ISM with  $H = 0.5$  plotted against  $\alpha$  (left) and  $\gamma$  (right). Analytical predictions are given by black curves, which are well consistent with the numerical data. Clear indication of the phase transition from small to large  $\Gamma$  is seen in the right panel, and thus finite-size effects seem to be not so large.

(iv) Continue (ii)-(iii) until the size of the observables set becomes equal to  $M$ . At the end, what is left in the pool is our set of possible MC-move directions.

We have to remove  $\mathbf{w}^0 = \mathbf{1}/\sqrt{2^N}$  from the beginning since it is connected to the normalization condition and can neither be a nontrivial observable nor a MC-move direction. The exponentially-decaying weight in (iii) corresponds to putting  $V_{\mathbf{k}}^{-1} = 2b|\mathbf{k}|$  in eq. (9).

As for the target model, we adopt the Boltzmann distribution  $\hat{p}_{\mathbf{s}} = e^{-\beta\mathcal{H}}/Z$  generated from the  $p$ -spin Hamiltonian

$$\mathcal{H} = - \sum_{i_1 < i_2 < \dots < i_p} J_{i_1 i_2 \dots i_p} s_{i_1} s_{i_2} \dots s_{i_p}, \quad (21)$$

where we assume  $J_{i_1 i_2 \dots i_p}$  is an i.i.d. random variable drawn from  $\mathcal{N}\left(0, \frac{p!}{N^{p-1}}\right)$ . Hereafter we call this model  $p$ -spin model (pSM). Parameter  $p$  measures the “smoothness” of the target distribution. In the large- $p$  limit (after the large- $N$  limit), this model converges to the random energy model (REM)<sup>1</sup>. The REM is a completely non-smooth model and we expect the smoothness of observables will not help inference of this model at all. On the other hand, if  $p$  is enough small, we expect that the ME distribution at large values of  $b$  preferentially has lower-spin multiplets in the effective Hamiltonian and matches to the target one, leading to an excellent performance in inference. Hence it is meaningful to compare different- $p$  results for examining the effect of smoothness in the observables.

We below see the MC results with changing the parameters  $N, M, b, p$  and  $\Gamma$ . The quality of inference is quantified by the order parameter  $q$  as introduced in the random case. Error bar is estimated from the standard deviation  $\sigma_{\text{sample}}$  of the objective quantity among different samples as

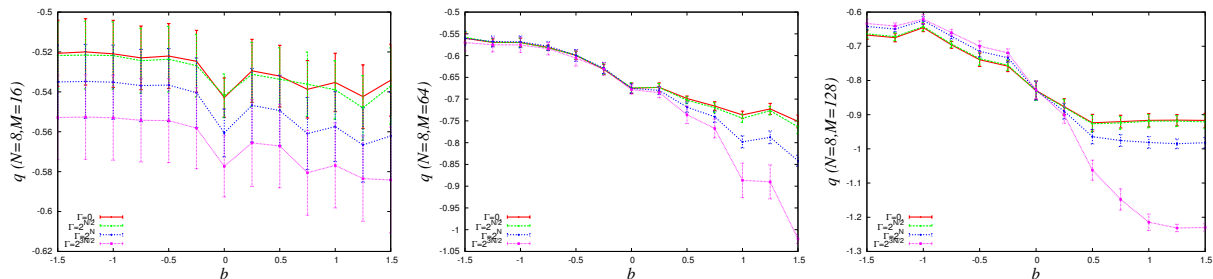
$$\text{Error bar} = \frac{\sigma_{\text{sample}}}{\sqrt{N_{\text{sample}} - 1}}. \quad (22)$$

where  $N_{\text{sample}}$  is the number of simulated samples and are chosen to be 10 in most of the below results.

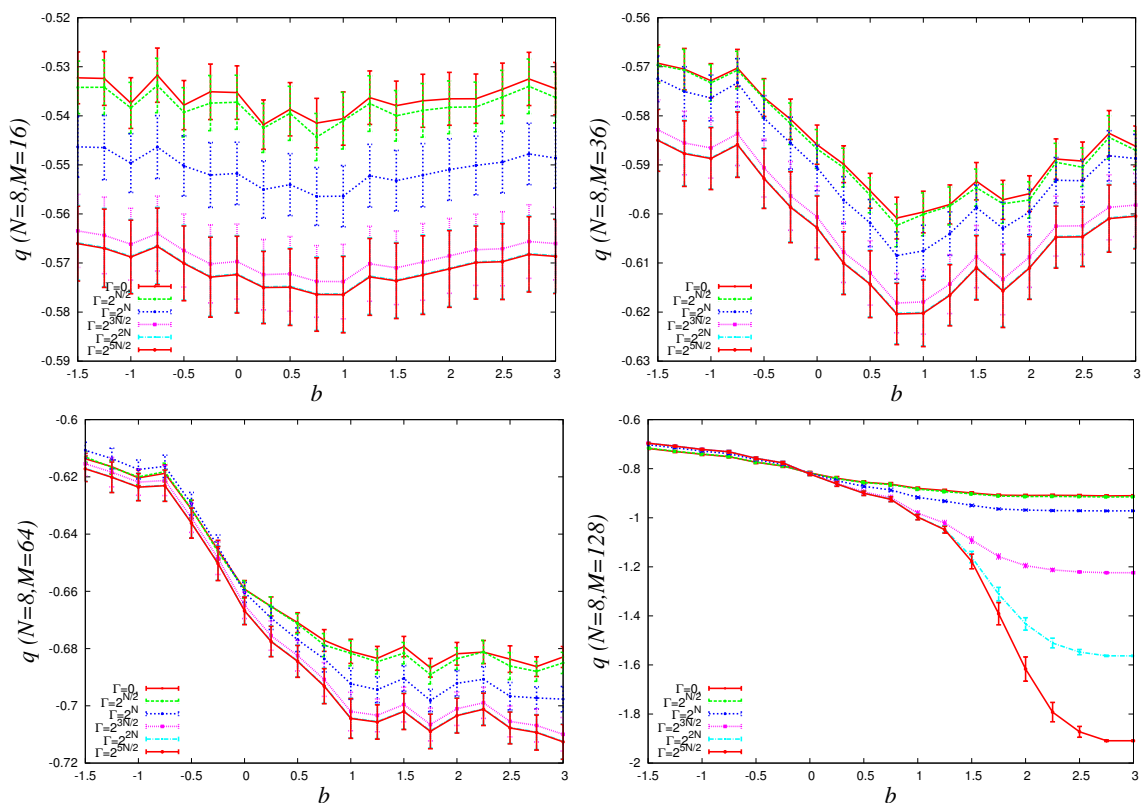
<sup>1</sup> The REM is defined by the Boltzmann distribution whose energies  $\{E_{\mathbf{s}}\}_{\mathbf{s}}$  are drawn from an i.i.d. normal distribution  $\mathcal{N}(0, 2/N)$  (see reprints in [12] for details).



**3.2.1. Result in the smooth case** We plot  $q$  of the 2SM in Fig. 5. We see that the larger  $b$  tends to give smaller  $q$  meaning better inference. Especially for  $M = 64$  and  $128$  with enough large  $b$ , we find that  $q$  becomes smaller and smaller as  $\Gamma$  grows, implying that the ME distribution given in  $\Gamma \rightarrow \infty$  realizes the perfect learning  $q \rightarrow \infty$ . Similar observation is found for the 3SM in Fig.

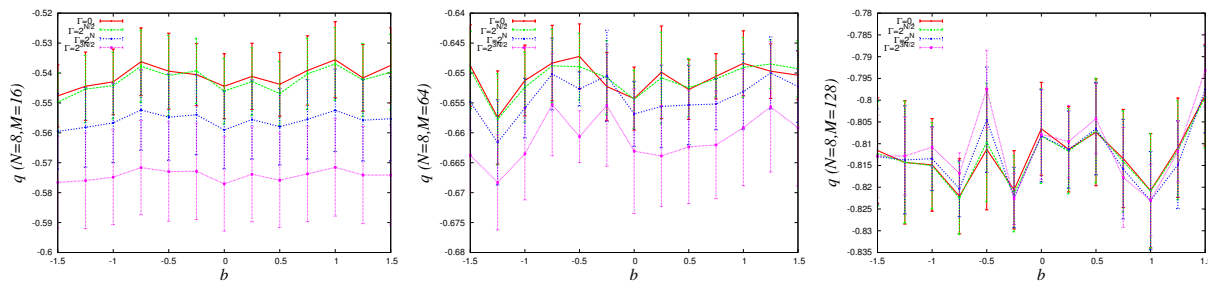


**Figure 5.** Plots of  $q$  for the  $N = 8$  2SM against  $b$ , for  $M = 16, 64$ , and  $128$  from left to right. Different colors correspond to different values of  $\Gamma$ .



**Figure 6.** Plots of  $q$  for the  $N = 8$  3SM vs.  $b$ , for  $M = 16, 36, 64$ , and  $128$  from left to right and from upper to lower. The curve for  $M = 36$  (right upper) shows a dip around  $b = 0.8$ , implying the existence of an optimal value of  $b$ .

6. In this case, we calculate a wider range values of  $b$  and  $\Gamma$ . The results of  $M = 16, 64$ , and  $128$  are very similar to the 2SM case. A characteristic dip is observed in  $M = 36$ , implying the existence of an optimal value of  $b$ , though the difference from other values of  $b$  is not so strong.



**Figure 7.** Plots of  $q$  for the  $N = 8$  REM vs.  $b$ , for  $M = 16, 64$ , and  $128$  from left to right.

Fig. 7 shows the result for the REM. We see no clear dependence of  $q$  on  $b$  for any value of  $M$ , meaning that the smoothness does not help the inference at all, as expected.

The above results for the 2SM, 3SM, and REM clearly exhibit a certain compatibility between the observables and the target distribution. For further examining this, the 3SM result, especially the existence of an optimal value of  $b$ , can be a good starting point. To interpret this result, we introduce  $N_k(\{\mathbf{v}_\mu\}_\mu)$  representing the number of modes with the wavenumber  $k$  in a given set of observables  $\{\mathbf{v}_\mu\}_\mu$ , and define the average occupation ratio of the  $k$ -modes by

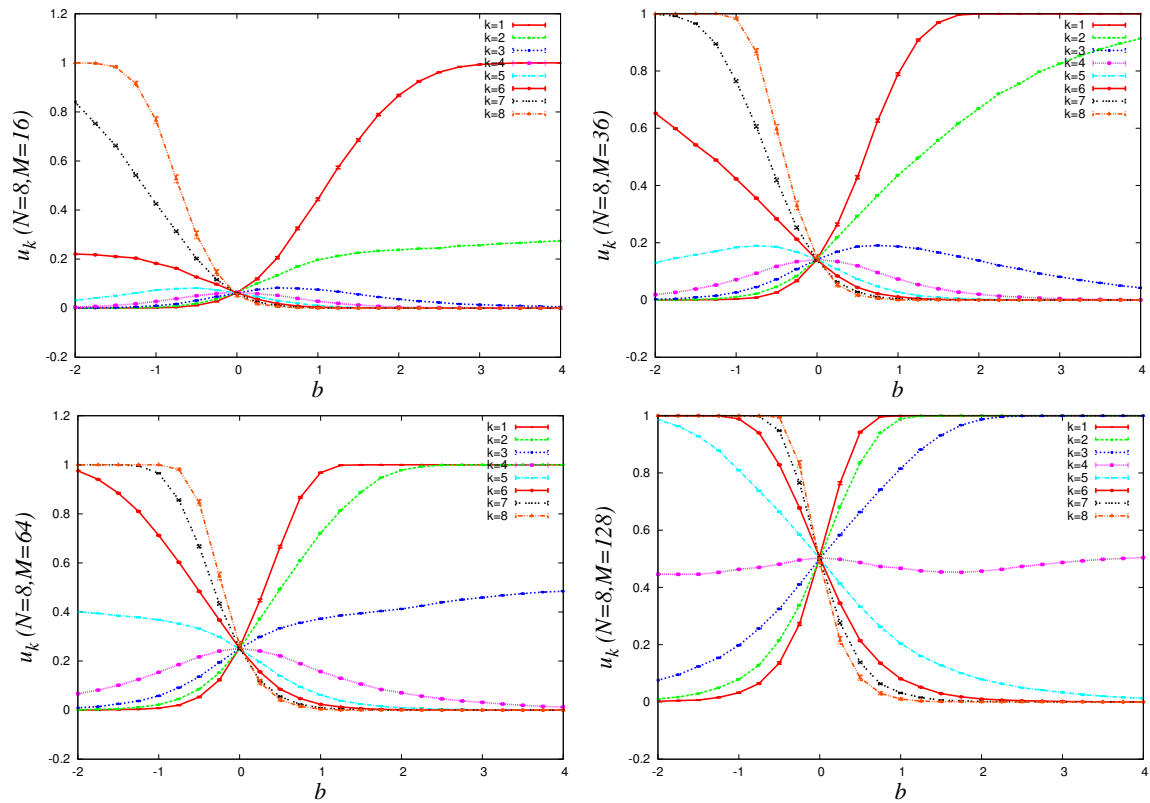
$$u_k = [N_k(\{\mathbf{v}_\mu\}_\mu)]_b \binom{N}{k}^{-1}, \quad (23)$$

where the brackets  $[\dots]_b$  denote the average over the distribution with the smoothness  $b$  stated above. This quantity is easy to calculate because this is independent of the target distribution and moreover the MC simulation is not needed. The occupation ratios for  $M = 16, 36, 64$ , and  $128$  are shown in Fig. 8. A remarkable quantity for examining the 3SM result is  $u_3$  which shows the following behavior

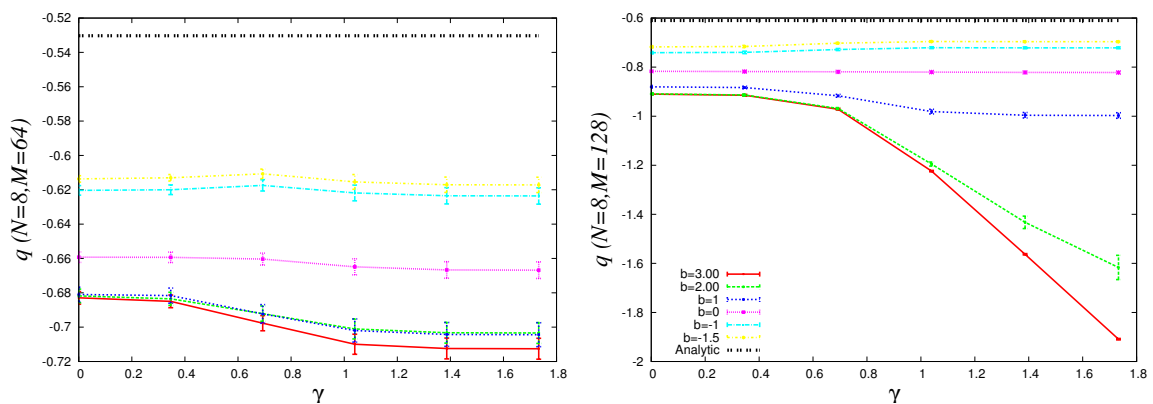
- For  $M = 36$ , it exhibits a clear peak around  $b = 0.8$ .
- For  $M = 64$ , it increases monotonically as  $b$  grows up to  $b = 4$ .
- For  $M = 128$ , it saturates to unity for  $b \gtrsim 2$ .

These observations are probably connected to the characteristic behavior of  $q$ : the dip for  $M = 36$ , the monotonic decrease of  $q$  with  $b$  for  $M = 64$ , and the significant decrease of  $q$  as  $\Gamma$  grows at large values of  $b$  for  $M = 128$ . It is also interesting to see how  $q$  behaves as  $\Gamma$  changes when  $u_3 = 1$  is satisfied, which is demonstrated in the right panel of Fig. 9. For  $b = 3$ , where  $u_3 = 1$ ,  $q$  seems to decrease as a linear function of  $\gamma$  in the large  $\gamma$  region, while at small values of  $b$  the dependence on  $\gamma$  is quite weak and can become positive (decreasing  $q$ ) and negative (increasing  $q$ ). For comparison, we also give the same plot of  $M = 64$  in the left panel, in which no significant decreasing of  $q$  occur.

**3.2.2. Consequences of the smooth-case results.** The above results certainly demonstrate some situations where the ME distribution, and more generally, the introduction of an entropic bias helps to infer the target distribution. However, the interpretation based on  $u_k$  implies the close connection between target distribution and observables. If we know the “smoothness” of a target distribution, we can determine the appropriate observables for applying to the target system, but in practical situations the nature of the target distribution is unclear.



**Figure 8.** The occupation ratio of Fourier modes with wavenumber  $k$  for  $M = 16, 36, 64$ , and  $128$  in  $N = 8$ . Different colors correspond to different  $k$ .

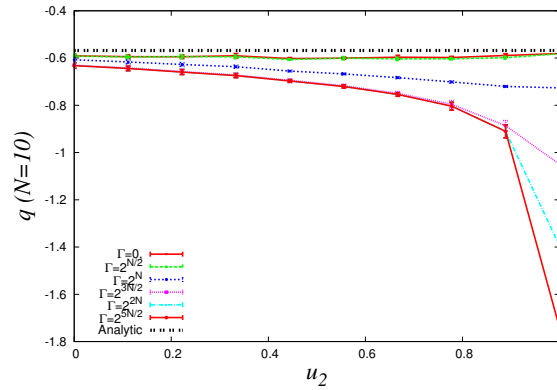


**Figure 9.** Plots of  $q$  for the  $N = 8$  3SM against  $\gamma$ , for  $M = 64$  (left), and  $128$  (right). Black dotted lines denote the analytic solution derived in the random case. For  $M = 128$  at  $b = 3$ ,  $q$  seems to decrease linearly as  $\gamma$  grows at large values of  $\gamma$ .

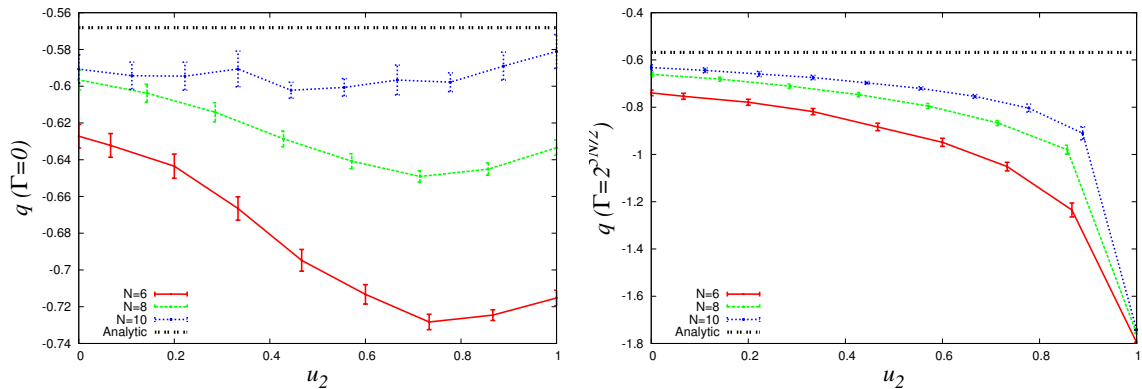
### 3.3. On control of the ratio of “correct” observables.

Here, we examine a quantitative relation between  $q$  and the ratio of “correct” observables. A “correct” observable here is such that the corresponding interaction in the target energy, see eq. (21), is nonzero. For example in the 2SM, only observables with the wavenumber  $k = 2$  are correct. Let us consider the 2SM and  $M = \binom{N}{2}$  observables. We assume that the observables are a mixture only of  $k = 2$  and 3 modes. We treat  $u_2$  as a control parameter on this setup, in contrast to the smooth case where  $u_k$  is the function of  $b$  and  $M$ .

In Fig. 10, we give the plot of  $q$  against  $u_2$  of the  $N = 10$  2SM. To examine the size dependence of  $q$ , we also show the  $u_2$ - $q$  plots for  $N = 6, 8$  and 10 at  $\Gamma = 0$  and  $2^{5N/2}$  in Fig. 11. In Fig.



**Figure 10.** The order parameter  $q$  is plotted against  $u_2$  in the  $N = 10$  2SM. Different colors correspond to different  $\Gamma$  values. At  $u_2=1$ ,  $q$  will go to  $-\infty$  in the limit  $\Gamma \rightarrow \infty$ .



**Figure 11.** The  $u_2$ - $q$  plots at  $\Gamma = 0$  (left) and  $2^{5N/2}$  (right). Different colors represent different system sizes. Black dotted lines denote the analytic solution in the random case.

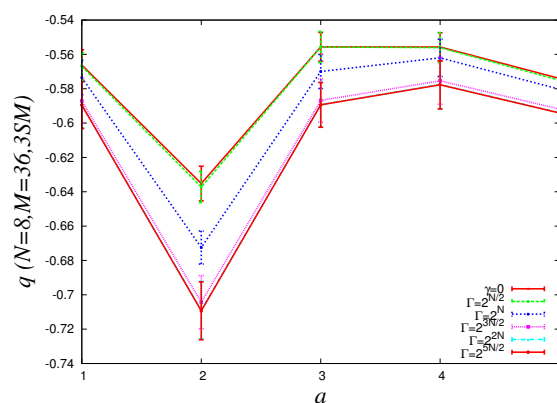
10, we see that  $q$  decreases as  $u_2$  increases at large values of  $\Gamma$ . Especially at  $u_2 = 1$ , the ME distribution appearing in the limit  $\Gamma \rightarrow \infty$  recovers the target distribution and thus  $q$  should go to  $-\infty$  in that limit. This tendency is actually observed in Fig. 10. On the other hand from Fig. 11, at  $\Gamma = 0$  the  $u_2$ - $q$  curves seem to converge to the analytical solution in the random case for any  $u_2$  as the system size grows. A similar behaviour is expected even for large  $\Gamma$  if  $u_2$  is not too large. These observations suggest to focus on  $u_2$  close to, but slightly smaller than 1.

Comparing the size dependence of  $q$  at  $u_2 = 1$  and around  $u_2 = 0.9$  in Fig. 11, we find that the size dependence at  $u_2 = 1$  seems to converge already for  $N = 10$  and  $q$  in the limit  $N \rightarrow \infty$  is expected to be different from the random-case one, while the one around  $u_2 = 0.9$  is still significant and  $q$  seems to converge to the random-case value as  $N$  grows. This observation yields a pessimistic description: in the  $N \rightarrow \infty$  limit, only the  $u_2 = 1$  point is special and for all other values of  $u_2 < 1$  the order parameter  $q$  converges to the random-case value irrespectively of the ME bias  $\Gamma$ .

These results imply that the ME principle does not work even as an approximation, at least in the present scale of the analysis, unless we exactly know what kind of interactions are needed to describe the target distribution. This consequence is only based on a limited result, and hence in the next subsection we reexamine it from another viewpoint.

### 3.4. Comparison of low- and high- $|k|$ modes

As stated above, in realistic situations low- $k$  modes are thought to be more important than high- $k$  modes, and  $k = 1$  and 2 modes are used in most applications. We now examine the possible difference in approximation performance between low- $k$  and high- $k$  modes. We choose 3SM as the target system. Due to computational limitations, we fix the system size to  $N = 8$ . The number of observables is fixed as  $M = 36$  which is equal to  $\binom{N}{1} + \binom{N}{2}$  with  $N = 8$ . Hence, we can include all  $k = 1$  and 2 modes, the common choice in the applications, in a set of observables. Let us label this set by  $a = 1$ . We also define sets of observables labelled by  $a = 2, 3, 4$ , and 5. The sets  $a = 2, 3$ , and 4 consists only of  $k = 3, 4$ , and 5 modes, respectively, and the  $a = 5$  one is a mixture of all the  $k = 6$  and 7 modes. We compare the values of  $q$  among those different sets. Since the target system is the 3SM, the set  $a = 2$  ( $k = 3$ ) gives a better performance than the others, which is trivial. An interesting possibility is whether the results for  $a = 1$  ( $k = 1, 2$ ) are better than the ones for  $a = 3, 4$ , and 5. If so, it could give support to why the ME distribution constructed with  $k = 1$ - and 2-observables works well in many applications. Unfortunately, we do not see any obvious difference in performance between  $a = 1$  and  $a \geq 3$  in Fig. 12.



**Figure 12.** The order parameter  $q$  is plotted against  $u_2$  in the  $N = 10$  2SM. Different colors correspond to different  $\Gamma$  values. At  $u_2=1$ ,  $q$  will go to  $-\infty$  in the limit  $\Gamma \rightarrow \infty$ .

## 4. Analytical treatment of the smooth-observables case: some results

In this section, we provide an analytical solution in some limited cases. We start from stating some assumptions and approximations for making the analysis amenable.

#### 4.1. Assumptions and approximations

4.1.1. *Prior over the target distribution and the annealed approximation.* We have introduced the prior distribution of observables (8) with the matrix  $M_{st}$  which controls the smoothness of observables. Similarly, we here introduce the prior over the target distribution. This was irrelevant in the random case due. In the presence of nontrivial  $M_{st}$ , however, we have to treat correlations among different spin configurations, which requires to specify the properties of the target distribution. The prior distribution of the target distribution is chosen as the following quadratic form

$$\hat{P}(\hat{\mathbf{p}}) = \hat{C} \prod_s \theta(\hat{p}_s) \int_{-\infty}^{\infty} d\hat{\Lambda} e^{\hat{\Lambda} \sum_s (\hat{p}_s - 1) - \frac{1}{2} \sum_{s,t} \hat{M}_{st} \hat{p}_s \hat{p}_t}. \quad (24)$$

The matrix  $\hat{M}_{st}$  controls the smoothness of the target distribution as  $M_{st}$  in eq. (8). The integration variable  $\Lambda$  keeps the normalization condition  $\sum_s \hat{p}_s = 1$ . The factor  $\hat{C}$  is just the normalization.

As eq. (10) in the random case, we calculate the averaged logarithm of the partition function,  $F$ , by using the replica method with the RS ansatz. After some calculations, we get

$$\begin{aligned} \log [V^n]_{\text{smooth}} = \text{E}_{\Omega}^{\text{tr}} \left\{ -\frac{M}{2} \left\{ \log \left( 1 + \frac{Q + (n-1)R}{E} \right) + (n-1) \log \left( 1 + \frac{Q-R}{E} \right) \right\} \right. \\ \left. + \frac{1}{2} n Q(Q' - R') - \frac{1}{2} n(n-1) R R' + \log \int \frac{\sqrt{\det M_{st}}}{\sqrt{2\pi}^{2^N}} \left( \prod_s dz_s \right) e^{-\frac{1}{2} \sum_{s,t} M_{st} z_s z_t} X, \right\} \quad (25) \end{aligned}$$

where the brackets  $[\cdots]_{\text{smooth}}$  denote the average over eqs. (8,24),  $\Omega$  represents the set of order parameters  $\Omega = \{Q, Q', R, R', \Lambda, \Lambda'\}$ , and

$$X = \int_0^\infty \prod_s d\hat{p}_s \hat{P}(\hat{\mathbf{p}}) Y^n, \quad (26)$$

$$\begin{aligned} Y = \int_0^\infty \prod_s dp_s \exp \left\{ \sqrt{R'} \sum_s z_s (p_s - \hat{p}_s) - \Lambda \sum_t (p_t - \hat{p}_t) \right. \\ \left. - \frac{1}{2} Q' \sum_{s,t} M_{st}^{-1} (p_s - \hat{p}_s) (p_t - \hat{p}_t) - \Gamma \sum_s p_s \log p_s \right\}. \quad (27) \end{aligned}$$

Now, we should in principle take the quenched limit  $n \rightarrow 0$ , but unfortunately, the computation in this limit is far from being an easy task. Instead, we work in the annealed case  $n = 1$ , and the resultant formula becomes much simpler

$$\log [V]_{\text{smooth}} = \text{E}_{\Omega'}^{\text{tr}} \left\{ -\frac{M}{2} \log \left( 1 + \frac{Q}{E} \right) + \frac{1}{2} Q Q' + \log X_{\text{anneal}} \right\}, \quad (28)$$

where  $\Omega' = \{Q, Q', \Lambda, \Lambda'\}$

$$\begin{aligned} X_{\text{anneal}} = \int_0^\infty \prod_s d\hat{p}_s \hat{P}(\hat{\mathbf{p}}) \int_0^\infty \prod_s dp_s \exp \left\{ -\Lambda \sum_t (p_t - \hat{p}_t) \right. \\ \left. - \frac{1}{2} Q' \sum_{s,t} M_{st}^{-1} (p_s - \hat{p}_s) (p_t - \hat{p}_t) - \Gamma \sum_s p_s \log p_s \right\}. \quad (29) \end{aligned}$$

4.1.2. *A simple choice of  $M_{st}$  and  $\hat{M}_{st}$ .* As noted above,  $M_{st}$  should only depend on the overlap  $\sum_i s_i t_i$ . Here we consider one of the simplest choice obeying this constraint:

$$M_{st}^{-1} = \frac{e^{B \sum_i s_i t_i}}{(2 \cosh B)^N}. \quad (30)$$

The parameter  $B \in \mathbb{R}$  controls the smoothness of observables and the random case is reproduced in the limit  $B \rightarrow \infty$ . Note that  $B$  is defined from the inverse of  $M_{st}$ . The matrix  $M_{st}$  itself can be written in the same form, but the parameter corresponding  $B$  becomes complex in general.

How about the matrix  $\hat{M}_{st}$ ? We do not want to bias a specific spin and thus  $\hat{M}_{st}$  should have the permutation symmetry of the spin index, meaning that  $\hat{M}_{st}$  is a function of the overlap  $\sum_i s_i t_i$  and the two magnetizations  $\sum_i s_i$  and  $\sum_i t_i$  only. The symmetry between the two magnetizations should be also maintained. Here we choose

$$\hat{M}_{st} = \hat{C} e^{\hat{B} \sum_i s_i t_i + \hat{H}(s_i + t_i)}. \quad (31)$$

As seen below, this functional form leads to the ISM as the target distribution.

4.1.3. *A saddle-point approximation of the integration over  $\mathbf{p}$ .* A notable point of the present formulation is that all the variables obey the exponential scale. Thanks to this, we can approximate the integration over  $\mathbf{p}$  by a saddle-point (SP) technique. Actually, this approximation gives the exact result, which was demonstrated in [11] in the random case. We expect the same is true in the present case. Hence, instead of integrating out over  $\mathbf{p}$  and  $\hat{\mathbf{p}}$ , we take the SP condition on the exponent of  $X_{\text{anneal}}$ . The result is

$$-\Lambda - Q' \sum_t M_{st}^{-1} p_t + Q' \sum_t M_{st}^{-1} \hat{p}_t - \Gamma(\log p_s + 1) = 0, \quad (32)$$

$$\hat{\Lambda} + Q' \sum_t M_{st}^{-1} (p_t - \hat{p}_t) - \sum_t \hat{M}_{st}^{-1} \hat{p}_t = 0, \quad (33)$$

Taking the extremization condition with respect to  $Q$  in eq. (28), we see  $Q'$  is proportional to  $M$ . Hence, for small  $M$ , we find that the SP solution of  $\hat{\mathbf{p}}$  gives back the ISM

$$\hat{p}_s^* = \hat{\Lambda} \sum_t \hat{M}_{st} = \frac{e^{H^* \sum_i s_i}}{(2 \cosh H^*)^N}, \quad (34)$$

where  $H^*$  is a function of  $\hat{B}$  and  $\hat{H}$  in eq. (31). We expect the smallness of the second term in eq. (33) to hold for general  $M$ . Therefore, we hereafter choose the ISM as the target distribution. This is natural since the target distribution should not be changed as  $M$  increases. Thanks to this assumption, we can concentrate on solving only eq. (32) for  $\hat{\mathbf{p}}$  defined by the ISM, see eq. (13).

Note that variable  $p_s$  is a probability value and cannot be negative, a fact ignored in deriving eq. (32). We give the name ‘‘SAT’’ to the set of spin configurations  $\mathbf{s}$  whose probability values derived from the SP equation becomes positive, and ‘‘UNS’’ to the other configurations. Thus the correct SP equation is

$$\begin{cases} -\Lambda - Q' \sum_{t \in \text{SAT}} M_{st}^{-1} p_t + Q' \sum_t M_{st}^{-1} \hat{p}_t - \Gamma(\log p_s + 1) = 0, & (\mathbf{s} \in \text{SAT}) \\ p_s = f_s(\Gamma), & (\mathbf{s} \in \text{UNS}) \end{cases} \quad (35)$$

This is the extension of the modified SP approximation used in [11]. For UNS, the function  $f_s(\Gamma)$  is expected to be a simple function of  $\Gamma$ . For example in the random case,  $f_s(\Gamma) = 0$  for  $0 \leq \Gamma < 2^N$  and  $f_s(\Gamma) = 1/2^N$  for  $\Gamma \geq 2^N$ .

The difficulty in solving this set of coupled equations is the classification of the SAT and UNS configurations. In the random case  $B \rightarrow \infty$ , target probabilities are learned in descending order of their values, and the structure in the spin-configuration space is completely irrelevant: the classification of the SAT and UNS configurations is reduced to the self-consistent determination of the boundary value, namely, the learning edge, between the learned, large probabilities and the unlearned, small ones. However in the present case, due to the presence of the nontrivial matrix  $M_{st}$ , we have to treat the spin-configuration space itself, and the classification of the SAT and UNS configurations is very complicated.

In the following subsection, we assume that the learning edge can be uniquely determined for finite but still enough large  $B$ , which enables us to conduct the analysis. The applicable range of this assumption is also discussed.

#### 4.2. Perturbation from large $B$ at $\Gamma = 0$

In the discussion below, we put  $\Gamma = 0$  in eq. (76) and use the following truncation with a fixed  $\Xi \in \mathbb{N}$

$$\sum_t M_{st}^{-1}(\dots) \approx \sum_{\{t|Q_{st}=N, N-2, \dots, N-2\Xi\}} M_{st}^{-1}(\dots), \quad (36)$$

where  $Q_{st} = \sum_i s_i t_i$  and  $\Xi$  is the truncation range. Let us write the solution of the SP equation as

$$p_s^* = \hat{p}_s - \frac{\Lambda}{Q'} + \Phi_s. \quad (37)$$

For  $s \in \text{SAT}$  the residual term  $\Phi_s$  is zero at the random case and we expect  $\Phi$  is small for enough large  $B$ , while for  $s \in \text{UNS}$   $\Phi_s = \Lambda/Q' - \hat{p}_s$  since  $p_s^* = 0$ .

4.2.1. *The  $\Xi = 1$  case.* Let us put  $T_B = e^{BN}/(2 \cosh B)^N$ . Remembering that our target distribution is the ISM and depends only on the magnetization  $M_s = \sum_i s_i = Nm_s$ , we get from the third term in the upper line of eq. (76)

$$\sum_{\{t|Q_{st}=N, N-2\}} M_{st}^{-1} \hat{p}_t = T_B \{ \hat{p}_{M_s} + e^{-2B} (D_s \hat{p}_{M_s+2} + U_s \hat{p}_{M_s-2}) \}, \quad (38)$$

where  $U_s(D_s)$  is the number of up (down) spins of the configuration  $s$ . The variables  $M_s, U_s$ , and  $D_s$  are extensive and proportional to  $N$ . Hence, to validate the above truncation, we need

$$\epsilon \equiv Ne^{-2B} \ll 1. \quad (39)$$

Inserting the solution form eq. (37) into eq. (76), we get

$$\epsilon \frac{(1 - m_s)}{2} \Phi_{M_s+2} + \Phi_{M_s} + \epsilon \frac{(1 + m_s)}{2} \Phi_{M_s-2} = \frac{\Lambda}{Q'} \frac{1 - T_B(1 + \epsilon)}{T_B}. \quad (40)$$

The term in the righthand side appears due to the modification from the normalization

$$0 = 1 - \sum_t M_{st}^{-1} \approx 1 - T_B(1 + \epsilon). \quad (41)$$

So, we can neglect the righthand side of eq. (40).



*$M_s$ -dependent solution.* As  $\Phi = \text{const}$  is an inappropriate solution to eq. (40), we look for a  $M_s$ -dependent solution. To lighten notations, we change the index of  $\hat{p}$  from the magnetization  $M_s$  into the number of up spins,  $U_s$ , and carry out the same change of notations on the other variables. The equation we want to solve is now

$$\epsilon(1 - u_s)\Phi_{U_s+1} + \Phi_{U_s} + \epsilon u_s \Phi_{U_s-1} = 0, \quad (42)$$

where  $u_s = U_s/N$ . The general solution of this equation is not simple due to the  $u_s$ -dependence of the coefficients. Thus we work on the large  $N$  limit and assume the scale  $\Phi(U_s) = e^{N\phi(u_s)}$ . Correspondingly,

$$\Phi_{U_s \pm 1} = e^{N\phi(u_s \pm 1/N)} \approx e^{N\phi(u_s)} e^{\pm \frac{\partial \phi}{\partial u} \big|_{u=u_s}}. \quad (43)$$

The SP equation becomes quadratic of  $y = e^{\frac{\partial \phi}{\partial u} \big|_{u=u_s}}$  and the solution is

$$y = \frac{-1 \pm \sqrt{1 - 4\epsilon^2 u_s(1 - u_s)}}{2\epsilon(1 - u_s)}. \quad (44)$$

In the limit  $\epsilon \rightarrow 0$ , we expect  $\Phi(U_s) = 0$  if  $U_s$ , implying  $y = 0$ . Thus we need to take the plus sign in eq. (44). Using this solution and integrating  $\log y$  with respect to  $u$ , we get

$$\begin{aligned} \phi(u) = \phi_0 + u \log \left\{ \frac{-1 + \sqrt{1 - 4\epsilon^2 u(1 - u)}}{2\epsilon(1 - u)} \right\} - \frac{1}{2\epsilon} \log \left\{ -\epsilon + 2\epsilon u + \sqrt{1 - 4\epsilon^2 u(1 - u)} \right\} \\ + \frac{1}{2} \log \left\{ 1 - 2\epsilon^2 + 2u\epsilon^2 + \sqrt{1 - 4\epsilon^2 u(1 - u)} \right\} \equiv \phi_0 + \phi_c(u). \end{aligned} \quad (45)$$

Since our target distribution is the ISM, there is the edge value of  $U$ ,  $U_E = Nu_E$ , corresponding to the learning edge  $\omega_E$ . The constant  $\Phi_0 = e^{N\phi_0}$  is determined from the boundary condition at  $U_s = U_E$ . We can put  $\Phi_{U_E-1} = \hat{p}_{U_E-1} - \frac{\Lambda}{Q'}$ , hence from the SP equation at  $U_s = U_E$ , we get

$$\Phi(u) = -y(u_E) \left( \hat{p}_{U_E-1} - \frac{\Lambda}{Q'} \right) e^{N(\phi_c(u) - \phi_c(u_E))}, \quad (46)$$

where we used a simple relation

$$\frac{2\epsilon u_E}{1 + \sqrt{1 - 4\epsilon^2 u_E(1 - u_E)}} = -y(u_E). \quad (47)$$

The corresponding  $p_U^*$  thus becomes

$$p_U^* = \hat{p}_U - \frac{\Lambda}{Q'} - y(u_E) \left( \hat{p}_{U_E-1} - \frac{\Lambda}{Q'} \right) e^{N(\phi_c(u) - \phi_c(u_E))}. \quad (48)$$

It is easy to confirm that  $\text{sign}(y(u)) = -1$  and  $0 < |y(u)| < 1$  for any  $0 < \epsilon < 1$  and  $0 < u < 1$ . Hence, eq. (48) has a good interpretation. The factor  $y(u_E)e^{N(\phi_c(u) - \phi_c(u_E))}$  controls the changing speed when we change  $U$  around  $U_E$ . Namely,

$$p_{U_E+K}^* \approx \hat{p}_{U_E+K} - \frac{\Lambda}{Q'} + (\text{sign}(y(u_E)))^{K+1} |y(u_E)| e^{K \log |y(u_E)|} \left( \frac{\Lambda}{Q'} - \hat{p}_{U_E-1} \right). \quad (49)$$

If the increasing speed  $e^{K \log |y(u_E)|}$  is slower than the one of  $\hat{p}_{U_E+K} = \hat{p}_{U_E} e^{-\frac{\partial \omega}{\partial u} K}$ , we can see all  $U \geq U_E$  are fine and the corresponding  $\{p_U^*\}_{U \geq U_E}$  take positive values, then the learning

edge is well-defined. This is the case for  $\epsilon < 1$  since  $\log |y|$  is always negative and smaller than  $-\frac{\partial \omega}{\partial u} = \frac{H^*}{2} > 0$ . Thus our analysis is consistent with the assumptions.

We should clarify the magnitude relation between  $\Lambda/Q'$  and  $\hat{p}_{U_E}$ . We expect the righthand side of eq. (49) to be positive if  $U \geq U_E$ , and zero or negative for  $U \leq U_E - 1$ . The definition of  $y$  leads to a relation  $e^{N\{\phi_c(u_E-1/N)-\phi_c(u_E)\}} \approx y^{-1}(u_E)$ , which simplifies the inequalities for  $U = U_E$  and  $U = U_E - 1$ :

$$\hat{p}_{U_E} \geq \frac{\Lambda}{Q'} \{1 - y(u_E)\} + y(u_E) \hat{p}_{U_E-1}, \quad (50)$$

$$\hat{p}_{U_E-1} \leq \frac{\Lambda}{Q'} \{1 - y(u_E)y^{-1}(u_E)\} + y(u_E)y^{-1}(u_E) \hat{p}_{U_E-1} = \hat{p}_{U_E-1}. \quad (51)$$

Hence the SP equation at  $U_s = U_E - 1$  is marginally satisfied, implying  $p_{U_E-2}^* < 0$ . These inequalities imply

$$\hat{p}_{U_E-2} < \frac{\Lambda}{Q'} < \hat{p}_{U_E}. \quad (52)$$

*Effective number of observables.* In the usual situation, we calculate the learning edge  $\omega_E$  for a given number of observables  $M$ , but here we perform the opposite analysis, namely we calculate the effective number of observables  $M_{\text{eff}}$  for a given  $\omega_E$ . If  $M_{\text{eff}}(\omega_E)$  is smaller than the one in the random case with the same value of  $\omega_E$ , it means that the inference improves as a result of the introduction of the nontrivial matrix  $M_{st}$ .

From the extremization conditions with respect to  $Q$  and  $Q'$  in eq. (28) with negligibly small  $E$ , we get

$$Q = \left\langle \sum_{s,t} M_{st}^{-1} (p_s - \hat{p}_s)(p_t - \hat{p}_t) \right\rangle, \quad (53)$$

$$Q'Q = M. \quad (54)$$

The summation over  $\sum_{s,t}$  in the first equation can be divided into three categories

$$\sum_{s,t} (\dots) = \left( \sum_{s,t \in \text{SAT}} + \sum_{s,t \in \text{UNS}} + 2 \sum_{s \in \text{SAT}, t \in \text{UNS}} \right) (\dots). \quad (55)$$

Let us assume the average over  $\rho$ ,  $\langle \dots \rangle$ , can be replaced by the SP values of  $\mathbf{p}$ . Then we get

$$\sum_{s,t \in \text{SAT}} M_{st}^{-1} (p_s^* - \hat{p}_s)(p_t^* - \hat{p}_t) = \sum_{s,t \in \text{SAT}} M_{st}^{-1} \left( \Phi_s - \frac{\Lambda}{Q'} \right) \left( \Phi_t - \frac{\Lambda}{Q'} \right) \doteq e^{N\hat{\sigma}} \left( \frac{\Lambda}{Q'} \right)^2, \quad (56)$$

where the symbol  $\doteq$  denotes the equality in the exponential scale. To derive this, we assume  $\Lambda/Q' > |\Phi_s|$  which is true in the random case and is strongly expected in the present case too. Similarly, we have

$$\sum_{s,t \in \text{UNS}} M_{st}^{-1} (p_s^* - \hat{p}_s)(p_t^* - \hat{p}_t) \approx \sum_{s,t \in \text{UNS}} M_{st}^{-1} \hat{p}_s \hat{p}_t \doteq \sum_{U_s < U_E} \hat{p}_{U_s}^2, \quad (57)$$

$$\sum_{s \in \text{SAT}, t \in \text{UNS}} M_{st}^{-1} (p_s^* - \hat{p}_s)(p_t^* - \hat{p}_t) \doteq e^{N\sigma_E} \left( \frac{\Lambda}{Q'} \right)^2, \quad (58)$$

where  $\sigma_E = \sigma(\omega_E)$ . Collecting and comparing all these terms, we see the dominant term is

$$Q \doteq \sum_{\mathbf{s} \in \text{UNS}} \hat{p}_{\mathbf{s}}^2 \doteq \begin{cases} e^{N(\sigma(\omega_2) - 2\omega_2)} & (\omega_E \leq \omega_2) \\ e^{N(\sigma_E - 2\omega_E)} & (\omega_E > \omega_2) \end{cases} = e^{Nq(\omega_E)}. \quad (59)$$

We expect  $\Lambda = 2^N$  in an interesting situation. Hence, we can calculate  $M_{\text{eff}}$  from a given value of  $\Lambda/Q' = \mathcal{L}$  and from eq. (54) as

$$M_{\text{eff}} = \frac{2^N}{\mathcal{L}} Q \doteq e^{N(\log 2 - \omega_E + q(\omega_E))}. \quad (60)$$

The last inequality comes from  $\mathcal{L} \doteq e^{-N\omega_E}$  derived from eq. (52). This effective  $M_{\text{eff}}$  is identical to the one in the random case for a given  $\omega_E$ , implying that the inference is not improved at least in the exponential scale, unfortunately.

*4.2.2. General  $\Xi = O(1)$ .* Next, we proceed to the case of larger values of  $\epsilon$ . Now, we need higher orders in the truncation. We assume  $\Xi$  is a  $O(1)$  constant and is not extensive. As in the  $\Xi = 1$  case, the SP equation is expressed as

$$\sum_{i=0}^{\Xi} \epsilon^i \sum_{j=0}^i \frac{u_{\mathbf{s}}^j (1 - u_{\mathbf{s}})^{i-j}}{j! (i-j)!} \Phi_{U_{\mathbf{s}} + i - 2j} = 0. \quad (61)$$

To derive this, we neglect sub-leading factors in the large-size limit as

$$\binom{U_{\mathbf{s}}}{j} \approx N^j \frac{u_{\mathbf{s}}^j}{j!}. \quad (62)$$

This truncation is justified if eq. (39) is satisfied. Fortunately, we can solve this equation for general  $\Xi = O(1)$ . We state the outline of deriving the solution.

First we need to solve the characteristic equation

$$\sum_{i=0}^{\Xi} \epsilon^i \sum_{j=0}^i \frac{u_{\mathbf{s}}^j (1 - u_{\mathbf{s}})^{i-j}}{j! (i-j)!} y(u_{\mathbf{s}})^{i-2j} = 0. \quad (63)$$

This has  $2\Xi$  different roots. The appropriate roots are the ones going to zero as  $\epsilon \rightarrow 0$ . As far as we searched, the appropriate roots consist of some real-valued ones and some complex-conjugate pairs. For example in  $\Xi = 3$ , we found that the appropriate roots are constituted by one real-valued root  $y_1$  and one pair of complex roots  $y_2$  and  $y_3$  which are in the complex conjugate relation. Using these roots, we can write

$$\Phi(u) = \sum_{i=1}^{\Xi} A_i e^{N\phi_i(u)}, \quad (64)$$

where

$$\phi_i(u) = \int^u dx \log y_i(x). \quad (65)$$

Now, the coefficients  $\{A_i\}$  are determined from the boundary conditions  $U_{\mathbf{s}} = U_E, \dots, U_{E+\Xi-1}$ , leading to

$$\frac{\Lambda}{Q'} - \hat{p}_{U_E - \xi} = \Phi(U_E - \xi) \approx \sum_{i=1}^R A_i y_i^{-\xi} e^{N\phi_i(u_E)} \quad (\xi = 1, \dots, \Xi). \quad (66)$$

Thus means  $p_{U_E-\xi}^* = 0$ , ( $\xi = 1, \dots, \Xi$ ) is maintained by the boundary conditions. From the boundary conditions, we also find the magnitude relation between  $\Lambda/Q'$  and  $\hat{p}_{U_E}$

$$\hat{p}_{U_E-\Xi-1} \leq \frac{\Lambda}{Q'} \leq \hat{p}_{U_E}. \quad (67)$$

These are the construction of the solution of eq. (61).

We confirmed up to a reasonable value of  $\Xi$  that the real part of the suitable  $\{y_i(u)\}$  are decreasing function with respect to  $u$  for a certain moderate range of  $\epsilon > 0$ , and thus  $\phi_i(u)$  too. Hence,  $\Phi(U_E + K)$  is expected to decrease as  $K$  grows by the factor  $e^{K \log |y_i(u_E)|}$ , and the learning edge is well-defined.

The effective number of observables can be also evaluated in this case. Unfortunately, as long as  $\Xi = O(1)$ ,  $M_{\text{eff}}(\omega_E)$  cannot be different in the exponential scale from  $M(\omega_E)$  in the random case, since both  $Q$  and  $\mathcal{L} = \Lambda/Q' \doteq e^{-N\omega_E}$  are essentially the same as the ones of the random case. This naturally motivates us to treat  $\Xi = O(N)$ , but that is not easy since it becomes even nontrivial to write down the SP equation. We leave this as a future work. Instead in the next subsection, we try to examine the effect of the ME bias  $\Gamma$  in the same setting.

### 4.3. Perturbation from large $B$ at finite $\Gamma$

*4.3.1. Review of the random case.* We start by reviewing the random case with relatively small  $M$ . The SP equation is

$$p_s = \hat{p}_s - \frac{\Lambda}{Q'} - \frac{\Gamma}{Q'} (1 + \log p_s). \quad (68)$$

For large  $\Gamma \geq 2^N$ , the order parameter  $\Lambda$  behaves as

$$\Lambda = (N \log 2 - 1)\Gamma, \quad (69)$$

and the solution of the SP equation is

$$p_s^* = \begin{cases} \hat{p}_s - \frac{\Gamma \log 2^N \hat{p}_s}{Q'}, & (s \in \text{SAT}) \\ p_0 + \Delta_s, & (s \in \text{UNS}) \end{cases}, \quad (70)$$

where

$$p_0 = e^{-1-\frac{\Lambda}{\Gamma}} = \frac{1}{2^N}, \quad \Delta_s = \frac{Q' p_0}{Q' p_0 + \Gamma} (\hat{p}_s - p_0). \quad (71)$$

The expression of  $\Delta_s$  can be derived from the perturbation with respect to  $\Delta$  in eq. (68). Inequalities

$$p_0 > \Delta_s, \quad \Gamma > Q' \hat{p}_s \quad (\forall s \in \text{UNS}) \quad (72)$$

are needed to justify this perturbation. To make the above expressions consistent, the learning edge should be

$$\omega_E = q' - \gamma = q' - \lambda. \quad (73)$$

For  $\Gamma < 2^N \doteq \Lambda(\Gamma = 0)$ , the situation is essentially identical to the case  $\Gamma = 0$ . This can be understood by comparing each term of the righthand side in eq. (68). For  $s \in \text{SAT}$  the dominant term is the first term  $\hat{p}_s$  and for  $s \in \text{UNS}$  the second one is dominating. The last term proportional to  $\Gamma/Q'$  is always smaller than the other two in the exponential scale and hence essentially does not change the SP equation.

*4.3.2. General  $\Xi = O(1)$ .* A lesson from sec. 4.2.2 tells that the introduction of the nontrivial  $M_{st}$  does not change the scale of the problem if  $\Xi = O(1)$ , and another lesson from sec. 4.3.1 says that the effect of the ME bias  $\Gamma$  appears when  $\Gamma$  becomes comparable with  $\Lambda$  in the scale. These two lessons imply that the ME bias and  $M_{st}$  cannot provide a meaningful effect on the inference in the present setting.

As an example, let us see the SP equation for the truncation range  $\Xi = 1$  with  $\Gamma$  which can be written as

$$\Phi(u_s) (\epsilon(1 - u_s) + 1 + \epsilon u_s) = -\frac{\Gamma}{Q'} (1 + \log p_s), \quad (74)$$

We found that  $\Phi(u_s) = p_{u_s}^* - \hat{p}_{u_s} + \frac{\Lambda}{Q'}$  is the decreasing function of  $u$ , and the leading-order terms in  $\Phi$  are  $p_s^* \doteq \hat{p}_s$  for  $s \in \text{SAT}$  and  $\Lambda/Q'$  for  $s \in \text{UNS}$ . The scale of  $\Phi$  itself is smaller than  $p_s^* \doteq \hat{p}_s$  for  $s \in \text{SAT}$  due to the cancellation in  $p_s^* - \hat{p}_s$ . If  $\Gamma/Q'$  is smaller than any of these two terms, the situation will be identical to  $\Gamma = 0$  as the random case with  $\Gamma < 2^N$ . This is consistent to the leading order of eq. (74). As  $\Gamma$  grows and  $\Gamma \doteq \Lambda(\Gamma = 0) \doteq 2^N$ , the situation is changed:  $\Lambda$  scales as  $\Gamma$  for  $\Gamma > 2^N$  and the learning edge determined by the balance between  $\Lambda/Q'$  and  $\hat{p}_s$  starts to decrease since  $Q'$  is kept to be constant against  $\Gamma$ . Thus introducing the ME bias worsens the inference performances, as found with the annealed approximation in the random case.

This situation seems to remain the same as long as  $\Xi = O(1)$  where the coefficient of  $\Phi$  is  $O(1)$ . Hence, from the annealed calculations, we cannot find any good effect on inference produced by the ME bias.

*4.3.3. Some comments on the quenched case.* A crucial difference between the quenched and annealed cases is the behavior of  $Q'$ . The extremization condition over  $Q'$  in eq. (25) after taking the  $n \rightarrow 0$  limit becomes

$$Q' = \frac{M}{D}, \quad (75)$$

and  $D$  is scaled as  $D = \Gamma^{-1}$  for  $\Gamma > 2^N$  in the random case. Hence,  $Q'$  increases with the growth of  $\Gamma$  as  $\Lambda$ . This is good for the ME principle since the learning edge  $\omega_E = q' - \lambda$  does not change as  $\Gamma$  grows.

On the other hand, the SP condition with respect to  $\mathbf{p}$  becomes

$$\begin{cases} -\Lambda - Q' \sum_{t \in \text{SAT}} M_{st}^{-1} p_t + Q' \sum_t M_{st}^{-1} \hat{p}_t + \frac{\sqrt{R'}}{Q'} z_s - \Gamma(\log p_s + 1) = 0, & (s \in \text{SAT}) \\ p_s = f_s(\Gamma), & (s \in \text{UNS}) \end{cases} \quad (76)$$

The only difference with the annealed case is the presence of the term  $\frac{\sqrt{R'}}{Q'} z_s$  and we expect that the nature of the SP equation does not change so much.

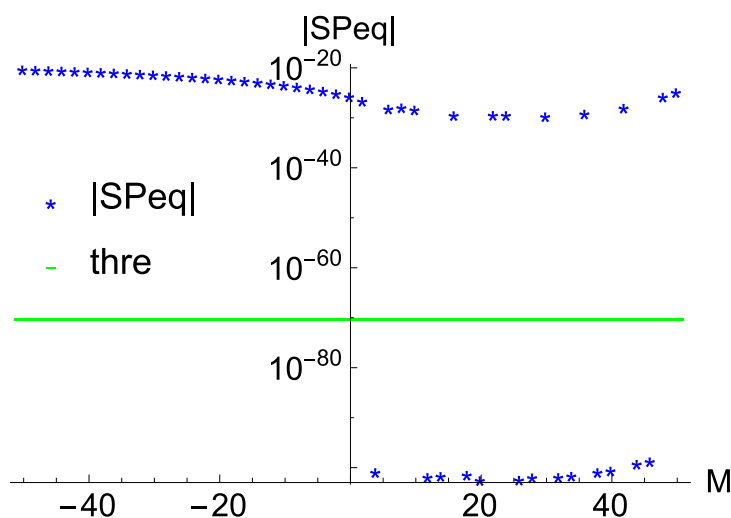
#### 4.4. A semi-analytical treatment: numerical evaluation of the SP equation

In the above subsections, we have concentrated on studying the truncated SP equation. For general  $B$ ,  $\epsilon = Ne^{-2B}$  can be large so the truncation cannot be justified. To examine this case, we numerically evaluate (76). Actually, directly evaluation of (76) is not easy, and instead of that, we minimize the following cost function

$$L = \frac{1}{2} \sum_{s,t} M_{st}^{-1} (p_s - \hat{p}_s)(p_t - \hat{p}_t) + \frac{\Lambda}{Q'} \sum_t (p_t - \hat{p}_t) + \frac{\Gamma}{Q'} \sum_s p_s \log p_s, \quad (77)$$

for given  $\Lambda/Q'$  and  $\Gamma/Q'$  under the condition  $p_s \geq 0$ . This is a convex optimization problem with constraints and can be solved by certain known algorithms such as the interior-point method. By this minimization, we can obtain the solution of the SP equation,  $\mathbf{p}^*$ , and thus related quantities can be evaluated. In solving this, we assumed the symmetry that  $p_s^*$  having the same magnetization value share the same probability value as well as  $\hat{p}_s$ , which reduces a lot of computational costs.

An interesting finding might be the absence of the unique learning edge, which is demonstrated in Fig. 13. In the figure, the absolute values of the lefthand side of eq. (32)



**Figure 13.** The absolute values of the lefthand side of eq. (32) of the upper line are plotted against the magnetization  $M_s = \sum_i s_i$ . The target model is the ISM with  $H = 0.5$ , and other parameters are  $N = 50, B = 0.2, H = 0.5, \Gamma = 0$  and  $\Lambda/Q' = 5.0 \times 10^{-20}$ . If  $\mathbf{s} \in \text{SAT}$  (UNS), the corresponding value denoted by a blue point becomes smaller (larger) than the chosen threshold value  $10^{-70}$  denoted by the green line. We see that the blue points locating below the threshold value are distributed in a patchy fashion, implying the absence of the learning edge.

of the upper line are plotted against the magnetization with the parameters:  $N = 50, B = 0.2, H = 0.5, \Gamma = 0$  and  $\Lambda/Q' = 5.0 \times 10^{-20}$ . Namely, if  $\mathbf{s} \in \text{SAT}$ , the corresponding value should be enough small (here we set the threshold value as  $10^{-70}$ ). In the random case  $B \rightarrow \infty$ , there exists an edge value of  $M$ ,  $M_E$  corresponding to the learning edge  $\omega_E$ , and all the SP-equation values for  $M \geq M_E$  are smaller than the threshold, while the ones for  $M < M_E$  are larger. In the present case  $B = 0.2$ , however, we see that the blue points locating below the threshold value are distributed in a patchy fashion and we cannot find such a clear edge value for  $M$ . This results implies that no learning edge can be defined, which makes the resolution of the SP equation directly (76) very difficult, as pointed out in sec. 4.1.3. Note that this calculation requires a very multiple precision because our variables obey the exponential scale and can be very large and small. For  $N = 50$ , we used 170 digits to express a real number.

## 5. Discussion and summary

In the present work, we have investigated the space of probability distributions constrained by smooth observables. We have first introduced a Monte Carlo procedure to sample the space of probability distribution over the set of configurations made of a small (in practice  $\leq 10$ ) number of spins. In this numerical study, we have used the order parameter  $q$  measuring the

distance between the target distribution and the averaged one to quantify the quality of inference. Our numerical investigations have revealed that the ME principle can work better than other typical distributions, in particular when the target distribution and the observables have comparable smoothnesses. As a consequence, once we know the nature of the interactions present in the target distribution, we can construct the ME distribution precisely describing the target system. However, our numerical simulations also show that a few errors in specifying the type of interactions leads to a significant worsening of the ME distribution. This suggests that the ME principle is not good in approximating a target distribution. In practical situations, it does not seem easy to precisely know the type of the interactions of the target system, and compatibility is a serious issue.

To get deeper insight, we have also carried out some analytical investigations. We have assigned a specific functional form to the matrix controlling the smoothness of the observables and introduced a prior distribution over the target distribution. In this setting, we have analyzed the perturbative regime from the random case, combining this approach with the annealed approximation and the saddle-point approximation. We were able to construct a solution in this case, but the solution was not informative. The annealed approximation, indeed, corresponds to ‘bad’ choices of the constraints, making the volume  $V$  increase compared to its typical value for quenched, random observables.

One of the main outcomes of the present work is that, in the presence of smooth observables, no learning edge, separating small from large target probabilities cannot be defined, see Fig. 13, as used to be the case in the non-smooth case [11]. This result highlights the complexity of learning in the smooth case, and its collective nature. While the probabilities of the different spin configuration were essentially learned independently of each other in the non-smooth case, this is clearly not the case anymore. How to describe in a mathematically controlled and precise way this collective phenomenon remains a challenge. We hope that the remarks presented here will be useful for further progresses in this direction.

## Acknowledgments

T. O. is grateful to E. Aurell for a fruitful discussion at KITPC in China. He also acknowledges the support by Grant-in-Aid for JSPS Fellows, as well as the JSPS Core-to-Core Program “Non-equilibrium dynamics of soft matter and information”. A part of numerical calculations were carried out on the TSUBAME2.5 supercomputer in the Tokyo Institute of Technology. R.M. was partly funded by the Agence Nationale de la Recherche Coevstat project (ANR-13-BS04-0012-01) and the [EU]-FP7 FET OPEN project Enlightenment 284801.

## References

- [1] Jaynes E T 1957 *Phys. Rev.* **106** 620
- [2] Shneidman E, Berry M J, Segev R, Bialek W 2006 *Nature* **440** 1007
- [3] Tkačik G, Marre O, Mora T, Amodei D, Berry II M J, Bialek W 2013 *J. Stat. Mech.* P03011
- [4] Cocco S, Leibler S, Monasson R 2009 *Proc. Nat. Acad. Sci. USA* **106** 14058
- [5] Seno F, Trovato A, Banavar J R, Maritan A 2008 *Phys. Rev. Lett.* **100** 078102
- [6] Weigt M, White R A, Szurmant H, Hoch J A, Hwa T 2009 *Proc. Nat. Acad. Sci. USA* **106** 67
- [7] Mora T, Walczak A M, Bialek W, Callan C G 2010 *Proc. Nat. Acad. Sci. USA* **107** 5405
- [8] Cocco S, Monasson R, Weigt M 2013 *PLOS Comput. Biol.* **9** e1003176; *J. Phys. Conference Series* **473** 012010
- [9] Bialek W, Cavagna A, Giardinà I, Mora T, Silvestri, E, Viale M, Walczak A 2012 *Proc. Nat. Acad. Sci. USA* **109** 4786
- [10] Kataoka S, Yasuda M, Furtlehner C, Tanaka K 2014 *Inverse Problems* **30** 025003
- [11] Obuchi T, Monasson R, and Cocco S, 2015 Learning probabilities from random observables in high dimensions: the maximum entropy distribution and others, submitted to *J. Stat. Phys.*
- [12] Mézard M, Parisi G, Virasoro M A 1987 *Spin Glass Theory and Beyond* (Singapore: World Scientific)
- [13] Bishop C M 2006 *Pattern Recognition and Machine Learning* (New York: Springer)