

# Parameter priors for Ising models

## research notes

Y. Roudi

<yasser.roudi@ntnu.no>

P.G.L. Porta Mana

<piero.mana@ntnu.no>

25 June 2018; updated 27 June 2018

Study of uniform priors in parameter space and in constraint space for Ising models

'Flat priors do not exist'  
(anonymous)

## 1 A two-unit model with sufficient statistics

Consider a population of two binary units  $s := (s_1, s_2)$  with values in  $\{0, 1\}$ . One observation of this population can thus give four results:  $s \in \{00, 01, 10, 11\}$ .

Assume that we have  $N$  observations  $(s^{(1)}, \dots, s^{(N)})$  of this or other populations prepared in similar conditions, so that knowledge of these observations is relevant for our forecast of a new observation  $s$ , again in similar conditions. Also assume that only the number, the mean, and the second moments of these past observations are relevant to forecast the new one; that is,

$$N, \quad \frac{1}{N}(s^{(1)} + \dots + s^{(N)}) =: \bar{s}, \quad \frac{1}{N}(s_1^{(1)} s_2^{(1)} + \dots + s_1^{(N)} s_2^{(N)}) =: \overline{s s} \quad (1)$$

are sufficient statistics. These assumptions are collectively denoted  $I$ .

There is a series of mathematical results, which we call the Koopman-Pitman-Lauritzen theorem, that says that our probabilistic forecasts must

assume this general form, for any  $N$ :

$$\begin{aligned} p(\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)} | I) &= \int \left[ \prod_{i=1}^N g(\mathbf{s}^{(i)}) \frac{\exp(\mu_1 s_1^{(i)} + \mu_2 s_2^{(i)} + \lambda s_1^{(i)} s_2^{(i)})}{Z(\mu_1, \mu_2, \lambda)} \right] \times \\ &\quad p(\mu_1, \mu_2, \lambda | I) d\mu_1 d\mu_2 d\lambda, \\ &= \int \left[ \prod_{i=1}^N g(\mathbf{s}^{(i)}) \right] \frac{\exp[N(\mu_1 \bar{s}_1 + \mu_2 \bar{s}_2 + \lambda \bar{s} \bar{s})]}{Z(\mu_1, \mu_2, \lambda)^N} \times \\ &\quad p(\mu_1, \mu_2, \lambda | I) d\mu_1 d\mu_2 d\lambda, \end{aligned}$$

$$\text{with } Z(\mu_1, \mu_2, \lambda) := 1 + \exp(\mu_1) + \exp(\mu_2) + \exp(\mu_1 + \mu_2 + \lambda). \quad (2)$$

Denote the three parameters that appear in this formula by  $\theta := (\mu_1, \mu_2, \lambda) \in \mathbf{R}^3$ .

The distribution  $g(\mathbf{s})$  and the density  $p(\mu_1, \mu_2, \lambda | I)$  in the formula above are not determined by the theorem: they need to be determined by additional assumptions. The distribution  $g$  is often determined by symmetry or combinatorial properties of the problem. From now on we assume it to be unity:  $g(\mathbf{s}) = 1$ . The density  $p(\theta | I)$  is called *prior parameter density*.

Formula (2) may appear deceptively specific in its dependence on the parameters  $\theta$ ; let's summarize in words the content of the theorem:

(a) our joint probability for the observations  $\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)}$  is given by a convex combination of joint probabilities:

$$p(\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)} | \theta, I) = \prod_{i=1}^N \frac{\exp(\mu_1 s_1^{(i)} + \mu_2 s_2^{(i)} + \lambda s_1^{(i)} s_2^{(i)})}{Z(\theta)}; \quad (3)$$

- (b) each joint probability in this convex combination factorizes into  $N$  independent probabilities for the  $N$  observations, as is clear from the formula above;
- (c) each joint probability in the convex combination is identified by a triplet of parameters  $\theta$ ; it therefore belongs to a three-dimensional submanifold of joint probabilities. Note that the full manifold of joint probabilities is  $(2^N - 1)$ -dimensional;
- (d) the weight assigned to the probability labelled by  $\theta$  is  $p(\theta | I) d\theta$ .

Point (c) shows that the Pitman-Koopman-Lauritzen theorem greatly reduces our freedom in specifying the joint probability. This is the

effect of assuming that the statistics (1) are sufficient; it's a very strong assumption.

Points (c) and (d) show that the theorem selects a particular three-dimensional submanifold within each of the  $(2^N - 1)$ -dimensional manifolds of probability distributions for  $N$  observations, for all  $N$ . But the theorem doesn't select any particular coordinate system within the submanifold: the parameters  $\theta$  are just coordinates, and there is nothing special about them, besides the fact that they appear as coefficients of the linear combination of statistics in the exponential (2). We could choose different coordinates  $t$ , with one-one coordinate transformations  $t = t(\theta)$ ,  $\theta = \theta(t)$ . In these new coordinates the mixed joint probabilities are

$$p(s^{(1)}, \dots, s^{(N)} | t, I) = \prod_{i=1}^N \frac{\exp[\mu_1(t) s_1^{(i)} + \mu_2(t) s_2^{(i)} + \lambda(t) s_1^{(i)} s_2^{(i)}]}{Z[\theta(t)]}; \quad (4)$$

and the weights of the convex combination are given by  $p(t | I) dt$ , the densities for  $\theta$  and for  $t$  being related by a Jacobian determinant:

$$p(\theta | I) = p[t(\theta) | I] \det\left(\frac{\partial t}{\partial \theta}\right). \quad (5)$$

This coordinate change will be central in the rest of this study.

## 2 New coordinates and their motivation

Assuming that (1) are sufficient statistics and therefore using formula (2), let's ask what's the limit probability of observing particular values of the statistics  $\bar{s}$ ,  $\bar{s}\bar{s}$  for very large  $N$ ; that is,  $p(\bar{s}, \bar{s}\bar{s} | I, \text{large } N)$ .

In this section we show that there is a particular coordinate system  $t := (m_1, m_2, I)$  of the three-dimensional manifold discussed above for which *the prior parameter density coincides, in the large- $N$  limit, with the probability of the observed statistics*:

$$p[(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s}) = x | I, \text{large } N] \approx p(t = x | I). \quad (6)$$

To see this, consider the parameterized, factorized joint probability  $p(s^{(1)}, \dots, s^{(N)} | \theta, I)$  of eq. (12). The expectation of the statistics

$(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s})$  is given by

$$E[(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s}) | \boldsymbol{\theta}, I] = \frac{1}{N} \sum_i E[(s_1^{(i)}, s_2^{(i)}, s_1^{(i)} s_2^{(i)}) | \boldsymbol{\theta}, I] = E[(s_1, s_2, s_1 s_2) | \boldsymbol{\theta}, I] \quad (7)$$

where  $(s_1, s_2, s_1 s_2)$  refer to any one of the  $N$  observations. The two equalities come from the properties of the expectation and the factorized form of the joint probability conditional on  $\boldsymbol{\theta}$ . From the properties of the variance we also have

$$V[(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s}) | \boldsymbol{\theta}, I] = \frac{1}{N} V[(s_1, s_2, s_1 s_2) | \boldsymbol{\theta}, I]. \quad (8)$$

This means that for a triplet  $\boldsymbol{\theta}$ , for large  $N$  we have a probability distribution for the statistics that is very peaked at particular values  $\mathbf{t} := (m_1, m_2, l)$  determined by the equations

$$\begin{aligned} m_1 = E(s_1 | \boldsymbol{\theta}, I) &\equiv \frac{\partial \ln Z(\boldsymbol{\theta})}{\partial \mu_1}, & m_2 = E(s_2 | \boldsymbol{\theta}, I) &\equiv \frac{\partial \ln Z(\boldsymbol{\theta})}{\partial \mu_2}, \\ l = E(s_1 s_2 | \boldsymbol{\theta}, I) &\equiv \frac{\partial \ln Z(\boldsymbol{\theta})}{\partial \lambda}. \end{aligned} \quad (9)$$

This system of equations actually puts the parameters  $\boldsymbol{\theta} := (\mu_1, \mu_2, \lambda)$  and  $\mathbf{t} := (m_1, m_2, l)$  into one-one correspondence (Mead et al. 1984). The former belong to  $\mathbf{R}^3$ ; the latter to the bounded domain

$$0 \leq m_1, m_2 \leq 1, \quad \max(0, m_1 + m_2 - 1) \leq l \leq \min(m_1, m_2) \quad (10)$$

shown in fig. 1.

Using these new parameters, the probability for statistics  $(\bar{s}, \bar{s}\bar{s})$  becomes for large  $N$

$$p(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s} | m_1, m_2, l, I) \approx \delta[(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s}) - (m_1, m_2, l)]. \quad (11)$$

Taking the convex combination of this expression in  $\mathbf{t}$  with weights  $p(\mathbf{t} | I) d\mathbf{t}$  we obtain eq. (6).

In the coordinates  $\mathbf{t}$ , the formula (2) given by the theorem can be interpreted in the following way.

- (1) We first assume to know that the limit statistics in a very large number of observations is  $\mathbf{t} := (m_1, m_2, l)$ . Given this knowledge

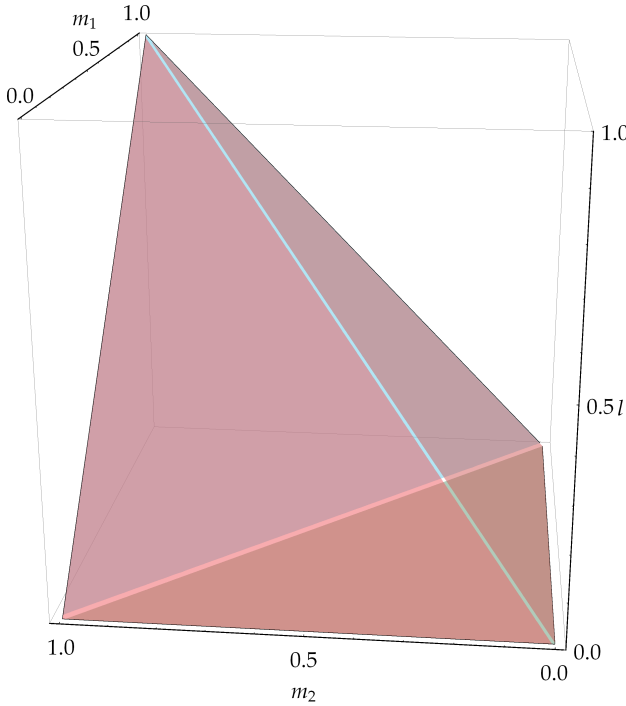


Figure 1

we can combinatorially calculate the probability of observing a finite sequence of  $N$  observations,  $p(s^{(1)}, \dots, s^{(N)} | t, I)$ , assuming that all sequences having given statistics are equally likely – this equiprobability corresponds to setting  $g(s) = 1$  in eq. (2). An example of this combinatorial calculation is given below.

- (2) We then express our uncertainty about the limit statistics with the probability density  $p(t | I) dt$ .
- (3) The two uncertainties above are finally combined in the usual way using the law of total probability.

Let's show, in the simplest case, that the probability conditional on the statistics  $t$ ,

$$p(s | t, I) = \frac{\exp[\mu_1(t) s_1 + \mu_2(t) s_2 + \lambda(t) s_1 s_2]}{Z[\theta(t)]} \quad (12)$$

is indeed given combinatorially assuming equiprobability of all sequences, as claimed above. Suppose that we know the limit statistics are  $(m_1, m_2, l)$ . Only the outcome  $s = 11$  gives a non-vanishing contribution to the second moment  $l$ , eq. (1). This number is therefore equal to the limit relative frequency of 11. Assuming equiprobability we set the probability of this outcome in the next observation equal to this frequency  $l$ . Only the outcomes 10 and 11 give non-vanishing contributions to the mean  $m_1$ ; this number is therefore equal to their joint limit relative frequencies. Since the frequency of 11 is given by  $l$ , the frequency of 10 must be given by  $m_1 - l$ , which is then our probability for this outcome in the next observation. Analogous reasoning holds for the outcome 01. Finally, the limit relative frequencies of all four outcomes must sum to 1; thus the limit frequency and probability of outcome 00 must be  $1 - l - (m_1 - l) - (m_2 - l)$ . Summarizing,

$$p(s|t, I) = \begin{cases} 1 + l - m_1 - m_2 & \text{for } s = 00 \\ m_1 - l & \text{for } s = 10 \\ m_2 - l & \text{for } s = 01 \\ l & \text{for } s = 11 \end{cases}$$

$$\equiv l^{s_1 s_2} (m_1 - l)^{s_1 - s_1 s_2} (m_2 - l)^{s_2 - s_1 s_2} (1 + l - m_1 - m_2)^{1 - s_1 - s_2 + s_1 s_2}. \quad (13)$$

This probability distribution is exactly eq. (12), as can be checked by finding  $\theta(t)$  with the inverse of the coordinate transformations (9),

$$\mu_1 = \ln \frac{m_1 - l}{1 + l - m_1 - m_2}, \quad \mu_2 = \ln \frac{m_2 - l}{1 + l - m_1 - m_2},$$

$$\lambda = \ln \frac{(1 + l - m_1 - m_2)l}{(m_1 - l)(m_2 - l)}, \quad (14)$$

and substituting it in the right side of eq. (12).

### 3 Scientifically motivated prior parameter densities

The interpretation of the Koopman-Pitman-Lauritzen formula (2) explained in the previous section gives us more intuitive grounds to choose the prior parameter density  $p(t|I) dt$ : given the interpretation of the observables  $s$  in a particular scientific context, which limit statistics  $t$  would we expect to observe?

If  $s$  represents the binned activity of a neural population in the brain, for example, from our research experience we consider more likely to find low mean values  $\bar{s}_1 = m_1$ ,  $\bar{s}_2 = m_2$  than high ones, close to 1. We may also have some vague expectations about the second moments  $\bar{s}\bar{s} = l$ . Even vague prior knowledge can be expressed by a probability density with particular features, and this leads to better predictions. Let's examine this possibility more concretely.

Using Bayes's theorem with formula (2) we find our probability for a new outcome  $s$  conditional on observations  $(s^{(1)}, \dots, s^{(N)})$ :

$$p(s | s^{(1)}, \dots, s^{(N)}, I) = \int \frac{\exp(\mu_1 s_1 + \mu_2 s_2 + \lambda s_1 s_2)}{Z(\theta)} p(\theta | s^{(1)}, \dots, s^{(N)} I) d\theta \quad (15a)$$

with

$$p(\theta | s^{(1)}, \dots, s^{(N)} I) \propto \left[ \prod_{i=1}^N \frac{\exp(\mu_1 s_1^{(i)} + \mu_2 s_2^{(i)} + \lambda s_1^{(i)} s_2^{(i)})}{Z(\theta)} \right] p(\theta | I) \equiv \exp\{N [\mu_1 a_1 + \mu_2 a_2 + \lambda \bar{s}\bar{s} - \ln Z(\theta)]\} p(\theta | I). \quad (15b)$$

The density  $p(\theta | s^{(1)}, \dots, s^{(N)} I)$  is called *posterior parameter density*.

The last expression shows that the  $N$  observations affect our forecast only through the averages  $\bar{s}$  and  $\bar{s}\bar{s}$ , eq. (1), as we assumed.

The proportionality relation of the last formula reminds us that we must perform an integral over  $\theta$  to calculate the posterior parameter density. We must also perform an integral over  $\theta$  to calculate the conditional probability for  $s$ . These integrals are difficult when we consider populations with many units. When the number  $N$  of known observations is large, the posterior parameter density is often approximated by a Dirac delta centred on the maximum of the posterior,

$$\theta_m := \arg \sup_{\theta} \{N [\mu_1 a_1 + \mu_2 a_2 + \lambda \bar{s}\bar{s} - \ln Z(\theta)] + \ln p(\theta | I)\}. \quad (16)$$

The probability for  $s$  then equals the exponential calculated at  $\theta_m$ . If the prior parameter density  $p(\theta | I)$  is constant or very broad, it can be dropped in the calculation of the maximum, as an approximation.

The literature indeed often assumes a prior parameter density  $p(\theta|I)$  that is constant in  $\theta$ . This is an ‘improper’, non-normalizable prior, because  $\theta \in \mathbf{R}^3$ . So we are properly considering a *sequence* of normalizable priors of increasing width – for example, normal distributions with increasing variance – and the resulting limit if it exists.

How reasonable is prior density constant in  $t$ ? Let’s find the equivalent density for the parameters  $t$  discussed in the previous section.

The Jacobian determinants of the transformations  $\theta(t)$  and  $t(\theta)$ , from eqs (14) and (9), are

$$\det\left(\frac{\partial \theta}{\partial t}\right) = \frac{1}{l(m_1 - l)(m_2 - l)(1 + l - m_1 - m_2)}, \quad (17a)$$

$$\det\left(\frac{\partial t}{\partial \theta}\right) = \det \frac{\partial^2 \ln Z(\theta)}{\partial \theta \partial \theta} = \frac{\exp(2\mu_1 + 2\mu_2 + \lambda)}{Z(\theta)^4}. \quad (17b)$$

These expressions are worthy of notice, because they can be uniquely written as

$$\det\left(\frac{\partial \theta}{\partial t}\right) = \frac{1}{p(00|t, I)p(10|t, I)p(01|t, I)p(11|t, I)} \equiv \frac{1}{\prod_s p(s|t, I)}, \quad (18a)$$

$$\det\left(\frac{\partial t}{\partial \theta}\right) = p(00|\theta, I)p(10|\theta, I)p(01|\theta, I)p(11|\theta, I) \equiv \prod_s p(s|\theta, I), \quad (18b)$$

as can be checked from eq. (12) for  $N = 1$ .

A prior density constant in  $\theta$  – denote this assumption by  $I_\theta$  – therefore corresponds to the density

$$p(t|I_\theta) dt \propto \det\left(\frac{\partial \theta}{\partial t}\right) dt \equiv \frac{dt}{l(m_1 - l)(m_2 - l)(1 + l - m_1 - m_2)}, \quad (19)$$

obtained from eqs (18) and (13). This density, besides being improper, gives very high probability to the extreme values of all three statistics. It doesn’t seem much appropriate in the context of brain activity discussed above, for example.

Thus two questions appear: What densities are more appropriate? Do they lead to more difficult computations than those used at present?

A simple parameter density that is normalizable and doesn’t give too high probability to extreme values of the statistics is that constant in the



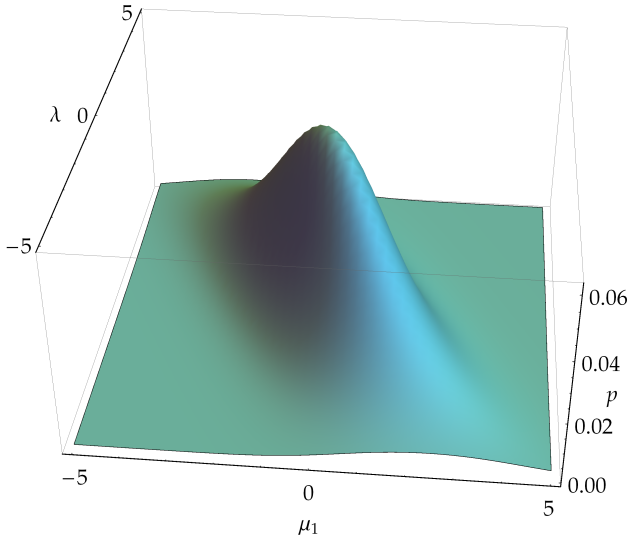


Figure 2

$t$  coordinates; denote this assumption by  $I_t$ :

$$p(t|I_t)dt = 6dt, \quad (20)$$

the normalization constant calculated from solid geometry looking at the pyramid of fig. 1, having volume  $1/6$ . In terms of  $\theta$  coordinates, using the Jacobian determinant (18), it is

$$p(\theta|I_t)d\theta = 6 \left[ \prod_s p(s|\theta, I_t) \right] d\theta \equiv 6 \frac{\exp(2\mu_1 + 2\mu_2 + \lambda)}{Z(\theta)^4} d\theta. \quad (21)$$

Its marginal for  $(\mu_1, \lambda)$  is shown in fig. 2.

Comparing the prior parameter density above with the general formula (15b) for the posterior density, we see that *the prior density  $I_t$  is equivalent to the posterior density  $I_\theta$  conditional on having observed all four possible outcomes once*:

$$p(\theta|I_t)d\theta = p(\theta|00,10,01,11,I_\theta)d\theta. \quad (22)$$

We can write this as  $I_t = I_\theta \wedge A$ , where  $A$  represents the observation of the four outcomes.

This result is computationally important. If we want to make inferences conditional on some data  $D$  using a density constant in  $t$ , we can use the same algorithms and approximations used for the density constant in  $\theta$ , but augmenting the data  $D$  with the ‘auxiliary data’  $A$ .

If it can be proven that the formula for the Jacobian determinant (18) holds for any number of units, then this method requires to add  $2^n$  auxiliary data if we consider  $n$  units. This should have a big influence on our predictions even when the observed data are numerous.

✚ To be continued

Conjecture:

$$\det \left[ \frac{\partial^2 \ln Z(\theta)}{\partial \theta_k \partial \theta_l} \right] = \prod_{s \in \{0,1\}^n} \frac{\exp(\sum_{i=1}^n h_i s_i + \sum_{1 \leq i < j \leq n} J_{ij} s_i s_j)}{Z(\theta)}$$

with  $Z(\theta) := \sum_s \exp(\dots)$ ,  $\theta := (\theta_k) := (h_i, J_{ij})$  (23)

## Bibliography

(‘de  $X$ ’ is listed under D, ‘van  $X$ ’ under V, and so on, regardless of national conventions.)

Mead, L. R., Papanicolaou, N. (1984): *Maximum entropy in the problem of moments*. J. Math. Phys. **25**<sup>8</sup>, 2404–2417. <http://bayes.wustl.edu/Manual/MeadPapanicolaou.pdf>.