# Parameter priors for Ising models

## research notes

Y. Roudi
<yasser.roudi@ntnu.no>

P.G.L. Porta Mana
<piero.mana@ntnu.no>

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 $m$ observations $(s^{(1)}, \ldots, s^{(m)})$ 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,

$$m, \qquad \tfrac{1}{m}\big(s^{(1)} + \cdots + s^{(m)}\big) =: \boldsymbol{a}, \qquad \tfrac{1}{m}\big(s_1^{(1)}s_2^{(1)} + \cdots + s_1^{(m)}s_2^{(m)}\big) =: c \quad (1)$$

are sufficient statistics; note that the second sum contains the first as its diagonal. These assumptions are collectively denoted $I$. Then the Koopman-Pitman theorem says that our probabilistic forecasts must assume this general form:

$$p(s^{(1)}, \ldots, s^{(n)} | I) =$$
$$\int \left[ \prod_{i=1}^{n} g(s^{(i)}) \frac{\exp\big(\mu_1 s_1^{(i)} + \mu_2 s_2^{(i)} + \lambda s_1^{(i)} s_2^{(i)}\big)}{Z(\mu_1, \mu_2, \lambda)} \right] p(\mu_1, \mu_2, \lambda | I)\, d\mu_1\, d\mu_2\, d\lambda,$$

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

Let's denote $\theta := (\mu_1, \mu_2, \lambda) \in \mathbf{R}^3$.

The distribution $g(s)$ and the density $p(\mu_1, \mu_2, \lambda \mid I)$ 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 system. In the present study we assume it to be unity: $g(s) = 1$. The density $p(\theta \mid I)$ is called *prior parameter density*.

Geometrically, the formula above says that our probability distribution is the convex combination of probability distributions from a three-dimensional family parameterized by $\theta$, with weights $p(\theta \mid I)\,d\theta$. The parameters $\theta$ are just coordinates of this three-dimensional manifold, and we can choose other coordinates $t$, the weights then being given by $p(t \mid I)\,dt$, with the densities for $\theta$ and for $t$ related by a Jacobian determinant:

$$p(\theta \mid I) = p[t(\theta) \mid I] \, \det\!\left(\frac{\partial t}{\partial \theta}\right). \tag{3}$$

Using Bayes's theorem with the probabilities (2) we find our forecast for a new observation $s$ conditional on observations $\big(s^{(1)}, \ldots, s^{(m)}\big)$:

$$p(s \mid s^{(1)}, \ldots, s^{(m)}, I) =$$
$$\int \frac{\exp\big(\mu_1 s_1 + \mu_2 s_2 + \lambda s_1 s_2\big)}{Z(\theta)} \, p(\theta \mid s^{(1)}, \ldots, s^{(m)} I)\,d\theta \tag{4a}$$

with

$$p(\theta \mid s^{(1)}, \ldots, s^{(m)} I) \propto \left[\prod_{i=1}^{m} \frac{\exp\big(\mu_1 s_1^{(i)} + \mu_2 s_2^{(i)} + \lambda s_1^{(i)} s_2^{(i)}\big)}{Z(\theta)}\right] p(\theta \mid I) =$$
$$\exp\{m\,[\mu_1 a_1 + \mu_2 a_2 + \lambda c - \ln Z(\theta)]\}\, p(\theta \mid I). \tag{4b}$$

The density $p(\theta \mid s^{(1)}, \ldots, s^{(m)} I)$ is called *posterior parameter density*.

The last expression shows that the $m$ observations affect our forecast only through the averages $a$ and $c$, 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 $m$ of known observations is large, the posterior parameter density is often approximated by

a Dirac delta centred on the maximum of the posterior,

$$\theta_m := \arg\sup_{\theta}\{m\left[\mu_1 a_1 + \mu_2 a_2 + \lambda c - \ln Z(\theta)\right] + \ln p(\theta \mid I)\}. \quad (5)$$

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

## 2   Other prior parameter densities

The literature often assumes a prior parameter density $p(\theta \mid I)$ that is constant in $\theta$. This is an 'improper', non-normalizable prior; we are really considering a sequence of normalizable priors of increasing width – for example, normal distributions with increasing variance – and the resulting limit if it exists.

As noted before, the parameters $\theta$ are just coordinates in a manifold of distributions for $s$. A constant density in these coordinates corresponds to an non-constant density in other coordinates. But are these coordinates 'special' in any way, to consider a constant density in them? Are there other coordinates in which it makes more sense to consider a constant density? How does such a different choice affect our inference about $s$?