

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 18 July 2018

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

‘Flat priors do not exist’
(anonymous)

1 Initial assumptions for models with second-order sufficient statistics

1.1 A two-unit model with sufficient statistics

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

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

$$N, \quad \frac{1}{N}(\mathbf{s}^{(1)} + \dots + \mathbf{s}^{(N)}) =: \bar{\mathbf{s}}, \quad \frac{1}{N}(\mathbf{s}_1^{(1)}\mathbf{s}_2^{(1)} + \dots + \mathbf{s}_1^{(N)}\mathbf{s}_2^{(N)}) =: \overline{\mathbf{s}\mathbf{s}} \quad (1)$$

These assumptions are collectively denoted I .

A series of mathematical results, which we call the Koopman-Pitman-Lauritzen theorem, says that our probabilistic forecasts must assume this general form, for any N :

$$\begin{aligned}
 p(s^{(1)}, \dots, s^{(N)} | I) \\
 &= \int \left[\prod_{i=1}^N g(s^{(i)}) \frac{\exp(\mu_1 s_1^{(i)} + \mu_2 s_2^{(i)} + \mu_{12} s_1^{(i)} s_2^{(i)})}{Z(\mu)} \right] p(\mu | I) d\mu \\
 &= \int \left[\prod_{i=1}^N g(s^{(i)}) \right] \frac{\exp[N(\mu_1 \bar{s}_1 + \mu_2 \bar{s}_2 + \mu_{12} \bar{s} \bar{s})]}{Z(\mu)^N} p(\mu | I) d\mu \\
 &\quad \text{with } \mu := (\mu_1, \mu_2, \mu_{12}) \in \mathbf{R}^3, \\
 &\quad Z(\mu) := g(00) + g(10) \exp(\mu_1) + g(01) \exp(\mu_2) + g(11) \exp(\mu_1 + \mu_2 + \mu_{12}).
 \end{aligned} \tag{2}$$

The distribution $g(s)$ and the density $p(\mu | I) d\mu$ 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: $g(s) = 1$. The density $p(\mu | I) d\mu$ is called *prior parameter density*.

The formula above says that our joint probability distribution for the N outcomes is given by a mixture of joint, factorizable distributions from a three-parameter family. This family is a submanifold in the space of all possible joint distributions, which has dimension $4^N - 1$. The parameters μ are coordinates in this submanifold, each triplet identifying a particular factorizable distribution

$$p(s^{(1)}, \dots, s^{(N)} | \mu, I) = \prod_{i=1}^N \frac{\exp(\mu_1 s_1^{(i)} + \mu_2 s_2^{(i)} + \mu_{12} s_1^{(i)} s_2^{(i)})}{Z(\mu)}. \tag{3}$$

The weights of the mixture are $p(\mu | I) d\mu$.

From the integral of the formula (6) and the geometric interpretation above it is clear that the theorem does not select the coordinates μ . We can choose another set of coordinates $m := (m_1, m_2, m_{12})$, related to the μ by a one-one transformation $m(\mu)$ with inverse $\mu(m)$. The

three-dimensional family is then labelled as

$$p(s^{(1)}, \dots, s^{(N)} | \mathbf{m}, I) = \prod_{i=1}^N \frac{\exp[\mu_1(\mathbf{m}) s_1^{(i)} + \mu_2(\mathbf{m}) s_2^{(i)} + \mu_{12}(\mathbf{m}) s_1^{(i)} s_2^{(i)}]}{Z[\boldsymbol{\mu}(\mathbf{m})]}, \quad (4)$$

the prior parameter density is

$$p(\mathbf{m} | I) d\mathbf{m} = \det\left(\frac{\partial \mathbf{m}}{\partial \boldsymbol{\mu}}\right) p(\boldsymbol{\mu} | I) d\boldsymbol{\mu}, \quad (5)$$

and the integral formula becomes, with $g(s) = 1$,

$$\begin{aligned} p(s^{(1)}, \dots, s^{(N)} | I) &= \int \left[\prod_{i=1}^N \frac{\exp[\mu_1(\mathbf{m}) s_1^{(i)} + \mu_2(\mathbf{m}) s_2^{(i)} + \mu_{12}(\mathbf{m}) s_1^{(i)} s_2^{(i)}]}{Z[\boldsymbol{\mu}(\mathbf{m})]} \right] p(\mathbf{m} | I) d\mathbf{m} \\ &= \int \frac{\exp\{N [\mu_1(\mathbf{m}) \bar{s}_1 + \mu_2(\mathbf{m}) \bar{s}_2 + \mu_{12}(\mathbf{m}) \bar{s}\bar{s}]\}}{Z[\boldsymbol{\mu}(\mathbf{m})]^N} p(\mathbf{m} | I) d\mathbf{m}, \end{aligned} \quad (6)$$

equivalent to (6).

This coordinate change is central to the rest of this study.

1.2 New coordinates and their motivation

Assuming that (1) are sufficient statistics and therefore using formula (6), 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 $\mathbf{m} := (m_1, m_2, m_{12})$ 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:

$$p[(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s}) = \mathbf{x} | I, \text{large } N] \approx p(\mathbf{m} = \mathbf{x} | I). \quad (7)$$

To see this, consider the parameterized, factorized joint probability $p(s^{(1)}, \dots, s^{(N)} | \boldsymbol{\mu}, I)$ of eq. (13). The expectation of the statistics $(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s})$ is given by

$$\begin{aligned} E[(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s}) | \boldsymbol{\mu}, I] &= \\ \frac{1}{N} \sum_i E[(s_1^{(i)}, s_2^{(i)}, s_1^{(i)} s_2^{(i)}) | \boldsymbol{\mu}, I] &= E[(s_1, s_2, s_1 s_2) | \boldsymbol{\mu}, I] \end{aligned} \quad (8)$$

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 μ . From the properties of the variance we also have

$$V[(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s}) | \mu, I] = \frac{1}{N} V[(s_1, s_2, s_1 s_2) | \mu, I]. \quad (9)$$

This means that for a triplet μ , for large N we have a probability distribution for the statistics that is very peaked at particular values $\mathbf{m} := (m_1, m_2, m_{12})$ determined by the equations

$$\begin{aligned} m_1 &= E(s_1 | \mu, I) \equiv \frac{\partial \ln Z(\mu)}{\partial \mu_1}, & m_2 &= E(s_2 | \mu, I) \equiv \frac{\partial \ln Z(\mu)}{\partial \mu_2}, \\ m_{12} &= E(s_1 s_2 | \mu, I) \equiv \frac{\partial \ln Z(\mu)}{\partial \mu_{12}}. \end{aligned} \quad (10)$$

This system of equations actually puts the parameters $\mu := (\mu_1, \mu_2, \mu_{12})$ and $\mathbf{m} := (m_1, m_2, m_{12})$ into one-one correspondence (Mead et al. 1984). The former belong to \mathbf{R}^3 ; the latter to the bounded domain

$$0 \leq m_1, m_2 \leq 1, \quad \max(0, m_1 + m_2 - 1) \leq m_{12} \leq \min(m_1, m_2) \quad (11)$$

shown in fig. 1.

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

$$p(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s} | m_1, m_2, m_{12}, I) \approx \delta[(\bar{s}_1, \bar{s}_2, \bar{s}\bar{s}) - (m_1, m_2, m_{12})]. \quad (12)$$

Taking the convex combination of this expression in \mathbf{m} with weights $p(\mathbf{m} | I) d\mathbf{m}$ we obtain eq. (7).

In the coordinates \mathbf{m} , the formula (6) 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 $\mathbf{m} := (m_1, m_2, m_{12})$. Given this knowledge we can combinatorially calculate the probability of observing a finite sequence of N observations, $p(\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)} | \mathbf{m}, I)$, assuming that all sequences having given statistics are equally likely – this equiprobability corresponds to setting $g(\mathbf{s}) = 1$ in eq. (6). An example of this combinatorial calculation is given below.

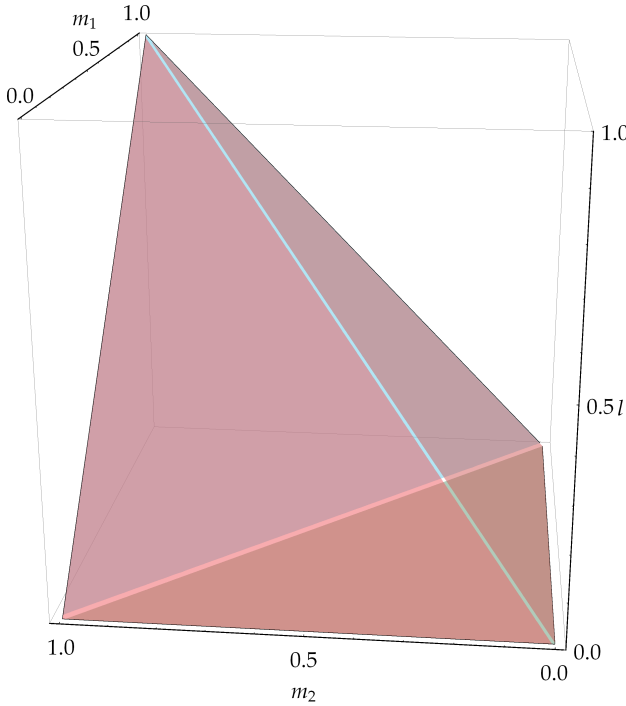


Figure 1

- (2) We then express our uncertainty about the limit statistics with the probability density $p(\mathbf{m} | I) d\mathbf{m}$.
- (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 \mathbf{m} ,

$$p(\mathbf{s} | \mathbf{m}, I) = \frac{\exp[\mu_1(\mathbf{m}) s_1 + \mu_2(\mathbf{m}) s_2 + \mu_{12}(\mathbf{m}) s_1 s_2]}{Z[\mu(\mathbf{m})]} \quad (13)$$

is indeed given combinatorially assuming equiprobability of all sequences, as claimed above. Suppose that we know the limit statistics are (m_1, m_2, m_{12}) . Only the outcome $\mathbf{s} = 11$ gives a non-vanishing contribution to the second moment m_{12} , 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 m_{12} . 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 m_{12} , the frequency of 10 must be given by $m_1 - m_{12}$, 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 $1 - m_{12} - (m_1 - 1) - (m_2 - m_{12})$. Summarizing,

$$p(s | \mathbf{m}, I) = \begin{cases} 1 + m_{12} - m_1 - m_2 & \text{for } s = 00 \\ m_1 - m_{12} & \text{for } s = 10 \\ m_2 - m_{12} & \text{for } s = 01 \\ m_{12} & \text{for } s = 11 \end{cases}$$

$$\equiv m_{12}^{s_1 s_2} (m_1 - m_{12})^{s_1 - s_1 s_2} (m_2 - m_{12})^{s_2 - s_1 s_2} (1 + m_{12} - m_1 - m_2)^{1 - s_1 - s_2 + s_1 s_2}. \quad (14)$$

This probability distribution is exactly eq. (13), as can be checked by finding $\mu(\mathbf{m})$ with the inverse of the coordinate transformations (10),

$$\mu_1 = \ln \frac{m_1 - m_{12}}{1 + m_{12} - m_1 - m_2}, \quad \mu_2 = \ln \frac{m_2 - m_{12}}{1 + m_{12} - m_1 - m_2},$$

$$\mu_{12} = \ln \frac{(1 + m_{12} - m_1 - m_2) m_{12}}{(m_1 - m_{12})(m_2 - m_{12})}, \quad (15)$$

and substituting it in the right side of eq. (13).

1.3 Scientifically motivated prior parameter densities

The interpretation of the Koopman-Pitman-Lauritzen formula (6) explained in the previous section gives us more intuitive grounds to choose the prior parameter density $p(\mathbf{m} | I) d\mathbf{m}$: given the interpretation of the observables \mathbf{s} in a particular scientific context, which limit statistics \mathbf{m} would we expect to observe?

If \mathbf{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

$\bar{s}\bar{s} = m_{12}$. 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 (6) we find our probability for a new outcome \mathbf{s} conditional on observations $(\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)})$:

$$p(\mathbf{s} | \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)}, I) = \int \frac{\exp(\mu_1 s_1 + \mu_2 s_2 + \mu_{12} s_1 s_2)}{Z(\boldsymbol{\mu})} p(\boldsymbol{\mu} | \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)} I) d\boldsymbol{\mu} \quad (16a)$$

with

$$p(\boldsymbol{\mu} | \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)} I) \propto \left[\prod_{i=1}^N \frac{\exp(\mu_1 s_1^{(i)} + \mu_2 s_2^{(i)} + \mu_{12} s_1^{(i)} s_2^{(i)})}{Z(\boldsymbol{\mu})} \right] p(\boldsymbol{\mu} | I) \equiv \exp\{N [\mu_1 \bar{s}_1 + \mu_2 \bar{s}_2 + \mu_{12} \bar{s}\bar{s} - \ln Z(\boldsymbol{\mu})]\} p(\boldsymbol{\mu} | I). \quad (16b)$$

The density $p(\boldsymbol{\mu} | \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)} I)$ 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 $\boldsymbol{\mu}$ to calculate the posterior parameter density. We must also perform an integral over $\boldsymbol{\mu}$ to calculate the conditional probability for \mathbf{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,

$$\boldsymbol{\mu}_m := \arg \sup_{\boldsymbol{\mu}} \{N [\mu_1 \bar{s}_1 + \mu_2 \bar{s}_2 + \mu_{12} \bar{s}\bar{s} - \ln Z(\boldsymbol{\mu})] + \ln p(\boldsymbol{\mu} | I)\}. \quad (17)$$

The probability for \mathbf{s} then equals the exponential calculated at $\boldsymbol{\mu}_m$. If the prior parameter density $p(\boldsymbol{\mu} | I)$ 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 $p(\boldsymbol{\mu} | I)$ that is constant in $\boldsymbol{\mu}$. This is an 'improper', non-normalizable prior, because $\boldsymbol{\mu} \in \mathbf{R}^3$. 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 \mathbf{m} ? Let's find the equivalent density for the parameters \mathbf{m} discussed in the previous section.

The Jacobian determinants of the transformations $\boldsymbol{\mu}(\mathbf{m})$ and $\mathbf{m}(\boldsymbol{\mu})$, from eqs (15) and (10), are

$$\det\left(\frac{\partial \boldsymbol{\mu}}{\partial \mathbf{m}}\right) = \frac{1}{m_{12}(m_1 - m_{12})(m_2 - m_{12})(1 + m_{12} - m_1 - m_2)}, \quad (18a)$$

$$\det\left(\frac{\partial \mathbf{m}}{\partial \boldsymbol{\mu}}\right) = \det \frac{\partial^2 \ln Z(\boldsymbol{\mu})}{\partial \boldsymbol{\mu} \partial \boldsymbol{\mu}} = \frac{\exp(2\mu_1 + 2\mu_2 + \mu_{12})}{Z(\boldsymbol{\mu})^4}. \quad (18b)$$

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

$$\det\left(\frac{\partial \boldsymbol{\mu}}{\partial \mathbf{m}}\right) = \frac{1}{p(00|\mathbf{m}, I)p(10|\mathbf{m}, I)p(01|\mathbf{m}, I)p(11|\mathbf{m}, I)} \equiv \frac{1}{\prod_s p(s|\mathbf{m}, I)}, \quad (19a)$$

$$\det\left(\frac{\partial \mathbf{m}}{\partial \boldsymbol{\mu}}\right) = p(00|\boldsymbol{\mu}, I)p(10|\boldsymbol{\mu}, I)p(01|\boldsymbol{\mu}, I)p(11|\boldsymbol{\mu}, I) \equiv \prod_s p(s|\boldsymbol{\mu}, I), \quad (19b)$$

as can be checked from eq. (13) for $N = 1$.

Consider an assumption, denoted I_μ , that lead us to assign a constant prior parameter density function $p(\boldsymbol{\mu}|I_\mu)$. For \mathbf{m} such an assumption corresponds to the density

$$p(\mathbf{m}|I_\mu) d\mathbf{m} \propto \det\left(\frac{\partial \boldsymbol{\mu}}{\partial \mathbf{m}}\right) d\mathbf{m} \equiv \frac{d\mathbf{m}}{m_{12}(m_1 - m_{12})(m_2 - m_{12})(1 + m_{12} - m_1 - m_2)}, \quad (20)$$

obtained from eqs (19) and (14). 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?

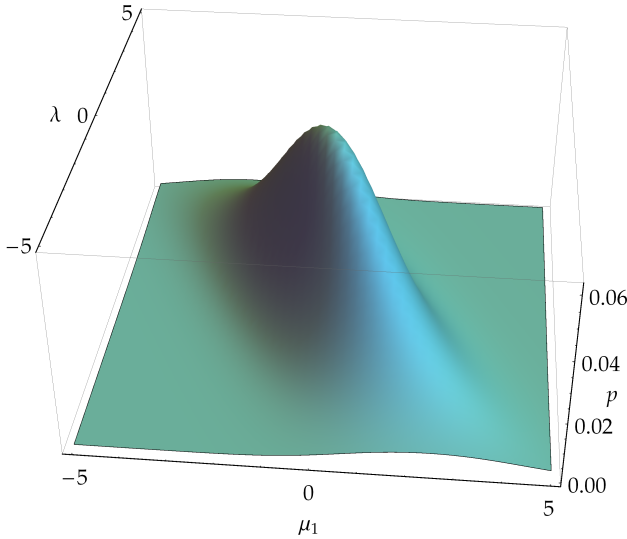


Figure 2

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 \mathbf{m} coordinates; denote this assumption by I_m :

$$p(\mathbf{m} | I_m) d\mathbf{m} = 6 d\mathbf{m}, \quad (21)$$

the normalization constant calculated from solid geometry looking at the pyramid of fig. 1, having volume $1/6$. In terms of μ coordinates, using the Jacobian determinant (19), it is

$$p(\mu | I_m) d\mu = 6 \left[\prod_s p(s | \mu, I_m) \right] d\mu \equiv 6 \frac{\exp(2\mu_1 + 2\mu_2 + \mu_{12})}{Z(\mu)^4} d\mu. \quad (22)$$

Its marginal for (μ_1, μ_{12}) is shown in fig. 2.

Comparing the prior parameter density above with the general formula (16b) for the posterior density, we see that *the prior density I_m is equivalent to the posterior density I_μ conditional on having observed all four possible outcomes once*:

$$p(\mu | I_m) d\mu = p(\mu | 00, 10, 01, 11, I_\mu) d\mu. \quad (23)$$

We can write this as $I_m = I_\mu \wedge A$, 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 m , we can use the same algorithms and approximations used for the density constant in μ , but augmenting the data D with the ‘auxiliary data’ A .

If it can be proven that the formula for the Jacobian determinant (19) 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.

2 [Luca:] Does it make sense to test against computer-generated distributions?

But we must note with sadness that, in much of the current Bayesian literature, very little of the orthodox baggage has been cast off. For example, it is rather typical to see a Bayesian article start with such phrases as: ‘Let X be a random variable with density function $p(x|\theta)$, where the value of the parameter θ is unknown. Suppose this parametric family contains the true distribution of $X \dots$ ’ Or, one describes a uniform prior $p(\theta|I)$ by saying: ‘ θ is supposed uniformly distributed’. The analytical solutions thus obtained will doubtless be a valid Bayesian result; but one is still clinging to the orthodox fiction of ‘random variables’ and ‘true distributions’. θ is simply an unknown constant; it is not ‘distributed’ at all. What is ‘distributed’ is our state of knowledge about θ : again there is that persistent mind projection fallacy that contaminates all of probability theory, leading inexperienced readers far astray as to what we are doing. Equally bad, those who commit this fallacy seem unaware that this is restricting the application to a small fraction of the real situations where the solution might be useful. In the vast majority of real applications there are no ‘random variables’ (What defines ‘randomness’?) and no ‘true distribution’ (What defines it? What test could we apply to decide whether some proposed distribution is or is not the ‘true’ one?); yet probability theory as logic applies to all of them.

Jaynes (2003 § 17.12)

✠ To be continued

Bibliography

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

Jaynes, E. T. (2003): *Probability Theory: The Logic of Science*. (Cambridge University Press, Cambridge). Ed. by G. Larry Bretthorst. First publ. 1994. <https://archive.org/details/XQUHIUXHIQUHIQXUIHX2>, <http://www-biba.inrialpes.fr/Jaynes/prob.html>, <http://omega.albany.edu:8008/JaynesBook.html>.

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