

---

# Representative samples and maximum-entropy distributions in neuroscience: a dilemma

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This note has three nested purposes. The first purpose is to show that the maximum-  
2 entropy method can be applied to a representative random sample of a population, to  
3 generate its probability distribution, along two different routes. Both routes appear  
4 legitimate, but they give inequivalent results. Which route should be chosen?  
5 Some arguments are presented in favour of one. The second more general purpose,  
6 motivated by the above dilemma, is to remind readers that models like maximum-  
7 entropy may contain hidden assumptions; in this case the hidden an unnatural  
8 assumption that the sample modelled is isolated from the rest of the population. The  
9 third purpose is to promote some old but possibly forgotten probability formulae  
10 that may be useful in neuroscientific sampling contexts.

## 11 1 Introduction: maximum-entropy and sampling in neuroscience

12 This note is mainly addressed to neuroscientists interested in maximum-entropy methods, but we  
13 would be pleased if its discussion of the probability of “sampling” were useful to neuroscientists that  
14 use other statistical methods to study the physical and dynamical characteristics of brain areas via  
15 neuronal recording.

16 Recent electrophysiological techniques [1] in fact allow experimenters to record samples comprising  
17 even a couple hundreds neurons from specific brain areas. From the observation of these samples  
18 we expect to learn something about all other neurons in the same brain area. That is, we assume that  
19 they are a “representative sample”. We don’t elaborate on the various important purposes of such  
20 recordings here, but stress that these sample sizes can be considered “large”, because their statistical  
21 analysis requires considerable computational power.

22 These computational costs are one, probably not the earliest, of the many reasons why maximum-  
23 entropy methods have been introduced in neuroscience. It would be useful to somehow compress the  
24 statistical wealth of large neuronal recordings into few quantities, like sample moments for example.  
25 This compression would also entail interesting physico-biological properties of neuronal activity.  
26 The standard maximum-entropy method [2–4] accomplishes this kind of compression: it associates a  
27 unique probability distribution with few experimental quantities. But this is only one of its uses. It is  
28 also used for various information-theoretic purposes or to generate reference probability distributions  
29 [5–11]. In all these uses the maximum-entropy distribution is chosen as the “maximally noncommittal”  
30 one [12]. This adjective means little without further technical characterizations. Different works give  
31 different characterizations, but in this paper we will use the quoted expression as an umbrella term  
32 for all of them. Our results will not depend on the specific characterization of “noncommittal”.

33 In this note we show that we must face a dilemma if we want to apply the maximum-entropy method  
 34 to a representative sample to find a maximally noncommittal distribution.

35 The dilemma is this. We can apply the maximum-entropy method to the sample, using a specified  
 36 set of experimental constraints, and generate a probability distribution for its state. But our sample  
 37 is representative of a larger population. We can apply the maximum-entropy method to the larger  
 38 population, using the same constraints, and generate a probability distribution for its larger state, and  
 39 then find the distribution for the sample by marginalization. Either application seems to have some  
 40 commendable features. However, *the distributions obtained by these two applications differ*. It goes  
 41 without saying that if one is “maximally noncommittal”, the other must be somehow “committal”.  
 42 Which choice is most meaningful?

43 In the final discussion we present some arguments in favour of one choice in our dilemma.

44 In the rest of the paper we mathematically formulate this dilemma. To this purpose we also present  
 45 some probability relations relevant to the sampling problem. The relations we present are well-known  
 46 in survey sampling and in the pedagogic problem of drawing from an urn without replacement, yet  
 47 they are somewhat hard to find explicitly written in the neuroscientific literature, so they may be  
 48 of interest on their own. We briefly discuss also two very different assumptions that lead to similar  
 49 initial probability assignments for sampling.

50 Our mathematical analysis pertains to a neurons modelled as binary units at one specific instant or  
 51 short window in time. It is possible to similarly discuss multi-state neuron models and population  
 52 dynamics; but for simplicity neurons are here assumed to be in a fixed, active “1” or inactive “0”  
 53 state; and evolution, change, time correlations, and similar concepts do not concern us. We consider  
 54 maximum-entropy models that have constraints based on some kinds of sample or population averages;  
 55 they are often called “homogeneous”. The final discussion touches upon “inhomogeneous” models as  
 56 well.

57 The notation in this note follows ISO and ANSI standards [13–15] but for the use of the comma “,”  
 58 to denote logical conjunction. Probability notation follows Jaynes [16]. By “probability” we mean  
 59 plausibility or the degree of belief which “would be agreed by all rational men if there were any  
 60 rational men” [17].

## 61 2 Setup

62 We have a population of  $N$  binary neurons. We assume that they can be distinguished, by their spike  
 63 shapes for example; but other details, like their locations, are unknown. The neurons have a joint  
 64 state  $(X_1, \dots, X_N) =: \mathbf{X}$  having fixed but unknown binary values  $(R_1, \dots, R_N) =: \mathbf{R} \in \{0, 1\}^N$ . A  
 65 particular sample of  $n$  neurons from this population has joint state  $(x_1, \dots, x_n) =: \mathbf{x}$  having fixed  
 66 binary values  $(r_1, \dots, r_n) =: \mathbf{r}$ . We will consider various averages of the population and the sample.  
 67 For this purpose we introduce a general averaging operator  $\bar{\cdot}$  defined by

$$\begin{aligned}\bar{X} &:= \frac{1}{N}(X_1 + X_2 + \dots + X_N), & \overline{X\bar{X}} &:= \binom{N}{2}^{-1}(X_1 X_2 + X_1 X_3 + \dots + X_{N-1} X_N), \\ \overline{X\bar{X}\bar{X}} &:= \binom{N}{3}^{-1}(X_1 X_2 X_3 + \dots + X_{N-2} X_{N-1} X_N),\end{aligned}\tag{1}$$

68 and so on. These formulae say that  $\bar{X}$  is the fraction of active neurons,  $\overline{X\bar{X}}$  the fraction of simultan-  
 69 eously active pairs out of all  $\binom{N}{2}$  pairs,  $\overline{X\bar{X}\bar{X}}$  the fraction of simultaneously active triplets, and so on.  
 70 Products of states like  $X_i \dots X_j$  also have values in  $\{0, 1\}$ ; from this we can combinatorially prove  
 71 that

$$\underbrace{\overline{X \dots X}}_{m \text{ factors}} = \binom{N}{m}^{-1} \binom{N\bar{X}}{m}.\tag{2}$$

72 Analogous formulae hold for quantities like  $\mathbf{x}$ ,  $\mathbf{R}$ ,  $\mathbf{r}$ .

73 Our uncertainty about the actual state of the population is completely expressed by the joint probability  
 74 distribution

$$P(X_1 = R_1, X_2 = R_2, \dots, X_N = R_N | K) \quad \text{or} \quad P(\mathbf{X} = \mathbf{R} | K), \quad \mathbf{R} \in \{0, 1\}^N,\tag{3}$$

75 where  $K$  denotes our state of knowledge, i.e. the evidence and assumptions backing this particular  
 76 probability assignment. Our uncertainty about the state of the sample is likewise expressed by

$$P(x_1 = r_1, x_2 = r_2, \dots, x_n = r_n | K) \quad \text{or} \quad P(\mathbf{x} = \mathbf{r} | K), \quad \mathbf{r} \in \{0, 1\}^n. \quad (4)$$

### 77 3 Initial assumptions: the probability of representative samples

78 We need to make an initial probability assignment before any experimental observations are made.  
 79 This initial assignment will be modified by our experimental observations. We would also like our  
 80 probability assignment to reflect that the sample is somehow “representative” of the population. We  
 81 consider here two states of knowledge that express this representativeness in different ways but lead  
 82 to identical *initial* probability assignments.

83 In the first, denoted  $I'$ , we know that the neurons in the population are physico-biologically similar,  
 84 for example in morphology and kind of input they receive. Knowledge of this similarity leads us to  
 85 assign a probability distribution for the population state  $\mathbf{X}$  that is symmetric, also called *exchangeable*,  
 86 under permutations of neuron identities.

87 In the second, denoted  $I''$ , we are completely ignorant about the physical details of the individual  
 88 neurons. Our ignorance is therefore symmetric under permutations of neuron identities. This also  
 89 leads to an exchangeable probability distribution for  $\mathbf{X}$ .

90 Let us use  $I$  to denote either of these two states of knowledge in those probabilities that are identical  
 91 in both.

92 The *representation theorem for finite exchangeability* states that the symmetric distribution of  $I$  must  
 93 obey

$$P(\mathbf{X} = \mathbf{R} | I) \equiv P(\mathbf{X} = \mathbf{R} | \bar{\mathbf{X}} = \bar{\mathbf{R}}, I) P(\bar{\mathbf{X}} = \bar{\mathbf{R}} | I) = \left( \frac{N}{N\bar{\mathbf{R}}} \right)^{-1} P(\bar{\mathbf{X}} = \bar{\mathbf{R}} | I), \quad (5)$$

94 the latter being the probability for the population average  $\bar{\mathbf{X}}$ . Proof of this theorem and generalizations  
 95 to non-binary and continuum cases are given by de Finetti [18], Kendall [19], Ericson [20], Diaconis  
 96 & Freedman [21; 22], Heath & Sudderth [23]. This theorem is intuitive: owing to symmetry, we must  
 97 assign equal probabilities to all states with  $N\bar{\mathbf{R}}$  active neurons.

98 By marginalization we obtain the probability for the state of the sample:

$$P(\mathbf{x} = \mathbf{r} | I) = \left( \frac{n}{n\bar{\mathbf{r}}} \right)^{-1} P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | I), \quad (6)$$

99 with

$$P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | I) = \sum_{N\bar{\mathbf{R}}=0}^N P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | \bar{\mathbf{X}} = \bar{\mathbf{R}}, I) P(\bar{\mathbf{X}} = \bar{\mathbf{R}} | I), \quad (7)$$

$$P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | \bar{\mathbf{X}} = \bar{\mathbf{R}}, I) = \binom{n}{n\bar{\mathbf{r}}} \binom{N-n}{N\bar{\mathbf{R}}-n\bar{\mathbf{r}}} \binom{N}{N\bar{\mathbf{R}}}^{-1} =: \Pi(\bar{\mathbf{r}} | \bar{\mathbf{R}}). \quad (8)$$

100 Our initial symmetric ignorance should intuitively also apply to the sample; indeed, the probability  
 101 for the state of the sample (7) automatically satisfies the representation theorem (5) as well. The  
 102 conditional probability in the last formula is a hypergeometric distribution  $\Pi(\bar{\mathbf{r}} | \bar{\mathbf{R}})$ , typical of  
 103 “drawing without replacement” problems. The combinatorial proof of the formulae above is in fact  
 104 the same as for this class of problems [16 ch. 3; 24 § 4.8.3; 25 § II.6].

105 How is it possible that the very different states of knowledge  $I'$  and  $I''$  lead to the same formulae  
 106 above? Their difference appears as soon as we make an experimental observation, say  $X_2 = R_2 \in$   
 107  $\{0, 1\}$  and update our probabilities (5),

$$P(\mathbf{X} = \mathbf{R} | X_2 = R_2, I) \equiv P(\mathbf{X} = \mathbf{R} | \bar{\mathbf{X}} = \bar{\mathbf{R}}, X_2 = R_2, I) P(\bar{\mathbf{X}} = \bar{\mathbf{R}} | X_2 = R_2, I) \quad (9)$$

108 Both the conditional probability and the one for the average on the right side will update in very  
 109 different ways for  $I'$  and  $I''$ . The discussion of this update process in the two cases is outside the  
 110 purposes of this note, and we hope to address it in full in a future work.

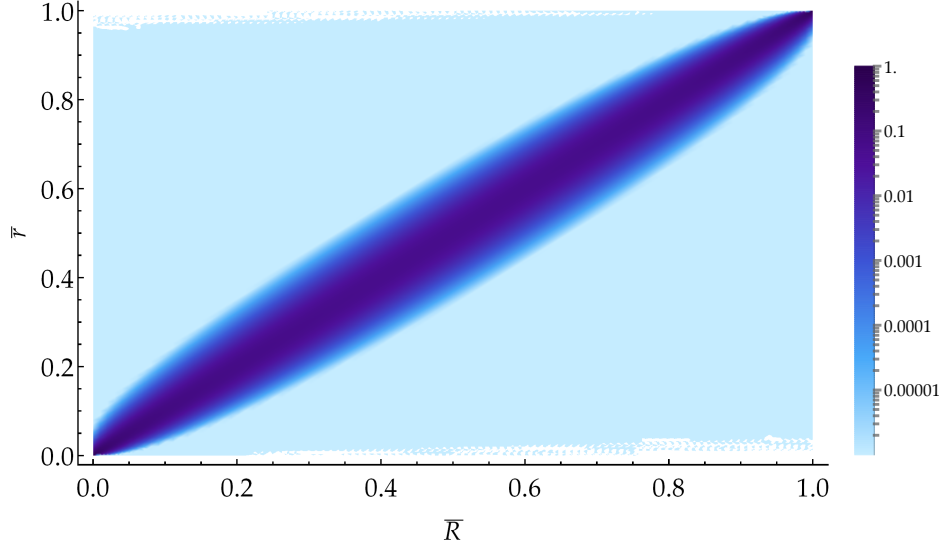


Figure 1: Log-plot of the hypergeometric distribution  $\Pi(\bar{r}|\bar{R}) := \binom{n}{\bar{r}} \binom{N-n}{N\bar{R}-n\bar{r}} \binom{N}{N\bar{R}}^{-1}$  for  $N = 5000$ ,  $n = 200$ . (Band artifacts may appear in the colourbar depending on your PDF viewer.)

111 The conditional probability  $\Pi(\bar{r}, \bar{R})$  relates the spaces of the sample average  $\bar{X} \in \{0, \dots, N\}$  and of  
 112 the population average  $\bar{x} \in \{0, \dots, n\}$  in a special way. It is a coarsening projector of any probability  
 113  $p$  for  $\bar{X}$  onto a marginal probability  $p_*$  for  $\bar{x}$ :

$$p_*(\bar{x} = \bar{r}) = \sum_{N\bar{R}=0}^N \Pi(\bar{r}|\bar{R}) p(\bar{X} = \bar{R}). \quad (10)$$

114 Conversely it also pulls back expectations of functions  $f$  of the sample average  $\bar{x}$  to expectations of  
 115 functions  $f^*$  of the population average  $\bar{X}$ :

$$f^*(\bar{X}) := \sum_{n\bar{r}=0}^n f(\bar{r}) \Pi(\bar{r}|\bar{X}),$$

$$\mathbb{E}[f(\bar{x})] = \mathbb{E}[f^*(\bar{X})] = \sum_{n\bar{r}=0}^n f(\bar{r}) \mathbb{P}(\bar{x} = \bar{r} | I) = \sum_{N\bar{R}=0}^N f^*(\bar{R}) \mathbb{P}(\bar{X} = \bar{R} | I). \quad (11)$$

116 A look at a plot of the hypergeometric distribution  $\Pi(\bar{r}|\bar{R})$ , see fig. 1, reveals that it is a sort of  
 117 “fuzzy identity matrix” between the  $\bar{X}$ -space  $\{0, \dots, N\}$  and  $\bar{x}$ -space  $\{0, \dots, n\}$ . When  $n = N$  it is  
 118 the identity matrix. We thus have that

$$\mathbb{P}(\bar{x} = a) \approx \mathbb{P}(\bar{X} = a), \quad \mathbb{E}[f(\bar{x})] \approx \mathbb{E}[f(\bar{X})]. \quad (12)$$

119 These are only very approximate equalities: they may miss important features of the two probability  
 120 distributions. In the next section we will in fact emphasize their differences. If the distribution for the  
 121 population average  $\bar{X}$  is bimodal, for example, the bimodality can be lost in the distribution for the  
 122 sample average  $\bar{x}$ , owing to the coarsening effect of  $\Pi(\bar{r}|\bar{R})$ .

123 The approximate equalities above express the fact that *our uncertainty about the sample is repres-*  
 124 *entative of our uncertainty about the population and about other samples*, and vice versa. This fact  
 125 comes about for very different reasons in the states of knowledge  $I'$  and  $I''$ . In  $I'$ , because our sample  
 126 is *physically* representative of the population and of other samples; we could write this as  $\bar{x} \approx \bar{X}$ . In  
 127  $I''$ , because we are ignorant about sample and population in similar symmetric ways; but this does  
 128 not imply that  $\bar{x} \approx \bar{X}$ . New observations may in fact break this symmetry by eq. (9).

129 Note that formulae (12) say more than the limits  $\mathbb{P}(\bar{x} = a) \rightarrow \mathbb{P}(\bar{X} = a)$ ,  $\mathbb{E}[f(\bar{x})] \rightarrow \mathbb{E}[f(\bar{X})]$  as  
 130  $n \rightarrow N$ . These limits are trivially valid because the sample becomes the full population as  $n \rightarrow N$ . In

particular, these limits hold even in cases where the conditional probability  $P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | \bar{\mathbf{X}} = \bar{\mathbf{R}})$  is not a fuzzy identity and our uncertainties about sample and about population can differ wildly.

For functions representing averaged products,  $f(\bar{\mathbf{x}}) := \overline{\mathbf{x} \cdots \mathbf{x}} = \binom{n\bar{\mathbf{x}}}{m} / \binom{n}{m}$ , formulae (11) have the useful form

$$\begin{aligned} \underbrace{(\bar{\mathbf{x}} \cdots \bar{\mathbf{x}})}_{m \text{ factors}}^* &= \underbrace{\bar{\mathbf{X}} \cdots \bar{\mathbf{X}}}_{m \text{ factors}}, \\ E(\overline{\mathbf{x} \cdots \mathbf{x}} | I) &= E(\overline{\mathbf{X} \cdots \mathbf{X}} | I) = \binom{n}{m}^{-1} \sum_{n\bar{\mathbf{r}}=0}^n \binom{n\bar{\mathbf{r}}}{m} P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | I) = \binom{N}{m}^{-1} \sum_{N\bar{\mathbf{R}}=0}^N \binom{N\bar{\mathbf{R}}}{m} P(\bar{\mathbf{X}} = \bar{\mathbf{R}} | I). \end{aligned} \quad (13)$$

The proof uses the expression for the  $m$ th factorial moment of the hypergeometric distribution [26]. Thus, in the states of knowledge  $I'$  and  $I''$  the averages of activity products *are initially the same for the sample and for the full population*. Similar relations can be found for the raw moments  $E(\bar{\mathbf{x}}^m)$  and  $E(\bar{\mathbf{X}}^m)$ , which can be written in terms of the product expectations via eq. (2).

#### 4 Enter maximum-entropy: dilemma

Formulae (5)–(8) are constraints on our initial probability assignment, but do not determine it numerically. The probability  $P(\bar{\mathbf{X}} = \bar{\mathbf{R}} | I)$  for the population average needs to be numerically specified, and by (7) it will determine that of the sample average,  $P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | I)$ . If we numerically specify the latter, the former is not completely specified, because eq. (7) linearly constrains  $N + 1$  unknowns by only  $n + 1$  equations in this case.

We may want to specify the probability by enforcing the sample expectations of several functions to have specific values, for example  $E(\bar{\mathbf{x}}) = c_1$ ,  $E(\bar{\mathbf{x}}\bar{\mathbf{x}}) = c_2$ . This is still an underdetermined linear problem: several distributions can have the same desired expectations, as clear from eqs (13).

The maximum-entropy method is brought into play to solve this indeterminacy. It selects one distribution, purported to be “maximally noncommittal”, among those that have the desired expectations. But here’s a dilemma: formulae (11) allow us to apply the method to find the probability of the population  $P(\bar{\mathbf{X}} = \bar{\mathbf{R}} | I)$ , or of the sample  $P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | I)$ . *The two applications, however, are inequivalent*. They lead to numerically different  $P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | I)$ .

Suppose we want to constrain the sample expectations of a vector function  $\mathbf{f} = (f_1, \dots, f_m)$  to the vector values  $\mathbf{c} = (c_1, \dots, c_m)$ , that is,  $E[\mathbf{f}(\bar{\mathbf{x}})] = \mathbf{c}$ . Application of maximum-entropy [3; 4] at the population level, denoted by  $I_p$ , gives

$$P(\bar{\mathbf{X}} = \bar{\mathbf{R}} | I_p) = K \binom{N}{N\bar{\mathbf{R}}} \exp \left[ \Lambda^\top \sum_{n\bar{\mathbf{r}}=0}^n \mathbf{f}(\bar{\mathbf{r}}) \Pi(\bar{\mathbf{r}} | \bar{\mathbf{R}}) \right], \quad (14)$$

and then by marginalization with eq. (6)

$$P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | I_p) = K \sum_{N\bar{\mathbf{R}}=0}^N \Pi(\bar{\mathbf{r}} | \bar{\mathbf{R}}) \binom{N}{N\bar{\mathbf{R}}} \exp \left[ \Lambda^\top \sum_{n\bar{\mathbf{r}}=0}^n \mathbf{f}(\bar{\mathbf{r}}) \Pi(\bar{\mathbf{r}} | \bar{\mathbf{R}}) \right], \quad (15)$$

where  $K$  is a normalization constant and  $\Lambda^\top = (\Lambda_1, \dots, \Lambda_m)^\top$  are Lagrange multipliers such that

$$\mathbf{c} = K \sum_{n\bar{\mathbf{r}}=0}^n \sum_{N\bar{\mathbf{R}}=0}^N \mathbf{f}(\bar{\mathbf{r}}) \Pi(\bar{\mathbf{r}} | \bar{\mathbf{R}}) \binom{N}{N\bar{\mathbf{R}}} \exp \left[ \Lambda^\top \sum_{n\bar{\mathbf{r}}=0}^n \mathbf{f}(\bar{\mathbf{r}}) \Pi(\bar{\mathbf{r}} | \bar{\mathbf{R}}) \right]. \quad (16)$$

Application of maximum-entropy at the sample level, denoted by  $I_s$ , gives

$$P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | I_s) = \kappa \binom{n}{n\bar{\mathbf{r}}} \exp[\lambda^\top \mathbf{f}(\bar{\mathbf{r}})] \quad (17)$$

where  $\kappa$  is a normalization constant and  $\lambda^\top$  are Lagrange multipliers such that

$$\mathbf{c} = \kappa \sum_{n\bar{\mathbf{r}}=0}^n \mathbf{f}(\bar{\mathbf{r}}) \binom{n}{n\bar{\mathbf{r}}} \exp[\lambda^\top \mathbf{f}(\bar{\mathbf{r}})]. \quad (18)$$

△ maxent at population level, marginalized    ○ maxent at sample level

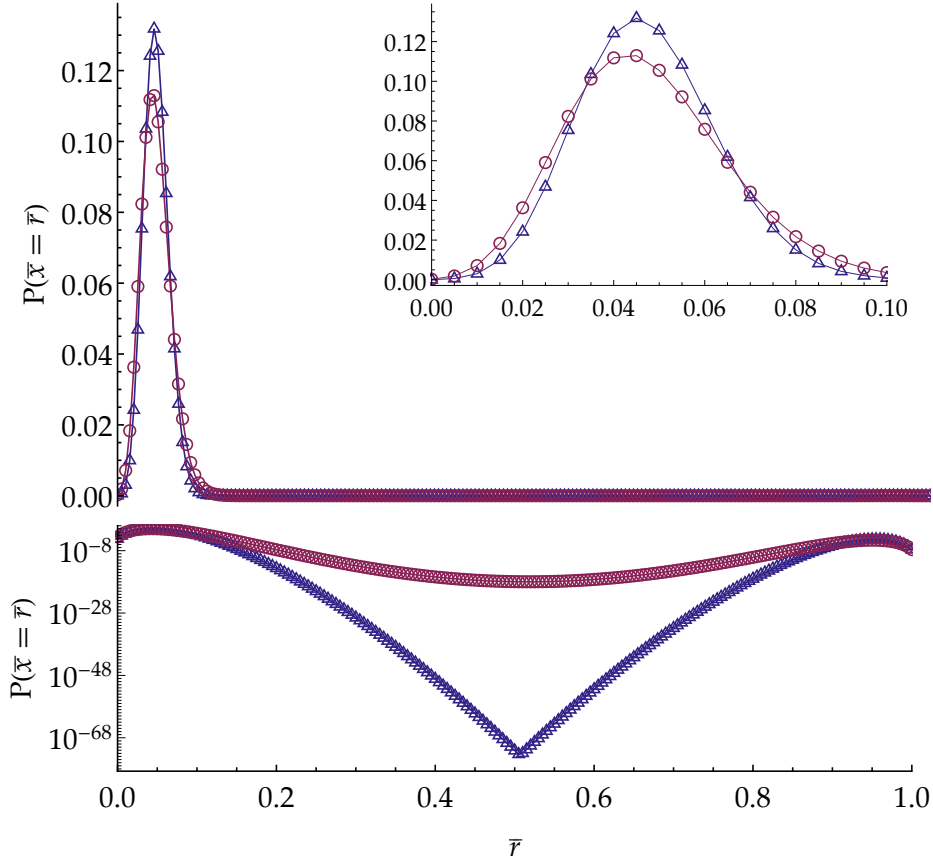


Figure 2: Linear and log-plots of  $P(\bar{x} = \bar{r})$  constructed by maximum-entropy at the population level followed by sample marginalization (blue triangles), eq. (15), and at the sample level (red circles), eq. (17), with  $N = 5000$ ,  $n = 200$ , constraints  $E(\bar{x}) = 0.0478$ ,  $E(\bar{x}\bar{x}) = 0.00257$ .

160 The probabilities for the sample average obtained from application at the population level (15)  
 161 and at the sample level (17) should be approximately equal, by our previous observation about  
 162 representativity (12) and also by the fact that they must satisfy the same expectations for  $f$ .

163 Yet they cannot be exactly equal, because their equality would require the Lagrange multipliers  $\Lambda$   
 164 and  $\lambda$  to satisfy equations (16), (18), and  $P(\bar{x} = \bar{r} | I_p) = P(\bar{x} = \bar{r} | I_s)$ ; that is,  $2m + n$  equations  
 165 (normalization is taken care of) in  $2m$  unknowns. A solution can exist, if at all, only for very special  
 166 choices of constraints functions  $f$  and values  $c$ .

167 The sample distribution obtained from maximum-entropy at the sample level will therefore likely  
 168 miss important features present in the one obtained at the population level, like additional modes or  
 169 particular tail behaviour.

170 We show two examples of this discrepancy in figs 2 and 3, for  $N = 5000$ ,  $n = 200$ , and constraint  
 171 functions of the form  $f(\bar{x}) = (\bar{x}, \bar{x}\bar{x}, \dots) \equiv (\bar{x}, \binom{n\bar{x}}{2}/\binom{n}{2}, \dots)$ . In the first example the constraint  
 172 functions are  $E(\bar{x})$  and  $E(\bar{x}\bar{x})$ . The distribution obtained at the sample level is broader than the one  
 173 obtained at the population level; the tails of the two distributions are very different. The second  
 174 example uses two additional constraint function  $E(\bar{x}\bar{x}\bar{x})$ ,  $E(\bar{x}\bar{x}\bar{x}\bar{x})$ . The distribution obtained at the  
 175 population level has two modes, replaced by only one in the distribution obtained at the sample level;  
 176 the tails are very different also in this case. The constraint values used in these examples, reported in  
 177 the figure captions, have neurobiologically realistic values [27].

△ maxent at population level, marginalized    ○ maxent at sample level

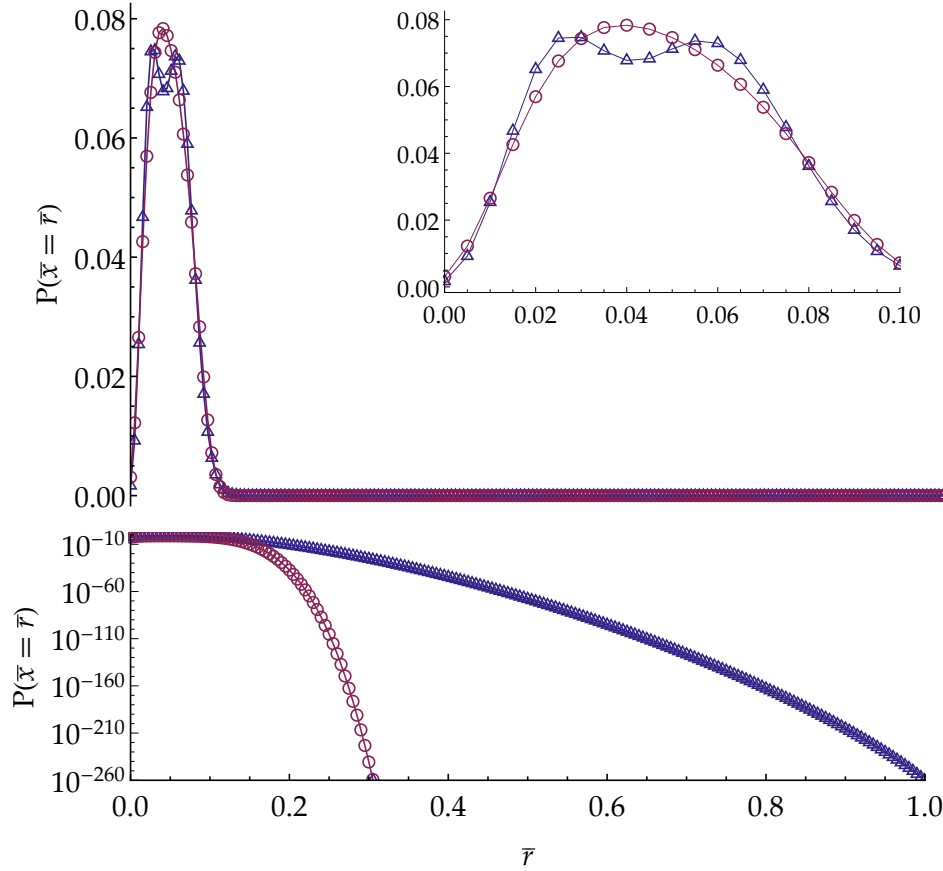


Figure 3: Linear and log-plots of  $P(\bar{x} = \bar{r})$  constructed by maximum-entropy at the population level followed by sample marginalization (blue triangles), eq. (15), and at the sample level (red circles), eq. (17), with  $N = 5000$ ,  $n = 200$ , constraints  $E(\bar{x}) = 0.0478$ ,  $E(\bar{x}\bar{x}) = 0.00257$ ,  $E(\bar{x}\bar{x}\bar{x}) = 1.48 \times 10^{-4}$ ,  $E(\bar{x}\bar{x}\bar{x}\bar{x}) = 8.81 \times 10^{-6}$ .

178 How should we apply the maximum-entropy method then? on the sample or on the population?  
 179 Which application is maximally noncommittal?

## 180 5 Discussion

181 The question that closed the preceding section cannot receive a categorical answer. An optimal answer  
 182 can only be given case by case, depending on the computational power available, on which inferences  
 183 we are trying to make, on which assumptions we need or want to make and those we wish to avoid.

184 The tricky point is this. Maximum-entropy applied at the population level and applied at the sample  
 185 level give different results; they are different statistical models. The former model clearly assumes,  
 186 by construction, the existence of a larger population from which the sample is taken. What does the  
 187 latter model assume in this respect? is it “unassuming”, as often claimed in the literature? or is it  
 188 actually assuming that *no* larger population exists? In the latter case it would not be the correct model  
 189 to use in our problem.

190 A perfunctory intuitive reasoning seems insufficient for clarifying this point. Let’s express it in  
 191 the language of the probability calculus. Suppose we do not know whether the sample is really  
 192 part of a larger population: we do not know whether  $N = n$  or how large is  $N$  otherwise. Call this  
 193 state of ignorance  $\gamma$ . In the probability calculus, this ignorance about  $N$  is expressed by assigning  
 194 a probability distribution  $P(N|\gamma)$  that vanishes if  $N < n$ , since we know that  $N \geq n$ ; see Good



[28; 29] and Rissanen [30] for examples of such distributions over the integers. Maintaining our assumption of symmetric ignorance, probability assignments that do not assume a specific value of  $N$  are then obtained by multiplication of all  $N$ -dependent probabilities by  $P(N|\gamma)$  and subsequent marginalization over  $N$ . Technically speaking,  $N$  becomes a *nuisance parameter* [16; 31; 32]. The probability obtained from maximum-entropy at the population level, eq. (15), then generalizes to

$$P(\bar{x} = \bar{r} | \gamma) = \sum_N \left\{ K_N \sum_{N\bar{R}=0}^N \Pi_N(\bar{r} | \bar{R}) \binom{N}{N\bar{R}} \exp \left[ \Lambda_N^\top \sum_{n\bar{r}=0}^n f(\bar{r}) \Pi_N(\bar{r} | \bar{R}) \right] \right\} P(N | \gamma), \quad (19)$$

where  $N$ -dependencies have been made explicit. This is a formidable expression. But the answer to our question, whether the usual maximum-entropy at the sample level (17) does not assume anything about a larger population, translates now into the precise mathematical question: are the distributions (19) and (17) equal for some choice of  $P(N|\gamma)$ , with  $P(N|\gamma) \neq 0$  for  $N > n$ ? We leave this mathematical problem for future work. Note, however, that this equality is satisfied if  $P(N = 1 | \gamma) = 1$ , which means that the usual maximum-entropy model can also be interpreted as assuming that *no* larger population exists.

Rough estimates of  $N$  may be available in neuroscientific contexts, so we find the maximum-entropy model constructed at the population level very natural and preferable. After all, *physical* neuronal-network models usually include some sort of external input to the neurons as well, mimicking their embedding in a larger network. The distribution of this maximum-entropy model, when used as a reference distribution for surprise analysis, may reveal features in a dataset that were unnoticed by the standard maximum-entropy model.

The possibility of using two different distributions is not a physical contradiction. Similar situations arise in statistical mechanics. It is known that if a system is described by a maximum-entropy Gibbs state, its subsystems need not be [33]. A dilemma quite similar to ours also appears in the statistical description of the final state of a non-equilibrium process starting and ending in two equilibrium states: we can describe our knowledge about the final state either by a Gibbs distribution, or by the distribution obtained from the Liouville evolution of the Gibbs distribution assigned to the initial state. The two descriptions differ – even though the final *physical* state is obviously exactly the same [34 § 4]. The difference in the two descriptions appears because in one case we can make sharper predictions about the state thanks to our knowledge of its preceding dynamics. In this example, though, both distributions are usually immensely sharp and practically lead to the same predictions. In the cases considered in this note the difference in predictions may be relevant instead.

Our analysis touched only constraints of the sample average. The corresponding models are usually called “homogeneous” in the literature. Purely “inhomogeneous” models have also been used [6–8], in which expectations for individual neurons or groups of neurons are constrained, for example  $E(x_2)$  or  $E(x_1 x_8 x_9)$ . A short computation shows that the maximum-entropy method with this kind of constraints gives the same result whether applied at the sample or at the population level: the states of any unconstrained neurons marginalize out. This is understandable: expressing different uncertainties about neurons 2 and 5 we are breaking the symmetry of our uncertainty, which thus cannot be representative of other neurons in the sample or in the population.

Inhomogeneous models, however, require enormous computational power for large sample sizes; homogeneous models therefore retain their importance. Our analysis and dilemma also persist for hybrid homogeneous-inhomogeneous models [10; 11].

## Acknowledgments

To be added after review.

## References

- [1] A. Berényi, Z. Somogyvári, A. J. Nagy, L. Roux, J. D. Long, S. Fujisawa, E. Stark, A. Leonardo et al.: *Large-scale, high-density (up to 512 channels) recording of local circuits in behaving animals*. J. Neurophysiol. **111**<sup>5</sup> (2014), 1132–1149. <http://www.buzsakilab.com/content/PDFs/Berenyi2013.pdf>.
- [2] E. T. Jaynes: *Information theory and statistical mechanics*. Phys. Rev. **106**<sup>4</sup> (1957), 620–630. <http://b.ayes.wustl.edu/etj/node1.html>, see also ref. [35].



- [3] D. S. Sivia: *Data Analysis: A Bayesian Tutorial*, 2nd ed. Oxford University Press, Oxford (2006). Written with J. Skilling. First publ. 1996.
- [4] L. R. Mead, N. Papanicolaou: *Maximum entropy in the problem of moments*. J. Math. Phys. **25**<sup>8</sup> (1984), 2404–2417. <http://bayes.wustl.edu/Manual/MeadPapanicolaou.pdf>.
- [5] S. M. Bohte, H. Spekreijse, P. R. Roelfsema: *The effects of pair-wise and higher-order correlations on the firing rate of a postsynaptic neuron*. Neural Comp. **12**<sup>1</sup> (2000), 153–179.
- [6] E. Schneidman, M. J. Berry II, R. Segev, W. Bialek: *Weak pairwise correlations imply strongly correlated network states in a neural population*. Nature **440**<sup>7087</sup> (2006), 1007–1012. [arXiv:q-bio/0512013, http://www.weizmann.ac.il/neurobiology/labs/schneidman/The\\_Schneidman\\_Lab/Publications.html](http://www.weizmann.ac.il/neurobiology/labs/schneidman/The_Schneidman_Lab/Publications.html).
- [7] J. Shlens, G. D. Field, J. L. Gauthier, M. I. Grivich, D. Petrusca, A. Sher, A. M. Litke, E. J. Chichilnisky: *The structure of multi-neuron firing patterns in primate retina*. J. Neurosci. **26**<sup>32</sup> (2006), 8254–8266. See also correction ref. [36].
- [8] Y. Roudi, S. Nirenberg, P. E. Latham: *Pairwise maximum entropy models for studying large biological systems: when they can work and when they can't*. PLoS Computational Biology **5**<sup>5</sup> (2009), e1000380. [arXiv:0811.0903](http://arxiv.org/abs/0811.0903).
- [9] J. H. Macke, M. Oppen, M. Bethge: *Common input explains higher-order correlations and entropy in a simple model of neural population activity*. Phys. Rev. Lett. **106**<sup>20</sup> (2011), 208102. [arXiv:1009.2855](http://arxiv.org/abs/1009.2855).
- [10] G. Tkačik, T. Mora, O. Marre, D. Amodè, S. E. Palmer, M. J. Berry II, W. Bialek: *Thermodynamics and signatures of criticality in a network of neurons*. Proc. Natl. Acad. Sci. (USA) **112**<sup>37</sup> (2014), 11508–11513. [arXiv:1407.5946](http://arxiv.org/abs/1407.5946).
- [11] H. Shimazaki, K. Sadeghi, T. Ishikawa, Y. Ikegaya, T. Toyozumi: *Simultaneous silence organizes structured higher-order interactions in neural populations*. Sci. Rep. **5** (2015), 9821.
- [12] E. T. Jaynes: *Information theory and statistical mechanics*. In: Ford [37] (1963), 181–218. Repr. in ref. [38 ch. 4, pp. 39–76]; <http://bayes.wustl.edu/etj/node1.html>.
- [13] *Quantities and units*, 3rd ed. International Organization for Standardization. Geneva (1993).
- [14] *ANSI/IEEE Std 260.3-1993: American National Standard: Mathematical signs and symbols for use in physical sciences and technology*. Institute of Electrical and Electronics Engineers. New York (1993).
- [15] *Guide for the Use of the International System of Units (SI): NIST special publication 811, 1995 edition*. National Institute of Standards and Technology. Washington, D.C. (1995). <http://physics.nist.gov/cuu/Uncertainty/bibliography.html>.
- [16] E. T. Jaynes: *Probability Theory: The Logic of Science*. Cambridge University Press, Cambridge (2003). Ed. by G. Larry Bretthorst; <http://omega.albany.edu:8008/JaynesBook.html>, <http://omega.albany.edu:8008/JaynesBookPdf.html>, <http://www-biba.inrialpes.fr/Jaynes/prob.html>; first publ. 1994.
- [17] I. J. Good: *How to estimate probabilities*. J. Inst. Maths. Applics **2**<sup>4</sup> (1966), 364–383.
- [18] B. de Finetti: *La probabilità e la statistica nei rapporti con l'induzione, secondo i diversi punti di vista*. In: de Finetti [39] (1959), 1–115. Transl. as ref. [40].
- [19] D. G. Kendall: *On finite and infinite sequences of exchangeable events*. Studia Sci. Math. Hung. **2** (1967), 319–327.
- [20] W. A. Ericson: *A Bayesian approach to two-stage sampling*. Tech. rep. AFFDL-TR-75-145. University of Michigan, Ann Arbor, USA (1976). <http://hdl.handle.net/2027.42/4819>.
- [21] P. Diaconis: *Finite forms of de Finetti's theorem on exchangeability*. Synthese **36**<sup>2</sup> (1977), 271–281. <http://statweb.stanford.edu/~cgates/PERSI/year.html>.
- [22] P. Diaconis, D. Freedman: *Finite exchangeable sequences*. Ann. Prob. **8**<sup>4</sup> (1980), 745–764.
- [23] D. Heath, W. Sudderth: *De Finetti's theorem on exchangeable variables*. American Statistician **30**<sup>4</sup> (1976), 188–189.
- [24] S. Ross: *A First Course in Probability*, 8th ed. Pearson, Upper Saddle River, USA (2010). First publ. 1976.
- [25] W. Feller: *An Introduction to Probability Theory and Its Applications. Vol. I*, 3rd ed. Wiley, New York (1968). First publ. 1950.
- [26] R. B. Potts: *Note on the factorial moments of standard distributions*. Aust. J. Phys. **6**<sup>4</sup> (1953), 498–499.
- [27] V. Rostami, P. G. L. Porta Mana, M. Helias: *Pairwise maximum-entropy models and their Glauber dynamics: bimodality, bistability, non-ergodicity problems, and their elimination via inhibition*. (2016). [arXiv:1605.04740](http://arxiv.org/abs/1605.04740).
- [28] I. J. Good: *The Estimation of Probabilities: An Essay on Modern Bayesian Methods*. MIT Press, Cambridge, USA (1965).
- [29] I. J. Good: *A Bayesian significance test for multinomial distributions*. J. Roy. Stat. Soc. B **29**<sup>3</sup> (1967), 399–431. With discussion.

- 303 [30] J. Rissanen: *A universal prior for integers and estimation by minimum description length*. Ann. Stat. **11**<sup>2</sup>  
304 (1983), 416–431.
- 305 [31] D. V. Lindley: *Introduction to Probability and Statistics from a Bayesian Viewpoint. Part 2: Inference*,  
306 reprint. Cambridge University Press, Cambridge (2008). First publ. 1965.
- 307 [32] J.-M. Bernardo, A. F. Smith: *Bayesian Theory*, reprint. Wiley, New York (2000). First publ. 1994.
- 308 [33] C. Maes, F. Redig, A. Van Moffaert: *The restriction of the Ising model to a layer*. J. Stat. Phys. **96**<sup>1</sup>  
309 (1999), 69–107. [arXiv:math/9810094](https://arxiv.org/abs/math/9810094).
- 310 [34] E. T. Jaynes: *Inferential scattering*. (1993). <http://bayes.wustl.edu/etj/node1.html>; extensively  
311 rewritten version of a paper first publ. 1985 in ref. [41], pp. 377–398.
- 312 [35] E. T. Jaynes: *Information theory and statistical mechanics. II*. Phys. Rev. **108**<sup>2</sup> (1957), 171–190. [http](http://bayes.wustl.edu/etj/node1.html)  
313 [://bayes.wustl.edu/etj/node1.html](http://bayes.wustl.edu/etj/node1.html), see also ref. [2].
- 314 [36] J. Shlens, G. D. Field, J. L. Gauthier, M. I. Grivich, D. Petrusca, A. Sher, A. M. Litke, E. J. Chichilnisky:  
315 *Correction, the structure of multi-neuron firing patterns in primate retina*. J. Neurosci. **28**<sup>5</sup> (2008), 1246.  
316 See ref. [7].
- 317 [37] K. W. Ford, ed.: *Statistical Physics*. Benjamin, New York (1963).
- 318 [38] E. T. Jaynes: *E. T. Jaynes: Papers on Probability, Statistics and Statistical Physics*, reprint. Kluwer,  
319 Dordrecht (1989). Ed. by R. D. Rosenkrantz. First publ. 1983.
- 320 [39] B. de Finetti, ed.: *Induzione e statistica*, reprint. Springer, Berlin (2011). First publ. 1959.
- 321 [40] B. de Finetti: *Probability, statistics and induction: their relationship according to the various points of*  
322 *view*. In: de Finetti [42] (1972), Ch. 9, 147–227. Transl. of ref. [18].
- 323 [41] C. R. Smith, W. T. Grandy Jr., eds.: *Maximum-Entropy and Bayesian Methods in Inverse Problems*.  
324 D. Reidel, Dordrecht (1985).
- 325 [42] B. de Finetti: *Probability, Induction and Statistics: The art of guessing*. Wiley, London (1972).

326 arXiv eprints available at <http://arxiv.org/>.  
327