

---

# Representative samples and maximum-entropy distributions: a dilemma

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This note shows that the maximum-entropy method can be applied to a representat-  
2 ive sample from a neuronal population along two different routes: (1) apply to the  
3 sample; or (2) apply to the population and marginalize to the sample. These two  
4 routes give inequivalent results. Which route should be chosen? Some arguments  
5 are presented in favour of the second. The note also touches upon probability formu-  
6 lae of representative sampling and discusses their possible meanings, a discussion  
7 that may be useful for sampling problems in neuroscience.

## 8 1 Introduction: maximum-entropy and sampling in neuroscience

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

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

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

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

32 The dilemma is this. We can apply the maximum-entropy method to the sample, using a specified  
33 set of experimental constraints, and generate a probability distribution for the state of the sample.

34 But our sample is representative of a larger population, in which it is biologically and functionally  
 35 embedded. We can apply the maximum-entropy method to the larger population, using the same  
 36 constraints, and generate a probability distribution for its larger state, and then find the distribution  
 37 for the sample by marginalization. Either application seems to have some commendable features.  
 38 However, *the distributions obtained from these two applications differ*. It goes without saying that if  
 39 one is “noncommittal”, the other must be “committal”. Which choice is most meaningful?

40 In the final discussion we present some arguments in favour of one choice.

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

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

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

## 57 2 Setup: population, sample, probabilities

58 We have a population of  $N$  binary neurons. We assume that they can be distinguished, by their spike  
 59 shapes for example; but other details, like their locations, are unknown. The neurons have a joint  
 60 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  
 61 particular sample of  $n$  neurons from this population has joint state  $(x_1, \dots, x_n) =: \mathbf{x}$  having fixed  
 62 binary values  $(r_1, \dots, r_n) =: \mathbf{r} \in \{0, 1\}^n$ . We will consider various averages of the population and the  
 63 sample. For this purpose we introduce a general averaging operator  $\overline{\phantom{x}}$  defined by

$$\begin{aligned}\overline{X} &:= \frac{1}{N}(X_1 + X_2 + \dots + X_N), & \overline{X\overline{X}} &:= \binom{N}{2}^{-1}(X_1X_2 + X_1X_3 + \dots + X_{N-1}X_N), \\ \overline{X\overline{X}\overline{X}} &:= \binom{N}{3}^{-1}(X_1X_2X_3 + \dots + X_{N-2}X_{N-1}X_N),\end{aligned}\tag{1}$$

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

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

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

69 Our uncertainty about the actual state of the population is completely expressed by the joint probability  
 70 distribution

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

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

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

### 73 3 Initial assumptions: the probability of representative samples

74 We need to make an initial probability assignment before any experimental observations are made.  
 75 This initial assignment will be modified by our experimental observations. Our probability assignment  
 76 should reflect that the sample is somehow “representative” of the population. We consider here two  
 77 states of knowledge that express this representativeness in different ways but lead to identical *initial*  
 78 probability assignments.

79 In the first state of knowledge, denoted  $I'$ , we know that the neurons in the population are biologically  
 80 or functionally similar, for example in morphology and kind of input or output they receive or give.  
 81 Knowledge of this similarity leads us to assign a probability distribution for the population state  $X$   
 82 that is symmetric under permutations of neuron identities, or *exchangeable* as it is usually called.

83 In the second state of knowledge or ignorance, denoted  $I''$ , we are completely ignorant about the  
 84 physical details of the individual neurons. Our ignorance is therefore symmetric under permutations  
 85 of neuron identities. This also leads to an exchangeable probability distribution for  $X$ .

86 Let us use  $I$  to denote either of these two states of knowledge, in those probabilities that are identical  
 87 for  $I'$  and  $I''$ .

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

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

90 the latter being the probability for the population average  $X$ . A sum is only apparently missing in  
 91 the central term: its summands  $\bar{X} = A$  are all zero except for  $A = \bar{\mathbf{R}}$ . Proof of this theorem and  
 92 generalizations to non-binary and continuum cases are given by de Finetti [18], Kendall [19], Ericson  
 93 [20], Diaconis & Freedman [21; 22], Heath & Sudderth [23]. This theorem is intuitive: owing to  
 94 symmetry, we must assign equal probabilities to all states with  $N\bar{\mathbf{R}}$  active neurons.

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

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

96 with

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

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

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

102 How is it possible that the very different states of knowledge  $I'$  and  $I''$  lead to the same formulae  
 103 above? Their difference appears as soon as we make an experimental observation, say  $X_2 = R_2 \in$   
 104  $\{0, 1\}$  and update our initial probabilities (5):

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

105 The conditional probability and the probability for the average on the right side will update in very  
 106 different ways for  $I'$  and  $I''$ . The discussion of this update process in the two cases is unnecessary  
 107 for the purposes of this note and outside their scope; we hope to address it in full in a future work.

108 The conditional probability  $\Pi(\bar{\mathbf{r}} | \bar{\mathbf{R}})$  relates the spaces of the sample average  $\bar{X} \in \{0, \dots, N\}$  and of  
 109 the population average  $\bar{x} \in \{0, \dots, n\}$  in a special way. It is a coarsening projector of any probability

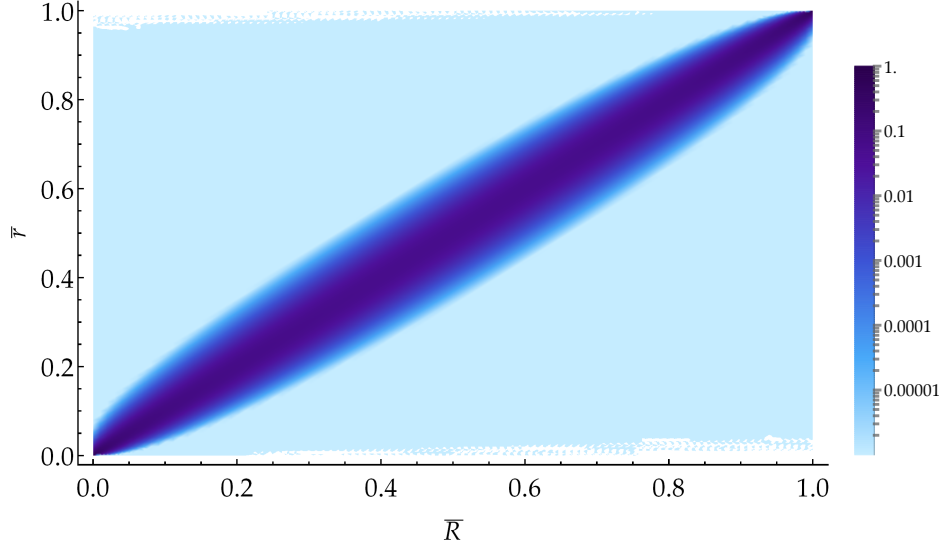


Figure 1: Log-plot of the hypergeometric distribution  $\Pi(\bar{r}|\bar{R}) := \binom{n}{\bar{r}} \binom{N-n}{N\bar{R}-n\bar{r}} \binom{N}{N\bar{R}}^{-1}$  for  $N = 5000$ ,  $n = 200$ . (Band artifacts may appear in the colourbar depending on your PDF viewer.)

110  $p$  for  $\bar{X}$  onto a marginal probability  $p_*$  for  $\bar{x}$ :

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

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

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

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

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

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

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

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

126 Note that formulae (12) say more than the limits  $\mathbb{P}(\bar{x} = a) \rightarrow \mathbb{P}(\bar{X} = a)$  and  $\mathbb{E}[f(\bar{x})] \rightarrow \mathbb{E}[f(\bar{X})]$ , as  
127  $n \rightarrow N$ , do. These limits are trivially valid because the sample becomes the full population as  $n \rightarrow N$ .  
128 In particular, these limits hold even in cases where the conditional probability  $\mathbb{P}(\bar{x} = \bar{r} | \bar{X} = \bar{R})$  is not  
129 a fuzzy identity and our uncertainties about sample and about population can differ wildly.

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

$$\underbrace{(\bar{\mathbf{x}} \cdots \bar{\mathbf{x}})}_{m \text{ factors}}^* = \underbrace{\bar{\mathbf{X}} \cdots \bar{\mathbf{X}}}_{m \text{ factors}}, \quad (13)$$

$$E(\bar{\mathbf{x}} \cdots \bar{\mathbf{x}} | I) = E(\bar{\mathbf{X}} \cdots \bar{\mathbf{X}} | I) = \binom{n}{m}^{-1} \sum_{n\bar{\mathbf{r}}=0}^n \binom{n\bar{\mathbf{r}}}{m} P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | I) = \binom{N}{m}^{-1} \sum_{N\bar{\mathbf{R}}=0}^N \binom{N\bar{\mathbf{R}}}{m} P(\bar{\mathbf{X}} = \bar{\mathbf{R}} | I).$$

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

The probability 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 marginalization (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.

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 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: the expectation 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 distributions for the sample average  $P(\bar{\mathbf{x}} = \bar{\mathbf{r}} | I)$ .

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

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

and then by marginalization (6)

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

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

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

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

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

where  $\zeta$  is a normalization constant and  $\boldsymbol{\lambda}^\top$  are Lagrange multipliers such that

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

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

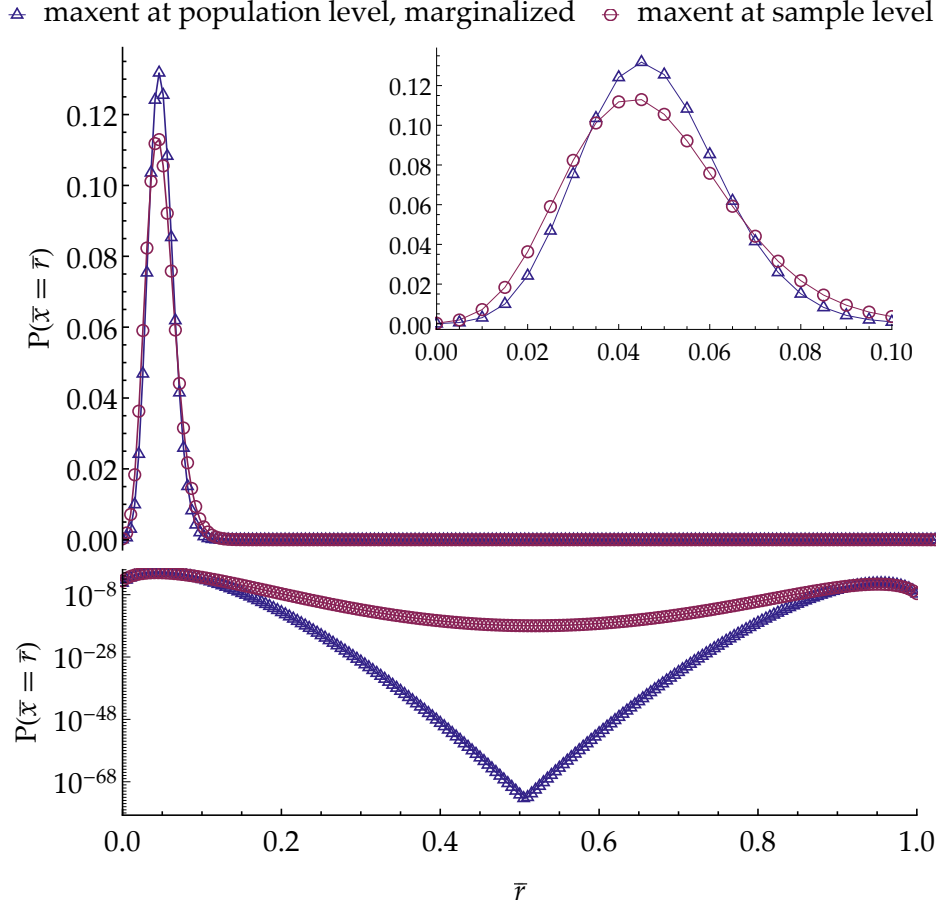


Figure 2: Linear and log-plots of  $P(\bar{x} = \bar{r})$  constructed by maximum-entropy at the population level followed by sample marginalization (blue triangles), eq. (15), and at the sample level (red circles), eq. (17), with  $N = 5000$ ,  $n = 200$ , constraints  $E(\bar{x}) = 0.0478$ ,  $E(\bar{x}\bar{x})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

160 Yet they cannot be exactly equal, because their equality would require the Lagrange multipliers  $\Lambda$   
 161 and  $\lambda$  to satisfy the constraint equations (16), (18), and also  $P(\bar{x} = \bar{r} | I_p) = P(\bar{x} = \bar{r} | I_s)$ ; that is,  
 162  $2m + n$  equations (one normalization is taken care of) in  $2m$  unknowns. A solution can exist, if at all,  
 163 only for very special choices of constraints functions  $f$  and values  $c$ .

164 The sample distribution obtained from maximum-entropy at the sample level will therefore likely  
 165 miss important features present in the one obtained at the population level, like additional modes or  
 166 particular tail behaviour. We show two examples of this discrepancy in figs 2 and 3, for  $N = 5000$ ,  
 167  $n = 200$ , and constraint functions of the form  $f(\bar{x}) = (\bar{x}, \bar{x}\bar{x}, \dots) \equiv (\bar{x}, \binom{n\bar{x}}{2}/\binom{n}{2}, \dots)$ , equivalent to  
 168 moments constraints. The constraint values used in these examples, reported in the figure captions,  
 169 have neurobiologically realistic values [27].

170 In the first example the constraint functions are  $E(\bar{x})$  and  $E(\bar{x}\bar{x})$ . The distribution obtained at the  
 171 sample level is broader than the one obtained at the population level; the tails of the two distributions  
 172 are very different.

173 The second example uses two additional constraint functions  $E(\bar{x}\bar{x}\bar{x})$ ,  $E(\bar{x}\bar{x}\bar{x}\bar{x})$ . The distribution  
 174 obtained at the population level has two modes, replaced by only one in the distribution obtained at  
 175 the sample level; the tails are very different also in this case.

176 How should we apply the maximum-entropy method then? on the sample or on the population?  
 177 Which application is “maximally noncommittal”?

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

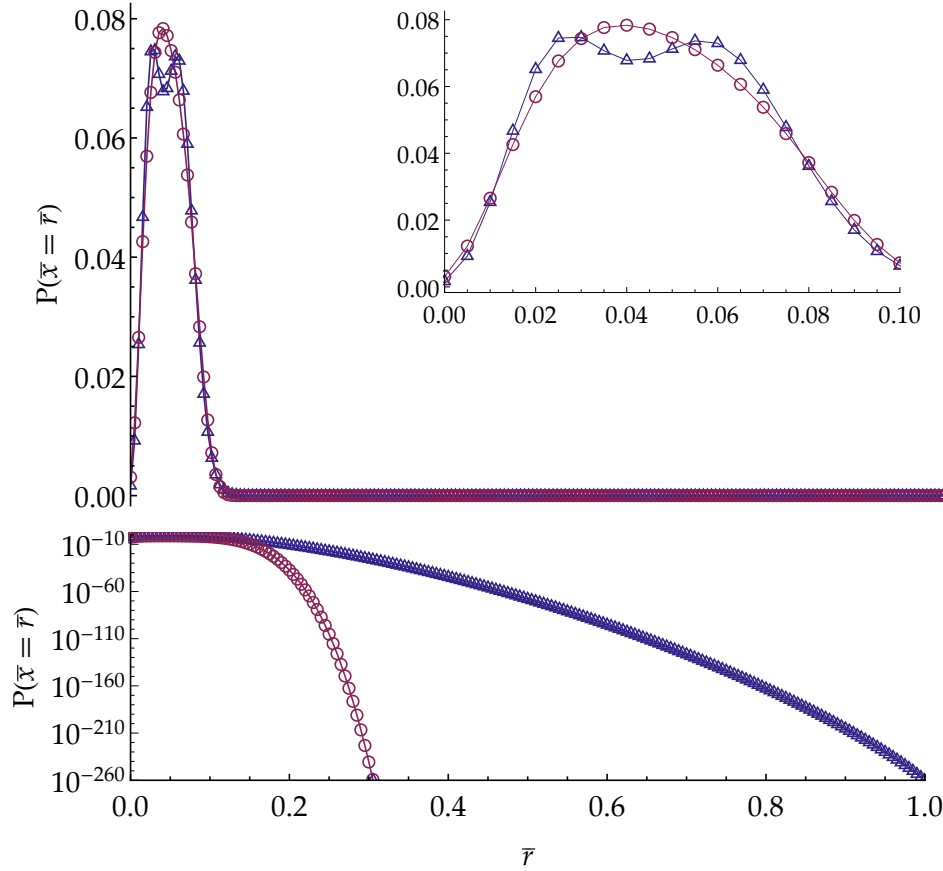


Figure 3: Linear and log-plots of  $P(\bar{x} = \bar{r})$  constructed by maximum-entropy at the population level followed by sample marginalization (blue triangles), eq. (15), and at the sample level (red circles), eq. (17), with  $N = 5000$ ,  $n = 200$ , constraints  $E(\bar{x}) = 0.0478$ ,  $E(\bar{x}\bar{x})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

## 178 5 Discussion

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

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

188 A perfunctory intuitive reasoning seems insufficient for clarifying this point. Let’s express it in  
 189 the language of the probability calculus. Suppose we do not know whether the sample is really  
 190 part of a larger population: we do not know whether  $N = n$  or how large  $N$  is otherwise. Call this  
 191 state of ignorance  $\gamma$ . In the probability calculus this ignorance about  $N$  is expressed by assigning  
 192 a probability distribution  $P(N | \gamma)$  that vanishes if  $N < n$ , since we know that  $N \geq n$ ; see Good  
 193 [28; 29] and Rissanen [30] for examples of such distributions over the integers. Maintaining our  
 194 assumption of symmetric ignorance, probability assignments that do not assume a specific value of  
 195  $N$  are then obtained via multiplication of all  $N$ -dependent probabilities by  $P(N | \gamma)$  and subsequent



196 marginalization over  $N$ . Technically speaking,  $N$  becomes a *nuisance parameter* [16; 31; 32]. The  
 197 probability obtained from maximum-entropy at the population level, eq. (15), then generalizes to

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

198 where  $N$ -dependencies have been made explicit. This is a formidable expression. But our question,  
 199 “is the usual maximum-entropy at the sample level (17) unassuming with regard to the existence of a  
 200 larger population?”, translates now into the precise mathematical question: “are the distributions (19)  
 201 and (17) equal for some choice of  $P(N | \gamma)$ , with  $P(N | \gamma) \neq 0$  for  $N > n$ ?”. We leave this mathematical  
 202 problem for future work. Note, however, that this equality is satisfied if  $P(N = 1 | \gamma) = 1$ , which  
 203 means that the usual maximum-entropy model can also be interpreted as assuming that *no* larger  
 204 population exists.

205 We find the maximum-entropy model constructed at the population level very natural and preferable.  
 206 After all, physical models of neuronal networks usually include some sort of external input to the  
 207 neurons as well, mimicking their embedding in a larger network. The sample distribution given by the  
 208 maximum-entropy model at the population level, when used as a reference distribution for surprise  
 209 analysis, may reveal features in a dataset that were unnoticed by the standard maximum-entropy  
 210 model. The question remains of how to specify  $N$ , though. We have tacitly intended  $N$  as the size of  
 211 the largest biologically or functionally homogeneous population from which our sample was recorded.  
 212 It could be the amount of neurons in a functional brain area, for example the primary visual cortex,  
 213 for which  $N \sim 10^8$  [33]. For large  $N$  – unfortunately we are not yet able to translate this “large” into  
 214 a numeric order of magnitude – the final distribution becomes independent of  $N$ , and continuous  
 215 approximations become available.

216 The possibility of using two different distributions is not a physical contradiction. Similar situations  
 217 arise in statistical mechanics. It is known that if a system is described by a maximum-entropy Gibbs  
 218 state, its subsystems need not be [34]. A dilemma quite similar to ours also appears in the statistical  
 219 description of the final state of a non-equilibrium process starting and ending in two equilibrium  
 220 states: we can describe our knowledge about the final state either by a Gibbs distribution, or by the  
 221 distribution obtained from the Liouville evolution of the Gibbs distribution assigned to the initial  
 222 state. The two descriptions differ – even though the final *physical* state is obviously exactly the same  
 223 [35 § 4]. The two descriptions differs because in one case we can make sharper predictions about the  
 224 state thanks to our knowledge of its preceding dynamics. In this example, though, both distributions  
 225 are usually immensely sharp and practically lead to the same predictions. In the neuroscientific  
 226 applications considered in this note the difference in predictions may be relevant instead.

227 Our analysis touched only constraints of the sample average,  $E[f(\bar{x})]$ . The corresponding models are  
 228 usually called “homogeneous” in the literature. Purely “inhomogeneous” models have also been used  
 229 [6–8], in which expectations for individual neurons or groups of neurons are constrained, for example  
 230  $E(x_2)$  or  $E(x_1 x_8 x_9)$ . A short computation shows that the maximum-entropy method with this kind of  
 231 constraints gives the same result whether applied at the sample or at the population level: the states of  
 232 any unconstrained neurons marginalize out. This is understandable: expressing different uncertainties  
 233 about, say, neurons 2 and 5 we are breaking the symmetry of our uncertainty, which thus cannot be  
 234 representative of other neurons in the sample or in the population. Inhomogeneous models, however,  
 235 require enormous computational power for large sample sizes; homogeneous models therefore retain  
 236 their importance. Our analysis and dilemma also persist for hybrid homogeneous-inhomogeneous  
 237 models [10; 11].

## 238 Acknowledgments

239 To be added after review.

## 240 References

- 241 [1] A. Berényi, Z. Somogyvári, A. J. Nagy, L. Roux, J. D. Long, S. Fujisawa, E. Stark, A. Leonardo et  
 242 al.: *Large-scale, high-density (up to 512 channels) recording of local circuits in behaving animals*. J.  
 243 Neurophysiol. **111**<sup>5</sup> (2014), 1132–1149. [http://www.buzsakilab.com/content/PDFs/Berenyi20](http://www.buzsakilab.com/content/PDFs/Berenyi2013.pdf)  
 244 [13.pdf](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://bayes.wustl.edu/etj/node1.html>, see also ref. [36].
- [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. [37].
- [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. Amodei, 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 [38] (1963), 181–218. Repr. in ref. [39 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 [40] (1959), 1–115. Transl. as ref. [41].
- [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.
- [30] J. Rissanen: *A universal prior for integers and estimation by minimum description length*. Ann. Stat. **11**<sup>2</sup> (1983), 416–431.
- [31] D. V. Lindley: *Introduction to Probability and Statistics from a Bayesian Viewpoint. Part 2: Inference*, reprint. Cambridge University Press, Cambridge (2008). First publ. 1965.
- [32] J.-M. Bernardo, A. F. Smith: *Bayesian Theory*, reprint. Wiley, New York (2000). First publ. 1994.
- [33] G. Leuba, R. Kraftsik: *Changes in volume, surface estimate, three-dimensional shape and total number of neurons of the human primary visual cortex from midgestation until old age*. Anat. Embryol. **190**<sup>4</sup> (1994), 351–366.
- [34] C. Maes, F. Redig, A. Van Moffaert: *The restriction of the Ising model to a layer*. J. Stat. Phys. **96**<sup>1</sup> (1999), 69–107. [arXiv:math/9810094](https://arxiv.org/abs/math/9810094).
- [35] E. T. Jaynes: *Inferential scattering*. (1993). <http://bayes.wustl.edu/etj/node1.html>; extensively rewritten version of a paper first publ. 1985 in ref. [42], pp. 377–398.
- [36] E. T. Jaynes: *Information theory and statistical mechanics. II*. Phys. Rev. **108**<sup>2</sup> (1957), 171–190. <http://bayes.wustl.edu/etj/node1.html>, see also ref. [2].
- [37] J. Shlens, G. D. Field, J. L. Gauthier, M. I. Grivich, D. Petrusca, A. Sher, A. M. Litke, E. J. Chichilnisky: *Correction, the structure of multi-neuron firing patterns in primate retina*. J. Neurosci. **28**<sup>5</sup> (2008), 1246. See ref. [7].
- [38] K. W. Ford, ed.: *Statistical Physics*. Benjamin, New York (1963).
- [39] E. T. Jaynes: *E. T. Jaynes: Papers on Probability, Statistics and Statistical Physics*, reprint. Kluwer, Dordrecht (1989). Ed. by R. D. Rosenkrantz. First publ. 1983.
- [40] B. de Finetti, ed.: *Induzione e statistica*, reprint. Springer, Berlin (2011). First publ. 1959.
- [41] B. de Finetti: *Probability, statistics and induction: their relationship according to the various points of view*. In: de Finetti [43] (1972), Ch. 9, 147–227. Transl. of ref. [18].
- [42] C. R. Smith, W. T. Grandy Jr., eds.: *Maximum-Entropy and Bayesian Methods in Inverse Problems*. D. Reidel, Dordrecht (1985).
- [43] B. de Finetti: *Probability, Induction and Statistics: The art of guessing*. Wiley, London (1972).

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