

# Posteriors for sufficiency hypotheses and maximum-entropy

P.G.L. Porta Mana  
<perio.mano@ntnu.no>

Draft of 7 April 2019 (first drafted 16 February 2019)

Assessing the probability of a hypothesis of sufficiency from the observation of a sample.

## 1 Hypotheses about sufficient statistics

We have a population of  $N$  neurons whose activities we imagine to have time-binned into  $T$  bins and binarized. Denote their total population activity at time bin  $t$  by  $S_t \in \{0, 1, \dots, N\}$ , their total activity at an unspecified bin by  $S$ , and the time sequence of total activities by  $S := (S_{t_1}, S_{t_2}, \dots, S_{t_T})$ .

We have recorded the activities of a sample of  $n$  neurons from the population above. Denote the total activity of the sample at time bin  $t$  by  $s_t \in \{0, 1, \dots, n\}$ , at an unspecified bin by  $s$ , and the time sequence by  $s := (s_{t_1}, s_{t_2}, \dots, s_{t_T})$ .

We don't know how these sampled neurons were chosen from the full population. This fact leads, for each time bin, to the following degree of belief about the activity of the sample if we knew the activity of the full population (Porta Mana et al. 2015 § 2.3; 2018 § 2):

$$p(s | S, I) = \binom{n}{s} \binom{N-n}{S-s} \binom{N}{S}^{-1} =: G_{sS}, \quad (1)$$

namely, a hypergeometric distribution.

Here and in the following  $I$  denotes the proposition stating our background information.

Now suppose that we knew the total activities  $S$  of the *full* population at some  $T$  time bins  $\{t\}$ , and we wanted to infer the total activities  $S'$  at  $T'$  *different* time bins  $\{t'\}$ :

$$p(S' | S, I). \quad (2)$$

We want to consider the hypotheses that *only a specific set of statistics* about our data  $S$  are *relevant* for our inference about  $S'$ ; that is, they

are a sufficient statistics. Any aspect of the data not contained in those statistics would be irrelevant for our inference. This inferential property could be the result of biological properties of the population.

Let's assume that there are  $R$  such statistics (besides  $T$ , which is always part of a set of sufficient statistics). Each statistic is the sum over time of a specific function of the total activity  $S$ . We can arrange these functions in an  $R$ -by- $(N + 1)$  matrix  $\mathbf{C} := (C_{rS})$ , where  $C_{rS}$  is the value of the function for the  $r$ th statistic when the total activity is  $S$ . The  $R$  sufficient statistics for the data  $\mathbf{S}$  would thus be

$$\bar{C}_r := \frac{1}{T} \sum_t C_{rS_t}, \quad r \in \{1, \dots, R\}. \quad (3)$$

Our goal is to quantify our uncertainty about these hypotheses of sufficient statistics, given the activity data from a sample of neurons. It's important to note that the hypotheses we must consider are not discrete or of a yes-or-no type: they form a continuum. This is because we have a continuum of degrees of relevance. Consider for example two statistics  $\bar{C}_1$  and  $\bar{C}_2$  from the bins  $\{t\}$ . Our degrees of belief about the activities at bins  $\{t'\}$  are

$$p(S' | \bar{C}_1, \bar{C}_2, I). \quad (4)$$

It may happen that lack of knowledge about  $\bar{C}_2$  doesn't change our degree of belief:

$$p(S' | \bar{C}_1, I) = p(S' | \bar{C}_1, \bar{C}_2, I), \quad (5)$$

in which case  $\bar{C}_2$  is irrelevant. It may also happen that our degree of belief is changed but in a negligible way, for all values of  $S'$  and  $\bar{C}_1$ :

$$p(S' | \bar{C}_1, I) \approx p(S' | \bar{C}_1, \bar{C}_2, I), \quad (6)$$

so that  $\bar{C}_2$  could be dropped in practice. We can imagine larger and larger changes to the point where dropping  $\bar{C}_2$  would lead to drastically different degrees of belief. The question of the relevance of  $\bar{C}_2$  is therefore not dichotomous. We will thus deal with a continuum of hypotheses, each representing a degree of relevance of some statistics. We shall shortly see how to mathematically represent this continuum of hypotheses.

## 2 The Koopman-Pitman theorem

How does a hypothesis about a sufficient statistic affect our degrees of belief? The answer comes from the Koopman-Pitman theorem (Koopman 1936; Pitman 1936; see also Darmais 1935; Barankin et al. 1963; Denny 1967; Hipp 1974; Lauritzen 1974a; 1984; 1988; for the discrete version: Fraser 1963; Andersen 1970), which says that the degree of belief (2) has a very specific mathematical expression if only some statistics of  $S$  are relevant. The main statement of the theorem is this: if  $R$  sufficient statistics are given by functions  $C_{rS}$ , then for any number of time bins  $T$

$$p(S | I) = \int d\lambda \, p(\lambda | I) \prod_t p(S_t | \lambda, I) \quad (7a)$$

with

$$p(S | \lambda, I) := \frac{g^S}{Z(\lambda)} \exp(\sum_r \lambda_r C_{rS}), \quad (7b)$$

$$\lambda := (\lambda_1, \dots, \lambda_R) \in \mathbf{R}^R, \quad Z(\lambda) := \sum_S g^S \exp(\sum_r \lambda_r C_{rS}), \quad (7c)$$

and  $g$  a positive function of  $S$ .

Some important remarks about the Pitman-Koopman formula (7):

- a. A hypothesis that only stated what the sufficient statistics are would not determine the density  $p(\lambda | I)$  or the function  $g$  in the formula above. The hypotheses we are going to compare thus contain additional information besides sufficiency.
- b. The formula can be interpreted this way: our degree of belief about  $S$  is given by the degree of belief we would have if we knew the values of the sufficient statistics  $\bar{C}$  for an unlimited number of time bins, mixed over our uncertainty about the values themselves:

$$p(S | I) = \int d\bar{C} \, p(\bar{C} | I) p(S | \bar{C}, I). \quad (8)$$

- c. Formula (7) is obtained from (8) by a one-to-one reparametrization:

$$\bar{C}_r(\lambda) = \sum_S C_{rS} \frac{g^S}{Z(\lambda)} \exp(\sum_r \lambda_r C_{rS}) \equiv \partial_{\lambda_r} \ln Z(\lambda), \quad (9)$$

$$p(\bar{C} | I) d\bar{C} = p(\lambda | I) d\lambda, \quad (10)$$

$$p[S | \bar{C}(\lambda), I] = p(S | \lambda, I) \equiv \frac{g^S}{Z(\lambda)} \exp(\sum_r \lambda_r C_{rS}). \quad (11)$$

Equation (9) cannot be solved explicitly for  $\lambda$  in terms of  $\bar{C}$  except for very simple cases. The parametrization in terms of  $\lambda$  has several special properties:

- c.1. The quantities  $(\lambda_r)$  can assume any values independently of one another, whereas the limit statistics  $(\bar{C}_r)$  have interdependent ranges.
- c.2. The expression for  $p(S \mid \bar{C}, I)$  can be written explicitly in terms of  $\lambda$ , but not in terms of  $\bar{C}$ .
- c.3. If some  $\lambda_r$  vanishes then the corresponding statistic is *irrelevant* – the corresponding term indeed disappears from the exponential in formula (7).

Remark c. suggests that the absolute value of each parameter,  $|\lambda_r|$ , might be used to quantify the degree of relevance of the corresponding statistic, zero meaning complete irrelevance. These absolute values can therefore be used as parameters for our continuum of hypotheses about the sufficient statistics. The degree of belief we ultimately want to quantify is thus

$$p(|\lambda| \mid \text{data}, I) \quad (12)$$

and its marginals for the various  $|\lambda_r|$ .

### 3 Conditionalization on full-population data

First let us recapitulate the expression for

$$p(\lambda \mid S, I) \propto p(\lambda \mid I) p(S \mid \lambda, I). \quad (13)$$

Introduce the relative frequencies  $F := (F_S)$  with which the activity  $S$  appears during the  $T$  time bins, given by

$$F_S = \frac{1}{T} \sum_t \delta(S_t - S). \quad (14)$$

From formula (7), bringing several terms within the exponential, we can rewrite, using the frequencies,

$$\begin{aligned} p(S \mid \lambda, I) &= \prod_t p(S_t \mid \lambda, I) \\ &= \prod_S [p(S \mid \lambda, I)]^{TF_S} \\ &= \exp\left(-T\{H[F, p(S \mid \lambda, I)] - H(F)\}\right), \end{aligned} \quad (15)$$

where  $H(a, b)$  is the relative entropy of  $a$  with respect to  $b$  and  $H(a)$  is the Shannon entropy of  $a$ .

The last expression shows that: (a) our degree of belief about the sequence  $S$  depends only on the frequencies  $F$  of that sequence; (b) if  $p(\lambda | I)$  is constant, the mode of the density (13) for  $\lambda$  can be calculated by minimizing with respect to  $\lambda$  the relative entropy between  $F$  and  $p(S | \lambda, I)$ ; (c) owing to the mathematical form (7b) of  $p(S | \lambda, I)$ , minimization of the relative entropy leads to the classical maximum-entropy equations

$$E(C_{rS} | \lambda, I) \equiv \partial_{\lambda_r} \ln Z(\lambda) = \bar{C}_r \equiv \sum_S C_{rS} F_S. \quad (16)$$

#### 4 Conditionalization on sample data

The data conditional on which we want to calculate the degree of belief (12) do not involve the full population, though, but only a sample of it: they are some statistics

$$\hat{c} := \frac{1}{T} \sum_t c_{s_t} \equiv \sum_s c_s f_s \quad (17)$$

of a sequence  $s$  of recorded sample activities with relative frequencies  $f$ . The statistic  $c_s$  can be vector-valued.

We use Bayes's theorem:

$$p(\lambda | \hat{c}, I) \propto p(\lambda | I) p(\hat{c} | \lambda, I) \quad (18)$$

and focus our attention on the last term. It can be written as the sum of our degrees of belief for the frequencies satisfying the statistics:

$$p(\hat{c} | \lambda, I) = \sum_f \delta(\sum_s c_s f_s = \hat{c}) p(f | \lambda, I). \quad (19)$$

To calculate our degree of belief about the frequencies, we note that all sequences  $s$  having frequencies  $f$  have the same degree of belief under

the hypothesis of sufficiency *for the full population* and assuming the long-run statistics are known – that is, with  $\lambda$  given:

$$p(s \mid \lambda, I) = \prod_t p(s_t \mid \lambda, I) = \prod_s p(s \mid \lambda, I)^{T f_s}. \quad (20)$$

The first, joint degree of belief factorizes because the degree of belief for each  $s_t$  is independent of the other  $t$  if the corresponding  $S_t$  is known, and the joint degree of belief for all the  $S_t$  factorizes if  $\lambda$  is known.

Our degree of belief  $p(s \mid \lambda, I)$  about the sample activity at any time, given  $\lambda$ , can be obtained by marginalizing the corresponding  $S$ :

$$p(s \mid \lambda, I) = \sum_S p(s \mid S, \lambda, I) p(S \mid \lambda, I) = \sum_S G_{sS} p(S \mid \lambda, I), \quad (21)$$

since the first degree of belief in the sum is independent on  $\lambda$ ; the second is given by eq. (7b).

Our degree of belief  $p(f \mid \lambda, I)$  about the relative frequencies of sample activities is then proportional to the degree of belief above, times the number of sequences having those frequencies:

$$\begin{aligned} p(f \mid \lambda, I) &= \binom{T}{f} p(s \mid \lambda, I) = \binom{T}{f} \prod_s \left[ \sum_S G_{sS} p(S \mid \lambda, I) \right]^{T f_s} \\ &\approx \exp \left\{ T H \left[ f, \left( \sum_S G_{sS} p(S \mid \lambda, I) \right)_s \right] \right\} \end{aligned} \quad (22)$$

The last expression valid for enough large  $T$ .

Combining eqs (19) and (22) we obtain

$$p(\hat{c} \mid \lambda, I) \approx \sum_f \delta \left( \sum_s c_s f_s = \hat{c} \right) \exp \left\{ T H \left[ f, \left( \sum_S G_{sS} p(S \mid \lambda, I) \right)_s \right] \right\}. \quad (23)$$

When  $T$  is very large compared to the number of possible values of  $s$ , the expression above becomes negligible for all  $\lambda$  such that  $\sum_S G_{sS} p(S \mid \lambda, I)$  equal some set of frequencies  $f_s$  that satisfy the constraints (17). It can therefore be approximated as

$$\begin{aligned} p(\hat{c} \mid \lambda, I) &\approx \delta \left[ \sum_s c_s \sum_S G_{sS} p(S \mid \lambda, I) = \hat{c} \right] \equiv \\ &\delta \left[ \sum_s \sum_S c_s G_{sS} \frac{g_s}{Z(\lambda)} \exp(\lambda C_s) = \hat{c} \right]. \end{aligned} \quad (24)$$

This equation differ from the constraint equation obtained with the usual maximum-entropy method: since

$$\sum_s c_s G_{sS} \neq C_S, \quad (25)$$

we are not seeking  $\lambda$  for which the expectation of *its conjugate statistics*  $C_S$  satisfies a constraint, but for which the expectation of a different statistics satisfies a constraint. The only exception is when

$$\sum_s c_s G_{sS} = C_S, \quad (26)$$

which happens, for example, for the factorial moments.

\*\*\*

As in the previous section, if the density  $p(\lambda | I)$  for  $\lambda$  is uniform then from eq. (18) the mode of  $p(\lambda | s, I)$  can be obtained by maximizing  $p(f | \lambda, I)$ , which amount to maximizing the relative entropy between  $f$  and  $(\sum_s G_{sS} p(S | \lambda, I))_s$  with respect to  $\lambda$ . We have the maximum of the latter for the  $\lambda$  satisfying

$$\sum_s G_{sS} \frac{g_s}{Z(\lambda)} \exp(\sum_r \lambda_r C_{rS}) = f_s \quad \forall s. \quad (27)$$

minimizing the relative entropy between  $f$  and  $p(s | \lambda, I)$ . In the present case, however, the expression for the latter probability, eq. (28b), does not lead to the classical maximum-entropy equations. We find instead

$$\partial_{\lambda_r} \ln Z(\lambda) = \hat{c}_r(\lambda) \quad (28a)$$

with

$$\hat{c}_r(\lambda) := \sum_s f_s c_{rs}(\lambda), \quad c_{rs}(\lambda) := \frac{\sum_s C_{rS} G_{sS} p(S | \lambda, I)}{\sum_s G_{sS} p(S | \lambda, I)}. \quad (28b)$$

We can calculate it by marginalizing with respect to all possible *sequences* of activities  $S$  for the full population, using the hypergeometric distribution (1) and the Koopman-Pitman formula (7):

$$p(s | I) = \sum_S p(s | S, I) p(S | I) = \int d\lambda p(\lambda | I) \prod_t p(s_t | \lambda, I) \quad (29a)$$

with

$$p(s | \lambda, I) := \sum_S G_{sS} \frac{g_s}{Z(\lambda)} \exp(\sum_r \lambda_r C_{rS}). \quad (29b)$$

This formula shows that our degree of belief about the sample activities does *not* have a sufficient statistics. This fact is mathematically similar to

what happens in statistical mechanics: if our uncertainty about a system of particles is expressed by a Gibbs distribution, then our uncertainty about a subsystem won't generally be of a Gibbsian type (Maes et al. 1999).

We can now calculate the density for  $\lambda$  given the sequence of sample activities  $s$  using Bayes's theorem. We can introduce the relative frequencies  $f$  for the sample activities and proceed as in § 3, obtaining

$$\begin{aligned} p(s \mid \lambda, I) &= \prod_t p(s_t \mid \lambda, I) \\ &= \prod_s [p(s \mid \lambda, I)]^{T f_s} \\ &= \exp\left(-T\{H[f, p(s \mid \lambda, I)] - H(f)\}\right). \end{aligned} \quad (30)$$

As in the previous section, if the density  $p(\lambda \mid I)$  for  $\lambda$  is uniform then the mode of  $p(\lambda \mid s, I)$  can be obtained by minimizing the relative entropy between  $f$  and  $p(s \mid \lambda, I)$ . In the present case, however, the expression for the latter probability, eq. (28b), does not lead to the classical maximum-entropy equations. We find instead

$$\partial_{\lambda_r} \ln Z(\lambda) = \hat{c}_r(\lambda) \quad (31a)$$

with

$$\hat{c}_r(\lambda) := \sum_s f_s c_{rs}(\lambda), \quad c_{rs}(\lambda) := \frac{\sum_S C_{rS} G_{sS} p(S \mid \lambda, I)}{\sum_S G_{sS} p(S \mid \lambda, I)}. \quad (31b)$$

or more explicitly

$$\frac{\sum_S C_{rS} g_S \exp(\sum_r \lambda_r C_{rS})}{\sum_S g_S \exp(\sum_r \lambda_r C_{rS})} = \sum_s f_s \frac{\sum_S C_{rS} G_{sS} g_S \exp(\sum_r \lambda_r C_{rS})}{\sum_S G_{sS} g_S \exp(\sum_r \lambda_r C_{rS})}. \quad (32)$$

The expressions (30) define  $r$  functions  $c_{rs}(\lambda)$  of the sample activity  $s$  and of  $\lambda$  that have a role analogous to the statistic  $C_{rS}$  of the full-population activity  $S$ , in the sense that

$$\begin{aligned} E[c_{rs}(\lambda) \mid \lambda, I] &= E(C_{rS} \mid \lambda, I), \\ \sum_s c_{rs}(\lambda) p(s \mid \lambda, I) &= \sum_S C_{rS} p(S \mid \lambda, I), \end{aligned} \quad (33)$$

as can be verified by substitution and the definition (1) of  $G_{sS}$ .



The mode  $\lambda$  is given by

$$\arg \inf_{\lambda} \left\{ \ln Z(\lambda) - \sum_s f_s \ln \left[ \sum_S G_{sS} g_S \exp(\sum_r \lambda_r C_{rS}) \right] \right\}. \quad (34)$$

Compare this with the standard maximum-entropy case (Mead et al. 1984), which can be written as

$$\arg \inf_{\lambda} \left[ \ln Z(\lambda) - \sum_S F_S \ln \exp(\sum_r \lambda_r C_{rS}) \right]. \quad (35)$$

### Important differences from the ‘dilemma’ paper:

Consider the following statistics:

$$C_{rS} := \binom{S}{r} / \binom{N}{r} \quad r \in 1, 2, \dots \quad (36)$$

In particular,  $C_{1S} = S/N$  and  $C_{rS}$  is the number of  $r$ -tuples of simultaneously active neurons divided by the number of possible ones. Analogous statistics  $\binom{s}{r} / \binom{n}{r}$  can be considered for the sample.

These statistics have this special property (Porta Mana et al. 2015; 2018):

$$\sum_s \binom{s}{r} / \binom{n}{r} p(s) = \sum_S \binom{S}{r} / \binom{N}{r} p(S), \quad (37)$$

no matter what the degrees of belief about full-population and sample might be, provided that they are related by  $p(s) = \sum_S G_{sS} p(S)$ .

Note, however, that given  $p(S)$  and any statistic  $C_{rS}$  it is always possible to create a function  $c_{rS}$  of the sample activity that satisfies

$$\sum_s c_{rs} p(s) = \sum_S C_{rS} p(S) \quad (38)$$

namely,

$$c_{rs} := \frac{\sum_S C_{rS} G_{sS} p(S)}{\sum_S G_{sS} p(S)}. \quad (39)$$

But this function *depends on the specific*  $p(S)$ .

Now consider the minimization of the relative entropy in eq. (29), corresponding to eqs (30), for statistics given by (35). Remember that the assumption of sufficiency asymptotically gives zero probability to observing relative frequencies of full population and sample that lie outside a particular  $R$ -dimensional set.

We have two possibilities when  $T$  is very large:

1. If the frequencies  $f$  observed for the sample are among those admitted by the sufficiency hypothesis, then there is a  $\lambda^*$  for which

$$f_s \approx \sum_S G_{sS} p(S | \lambda^*, I). \quad (40)$$

Moreover, for that  $\lambda$  we also have

$$\sum_s f_s \binom{s}{r} / \binom{n}{r} \approx \sum_s f_s c_{rs}(\lambda^*). \quad (41)$$

and the result of the relative-entropy minimization is equivalent to finding the maximum-entropy distribution for the full population with expectations for the statistics (35) constrained to equal the empirical ones observed in the sample.

2. If the frequencies  $f$  observed for the sample are *not* among those admitted by the sufficiency hypothesis, then the minimization of the relative entropy yields a distribution for the full-population that does *not* satisfy those constraints.

## Bibliography

- (‘de  $X$ ’ is listed under D, ‘van  $X$ ’ under V, and so on, regardless of national conventions.)
- Andersen, E. B. (1970): *Sufficiency and exponential families for discrete sample spaces*. J. Am. Stat. Assoc. **65**<sup>331</sup>, 1248–1255.
- Barankin, E. W., Maitra, A. P. (1963): *Generalization of the Fisher-Darmois-Koopman-Pitman theorem on sufficient statistics*. Sankhyā A **25**<sup>3</sup>, 217–244.
- Barndorff-Nielsen, O. E., Blæsild, P., Schou, G., eds. (1974): *Proceedings of Conference on Foundational Questions in Statistical Inference: Aarhus, May 7–12, 1973*. (University of Aarhus, Aarhus).
- Barndorff-Nielsen, O. E., Dawid, A. P., Diaconis, P., Johansen, S., Lauritzen, S. L. (1984): *Discussion of Steffen Lauritzen’s paper [‘Extreme Point Models in Statistics’]*. Scand. J. Statist. **11**<sup>2</sup>, 83–91. See Lauritzen (1984).
- Darmois, G. (1935): *Sur les lois de probabilité à estimation exhaustive*. Comptes rendus hebdomadaires des séances de l’Académie des sciences **200**, 1265–1266.
- Denny, J. L. (1967): *Sufficient conditions for a family of probabilities to be exponential*. Proc. Natl. Acad. Sci. (USA) **57**<sup>5</sup>, 1184–1187.
- Fraser, D. A. S. (1963): *On sufficiency and the exponential family*. J. Roy. Stat. Soc. B **25**<sup>1</sup>, 115–123.
- Hipp, C. (1974): *Sufficient statistics and exponential families*. Ann. Stat. **2**<sup>6</sup>, 1283–1292.
- Koopman, B. O. (1936): *On distributions admitting a sufficient statistic*. Trans. Am. Math. Soc. **39**<sup>3</sup>, 399–409.
- Lauritzen, S. L. (1974a): *Sufficiency, prediction and extreme models*. In: Barndorff-Nielsen, Blæsild, Schou (1974), 249–269. With discussion. Repr. without discussion in Lauritzen (1974b).
- (1974b): *Sufficiency, prediction and extreme models*. Scand. J. Statist. **1**<sup>3</sup>, 128–134.
- (1984): *Extreme point models in statistics*. Scand. J. Statist. **11**<sup>2</sup>, 65–83. See also discussion and reply in Barndorff-Nielsen, Dawid, Diaconis, Johansen, Lauritzen (1984).
- (1988): *Extremal Families and Systems of Sufficient Statistics*. (Springer, Berlin). First publ. 1982.
- Maes, C., Redig, F., Van Moffaert, A. (1999): *The restriction of the Ising model to a layer*. J. Stat. Phys. **96**<sup>1</sup>, 69–107.
- Mead, L. R., Papanicolaou, N. (1984): *Maximum entropy in the problem of moments*. J. Math. Phys. **25**<sup>8</sup>, 2404–2417. <http://bayes.wustl.edu/Manual/MeadPapanicolaou.pdf>.
- Pitman, E. J. G. (1936): *Sufficient statistics and intrinsic accuracy*. Math. Proc. Camb. Phil. Soc. **32**<sup>4</sup>, 567–579.
- Porta Mana, P. G. L., Rostami, V., Torre, E., Roudi, Y. (2018): *Maximum-entropy and representative samples of neuronal activity: a dilemma*. Open Science Framework doi:10.17605/osf.io/uz29n, bioRxiv doi:10.1101/329193, arXiv:1805.09084.
- Porta Mana, P. G. L., Torre, E., Rostami, V. (2015): *Inferences from a network to a subnetwork and vice versa under an assumption of symmetry*. bioRxiv doi:10.1101/034199.