

Posteriors for sufficiency hypotheses and maximum-entropy

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

Draft of 23 February 2019 (first drafted 16 February 2019)

Note: Dear Reader & Peer, this manuscript is being peer-reviewed by you. Thank you.

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 $S(\cdot) := (S(1), S(2), \dots, S(T))$.

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 $s(\cdot) := (s(1), s(2), \dots, s(T))$.

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):

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

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

$$p[\{S(t')\} | \{S(t)\}, I]. \quad (2)$$

We want to consider the hypotheses that *only a specific set of statistics* about our data at $\{t\}$ are *relevant* for our inference about $\{t'\}$; 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). 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 data at $\{t\}$ 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. Let's make this goal more precise and note some important points:

1. 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

$$p[\{S(t')\} | \bar{C}_1, \bar{C}_2, I]. \quad (4)$$

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

$$p[\{S(t')\} | \bar{C}_1, I] = p[\{S(t')\} | \bar{C}_1, \bar{C}_2, I], \quad (5)$$

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 $\{S(t')\}$ and \bar{C}_1 :

$$p[\{S(t')\} | \bar{C}_1, I] \approx p[\{S(t')\} | \bar{C}_1, \bar{C}_2, I], \quad (6)$$

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.

2. 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 Darrois 1935; Barankin et al.

1963; Denny 1967; Hipp 1974; Lauritzen 1974; 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 $\{S(t')\}$ are relevant. The main statement of the theorem is this: if R sufficient statistics are given by functions C_{rS} , then

$$p[\{S(t)\} | I] = \int d\lambda \, p(\lambda | I) \prod_t \frac{g_{S(t)}}{Z(\lambda)} \exp[\sum_r \lambda_r C_{rS(t)}] \quad (7a)$$

with

$$\lambda := (\lambda_1, \dots, \lambda_R), \quad Z(\lambda) := \sum_S g_S \exp(\sum_r \lambda_r C_{rS}). \quad (7b)$$

It is important to note that a hypothesis solely about a sufficient statistics does not determine the density $p(\lambda | I)$ or the function g in the formula above.

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**³³¹, 1248–1255.
- Barankin, E. W., Maitra, A. P. (1963): *Generalization of the Fisher-Darmois-Koopman-Pitman theorem on sufficient statistics*. Sankhyā A **25**³, 217–244.
- 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**⁵, 1184–1187.
- Fraser, D. A. S. (1963): *On sufficiency and the exponential family*. J. Roy. Stat. Soc. B **25**¹, 115–123.
- Hipp, C. (1974): *Sufficient statistics and exponential families*. Ann. Stat. **2**⁶, 1283–1292.
- Koopman, B. O. (1936): *On distributions admitting a sufficient statistic*. Trans. Am. Math. Soc. **39**³, 399–409.
- Lauritzen, S. L. (1974): *Sufficiency, prediction and extreme models*. In: **barndorffnielsenetal1974**, 249–269. With discussion. Repr. without discussion in **lauritzen1974_r1974**.
- (1984): *Extreme point models in statistics*. Scand. J. Statist. **11**², 65–83. See also discussion and reply in **barndorffnielsenetal1984**.
- (1988): *Extremal Families and Systems of Sufficient Statistics*. (Springer, Berlin). First publ. 1982.
- Pitman, E. J. G. (1936): *Sufficient statistics and intrinsic accuracy*. Math. Proc. Camb. Phil. Soc. **32**⁴, 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.