# A  Derivation of the probability distribution and of its approximations

## A.1  Notation

$N = 1000$: size of the full network.

$n = 65$: size of the sample.

$T = 417\,641$: number of bins in the recording.

Bins are indexed by $t \in \{1, \dots, T\}$.

$A \in \{0, 1, \dots, N\}$: levels of total activity of the full network.

$A_t \in \{0, 1, \dots, N\}$: total activity of the full network at bin $t$.

$\overline{A} := (A_1, A_2, \dots, A_T)$: sequence of activities of full network during the recording.

$F_A \in \{0, 1/T, \dots, 1\}$: relative frequency of activity level $A$ of full network during the recording. That is, activity level $A$ appears in $TF_A$ out of $T$ bins: $TF_A = \sum_t \delta(A_t = A)$.

$\boldsymbol{F} := (F_0, F_1, \dots, F_N)$: frequency distribution.

Note that $\overline{A}$ completely determines $\boldsymbol{F}$. We write the $\boldsymbol{F}$ determined by $\overline{A}$ as $\boldsymbol{F}(\overline{A})$.

$\nu_A \in [0, 1]$: long-run relative frequency of activity level $A$ of full network during all hypothetical repetitions of the same recording in the same conditions.

$a \in \{0, 1, \dots, n\}$: levels of total activity of the sample.

$a_t \in \{0, 1, \dots, n\}$: total activity of the sample at bin $t$.

$\bar{a} := (a_1, a_2, \dots, a_T)$: sequence of activities of sample during the recording.

$f_a \in \{0, 1/T, \dots, 1\}$: relative frequency of activity level $a$ of sample during the recording.

$\boldsymbol{f} := (f_0, f_1, \dots, f_n)$: frequency distribution.

Note that $\bar{a}$ completely determines $\boldsymbol{f}$. We write the $\boldsymbol{f}$ determined by $\bar{a}$ as $\boldsymbol{f}(\bar{a})$.

$J_{aA} \in \{0, 1/T, \dots, 1\}$: joint relative frequency of activity levels $a$ and $A$ for sample and full network during the recording. That is, the pair $(a, A)$ appears in $TJ_{aA}$ out of $T$ bins: $TJ_{aA} = \sum_t \delta(a_t = a)\,\delta(A_t = A)$.

Note that the full network can't have fewer active neurons or fewer silent neurons than its sample (for example, if $A_t = 0$ then $a_t = 0$, and if $A_t = N$ then $a_t = n$), so $J_{a\,A} \equiv 0$ for $a > A$ or $n - a > N - A$.

Obviously $F_A \equiv \sum_a J_{a\,A}$ and $f_a \equiv \sum_A J_{a\,A}$.

$\boldsymbol{J} := (J_{00}, J_{01}, \dots, J_{n\,N})$: joint frequency distribution.

Note that $(\bar{a}, \bar{A})$ together completely determine $\boldsymbol{J}$. We write the $\boldsymbol{J}$ determined by $(\bar{a}, \bar{A})$ as $\boldsymbol{J}(\bar{a}, \bar{A})$.

$G_{a\,A} := \binom{n}{a}\binom{N-n}{A-a}\binom{N}{A}^{-1} \equiv \binom{A}{a}\binom{N-A}{n-a}\binom{N}{n}^{-1}$: hypergeometric distribution.

$I$: background information and assumptions; includes knowledge of $N$ and $n$.

## A.2   Derivation

We want $p(\boldsymbol{F} \mid \bar{a}, I)$. This can be obtained as the marginal of $p(\boldsymbol{J} \mid \bar{a}, I)$. If this probability distribution is represented by a set of Monte Carlo samples $\{\boldsymbol{J}^{(i)}\}$, then $\{\boldsymbol{F}^{(i)} \equiv \sum_a J_{a\,A}^{(i)}\}$ are automatically samples of $p(\boldsymbol{F} \mid \bar{a}, I)$.

By the theorem of total probability,

$$p(\boldsymbol{J} \mid \bar{a}, I) = \sum_{\bar{A}} p(\boldsymbol{J} \mid \bar{a}, \bar{A}, I)\, p(\bar{A} \mid \bar{a}, I), \tag{35}$$

the sum being over all possible sequences $\{\bar{A}\}$ of total activities.

The first factor is a singular distribution:

$$\begin{aligned} p(\boldsymbol{J} \mid \bar{a}, \bar{A}, I) &= \delta\big[\boldsymbol{J} = \boldsymbol{J}(\bar{a}, \bar{A})\big] \\ &\equiv \prod_{a\,A} \delta\big[T J_{a\,A} = \textstyle\sum_t \delta(a_t = a)\,\delta(A_t = A)\big]. \end{aligned} \tag{36}$$

The second factor is derived from Bayes's theorem:

$$p(\bar{A} \mid \bar{a}, I) = \frac{p(\bar{a} \mid \bar{A}, I)\, p(\bar{A} \mid I)}{\sum_{\bar{A}} p(\bar{a} \mid \bar{A}, I)\, p(\bar{A} \mid I)}. \tag{37}$$

In the last formula, let's first derive $p(\bar{a} \mid \bar{A}, I)$. We make this

**assumption**: if $A_t$ is known, then knowledge of $a_{t'}$ or $A_{t'}$ with $t' \neq t$ is irrelevant for our inferences about $a_t$. $\tag{38}$

It isn't a realistic assumption, but it's the one behind the maximum-entropy method in this particular application. From this assumption we have, using the product rule,

$$p(\bar{a} \mid \bar{A}, I) = \prod_t p(a_t \mid A_t, I), \tag{39}$$

and from simple sampling theory $p(a_t \mid A_t, I) = G_{a_t A_t}$ for each $t$, so that

$$p(\bar{a} \mid \bar{A}, I) = \prod_t G_{a_t A_t}. \tag{40}$$

The latter product can be rewritten by grouping together all $t$ in which the same $(a, A)$ pair appears: there are $T J_{a A}(\bar{a}, \bar{A})$ such $t$. Then we consider all such pairs:

$$p(\bar{a} \mid \bar{A}, I) = \prod_{a,A} G_{a A}{}^{T J_{a A}(\bar{a}, \bar{A})} \tag{41}$$

where $\prod_{a,A} \coloneqq \prod_a \prod_A$.

Now let's derive $p(\bar{A} \mid I)$ in formula (37). We assume that the value of this probability has the same value for all sequences $\bar{A}$ having the same frequency distribution $F(\bar{A})$; that is, it functionally depends on $\bar{A}$ only through $F(\bar{A})$. More specifically we assume it can be written hierarchically as

$$p(\bar{A} \mid I) = \pi(F) = \int d\boldsymbol{v} \, q(\boldsymbol{v}) \prod_A v_A{}^{T F_A} \tag{42}$$

for some density function $q(\boldsymbol{v}) \, d\boldsymbol{v}$. The integral expression is simply a mathematical way to write $\pi(F)$, and doesn't need to be further interpreted (if we can solve the integral analytically, it just disappears). But it's also the formula for infinite exchangeability, and from this point of view $\boldsymbol{v}$ can be interpreted as the long-run frequency distribution of the activities in all experiments performed in the same conditions (imagining to joint together their recording times). This is the point of view underlying the maximum-entropy method.

Replacing (41) and (42) into (37) and rearranging we find

$$
\begin{aligned}
p(\bar{A} \mid \bar{a}, I) &= \frac{1}{Z} \int d\boldsymbol{v} \, q(\boldsymbol{v}) \prod_{a,A} \left[ G_{a A}{}^{T J_{a A}(\bar{a}, \bar{A})} \right] \prod_A v_A{}^{T F_A} \\
&\equiv \frac{1}{Z} \int d\boldsymbol{v} \, q(\boldsymbol{v}) \prod_{a,A} (G_{a A} \, v_A)^{T J_{a A}(\bar{a}, \bar{A})},
\end{aligned}
\tag{43}
$$

where $Z := \sum_{\overline{A}} \mathrm{p}(\bar{a} \mid \overline{A}, I)\, \mathrm{p}(\overline{A} \mid I)$ is the normalization factor and the second equality comes from rewriting

$$\nu_A{}^{TF_A} \equiv \nu_A{}^{T \sum_a J_{aA}} \equiv \prod_a \nu_A{}^{T J_{aA}}. \tag{44}$$

Note that the normalization factor is independent of $\overline{A}$ and only depends on $\bar{a}$, which is given and fixed in our inference.

We can now return to our initial probability distribution (35). Replacing (36) and (43) into it and rearranging we have

$$\mathrm{p}(J \mid \bar{a}, I) = \frac{1}{Z} \int \mathrm{d}\nu\, q(\nu) \sum_{\overline{A}} \delta[J = J(\bar{a}, \overline{A})] \prod_{a,A} \left( G_{aA}\, \nu_A \right)^{T J_{aA}(\bar{a}, \overline{A})}. \tag{45}$$

This expression can be considered as the marginal distribution of the joint density function

$$\mathrm{p}(J, \nu \mid \bar{a}, I) = \frac{1}{Z}\, q(\nu) \sum_{\overline{A}} \delta[J = J(\bar{a}, \overline{A})] \prod_{a,A} \left( G_{aA}\, \nu_A \right)^{T J_{aA}(\bar{a}, \overline{A})}. \tag{46}$$

If we represent this density by a set of Monte Carlo samples for $(J, \nu)$, we automatically also have samples for the distribution for $J$, that for $F$, and that for $\nu$.

Let's consider the sum over all possible sequences, $\sum_{\overline{A}}$, in the joint density (46). Owing to the delta in the summand, the only $\overline{A}$s that contribute to the sum are those for which $J(\bar{a}, \overline{A}) = J$. This also means that all terms in the sum have the same numerical value, because their values depend on $\overline{A}$ only through $J(\bar{a}, \overline{A})$, which is fixed. We must therefore only find how many terms there are in the sum, and multiply the value of one term for the number of terms.

With $\bar{a}$ and $J$ fixed, we are asking how many $\overline{A}$s satisfy $J(\bar{a}, \overline{A}) = J$. Consider the problem from the following point of view. We have a grid of boxes with $(n + 1)$ rows and $(N + 1)$ columns, indexed by $(a, A)$. A sequence of $T$ balls go into the boxes: if the $t$th ball goes into the $(a, A)$ box, it means that at the $t$th time bin the activities of sample and full network were $a_t = a$ and $A_t = A$. In the typical combinatorial problem we would ask in how many ways we can fill the boxes, with $T J_{aA}$ balls in the box $(a, A)$, by throwing the $T$ balls. The number would be given by all $T!$ possible permutations of the balls, but considering permutations

of two or more balls within the same box as equivalent, thus finding the multinomial coefficient

$$\binom{T}{T\boldsymbol{J}} \equiv \binom{T}{TJ_{00}, \dots, TJ_{nN}} \coloneqq \frac{T!}{\prod_{a,A}(TJ_{aA})!}. \tag{47}$$

In our case, however, we have one constraint: $\bar{a}$ is fixed, which means that the *row* in which each ball must fall is determined and fixed. The $t$th ball must perforce fall into row $a_t$. In considering all possible permutations we must therefore exclude those that change the rows of the balls. This means that only within-row permutations are allowed. Within each row the counting proceeds as usual, so that for row $a$, which has a total of $T \sum_A J_{aA} \equiv T f_a$ balls, we have

$$\binom{T f_a}{TJ_{a0}, \dots, TJ_{aN}} = \frac{(T f_a)!}{\prod_A (TJ_{aA})!}. \tag{48}$$

We therefore find that the number of terms in the sum of (46) is

$$\prod_a \binom{T f_a}{TJ_{a0}, \dots, TJ_{aN}}. \tag{49}$$

If $T$ is large the logarithm of the multinomial coefficient can be approximated using the Shannon entropy $H$, which for a generic distribution $(x_i)$ satisfies the bounds[44]

$$TH(x_0, \dots, x_N) - \ln\binom{T+N}{T} \leqslant \ln\binom{T}{Tx_0, \dots, Tx_N} \leqslant TH(x_0, \dots, x_N). \tag{50}$$

We therefore approximate the multiplicity (49) as

$$\exp\left[T \sum_a f_a\, H\left(\frac{J_{a0}}{f_a}, \dots, \frac{J_{aN}}{f_a}\right)\right] \equiv \exp\left[-T \sum_a f_a \sum_A \frac{J_{aA}}{f_a} \ln \frac{J_{aA}}{f_a}\right] \equiv$$

$$\exp\left[-T \sum_{a,A} J_{aA} \ln J_{aA} + T \sum_a f_a \ln f_a\right] \equiv \exp[T\,H(\boldsymbol{J}) - T\,H(\boldsymbol{f})], \tag{51}$$

where the second line is obtained by combining the sums, simplifying, and using the definition of Shannon entropy. ✚ Is it possible that this approximation affects the result? Check this!

[44]Lemma 2.2 pp. 429–430 in I. Csiszár et al. (2004): *Information theory and statistics: a tutorial.* Foundations and Trends in Communications and Information Theory **1**[4], 417–528. http://www.renyi.hu/~csiszar/.

The sum over $\bar{A}$ in formula (46) can therefore be replaced by any single summand multiplied by the multiplicity (51) above. The marginal frequency distribution $f$ is still fixed, determined by $\bar{a}$, so we still have the constraint $\sum_A J_{aA} = f_a$ for each $a$, which we compactly write $\delta(\sum J = f)$. We find

$$p(J, \nu \mid \bar{a}, I)$$

$$= \frac{1}{Z} \, \delta(\textstyle\sum J = f) \, q(\nu) \, \exp[T \, H(J) - T \, H(f)] \prod_{a,A} (G_{aA} \, \nu_A)^{T J_{aA}(\bar{a},\bar{A})}$$

$$= \frac{1}{Z} \, \delta(\textstyle\sum J = f) \, q(\nu) \, \exp\{T[H(J) - H(f) + \textstyle\sum_{a,A} J_{aA} \ln(G_{aA} \, \nu_A)]\},$$

$$(52)$$

where the second equality comes from re-expressing the product over $(a, A)$ in exponential-logarithm form.

Now we note that $G_{aA} \, \nu_A$ is a normalized distribution in $(a, A)$ owing to the properties of the hypergeometric distribution and of $\nu$. We denote it $G \cdot \nu$. We can also use the definition of relative entropy

$$H[(x_i); (y_i)] \coloneqq \textstyle\sum_i x_i \ln \frac{x_i}{y_i} \equiv -H[(x_i)] - \textstyle\sum_i x_i \ln y_i, \qquad (53)$$

to combine the first and last terms within the exponential. The logarithm of our joint density (46) can then finally be written as

$$\boxed{\begin{aligned} \ln p(J, \nu \mid \bar{a}, I) = \ln \delta(\textstyle\sum J = f) + \ln q(\nu) - T \, H(J; G \cdot \nu) \\ - T \, H(f) - \ln Z(\bar{a}) \end{aligned}}$$

$$(54)$$

The last two terms are constants: they only depend on $\bar{a}$ and $f(\bar{a})$, which are given and fixed in our problem. They can therefore be discarded in Monte Carlo sampling. In sampling, the delta term is taken care of by restricting the sampling to the set of allowed $J$.

The density $q(\nu)$ remains to be specified. We use an entropic density

$$q(\nu) \propto \exp[L \, H(\nu; r)] \qquad (55)$$

where $r$ is a reference distribution – we take the uniform one for simplicity – and $L$ is of order unity, say less than 10. This density is motivated my

the consideration of multiplicities: a frequency distribution, like $\boldsymbol{v}$, can be realized in a number of ways given by a multinomial coefficient. This coefficient, by formula (50), approximately equal to the exponential of the Shannon entropy of the distribution. So the density (55) keeps track of this multiplicity, but regulates its importance via the parameter $L$. We can also choose $L = 0$, which leads to $q(\boldsymbol{v})$ being uniform in $\boldsymbol{v}$.

## A.3   Approximations

We consider two successive approximations of the log-density (54), based on the fact that $T$ is large.

The relative entropy $H(J; G \cdot \boldsymbol{v})$ is always positive, and vanishes if only if $J = G \cdot \boldsymbol{v}$. This term is multiplied by $-T$ and appears in an exponential. We therefore approximate it with a delta.

This delta restricts $J$ to the form $G \cdot \boldsymbol{v}$, making it completely determined by $\boldsymbol{v}$. The density (54) thus collapses onto the $\boldsymbol{v}$-space. More precisely a subspace, since $J$ must also satisfy the constraint $\sum_a J = f$, which becomes $\sum_a G \cdot \boldsymbol{v} = f$, constraining $\boldsymbol{v}$.

Such a constraint is actually too strong, in the sense that no distribution $\boldsymbol{v}$ can satisfy it – it leads to negative values $nu_a$ for some $a$s. This problem is mitigated by considering a softer constraint between $J$ and $f$, for example taking into account errors in spike sorting or just considering the first $m$ moments of $f$. We express the softer constraint by replacing the delta with a normal N whose covariance matrix $\boldsymbol{\sigma}$ is determined from the spike-sorting error. With this approximation we reduce our problem to the density

$$\ln \mathrm{p}(\boldsymbol{v} \mid \bar{a}, I) \propto \ln \mathrm{N}(\textstyle\sum_a G \cdot \boldsymbol{v} \mid f, \boldsymbol{\sigma}) + \ln q(\boldsymbol{v}), \qquad (56)$$

with $q(\boldsymbol{v})$ given by the entropic density (55). We can call this the 'maximum-entropy approximation with Bayesian correction'. It can also be represented by a set of Monte Carlo samples $\{\boldsymbol{v}^{(i)}\}$.

As a further approximation we can simply consider the mode of (56). This is found by maximizing the relative entropy (55) under the constraints about $f$. The solution is the one given by the maximum-entropy method applied to the full network with subnetwork constraints.

Figure A shows comparison of the expected values (means of the Monte Carlo samples) for the density (54) (for $F$), the density (56) (for $\boldsymbol{v}$

and using the first six moments of $f$ as constraints), and the maximum-entropy solution (six moments as constraints).
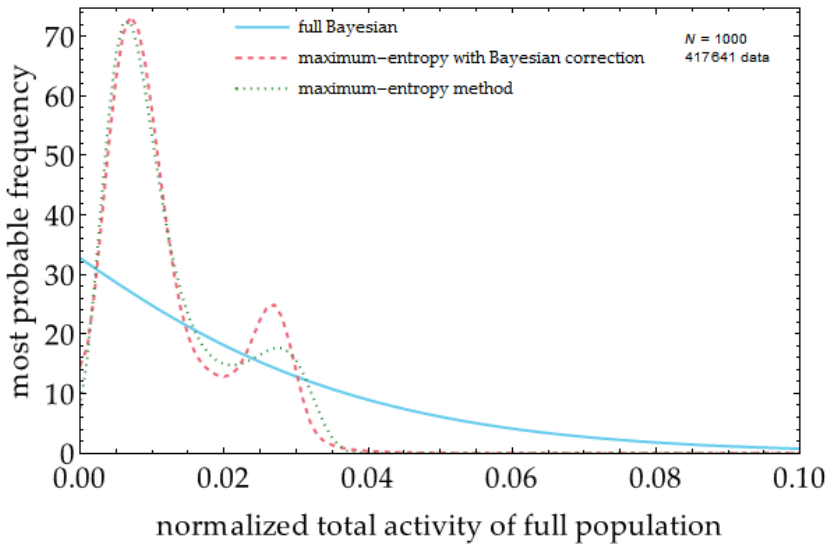


Figure 8   Comparison of (54), (56), and the maximum-entropy approximation (the vertical-axis label is wrong).