

# Notes on inferring connectivity

## (Bente's problem) [draft]

P.G.L. Porta Mana <[piero.mana@ntnu.no](mailto:piero.mana@ntnu.no)>

17 October 2019; updated 3 November 2019

Notes on Bente's problem about inferring connectivity. This is just a memo, not meant to be detailed or precise.

### 1 Synopsis of the problem

We're considering neurons in a specific region  $A$  of the mouse brain, and neurons in a specific region  $B$  which connects to those in region  $A$ . Our question is how many neurons in region  $B$  project, on average, to a single neuron in region  $A$ .

First of all let's make our question more precise. We can rephrase it in two different ways; it's the second that probably interests us:

- Q1.** If we examine a new neuron in region  $A$  in a new mouse, how many neurons in region  $B$  will project to it? The possible answers to this question are: 0, 1, 2, and so on.
- Q2.** If we examine all neurons in region  $A$  in a very large sample of mice, and for each neuron we count how many neurons from region  $B$  project to it, then what are the frequencies with which we'll observe 0 connections, 1 connection, 2 connections, and so on? One possible answer could be this distribution:

0: 0.1%,    1: 0.3%,    ...,    1 000: 24%,    ...

another answer could be this distribution:

0: 0.02%,    1: 0.07%,    ...,    1 000: 31%,    ...

and so on, for all possible distributions of frequencies.

For the first question we must give a distribution of probability over all possible answers 0, 1, and so on. For the second the probability is distributed over all possible frequency distributions.

The answers to these two questions and their probabilities are connected; in particular, from the second we can obtain the first. The second is more informative, because it gives us also a range of variability.

Our analysis will also produce probabilities for other another quantity: how many neurons are there in region  $B$ , in the next mouse we observe?

Our question actually concerns several input regions at once. For the moment we focus on setting up the calculation for one input region only, but the calculation will be generalized to several input regions later. Considering several at once will actually improve the results, even if just a little.

**The observations and data we use:** in several mice a small number of neurons in region  $A$  – *starter neurons* – are infected with a virus. This virus spreads, with time, to neurons – *input neurons* – that projects to the starter neurons. The virus can't spread further than that. Some of the infected neurons will be in region  $B$ . The virus also leads to colouring of starter and input neurons. Each mouse is killed after several days (different for each mouse) and the starter and input neurons are counted. So our data are the number of starter and input neurons for several mice.

There's a complication: if the mouse is killed too early, the virus hasn't had the time to spread to all input neurons. If the mouse is killed too late, the starter cells decay and are no longer visible. So we also need to guess, for each observed mouse, whether there were additional starter neurons, no longer visible, and what proportion of input neurons have been reached by the virus. A data, we have the number of days between virus injection and killing, for each mouse. This complication for the moment is set aside, to simplify the initial calculation, but will be taken into account later.

## 2 Strategy for solving the problem

We must consider the set of mice that have been observed and those that will be observed in the future. Let's label each mouse with an index  $m$ . For our inferences to be valid, these mice must be judged 'similar'.

For each mouse consider all neurons in the input region, label them with  $i = 1, 2, \dots$ ; and all neurons in the target region, label them  $j = 1, 2, \dots$ . Mouse  $m$  has  $K_m$  neurons in the input region  $B$  and  $L_m$  neurons in the target region  $A$ .

Now that we've labelled the mice and their neurons in the two regions, we can make several concrete statements, such as

*In mouse #64, the neuron #931 in the input region projects to the neuron #1488 in the target region.*

which can be either true or false. Let's denote such kind of statement with

$$C_{mij} := \text{In mouse } m, \text{ neuron } i \text{ in input region projects to neuron } j \text{ in target region} \quad (1)$$

and let's introduce the quantity

$$c_{mij} = \begin{cases} 1 & \text{if } C_{mij} \text{ is true,} \\ 0 & \text{if } C_{mij} \text{ is false.} \end{cases} \quad (2)$$

Obviously for most of these statements we'll never be able to ascertain whether they're true or false. But the answers to our question Q2 and the statement of our experimental observations are built up from the specific statements (1). For example, the statement *In mouse # 3, neuron # 1 in the target region receives input from 2 neurons in the input region* is equivalent to the composite statement *In mouse # 3, neuron # 1 in the target region receives input from neurons # 1 and # 2, or # 1 and # 3, or # 1 and # 4, or ... in the input region*, which we can express as

$$(C_{311} \text{ and } C_{312}) \text{ or } (C_{311} \text{ and } C_{313}) \text{ or } \dots \\ \text{or } (C_{312} \text{ and } C_{313}) \text{ or } (C_{312} \text{ and } C_{314}) \text{ or } \dots \quad (3)$$

The probabilities of such combinations are completely determined once we give an initial joint probability distribution for all statements (1) and their negations, which is equivalent to the joint distribution of

$$p[(c_{mij}) | I] \equiv p(c_{111}, c_{112}, c_{113}, \dots | I) \quad (4)$$

conditional on the information  $I$  that we have before doing the experiments and the observations.

Now, we'll want the probability of a specific answer  $A$  to our question, conditional on the data  $D$  we observed:  $p(A | D, I)$ . By the rule of conditional probability we obtain it as

$$p(A | D \text{ and } I) = \frac{p(A \text{ and } D | I)}{p(D | I)}. \quad (5)$$

The probabilities in the numerator and denominator can be obtained from our main joint probability (4) using the basic rules of the probability calculus.

So our strategy is as follows:

1. Quantitatively determine the initial joint probability  $p[(c_{mij}) | I]$ .
2. Resolve the answers to question Q2 in terms of the basic statements  $C_{mij}$ .
3. Resolve the experimental observations in terms of the basic statements  $C_{mij}$ .
4. Calculate the probabilities in (5).