

# Do genes keep us awake?

## Research notes

Cüneyt

<cuneyt.guzey@ntnu.no>

Daniela

<daniela.bragantini@ntnu.no>

Luca

<piero.mana@ntnu.no>

Yasser

<yasser.roudi@ntnu.no>

Draft of 23 September 2018 (first drafted 22 August 2018)

Research notes

## 1 Introductory notes

### 1.1 Preliminary remarks about Bayesian probability theory

Bayesian probability theory is not just a set of new, better recipes meant to replace old ones. It also requires a different – and simpler – mindset about problems of inference. Three points are especially important:

1. The only purpose of Bayesian theory is to give the probability of some statements – more exactly, ‘propositions’ (Copi et al. 2014; Barwise et al. 2003) – given other statements that may concern data, facts, hypotheses. For example, Bayesian theory can tell us that hypothesis  $A$  has probability  $x$  given some data  $C$  and initial information  $I$ , while hypothesis  $B$  has probability  $y$  given the same conditions:

$$P(A|C, I) = x, \quad P(B|C, I) = y.$$

That’s all there is to it. We can then use these probabilities as we like; in particular, we can use them within decision theory to choose courses of action (Raiffa et al. 2000; Pratt et al. 1996; Sox et al. 2013). But notions like ‘statistical significance’, ‘acceptance level’, ‘confidence’, and similar are foreign to Bayesian theory; or at best they’re just secondary notions.

2. Bayesian theory is an extension of formal logic, the truth calculus. In fact we’ll call it *probability calculus* from now on.

In formal logic, to prove a theorem we need some axioms to start from. These may partly include experimental facts or data, but they always also include assumptions that are purely conjectural. It’s impossible to avoid

this conjectural element (see for example Harding 1976).<sup>1</sup> Likewise, in the probability calculus we need to specify initial probabilities. These may originate in data, but they always also include additional assumptions. The motto ‘let the data speak for themselves’ is simply impossible.

The difference between Bayesian methods and traditional methods is *not* that the former need additional assumptions while the latter don’t. Rather, Bayesian methods make these assumptions explicit, while traditional methods hide them. This is the reason why many traditional results can be obtained as special cases of Bayesian ones.

3. Conditional probabilities like  $P(A|B)$  do not express a causal connection between  $A$  and  $B$ , but an *informational* connection. In that conditional probability,  $A$  could be the cause of  $B$ , or  $B$  of  $A$ , or neither could be the cause of the other. The classical example of this is

$$P(\text{clouds in the sky} | \text{rain on the pavement}, I) > 0.5, \quad (1)$$

not because the rain is the cause of the clouds, but because its presence gives us *relevant information* about the cloudiness of the sky.

The previous remarks may appear pedantic, but they’re important lest we misuse Bayesian methods.

## 1.2 What is the question?

✂ Luca: the following thoughts may be naive; I must still read (Stingo et al. 2015) and (Bush et al. 2012)

---

<sup>1</sup>This impossibility is well known in modern science; we can quote Poincaré (1992): ‘But upon more reflection we realize the position held by hypothesis; we see that the mathematician wouldn’t know how to do without it, and the experimenter can’t do without it at all’ (Introduction); ‘Every generalization is a hypothesis’ (ch. IX, p. 176). Duhem (1991): ‘In sum, the physicist can never subject an isolated hypothesis to experimental test, but only a whole group of hypotheses; when the experiment is in disagreement with his predictions, what he learns is that at least one of the hypotheses constituting this group is unacceptable and ought to be modified; but the experiment does not designate which one should be changed’ (§ VI.2, p. 187); ‘Unlike the reduction to absurdity employed by geometers, experimental contradiction does not have the power to transform a physical hypothesis into an indisputable truth; in order to confer this power on it, it would be necessary to enumerate completely the various hypotheses which may cover a determinate group of phenomena; but the physicist is never sure he has exhausted all the imaginable assumptions’ (§ VI.3, p. 190); ‘the realization and interpretation of no matter what experiment in physics imply adherence to a whole set of theoretical propositions’ (§ VI.5, p. 200). Medawar (1963): ‘the starting point of induction, naive observation, innocent observation, is a mere philosophic fiction. There is no such thing as unprejudiced observation’ [Jeffreys [quote]][ref].

We want to assess the informational relevance between some genetic variations  $\{G\}$  and (combinations of) insomnia symptoms  $\{S\}$  in the Norwegian or European population. To assess this relevance we use data  $D$  from a population sample. Some assumptions or background information  $I$  are also always present in our assessment.

The traditional approach to this kind of problems is to set two hypotheses against each other: ‘there is a correlation’ vs ‘there isn’t a correlation’, and to assess which is more ‘significant’ in view of the data. Mathematically this corresponds to a dichotomy between an exactly zero value and non-zero values of a correlation-like quantity.

Here we approach this problem differently. Rather than contrasting zero against non-zero values, we simply calculate how *relevant* one quantity is to make inferences about the other quantity. ✚check (Stephens et al. 2009): they seem to have a reference with a similar philosophy One quantity, for example  $G$ , is *inferentially* or *informationally* relevant when the knowledge of its value leads us to have a different degree of belief regarding the value of the other, for example  $S$ , compared to when we don’t know its value. In other words, the probabilities  $P(S|G DI)$  and  $P(S|DI)$  are numerically different. If these two probabilities are approximately equal then the particular genetic variation  $G$  are *irrelevant* for our prediction of the insomnia symptom  $S$ . The same conclusion holds with  $G$  and  $S$  exchanged: the probability calculus says that

$$P(S|G DI) = P(S|DI) \iff P(G|S DI) = P(G|DI) \quad (2)$$

if  $P(S|DI)$ ,  $P(G|DI)$  aren’t zero.

This measure of relevance can be extended to sets of (combinations of) symptoms  $\{S\}$  and of genetic variations  $\{G\}$  by using the conditional entropy (Shannon 1948; Kelly 1956; Press et al. 2007 § 14.7; Cover et al. 2006 ch. 2)

$$H(\{S\}|\{G\}, DI) := - \sum_G P(G|DI) \sum_S P(S|G DI) \ln P(S|G DI), \quad (3)$$

which is zero only if  $G$  gives us certainty about  $S$ , and is equal to the entropy

$$H(\{S\}|DI) := - \sum_S P(S|DI) \ln P(S|DI) \quad (4)$$

if  $G$  is irrelevant for predicting  $S$  (Press et al. 2007 § 14.7; Cover et al. 2006 ch. 2). Another, symmetric measure of relevance is the mutual information, discussed in § 2.

If we find that there is mutual informational relevance between genetic variations and insomnia symptoms, we can conclude from biologic reasons that those variations must have a direct or indirect influence on the symptoms, for example they may give susceptibility to insomnia.

There are three main ways to calculate the conditional probabilities of above: By calculating first  $P(SG|DI)$ , or  $P(S|GDI)$ , or  $P(G|SDI)$ . Also, we can consider all possible combinations of genetic variations from the outset, or consider combination of few variations, gradually increasing the numbers. We shall try all these approaches and see whether their results are mutually consistent.

### 1.3 Why exchangeability

The sentence ‘the probability for an insomnia symptom given a genetic variation’ is vague. What we mean is our guess that a given individual with that variation presents that symptom. Our guess about a particular individuals in the full population is updated from our knowledge about individuals we have sampled. Our guess about a new individual is affected by the sample data only insofar we believe those data to be representative for that individual.

The notion of exchangeability expresses this representativeness in probabilistic terms, as explained at length by de Finetti (1931; 1937; 1938). Denote by  $S_s^{(i)}$  the statement that individual  $i$  has symptom  $S_s$ , and likewise with  $G_g^{(i)}$  for the combination of genetic variations  $g$ . The fact that we believe, for inferential purposes, that the individuals from a population having the same genetic variation  $g$  are representative of one another is expressed by

$$P(S_{s_1}^{(1)} S_{s_2}^{(2)} S_{s_3}^{(3)} \dots | G_g^{(1)} G_g^{(2)} G_g^{(3)} \dots I) = P(S_{s_{\pi(1)}}^{(1)} S_{s_{\pi(2)}}^{(2)} S_{s_{\pi(3)}}^{(3)} \dots | G_g^{(1)} G_g^{(2)} G_g^{(3)} \dots I) \quad (5)$$

where  $\pi$  is an arbitrary permutation of the individuals’ labels  $i$ . For example,

$$P(S_c^{(1)} S_a^{(2)} S_b^{(3)} \dots | G_g^{(1)} G_g^{(2)} G_g^{(3)} \dots I) = P(S_b^{(1)} S_c^{(2)} S_a^{(3)} \dots | G_g^{(1)} G_g^{(2)} G_g^{(3)} \dots I). \quad (6)$$

This mathematical property is called *exchangeability*. If the number of individuals is finite it is called *finite* exchangeability; letting this number increase indefinitely we reach *infinite* exchangeability as a limit (Heath et al. 1976).

Assuming exchangeability for each group of individuals sharing the same genetic variation, the mathematical relations above generalize as follows. Our probability assignment is the same if we exchange symptom labels among individuals *having the same genetic variation*; but it may be different if we exchange symptom labels among individuals with different genetic variations. The exact mathematical expression may look complicated; a concrete example is

$$P(S_a^{(1)} S_c^{(2)} S_d^{(3)} S_a^{(4)} S_b^{(5)} S_d^{(6)} \dots | G_\alpha^{(1)} G_\beta^{(2)} G_\beta^{(3)} G_\gamma^{(4)} G_\alpha^{(5)} G_\gamma^{(6)} \dots I) = \\ P(S_b^{(1)} S_d^{(2)} S_c^{(3)} S_d^{(4)} S_a^{(5)} S_a^{(6)} \dots | G_\alpha^{(1)} G_\beta^{(2)} G_\beta^{(3)} G_\gamma^{(4)} G_\alpha^{(5)} G_\gamma^{(6)} \dots I) \quad (7)$$

where the probability remains the same as we exchange symptoms a and b between individuals 1 and 5, both having genetic variation  $\alpha$ ; symptoms c and d between individuals 2 and 3, both having genetic variation  $\beta$ ; symptoms a and d between individuals 4 and 6, both having genetic variation  $\gamma$ . These exchanges can involve an arbitrary number of individuals, symptoms, genetic variations. This general property is called *partial* exchangeability (de Finetti 1938; Diaconis et al. 1980; Diaconis 1988; for a connection with sampling theory see Sugden 1982; 1993).

Note that the property exemplified by eq. (7) is more general than just separately stating exchangeability for the probability distributions for the individuals sharing the same genetic variations. Property (7) allows data about individuals with a genetic variation to be *relevant* for prediction of data about individuals with *another* genetic variation, as we'll see shortly.

✚ add: de Finetti's representation theorem for probability distributions with the property above

## 1.4 Selection of variables and robustness

Denote the presence of the genetic variation labelled  $i$  by  $G_i$  and its absence by  $\neg G_i$ . We can consider the relevance of each variation individually, say

$$P(S | G_1 DI), \quad (8)$$

or of the combination of any number of variations, say

$$P(S | G_1 \neg G_2 \neg G_3 G_4 DI). \quad (9)$$

The probability calculus allows us to assign all these probabilities for any amount of data  $D$  – since they represent beliefs. If the number of combinations is high compared with the number of data, however, these probabilities will usually change noticeably when updated with new data; we can say that they are less ‘robust’ to the acquisition of new data. This robustness can be quantified in various ways to be discussed later.

From this point of view it makes sense to first consider each genetic variation individually and then larger and larger combinations of variations, as long as we see that our probabilities conditional on data  $D$  are robust.

✠ Jeffreys (1983) § 3.2 *very* relevant to our problem! Also Broad (1918)

✠ See Jeffreys (1983 § 3.1, p. 124) on the probability to be given to the ratio values  $\{0, 1\}$ : ‘In genetics the suggested values are usually intermediate, such as  $1/2$ ,  $1/4$ , and  $3/8$ ’. Also, ‘we cannot give a universal rule for them beyond the common-sense one, that if anybody does not know what his suggested value is, or whether there is one, he does not know what question he is asking and consequently does not know what his answer means.’

## 2 First approach: joint probability and mutual information

### 2.1 Notation

The following notation produces compact but readily understandable formulae. Functions and operations on tuples  $x := (x_1, \dots, x_C)$ ,  $y := (y_1, \dots, y_C)$ , and numbers  $a$  operate component-wise. For example:

$$\begin{aligned} \exp x &:= (\exp x_1, \dots, \exp x_C) & xy &:= (x_1 y_1, \dots, x_C y_C) \\ ax &:= (ax_1, \dots, ax_C) & x^a &:= (x_1^a, \dots, x_C^a) \quad \text{and so on.} \end{aligned} \quad (10a)$$

The exception are the sum and multiplication operators  $\sum, \prod$ :

$$\sum x := x_1 + \dots + x_C \quad \prod x := x_1 \cdots x_C \quad (10b)$$

so that, for example,

$$\sum \ln(x/y) := \sum_{i=1}^C \ln(x_i/y_i).$$

Note also the conventions

$$\binom{a}{x} := \binom{a}{x_1, \dots, x_C} := \frac{a!}{x_1! \cdots x_C!} \quad \prod \binom{y}{x} := \binom{y_1}{x_1} \cdots \binom{y_C}{x_C}, \quad (10c)$$

where the first expression is the multinomial coefficient.

## 2.2 Scheme of this approach

Here is the way of thinking and general form of the calculations for this approach. We'll see later if these calculations are practically feasible.

Denote by  $N$  the amount of Norwegian or European population, roughly equal to  $5.3 \times 10^6$  or  $740 \times 10^6$ . Our initial information  $I$  says that each individual is characterized by an insomnia parameter  $\sigma \in \{1, \dots, 8\}$  with  $C_\sigma := 8$  possible values, and by a combination of 94 gene allele pairs, denoted by  $\gamma := (g_1, \dots, g_{94}) \in \{0, 1\}^{94}$ , with  $C_\gamma := 2^{94}$  possible values. The combined variate  $\xi := (\sigma, \gamma)$  can thus assume  $C := C_\sigma \times C_\gamma \approx 1.6 \times 10^{29}$  possible values.

We have data  $D$  consisting in the values of the variate for a sample of  $n := 6029$  individuals:  $\xi^{(i)} := (\sigma^{(i)}, \gamma^{(i)}), i \in \{1, \dots, n\}$ .

Denote the relative frequency for  $\sigma$  in the full population by  $S$ , that for  $\gamma$  by  $G$ , and their joint frequency distribution by  $X$ . Note that  $\sum X = 1$  and so on for the other frequency distributions, and that  $S_\sigma = \sum_\gamma X_{\sigma, \gamma}$  and  $G_\gamma = \sum_\sigma X_{\sigma, \gamma}$ . Denote the corresponding relative-frequency distributions in the sampled population by  $s$ ,  $g$ , and  $x$ ; analogous equalities hold for them.

The  $C_\gamma$  possible gene variations are much more numerous than the data  $n$ ; not all of them can therefore appear in the data, that is, have frequencies  $g_\gamma^{(0)} > 0$ . Denote by  $C_\gamma^+$  the number of distinct variations present in the data, and by  $C_\gamma^-$  the number of those not present. Let  $C^+$  and  $C^-$  have the same meaning regarding the joint variate  $(\sigma^{(0)}, \gamma^{(0)})$ . Obviously  $C_\gamma = C_\gamma^+ + C_\gamma^-$ ,  $C = C^+ + C^-$ , and  $C^+ \leq C_\gamma^+$ .

We assess the mutual relevance of insomnia symptoms and genetic data by considering their joint probability for an individual '0' chosen at random from the full population, given our data:

$$p(\sigma^{(0)}, \gamma^{(0)} | D, I), \quad (11)$$

and the marginal probabilities

$$p(\gamma^{(0)} | D, I), \quad p(\sigma^{(0)} | D, I). \quad (12)$$

From these we compute the mutual information (Shannon 1948; Kelly 1956 in these called ‘rate of transmission’; Press et al. 2007 § 14.7; Cover et al. 2006 ch. 2):

$$I(\sigma^{(0)} : \gamma^{(0)} | D, I) := \sum_{\sigma^{(0)}, \gamma^{(0)}} p(\sigma^{(0)}, \gamma^{(0)} | D, I) \ln \frac{p(\sigma^{(0)}, \gamma^{(0)} | D, I)}{p(\sigma^{(0)} | D, I) p(\gamma^{(0)} | D, I)}. \quad (13)$$

The mutual information vanishes if the symptoms and genetic variations are completely irrelevant to one another, and is equal to the least of the entropies of either marginal probability distribution if they are deterministically related. Thus, instead of asking whether they are ‘uncorrelated’ or ‘correlated’, we’re quantifying their mutual relevance.

In the rest of this section we’ll denote  $\sigma^{(0)}, \gamma^{(0)}$  simply by  $\sigma, \gamma$  for brevity.

Our first step is to calculate the joint probability distribution (11). It can be calculated if we know the joint relative-frequency distribution  $X$  of the variates in the full population. The joint probability (11) *when  $X$  is known* is given as a special case of the multivariate hypergeometric distribution:

$$p(\sigma, \gamma | x, n, X, N, I) \equiv p(\sigma, \gamma | X, N, I) = X_{\sigma, \gamma}. \quad (14)$$

If the frequency distribution  $X$  isn’t known, from the theorem of total probability we obtain

$$p(\sigma, \gamma | x, n, N, I) = \sum_X p(\sigma, \gamma | x, n, X, N, I) p(X | x, n, N, I) \equiv \sum_X X_{\sigma, \gamma} p(X | x, n, N, I). \quad (15)$$

We must therefore find the probability for  $X$  given the data. It can be found using Bayes’s theorem:

$$p(X | x, n, N, I) = \frac{p(x | n, X, N, I) p(X | N, I)}{\sum_X p(x | n, X, N, I) p(X | N, I)}. \quad (16)$$

Bayes’s theorem requires two probabilities: the probability for the frequency distribution  $x$  of the sampled population, given  $X$ ; and our initial



probability for  $\mathbf{X}$ . The first is given by the multivariate hypergeometric distribution ✂ add refs:

$$p(\mathbf{x}|n, \mathbf{X}, N, I) = \binom{N}{n}^{-1} \prod \binom{N\mathbf{X}}{n\mathbf{x}}. \quad (17)$$

The second will be discussed later.

Combining eqs (15), (16), (17), and simplifying we obtain

$$p(\sigma, \gamma|\mathbf{x}, n, N, I) = \frac{\sum_{\mathbf{X}} X_{\sigma, \gamma} \prod \binom{N\mathbf{X}}{n\mathbf{x}} p(\mathbf{X}|N, I)}{\sum_{\mathbf{X}} \prod \binom{N\mathbf{X}}{n\mathbf{x}} p(\mathbf{X}|N, I)}. \quad (18)$$

With an analogous reasoning we find analogous formulae for  $p(\sigma|\mathbf{x}, n, N, I)$  by replacing  $\mathbf{X}$ ,  $\mathbf{x}$  with  $\mathbf{S}$ ,  $\mathbf{s}$ ; and for  $p(\gamma|\mathbf{x}, n, N, I)$  by replacing  $\mathbf{X}$ ,  $\mathbf{x}$  with  $\mathbf{G}$ ,  $\mathbf{g}$ . From these we can calculate the mutual information (13).

### 2.3 First calculation: constant initial probability

We need to assess our initial probability for  $\mathbf{X}$  according to some background knowledge. As a first tentative let's consider background knowledge  $I_{C_\sigma}$  such that the probability is constant for all possible distributions  $\mathbf{X}$ . There are  $\binom{N+C-1}{C-1}$  possible distributions  $\mathbf{X}$  (Csiszár et al. 2004 § 2.1); hence

$$p(\mathbf{X}|N, I_{C_\sigma}) = \binom{N+C-1}{C-1}^{-1}. \quad (19)$$

This yields a constant initial probability distribution for  $(\sigma, \gamma)$ :

$$\begin{aligned} p(\sigma, \gamma|N, I_{C_\sigma}) &= \sum_{\mathbf{X}} p(\sigma, \gamma|\mathbf{X}, N, I_{C_\sigma}) p(\mathbf{X}|N, I_{C_\sigma}) = \\ &= \sum_{\mathbf{X}} X_{\sigma, \gamma} \binom{N+C-1}{C-1}^{-1} = 1/C, \end{aligned} \quad (20)$$

because the last sum is a convex combination, with equal weights, of all points in a simplex of dimension  $C-1$ , giving its centre of mass  $(1/C, 1/C, \dots)$ .

With the initial knowledge  $I_{C_\sigma}$ , eq. (18) simplifies to

$$p(\sigma, \gamma | x, n, N, I_{C_\sigma}) = \frac{\sum_{\mathbf{x}} X_{\sigma, \gamma} \prod \binom{N\mathbf{x}}{n\mathbf{x}}}{\sum_{\mathbf{x}} \prod \binom{N\mathbf{x}}{n\mathbf{x}}}. \quad (21)$$

The sum in the denominator can be calculated with an identity for multinomial coefficients ✂ add refs for summation formula:

$$\sum_{\mathbf{x}} \prod \binom{N\mathbf{x}}{n\mathbf{x}} = \binom{N+C-1}{n+C-1}, \quad (22)$$

which leads also to

$$p(\mathbf{X} | x, n, N, I_{C_\sigma}) = \binom{N+C-1}{n+C-1}^{-1} \prod \binom{N\mathbf{x}}{n\mathbf{x}}. \quad (23)$$

The sum in the numerator can be rewritten ✂ explain the steps obtaining

$$\sum_{\mathbf{x}} X_{\sigma, \gamma} \prod \binom{N\mathbf{x}}{n\mathbf{x}} = \frac{nx_{\sigma, \gamma} + 1}{N} \binom{N+C}{n+C} - \frac{1}{N} \binom{N+C-1}{n+C-1}. \quad (24)$$

Simplifying we finally find

$$p(\sigma, \gamma | x, n, N, I_{C_\sigma}) = \frac{(N+C)nx_{\sigma, \gamma} + N - n}{N(n+C)} \quad (25)$$

This expression can be interpreted as a weighted sum of the observed frequency  $x_{\sigma, \gamma}$  in the data and the initial probability  $1/C$ , eq. (20):

$$p(\sigma, \gamma | x, n, N, I_{C_\sigma}) \propto x_{\sigma, \gamma} + \frac{N-n}{n} \frac{C}{N+C} \frac{1}{C}, \quad (26)$$

the ratio of the second to the first weight being  $(N-n)/n \times C/(N+C)$  ✂ Luca: I don't find this intuitively satisfying, though I don't know why. It'd be good to try another initial state of knowledge. When  $n = N$ , that is, when we've sampled the full population, the second weight is zero and we're left with  $x_{\sigma, \gamma}$ , as intuition demands.

For the marginal distributions we obtain, by summation,

$$p(\sigma | x, n, N, I_{C_\sigma}) = \frac{(N+C)ns_\sigma + (N-n)C_\gamma}{N(n+C)}, \quad (27)$$

$$p(\gamma | x, n, N, I_{C_\sigma}) = \frac{(N+C)ng_\gamma + (N-n)C_\sigma}{N(n+C)}. \quad (28)$$

Note that if the initial probability distribution for  $(\sigma, \gamma)$  is constant, then the corresponding marginals for  $\sigma$  and  $\gamma$  are *not* constant.

Using the distributions (25) and (27) in the formula for the mutual information (13) and simplifying we find

$$I(\sigma : \gamma | D, I_{C_\sigma}) = \ln[N(n+C)] + \sum_{\sigma, \gamma} \frac{(N+C)nx_{\sigma, \gamma} + N-n}{N(n+C)} \times \ln \frac{(N+C)nx_{\sigma, \gamma} + N-n}{[(N+C)ns_\sigma + (N-n)C_\gamma][(N+C)ng_\gamma + (N-n)C_\sigma]} \quad (29)$$

We can divide this sum into two parts: one sum over the values of  $\sigma$  and  $\gamma$  which appear in our data, for which  $g_\gamma > 0$ ; and one sum over the  $C_\gamma^-$  remaining  $\gamma$  values, for which  $g_\gamma = x_{\sigma, \gamma} = 0$ :

$$\begin{aligned} I(\sigma : \gamma | D, I_{C_\sigma}) &= \ln[N(n+C)] \\ &+ \sum_{\sigma} \sum_{\gamma \in D} \frac{(N+C)nx_{\sigma, \gamma} + N-n}{N(n+C)} \times \\ &\quad \ln \frac{(N+C)nx_{\sigma, \gamma} + N-n}{[(N+C)ns_\sigma + (N-n)C_\gamma][(N+C)ng_\gamma + (N-n)C_\sigma]} \\ &- C_\gamma^- \frac{N-n}{N(n+C)} \sum_{\sigma} \ln\{C_\sigma [(N+C)ns_\sigma + (N-n)C_\gamma]\} \end{aligned} \quad (30)$$

## 2.4 Second calculation: uniform marginals

The initial state of knowledge  $I_{C_\sigma}$  used in the previous section was denoted with ' $C_\gamma$ ' because it depends on the number of gene variations we consider. Two different values  $C'_\gamma$  and  $C''_\gamma$  correspond to two different states of knowledge,  $I_{C'_\gamma} \neq I_{C''_\gamma}$ , which yield different initial probabilities for the marginal frequency distribution of a single gene. For example, for very large  $N$ , if we consider a single gene,  $C_\gamma = 2$ , the probability density for its marginal frequency  $f$  is constant in  $df$ . If we consider an additional gene,  $C_\gamma = 4$ , the probability density for the marginal frequency of the previous gene becomes  $6f(1-f)df$ . If we consider an additional third gene, this probability density will be again different. As a consequence, if we decide to consider additional genes we cannot compare our inference with the ones previously made, because they come from contradictory initial states of knowledge.

Another bothering feature of the state of knowledge  $I_{C_\sigma}$  is that the marginal probability for the insomnia symptoms alone given the full set

of data, eq. (27), depends on the number of gene variations considered in the data. Likewise, the probability for the gene variations depend on the number of insomnia symptoms.

It would be advantageous to consider an initial state of knowledge that doesn't lead to different marginal probability distributions for the frequencies when we want to consider an additional gene variation, and that doesn't present the counter-intuitive features above.

One such a state of knowledge  $I_0$  exists (perhaps it isn't unique) and is characterized by a Dirichlet probability distribution for  $X$  in the  $N \rightarrow \infty$  limit:

$$p(X|I_0) dX = \frac{1}{\prod \Gamma(\frac{2}{C})} \prod X^{\frac{2}{C}-1} \delta(\sum X - 1) dX \quad (31)$$

This distribution has the property of leaving the probability distribution for the marginal frequencies of any number of gene variations invariant if we consider one more gene, thanks to the marginalization properties of the Dirichlet distribution (Ferguson 1973 § 2 p. 211). It moreover assigns a uniform distribution to the frequency of each gene variation. It doesn't assign equal probabilities to the possible frequency distributions for the  $C_\sigma \equiv 8$  symptom combinations, but rather

$$p(S|I_0) dS = \frac{1}{\prod \Gamma(\frac{2}{C_\sigma})} \prod S^{\frac{2}{C_\sigma}-1} \delta(\sum S - 1) dS = \frac{1}{8\Gamma(\frac{1}{4})} \prod S^{-3/4} \delta(\sum S - 1) dS. \quad (32)$$

Note, however, that the probability density for the frequency of controls vs cases is uniform.

The approximation  $N \rightarrow \infty$  needs to be justified. If we consider a limited amount of gene variations, so that  $N/C$  is large (say, less than 12 gene variations), this approximation should be valid. If we consider a large amount of gene variations, so that  $C/N$  is large (say, more than 25 gene variations), the possible frequencies lie on the faces of the  $(C - 1)$ -dimensional simplex; but the density (31) is concentrated in regions at the boundary of the simplex, and thus the approximation may still be reasonable.

Result:

$$p(\sigma, \gamma|x, nI_0) = \frac{nx_{\sigma, \gamma} + 2/C}{n + 2} \quad (33)$$

[Luca's memoranda:]

- Use of partial exchangeability *has to* distinguish also between men and women: see Gehrman et al. (2013 p. 327).
- This study could also be used to detect most relevant genes, by eliminating them in turn (and in pairs etc) and checking the ensuing predictions.
- Is it computationally possible to use a 'nonparametric model'? It would avoid unwarranted assumptions and phenomena like over-training.

## Bibliography

('de X' is listed under D, 'van X' under V, and so on, regardless of national conventions.)

- Barwise, J., Etchemendy, J. (2003): *Language, Proof and Logic*. (CSLI, Stanford). Written in collaboration with Gerard Allwein, Dave Barker-Plummer, Albert Liu. First publ. 1999.
- Bernardo, J.-M., DeGroot, M. H., Lindley, D. V., Smith, A. F. M., eds. (1988): *Bayesian Statistics 3*. (Oxford University Press, Oxford).
- Broad, C. D. (1918): *On the relation between induction and probability*. – (Part I.) *Mind* 27<sup>108</sup>, 389–404. See also Broad (1920).
- (1920): *On the relation between induction and probability*. – (Part II.) *Mind* 29<sup>113</sup>, 11–45. See also Broad (1918).
- Bush, W. S., Moore, J. H. (2012): *Genome-wide association studies*. *PLoS Comput. Biol.* 8<sup>12</sup>, e1002822.
- Copi, I. M., Cohen, C., McMahon, K. (2014): *Introduction to Logic*, 14th ed. (Pearson, Harlow, UK). First publ. 1953.
- Cover, T. M., Thomas, J. A. (2006): *Elements of Information Theory*, 2nd ed. (Wiley, Hoboken, USA). First publ. 1991.
- Csiszár, I., Shields, P. C. (2004): *Information theory and statistics: a tutorial*. *Foundations and Trends in Communications and Information Theory* 1<sup>4</sup>, 417–528. <http://www.renyi.hu/~csiszar/>.
- de Finetti, B. (1931): *Probabilismo*. *Logos* 14, 163–219. Transl. as de Finetti (1989). See also Jeffrey (1989).
- (1937): *La prévision : ses lois logiques, ses sources subjectives*. *Ann. Inst. Henri Poincaré* 7<sup>1</sup>, 1–68. Transl. in Kyburg, Smokler (1980), pp. 53–118, by Henry E. Kyburg, Jr.
- (1938): *Sur la condition d'équivalence partielle*. In: *Colloque consacré à la théorie des probabilités. VI : Conceptions diverses*. Ed. by B. de Finetti, V. Glivenko, G. Neymann (Hermann, Paris), 5–18. Transl. in Jeffrey (1980), pp. 193–205, by P. Benacerraf and R. Jeffrey.
- (1989): *Probabilism: A critical essay on the theory of probability and on the value of science*. *Erkenntnis* 31<sup>2-3</sup>, 169–223. Transl. of de Finetti (1931) by Maria Concetta Di Maio, Maria Carla Galavotti, and Richard C. Jeffrey.

- Diaconis, P. (1988): *Recent progress on de Finetti's notions of exchangeability*. In: Bernardo, DeGroot, Lindley, Smith (1988), 111–125. With discussion by D. Blackwell, Simon French, and author's reply. <http://statweb.stanford.edu/~cgates/PERSI/year.html>, <https://statistics.stanford.edu/research/recent-progress-de-finettis-notions-exchangeability>.
- Diaconis, P., Freedman, D. (1980): *De Finetti's generalizations of exchangeability*. In: Jeffrey (1980), 233–249.
- Duhem, P. (1914): *La Théorie Physique : son objet – sa structure*, 2nd ed. (Marcel Rivière, Éditeur, Paris). [http://virtualbooks.terra.com.br/freebook/fran/la\\_theorie\\_physique.htm](http://virtualbooks.terra.com.br/freebook/fran/la_theorie_physique.htm). First publ. 1906. Transl. as Duhem (1991).
- (1991): *The Aim and Structure of Physical Theory*, Transl. of the 2nd ed. (Princeton University Press, Princeton). Transl. of Duhem (1914) by P. P. Wiener. First publ. in French 1906.
- Ferguson, T. S. (1973): *A Bayesian analysis of some nonparametric problems*. *Ann. Stat.* **1**<sup>2</sup>, 209–230.
- Gehrman, P. R., Pfeiffenberger, C., Byrne, E. M. (2013): *The role of genes in the insomnia phenotype*. *Sleep Med. Clin.* **8**<sup>3</sup>, 323–331.
- Harding, S. G., ed. (1976): *Can Theories Be Refuted?* (D. Reidel, Dordrecht).
- Heath, D., Sudderth, W. (1976): *De Finetti's theorem on exchangeable variables*. *American Statistician* **30**<sup>4</sup>, 188–189.
- Jeffrey, R. (1989): *Reading Probabilismo*. *Erkenntnis* **31**<sup>2–3</sup>, 225–237. See de Finetti (1931).
- Jeffrey, R. C., ed. (1980): *Studies in inductive logic and probability*. Vol. II. (University of California Press, Berkeley).
- Jeffreys, H. (1983): *Theory of Probability*, third ed. with corrections. (Oxford University Press, London). First publ. 1939.
- Kelly Jr., J. L. (1956): *A new interpretation of information rate*. *Bell Syst. Tech. J.* **35**<sup>4</sup>, 917–926. <http://turtletrader.com/kelly.pdf>, <https://archive.org/details/bstj35-4-917>.
- Koch, G., Spizzichino, F., eds. (1982): *Exchangeability in Probability and Statistics*. (North-Holland, Amsterdam).
- Kyburg Jr., H. E., Smokler, H. E., eds. (1980): *Studies in Subjective Probability*, 2nd ed. (Robert E. Krieger, Huntington, USA). First publ. 1964.
- Medawar, P. B. (1963): *Is the scientific paper a fraud?* *Listener* **70**, 377–378.
- Poincaré, H. (1905): *Science and Hypothesis*. (Walter Scott, London). Transl. of Poincaré (1992) by W. J. Greenstreet; with a Preface by J. Larmor. First publ. 1902. Partly repr. in Poincaré (1958).
- (1958): *The Value of Science*. (Dover, New York). Authorized transl. with an introduction by G. B. Halsted. First publ. 1913.
- (1992): *La science et l'hypothèse*. (Éditions de la Bohème, Rueil-Malmaison, France). <http://gallica.bnf.fr/document?0=N026745>. First publ. 1902; transl. as Poincaré (1905).
- Pratt, J. W., Raiffa, H., Schlaifer, R. (1996): *Introduction to Statistical Decision Theory*, 2nd pr. (MIT Press, Cambridge, USA). First publ. 1995.
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., Flannery, B. P. (2007): *Numerical Recipes: The Art of Scientific Computing*, 3rd ed. (Cambridge University Press, Cambridge). First publ. 1988.
- Raiffa, H., Schlaifer, R. (2000): *Applied Statistical Decision Theory*, reprint. (Wiley, New York). First publ. 1961.

- Shannon, C. E. (1948): *A mathematical theory of communication*. Bell Syst. Tech. J. **27**<sup>3, 4</sup>, 379–423, 623–656. <https://archive.org/details/bstj27-3-379>, <https://archive.org/details/bstj27-4-623>, <http://math.harvard.edu/~ctm/home/text/others/shannon/entropy/entropy.pdf>.
- Sox, H. C., Higgins, M. C., Owens, D. K. (2013): *Medical Decision Making*, 2nd ed. (Wiley, New York). First publ. 1988.
- Stephens, M., Balding, D. J. (2009): *Bayesian statistical methods for genetic association studies*. Nat. Rev. Genet. **10**, 681–690.
- Stingo, F. C., Swartz, M. D., Vannucci, M. (2015): *A bayesian approach to identify genes and gene-level SNP aggregates in a genetic analysis of cancer data*. Stat. Interface **8**<sup>2</sup>, 137–151.
- Sugden, R. A. (1982): *Exchangeability and survey sampling inference*. In: Koch, Spizzichino (1982), 321–330.
- (1993): *Partial exchangeability and survey sampling inference*. Biometrika **80**<sup>2</sup>, 451–455.