

# Statistical relations between SNPs and insomnia: a Bayesian study [draft]

D. Bragantini

[<daniela.bragantini@ntnu.no>](mailto:daniela.bragantini@ntnu.no)

P.G.L. Porta Mana

[<piero.mana@ntnu.no>](mailto:piero.mana@ntnu.no)

C. Güzey

[<cuneyt.guzey@ntnu.no>](mailto:cuneyt.guzey@ntnu.no)

Y. Roudi

[<yasser.roudi@ntnu.no>](mailto:yasser.roudi@ntnu.no)

Dept of Mental Health, NTNU, Trondheim

Kavli Institute, Trondheim

11 November 2018; updated 25 August 2019

\*\*\*

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

## 1 SNPs and insomnia: introduction and goals

The purpose of the present work is to investigate the statistical link between some insomnia symptoms and a specific set of SNPs from HUNT 3 data. Our approach is characterized by two main features:

1. We quantify the probable link between the two variates in a continuous way, instead of asking the more traditional null-hypothesis question ‘is there an association between this SNP and this symptom – yes or no?’ The complexity of the human organism makes this kind of dichotomous questions very artificial, and their yes/no answers suffer from an arbitrary choice of cut-off. Our results (fig. 5) indeed show that there’s a continuum of statistical links, with no clear ‘significant vs non-significant’ division. Any such cut-off should be made case by case depending on further research questions of interest. For example, one research group might want to further investigate some specific SNPs, but have budget and resources only for two; in this case they’d choose the two most strongly linked SNPs. Another research group might have enough budget and resources for five; they’d choose the five most strongly linked SNPs.



Our perspective reflects old and recent foundational reevaluations in the Probability & Statistics communities<sup>1</sup>.

---

<sup>1</sup> ASA 2016, 2019, Amrhein et al. 2019, see also Kadane in Cox et al. 1987 pp. 347–348, Berger et al. 1988, Johnson 1999, Stephens et al. 2003 Box 3 p. 687.

2. We make full use of the principles of the (Bayesian) probability calculus. Derivations similar to ours have appeared before in the literature of genetic and related studies<sup>2</sup>. These principles allow us in particular to take care of SNPs with very low minor-allele counts: we show (fig. 5) that their probable link strength with respect to other SNPs are qualitatively the same whether we adopt a conservative or a more unrestrained pre-data guess.


We also strive to make our probability calculations intuitively understandable.

Our study reveals a very probable strong link between individual SNPs located in  [recheck location](#) and the three main insomnia symptoms. We also study the link between *pairs* of SNPs and each symptom; this interesting interactions between pairs of SNPs  [figure to be added](#).

Section 3.3 presents our variates, data, and the way we measure the link strength between variates. Section 4 summarizes our methods. Section 5 presents and discusses the results.

## 2 Variates, data, and link measure

### Description of the HUNT data and the locations of the SNPs

We consider three insomnia symptoms: onset insomnia O, maintenance insomnia M, and terminal insomnia T  [explain them](#).

Let's discuss the measure to quantify the link between the alleles  $a$  and  $b$  of a specific SNP on one side, and the presence of a specific insomnia symptom on the other.

Consider a hypothetical, arbitrarily large population from the same genetic pool as our sample. Each individual belongs to one of four mutually exclusive classes: (i) allele  $a$  and symptom, (ii) allele  $b$  and symptom, (iii) allele  $a$ , no symptom, (iv) allele  $b$ , no symptom. Statistical relations between the two variates in the population are fully contained in the *conditional relative frequency* of one variate given the other. Of course these conditional frequencies don't necessarily indicate a causal connection – there can be confounders – but they are our first and most available handle in the investigation of such connection. Since any causal link should go in the direction  $\text{SNP} \rightarrow \text{symptom}$ , we consider the conditional frequencies of the symptom given the two alleles:  $f_{|a}$  and

<sup>2</sup> Lange 1995, 2003, Lewinger et al. 2007, Stingo et al. 2015, see also refs in Stephens et al. 2009.

$f_{|b}$ , which should be more robust to contextual changes<sup>3</sup> (the bar in the notation wants to remind us that these are conditional frequencies, so that  $f_{|a} + f_{|b} \neq 1$  in general).

Any statistical difference in the links between the two alleles and the symptom is then reflected in the difference between the two conditional frequencies:  $\Delta f := f_{|a} - f_{|b}$ . A large positive difference can indicate some biologic association between allele  $a$  and the symptom, and analogously for a negative difference and allele  $b$ . A difference around zero can indicate a weak association or lack thereof. If we are only interested in the link strength, and not in its direction towards one or the other allele, we can focus on the absolute difference  $|\Delta f|$ .


Our approach, as discussed in the Introduction, is to assess the plausible value of the difference  $|\Delta f|$  in our hypothetical population, without choosing any ‘significance’ threshold a priori. Such threshold can always be chosen a posteriori by the readers according to their research goals and resources. The plausibility of the possible differences is then expressed by the probability distribution  $p(|\Delta f|)$ , such as that shown in fig. 1. This distribution is conditional on: (a) the data we’ve collected, (b) our pre-data information or guesses about the frequencies of the hypothetical population.

This distribution contains all probabilistic information we need, but we can also try to summarize it with a single number. We can for example check what is the minimal frequency difference we are 90% sure about, that is, the lowest 10-quantile:

$$x \text{ such that } p(|\Delta f| > x) = 0.9, \quad (1)$$

or simply the expectation

$$E(|\Delta f|). \quad (2)$$

For the plot of fig. 1 such measures are 0.0112 and 0.0262. In this work we use the quantile measure (7), but the results are the same if we use the expectation.  add maybe something about dependence of broadness on sample size, and ‘smoothing’ as discussed by MacKay & Bauman Peto<sup>4</sup>.

A more complete quantification of our plausible guesses is the joint probability distribution  $p(f_{|a}, f_{|b})$  for the two conditional frequencies given the data and our pre-data assumptions, from which  $p(\Delta f)$  can be

<sup>3</sup> Pearl 2009 § 1.3.2  check also pearl2004pearl2003. <sup>4</sup> MacKay et al. 1995 § 2.6.

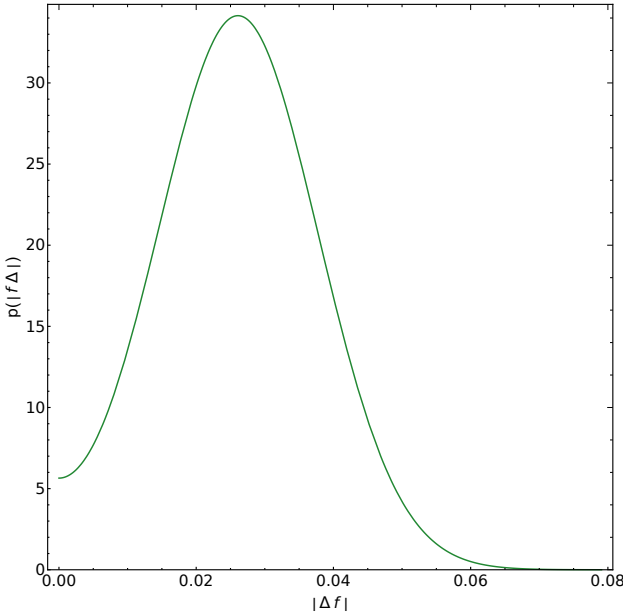


Figure 1 Example of probability distributions for  $|\Delta f|$ , the absolute difference between conditional frequencies. [add lines indicating the 10-quantile and expectation](#)

calculated. This joint distribution is therefore the pivot of our method in the next section.

### 3 Methods

#### 3.1 General plan

To explain our method, consider a specific SNP with alleles  $a$  and  $b$ , and a specific insomnia symptom, for example O. As announced in the previous section, our calculations hinge on the joint probability for the conditional frequencies of the symptom given the two alleles in a hypothetical infinite population,

$$p(f_{|a}, f_{|b} \mid \text{data, initial info}), \quad (3)$$

conditional on the collected data and on our pre-data information or assumptions.

The first step is to use Bayes's theorem:

$$p(f_{|a}, f_{|b} \mid \text{data, pre-data info}) \propto$$

$$p(\text{data} \mid f_{|a}, f_{|b}, \text{pre-data info}) \times p(f_{|a}, f_{|b} \mid \text{pre-data info}), \quad (4)$$

which says that we need to provide a probability and a probability distribution: (i) the probability of observing our sampled data if we had known the conditional frequencies in the larger population; (ii) our pre-data probability distribution for the conditional frequencies. The first is given by a sampling formula. The second can be modelled in several reasonable ways, and we'll see that, owing to the large size of our sample, they all lead to very similar conclusions about the conditional frequencies and their difference.

First of all, introduce some notation. The conditional frequencies are denoted  $f_{|a}, f_{|b}$  as above. The data  $D$  consist in the conditional frequencies observed in the sample:  $F_{|a}$  is the fraction showing the symptom among the sampled individuals having allele  $a$ , and  $F_{|b}$  is the analogous fraction for allele  $b$ . Finally, our initial information  $I$  consists in our initial beliefs and in the numbers  $N_a, N_b$  of sampled individuals having allele  $a$  and  $b$ ; although these numbers are part of  $I$  we will sometimes write them explicitly besides  $I$ . (Note that  $N_a, N_b$  are not part of the 'data' in the sense that they don't help us to update our belief about the frequencies  $f_{|a}, f_{|b}$  and therefore remain on the right side of the conditional.)

With this notation, Bayes's theorem above is written

$$p(f_{|a}, f_{|b} \mid D, I) \propto p(D \mid f_{|a}, f_{|b}, I) p(f_{|a}, f_{|b} \mid I) \equiv$$

$$p(F_{|a}, F_{|b} \mid f_{|a}, f_{|b}, N_a, N_b, I) p(f_{|a}, f_{|b} \mid N_a, N_b, I). \quad (10)_r$$

### 3.2 Plausibility for the data given the frequencies

The plausibility of obtaining the data given the frequencies is given by simple counting and symmetry. We have an arbitrarily large population where a fraction  $f_{|a}$  of individuals having allele  $a$  show the symptom, and a fraction  $1 - f_{|a}$  therefore don't show the symptom. We casually observe  $N_a$  individuals with allele  $a$ . The plausibility that a fraction  $F_{|a}$  of these show the symptom is then counted as a 'drawing with replacement' because of the arbitrarily large size of the full population, given by a

(rescaled) binomial distribution Jaynes 2003 ch. 3, Ross 2010 § 4.6, Feller 1968 § VI.2:

$$p(F_{|a} \mid f_{|a}, N_a, I) = \binom{N_a}{N_a F_{|a}} f_{|a}^{N_a F_{|a}} (1 - f_{|a})^{N_a (1 - F_{|a})}. \quad (5)$$

An analogous reasoning holds for allele  $b$ . Our belief about obtaining the data is therefore

$$p(D \mid f_{|a}, f_{|b}, I) = \prod_{x=a,b} \binom{N_x}{N_x F_{|x}} f_{|x}^{N_x F_{|x}} (1 - f_{|x})^{N_x (1 - F_{|x})}. \quad (6)$$

---

old text below

---

### ✚ Some intro about insomnia and its symptoms here

Every single-nucleotide polymorphism (SNP), together with a huge variety of external factors, can in principle affect the appearance of insomnia symptoms. It is extremely complicated to identify and untangle these causal mechanisms and interactions and to ascertain their degrees, and thus to build a causal graph<sup>5</sup> for them.

An indication of the causal strength of one or more SNPs on one or more insomnia symptoms can be obtained by replacing the causal graph with a corresponding simplified Bayesian network<sup>6</sup> of *conditional probabilities*. In other words, we can study the probability of observing a specific insomnia symptom among the individuals having a particular SNP allele, as a proxy for the possible influence of the latter on the former. This is the approach of genetic association studies ✚ \*\*\*refs here.

In this work we present evidence of a strong association between some SNPs located in ✚ name the genes here and the three main insomnia symptoms ✚ name the relevant SNPs here?, within population sampled in ✚ ref to HUNT study. Our results involve associations between each symptom and several SNPs individually, and also associations between each symptom and *pairs* of SNPs; the latter result shows different kinds of interaction between the alleles of a pair of SNPs.

We study this probabilistic association using Bayesian methods. The basic idea is simple: the conditional frequencies – of each symptom given each allele – in our sample are a rough estimate of the conditional frequencies that we would observe in a hypothetical infinite population. We just need to quantify more precisely our uncertainty about the latter hypothetical frequencies, in the form of a probability distribution for them.

An important point in this calculation is that we may have large samples for specific symptom-allele pairs, and small samples for others. A large samples intuitively gives us a good estimate and lead to a narrower probability distribution. The estimate suggested by a small sample could be deceiving, instead; such estimate is therefore corrected in a conservative way, considering the mean of all other samples. It also leads to a broader probability distribution. This approach is often called *shrinkage* in the probability literature, and has been used in a

---

<sup>5</sup> Pearl 2009. <sup>6</sup> Pearl 2009.

broad variety of association studies, from genetics<sup>7</sup> to contextual text prediction<sup>8</sup> and even baseball<sup>9</sup>. It gives simple, clear, and intuitively understandable inferences about symptom-SNP associations; it is computationally fast; and it is easy to generalize to association studies of symptom combinations vs multiple SNPs, as we'll show in later sections.

The Method section of this work focus on explaining our Bayesian approach, but use the real data from which our results are derived. In the subsequent Result section we discuss the different relevant associations found.


### 3.3 The data

 [description of our data here](#)

## 4 Methods

### 4.1 Outline

We focus on the simplest kind of association: between one particular SNP with two alleles  $a$ ,  $b$ , and one insomnia symptom. For example, we could be speaking of the SNP rs875994 with alleles  $a = C$ ,  $b = T$ , and onset insomnia O.

We imagine to have an arbitrarily large population from the same genetic pool as our sample. In this population we would like to know how large are the fraction  $f_{|a}$  of individuals that show the symptom among those having allele  $a$ , and the fraction  $f_{|b}$  that show the symptom among those having allele  $b$ . These two fractions are *conditional relative frequencies* (note that  $f_{|a} + f_{|b} \neq 1$  in general, since we are not speaking about the frequencies of two mutually exclusive and exhaustive alternatives, but of frequencies *conditional* on such alternatives).  [Add some remarks about the role of 'limit frequencies' and 'arbitrarily large population'? Relation to partial exchangeability and belief about next individual in an endless sequence.](#)

We are in particular interested in knowing how much these two conditional frequencies differ, that is, in  $f_{|a} - f_{|b} =: \Delta f$ . A larger difference may indicate some sort of biologic association between the SNP (or another SNP linked to it) and the symptom. This difference in the large population is unknown, however; we can only make a plausible inference about its

<sup>7</sup> Lange 1995, Lockwood et al. 2001. <sup>8</sup> MacKay et al. 1995. <sup>9</sup> Jiang et al. 2010.



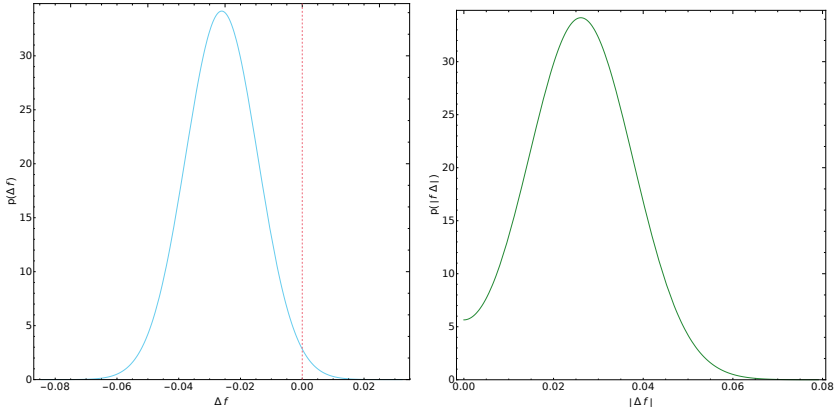


Figure 2 Example of uncertainty distributions about the conditional-frequency difference  $\Delta f$  and its absolute value  $|\Delta f|$ .

value from sampled data and other initial information. The most detailed statement we can make about it is by quantifying the distribution of our degree of belief about  $\Delta f$  or its absolute values,  $p(\Delta f)$  and  $p(|\Delta f|)$ , like the ones plotted in fig. 1. These distributions give us many pieces of information: for example, a negative difference  $\Delta f$  is more plausible than a positive one; it is 8.3% plausible that  $|\Delta f| < 0.01$ , and therefore 91.7% plausible that  $|\Delta f| > 0.01$ ; it is 90% plausible that  $-0.045 < \Delta f < -0.007$ ; and so on. Several measures can be chosen to summarize the distribution with one number. For example, we could use some minimal frequency difference we are highly sure about, for example the lowest 10-quantile:

$$x \text{ such that } p(|\Delta f| > x) = 0.9, \quad (7)$$

or simply the expectation of the absolute value of the difference<sup>10</sup>:

$$E(|\Delta f|) \equiv \sqrt{\frac{2}{\pi}} \sigma \exp\left(-\frac{\mu^2}{2\sigma^2}\right) + \mu \operatorname{erf}\left(\frac{\mu}{\sqrt{2}\sigma}\right), \quad (8)$$

where  $\mu$  and  $\sigma$  are the mean and standard deviation of  $p(\Delta f)$ . For the plot of fig. 1 such measures are 0.0112 and 0.0262. ✚ add maybe something about dependence of broadness on sample size, and ‘smoothing’ as discussed by MacKay & Bauman Peto<sup>11</sup>.

<sup>10</sup> Leone et al. 1961. <sup>11</sup> MacKay et al. 1995 § 2.6.

A more complete quantification of our uncertainty is the distribution  $p(f_{|a}, f_{|b})$  among the possible joint values of the two conditional frequencies, from which  $p(\Delta f)$  can be calculated. This joint distribution is also the optimal starting point when frequencies conditional on allele combinations of several SNP are considered. Our goal is therefore to quantify this joint degree of belief, given: (1) the conditional frequencies in a population sample, (2) our initial information or guesses about such frequencies. In formulae, we want to assign a numerical value to

$$p(\text{conditional frequencies} \mid \text{sample data, initial information}). \quad (9)$$

In the rest of this section we shall calculate this degree of belief by methodically applying the probability calculus.

First let's observe that the approach just outlined is not dichotomous, unlike a classical significance test. We are not asking whether the frequency difference is 'significant' or not. Rather, we find a gradation of cases: from conditional frequencies likely to be very distinct, to conditional frequencies likely to be very similar. These cases can be sorted, for example with the measure of formula (7), obtaining a sequence of SNPs with a decreasing belief of causal association with the symptom. How many of these SNPs are to be selected for further study depends on one's experimental and computational resources.

## 4.2 Inference: concrete calculation

Once we know the quantity we want to quantify our uncertainty about, the calculation of our degrees of belief follows almost mechanically from the rules of the probability calculus<sup>12</sup>. This calculus also shows which degrees of belief we must provide at the start, to arrive at the desired ones.

The first step is to use Bayes's theorem:

$$p(\text{frequencies} \mid \text{data, initial info}) \propto$$

$$p(\text{data} \mid \text{frequencies, initial info}) \times p(\text{frequencies} \mid \text{initial info}), \quad (10)$$

which says that we need to provide two kinds of degree of belief: the plausibility of obtaining our sampled data if we had known the

<sup>12</sup> Jeffreys 1983, Cox 1946, Jaynes 2003, Hailperin 1996.

conditional frequencies in the larger population; and our initial guess about the conditional frequencies, before we observed the data. The first is given by a simple sampling formula. The second can be modelled in several reasonable ways; we'll see, however, that they all lead to very similar conclusions about the conditional frequencies and their difference, owing to the large size of our sample. Let's give explicit expressions for both.

First of all, introduce some notation. The conditional frequencies are denoted  $f_{|a}, f_{|b}$  as above. The data  $D$  consist in the conditional frequencies observed in the sample:  $F_{|a}$  is the fraction showing the symptom among the sampled individuals having allele  $a$ , and  $F_{|b}$  is the analogous fraction for allele  $b$ . Finally, our initial information  $I$  consists in our initial beliefs and in the numbers  $N_a, N_b$  of sampled individuals having allele  $a$  and  $b$ ; although these numbers are part of  $I$  we will sometimes write them explicitly besides  $I$ . (Note that  $N_a, N_b$  are not part of the 'data' in the sense that they don't help us to update our belief about the frequencies  $f_{|a}, f_{|b}$  and therefore remain on the right side of the conditional.)

With this notation, Bayes's theorem above is written

$$p(f_{|a}, f_{|b} \mid D, I) \propto p(D \mid f_{|a}, f_{|b}, I) p(f_{|a}, f_{|b} \mid I) \equiv \\ p(F_{|a}, F_{|b} \mid f_{|a}, f_{|b}, N_a, N_b, I) p(f_{|a}, f_{|b} \mid N_a, N_b, I). \quad (10)_r$$

### 4.3 Plausibility for the data given the frequencies

The plausibility of obtaining the data given the frequencies is given by simple counting and symmetry. We have an arbitrarily large population where a fraction  $f_{|a}$  of individuals having allele  $a$  show the symptom, and a fraction  $1 - f_{|a}$  therefore don't show the symptom. We casually observe  $N_a$  individuals with allele  $a$ . The plausibility that a fraction  $F_{|a}$  of these show the symptom is then counted as a 'drawing with replacement' because of the arbitrarily large size of the full population, given by a (rescaled) binomial distribution Jaynes 2003 ch. 3, Ross 2010 § 4.6, Feller 1968 § VI.2:

$$p(F_{|a} \mid f_{|a}, N_a, I) = \binom{N_a}{N_a F_{|a}} f_{|a}^{N_a F_{|a}} (1 - f_{|a})^{N_a (1 - F_{|a})}. \quad (11)$$

An analogous reasoning holds for allele  $b$ . Our belief about obtaining the data is therefore

$$p(D | f_{|a}, f_{|b}, I) = \prod_{x=a,b} \binom{N_x}{N_x F_{|x}} f_{|x}^{N_x F_{|x}} (1 - f_{|x})^{N_x (1-F_{|x})}. \quad (12)$$

#### 4.4 Initial plausibility for the frequencies

We consider two different initial assumptions about the conditional frequencies. Their densities are represented in fig. 3. The first (left plot) assigns equal plausibility to equal areas in the  $(f_{|a}, f_{|b})$ -coordinate plane; we call this the ‘uniform’ state of knowledge and denote it  $I_u$ . The second (right plot) expresses a strong belief that the two frequencies should be equal, as clear from the steep density increase towards the diagonal of the  $(f_{|a}, f_{|b})$ -coordinate plane; we call this a ‘conservative’ state of knowledge and denote it  $I_c$ .

The uniform state of knowledge is given by the simple density

$$p(f_{|a}, f_{|b} | I_u) = 1. \quad (13)$$

The conservative state of knowledge is given by the following integral:

$$p(f_{|a}, f_{|b} | I_c) = \int_0^\infty du \int_0^\infty dv p(f_{|a}, f_{|b} | u, v, I_c) p(u, v | I_c) \quad (14a)$$

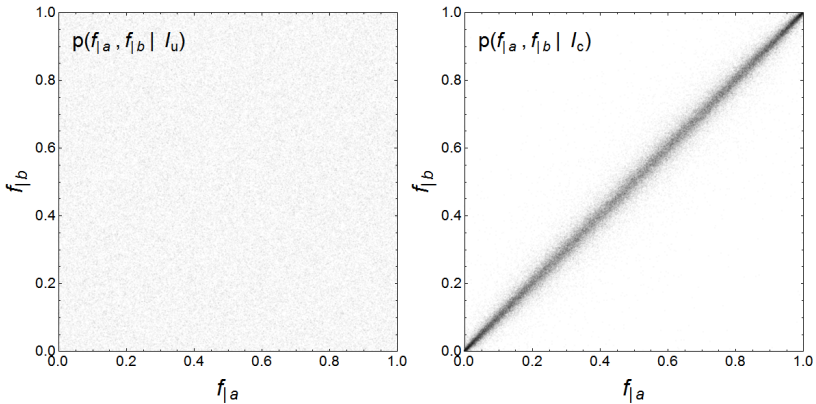


Figure 3 Sampling plots of our two initial states of knowledge, 100 000 samples each. Left: uniform belief (13); right: conservative belief (14).

with

$$p(f_{|a}, f_{|b} \mid u, v, I_c) := \beta(f_{|a} \mid u, v) \beta(f_{|b} \mid u, v), \quad (14b)$$

$$p(u, v \mid I_c) := \frac{1}{u+v} \gamma(u+v \mid 1, 1000) \quad (14c)$$

where  $\beta$  and  $\gamma$  are beta and gamma densities:

$$\beta(f \mid u, v) := \frac{\Gamma(u+v)}{\Gamma(u)\Gamma(v)} f^{u-1} (1-f)^{v-1}, \quad u, v > 0, \quad (14d)$$

$$\gamma(\alpha \mid 1, 1000) := \frac{1}{1000} \exp(-\alpha/1000). \quad (14e)$$

The integral expression (14) can be interpreted as follows: we consider several independent beliefs about the two frequencies – the product of beta densities, depending on identical shape parameters – and then weigh and mix them together to arrive at a belief that is not independent about the two frequencies. This also corresponds to defining a so-called hierarchic model<sup>13</sup>. The particular shape and scale values of the gamma density are chosen so as to concentrate our initial belief along the  $(f_{|a}, f_{|b})$  diagonal. The use of beta densities in this problem seems to have been first endorsed by G. F. Hardy<sup>14</sup>; Good<sup>15</sup> motivates the use of their mixtures. We discuss further interpretations of this initial state of knowledge in appendix\*\*\*.

We will examine how much our conclusions about the frequencies differ between the two initial states of knowledge. Note that the second is indeed very conservative. Therefore, if our updated belief conditional on the data will show clearly distinct conditional frequencies, it will be because the data have given enough evidence to overwhelm this conservative belief.

## 4.5 Final plausibilities

Now we have the plausibility (12) for the sampled data, and two choices of initial plausibilities (13), (14), and we can multiply and renormalize these in Bayes's theorem (10) to obtain our final degrees of belief.

In the case of the uniform initial state of knowledge, multiplication of the sampling formula (12) and initial degree of belief (13) leads to

<sup>13</sup> Good 1980. <sup>14</sup> Hardy 1889. <sup>15</sup> Good 1965 § 4.1, 1980 § 4.

a density for the frequencies proportional to the product of two beta densities. Our final degree of belief in this case is

$$p(f_{|a}, f_{|b} \mid D, I_u) = \prod_{x=a,b} \beta[f_{|x} \mid N_x F_{|x} + 1, N_x (1 - F_{|x}) + 1] \quad (15)$$

This is the product of two independent beliefs about the frequencies.

In the case of the conservative initial state of knowledge, multiplication of (12) and the beta densities within the integral of eq. (14) leads to two new unnormalized beta densities. Grouping the ensuing normalization constant with the density for the parameters (14c) yields a final degree of belief in integral form:

$$p(f_{|a}, f_{|b} \mid D, I_c) = \int_0^\infty du \int_0^\infty dv \prod_{x=a,b} \{\beta[f_{|x} \mid N_x F_{|x} + u, N_x (1 - F_{|x}) + v]\} p(u, v \mid D, I_c) \quad (16a)$$

with

$$p(u, v \mid D, I_c) \propto \prod_{x=a,b} \left\{ \frac{\Gamma(N_x F_{|x} + u) \Gamma[N_x (1 - F_{|x}) + v]}{\Gamma(N_x + u + v)} \frac{\Gamma(u + v)}{\Gamma(u) \Gamma(v)} \right\} \frac{\gamma(u + v \mid 1, 1000)}{u + v}. \quad (16b)$$

This is a hierarchic model. The integral doesn't have an analytic form, but the final density for the frequencies can be sampled by generating samples of the shape parameters  $u, v$  from (16b) via Markov-Chain Monte Carlo methods, and for each generating samples of the frequencies from the beta densities with those shape parameters (alternatively we can write an averaged sum of beta distributions having the sampled shape parameters).

Thanks to the large sample sizes  $N_a, N_b$  the final densities above can be well approximated by normal densities. We only need to determine their means and covariance matrices. These can be easily calculated using the first and second raw moments of a beta distribution: [↗ refs for this](#)

$$E(f \mid u, v) = \frac{u}{u + v}, \quad E(f^2 \mid u, v) = \frac{u(u + 1)}{(u + v)(u + v + 1)}. \quad (17)$$

For the uniform initial state of knowledge (15) the means, variances, and covariance are readily calculated:

$$E(f_{|x} | D, I_u) = \frac{N_x F_{|x} + 1}{N_x + 2}, \frac{N_b F_{|b} + 1}{N_b + 2} \quad x = a, b, \quad (18a)$$

$$\text{var}(f_{|x} | D, I_u) = \frac{(N_x F_{|x} + 1)[N_x(1 - F_{|x}) + 1]}{(N_x + 2)^2(N_x + 3)} \quad x = a, b, \quad (18b)$$

$$\text{cov}(f_{|a}, f_{|b} | D, I_u) = 0. \quad (18c)$$

The covariance vanishes owing to the product form of the final density.

For the conservative initial state of knowledge (16) let us first introduce the notation

$$\langle X(u, v) \rangle := \int_0^\infty du \int_0^\infty dv X(u, v) p(u, v | D, I_c); \quad (19)$$

such an average can approximately be computed by averaging from a Monte Carlo sample. The means, variances, and covariance of the final density are

$$E(f_{|x} | D, I_c) = \left\langle \frac{N_x F_{|x} + u}{N_x + u + v} \right\rangle \quad x = a, b, \quad (20a)$$

$$\text{var}(f_{|x} | D, I_c) = \left\langle \frac{(N_x F_{|x} + u)(N_x F_{|x} + u + 1)}{(N_x + u + v)(N_x + u + v + 1)} \right\rangle - \left\langle \frac{N_x F_{|x} + u}{N_x + u + v} \right\rangle^2 \quad x = a, b \quad (20b)$$

$$\text{cov}(f_{|a}, f_{|b} | D, I_c) = \left\langle \prod_{x=a,b} \frac{N_x F_{|x} + u}{N_x + u + v} \right\rangle - \prod_{x=a,b} \left\langle \frac{N_x F_{|x} + u}{N_x + u + v} \right\rangle. \quad (20c)$$

The covariance does not vanish owing to our initial correlated beliefs about the two conditional frequencies.

It is worth remarking that for our data sample the typical average values of the shape parameters  $\langle u \rangle$ ,  $\langle v \rangle$  are comparable to the sample size  $N_a + N_b$ . Trying an expansion of formulae (20) around the uniform formulae (18) in powers of  $\langle u \rangle$ ,  $\langle v \rangle$  is therefore prone to produce bad approximations.

#### 4.6 Plausibilities for the frequency differences

We finally find an approximate expression for our final degree of belief about the conditional-frequency difference  $\Delta f$ . The density is also

approximately normal, and for each initial state of knowledge  $I = I_u, I_c$  its mean and variance can be found by the standard formulae [give ref for these](#)

$$E(f_{|a} - f_{|b} \mid D, I) = E(f_{|a} \mid D, I) - E(f_{|b} \mid D, I), \quad (21)$$

$$\text{var}(f_{|a} - f_{|b} \mid D, I) = \text{var}(f_{|a} \mid D, I) + \text{var}(f_{|b} \mid D, I) - 2 \text{cov}(f_{|a}, f_{|b} \mid D, I), \quad (22)$$

from the formulae (18), (20) given in the previous section.

## 5 Results



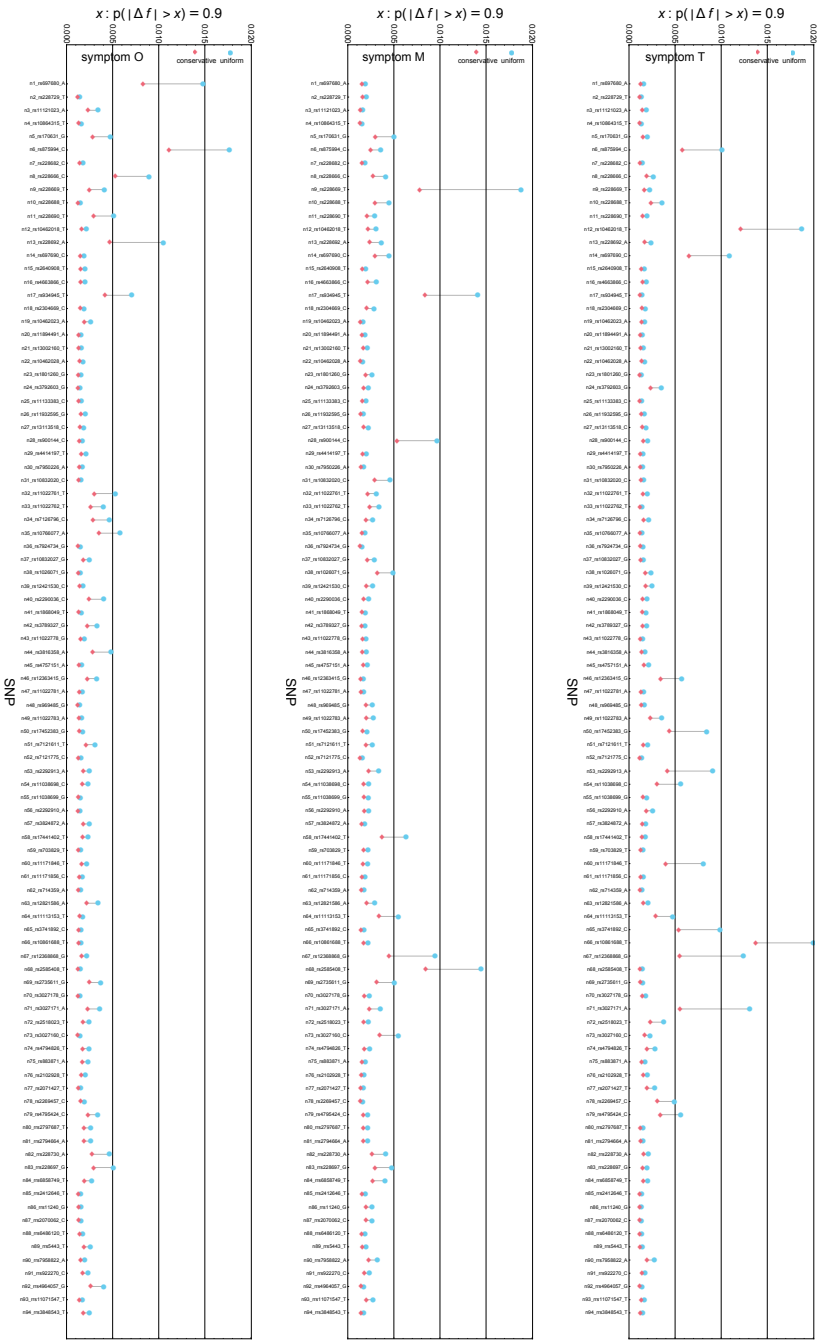


Figure 4 'Relevance' of the 94 SNPs for the three symptoms.

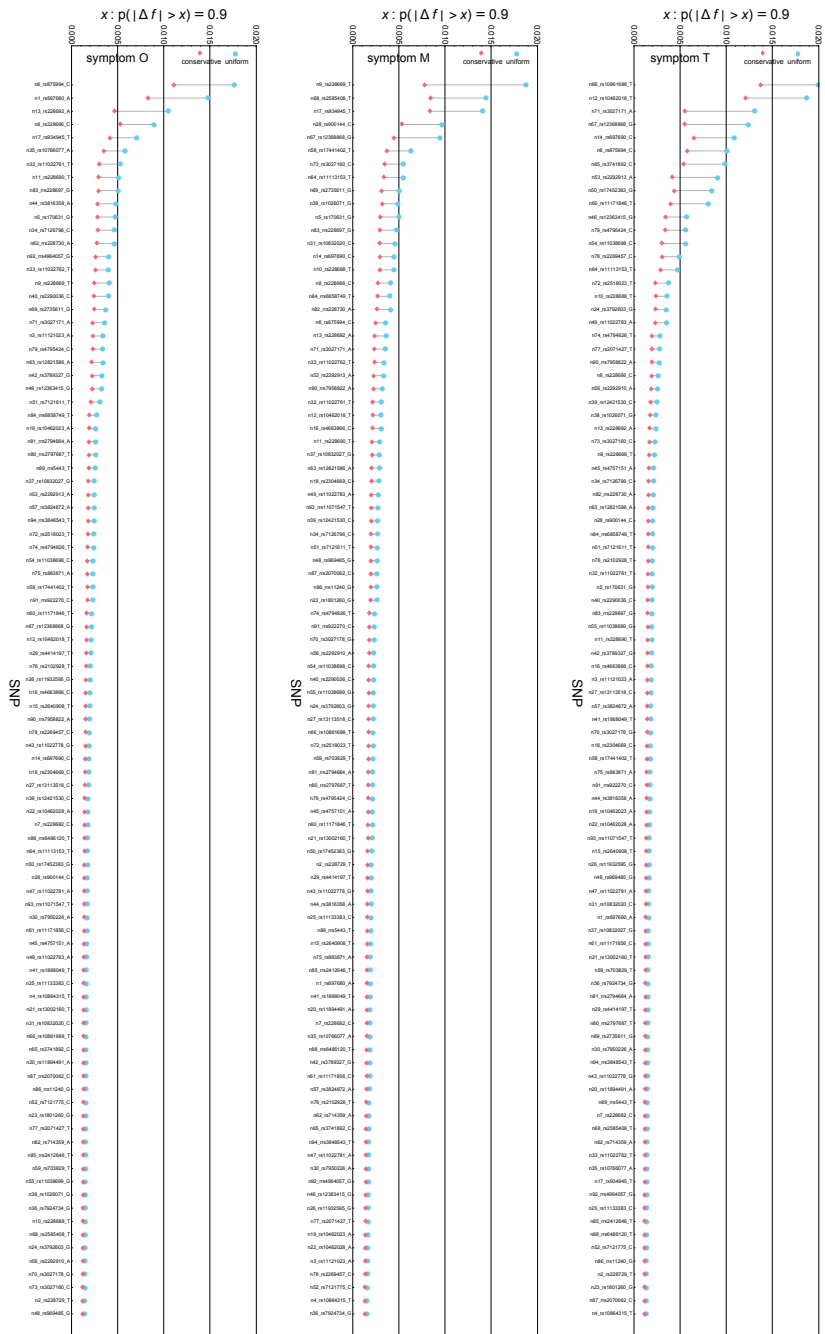


Figure 5 'Relevance' of the 94 SNPs for the three symptoms, sorted.

## Appendices

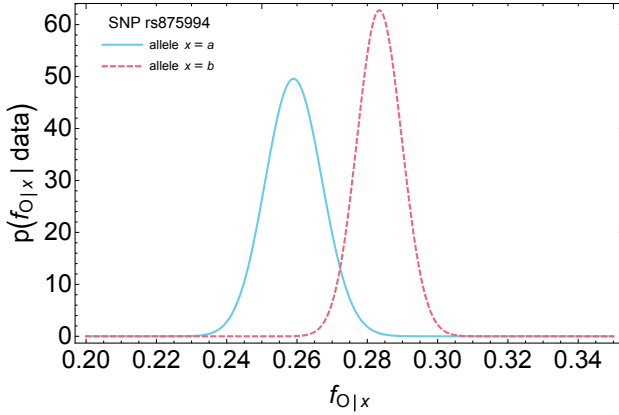


Figure 6 Example of distributions of belief

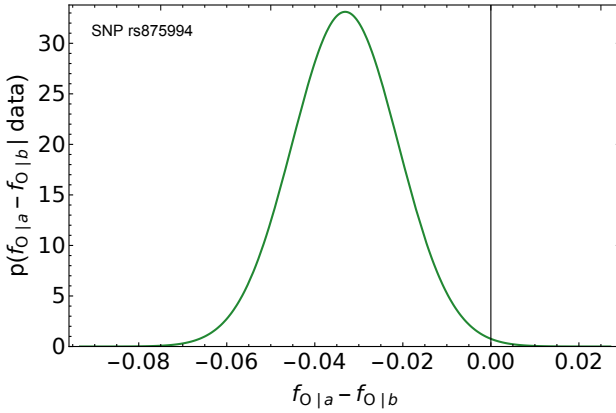


Figure 7 Distribution for the difference between the frequencies  $f_{O|a}$ ,  $f_{O|b}$  of onset insomnia (O) conditional on the two alleles of SNP rs875994

## A Generalizations

### A.1 Generalization to SNP combinations

The ideas and mathematical results discussed in the Method section are easily generalized to the frequencies of a symptom conditional on the combination of the alleles of two or more SNPs. Consider the simplest case of two; denote the two alleles of the first SNP by  $a, b$ , and of the second  $c, d$ . We then have  $2^2$  conditional frequencies  $f_{|ac}, f_{|ad}, f_{|bc}, f_{|bd}$  and we want to quantify the plausibility density  $p(f_{|ac}, f_{|ad}, f_{|bc}, f_{|bd} \mid D, I)$  given the conditional frequencies  $F_{|ac}, F_{|ad}, F_{|bc}, F_{|bd}$  observed in a sample of size  $N_{ac} + N_{ad} + N_{bc} + N_{bd}$ .

Bayes's theorem (10) applies as before. The sampling plausibility (12), obtained by the same reasoning, generalizes to the product over  $x = ac, ad, bc, bd$ . The uniform initial state of knowledge has still unit density as in (13). The conservative initial state of knowledge (14) generalizes to the mixture of the product of  $2^2$  beta distributions, one for each combination of alleles, conditional on the same parameters  $u, v$ . The rest of the calculations proceed analogously.

In this case we have  $\binom{4}{2} \equiv 6$  conditional-frequency differences of interest:  $f_{|ac} - f_{|ad}, f_{|ac} - f_{|bc}, f_{|ac} - f_{|bd}$ , and so on. The analysis of the results therefore offers more interesting possibilities.

### A.2 Generalization to symptom combinations

The Method section is also easily generalized to the conditional frequencies of non-dichotomous symptoms, or of symptom combinations. For example, for our study we could consider the  $2^3$  combinations: individual with no symptoms, or with only symptom O present, or with only symptoms O and M present, and so on up to all three symptoms O, M, T present. We can denote these cases by O, O, OM, and so on. Each of these frequencies would be conditional on whether the individual has allele  $a$  or  $b$  of a particular SNP. We thus have two conditional-frequency distributions:

$$\begin{aligned} (f_{O|a}, f_{O|a}, \dots, f_{OM|a}, \dots, f_{OMT|a}), & \quad \sum_y f_{y|a} = 1, \\ (f_{O|b}, f_{O|b}, \dots, f_{OM|b}, \dots, f_{OMT|b}), & \quad \sum_y f_{y|b} = 1. \end{aligned} \tag{23}$$

### A.3 Choices of initial beliefs and meaning of the parameters

The parameters  $\mathbf{u}$  of the Dirichlet density for the frequency distribution  $\mathbf{f}$  have very intuitive meanings, even more evident if we rewrite them as their sum and a pair of normalized parameters:

$$\alpha := \sum \mathbf{u}, \quad \mathbf{v} := \mathbf{u} / \sum \mathbf{u}. \quad (24)$$

The normalized parameters  $\mathbf{v}$  are the expected frequency distribution:

$$\mathbb{E}(\mathbf{f} \mid \mathbf{u}) = \mathbf{v} \quad (25)$$

The sum of the parameters  $\alpha$  expresses the sharpness of the distribution, as seen from the covariance matrix:

$$\text{var}(f_i \mid \mathbf{u}) = \frac{v_i (1 - v_i)}{\alpha + 1}, \quad \text{cov}(f_i f_j \mid \mathbf{u}) = -\frac{v_i v_j}{\alpha + 1}. \quad (26)$$

In fact, from the update formula\*\*\* we see that  $\alpha$  quantifies the amount of data necessary to modify our initial belief.

$$\frac{\mathbb{E}(f_{O|a} - f_{O|b} \mid D, I)}{\sigma(f_{O|a} - f_{O|b} \mid D, I)}$$

\*\*\*\*\*

This construction has the interesting consequence for our updated belief  $p(f_{|a}, f_{|b} \mid D, I)$ . According to Bayes's theorem (10) combined with our initial belief (??) it is given by

$$p(f_{|a}, f_{|b} \mid D, I) \propto p(D \mid f_{|a}, f_{|b} \mid I) p(f_{|a}, f_{|b} \mid I), \quad (27)$$

but, as shown in appendix B, it can also be written in the following way:

$$p(f_{|a}, f_{|b} \mid D, I) = \int d\mathbf{u} \int d\mathbf{v} p(f_{|a}, f_{|b} \mid D, \mathbf{u}, \mathbf{v}, I) \pi(\mathbf{u}, \mathbf{v} \mid D, I), \quad (28)$$

with the normalized density  $\pi(\mathbf{u}, \mathbf{v} \mid D, I)$  defined by

$$\pi(\mathbf{u}, \mathbf{v} \mid D, I) \propto \pi(\mathbf{u}, \mathbf{v} \mid I) \times$$

$$\underbrace{\int d\mathbf{f}_{|a} \int d\mathbf{f}_{|b} p(D \mid f_{|a}, f_{|b}, I) p(f_{|a}, f_{|b} \mid \mathbf{u}, \mathbf{v}, I)}_{\equiv: p(D \mid \mathbf{u}, \mathbf{v}, I)}. \quad (29)$$

It is as if  $\mathbf{u}, \mathbf{v}$  were unknown parameters with initial belief density  $\pi(\mathbf{u}, \mathbf{v} \mid I)$ , and our update for the frequencies proceeded by first updating

our belief about the parameters, eq. (29), and then marginalizing them out, eq. (28). This is a so-called *hierarchical* model<sup>16</sup>. This hierarchical way of thinking often helps in constructing densities that better represent our initial beliefs, and also leads to formulae that can be better approximated when exact computation is unfeasible. A point rarely emphasized in the literature, though, is that there is no mathematical difference between a hierarchic and a non-hierarchic model: we could forget about the integrals in formula (??) and about the update formula (28), and simply treat  $p(f_{|a}, f_b | I)$  as the density depicted in fig. ??, with update (10). The results would be the same. Further discussion about the formulae above is given in §\*\*\*.

With large sample sizes the density  $\pi(u, v | D, I)$  turns out to be so peaked with respect to  $p(f_{|a}, f_b | u, v, I)$  that it can be considered as a Dirac delta centred on the parameters  $u_M, v_M$  that maximize it. We thus obtain a good approximation of the updated belief (28) that doesn't involve parameter integration:

$$p(f_{|a}, f_b | D, I) \approx p(f_{|a}, f_b | u_M, v_M, I) \quad \text{with} \quad (u_M, v_M) := \arg \max_{u, v} \pi(u, v | D, I). \quad (30)$$

The maximum of  $\pi(u, v | D, I)$ , or better of its logarithm, can easily be found with most optimization methods. The explicit expression to be optimized, discussed in appendix\*\*\*,

$$k [\ln \Gamma(\sum_s \mathbf{u}_s) - \sum_s \ln \Gamma(\mathbf{u}_s)] - \sum_{x=a, b} [\ln \Gamma(N_x + \sum_s \mathbf{u}_s) - \sum_s \ln \Gamma(N_x F_{s|x} + \mathbf{u}_s)] \quad (31)$$

## B Derivation of Bayes's theorem in hierarchic form

We write Bayes's theorem (27) with our initial belief (??) written in full:

$$p(f_{O|a}, f_{O|b} | D, I) \propto \int d\alpha \int d\mathbf{v} \, p(D | f_{O|a}, f_{O|b} | I) p(f_{O|a}, f_{O|b} | \alpha, \mathbf{v}, I) \pi(\alpha, \mathbf{v} | I). \quad (32)$$

Multiplying and dividing within the integral with the expression

$$p(D | \alpha, \mathbf{v}, I) := \int df_{O|a} \int df_{O|b} \, p(D | f_{O|a}, f_{O|b}, I) p(f_{O|a}, f_{O|b} | \alpha, \mathbf{v}, I) \quad (33)$$

we obtain the alternative form (28)

---

<sup>16</sup> Good 1980.

Combining together the sampling formula (12), the expression of the beta density (14d), and the update formula (29) we obtain

$$\pi(\alpha, \mathbf{v} \mid D, I) \propto \pi(\alpha, \mathbf{v}) \times \left[ \int \mathrm{d}f_{O|a} \int \mathrm{d}f_{O|b} \beta\left(f_{O|a} \mid \alpha + N_a, \frac{\alpha \mathbf{v} + N_a F_{O|a}}{\alpha + N_a}\right) \beta\left(f_{O|b} \mid \alpha + N_b, \frac{\alpha \mathbf{v} + N_b F_{O|b}}{\alpha + N_b}\right) \right] \quad (34)$$

## C Summary of the main formulae

We have a sample of size  $n$ . We check the subsample of individuals that have a particular allele, say Bx, for a particular gene, say rs697680\_A. Suppose that in this subsample  $n_0$  individuals *don't* show symptom A and  $n_1$  *do* show symptom A. This also means that the size of our subsample (individuals with allele Bx) is  $n := n_0 + n_1$ .

Our degree of belief about the frequency  $f_1$  of symptom A among the individuals with allele Bx in an *infinite* population is a Beta distribution with parameters  $n_0 + \theta_0$ ,  $n_1 + \theta_1$ , with  $\theta := \theta_0 + \theta_1$ :

$$p(f_1 \mid n_0, n_1, \theta_0, \theta_1) \mathrm{d}f_1 = \frac{\Gamma(n + \theta)}{\Gamma(n_0 + \theta_0) \Gamma(n_1 + \theta_1)} (1 - f_1)^{n_0 + \theta_0 - 1} f_1^{n_1 + \theta_1 - 1} \mathrm{d}f_1 \quad (35)$$

This distribution has expected value and variance

$$\begin{aligned} E(f_1 \mid n_0, n_1, \theta_0, \theta_1) &= \frac{n_1 + \theta_1}{n + \theta}, \\ \text{var}(f_1 \mid n_0, n_1, \theta_0, \theta_1) &= \frac{(n_0 + \theta_0)(n_1 + \theta_1)}{(n + \theta)^2 (n + \theta + 1)}. \end{aligned} \quad (36)$$

✚ Possible further developments: use of hyper-Dirichlet priors, use of graphical models to infer causal relationships<sup>17</sup>

---

<sup>17</sup> Pearl 2009.

$$\begin{aligned}
p(f | D) &\propto p(D | f) p(f | I) \\
&\propto p(D | f) \int du p(f | u) p(u | I) \\
&\propto \int du p(D | f) p(f | u) p(u | I) \\
&\propto \int du p(D | f, u) p(f | u) p(D | u) p(u | I) \\
&\propto \int du p(f | D, u) p(u | D) \\
p(u | D) &\propto \int df p(D | f) p(f | u) p(u | I) \\
&\propto \prod_{x=a,b} \left[ \frac{\Gamma(\sum_s u_s)}{\prod_s \Gamma(u_s)} \frac{\prod_s \Gamma(N_x F_{s|x} + u_s)}{\Gamma(N_x + \sum_s u_s)} \right] \\
p(f_{O|a}, f_{O|b} | D, I) &\approx p(f_{O|a}, f_{O|b} | u_M, v_M, I)
\end{aligned}$$

## Thanks

PGLPM thanks Mari & Miri for continuous encouragement and affection; Buster Keaton and Saitama for filling life with awe and inspiration; and the developers and maintainers of L<sup>A</sup>T<sub>E</sub>X, Emacs, AUC<sub>T</sub>E<sub>X</sub>, Open Science Framework, Python, Inkscape, Sci-Hub for making a free and impartial scientific exchange possible.

## Bibliography

- (‘de X’ is listed under D, ‘van X’ under V, and so on, regardless of national conventions.)
- Amrhein, V., Greenland, S., McShane, B. (2019): *Retire statistical significance*. *Nature* **567**<sup>7748</sup>, 305–307.
- ASA (American Statistical Association) (2016): *ASA statement on statistical significance and p-values*. *American Statistician* **70**<sup>2</sup>, 131–133. Ed. by R. L. Wasserstein. See also introductory editorial in Wasserstein, Lazar (2016) and discussion in Greenland, Senn, Rothman, Carlin, Poole, Goodman, Altman, Altman, et al. (2016).
- (2019): *Moving to a world beyond “p < 0.05”*. *American Statistician* **73**<sup>S1</sup>, 1–19. Ed. by R. L. Wasserstein, A. L. Schirm, N. A. Lazar.
- Berger, J. O., Berry, D. A. (1988): *Statistical analysis and the illusion of objectivity*. *Am. Sci.* **76**<sup>2</sup>, 159–165. <http://drsmorey.org/research/rdmorey/bibtex/upload/Berger:Berry:1988.pdf>.



- Berger, J. O., Delampady, M. (1987): *Testing precise hypotheses*. Stat. Sci. **2**<sup>3</sup>, 317–335. See also comments and rejoinder in Cox, Eaton, Zellner, Bayarri, Casella, Berger, Kadane, Berger, et al. (1987).
- Cox, D. R., Eaton, M. L., Zellner, A., Bayarri, M. J., Casella, G., Berger, R. L., Kadane, J. B., Berger, J. O., et al. (1987): *[Testing precise hypotheses:] Comments and rejoinder*. Stat. Sci. **2**<sup>3</sup>, 335–352. See Berger, Delampady (1987).
- Cox, R. T. (1946): *Probability, frequency, and reasonable expectation*. Am. J. Phys. **14**<sup>1</sup>, 1–13. <http://jimbeck.caltech.edu/summerlectures/references/ProbabilityFrequencyReasonableExpectation.pdf>, [https://wwwusers.ts.infn.it/~milotti/Didattica/Bayes/Cox\\_1946.pdf](https://wwwusers.ts.infn.it/~milotti/Didattica/Bayes/Cox_1946.pdf).
- Feller, W. (1968): *An Introduction to Probability Theory and Its Applications*. Vol. I, 3rd ed. (Wiley, New York). First publ. 1950.
- Good, I. J. (1965): *The Estimation of Probabilities: An Essay on Modern Bayesian Methods*. (MIT Press, Cambridge, USA).
- (1980): *Some history of the hierarchical Bayesian methodology*. Trabajos de Estadística y de Investigación Operativa **31**<sup>1</sup>, 489–519. Repr. in Good (1983) ch. 9 pp. 95–105.
- (1983): *Good Thinking: The Foundations of Probability and Its Applications*. (University of Minnesota Press, Minneapolis, USA).
- Greenland, S., Senn, S. J., Rothman, K. J., Carlin, J. B., Poole, C., Goodman, S. N., Altman, D. G., Altman, N. S., et al. (2016): *Online supplement and discussion: ASA statement on statistical significance and p-values*. American Statistician **70**<sup>2</sup>, 129. <http://www.tandfonline.com/doi/suppl/10.1080/00031305.2016.1154108>. See ASA (2016) and Wasserstein, Lazar (2016).
- Hailperin, T. (1996): *Sentential Probability Logic: Origins, Development, Current Status, and Technical Applications*. (Associated University Presses, London).
- Hardy, G. F. (1889): *[Correspondence]*. Summarized in Whittaker (1921), pp. 174–182.
- Jaynes, E. T. (2003): *Probability Theory: The Logic of Science*. (Cambridge University Press, Cambridge). Ed. by G. Larry Bretthorst. First publ. 1994. <https://archive.org/details/XQUHIUXHIQUHIQUXUIHX2>, <http://www-biba.inrialpes.fr/Jaynes/prob.html>.
- Jeffreys, H. (1983): *Theory of Probability*, third ed. with corrections. (Oxford University Press, London). First publ. 1939.
- Jiang, W., Zhang, C.-H. (2010): *Empirical Bayes in-season prediction of baseball batting averages*. IMS Coll. **6**, 263–273.
- Johnson, D. H. (1999): *The insignificance of statistical significance testing*. J. Wildl. Manage. **63**<sup>3</sup>, 763–772. [http://www.ecologia.ufrgs.br/~adrimelo/lm/apostilas/critic\\_to\\_p-value.pdf](http://www.ecologia.ufrgs.br/~adrimelo/lm/apostilas/critic_to_p-value.pdf).
- Lange, K. (1995): *Applications of the Dirichlet distribution to forensic match probabilities*. Genetica **96**<sup>1–2</sup>, 107–117.
- (2003): *Mathematical and Statistical Methods for Genetic Analysis*, corr. pr. of 2nd ed. (Springer, New York). First publ. 1997.
- Leone, F. C., Nelson, L. S., Nottingham, R. B. (1961): *The folded normal distribution*. Technometrics **3**<sup>4</sup>, 543–550.
- Lewinger, J. P., Conti, D. V., Baurley, J. W., Triche, T. J., Thomas, D. C. (2007): *Hierarchical Bayes prioritization of marker associations from a genome-wide association scan for further investigation*. Genet. Epidemiol. **31**<sup>8</sup>, 871–882.
- Lidstone, G. J. (1921): *Note on the general case of the Bayes-Laplace formula for inductive or a posteriori probabilities*. Trans. Faculty Actuaries **8**<sup>77</sup>, 182–192. See Whittaker (1921) and also Low, Lidstone, Armstrong, Thomson, Sprague, Nicholl, Whittaker (1921).

- Lockwood, J. R., Roeder, K., Devlin, B. (2001): *A Bayesian hierarchical model for allele frequencies*. Genet. Epidemiol. **20**<sup>1</sup>, 17–33. <http://www.stat.cmu.edu/~roeder/publications/lrd2001.pdf>.
- Low, G. M., Lidstone, G. J., Armstrong, J. R., Thomson, W. L., Sprague, A. E., Nicholl, C. C., Whittaker, E. T. (1921): *Discussion [On some disputed questions of probability]*. Trans. Faculty Actuaries **8**<sup>77</sup>, 193–206. See Whittaker (1921) and Lidstone (1921).
- MacKay, D. J. C., Bauman Peto, L. C. (1995): *A hierarchical Dirichlet language model*. Nat. Lang. Eng. **1**<sup>3</sup>, 289–307.
- Pearl, J. (2009): *Causality: Models, Reasoning, and Inference*, 2nd ed. (Cambridge University Press, Cambridge). First publ. 2000.
- Ross, S. (2010): *A First Course in Probability*, 8th ed. (Pearson, Upper Saddle River, USA). First publ. 1976.
- Stephens, M., Balding, D. J. (2009): *Bayesian statistical methods for genetic association studies*. Nat. Rev. Genet. **10**, 681–690.
- Stephens, M., Donnelly, P. (2003): *A comparison of Bayesian methods for haplotype reconstruction from population genotype data*. Am. J. Hum. Genet. **73**<sup>5</sup>, 1162–1169.
- 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.
- Wasserstein, R. L., Lazar, N. A. (2016): *The ASA’s statement on p-values: context, process, and purpose*. American Statistician **70**<sup>2</sup>, 129–133. See ASA (2016) and discussion in Greenland, Senn, Rothman, Carlin, Poole, Goodman, Altman, Altman, et al. (2016). [https://catalyst.harvard.edu/pdf/biostatseminar/ASA\\_s\\_statement\\_on\\_p\\_values\\_context\\_process\\_and\\_purpose.pdf](https://catalyst.harvard.edu/pdf/biostatseminar/ASA_s_statement_on_p_values_context_process_and_purpose.pdf).
- Whittaker, E. T. (1921): *On some disputed questions of probability*. Trans. Faculty Actuaries **8**<sup>77</sup>, 163–182. See also Lidstone (1921) and Low, Lidstone, Armstrong, Thomson, Sprague, Nicholl, Whittaker (1921).