

# Data don't speak for themselves

## a Bayesian course

JOÃO P. FARIA

*Institute of Astrophysics and Space Sciences, Porto*

February 5, 2016

### Abstract

Bayesian statistics is rising in popularity in the astrophysical literature. It is no longer a debate: “work in Bayesian statistics now focuses on applications, computations, and models. Philosophical debates [...] are fading to the background” [1]. This is happening for two main reasons: faster computers and more complex models. In order to keep up, it is important to understand the fundamentals of Bayesian statistics, but it is as important to know how to deal with data analysis applications. In this course I want to provide a brief introduction to advanced concepts in Bayesian statistics. Emphasis will be on *intuition* and *computation*. No coin tossing, only real applications that relate to our day-to-day problems.

## Contents

<b>A (tiny) primer on philosophy</b>	<b>3</b>
<b>1 The basics of probability theory</b>	<b>3</b>
1.1 Some familiarity . . . . .	4
<b>2 Assigning probabilities</b>	<b>6</b>
2.1 The prior . . . . .	6
2.2 The likelihood . . . . .	9
<b>3 MCMC</b>	<b>9</b>
3.1 A basic MCMC implementation . . . . .	9

## **Before you read**

This document is not finished!

What I have written here is not completely original work: I follow closely, and base the discussion on, many different sources. Because it's not finished, I have not yet made fair attribution to all these sources.

I mean no infringement of copyright and no plagiarism.

draft

## A (tiny) primer on philosophy

While preparing for this course, I asked people what it was they wanted to learn about. Much to my dismay, no one answered with “the philosophy of Bayesian statistics”. Most want to learn about the “theory” and “applications” (read MCMC). They want to know how they can *use* Bayesian statistics in their work<sup>1</sup>.

Fair enough. But bypassing philosophy means making assumptions. Therefore, I shall introduce here Jaynes’ *robot*, an imaginary being whose brain reasons according to certain definite rules [3]. These rules will simply be stated, not derived and not defended (Sect. 1). They will be taken as being true. Then, everything that follows logically from them, everything the robot does, will be true.

In the words of [2], these rules make Bayesian inference “more like integration” in that it provides a “collection of tools for finding definitive answers to well-posed problems”. Posing the problem will require subjective knowledge; solving it will not.

### 1 The basics of probability theory

Here we present the rules of the game, following very closely (too closely!) the presentation in [4]. Using Bayesian statistics to analyse data (some might call it doing probabilistic inference) means working with likelihoods, prior probabilities, posterior probabilities and marginalization of nuisance parameters. All this will be explained later but the options we have when working with these mathematical objects are strongly constrained by the rules of probability calculus.

*On voit, par cet Essai, que la théorie des probabilités n’est, au fond,  
que le bon sens réduit au calcul;*

– Pierre-Simon Laplace

If we have a continuous parameter  $a$ , and a probability distribution function  $p(a)$  for<sup>2</sup>  $a$ , it must obey the normalization condition

$$1 = \int p(a) da \quad (1)$$

where the integral is limited by the domain of  $a$ .

Often times we will have more than one parameter. Even if we *condition*  $p(a)$  on some particular value of another parameter  $b$ , that is, we ask for

---

<sup>1</sup>I later read [2] and realised this is not the case for our institute only.

<sup>2</sup>I say *for* instead of *of*, but I promised no philosophy...

$p(a|b)$  (read “the pdf for  $a$  given  $b$ ”), it must obey the same normalization

$$1 = \int p(a|b) da \quad (2)$$

If we have a probability distribution for two things,  $p(a, b)$  (read “the pdf for  $a$  and  $b$ ”), you can always factorize it into two distributions, one for  $a$ , and one for  $b$  given  $a$  or the other way around:

$$p(a, b) = p(a) p(b|a) \quad (3)$$

$$p(a, b) = p(b) p(a|b) \quad (4)$$

These two factorisations together lead<sup>3</sup> to Bayes’ theorem:

$$p(a|b) = \frac{p(b|a) p(a)}{p(b)}. \quad (5)$$

Conditional probabilities factor just the same as unconditional ones

$$p(a, b|c) = p(a|c) p(b|a, c) \quad (6)$$

$$p(a, b|c) = p(b|c) p(a|b, c) \quad (7)$$

$$p(a|b, c) = \frac{p(b|a, c) p(a|c)}{p(b|c)} \quad (8)$$

where we just carried the condition  $c$  through all the terms.

You can integrate out or *marginalize* variables you want to get rid of (or *not* infer) by integrals like

$$p(a|c) = \int p(a, b|c) db \quad (9)$$

$$p(a|c) = \int p(a|b, c) p(b|c) db \quad (10)$$

where the second is a factorized version of the first. Remember that, since  $b$  can be a very high-dimensional mathematical object (a set of parameters), integrals like these can be extremely difficult to calculate in practice.

## 1.1 Some familiarity

Let us write some of the preceding equations with more familiar terms<sup>4</sup>. We have data  $D$  and a set of parameters  $\theta$  that we are interested in learning about. In all our analyses we condition on information  $\mathcal{I}$  that we have about the world. Then we can write Eq. (8) as

$$p(\theta|D, \mathcal{I}) = \frac{p(\theta|\mathcal{I}) p(D|\theta, \mathcal{I})}{p(D|\mathcal{I})} \quad (11)$$

<sup>3</sup>Bayes’ theorem is a consequence of the previous equations, it is not assumed as a rule. That’s why it’s called a theorem, actually.

<sup>4</sup>Only if you have heard about Bayesian statistics, otherwise just different terms.

The terms in this equation are usually called

$p(\theta|D, \mathcal{I})$  the posterior distribution

$p(\theta|\mathcal{I})$  the prior distribution

$p(D|\theta, \mathcal{I})$  the likelihood

$p(D|\mathcal{I})$  the evidence, sometimes denoted with a  $\mathcal{Z}$ .

but calling them this may hide something important (see Section 2.2).

From Eq. (2) we can derive<sup>5</sup> that

$$p(D|\mathcal{I}) = \mathcal{Z} = \int p(\theta|\mathcal{I}) p(D|\theta, \mathcal{I}) d\theta \quad (12)$$

Because  $\mathcal{Z}$  does not depend on  $\theta$ , you might see many times Eq. (11) written as

$$p(\theta|D, \mathcal{I}) \propto p(\theta|\mathcal{I}) p(D|\theta, \mathcal{I}), \quad (13)$$

so **the posterior is proportional to the likelihood times the prior**. Nevertheless, try to stick with Eq. (11);  $\mathcal{Z}$  contains in it a big deal of information.



#### think about it

Eq. (13) means that the likelihood and the prior can be defined up to an arbitrary multiplicative constant (i.e. not a function of  $\theta$ ).

For example, imagine we are interested in comparing two models  $M_1$  and  $M_2$  which have parameters  $\theta_1$  and  $\theta_2$ , respectively. We still only have data  $D$ . Bayes' theorem works the same:

$$p(M_i|D, \mathcal{I}) = \frac{p(M_i|\mathcal{I}) p(D|M_i, \mathcal{I})}{p(D|\mathcal{I})} \quad (14)$$

and the ratio of model probabilities is

$$\frac{p(M_1|D, \mathcal{I})}{p(M_2|D, \mathcal{I})} = \frac{p(M_1|\mathcal{I}) p(D|M_1, \mathcal{I})}{p(M_2|\mathcal{I}) p(D|M_2, \mathcal{I})} \quad (15)$$

$$= \frac{p(M_1|\mathcal{I}) \int p(D, \theta_1|M_1, \mathcal{I}) d\theta_1}{p(M_2|\mathcal{I}) \int p(D, \theta_2|M_2, \mathcal{I}) d\theta_2} \quad (16)$$

$$= \frac{p(M_1|\mathcal{I}) \int p(\theta_1|M_1, \mathcal{I}) p(D|\theta_1, M_1, \mathcal{I}) d\theta_1}{p(M_2|\mathcal{I}) \int p(\theta_2|M_2, \mathcal{I}) p(D|\theta_2, M_2, \mathcal{I}) d\theta_2} . \quad (17)$$

See how the evidence (of both models) turns out to be important!

<sup>5</sup>Multiply both sides of Eq. (8) by  $p(D|\mathcal{I})$  and integrate over  $\theta$ , noting that  $p(D|\mathcal{I})$  does not depend on  $\theta$  and that the posterior obeys Eq. (2).

The term

$$\frac{p(M_1|\mathcal{I})}{p(M_2|\mathcal{I})}$$

is the ratio of prior probabilities for the two models. If we believe they are equally likely at the start, then this is equal to 1.



### think about it

$\theta_1$  and  $\theta_2$  can be any set of parameters and, in particular, they can have different sizes. Say model  $M_1$  has 2 parameters  $\theta_1 = (a, b)$  and model  $M_2$  only has one parameter  $\theta_2 = (c)$ . Then the integrals in Eq. (17) are of different dimensionality. That's fine, probability theory doesn't care. Enough of that "divide by the degrees of freedom" nonsense! [5]

## 2 Assigning probabilities

Now that we have a set of rules with which we can manipulate our distributions, it is useful to learn how we can assign values to all these terms and actually start calculating things.

### 2.1 The prior

The dreadful prior, we wish we didn't have to deal with it. But having to specify priors is a good excuse to learn about different probability distributions. See Fig. 1, which shows the general form of many continuous probability distributions.

We can start with the Normal distribution. A variable  $x$  is normally distributed if its pdf is

$$\text{pdf}(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[ -\frac{(x - \mu)^2}{2\sigma^2} \right] \quad (18)$$

where  $\mu$  and  $\sigma$  are the two parameters of the Normal distribution. We will use the following notation to say that  $x$  follows a Normal distribution:

$$x \sim \mathcal{N}(\mu, \sigma)$$

The Normal distribution is defined on the real line, so  $x$  can take any value between  $-\infty$  and  $\infty$ . The parameter  $\mu$  is the mean of the distribution and  $\sigma$  is the standard deviation. Values of  $x$  closer to  $\mu$  are more probable than those far from  $\mu$ .

The uniform distribution is also quite common

$$x \sim \mathcal{U}(a, b) \quad \text{if} \quad \text{pdf}(x|a, b) = \frac{1}{a - b}$$

For the uniform distribution, all values between  $a$  and  $b$  are equally probable. The probability that  $x < a$  is 0 and the same for  $x > b$ . These values of  $x$  are *impossible*.

Let us introduce briefly one more distribution, the  $t$ -distribution:

$$x \sim t(\nu, \mu, \sigma) \quad \text{if} \quad \text{pdf}(x|\nu, \mu, \sigma) = \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})\sqrt{\pi\nu}\sigma} \left[ 1 + \frac{1}{\nu} \left( \frac{x-\mu}{\sigma} \right)^2 \right]^{-\frac{\nu+1}{2}} \quad (19)$$

with the degrees of freedom  $\nu$ , the mean  $\mu$  and the scale  $\sigma$  as parameters. In Eq. (19),  $\Gamma$  is the gamma function.

This distribution looks complicated but I mention it here for one reason: note how in Fig. 1, the normal and  $t$  distributions look very similar, they are both bell-shaped. The distinguishing feature is that the  $t$ -distribution has *heavier tails*, that is, values of  $x$  far from the mean are more probable if  $x$  is  $t$  distributed than if  $x$  is normally distributed. How heavy the tails are depends on the parameter  $\nu$ . For example, imagine that

$$x_1 \sim \mathcal{N}(0, 1) \quad (20)$$

$$x_2 \sim t(2, 0, 1) \quad (21)$$

and you obtain 100 000 random samples from these two distributions. Then you calculate the minimum value of all the samples for  $x_1$  and all the samples for  $x_2$ . Will these values be very different?

---

```

1 from scipy.stats import norm, t
2
3 print min( t(1., loc=0, scale=1).rvs(100000) )
4 print min( norm(loc=0, scale=1).rvs(100000) )

```

---

They are completely different! By orders of magnitude!

The  $t$ -distribution allows for more extreme values. This will come in handy when we deal with outliers.

Other distributions have different pdfs and different parameters. Wikipedia usually has a lot of information about the most common ones, just search for the names you see in Fig. 1.

The key concept about priors is to use (some of) these probability distributions to describe our beliefs about the parameters in our models.

- if any value of the parameter seems equally likely but you have an idea of its order of magnitude, try a uniform distribution in a sensible range.

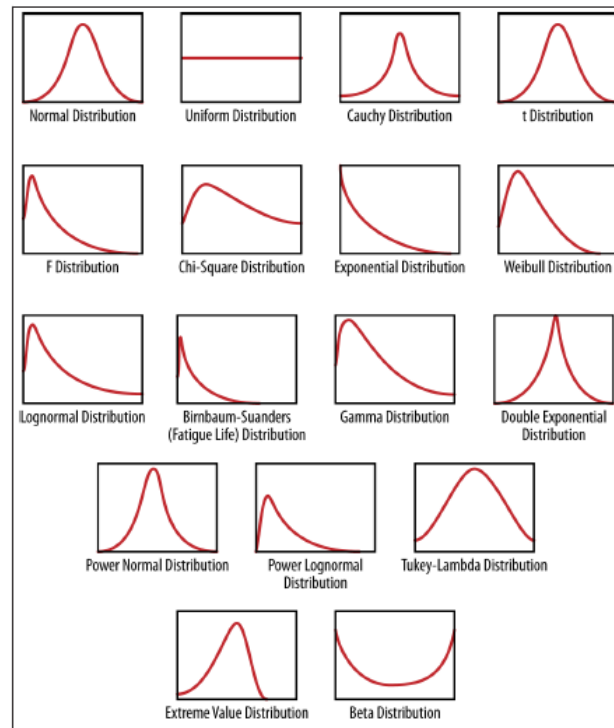


Figure 1: Continuous probability distributions. Adapted from here.

- if one value is somewhat more likely and, again, you know the order of magnitude, try a normal distribution centred in that value.
- if any value is equally likely within a few orders of magnitude, use a log-normal or an exponential distribution. [check!!!]
- if the parameter takes values between 0 and 1, the beta distribution is defined in this interval and can take a wide range of forms, depending on its parameters.

Every time you come up with a prior, you are using your expert knowledge about the problem. Don't underestimate this knowledge by trying to set an "uninformative" prior at all costs. If you are not sure, try a few different priors and see if the results make sense and are consistent.

The only way in which you can go completely wrong and change the results dramatically is by setting the prior to 0. No amount of data can change a prior which is zero.



## 2.2 The likelihood

The likelihood is a prior. I will repeat so you know this is not a typo: the likelihood is a prior. The term  $p(D|\theta, I)$  represents your beliefs on what the data will be like, given parameters  $\theta$  and information  $I$ . Therefore, the likelihood encodes prior information. Following [6], the likelihood is *not*

- the process that generated the data
- the pdf that your data kind of looks like when you plot it in a histogram

It is, instead, what we usually call *the model*. And the model is nothing else than a set of assumptions about the data and how they relate to the parameters. **The likelihood provides a (*the*) connection from  $\theta$  to  $D$ .**

## 3 MCMC

Markov chain Monte Carlo (MCMC) is the workhorse of Bayesian statistics. But how does it work? And what does it actually do? We will write the code for our own MCMC in the practical classes. If that is enough for you to *understand* it, skip the next section. But if you want to know *why* it works, carry on reading.

### 3.1 A basic MCMC implementation

We will only talk about the Metropolis-Hastings algorithm to do MCMC. It is probably not the simplest, but it is definitely the most often used. The following is adapted slightly from Wikipedia:

Let  $f(x)$  be a function that is proportional to the desired probability distribution  $P(x)$  (a.k.a. the target distribution).

#### Initialization:

Choose an arbitrary point  $x_0$  to be the first sample, and choose an arbitrary probability density  $Q(x|y)$  which suggests a candidate for the next sample value  $x$ , given the previous sample value  $y$ . A usual choice is to let  $Q(x|y)$  be a Gaussian distribution centered at  $y$ . The function  $Q$  is referred to as the proposal density or jumping distribution.

#### For each iteration $t$ :

- Generate a candidate  $x'$  for the next sample by picking from the distribution  $Q(x'|x_t)$ .

- Calculate the acceptance ratio  $\alpha = f(x')/f(x_t)$ . Because  $f$  is proportional to  $P$ , we have that  $\alpha = f(x')/f(x_t) = P(x')/P(x_t)$ .
- If  $\alpha \geq 1$ , then the candidate is more likely than  $x_t$ : automatically accept the candidate by setting  $x_{t+1} = x'$ . Otherwise, accept the candidate with probability  $\alpha$ ; if the candidate is rejected, set  $x_{t+1} = x_t$ , instead.

Let's see this in Python code

---

```

1 import numpy as np
2 from scipy.stats import norm
3
4 niter = 500 # number of iterations
5
6 # posterior distribution (the distribution we want to sample from)
7 p = lambda x: norm.pdf(x)
8
9 x = [0.5] # starting point
10
11 for step in range(niter):
12     d = np.random.randn() # propose a step
13
14     # calculate the ratio of posterior probabilities
15     alpha = p(x[-1]+d) / p(x[-1])
16
17     if alpha > 1.:
18         x.append(x[-1]+d) # accept the new step
19     else:
20         u = np.random.uniform()
21         if ratio > u:
22             x.append(x[-1]+d) # accept the new step
23         else:
24             x.append(x[-1]) # reject the new step (stay where we are)

```

---

## References

1. Andrew Gelman. *Bayesian data analysis*. Chapman & Hall/CRC texts in statistical science. CRC Press, Boca Raton, third edition edition, 2014.
2. Thomas J. Loredo. *The Return of the Prodigal: Bayesian Inference in Astrophysics*. Valencia, 1994.
3. E. T. Jaynes and G. Larry Bretthorst. *Probability theory: the logic of science*. Cambridge University Press, Cambridge, UK ; New York, NY, 2003.
4. David W. Hogg. Data analysis recipes: Probability calculus for inference. *ArXiv e-prints*, 1205:4446, May 2012.
5. Rene Andrae, Tim Schulze-Hartung, and Peter Melchior. Dos and don'ts of reduced chi-squared. *ArXiv10123754 Astro-Ph Physicsphysics Stat*, December 2010.
6. B. J. Brewer. The prior isn't the only prior, 2013.