

1 The Basics of Bayesian Statistics

In this course, we will apply statistics to draw conclusions on astrophysical data. For the first few lectures, we therefore study Bayesian statistics to equip ourselves with the tools and abilities that we later apply to cosmology.

First of all let's introduce notations. We denote the probability that some hypothesis "X" is true given some information "I" with

$\text{prob}(X|I)$
probability that "X" is true "given" I

So the bar "|" is read as "given".

The exact opposite of X is denoted by "not-X": \overline{X}

Without loss of generality, we can normalize

$$\text{prob}(X|I) \in [0, 1]$$

There are two axioms that are needed for logical consistency. They are called sum-rule and product rule respectively.

The sum-rule says that the probability of X being true plus the probability of X being false is 100%.

Sum rule: $\text{prob}(X | I) + \text{prob}(\bar{X} | I) = 1$

The second axiom, the product rule says:

product rule: $\text{prob}(X, Y | I) = \text{prob}(X | Y, I) \text{prob}(Y | I)$

As an example, consider:

X: "it's raining"

Y: "there are clouds on the sky" and you are

I: It's the 15th of October in Germany

Then the probability

$\text{prob}(X, Y | I)$: "It's raining and there are clouds on the sky"

is the product of

$\text{prob}(X | Y, I)$: "It's raining given that there are clouds and I"

and

$\text{prob}(Y | I)$: "There are clouds given I"

1.1 Bayes' Theorem and Marginalization

Bayes was the first to state an in principle simple corollary that follows from the product rule and the observation that

$$\text{prob}(X, Y | I) = \text{prob}(Y, X | I)$$

So:

$$\text{prob}(X | I, Y) \text{prob}(Y | I) = \text{prob}(Y | X, I) \text{prob}(X | I)$$

Which leads to Bayes Theorem:

$$\text{prob}(X | Y, I) = \frac{\text{prob}(Y | X, I) \text{prob}(X | I)}{\text{prob}(Y | I)}$$

Bayes theorem is extremely powerful! You'll immediately see this, when we use "hypothesis" and "data" instead of X and Y :

$$\text{prob}(\text{hypothesis} | \text{data}, I) \propto \text{prob}(\text{data} | \text{hypothesis}, I) \times \text{prob}(\text{hypothesis} | I)$$

So you learn about a hypothesis (model) you have given the data and additional information by computing how probable that it is to measure the data you took given the hypothesis (and I) times the probability that the hypothesis is valid given I .

Bayes' Theorem is so important, that each term has its own name:

$$\begin{array}{c} \text{"posterior"} \\ \downarrow \\ \text{prob}(X|Y, I) = \end{array} \frac{\begin{array}{c} \text{"likelihood"} \\ \downarrow \\ \text{prob}(Y|X, I) \end{array} \begin{array}{c} \text{"prior"} \\ \downarrow \\ \text{prob}(X|I) \end{array}}{\begin{array}{c} \text{prob}(Y|I) \\ \uparrow \text{"evidence"} \end{array}}$$

So you infer the "posterior" probability that your hypothesis X is right by the measurement yielding the "likelihood" that your hypothesis explains the (data) Y times your "prior" knowledge or prejudice that your hypothesis X is right divided by the "evidence" (more of this strange name later) that the (data) Y ~~is~~ exists given the Information I .

Very often, the data Y is measured already and you want to test a hypothesis X . So the evidence $\text{prob}(Y|I)$, which is unchanged by your hypothesis ~~can~~ is treated as a rather uninteresting nuisance. ~~that is to be~~

A very useful tool is marginalization. Suppose first that you have a discrete set of propositions $\{Y_k\} = Y_1, Y_2, \dots, Y_M$

For instance with 5 presidential candidates, Y_1 could be the ~~probabilist~~ proposition that candidate "1" wins etc.

The probability that X is true, for instance that unemployment is lower in a year's time irrespective of who is president is then

$$\text{prob}(X | I) = \sum_{k=1}^M \text{prob}(X, Y_k | I)$$

While this seems sensible, let us prove it by using the ~~sum~~^{product} rule. First, we have

$$\text{prob}(X, Y_1 | I) = \text{prob}(Y_1 | X, I) \text{prob}(X | I)$$

Now add all k propositions of Y_k :

$$\begin{aligned} & \text{prob}(X, Y_1 | I) + \text{prob}(X, Y_2 | I) + \dots + \text{prob}(X, Y_k | I) \\ &= \left[\text{prob}(Y_1 | X, I) + \text{prob}(Y_2 | X, I) + \dots + \text{prob}(Y_k | X, I) \right] \\ & \quad \times \text{prob}(X | I) \end{aligned}$$

But as $\{Y_k\}$ is a mutually exclusive and exhaustive set, i.e.

$$\sum_{k=1}^M \text{prob}(Y_k | X, I) = 1,$$

the quantity in square brackets above is 1 and we obtain our desired result.

Now suppose that instead of 5 presidential candidates, we'd have a continuous set of propositions, for instance if Y were the amplitude of some nuisance noise of a telescope.

in this case, it is obvious to generalize the concept of probability to probability density functions (pdf's):

$$\text{pdf}(X, Y=y | \mathcal{I}) = \lim_{\delta y \rightarrow 0} \frac{\text{prob}(X, y \in Y < y + \delta y | \mathcal{I})}{\delta y}$$

Yet, we'll use the notation $\text{prob}()$ instead of $\text{pdf}()$ for notational convenience!

So finally, we can generalize the sum we computed above to the important statement of marginalization over some pa. proportion (parameter) Y which we are not interested in:

Marginalization: $\text{prob}(X | \mathcal{I}) = \int_{-\infty}^{\infty} \text{prob}(X, Y | \mathcal{I}) dy$

2 Parameter Estimation I

Suppose we would like to find out, if a certain coin is fair in the sense that flipping it yields the same chance for heads as it does for tails. Let us denote the bias of the coin (if any) by a parameter $\theta \in [0, 1]$.

If $\theta = 0 \Rightarrow$ all tails

" $\theta = 1 \Rightarrow$ all heads

and an unbiased coin would have $\theta = 1/2$.

Suppose further that we really have no opinion on whether it is biased and on how much it is.

So our prior could be:

$$\text{prob}(\theta | \mathcal{I}) = \begin{cases} 1 & ; 0 \leq \theta \leq 1 \\ 0 & ; \text{otherwise} \end{cases}$$

and Baye's theorem says

$$\text{prob}(\theta | \{\text{data}\}, \mathcal{I}) = \frac{\text{prob}(\{\text{data}\} | \theta, \mathcal{I}) \text{prob}(\theta | \mathcal{I})}{\text{prob}(\{\text{data}\} | \mathcal{I})}$$

Please note that since our posterior will be normalized,

$$\int_0^1 \text{prob}(\theta | \{\text{data}\}, \mathcal{I}) d\theta = 1$$

we can neglect the evidence $\text{prob}(\{\text{data}\} | \mathcal{I})$, because it does not depend on θ and can be retrieved from the normalization condition later, if we really want to know it.

So Baye's theorem says

$$\text{prob}(H | \{data\}, I) \propto \text{prob}(\{data\} | H, I) \text{prob}(H | I)$$

and the likelihood to obtain 'R' heads in 'N' tosses is given by

$$\text{prob}(\{data\} | H, I) \propto H^R (1-H)^{N-R}$$

So

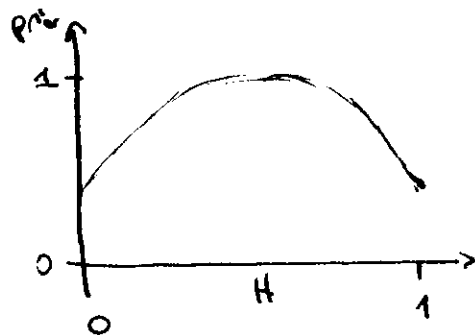
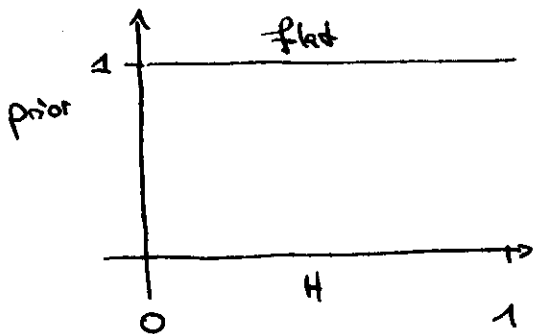
$$\text{prob}(H | \{data\}, I) \propto \text{prob} H^R (1-H)^{N-R} ; 0 \leq H \leq 1$$

and 0 otherwise.

Figure!

The choice of priors

What if we chose a different prior? Say one which gives more weight to $H = 1/2$ than to $H = 0$ or $H = 1$:



As the prior still has $\text{prob}(H | I) \neq 0$ everywhere, it does not exclude the case of $H = 0$ or $H = 1$.

In the long run, the data-taking will take over the prior, for as long as the prior is not too narrow or excludes values of H that are actually the coin's true value.

On the other hand, a sensible prior can improve the quality of a deduction: suppose you measured the length of some sort of fish in a river by taking a random sample of 10 such fishes. If you'd look up the usual length (plus deviations) of that sort of fish, you could improve your estimate of the average length by this prior knowledge, i.e. if you denote the average length as L , then

$$\text{prob}(L | \{\text{data}\}, I) \propto \text{prob}(\{\text{data}\} | L, I) \text{prob}(L | I)$$

and your prior knowledge of $\text{prob}(L | I)$ could help you here? You can also use the prior like here ~~to~~ to incorporate prior experiments.

Say you'd like to infer the Hubble parameter H from the WMAP measurements of the CMB. You could then add the HST measurement of H , to improve on your deduction and account for the prior measurement of H .