

Exploring Monte-Carlo-integration techniques in Bayesian model selection

Jakob Krause* and Dominic Schüchter†

(Dated: March 1, 2021)

An article usually includes an abstract, a concise summary of the work covered at length in the main body of the article.

Usage: Secondary publications and information retrieval purposes.

Structure: You may use the `description` environment to structure your abstract; use the optional argument of the `\item` command to give the category of each item.

I. INTRODUCTION

In physics, one is often faced with the problem of *Model Selection* for a given data set. That means finding a mathematical description of the data, which on the one hand sufficiently characterizes the data structure and on the other hand satisfies the expected dependencies. This is by no means a trivial task; one has to understand the underlying physical model beforehand to not make the mistake of choosing a too complicated model although it may seemingly fit the data. At the same time too trivial assumptions can also lead in the wrong direction. Compactly this problem can be formulated in the following way (adapted from [3, Chap. 4]):

Alice has a theory; Bob also has a theory, but with an adjustable parameter λ . Whose theory should we prefer on the basis of data D ?

BAYESian inference provides quantitative measures for this model-selection problem, e.g. the BAYES-factor and the BAYES-complexity, these have among others been successfully used in astronomy as can be found in [2, 5, 6] respectively. In this paper we will investigate two simulated example problems and apply various measures of bayesian model selection as to find out the true underlying model which was used to generate the simulated data. We will focus on the numerical evaluation of such problems, especially MONTE-CARLO techniques.

II. THEORY

In the following we will give a short introduction into BAYES THEORY and the underlying concepts of model selection.

A. Bayes' Theorem

The fundamental equation of BAYESIAN statistics – for a dataset y and parameters θ – is given by BAYES The-

orem.

$$\text{prob}(\theta|y) = p(\theta|y) = \frac{p(y|\theta) \cdot p(\theta)}{p(y)} \quad (1)$$

Where $p(\theta|y)$ is the posterior probability for parameters θ given the data y , $p(y|\theta)$ is the likelihood that the data fits a model with parameters θ , $p(\theta)$ is the prior probability of θ and $p(y) = \int_{-\infty}^{+\infty} d\theta p(y|\theta)p(\theta)$ is the marginal likelihood which acts as a normalization. In the case of parameter selection, the normalization can often be neglected since it is only a constant. In the case of the model comparison it is a crucial quantity, as we will discuss in section (II C) [3, Chap. 2].

B. Parameter estimation

Assume we want to find the best parameters for a given dataset y and model M . Eq. (1) then can be written as

$$p(\theta|y, M) = \frac{p(y|\theta, M) \cdot p(\theta|M)}{p(y|M)}. \quad (2)$$

Evaluating this equation will give probability density functions (PDFs)

$$p(\theta|D, M) = \int d\varphi p(\theta, \varphi|D, M) \quad (3)$$

for each parameter θ , where $\theta = (\varphi, \theta)$. Because of equation (3) $p(\theta|D, M)$ is also called *marginal posterior*. From this one can find the best fit value either by finding the value of θ for which the marginal posterior is maximized or by calculating the mean with respect to the marginal posterior.

$$p(\theta|D, M) = \max \Leftrightarrow \theta = \hat{\theta} \\ \langle \theta \rangle = \int_{-\infty}^{\infty} d\theta p(\theta|D, M) \cdot \theta \quad (4)$$

Throughout this paper we will use $\langle \theta \rangle$ as our best fit estimate. [3]

C. Model comparison

To compare two (or more) competing models M_i that describe a dataset D let us write BAYES theorem once

* <http://www.github.com/krausejm>; krause@hiskp.uni-bonn.de

† <http://www.github.com/dschuechter>; dschuechter@uni-bonn.de

again

$$p(M_i|D) = \frac{p(M_i) \cdot p(D|M_i)}{p(D)}. \quad (5)$$

We recognize $p(D|M_i)$ as the marginal likelihood of the model M_i . Let now the two competing models be those of Alice and Bob introduced in section (I). Thus,

$$p(D|M_i) = \int_{-\infty}^{\infty} d\lambda p(D|\lambda, M) \cdot p(\lambda|M)$$

follows [3, Chap. 3].

1. Bayes Factor

One way of comparing two models is to compare their posterior odds, that is the ratio **quote lecture A. Wirzba?**

$$O_{ij} := \underbrace{\frac{p(M_i|D)}{p(M_j|D)}}_{\text{posterior odds}} = \underbrace{\frac{p(D|M_i)}{p(D|M_j)}}_{\text{BAYES Factor}} \cdot \underbrace{\frac{p(M_i)}{p(M_j)}}_{\text{prior odds}} = B_{ij} \cdot \frac{p(M_i)}{p(M_j)}. \quad (6)$$

Usually we can assume the prior odds to be ~ 1 because we do not favor one model over the other just on prior beliefs. With this equation (6) simply becomes

$$O_{ij} = B_{ij}.$$

Calculating B_{ij} will be referred to as calculating the BAYES factor *in favor of* model M_i . We can then evaluate the result of B_{ij} as follows

$$B_{ij} = B_{ji}^{-1} = \begin{cases} \geq 1 & \text{model } M_i \text{ favored} \\ < 1 & \text{model } M_j \text{ favored} \end{cases}$$

Nevertheless a BAYES Factor ≥ 1 does not mean that the favored model is far superior. There is an empirical scale which points to the strength of evidence for a given B_{ij} . This is displayed in table (I).

$ \ln B_{ij} $	Odds	Probability	Strength of evidence
< 1.0	$\lesssim 3 : 1$	< 0.750	Inconclusive
1.0	$\sim 3 : 1$	0.750	Weak evidence
2.5	$\sim 12 : 1$	0.923	Moderate evidence
5.0	$\sim 150 : 1$	0.993	Strong evidence

TABLE I: Empirical scale for evaluating the strength of evidence when comparing two models M_i vs. M_j , adapted from [6]

Bayes factor vor allem attraktiv für genestete modelle

2. Bayes Complexity

When comparing models with a different amount of parameters, we obviously denote more parameters as more

complex. A fundamental question of model-selection is then if the given data supports more or less parameters. A naïve measure of complexity would be the number of free parameters a model has, we denote this as \mathcal{C}_0 . A more sophisticated measure of complexity is given by the so called BAYESIAN complexity, first introduced by Spiegelhalter et al [4], which can be written as [2]

$$\mathcal{C}_b = -2 \int d\theta p(\theta|D, M) \log(\mathcal{L}(\theta)) + 2 \log(\mathcal{L}(\tilde{\theta})), \quad (7)$$

with the likelihood $\mathcal{L}(\theta) = p(D|\theta, M)$ and $\tilde{\theta} = \langle \theta \rangle$. \mathcal{C}_b describes how many model parameters the data is able to constrain [2] and is thus a useful tool for examining models with an increasing number of parameters. If we define a χ^2 as $\mathcal{L}(\theta) \propto \exp(-\chi^2/2)$, we can write

$$\mathcal{C}_b = \overline{\chi^2(\theta)} - \chi^2(\tilde{\theta}), \quad (8)$$

where $\bar{\chi}$ denotes the mean taken over the posterior PDF.

The definition of \mathcal{C}_b is chosen such that $\mathcal{C}_b \rightarrow \mathcal{C}_0$ for highly informative data [2].

III. METHODS

We will now describe the functions of the used algorithms for the model selection. Since the implementation of the so called nested sampling is sufficiently dealt with by **Gruppe - Bayesian parameter fitting** our implementation uses the PyMC3-Python Library [1].

A. Monte-Carlo integration

Here we will explain Monte-Carlo sampling, that is *Sequential Monte Carlo* and therein METROPOLIS-HASTINGS. its probably better to put these two sub-sections in separate sections.

B. Savage Dickey Density Ration (SDDR)

IV. EXAMPLES

A. Betabinomial example (coin flip)

Let us now consider as a starting example, the flipping of a two-sided coin, i.e. an experiment where we can measure either heads (H) or tails (T) with 50% probability, respectively. This, while simple, allows us an intuitive approach to Bayesian inference and model selection as well as to the MCMC techniques discussed before. Furthermore is this example easily altered to many real-life problems, such as birth rates, ..., or anything with the option of either success or failure.

1. Analytical approach?

Assume we throw a coin 20 times. We observe 6 H and 14 T. "Is this a fair coin?" might be a question to ask yourself since the bias in outcome is quite large. Naively expecting a fair coin we could assign a *prior* to the probability of heads θ as centred around 0.5, so for example a gaussian with $\mu = 0.5, \sigma = 0.1$.

2. Numerical approach

B. Fitting a polynomial of unknown degree

1. Analytical approach?

2. Numerical approach

V. DISCUSSION

VI. SUMMARY

-
- [1] Pymc3 documentation — pymc3 3.10.0 documentation. URL <https://docs.pymc.io/>.
 - [2] M. Kunz, R. Trotta, and D. R. Parkinson. Measuring the effective complexity of cosmological models. *Phys. Rev. D*, 74:023503, Jul 2006. doi: 10.1103/PhysRevD.74.023503. URL <https://link.aps.org/doi/10.1103/PhysRevD.74.023503>.
 - [3] D. Sivia and J. Skilling. *Data Analysis - A Bayesian tutorial*, volume 2. Oxford University Press, 2006.
 - [4] D. J. Spiegelhalter, N. G. Best, B. P. Carlin, and A. van der Linde. Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 64(4):583–639, 2002. ISSN 13697412, 14679868. URL <http://www.jstor.org/stable/3088806>.
 - [5] R. Trotta. Applications of Bayesian model selection to cosmological parameters. *Monthly Notices of the Royal Astronomical Society*, 378(1):72–82, 05 2007. ISSN 0035-8711. doi:10.1111/j.1365-2966.2007.11738.x. URL <https://doi.org/10.1111/j.1365-2966.2007.11738.x>.
 - [6] R. Trotta. Bayes in the sky: Bayesian inference and model selection in cosmology. *Contemporary Physics*, 49(2):71–104, Mar 2008. ISSN 1366-5812. doi: 10.1080/00107510802066753. URL <http://dx.doi.org/10.1080/00107510802066753>.