

Bayesian modelling

Léo Belzile

Table of contents

Welcome	1
1 Introduction	3
1.1 Probability and frequency	3
1.1.1 Bayesian updating	5
2 Priors	7
References	9

Welcome

This book is a web complement to MATH 80601A *Bayesian modelling*, a graduate course offered at HEC Montréal.

These notes are licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License and were last compiled on Sunday, July 16 2023.

The objective of the course is to provide a hands on introduction to Bayesian data analysis. The course will cover the formulation, evaluation and comparison of Bayesian models through examples and real-data applications.

1 Introduction

1.1 Probability and frequency

In classical (frequentist) parametric statistic, we treat observations Y as realizations of a distribution whose parameters θ are unknown. The *likelihood principle* states that all information about parameters is encoded by the likelihood function, which is optimized numerically or analytically to find the maximum likelihood estimator. This gives a single value for the parameter, and large-sample theory shows that the resulting estimator is asymptotically normal under regularity conditions.

The interpretation of probability in the classical statistic is somewhat counterintuitive and is understood in terms of long run frequency, which is why we call this approach frequentist statistic. Think of a fair die: when we state that values $\{1, \dots, 6\}$ are equiprobable, what we mean is that repeatedly tossing the die should result, in large sample, in each outcome being realized roughly $1/6$ of the time (the symmetry of the object also implies they should equally likely). This interpretation also carries over to confidence intervals. A $(1 - \alpha)$ confidence interval either contains the true parameter value or it doesn't, so the probability level is only the long-run proportion of intervals created by the procedure that should contain the true fixed value, not the probability that a single interval contains the true value. This is counterintuitive to most.

In practice, the true value of the parameter θ vector is unknown to the practitioner, thus uncertain: Bayesians would argue that we should treat the latter as a random quantity rather than a fixed constant to reflect this lack of knowledge. Since different people may have different knowledge about these potential values, the prior knowledge is a form of subjective probabilities, meaning they are individual specific. For example, if you play cards, one person may have recorded the previous cards that were played, whereas other may not. They then assign different probability of certain cards being played.

In Bayesian inference, we consider θ as random variables to reflect our lack of knowledge about potential values taken. Italian scientist Bruno de Finetti, who is famous for the claim "Probability does not exist", stated in the preface of Finetti (1974):

Probabilistic reasoning — always to be understood as subjective — merely stems from our being uncertain about something. It makes no difference whether the

1 Introduction

uncertainty relates to an unforeseeable future, or to an unnoticed past, or to a past doubtfully reported or forgotten: it may even relate to something more or less knowable (by means of a computation, a logical deduction, etc.) but for which we are not willing or able to make the effort; and so on [...]. The only relevant thing is uncertainty — the extent of our knowledge and ignorance. The actual fact of whether or not the events considered are in some sense *determined*, or known by other people, and so on, is of no consequence.

On page 3, de Finetti continues (Finetti 1974)

only subjective probabilities exist — i.e., the degree of belief in the occurrence of an event attributed by a given person at a given instant and with a given set of information.

The likelihood $\mathcal{L}(\boldsymbol{\theta}; \mathbf{y}) \equiv p(\mathbf{y} | \boldsymbol{\theta})$ is the starting point for Bayesian inference. However, we adjoin to it a **prior** distribution $p(\boldsymbol{\theta})$ that reflects the prior knowledge about potential values taken by the p -dimensional parameter vector, before observing the data \mathbf{y} . We thus seek $p(\boldsymbol{\theta} | \mathbf{y})$: the observations are random variables but inference is performed conditional on the observed sample. By Bayes' theorem, the posterior distribution $p(\boldsymbol{\Theta} | \mathbf{Y})$ is

$$p(\boldsymbol{\Theta} | \mathbf{Y}) = \frac{p(\mathbf{Y} | \boldsymbol{\Theta})p(\boldsymbol{\Theta})}{\int p(\mathbf{Y} | \boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}}, \quad (1.1)$$

so the posterior $p(\boldsymbol{\theta} | \mathbf{y})$ is proportional, as a function of $\boldsymbol{\theta}$, to the product of the likelihood and the prior function. The integral in the denominator, termed marginal likelihood and denoted $p(\mathbf{Y}) = E_{\boldsymbol{\theta}}\{p(\mathbf{Y} | \boldsymbol{\theta})\}$, is a normalizing constant that makes the right hand side integrate to unity.

For the posterior to be **proper**, we need the product on the right hand side to be integrable. The denominator of Equation 1.1 is a normalizing constant so that the posterior is a distribution. If $\boldsymbol{\theta}$ is low dimensional, numerical integration such as quadrature methods can be used to compute the latter. To obtain the marginal posterior $p(\theta_j | \mathbf{y}) = \int p(\boldsymbol{\theta} | \mathbf{y})d\boldsymbol{\theta}_{-j}$, additional integration is needed.

When $\boldsymbol{\theta}$ is high-dimensional, the marginal likelihood is untractable. This is one of the main challenges of Bayesian statistics and the popularity and applicability has grown drastically with the development and popularity of numerical algorithms Gelfand and Smith (1990). Markov chain Monte Carlo methods circumvent the calculation of the denominator by drawing approximate samples from the posterior.

1.1.1 Bayesian updating

Subjective probabilities imply that different people with different prior beliefs would arrive at different conclusions. However, as more data are gathered, we can use Bayes theorem to update these prior beliefs and update the posterior. In most instances, the relative weight of the prior relative to the likelihood becomes negligible: if we consider independent data $\mathbf{y}_1, \mathbf{y}_n$ observed sequentially, then

$$\begin{aligned} p(\boldsymbol{\theta} \mid \mathbf{y}_1, \dots, \mathbf{y}_k) &\propto p(\mathbf{y}_k \mid \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathbf{y}_1, \dots, \mathbf{y}_{k-1}) \\ &\propto \prod_{i=1}^k p(\mathbf{y}_i \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) \end{aligned}$$

If data are exchangeable, the order in which observations are collected and the order of the belief updating is irrelevant to the full posterior $p(\boldsymbol{\theta} \mid \mathbf{y}_1, \dots, \mathbf{y}_n)$.

Example 1.1. Back in January 2021, the Quebec government was debating whether or not to distribute antigen rapid test, with strong reluctance from authorities given the paucity of available resources and the poor sensitivity.

A Swiss study analyse the efficiency of rapid antigen tests, comparing them to repeated polymerase chain reaction (PCR) test output, taken as benchmark (Jegerlehner et al. 2021). The results are presented in Table 1.1

Table 1.1: Confusion matrix of Covid test results for PCR tests versus rapid antigen tests, from Jegerlehner et al. (2021).

	PCR +	PCR –
rapid +	92	2
rapid –	49	1319
total	141	1321

Estimated seropositivity at the end of January 2021 according to projections of the Institute for Health Metrics and Evaluation (IHME) of 8.18M out of 38M inhabitants (Mathieu et al. 2020), a prevalence of 21.4%. Assuming the latter holds uniformly over the country, what is the probability of having Covid if I get a negative result to a rapid test?

1 Introduction

Let rapid− (rapid+) denote a negative (positive) rapid test result and C+ (C−) Covid positivity (negativity). Bayes' formula gives

$$\begin{aligned}\Pr(C+ \mid \text{rapid}-) &= \frac{\Pr(\text{rapid}- \mid C+) \Pr(C+)}{\Pr(\text{rapid}- \mid C+) \Pr(C+) + \Pr(\text{rapid}- \mid C-) \Pr(C-)} \\ &= \frac{49/141 \cdot 0.214}{49/141 \cdot 0.214 + 1319/1321 \cdot 0.786} \\ &\approx 0.0866\end{aligned}$$

so there is a small, but non-negligeable probability that the rapid test result is misleading. Jegerlehner et al. (2021) indeed found that the sensitivity was 65.3% among symptomatic individuals, but dropped down to 44% for asymptomatic cases. This may have fueled government experts scepticism.

2 Priors

References

- Finetti, Bruno de. 1974. *Theory of Probability: A Critical Introductory Treatment*. Vol. 1. New York: Wiley.
- Gelfand, Alan E., and Adrian F. M. Smith. 1990. "Sampling-Based Approaches to Calculating Marginal Densities." *Journal of the American Statistical Association* 85 (410): 398–409. <https://doi.org/10.1080/01621459.1990.10476213>.
- Geman, Stuart, and Donald Geman. 1984. "Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images." *IEEE Transactions on Pattern Analysis and Machine Intelligence* PAMI-6 (6): 721–41. <https://doi.org/10.1109/TPAMI.1984.4767596>.
- Jegerlehner, Sabrina, Franziska Suter-Riniker, Philipp Jent, Pascal Bittel, and Michael Nagler. 2021. "Diagnostic Accuracy of a SARS-CoV-2 Rapid Antigen Test in Real-Life Clinical Settings." *International Journal of Infectious Diseases* 109 (August): 118–22. <https://doi.org/10.1016/j.ijid.2021.07.010>.
- Mathieu, Edouard, Hannah Ritchie, Lucas Rod  s-Guirao, Cameron Appel, Charlie Giattino, Joe Hasell, Bobbie Macdonald, et al. 2020. "Coronavirus Pandemic (COVID-19)." *Our World in Data*.

