



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Bayesian statistics and modelling

Citation for published version:

van de Schoot, R, Depaoli, S, Gelman, A, King, R, Kramer, B, Märtens, K, Tadesse, MG, Vannucci, M, Willemsen, J & Yau, C 2021, 'Bayesian statistics and modelling', *Nature Reviews Methods Primers*, vol. 1, 3. <https://doi.org/10.1038/s43586-020-00003-0>

Digital Object Identifier (DOI):

[10.1038/s43586-020-00003-0](https://doi.org/10.1038/s43586-020-00003-0)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Nature Reviews Methods Primers

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Title: Bayesian Statistics and Modelling

Authors: Rens van de Schoot^{1*}, Sarah Depaoli², Andrew Gelman³, Ruth King⁴, Bianca Kramer⁵, Kaspar Märtens⁶, Mahlet G. Tadesse⁷, Marina Vannucci⁸, Duco Veen¹, Joukje Willemsen¹, Christopher Yau^{9, 10}

Affiliations

¹ Department of Methods and Statistics, Utrecht University, Utrecht, The Netherlands

² Department of Quantitative Psychology, University of California Merced, Merced, CA, USA

³ Department of Statistics, Columbia University, New York, USA

⁴ School of Mathematics, University of Edinburgh, Edinburgh, UK

⁵ Utrecht University Library, Utrecht University, Utrecht, The Netherlands

⁶ Department of Statistics, University of Oxford, Oxford, UK

⁷ Department of Mathematics and Statistics, Georgetown University, Washington DC, USA

⁸ Department of Statistics, Rice University, Houston, TX, USA

⁹ Division of Informatics, Imaging & Data Sciences, University of Manchester, Manchester, UK

¹⁰ The Alan Turing Institute, British Library, 96 Euston Road, London

Corresponding author: Rens van de Schoot: Department of Methods and Statistics, Utrecht University, P.O. Box 80.140, 3508TC, Utrecht, The Netherlands; Tel.: +31 302534468; E-mail address: a.g.i.vandeschoot@uu.nl.

Acknowledgements [AU: do any of the other authors want to add funding information?]

The first author (RvdS) was supported by a grant from the Netherlands organization for scientific research: NWO-VIDI-452-14-006. RK was supported by a Leverhulme research fellowship grant reference RF-2019-299.

Author contributions

Introduction (R.v.d.S., A.G.); Experimentation (R.v.d.S., S.D., J.W.); Results (R.v.d.S., R.K., M.G.T., M.V., D.V., K.M., C.Y.); Applications (S.D., R.K., K.M., M.G.T., M.V., C.Y.); Reproducibility and data deposition (B.K., D.V., S.D., R.v.d.S.); Limitations and optimizations (A.G.); Outlook (K.M., C.Y.); Overview of the Primer (R.v.d.S.).

Competing interests

The authors declare no competing interests.

33

34 **ORCID :**

35

36 RvdS: <https://orcid.org/0000-0001-7736-2091>

37 SD: <https://orcid.org/0000-0002-1277-0462>

38 AG: <https://orcid.org/0000-0002-6975-2601>

39 RK: <https://orcid.org/0000-0002-5174-8727>

40 BK: <https://orcid.org/0000-0002-5965-6560>

41 KM: <https://orcid.org/0000-0002-7631-727X>

42 MGT: <https://orcid.org/0000-0003-2671-1663>

43 MV: <https://orcid.org/0000-0002-7360-5321>

44 DV: <https://orcid.org/0000-0002-8352-7574>

45 JW: <https://orcid.org/0000-0002-7260-0828>

46 CY: <https://orcid.org/0000-0001-7615-8523>

47 **Abstract**

48

49 Bayesian statistics is an approach to data analysis and parameter estimation based on Bayes' Theorem.
50 This Primer describes the stages involved in Bayesian analysis, from specifying the prior and data models,
51 to deriving inference, model checking and refinement. Bayesian analysis has been successfully employed
52 across a variety of research fields, including social sciences, ecology, genetics, medicine, and more. We
53 discuss these applications and propose strategies for reproducibility and reporting standards. Finally, we
54 outline the impact of Bayesian analysis in artificial intelligence, a major goal in the next decade.

55

56

57 [H1] Introduction

58

59 It all started with an essay written by Reverend Thomas Bayes, published by Richard Price¹, on inverse
60 probability: how to determine the probability of a future event solely based on past events? It was Pierre
61 Simon Laplace² who actually published the theorem we now know as Bayes' theorem (Box 1). The typical
62 Bayesian workflow consists of three main steps (Figure 1). (1) The first ingredient has to do with
63 knowledge available about the parameter in a statistical model without the data itself and is captured in
64 the so-called **prior distribution [G]**. (2) The second ingredient is the information about the same
65 parameters in the data; it is the observed evidence expressed in terms of the **likelihood function [G]** of
66 the data given the parameters. Both prior distribution and likelihood function are combined via Bayes'
67 Theorem and are summarized by (3) the so-called **posterior distribution [G]**, which is a compromise of the
68 prior knowledge and the observed evidence. This joint distribution is also called a generative model. The
69 posterior distribution reflects one's updated knowledge, balancing prior knowledge with observed data.

70 Although the idea of inverse probability and Bayes' theorem have been longstanding within mathematics,
71 these tools have only become prominent in applied statistics in the past fifty years³⁻¹⁰. There are many
72 reasons for using Bayesian methods: Sometimes researchers may be "forced into" the use of Bayes'
73 theorem some models, for example mixture or multilevel models, require Bayesian methods to improve
74 convergence issues¹¹, exact quantification of uncertainty, aid in model identification¹², produce more
75 accurate parameter estimates¹³, data augmentation or data fusion. We will describe much more
76 advantages and disadvantages throughout the manuscript.

77 The goal of this primer is to provide an overview of the current and future use of Bayesian statistics across
78 different fields of science and to provide an overview of literature that can be used for further study.
79 Moreover, we use many examples how to actually implement a Bayesian model on real data, with all
80 data and code is available for teaching purposes. We aim at a broader group of quantitative researchers
81 working in science-related areas with at least some knowledge of regression modelling. In order to keep
82 the current paper as general as possible with respect to implementing Bayesian methods, there are
83 several concepts listed in Figure 1 that we will be focusing on, like priors and posteriors, and several that
84 we will not specifically address, see the left part in the Figure. We also only briefly touch upon topics like
85 model averaging, network analyses, utility functions/ loss functions without giving a full introduction and
86 do not discuss topics like nonparametric methods. For the non-Bayesian parts we do not discuss we refer
87 the interested reader to classical textbooks.^{14,15} This Primer discusses the general framework, algorithms,
88 and a Bayesian research cycle with a special focus on prior specifications (Experimentation). We discuss
89 model fitting, a thorough example of variable selection and we provide an example calculation with
90 posterior predictive checking (Results). Then, we describe how Bayesian statistics is being used in different
91 fields of science (Applications), followed by guidelines for data sharing, reproducibility, and reporting
92 standards (Reproducibility and Data Deposition). We conclude with a discussion of avoiding bias with
93 incorrect models (Limitations and optimizations), and provide a look into the future with Bayesian
94 Artificial Intelligence (Outlook).

95

96 [H1]Experimentation

97 There are several main issues included in this section. First, prior distributions are detailed, highlighting
98 different levels of informativeness (informative, weakly informative, and diffuse priors). The selection of
99 priors is often viewed as one of the more important choices that a researcher makes when implementing
100 Bayesian methods since the priors can have a substantial impact on final model results. This is followed
101 by a description of the prior predictive checking process, which can be used to assess whether the prior
102 settings being implemented are viable. This section concludes with a description of how to determine the
103 likelihood, which is combined with the prior to form the posterior. Given the important roles that the prior
104 and the likelihood have in determining the posterior, it is imperative that prior and model selection be
105 conducted with care.

106 **H3: An Empirical Example - Predicting PhD delays**

107 To illustrate many aspects of Bayesian statistics we provide an example based on real-life data. Note that
108 we simplified the statistical model and the results are only meant for instructional purposes. Instructions
109 for running the code is available for different software including additional data exploration steps¹⁶.
110 Consider an empirical example of a study predicting PhD delays¹⁷ in which the researchers asked 333 PhD
111 recipients in The Netherlands how long it had taken them to finish their PhD thesis. Based on this
112 information they computed the amount of delay as defined as the difference between planned and actual
113 project time in months ($M = 9.97$, $min/max = -31/91$, $SD = 14.43$). Suppose we are interested in
114 predicting PhD delay (y) using a simple regression model, $y = \beta_{age} + \beta_{age^2} + \epsilon$, with age (in years) as
115 a predictor, denoted by β_{age} , and we expect this relation to be quadratic, denoted by β_{age^2} . Also, the
116 model contains an intercept, $\beta_{intercept}$ and we assume the residuals, ϵ , are normally distributed with an
117 unknown variance, σ_{ϵ}^2 . We will refer to this example throughout the following sections to illustrate key
118 concepts.

119

120 **[H2] Formalizing Prior Distributions**

121 Prior distributions— play a defining role in Bayesian statistics. Prior distributions, or priors, can come in
122 many different distributional forms such as a normal, uniform, Poisson distribution, among others, see
123 also the section Variable Selection for some examples of, so-called, Shrinkage priors. They can also
124 represent different levels of informativeness; the information reflected in a prior distribution can be
125 anywhere along the continuum of complete uncertainty to relative certainty. Although it is important to
126 remember that priors can fall along this continuum, there are three main classifications of priors that are
127 often used in the literature to capture the degree of (un)certainty surrounding the population parameter
128 value: (1) informative, (2) weakly informative, and (3) diffuse. These classifications can be made based on
129 the researcher’s personal judgment. For example, a normal prior with a variance of 1000 may be
130 considered diffuse in one setting and informative in another—it depends on the values of the parameter,
131 as well as parameterization or scaling for the parameter.

132

133 Figure 2 illustrates the relationship between the likelihood, prior, and posterior for different prior settings
134 for β_{age} . In this figure, the first column represents the prior distribution, which is normally distributed for
135 the sake of this example. Notice that there are five different rows of priors, representing different prior
136 settings (some varying in the level of informativeness). The second column represents the likelihood. The
137 prior and the likelihood form together to create the posterior according to Bayes’ rule. The third column

138 illustrates the prior, likelihood, and the resulting posterior, which is derived for illustrative purposes in the
139 current section. In the next section Results we demonstrate how to obtain the posterior.

140 The individual parameters that control the amount of (un)certainty in the priors are called
141 **hyperparameters** [G]. Take the normal distribution as an example. This distribution is defined by a mean
142 and a variance which are the hyperparameters for the normal prior, and we can write this distribution as:
143 $N(\mu_0, \sigma_0^2)$, where the hyperparameters represent the mean (μ_0) and variance (σ_0^2) for the prior,
144 respectively. If the variance is relatively large, then it represents more uncertainty surrounding the mean,
145 vice versa. For example, Figure 2 illustrates five prior settings in the first column with different values for
146 μ_0 and σ_0^2 . The diffuse and weakly informative priors (first three rows) show more spread, that is, a
147 larger variance, compared to the informative priors (last two rows). The mean hyperparameter can be
148 seen as the peak in the distribution.

149

150 An **informative prior** [G] is one that reflects a high degree of certainty surrounding the population
151 parameter. Specifically, the hyperparameters for these priors are specified to express particular
152 information reflecting a greater degree of certainty about the model parameters being estimated. In the
153 case of a normal probability distribution, this would indicate that the prior would have a very small, or
154 narrowed, variance. A researcher may want to use an informative prior when existing information
155 suggests restrictions on the viable range of a particular parameter, or a relationship between parameters,
156 like a positive but imperfect (population) correlation between susceptibility to various medical
157 problem^{18,19}. The information embedded in the informative prior can come from a variety of places, which
158 is referred to as prior elicitation. Strategies for prior elicitation can be to ask an expert or a panel of experts
159 to provide an estimate for the hyperparameters based on knowledge of the field²⁰⁻²³, use the results of a
160 previous publication or meta-analysis^{24,25}, or a combination thereof²⁶. Consider the prior
161 $\beta_{age} \sim N(2.5, 5)$, which was derived from a ShinyApp containing a visualization of how the different
162 priors interact²⁷.

163

164 Finally, another method that can be used for prior elicitation involves implementing data-based priors,
165 which are derived based on a variety of methods including maximum likelihood²⁸⁻³¹ or sample statistics³²⁻
166 ³⁴. Although data-based priors are relatively common, we do not recommend use of so-called “double-
167 dipping” procedures, where estimation occurs based on the sample data and then results are used to
168 derive priors implemented (with the same sample data) for final model estimation. We refer the reader
169 elsewhere³² for more details on this topic. Instead, a hierarchical modelling strategy can be implemented,
170 where priors can depend on hyperparameter values that are data-driven, for example sample statistics
171 pulled from the data, thus avoiding the direct problems linked to “double-dipping.” In some cases, an
172 informative prior can produce a posterior that is not reflective of the population model parameter. There
173 are circumstances when informative priors are needed, but it is also important to assess the impact these
174 priors have on the posterior through a sensitivity analysis as discussed below.

175

176 A **weakly informative prior** [G] is typically not too diffuse, and it is not too restrictive either. In the case of
177 a normal prior, a weakly informative prior would have a variance hyperparameter that exhibits wider

178 variance compared to an informative prior. Such priors will have a small impact on the posterior,
179 depending on the scale of the variables, and the posterior results are still data driven.

180 Some researchers find this to be a nice middle ground regarding the informativeness of the prior. A
181 researcher may want to use a weakly informative prior when some information is assumed about a
182 parameter, but there is still a desired degree of uncertainty. For example, a weakly informative normal
183 prior for the regression coefficient could allow 95% of the prior density mass to fall within values between
184 -10 and 10 or between 0 and 10, see the two different examples in Figure 2, respectively. Essentially,
185 weakly-informative priors do not supply any strict information, but yet are still strong enough to avoid
186 inappropriate inferences that can be produced from a diffuse prior^{35,36}. For this purpose a plausible
187 parameter space should be specified capturing a range of plausible parameter values that is considered
188 to be a reasonable range, thereby excluding improbable values and attaining only a limited density mass
189 to implausible values. For example, if a regression coefficient is known to be near zero, then a weakly
190 informative prior can be specified to reduce the plausible range between, for example, ± 5 . This prior
191 would reduce the probability of observing out-of-bound values (e.g., a regression coefficient of 100)
192 without being too informative.

193 Finally, *diffuse priors* [G] reflect a great deal of uncertainty about the model parameter. This form of priors
194 represents a decision to not include knowledge about the value of the parameter being estimated. Such
195 a prior would be represented by a distribution with a relatively flat density (Figure 2). A researcher may
196 want to use a diffuse prior when there is a complete lack of certainty surrounding the parameter. In this
197 case, the data will largely determine the posterior. Sometimes researchers will use the term “non-
198 informative prior” as a synonym to “diffuse”³⁷. However, we refrain from using this term because we
199 argue that even a completely flat prior, for example, a so-called Jeffreys prior³⁸, is still providing
200 information about the degree of uncertainty³⁹. Therefore, no prior is really non-informative. Diffuse priors
201 can be useful for expressing a complete lack of certainty surrounding parameters, but they can also have
202 unintended consequences on the posterior⁴⁰. For example, diffuse priors can have an adverse impact on
203 parameter estimates via the posterior when sample sizes are small, especially under complex modelling
204 situations involving meta-analytic models⁴¹, logistic regression models³⁹, or mixture models¹³. In addition,
205 improper priors are sometimes used with the intention of using them as diffuse priors. Although improper
206 priors are common, and they can be implemented with relative ease within a variety of Bayesian
207 programs, it is important to note that improper priors can lead to improper posteriors. We mention this
208 caveat here because obtaining an improper posterior can impact the degree to which results can be
209 substantively interpreted. Overall, we note that a diffuse prior can be used as a placeholder, in the same
210 way that we might start with a simple statistical model with the intent to improve it as necessary. It may
211 be that future analyses (e.g., with subsequent data) are conducted with more informative priors.

212

213 Overall, there is no right or wrong prior setting. Many times, diffuse priors can produce results that are
214 aligned with the likelihood, whereas sometimes inaccurate (e.g., biased) results can be obtained with
215 relatively flat priors¹³. Likewise, and as described above in the context of informative priors, an informative
216 prior that is not centered in the same place as the likelihood can pull the posterior away from the
217 likelihood. Because there can be an unintended impact of the priors - despite the level of informativeness
218 - it is always important to conduct a prior sensitivity analysis in order to fully understand the influence
219 that the prior settings have on posterior estimates. Especially when sample size is small, Bayesian

220 estimation with mildly informative priors is often used^{9,42,43}, but the prior specification might have a huge
221 effect on the posterior results.

222 In addition, it is important to note that when priors do not conform with the likelihood, it is not necessarily
223 evidence that there is an issue with the prior. It may be that the likelihood is at fault due to a misspecified
224 model or biased data. In turn, the difference between the prior and the likelihood may be reflective of
225 variation that is not captured by the prior or likelihood alone. These issues can be identified through a
226 sensitivity analysis of the likelihood - for example, by modifying the model - in order to assess how the
227 priors and the likelihood align.

228 Although it is important to distinguish between these different types of priors, there is an overarching
229 issue that needs addressing. We would like to conclude this section with a final thought about the impact
230 of priors. It is common for critics of Bayesian methods to point toward the subjectivity of priors as a
231 potential downfall of the approach. We argue two distinct points here. First, many elements of the
232 estimation process are subjective, including the model itself or the error assumptions. To place the notion
233 of subjectivity solely on the priors is a misleading distraction from the fact that many other elements in
234 the process are inherently subjective by nature. Second, priors are not necessarily a point of subjectivity.
235 They can be used as tools to allow for data-informed shrinkage, enact regularization, or influence
236 algorithms toward a likely high-density region and improve estimation efficiency. In turn, priors are
237 typically defined through previous beliefs, information, or knowledge. Although beliefs can be
238 characterized as subjective points of view from the researcher, information is typically defined as being
239 outside of the researcher and something that can be rigorously quantified, and knowledge can be defined
240 as objective and consensus-based. Therefore, we urge the reader to consider priors in this broader sense,
241 and not simply as a means of incorporating subjectivity into the estimation process.

242

243 Lastly, the current section on informative, weakly informative, and diffuse priors was written in a general
244 sense in that these terms can be used to help define univariate and multivariate priors. The majority of
245 discussion presented in the current paper surrounds univariate priors placed on individual model
246 parameters. However, these concepts can be extended to the multivariate sense, where priors are placed
247 on, for example, an entire covariance matrix rather than a single element from a matrix. For more
248 information on multivariate priors, see^{44,45}.

249

250 **[H2]Prior Predictive Checking**

251 Because the inference based on a Bayesian analysis is subject to the “correctness” of the prior, it is of
252 importance to carefully check whether the specified model can be considered to be generating the actual
253 data^{46,47}. Note that priors are based on background knowledge and cannot be inherently wrong if the prior
254 elicitation procedure is valid. There is an extensive history of expert elicitation across many different
255 disciplines. MATCH⁴⁸ is a generic elicitation tool, but many elicitation problems require custom elicitation
256 procedures and tools, see for instance⁴⁹⁻⁵³ as examples of elicitation procedures designed for specific
257 models. For an abundance of elicitation examples and methods, see the data base of over 67,000 elicited
258 judgements⁵⁴, or the following collections^{20,55,56}. However, even in case of a valid prior elicitation
259 procedure, it is extremely important to understand the exact specification of the priors. This holds

260 especially for smaller sample sizes *in relation to* the complexity of the model, for numerous examples⁹ In
261 the case of smaller sample sizes, priors will exhibit a strong influence on the posteriors. The step of prior
262 prediction is an exercise to improve the understanding of the priors specified and not a method for
263 changing the original prior, unless the prior explicitly generates data that are incorrect.

264 Box⁵⁷ suggested deriving a **prior predictive distribution [G]** from the specified prior. The prior predictive
265 distribution is a distribution of all possible samples that could occur if the model is true. In theory, a
266 “correct” prior provides a prior predictive distribution similar to the true data generating distribution⁴⁶.
267 The prior predictive checking approach compares the observed data to the prior predictive distribution,
268 and checks their compatibility⁴⁷. The compatibility can be summarized by a p-value, describing how far
269 out in the tails of the reference prior predictive distribution the observed data lie⁵⁸. When the **prior**
270 **predictive-value [G]** is “small”, say 0.05, it would indicate that the observed data is unlikely to be
271 generated by the model, and thus call it into question⁴⁷. Evans and Moshonov⁵⁹ suggested restricting the
272 approach of Box to minimal sufficient statistics, i.e. statistics that are as efficient as possible in relaying
273 information about the value of a certain parameter from a sample⁶⁰.

274 Young and Pettit⁶¹ argue that measures being based on a tail area, such as the approaches of Box and
275 Evans and Moshonov, do not produce the required behaviour; favouring the more precise prior if two
276 priors are both specified at the correct value. They propose to use a **Bayes factor [G]**⁶² to compare two
277 priors, see also Box 3. All aforementioned methods leave the determination of the existence of prior-data
278 conflict up to debate depending on an arbitrary cut-off value. The data agreement criterion⁶³ tries to
279 resolve this issue by introducing a clear classification of prior-data conflict, removing the subjective
280 element of the decision⁶⁴. This is done at the expense of selecting an arbitrary divergence based criterion.
281 An alternative has been developed⁶⁵ which computes whether the distance is surprising in relation to the
282 expert’s prior predictive distribution, see for a comparison of both criterion Lek et al⁶⁶

284 **H3: An Empirical Example - Predicting PhD delays - continued**

285 Prior predictive checks can help prevent mistakes from being made. For instance, various
286 software packages can notate the same distribution differently. The normal distribution can be specified
287 by the hyperparameters mean and variance, mean and standard deviation or mean and precision. The
288 precision is the inverse of the variance. For the last prior shown in Figure 2, we have mis-specified the
289 prior variance, that is instead of using a variance of 5 we mis-specified the variance and used the inverse
290 of the variance (i.e., a precision) instead ($1/5=0.2$), $\beta_{age} \sim N(2.5, 0.2)$. If a user is not aware of such
291 differences, a prior which was intended to be weakly informative can easily turn into an informative prior
292 distribution. The prior predictive checks in Figure 3 help to avoid misspecifications like this. Panel A
293 displays a scenario in which precision was mistakenly used instead of variance for β_{age} , and displays an
294 unexpected pattern for the prior predictive distribution. Panel B shows reasonable results for the prior
295 predictive distribution for the correct implementation of the hyperparameters. Additionally, in panel C,
296 the **kernel density estimate (i.e., the estimate of the probability density function) [G]**⁶⁷ of the observed
297 data is displayed (y - in dark blue) which fall neatly in the distribution of the simulated data (y_{rep} - in light
298 blue). The kernel densities for the prior predictive data are based on combinations of possible values of
299 the different priors. Because of the combinations of uncertainty in the priors, the prior predictive kernel
300 density estimates can be quite different from the observed data. The main focus for Panel C is to check

301 that the prior predictive kernel distributions are not order-of-magnitudes different from the observed
302 data.

303 The scripts to reproduce the results are available at the Open Science Framework: [https://osf.io/ja859/](https://osf.io/ja859/DOI.10.17605/OSF.IO/JA859)
304 [DOI.10.17605/OSF.IO/JA859](https://osf.io/ja859/DOI.10.17605/OSF.IO/JA859). Note that in this example the prior predictive distribution and the data are
305 compared on the test statistics mean and standard deviation(sd). It is common to desire descent prior
306 predictive performance on these simple statistics at least. The test statistic can however be chosen to
307 reflect important characteristics of the data, e.g. skewness. It is common to desire descent prior predictive
308 performance on these simple statistics at least. The test statistic can however be chosen to reflect
309 important characteristics of the data, e.g. skewness.

310

311 **[H2] Determining the Likelihood Function**

312 The likelihood, which is used in both Bayesian and frequentist inference⁶⁸, is the conditional probability
313 distribution $p(y|\theta)$ of the data y given parameters θ . In Bayesian inference, the likelihood $p(\mathbf{y}|\boldsymbol{\theta})$ comes
314 into the posterior as a function of $\boldsymbol{\theta}$ for observed data \mathbf{y} . The likelihood function summarises the
315 information of the following elements: a statistical model that stochastically generates all the data, a
316 range of possible values for $\boldsymbol{\theta}$, and the observed data. In a Bayesian model, the likelihood function is part
317 of the generative model, the joint distribution of \mathbf{y} and $\boldsymbol{\theta}$. Because the concept of likelihood is not specific
318 to Bayesian methods, we do not provide a more elaborate introduction of the statistical concept here.
319 Instead, the interested reader is directed to the paper by Etz⁶⁹ for an introduction of how likelihood
320 underlies common frequentist and Bayesian statistical methods and to the work of Pawitan⁷⁰ for a
321 complete mathematical explanation on this topic.

322

323 Much of the discussion surrounding Bayesian inference focuses on the choice of priors, and there is a vast
324 literature on potential defaults^{71,72}. The inclusion of prior knowledge in the form of a prior is the most
325 noticeable difference between frequentist and Bayesian methods and a source of controversy. However,
326 as argued by Gelman, Simpson and Betancourt⁷¹, a prior can in general only be interpreted in the context
327 of the likelihood with which it will be paired. The importance of the likelihood often gets left out of the
328 discussion, even though the specified model for the data - instantiated by the likelihood function - is the
329 foundation for the analysis⁷³.

330 In some cases, specifying a likelihood function can be very straightforward, see Box 2 for an example.
331 However, in practice the underlying data-generating model is not always known. Researchers often
332 naively choose a certain distribution out of habit or because they cannot change it (easily) in the software.
333 The choice of the statistical data-generating model is subjective (based on background knowledge) and
334 should therefore be well understood and described in detail. Robustness checks should be performed to
335 verify the influence of the choice of the likelihood function on the posterior results⁷². Although most
336 research in the theory of Bayesian robustness has concerned the sensitivity of the posterior to imprecision
337 solely in the prior, a few contributions have focussed on the problem of robustness with respect to the
338 likelihood, see for instance⁷⁴⁻⁷⁶ and references therein.

339 [H1] Results

340 After specifying the prior and the likelihood, in this section we assume the data has been collected and
341 we describe the posterior parts of Figure 1. That is, we explain how a model can be fitted to data with the
342 goal of obtaining a posterior distribution, how to select variables, and why posterior predictive checking
343 would be needed. In practice, model building is an iterative process. Any Bayesian model (which includes
344 both the prior distribution and the probability model for data given parameters, which serves also as the
345 likelihood function) can be viewed as a placeholder which can later be improved, in response to the
346 availability of new data, lack of fit to existing data, or simply a process of refinement of the model. Box⁵⁷,
347 Rubin⁷⁷, and chapter 6 of Gelman et al.⁷³ discuss the fluidity of Bayesian model building, inference,
348 diagnostics, and model improvement.

349 [H2] Model Fitting

350 Once the general model structure has been formulated to describe the data, and the associated likelihood
351 function derived, the next step is to fit the model to the observed data to estimate the model parameters.
352 Although the statistical models necessarily simplify reality, they aim to capture the main processes driving
353 the data. Models may differ substantially in their complexity, taking into account the different
354 mechanisms acting on the system and sources of stochasticity and variability. Some examples of the types
355 of data and associated models are provided in *Applications*. Fitting the models to the observed data
356 permits the estimation of the model parameters, or functions of these, leading to an improved
357 understanding of the system, and associated underlying factors via relevant interpretable quantities given
358 the data.

359 There are two main paradigms for model fitting and parameter estimation: Bayesian and frequentist.
360 These approaches differ fundamentally. Within the Bayesian framework probabilities are assigned to the
361 model parameters, describing the associated uncertainties; whereas the frequentist framework focuses
362 on the expected long-term outcomes of an experiment. The corresponding implication is that frequentist
363 methods focus on producing a single point estimate for each model parameter, such as the maximum
364 likelihood estimate, (with an associated uncertainty interval: the confidence interval); whereas in
365 Bayesian statistics, the focus is on estimating the entire posterior distribution of the model parameters.
366 This posterior distribution is often summarised, for simplicity, via associated point estimates (such as the
367 posterior mean or median) and an interval estimate in the form of a credible interval (i.e. an interval that
368 contains a given % of the posterior distribution). Direct inference on the posterior distribution is typically
369 not possible as the mathematical equation describing the posterior distribution is typically both high-
370 dimensional (the number of dimensions is equal to the number of parameters) and of a very complex
371 form. In particular, the expression for the posterior distribution is typically only known up to a constant
372 of proportionality, with the denominator expressible as a function of only the data, where this function is
373 not available in closed form but expressible as an analytically intractable integral. We note that this
374 intractability of the posterior distribution was the primary practical reason that Bayesian statistics was
375 discarded by many scientists for the alternative frequentist statistics. However, the seminal paper by
376 Gelfand and Smith⁷⁸ transformed the data analytic world, describing how **Markov chain Monte Carlo**
377 **(MCMC)** [G], a technique for sampling from a probability distribution, can be used to fit models to data
378 within the Bayesian paradigm.⁷⁹

379 MCMC is able to indirectly obtain inference on the posterior distribution via simulation⁷⁹. In particular,
380 MCMC permits a set of sampled parameter values of arbitrary size to be obtained from the posterior
381 distribution of interest, despite the posterior distribution being high dimensional and only known up to a
382 constant of proportionality. These sample values are used to obtain empirical estimates of the posterior
383 distribution of interest, which can be estimated up to the desired accuracy by increasing the number of
384 sampled parameter values, if necessary. We note that due to the high dimensionality of the posterior
385 distribution it is often useful to focus on the marginal posterior distribution of each parameter, defined
386 by marginalising (or integrating) out over the other parameters (i.e. dimensions). Marginal distributions
387 are useful for focusing on individual parameters but by definition do not provide any information on the
388 relationship between the parameters.

389 Whilst MCMC is the most common algorithm used in Bayesian analyses, there are other model-fitting
390 algorithms, see Table 1 for a non-exhaustive overview of MCMC techniques of sampling and
391 approximation techniques. We refer the interested reader for running the PhD-example with different
392 estimators to^{80,81}. In this article for posterior inference, we focus on MCMC which combines two concepts:
393 (i) obtain a set of parameter values from the posterior distribution (using the **Markov chain [G]** , or the
394 first “MC”); and (ii) given sampled parameter values obtain a distributional estimate of the posterior and
395 associated posterior statistics of interest (using **Monte Carlo [G]** , or the second “MC”). We discuss each
396 of these “MC” components in turn, in reverse order.

397 Consider concept (ii) “Monte Carlo”. Suppose we have a set of parameter values from some distribution.
398 Monte Carlo integration permits estimation of this distribution using associated empirical estimates⁸². For
399 example, to estimate distributional summary statistics, such as the mean, variance or symmetric 95%
400 credible interval of a parameter we use the corresponding sample mean, variance and 2.5% and 97.5%
401 quantile parameter values. Similarly, probability statements can be estimated (such as the probability that
402 a parameter is positive/negative; or lies in the range [a,b]) as the proportion of the sampled values that
403 satisfy the given statement; while the posterior marginal density of any given parameter can be obtained
404 via kernel density estimation, which uses a non-parametric approach for estimating the associated density
405 from which sampled values have been drawn⁸³.

406 However, in general, it is not possible to directly and independently sample parameter values from the
407 posterior distribution. This leads to concept (i) the “Markov chain”. The idea is to obtain a sample from
408 the posterior distribution by constructing a Markov chain with some specified first-order transition kernel
409 which defines the distribution of the parameters at iteration $t+1$, given their state at time t , such that the
410 resulting stationary/equilibrium distribution of the Markov chain is equal to this posterior distribution of
411 interest. Thus, if we run the Markov chain long enough so that it has reached its stationary distribution,
412 subsequent realisations of the chain can be regarded as a (dependent) sample from the posterior
413 distribution and used to obtain the corresponding Monte Carlo estimates, see for an example Figure 4A.
414 We emphasise that the sampled parameter values obtained from the Markov chain are auto-correlated,
415 in that the parameter values are dependent on their previous values in the chain, and generated via the
416 first order Markov chain. The Markov chain is defined by the specification of the initial parameter values
417 and **transition kernel [G]**. There are standard approaches for defining the transition kernel so that the
418 corresponding stationary distribution is the correct posterior distribution: such as the Gibbs sampler⁸⁴;
419 Metropolis-Hastings algorithm^{85,86}; and Hamiltonian Monte Carlo⁸⁷.

420 Obtaining posterior inference, by fitting models to observed data can be complicated due to model
421 complexities or data collection processes. For example, for random effect models or in the presence of
422 latent variables, the likelihood may not be available in closed form, but only expressible as an analytically
423 intractable integral (over the random effect terms or latent variables). Alternatively, the likelihood may
424 be available in closed form, for example, for a finite mixture model (or discrete latent variable model), but
425 where the likelihood is multimodal leading to slow mixing within a standard MCMC approach. In such
426 circumstances data augmentation is often used⁸⁸, where we define additional variables, or **auxiliary**
427 **variables [G]**, such that the joint distribution of the data and auxiliary variables (often referred to as the
428 “complete data” likelihood) is now available in closed form and quick to evaluate. For example, for a
429 random effects model, the auxiliary variables correspond to the individual random effect terms (that
430 would previously have been integrated out); for a finite mixture model the auxiliary variables correspond
431 to the mixture component that each observation belongs to. A new joint posterior distribution is then
432 constructed over both the model parameters and auxiliary variables, which is defined to be proportional
433 to the complete data likelihood and associated parameter priors. A standard MCMC algorithm can then
434 be applied that obtains a set of sampled parameter values over both the model parameters and auxiliary
435 variables. Considering the values of only the model parameters of interest within the Markov chain,
436 essentially discarding the auxiliary variables, provides a sample from the original (marginal) posterior
437 distribution of the model parameters given the observed data. Finally we note that the auxiliary variables
438 may themselves be of interest themselves in some cases, and inference on these can be easily obtained
439 via the sampled values.

440 The transition kernel determines the MCMC algorithm, describing how the parameter values (and any
441 other additional auxiliary variables) are updated at each iteration of the Markov chain. In order for the
442 stationary distribution of the Markov chain to be the posterior distribution of interest, the transition
443 kernel is specified such that it satisfies some relatively straightforward rules. The transition kernel is
444 typically defined via some proposal distribution – this name arises as the process of updating the
445 parameter values involves proposing a set of new parameter values from some distribution which, in the
446 general case, are subsequently either accepted or rejected with some probability, where this acceptance
447 probability is a function of the proposal distribution. If the proposed values are accepted the Markov chain
448 moves to this new state; if the values are rejected the Markov chains remains in the same state at the
449 next iteration. Thus, the transition kernel is non-unique with many general choices for the proposal
450 distribution. For example these include the posterior conditional distribution (i.e. the Gibbs sampler;
451 where the acceptance probability in the updating step is equal to unity), Metropolis-Hastings random walk
452 sampler (randomly perturbing the parameter values from their current values), slice sampler and no-U-
453 turn sampler, amongst many others. We do not focus further on the internal mechanics of the MCMC
454 algorithm here as there is a wealth of literature on this topic and also associated computational tools and
455 programs for performing a Bayesian analysis via an MCMC approach (see later in this section).

456 Beyond the necessity of specifying a transition kernel, such that the corresponding stationary distribution
457 is the posterior distribution of interest, the choice of transition kernel defines the performance of the
458 MCMC algorithm in terms of how long the Markov chain needs to be run to obtain reliable inference on
459 the posterior distribution of interest. **Trace plots [G]** of the parameters display the value of the parameters
460 over iteration number. One-dimensional trace plots are most commonly plotted that describe the
461 parameter value at each iteration of the Markov chain (on the y-axis) against iteration number (on the x-
462 axis) and are often a useful exploratory tool (Figure 4A). They provide a visualisation of the chain in terms

463 of how each parameter is exploring the parameter space, often referred to as mixing, which, if poor,
464 require changes to the specified transition kernel; and also for identifying when the Markov chain has
465 reached its stationary distribution. Recall that the Markov chain only converges to the posterior
466 distribution, so that realisations of the chain prior to convergence to its stationary distribution are
467 discarded – this was originally called the burn-in but we prefer the term warm-up.⁸⁹ The most common
468 technique applied to assess convergence is the \hat{R} statistic [G]^{90,91} where multiple independent runs of the
469 MCMC algorithms are run and the within-chain variability and between-chain variability compared (Figure
470 4B). Ideally, each of the multiple chains should be started from different (over-dispersed) starting values
471 (and using different random seeds) to provide greater initial variability across the Markov chains, to make
472 it more likely that non-convergence of the chain to the stationary distribution will be identified, for
473 example, if different sub-modes of the posterior distribution are being explored. \hat{R} is defined to be the
474 ratio of the within- and between-chain variability. Values close to 1 for all parameters and quantities of
475 interest suggest the chain has sufficiently converged to the stationary distribution, so that future
476 realisations can be regarded as a sample from the posterior distribution of interest (Figure 4B). Once the
477 stationary distribution is reached, a further question relates to how many iterations are needed to obtain
478 reliable Monte Carlo estimates (i.e. for sufficiently small Monte Carlo error). To assess this, batching the
479 sampled values is often used which involves sub-dividing the sampled values into non-overlapping
480 “batches” of consecutive iterations and considering the variability of the estimated statistic using the
481 sampled values in each batch⁹².

482 Additionally, to determine if the entire posterior parameter space has been explored the effective sample
483 size (ESS) of the sampled parameter values may be obtained. The ESS roughly expresses how many
484 independent sampled parameter values contain the same information as the autocorrelated MCMC
485 samples- recall that the sampled MCMC values are not independent as they are generated via a first-order
486 Markov chain. Note that ‘sample size’ in the ESS does not refer to sample size of the data but can be seen
487 as the effective length of the MCMC chain instead of the actual length of the chain. Low sampling
488 efficiency is related to high autocorrelation (so that the variability of the parameter values is small over
489 successive iterations) and non-smooth histograms of posteriors, which in turn could point towards
490 potential problems in the model estimation or weak identifiability of the parameters⁵¹. Therefore, when
491 problems occur in obtaining reliable Monte Carlo estimates, a good starting point is to sort all variables
492 based on ESS and investigating the ones with the lowest ESS first. ESS is also useful for diagnosing the
493 sampling efficiency for a large number of variables⁹³.

494 For further discussion of MCMC-related issues, see for example^{73,94,95}. There are now many standard
495 computer packages for implementing Bayesian analyses, and a summary of the main packages are given
496 in Table 2 (see also *Reproducibility and data deposition*), which have subsequently led to the explosion of
497 Bayesian inference across many scientific fields (for examples, see *Applications*). Many of the available
498 packages perform the MCMC algorithm as a black-box (though often with options to change default
499 settings), permitting the analyst to focus on the prior and model specification, and avoid any technical
500 coding. Note there are many additional packages that make it easier to work with the sometimes heavily
501 code-based software, for example the packages BRMS⁹⁶ and Blavaan⁹⁷ in R for making it easy to use Stan.⁹⁸

502

503 H3: Empirical Example - Continued

504 The priors for the PhD delay example were updated with the data and posteriors were computed in Stan⁹⁸.
505 All scripts to reproduce the results are available at the Open Science Framework: DOI
506 10.17605/OSF.IO/JA859_The trace plot of four independent runs of the MCMC algorithms for $\beta_{intercept}$
507 is shown in Figure 4A and displays stability post-burn in. Also, the associated \hat{R} statistic stabilizes after
508 approximately 2,000 iterations, see Figure 4B. The prior and posterior distributions are displayed in Panels
509 4C-E. The posterior parameter estimates can be summarized using, for example, the median of the
510 posterior distributions. Based on these point summaries, it appears the delay peaks at around the age of
511 50, with an explained variance of only of 6%. If we compare our prior and posterior predictive
512 distributions, we are less uncertain and more consistent in what we expect after observing the data. So,
513 accurate predictions of delay for individual cases may not be possible, but we can predict general trends
514 at group level.

515

516 [H2] Variational inference

517 As we have outlined, Bayesian analysis consists of a number of stages including detailed model
518 development, including specifying the prior and data models, the derivation of exact inference
519 approaches based on MCMC, and model checking and refinement (Figure 1). Each is ideally treated
520 independently, separating model construction from its computational implementation. The focus on exact
521 inference techniques has spurred considerable activity in developing Monte Carlo methods which are
522 considered as a gold standard for Bayesian inference. Monte Carlo methods for Bayesian inference adopt
523 a simulation-based strategy for approximating the high-dimensional integrals required to compute
524 posterior quantities. An entirely alternative approach is to produce functional approximations of the
525 posterior using approaches including **Variational Inference [G]** (VI)⁹⁹ or Expectation Propagation¹⁰⁰. In the
526 following, we describe the variational approach, also known as variational methods or variational Bayes,
527 due to its popularity and prevalence of use in machine learning.

528 Variational inference begins by constructing an approximating distribution to approximate the desired,
529 but intractable, posterior distribution. Typically, the approximating distribution is chosen from a family of
530 standard probability distributions, e.g. multivariate Normal, and further assumes that some of the
531 dependencies between the variables in our model are broken. In the case, where the approximating
532 distribution assumes all variables are independent, this gives us the well-known “mean-field
533 approximation”. The approximating distribution will be specified up to a set of “variational parameters”
534 that we optimise to find the best posterior approximation by minimising the Kullback-Leibler divergence
535 to the true posterior. As a consequence, variational inference reposes Bayesian inference problems as
536 optimisation rather than as sampling problems and can be solved using numerical optimisation, i.e.
537 gradient descent. When combined with subsampling-based optimisation techniques such as stochastic
538 gradient descent, variational inference makes approximate Bayesian inference possible for complex large-
539 scale problems.

540 Variational methods therefore transform the inference problem into an optimisation task to identify the
541 parameters of the approximation that minimise its discrepancy with respect to the true posterior. In
542 Bayesian machine learning (see also the Outlook section), coordinate descent approaches for
543 optimisation, have generally given way to stochastic optimisation approaches which provide further
544 scalability benefits in the presence of large data sets¹⁰¹⁻¹⁰³. Stochastic gradient descent uses only subsets

545 of the data (mini-batches) to compute noisy estimates of the gradients whilst still retaining convergence
546 guarantees. However, there is no free lunch, unless the true posterior belongs to the pre-specified family
547 of approximating distributions, it is often difficult to determine how good the variational approximation
548 represents the true posterior.

549

550

551 [H2] Variable Selection

552 Variable selection is the process of identifying the subset of predictors to include in a model. It is a major
553 component of model building along with determining the functional form of the model. Variable selection
554 is especially important in situations where a large number of potential predictors is available. The inclusion
555 of unnecessary variables in a model has several disadvantages, such as increasing the risk of
556 multicollinearity, lacking enough samples to estimate all model parameters, overfitting the current data
557 thus leading to poor predictive performance on new data, and making the model interpretation more
558 difficult. For example, in genomic studies where high-throughput technologies are used to profile
559 thousands of genetic markers, only a few of those predictors are expected to be associated with the
560 phenotype or outcome under investigation. Methods for variable selection can be categorized into those
561 based on hypothesis testing and those that perform penalized parameter estimation. In the Bayesian
562 framework, hypothesis testing approaches use Bayes factors and posterior model probabilities, while
563 penalized parameter estimation approaches specify **shrinkage priors [G]** that induce **sparsity [G]**, as
564 discussed below. Bayes factors are often used when dealing with a small number of potential predictors
565 as they involve fitting all candidate models and choosing between them, whereas penalization methods
566 fit a single model and thus can scale up to larger dimensions.

567 We provide a brief review of these approaches in the context of a classical linear regression model, where
568 the response variable from n independent observations, \mathbf{y} , are related to p potential predictors defined
569 in an $n \times p$ covariate matrix \mathbf{X} via the model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$. The regression coefficients $\boldsymbol{\beta}$ capture the effect
570 of each covariate on the response and $\boldsymbol{\varepsilon}$ are the residuals assumed to follow a Normal distribution with
571 mean 0 and variance σ^2 . Bayes factors⁶² (Box 3) can be used to compare and choose between candidate
572 models, where each candidate model would correspond to a hypothesis. Unlike frequentist hypothesis
573 testing methods, Bayes factors do not require the models to be nested. In the context of variable selection,
574 each candidate model corresponds to a distinct subset of the p potential explanatory variables^{104,105}.
575 These 2^p possible models can be indexed by a binary vector $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)'$, where $\gamma_j = 1$ if covariate
576 X_j is included in the model, that is $\beta_j \neq 0$, and $\gamma_j = 0$ otherwise. Let $M_{\boldsymbol{\gamma}}$ be the model that includes the
577 X_j 's with $\gamma_j = 1$. Prior distributions for each model $p(M_{\boldsymbol{\gamma}})$ and for the parameters under each model
578 $p(\boldsymbol{\beta}_{\boldsymbol{\gamma}}, \sigma^2 | M_{\boldsymbol{\gamma}})$ are specified, and Bayes factors $BF_{\boldsymbol{\gamma}b}$ are evaluated to compare each model $M_{\boldsymbol{\gamma}}$ to one of
579 the models taken as a baseline, M_b . The posterior probability, $p(M_{\boldsymbol{\gamma}} | \mathbf{y})$, for each model can be expressed
580 in terms of the Bayes factors as

$$581 \quad p(M_{\boldsymbol{\gamma}} | \mathbf{y}) = \frac{BF_{\boldsymbol{\gamma}b} p(M_{\boldsymbol{\gamma}})}{\sum_{\boldsymbol{\gamma}'} BF_{\boldsymbol{\gamma}'b} p(M_{\boldsymbol{\gamma}'})}$$

582 where the denominator sums over all considered models $M_{\boldsymbol{\gamma}'}$. The models with largest posterior
583 probabilities would correspond to those with the highest amount of evidence in their favor among the

584 ones under consideration. When p is relatively small (say $p < 20$), all 2^p variable subsets and their posterior
585 probabilities can be evaluated. The model with highest posterior probability (the maximum *a posteriori*
586 model) may be selected as the one most supported by the data. Alternatively, the covariates with high
587 marginal posterior inclusion probabilities, $p(\gamma_j = 1|\mathbf{y}) = \sum_{X_j \in M_\gamma} p(M_\gamma|\mathbf{y})$, may be selected. For
588 moderate to large p , this strategy is not practically feasible as an exhaustive evaluation of all 2^p possible
589 models becomes computationally expensive. Instead, shrinkage priors that induce sparsity, either by
590 setting the regression coefficients of non-relevant covariates to zero or by shrinking them towards zero,
591 are specified and MCMC techniques are used to sample from the posterior distribution..

592 Various shrinkage priors have been proposed over the years. A widely used **shrinkage prior [G]** is the **spike-**
593 **and-slab prior [G]**, which uses the latent binary indicator vector $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p) \in \{0,1\}^p$ to induce a
594 mixture of two distributions on β_j , one peaked around zero (spike) to identify the zero elements and the
595 other a flat distribution (slab) to capture the non-zero coefficients^{106,107}. The discrete spike-and-slab
596 formulation¹⁰⁶ uses a mixture of a point mass at zero and a flat prior (see Figure 5A), while the continuous
597 spike-and-slab prior¹⁰⁷ uses a mixture of two normal distributions (see Figure 5B). Another widely used
598 formulation puts the spike-and-slab prior on the variance of the regression coefficients¹⁰⁸. After specifying
599 prior distributions for the other model parameters, MCMC algorithms are used to explore the large model
600 space and yield a chain of visited models. Variable selection is then achieved through the marginal
601 posterior inclusion probabilities, $P(\gamma_j = 1|\mathbf{y})$. Integrating out the parameters $\boldsymbol{\beta}$ and σ^2 can accelerate
602 the MCMC implementation, while speeding up its convergence and mixing. Various computational
603 methods have also been proposed to rapidly identify promising high posterior probability models, by
604 combining variable selection methods with modern Monte Carlo sampling techniques^{109,110}(see also Table
605 1).

606 Another class of regularization priors that have received a lot of attention in recent years are continuous
607 shrinkage priors¹¹¹⁻¹¹³. These are unimodal distributions on β_j that promote shrinkage of small regression
608 coefficients towards zero, similarly to frequentist penalized regression methods that accomplish
609 regularization by maximizing the log-likelihood function subject to a penalty¹¹⁴. The least absolute
610 shrinkage and selection operator (lasso)¹¹⁴, for instance, uses the penalty function $\lambda \sum_{j=1}^p |\beta_j|$ with λ
611 controlling the level of sparsity. The lasso estimate of β_j can be interpreted as a Bayesian posterior mode
612 estimate using independent Laplace priors for the regression coefficients. Motivated by this connection,
613 the Bayesian lasso¹¹¹ specifies conditional Laplace priors on $\beta_j|\sigma^2$. It should be noted that Bayesian
614 penalization methods do not shrink regression coefficients to be exactly zero, as the lasso penalization
615 does. Instead, the variable selection is carried out using credible intervals for β_j or by defining a selection
616 criterion on the posterior samples. Many continuous shrinkage priors can be parametrized as a scale
617 mixture of normal distributions, which facilitates the MCMC implementation. For example, the Laplace
618 prior in the Bayesian lasso can be obtained as a scale mixture of normals with an exponential mixing
619 density. The exponential mixing distribution has a single hyperparameter, which limits its flexibility in
620 differentially shrinking small and large effects (see Figure 5C). This limitation can be overcome by using a
621 class of shrinkage priors that introduce two shrinkage parameters, which respectively control the global
622 sparsity and the amount of shrinkage for each regression coefficient. The resulting marginalized priors for
623 β_j are characterized by a tight peak around zero that shrinks small coefficients to zero, and heavy tails
624 that prevent excessive shrinkage of large coefficients. These priors are known as global-local shrinkage
625 priors¹¹³. The **horseshoe prior [G]**, for example, achieves this by specifying a normal distribution for the

626 regression coefficient, β_j , conditional on its scale parameters, which in turn, follow half-Cauchy
627 distributions¹¹²(see Figure 5D). A comprehensive review and thorough comparison of the characteristics
628 and performance of different shrinkage priors can be found in ¹¹⁵. Bayesian variable selection methods
629 have been extended to a wide variety of models. Extensions to multivariate regression models include
630 spike-and-slab priors that select variables as relevant to either all or none of the responses¹¹⁶, as well as
631 multivariate constructions that allow each covariate to be relevant for subsets and/or individual response
632 variables¹¹⁷. Other extensions include generalized linear models, random effects and time-varying
633 coefficient models^{118,119}, mixture models for unsupervised clustering¹²⁰, and estimation of single and
634 multiple Gaussian graphical models^{121,122}. The forthcoming Handbook of Bayesian Variable Selection¹²³
635 presents a comprehensive review and highlights recent developments.

636

637 **[H3] Examples of Recent Applications of Bayesian Variable Selection in** 638 **Biomedical studies**

639 The variable selection priors for linear models described in the Results section have found important
640 applications in biomedical studies. We briefly discuss some examples of recent applications of Bayesian
641 variable selection methods.

642 The advent of high-throughput technologies has made it possible to measure thousands of genetic
643 markers on individual samples. Linear models are routinely used to relate large sets of biomarkers to
644 disease-related outcomes, and variable selection methods are employed to identify the significant
645 predictors. In Bayesian approaches, additional knowledge about correlation structure among the variables
646 can be easily incorporated into the analysis. For example, in models with gene expression data, spike-and-
647 slab variable selection priors incorporating knowledge on gene-to-gene interaction networks have been
648 employed to aid the identification of predictive genes¹²⁴, as well as the identification of both relevant
649 pathways and subsets of genes¹²⁵. Other successful applications of Bayesian variable selection priors have
650 been in genome-wide association studies (GWAS), where hundreds of thousands of single nucleotide
651 polymorphisms (SNPs) are measured in thousands or tens of thousands of individuals, with the goal of
652 identifying genetic variants that are associated with a single phenotype or a group of correlated
653 traits.^{126,127}

654 Air pollution is a major environmental risk factor for morbidity and mortality. Small particles produced by
655 traffic and industrial pollution can enter the respiratory tract and have adverse health effects. Particulate
656 matter exposure and their health effects exhibit both spatial and temporal variability. For a treatment of
657 Bayesian hierarchical models for spatial data we refer readers to¹²⁸. Spatially varying coefficients models
658 with spike-and-slab priors inducing spatial correlation have been proposed to identify pollutants
659 associated to adverse health outcomes over a whole region, as well as in different subregions¹²⁹. Over the
660 past couple of decades, a number of -omic studies have been conducted to investigate the effects of
661 environmental exposures on genomic markers and gain a better understanding of the mechanisms
662 underlying lung injury from exposure to air pollutants. Multivariate response models with structured
663 spike-and-slab priors that leverage the dependence across markers have been proposed to identify and
664 estimate the joint effect of pollutants on DNA methylation outcomes¹¹⁷.

665 In neuroscience, neuroimaging studies often employ functional magnetic resonance imaging (fMRI), a
666 non-invasive technique that provides an indirect measure of neuronal activity by detecting blood flow

667 changes. These studies produce massive collections of time series data, arising from spatially distinct
668 locations of the brain, on one or multiple subjects. In a typical task-based experiment, the whole brain is
669 scanned at multiple times while the subject performs a series of tasks. The objective of the analysis is to
670 detect those brain regions that get activated by the external stimulus. Bayesian approaches to general
671 linear models that employ spatial priors have played an important role in the analysis of such data, as they
672 allow a flexible modelling of the correlation structure of the data¹³⁰. Spike-and-slab variable selection
673 priors that incorporate structural information on the brain have been investigated within a wide class of
674 spatio-temporal hierarchical models for the detection of the activation patterns^{131,132}. Other applications
675 of Bayesian variable selection priors in fMRI analysis have been in brain connectivity studies. Here, fMRI
676 data are measured, on subjects typically at rest, with the aim of inferring how brain regions interact with
677 each other and how information is transmitted between them. Among other approaches, multivariate
678 vector autoregressive linear models have been investigated as a way to infer effective (i.e., directed)
679 connectivity. Continuous shrinkage priors as well as structured spike-and-slab prior constructions have
680 been employed for the selection of the active connections^{133,134}. Bayesian variable selection methods have
681 been successfully applied to a number of other biomedical areas, involving longitudinal data, functional
682 data, survival outcomes and case-control studies, to mention a few.

683

684

685 [H2] Posterior Predictive Checking

686 Once a posterior distribution for a particular model is obtained, it can be used to simulate new data
687 conditional on this distribution. Those simulations can be used for, at least, three purposes: First, to check
688 if the simulated data from the model resemble the observed data. To this end, one could compare kernel
689 density estimate of the observed data to density estimates for the simulated data⁶⁷. Second, a more
690 formal posterior predictive checking approach can be taken to evaluate if the model can be considered a
691 good fit with the data generating mechanism^{67,77,135-137}. Any parameter-dependent statistic or discrepancy
692 can be used for the posterior predictive check¹³⁶. This is similar to how prior predictive checks can be used
693 but much more stringent in the comparison⁶⁷. Because posterior distributions are usually more
694 concentrated on the parameter space compared to prior distributions, the tails of the predictive
695 distributions are more concentrated and tail-area probabilities for any observed statistic or discrepancy
696 are hence more sensitive. The sensitivity of the posterior predictive checks is useful because if realistic
697 models are used, the expectation is that these are well calibrated in the long-term average⁷⁷, for more
698 details see *Limitations and optimizations*. Third, posterior predictive distributions can be used to
699 extrapolate beyond the observed data to predict what data we would expect for new situations based
700 upon our model, e.g. in time series. The first two uses of posterior predictive checking should be used
701 with care. There is a risk of over adjusting and refining models to much to the details of a specific data set.
702 An example of this third kind of use of posterior predictive distributions can be found in the time series of
703 Figure 6. The analysis highlights how daily webpage views can be decomposed into non-periodic changes,
704 holiday effects, weekly seasonality, and yearly seasonality effects. Based on the posterior distributions for
705 the particular model, posterior predictive distributions were simulated for the observed and future data,
706 naturally becoming more uncertain when they are further ahead due to accumulated uncertainty. It is
707 also important to be aware that in temporal models some challenges in terms of posterior inference that
708 are inherent to the spatial and/or temporal dependencies^{44,138-140}.

709

710 **H3: Empirical Example – Time Series Wikipedia page views**

711 To illustrate the use of posterior predictive distributions suppose that it is of interest to know how many
712 pageviews a webpage has, and what time related factors might be relevant. Consider the Wikipedia page
713 views for the premier league, the highest English professional soccer league, obtained using the
714 ‘wikipediatrend’¹⁴¹ R package. The scripts are available at the Open Science Framework:
715 <https://osf.io/7yrud/> - DOI [10.17605/OSF.IO/7YRUD](https://doi.org/10.17605/OSF.IO/7YRUD). The decomposable time series model¹⁴²
716 implemented in the ‘prophet’¹⁴³ R package, allows the estimation of trends with non-periodic changes
717 (Figure 6A), holiday effects (B), weekly seasonality (C), and yearly seasonality effects (D). Notable effects
718 in this time series are the peaks of interest surrounding the start of the seasons in August, the end of the
719 seasons in May, and the dip on 29-04-2011 – the wedding day of Prince William and Catherine Middleton.
720 Additionally, a decrease in webpage views occur on each Christmas day, and notable increases occur on
721 Boxing day and at the start of the year when traditionally matches are played during the Christmas break.
722 The model is estimated using observed data in the period between January 1st 2010 and January 1st 2018.
723 Based on the posterior distributions for the particular model, posterior predictive distributions can be
724 simulated for the observed and future data. In panels E and F posterior predictive distributions at each
725 time point can be seen. In general, the simulated data from the model resembles the observed data for
726 the observed time frame. The posterior predictive distributions for future time points are more uncertain
727 when they are further ahead due to accumulated uncertainty. Notice that increases and decreases in page
728 views are accurately predicted for future page views, with the exception of increased interest in July 2018
729 which might relate to the final stage of the World cup Soccer at that time.

730

731 **[H1] Applications**

732 Bayesian inference has been used across all fields of science. We describe a few examples here but there
733 are many other areas of application such as philosophy, pharmacology, economics, physics, political
734 science and beyond.

735 **[H2] Social and Behavioural sciences**

736 A recent systematic review examining the use of Bayesian statistics found that the social and behavioural
737 sciences (e.g., psychology, sociology, and political sciences) have experienced an increase in empirical
738 Bayesian work⁴. The number of Bayesian publications has been steadily rising since about 2004, with more
739 notable increases in the last decade. In part, this focus on Bayesian methods has been due to the
740 development of more accessible software, as well as a focus on publishing tutorials aimed at applied social
741 and behavioural scientist researchers. The increase in prevalence of Bayesian methods is also due to the
742 continued use of Bayes’ rule as a theory for developmental processes.

743

744 Specifically, there have been two parallel uses of Bayesian methods within the social and behavioural
745 sciences: theory development and estimation. The field has experienced an increase in use with respect
746 to each of these two perspectives.

747 Bayes' rule has been used as an underlying theory for understanding reasoning, decision-making,
748 cognition, and theories of mind. This implementation has been especially prevalent within developmental
749 psychology and related fields. For example, Bayes' rule was used as a conceptual framework for cognitive
750 development in young children, capturing how children develop an understanding of the world around
751 them¹⁴⁴. Bayesian methodology has also been discussed in terms of enhancing cognitive algorithms used
752 for learning. Specifically, Gigerenzer and Hoffrage¹⁴⁵ discussed the use of frequencies, opposed to
753 probabilities, as a method to improve upon Bayesian reasoning. In another seminal paper, Slovic and
754 Lichtenstein¹⁴⁶ discussed how Bayesian methods can be used for judgement and decision-making
755 processes. Within this area of the social and behavioural sciences, Bayes' rule has been used as an
756 important conceptual tool for developing theories and understanding developmental processes.

757 The second way that Bayes' rule is used within the social and behavioural sciences, and the focus of much
758 of the current paper, is as a tool for estimation.

759 The social and behavioural sciences are a terrific setting for implementing Bayesian inference. The
760 literature is rich with information that can be used to derive prior knowledge. In turn, informative priors
761 are useful in complex modelling situations, which are common in the social sciences, as well as in cases of
762 small sample sizes. Likewise, certain models (e.g., some multidimensional item response theory models)
763 used to explore education outcomes and standardized tests are intractable using frequentist methods and
764 require the use of Bayesian methods.

765 There have been many tutorials aimed at explaining Bayesian methods to empirical researchers in a
766 variety of subsections of the social and behavioural sciences. To highlight the scope of tutorials, a
767 systematic review of Bayesian methods in the field of psychology uncovered 740 eligible regression-based
768 papers using this approach. Of these, 100 papers (13.5%) were tutorials for implementing Bayesian
769 methods, and an additional 225 papers (30.4%) were either technical papers or commentaries on Bayesian
770 statistics. Some examples of tutorials within this field are as follows. Hoijtink et al.¹⁴⁷ discussed the use of
771 Bayes factors for informative hypotheses within cognitive diagnostic assessment. They illustrated how
772 Bayesian evaluation of informative diagnostics hypotheses can be used as an alternative approach to the
773 traditional diagnostic methods. There is added flexibility with the Bayesian approach since informative
774 diagnostic hypotheses can be evaluated using the Bayes factor using only data from the individual person
775 being diagnosed. Lee¹⁴⁸ published an overview of how Bayes' theorem can be used within the field of
776 cognitive psychology. They discuss how Bayesian methods can be used to develop more complete theories
777 of cognitive psychology, account for observed behaviour in terms of different cognitive processes, explain
778 behaviour on a wide range of cognitive tasks, and provide a conceptual unification of different cognitive
779 models. Depaoli et al.¹⁴⁹ showed how Bayesian methods can benefit health-based research being
780 conducted within psychology. Specifically, they highlighted how informative priors via expert knowledge
781 and previous research can be used to better understand the physiological impact of a health-based
782 stressor. In this research scenario, frequentist methods would not have produced viable results because
783 the sample size was relatively small for the model being estimated (data were expensive to collect and
784 analyse and the population was difficult to access for sampling). Finally, Kruschke¹⁵⁰ presented the
785 simplest example using a *t*-test geared toward experimental psychologists, showing how Bayesian
786 methods can benefit the interpretation of any model parameter. This paper highlights the Bayesian way
787 of interpreting results, focusing on the interpretation of the entire posterior rather than a point estimate.

788 Methodologists have been attempting to guide applied researchers toward using Bayesian methods
789 within the social and behavioural sciences. Although the implementation has been slower to catch on
790 (e.g., the systematic review found only 167 regression-based papers (22.6%) were empirical applications
791 using human samples), some subfields are regularly publishing work implementing Bayesian methods.

792 The field has gained many interesting insights to psychological and social behaviour through Bayesian
793 methods, and the substantive areas where this work has been conducted are quite diverse. For example,
794 Bayesian statistics helped to: uncover the role that craving suppression has in smoking cessation¹⁵¹, make
795 population forecasts based on expert opinions¹⁵², examine the role that stress related to infant care has
796 in divorce¹⁵³, examine the impact of the President of the United States' ideology on U.S. Supreme Court
797 rulings¹⁵⁴, and predict behaviours that limit the intake of "free sugars" in one's diet¹⁵⁵.

798

799 These examples all represent different ways in which Bayesian methodology is captured in the literature.
800 It is common to find papers that highlight Bayes' rule as a mechanism to explain theories of development
801 and critical thinking¹⁴⁴, are expository^{149,150}, focus on how Bayesian reasoning can inform theory through
802 use of Bayesian inference¹⁴⁸, and papers using Bayesian modelling to extract findings that would have
803 been difficult using frequentist methods¹⁵¹. Overall, there is broad use of Bayes' rule within the social and
804 behavioural sciences.

805

806 We argue that the increased use of Bayesian methods in the social and behavioural sciences is a great
807 benefit to improving substantive knowledge. However, we also feel that the field needs to continue to
808 develop strict implementation and reporting standards so that results are replicable and transparent, as
809 discussed in the next section. We believe that there are important benefits to implementing Bayesian
810 methods within the social sciences, and we are optimistic that a strong focus on reporting standards can
811 make the methods optimally useful for gaining substantive knowledge.

812

813 **[H2] Ecology**

814 Applying Bayesian analyses to ecological applications has become increasingly widespread due to both
815 philosophical arguments and practical model-fitting advantages. This is combined with readily available
816 software, see Table 2, and numerous publications describing Bayesian ecological applications using a
817 range of software packages (see for example¹⁵⁶⁻¹⁶² amongst many others). The underlying Bayesian
818 philosophy is attractive in many ways within ecology¹⁶³ as it permits: the incorporation of external,
819 independent, prior information within a rigorous framework (such information may be from previous
820 studies on the same/similar species or from using inherent knowledge of the biological processes)^{164,165};
821 the ability to make direct probabilistic statements on parameters of interest (such as survival probabilities,
822 reproductive rates, population sizes and future predictions)¹⁵⁸; the calculation of relative probabilities of
823 competing models (for example, the presence/absence of density dependence or environmental factors
824 in driving the dynamics of the ecosystem) which in turn permit model-averaged estimates incorporating
825 both parameter and model uncertainty. The ability to provide probabilistic statements is particularly
826 useful in relation to wildlife management and conservation. For example, King et al¹⁶⁶ provide probability

827 statements in relation to the level of population decline over a given time period, which in turn provides
828 probabilities associated with species' conservation status.

829

830 A Bayesian approach is also often applied in practice for pragmatic reasons. Many ecological models are
831 complex (for example, they may be spatio-temporal in nature, high-dimensional and/or involving multiple
832 interacting biological processes) leading to computationally expensive likelihoods that are slow to
833 evaluate; while imperfect or limited data collection processes often lead to missing data and associated
834 intractable likelihoods. In such circumstances standard Bayesian model-fitting tools, such as data
835 augmentation, may permit the models to be fitted; whereas in the alternative frequentist framework,
836 additional model simplifications or approximations may be required. The application of Bayesian statistics
837 in ecology is vast and encompasses a range of spatio-temporal scales from an individual organism level to
838 ecosystem level, from understanding the population dynamics of the given system¹⁶⁷, modelling spatial
839 point pattern data¹⁶⁸, to population genetics, to estimating abundance¹⁶⁹ or assessing conservation
840 management¹⁷⁰.

841 Ecological data collection processes are generally from observational studies, where a sample is observed
842 from the population of interest using some given data survey protocol. In general, the survey should be
843 carefully designed, taking into account the ecological question(s) of interest and so that it minimises the
844 complexity of the model required to fit to the data to be able to answer the given question with a high
845 degree of accuracy. Nevertheless, due to data collection problems (which may, for example, be as a result
846 of equipment failure or due to poor weather conditions), or inherent data collection problems (for
847 example it is not possible to record any individual level information, such as breeding status, if an
848 individual is unobserved), associated model-fitting challenges may arise. Such challenges may include (but
849 are far from limited to) irregularly spaced observations in time (possibly due to equipment failure or
850 motion sensor detections), measurement error (for example, in relation to population counts or
851 disease/breeding status of individuals made from visual observations), missing information (such as
852 individual covariate information or global environmental factors) and multi-temporal and/or spatial scales
853 where different aspects of data are recorded at different temporal scales (for example, hourly GPS
854 location data of individuals; daily environmental data collected at fixed locations; monthly aerial/satellite
855 photographs and annual censuses). The data complexities that arise, combined with associated modelling
856 choices, may lead to a range of model-fitting challenges which can often be more easily addressed within
857 the Bayesian paradigm.

858 For a given ecological study, separating out the individual processes acting on the ecosystem is an
859 attractive mechanism for simplifying the model specification process.¹⁶⁷ For example, state-space models
860 provide a general and flexible modelling framework that describe two distinct types of processes: (i) the
861 system process and (ii) the observation process. The system process describes the true underlying state
862 of the system and how this changes over time. These states may be univariate (such as population size)
863 or multivariate (such as location data); and the system process may describe multiple processes acting on
864 the system (such as birth/reproduction/dispersal/death). However, we are typically not able to observe
865 the true states without some associated error: the observation process describes how the observed data
866 relate to the true (unknown) states. These general state-space models span many applications, including
867 for example, animal movement¹⁷¹; population count data¹⁷²; capture-recapture-type data¹⁶⁶; fisheries
868 stock assessment¹⁷³; and biodiversity¹⁷⁴ (for a review and further applications, see for example^{167,175,176}).

869 Bayesian model-fitting tools, such as MCMC with data augmentation¹⁷⁷, sequential Monte Carlo or particle
870 (P)MCMC,¹⁷⁸⁻¹⁸⁰ permit general state-space models to be fitted to the observed data without the need to
871 specify further restrictions on the model specification (such as distributional assumptions) or make
872 additional likelihood approximations.

873 The process of collecting data continues to evolve with advances in technology, for example, use of GPS
874 geo-location tags and associated additional accelerometers; remote sensing; use of drones for localised
875 aerial photographs; unmanned underwater vehicles; motion-sensor camera traps; citizen science etc. The
876 use of these technological devices has led to new forms of data, and in greater quantity, and associated
877 model-fitting challenges, providing a fertile ground for Bayesian analyses.

878 [H2] Genetics

879 Genetics and genomics have been a popular application of Bayesian methods. In genome-wide association
880 studies (GWAS), Bayesian approaches have provided a powerful alternative to frequentist approaches for
881 assessing the evidence of population associations between genetic variants and a phenotype of
882 interest¹⁸¹. These include approaches for incorporating genetic diversity (e.g. admixture¹⁸²), fine-mapping
883 to identify causal genetic variants¹⁸³, imputation of genetic markers not directly measured using reference
884 populations¹⁸⁴ and meta-analysis for combining information across studies. These applications further
885 benefit from the use of marginalisation in order to account for modelling uncertainties when drawing
886 inferences. More recently, large cohort studies such as the UK Biobank (UKBB)¹⁸⁵ have collated
887 heterogeneous datasets (e.g. imaging, lifestyle, routinely collected health data) alongside genetic
888 information that have expanded the methodological requirements for identifying genetic associations
889 with complex (sub)phenotypes. For example, a Bayesian analysis framework TreeWAS¹⁸⁶ has extended
890 genetic association methods to allow for the incorporation of tree-structured disease diagnosis
891 classifications by modelling the correlation structure of genetic effects across observed clinical
892 phenotypes. This approach incorporates prior knowledge of phenotype relationships that can be derived
893 from a diagnosis classification tree (e.g. ICD-10).

894 Beyond genetics, the availability of multiple molecular data types (“multi-omics”) has also attracted
895 Bayesian solutions to the problem of multimodal data integration. Bayesian latent variable models can be
896 used as an unsupervised learning approach to identify latent structures that correspond to known or
897 previously uncharacterised biological processes across different molecular scales. Multi-Omics Factor
898 Analysis (MOFA)¹⁸⁷ uses a Bayesian linear factor model to disentangle sources of heterogeneity that are
899 common across multiple modalities from those specific to individual data modalities.

900 In recent years, high-throughput molecular profiling technologies have advanced to allow the routine -
901 omics analysis of individual cells¹⁸⁸. This has led to a methodological revolution with an explosion of novel
902 approaches to account for the challenges of modelling single cell measurement noise, cell-to-cell
903 heterogeneity, high-dimensionality, large sample sizes (millions of cells) and perturbation effects from,
904 for instance, genome editing¹⁸⁹. Cellular heterogeneity lends itself naturally to Bayesian hierarchical
905 modelling and formal uncertainty propagation and quantification due to the layers of variability induced
906 by tissue-specific activity, heterogenous cellular phenotypes within a given tissue and stochastic
907 molecular expression at the level of the single cell. In BASiCS¹⁹⁰ this approach is used to account for cell-
908 specific normalisation constants, technical variability to decompose total gene expression variability into
909 technical and biological components.

910 Deep neural networks have also been utilised to specify flexible, non-linear conditional dependencies
911 within hierarchical models for single cell -omics. SAVER-X¹⁹¹ couples a Bayesian hierarchical model with
912 a pretrainable deep autoencoder to extract transferable gene–gene relationships across datasets from
913 different laboratories, variable experimental conditions and divergent species to denoise novel target
914 datasets. While in scVI¹⁹², hierarchical modelling is used to aggregate information across similar cells and
915 genes to infer the distributions that underlie observed expression values. Approximate and scalable
916 inference in both applications is enabled through the use of mini-batch **stochastic gradient descent [G]**
917 (the latter within a variational setting) - a standard technique with modern use of deep neural networks
918 - that allow these models to be fitted to hundreds of thousands to millions of cells (see also the outlook
919 section).

920 Bayesian approaches have also been popular for cancer genomics where large-scale cancer genomic
921 datasets¹⁹³ have enabled a data-driven approach to identifying novel molecular changes that drive cancer
922 initiation and progression. Bayesian network models¹⁹⁴ have been developed to identify the interactions
923 between mutated genes and capture mutational patterns (signatures) that highlight key genetic
924 interactions that potentially allow for genomic-based patient stratification for clinical trials and the
925 personalised use of therapeutics.

926 Bayesian methods have been important in answering questions about evolutionary processes in cancer.
927 Several Bayesian approaches for phylogenetic analysis of heterogeneous cancers enable the identification
928 of the distinct subpopulations that can exist with tumours and the ancestral relationships between these
929 through the analysis of single cell and bulk tissue sequencing data¹⁹⁵. These models therefore consider
930 the joint problem of learning a mixture model (number and identity of the subpopulations) and graph
931 inference (phylogenetic tree).

932

933 **[H1] Reproducibility and Data Deposition**

934 Proper reporting on statistics, including sharing of data and scripts, is a crucial element in the verification
935 and reproducibility of research¹⁹⁶. A typical workflow for good research practices across the research
936 workflow that can contribute to reproducibility is displayed in Figure 7. We demonstrate where the
937 Bayesian research cycle (Figure 1) and the *When to Worry, and how to Avoid the Misuse of Bayesian*
938 *Statistics* checklist¹⁴⁹ (Box 4) fit in the wider context of transparency in research. In this section we
939 highlight some important aspects of reproducibility and data /script deposition.

940 Allowing others to assess the statistical methods used, including access to the underlying data if possible,
941 can help in interpreting the results, assess the suitability of the parameters used, and detect and fix errors.
942 Reporting practices are not yet consistent across many fields, nor across journals in individual fields.
943 Within the systematic review on Bayesian statistics in psychology⁴, huge discrepancies within reporting
944 practices and standards were uncovered in the social sciences. For example, of the 167 regression-based
945 Bayesian papers using human samples in Psychology, 31% did not mention the priors that were
946 implemented, 43.1% did not report on chain convergence, and only 40% of those implementing
947 informative priors conducted a sensitivity analysis. We view this as a major impediment to the
948 implementation of Bayesian statistics within the social and behavioural sciences, as well as other fields of
949 research.

950 Specifically, for Bayesian methods there are many dangers in naïvely using priors. That is, the exact
951 influence of the priors is often not well understood, and priors might have a huge, sometimes unwanted,
952 impact on the study results. Therefore, one might want to pre-register the specification of the priors (and
953 likelihood) when possible, e.g. in a confirmatory study when the actual statistical model is known
954 beforehand. Moreover, akin to many elements of frequentist statistics, some Bayesian features can be
955 easily misused. For example, the impact of priors on final model estimates can be easily overlooked. A
956 researcher may estimate a model with certain priors and be unaware that using different priors with the
957 same model and data can result in substantively different results. In both cases, the results could look
958 completely viable, for example, chains appeared to be converged, posteriors appear viable and
959 informative. Without examining the impact of priors through a sensitivity analysis and prior predictive
960 checking, the researcher would not be aware of how sensitive results are to changes in the priors. Consider
961 the prior variance in the PhD delay example for β_{age} which was mis-specified as being a precision instead
962 of a variance.

963 Also, reporting on Bayesian statistics is not consistent with reporting on frequentist statistics, since there
964 are elements included in the Bayesian framework that are fundamentally different from frequentist
965 settings. Therefore, the WAMBS-checklist¹⁴⁹ was developed to promote proper use and reporting of
966 Bayesian methods. We offer an updated version (WAMBS, version 2) here (Box 4).

967 To enable reproducibility and allow others to rerun Bayesian statistics on the same data with, e.g. other
968 priors, model or likelihood functions for sensitivity analyses¹⁹⁷, it is important that the underlying data and
969 code used are properly documented and shared, following the FAIR principles^{198,199}: Findable, Accessible,
970 Interoperable and Reusable. Preferably, data and code are shared in a trusted repository²⁰⁰ rather than as
971 supplemental information in a journal, with their own persistent identifier (such as a doi) and tagged with
972 metadata describing the dataset or codebase. This also allows the dataset and code to be recognized as
973 separate research outputs and allows other to cite them accordingly²⁰¹. Repositories can be general (such
974 as Zenodo), language-specific such as CRAN for R packages, and PyPI for Python code, or domain-
975 specific²⁰¹. As data and code require different license options, metadata, and other attributes, data are
976 generally best stored in dedicated data repositories, which can be general or discipline-specific²⁰². Some
977 journals, like Nature Research' Scientific Data, have their own list of recommended data repositories
978 (<https://www.nature.com/sdata/policies/repositories>). To make depositing data and code easier for
979 researchers, two repositories (Zenodo and Dryad) are exploring collaboration to allow deposition of code
980 and data through one interface, with data stored in Dryad and code in Zenodo
981 (<https://blog.datadryad.org/2020/03/10/dryad-zenodo-our-path-ahead/>). Many scientific journals
982 adhere to TOP guidelines²⁰³ for transparency and openness in research, which specify requirements for
983 code and data sharing.

984 Verification and reproducibility do not only require access to the data, but also to the code used in
985 Bayesian modelling, ideally replicating the original environment the code was run in, with all
986 dependencies documented either in a dependency file accompanying the code or by creating a static
987 container image that provides a virtual environment to run the code in²⁰². Open source software should
988 be used as much as possible, as open sources reduce the monetary and accessibility threshold to
989 replicating scientific results. Moreover, it can be argued that closed source software keeps part of the
990 academic process hidden, including from the researchers who use the software. However, open-source
991 software is only truly accessible with proper documentation (e.g. listing dependencies and configuration

992 instructions in Readme files, commenting code to explain functionality, and including a comprehensive
993 reference manual when releasing packages).

994 [H1] Limitations and Optimizations

995 Bayesian inference is optimal conditional on the assumed model. That is, Bayesian posterior probabilities
996 are calibrated in long-term average, if parameters are drawn from the prior distribution and data are
997 drawn from the data distribution. That is, events with stated probability occur with that frequency in the
998 long term, when averaging over the generative model. In practice, our models are never correct; this is
999 where the limitations come from. There are two ways we would like to overcome these limitations: by
1000 identifying and fixing problems with the model, and by demonstrating that certain inferences are robust
1001 to reasonable departures from the model. There are many examples of model checks, see the sections on
1002 prior and posterior predictive checking, and robustness checks, like sensitivity analyses and checklists like
1003 the WAMBS (see Box 4), in the Bayesian literature.

1004
1005 Even the simplest and most accepted Bayesian inferences can have serious limitations. For example,
1006 suppose an experiment is conducted yielding an unbiased estimate z of a parameter θ which represents
1007 the effect of some treatment. If this estimate z is normally distributed with standard error s , we can write
1008 $z \sim \text{Normal}(\theta, s)$, a normal distribution parameterized by its location and scale parameter. Suppose that
1009 θ has a flat uniform prior distribution, then the posterior distribution is $\theta \sim N(z, s)$. These are all familiar
1010 calculations. Now suppose we observe $z = s$; that is, the estimate of θ is 1 standard error from zero. In
1011 practice, this would be considered statistically indistinguishable from noise, in the sense that such an
1012 estimate could occur by chance, even if the true parameter value were zero. But the Bayesian calculation
1013 gives a posterior probability $\Pr(\theta > 0|z) = 0.84$. Would you really be willing to offer 5-to-1 odds on a
1014 bet that $\theta > 0$, given these data? If not, in what sense can we say this probability is calibrated?

1015 The answer is that the probability is calibrated if you average over the prior. You can't average over a
1016 uniform distribution on an infinite range, so let's consider a very diffuse prior, for example $\theta \sim N(0, 1000)$,
1017 where we are assuming that s is roughly on unit scale. Under this model, when z is observed to equal s ,
1018 the parameter θ will be positive approximately 84% of the time. The reason why the 84% probability
1019 doesn't seem correct is that the uniform, or very diffuse, prior does not generally seem appropriate. In
1020 practice, studies are designed to estimate treatment effects with a reasonable level of precision. True
1021 effects may be one or two standard errors from zero, but they are rarely 5 or 10 or 100 standard errors
1022 away. In this example, Bayesian inference if taken literally would lead to over-certainty: an 84% posterior
1023 probability corresponds to the willingness to bet at 5-to-1 odds. There is a positive way to look at this
1024 story, though: the evident problem with the bet allowed us to recognize that prior information was
1025 available that we had not included in our model. Moreover, a weakly informative prior such as
1026 $\theta \sim \text{Normal}(0, s)$ does not change the posterior by much, as then the posterior becomes normal
1027 $\text{Normal}(0, 5s, 1/\sqrt{2}s)$, so $\Pr(\theta > 0|z) = 0.76$, and the betting odds only change to roughly 4:1.
1028 Ultimately, only a strong prior will make a big difference. Bayesian probabilities are only calibrated when
1029 averaging over the true prior or population distribution of the parameters.

1030 More generally, Bayesian models can be checked by comparing posterior predictive simulations to data¹³⁶
1031 and by estimating out-of-sample predictive error²⁰⁴. There is a benefit to strong prior distributions

1032 that regularize (constrain parameters to reasonable values) to allow the inclusion of more data while
1033 avoiding overfitting. More data can come from various sources, including additional data points,
1034 additional measurements on existing data, and prior information summarizing other data or theories. All
1035 methods, Bayesian and otherwise, require subjective interpretation in order to tell a plausible story, and
1036 all models come from researcher decisions. The point is that any choice of model has implications. For
1037 example, the flat prior is weak in the sense of providing no shrinkage of the estimate, but it is strong in
1038 the sense of leading to an inappropriate level of certainty about the sign of theta.

1039

1040 [H1] Outlook

1041 The widespread adoption of Bayesian Statistics across disciplines is a testament to the power of the
1042 Bayesian paradigm for the construction of powerful and flexible statistical models within a rigorous and
1043 coherent probability framework. Modern Bayesian practitioners have access to a wealth of knowledge
1044 and techniques that allows the creation of bespoke models and computational approaches for particular
1045 problems. While probabilistic programming languages, such as Stan, can take away much of the
1046 implementation details for many applications allowing the focus to remain on the fundamentals of
1047 modelling and design.

1048 Nevertheless, an ongoing challenge for Bayesian Statistics is the ever-growing demands posed by
1049 increasingly complex real-world applications. These are often associated with issues such as large datasets
1050 and uncertainties regarding model specification. All of this occurs within the context of rapid advances in
1051 computing hardware, the emergence of novel software development approaches and the growth of “data
1052 sciences” which has attracted a larger and more heterogeneous scientific audience than ever before.

1053 In particular, in recent years, the revision and popularisation of the term “artificial intelligence” (AI) to
1054 encompass a broad range of ideas including Statistics and Computation has blurred the traditional
1055 boundaries between disciplines. This has been hugely successful in popularising probabilistic modelling
1056 and Bayesian concepts outside of its traditional roots in Statistics but has also seen transformations in the
1057 way Bayesian inference is being carried out and new questions about how Bayesian approaches can
1058 continue to be right at the innovative forefront of AI research.

1059 Driven by the need to support large-scale applications involving datasets of increasing dimensionality and
1060 sample numbers, Bayesians have exploited the growth of new technologies centred around Deep Learning
1061 (DL). This includes deep learning programming frameworks (e.g. TensorFlow,²⁰⁵ PyTorch²⁰⁶) that
1062 greatly simplify the use of and computations with deep neural networks (DNN) that permit the
1063 construction of more expressive, data-driven models that are immediately amenable to inference
1064 techniques using off-the-shelf optimisation algorithms and state-of-the-art hardware (multicores, GPUs,
1065 TPUs). In addition to providing a powerful tool to specify flexible and modular generative models, DNNs
1066 have also been employed to develop new approaches for approximate inference and stimulated a new
1067 paradigm for Bayesian practice that sees the integration (not separation) of statistical modelling and
1068 computation at its core.

1069 An archetypal example is the “Variational Autoencoder” (VAE)²⁰⁷. VAEs have been successfully used in a
1070 variety of applications, including single cell genomics^{191,192}, and they provide a general modelling

1071 framework that has led to a number of extensions including latent factor disentanglement²⁰⁸⁻²¹⁰. The
1072 underlying statistical model is actually a simple Bayesian hierarchical latent variable model. This model
1073 maps high-dimensional observations to low-dimensional latent variables that are assumed to be normally
1074 distributed through functions defined by DNNs. Variational inference (VI) is used to approximate the
1075 posterior distribution over the latent variables. However, in standard VI we would introduce a local
1076 variational parameter for each latent variable, in which case the computational requirements would scale
1077 linearly with the number of data samples. VAEs use a further approximation process known as
1078 *amortization* to replace inference over the many individual variational parameters with a single global set
1079 of parameters that are used to parameterise a DNN (known as a *recognition network*) that outputs the
1080 local variational parameters for each data point.

1081 Remarkably, when the model and inference are combined and interpreted together, the VAE has an
1082 elegant interpretation as an encoding-decoding algorithm: It consists of a probabilistic *encoder* - a DNN
1083 that maps every observation to a distribution in the latent space - and a probabilistic *decoder* - a
1084 complementary DNN that maps each point in the latent space to a distribution in the observation space.
1085 Thus, model specification and inference have become entangled within the VAE, demonstrating the
1086 increasingly blurry boundary between principled Bayesian modelling and algorithmic DL techniques.
1087 Other recent examples include the use of DNNs to construct probabilistic models that define distributions
1088 over possible functions²¹¹⁻²¹³, build complex probability distributions by applying a sequence of invertible
1089 transformations (normalizing flows)^{214,215} and define models for exchangeable sequence data²¹⁶.

1090 The expressive power of DNNs and their utility within model construction and inference algorithms come
1091 with compromises that are fertile ground for further Bayesian research. The trend toward entangling
1092 models and inference has popularised these techniques for large-scale data problems but fundamental
1093 Bayesian concepts remain to be fully incorporated within this paradigm. Marginalisation, model
1094 averaging, decision theoretic approaches rely on accurate posterior characterisation which remains
1095 elusive due to the challenge posed by high-dimensional neural network parameter spaces²¹⁷. While
1096 Bayesian approaches to neural network learning have been around for decades²¹⁸⁻²²¹, further investigation
1097 into prior specifications for modern Bayesian deep learning models which involve complex network
1098 structures is required to understand how priors translate to specific functional properties²²².

1099 Recent debates within the field of artificial intelligence have questioned the requirement for Bayesian
1100 approaches and highlighted potential alternatives. For instance, Deep Ensembles²²³ have been shown to
1101 be alternatives to Bayesian methods for dealing with model uncertainty. However, more recent work has
1102 shown that “Deep Ensembles” can actually be reinterpreted as approximate Bayesian model
1103 averaging²²⁴. Similarly, “Dropout” is a regularization approach popularised for use in the training of deep
1104 neural networks to improve robustness by randomly dropping out nodes during the training of the
1105 network²²⁵. Dropout has been empirically shown to improve generalizability and reduce overfitting.
1106 Bayesian interpretations of dropout have emerged linking it to forms of Bayesian approximation of
1107 probabilistic deep Gaussian processes²²⁶. While the full extent of Bayesian principles have not yet been
1108 generalised to all recent developments in artificial intelligence, it is nonetheless a success that Bayesian
1109 thinking is deeply embedded and crucial to a number of innovations that have arisen. The next decade is
1110 sure to bring a new wave of exciting innovative developments for Bayesian Intelligence.

1111

Tables

Table 1. A non-exhaustive overview of sampling and approximation techniques

Name	Short description
MCMC	Markov chain Monte Carlo
Metropolis-Hastings (MH)	Updating algorithm uses general proposal distribution, with an associated accept/reject step for the proposed parameter value(s). ^{85,86}
Reversible jump (RJ)MCMC	Extension of MH algorithm to permit trans-dimensional moves within parameter space – most often applied in presence of model uncertainty. ^{33,227}
Hamiltonian Monte Carlo	Special case of MH algorithm based on Hamiltonian dynamics. ⁸⁷
No-U-Turn sampler (NUTS)	An extension to Hamiltonian Monte Carlo that optimizes the generation of candidate points. ²²⁸
Gibbs sampler	Special case of MH algorithm where the proposal distribution is the corresponding posterior conditional distribution, with an associated acceptance probability of 1. ⁸⁴
Particle (P)MCMC	Combined sequential Monte Carlo algorithm and MCMC used when the likelihood is analytically intractable ¹⁷⁸
Evolutionary Monte Carlo	MCMC algorithm that incorporates features of genetic algorithms and simulated annealing. ²²⁹
Other	
Sequential Monte Carlo	Algorithm based on multiple importance sampling steps for each observed data point - often used for on-line or real-time processing of data arrivals. ²³⁰
Approximate Bayesian Computation	Approximate approach, typically used when the likelihood function is analytically intractable or very computationally expensive. ²³¹
Integrated nested Laplace approximations (INLA)	Approximate approach developed for the large class of latent Gaussian models, which includes, for example, generalized additive spline models, Gaussian Markov processes and random fields. ²³²
Variational Bayes	Variational Inference describes a technique to approximate posterior distributions via simpler approximating distributions. Optimisation is used to adapt the variational parameters within these approximating distributions to make them as close to the true posterior distribution as possible using the KL-divergence as a measure of discrepancy ⁹⁹ .

1116 Table 2. A non-exhaustive summary of commonly used and open Bayesian software programs.

1117

Software Package	Summary	Type of sampling	System specifications
General-purpose Bayesian inference software			
BUGS ²³³⁻²³⁵ (Bayesian Inference Using Gibbs Sampler) / JAGS ²³⁶ (Just Another Gibbs Sampler)	The original general-purpose Bayesian inference engine, in different incarnations. Uses Gibbs and Metropolis sampling. Windows based software (<i>WinBUGS</i> ²³³), with user-specified model and black-box MCMC algorithm. Developments include an open source version (<i>OpenBUGS</i> ²³⁵) also available on Linux and Mac (using <i>WINE</i>); and parallel algorithm version (<i>MultiBUGS</i> ²³⁷). R packages are available for calling BUGS from R (such as <i>R2WinBUGS</i> ²³⁸ , <i>R2OpenBUGS</i> ²³⁸ and <i>BRugs</i> ²³⁹). <i>JAGS</i> ²³⁶ (Just Another Gibbs Sampler) is an open source variation of BUGS which can run cross-platform and can run from R via <i>rjags</i> ²³⁶ .	MCMC	OpenBUGS = Windows, Linux, Mac (using <i>WINE</i>) MultiBUGS = Windows JAGS = all platforms.
PyMC3 ²⁴⁰	Framework for Bayesian modeling and inference entirely within Python; includes Gibbs sampling and Hamiltonian Monte Carlo		
Stan ⁹⁸	General-purpose Bayesian inference engine using Hamiltonian Monte Carlo; can be run from R, Python, Julia, Matlab, and Stata. Open source software that implements efficient Hamiltonian Monte Carlo (HMC). Versions available for R, Python, MATLAB, Julia and Stata.	MCMC (Hamiltonian Monte Carlo)	All platforms
NIMBLE ²⁴¹	Generalization of the Bugs language in R; includes sequential Monte Carlo as well as MCMC. Open source R package using BUGS/JAGS-model	MCMC and sequential Monte carlo	All platfoms

	language to develop a model; and different algorithms for model fitting including MCMC and sequential Monte Carlo approaches including the ability to write novel algorithms.		
<i>Programming languages that can be used for Bayesian inference</i>			
TensorFlow Probability ^{242, 243}	A Python library for probabilistic modelling built on Tensorflow ²⁰⁵ from Google.	MCMC	Python 3.5 – 3.8 Ubuntu 16.04 or later Windows 7 or later (with C++ redistributable) macOS 10.12.6 (Sierra) or later (no GPU support) Raspbian 9.0 or later
Pyro ²⁴⁴	Probabilistic programming language built on Python and PyTorch ²⁰⁶ .	MCMC	
Julia ²⁴⁵	In addition to Stan, numerous other probabilistic programming libraries are available for the Julia programming language including Turing.jl ²⁴⁶ and Mamba.jl ²⁴⁷ .	MCMC	Windows macOS Linux FreeBSD
<i>Specialized software doing Bayesian inference for particular classes of models</i>			
JASP ²⁴⁸ (Jeffreys's Amazing Statistics Program)	JASP is a user friendly higher-level interface, offering standard analysis procedures in both their classical and Bayesian form. It is open source and relies upon a collection of open-source R packages.		Windows MAC Linux
R-INLA ²³²	Open source R package for implementing INLA. ²⁴⁹ Fast inference in R for a certain set of hierarchical models using nested Laplace approximations.	INLA	All platforms

GPstuff ²⁵⁰	Fast approximate Bayesian inference for Gaussian processes using expectation propagation; runs in Matlab, Octave, and R.		Unix and Windows Matlab
------------------------	--	--	-------------------------

1118

1119

1120 **Figures Headings**

1121

1122 **Figure 1. The Bayesian Research Cycle.**

1123 Typical steps needed for a research cycle using Bayesian statistics. The first part of the Bayesian Research
1124 Cycle, indicated with (A) is identical to any research cycle: starting with reading literature, defining a
1125 problem, specifying the research question and hypothesis^{14,15}. The analytic strategy should be pre-
1126 registered to enhance transparency. The second part of the Bayesian Research Cycle, indicated with (B) is
1127 specifically for a Bayesian workflow. It includes formalizing prior distributions based on background
1128 knowledge and prior elicitation, determining the likelihood function by specifying a data generating model
1129 and including observed data, and obtaining the posterior distribution as a function of both the specified
1130 prior and likelihood function^{135,251}. To probe the consequences of the specified model, it is important to
1131 perform robustness checks along the way and after. All concepts are briefly discussed in the primer with
1132 references for the interested user.

1133

1134 **Figure 2. Illustration of the Key Ingredients of Bayes' Theorem.**

1135 This figure displays how the likelihood and prior work together to form the posterior distribution. Notice
1136 that the likelihood remains constant across all rows. Each row only differs in the prior distribution
1137 specified. Priors are typically deemed to be informative, weakly informative, or diffuse, each defined
1138 through different degrees of (un)certainty—in this case, through the variance (or spread) of the prior. The
1139 posterior distribution is a compromise between the prior and the likelihood.

1140

1141 **Figure 3: Prior Predictive Checks.**

1142 Prior predictive checks for the PhD-delay example, computed via Stan⁹⁸ – the scripts are available at the
1143 Open Science Framework: <https://osf.io/ja859/> - DOI 10.17605/OSF.IO/JA859 (A) displays a scenario in
1144 which precision was mistakenly used instead of variance for β_{age} and displays an unexpected pattern for
1145 the prior predictive distribution. Note, in dark blue the observed mean and SD are presented, in light blue
1146 samples of the prior predictive distribution. (B) shows the prior predictive distribution for the correct
1147 implementation of the hyperparameters. The prior predictive checks for the correct implementation of
1148 the priors seem reasonable given the data. Additionally, in panel C, a kernel density estimate of the

1149 observed data is displayed (y - in dark blue), and kernel density estimates for the simulated data (y_{rep} -
1150 in light blue)⁶⁷.⁸³As can be seen the priors cover the entire plausible parameter space with the observed
1151 data in the center.

1152 **Figure 4. Posterior mean and SD estimation using MCMC**

1153 In panel (A) trace plots (iteration number against parameter value) for the PhD delay data, computed in
1154 Stan⁹⁸ of four independent MCMC algorithms are shown for exploring the same posterior distribution of
1155 $\beta_{intercept}$, with the first part omitted for constructing the posterior distribution (i.e, warm-up phase); In
1156 panel (B) the associated \hat{R} statistic is shown which appears to settle down around the value of 1 after
1157 approximately 2,000 iterations; and (C, D and E) prior and posterior distributions for the in the model, the
1158 intercept (Panel C, $\beta_{intercept}$), the linear effect of age on PhD delay (Panel D, β_{age}), and the quadratic
1159 effect of age on PhD delay (Panel E, β_{age^2}). For each chain, the first 2,000 iterations are discarded as
1160 warm-up. The scripts are available at the Open Science Framework: <https://osf.io/ja859/> - DOI
1161 [10.17605/OSF.IO/JA859](https://doi.org/10.17605/OSF.IO/JA859).

1162

1163 **Figure 5. Examples of shrinkage priors for Bayesian variable selection.**

1164 In Panel A, the discrete spike-and-slab prior for β_j (solid blue line) is specified as a mixture of a point mass
1165 at 0 (spike; dashed black line) and a flat prior (slab; dotted red line). In panel B, the continuous spike-and-
1166 slab prior for β_j (solid blue line) is specified as a mixture of two normal distributions, one peaked around
1167 0 (dashed black line) and the other with a large variance (dotted red line). In panel C, the Bayesian lasso
1168 specifies a conditional Laplace prior, which can be obtained as a scale mixture of normal distributions with
1169 an exponential mixing density. This prior does not offer enough flexibility to allow simultaneously a lot of
1170 mass around zero and heavy tails. In panel D, the horseshoe prior falls in the class of global-local shrinkage
1171 priors, which are characterized by a high concentration around zero to shrink small coefficients and heavy
1172 tails to avoid excessive shrinkage of large coefficients.

1173

1174 **Figure 6. Posterior Predictive Checking**

1175 Wikipedia page views for the premier league as obtained using the ‘wikipediatrend’¹⁴¹ R package and
1176 analyzed with the ‘prophet’¹⁴³ R package. The scripts are available at the Open Science Framework:
1177 <https://osf.io/7yrud/> - DOI [10.17605/OSF.IO/7YRUD](https://doi.org/10.17605/OSF.IO/7YRUD). Panels show posterior means for the following

1178 parameters along with 95% CIs for non-periodic changes (A), holiday effects (B), weekly seasonality (C),
1179 and yearly seasonality effects (D). In panels E and F posterior predictive distributions at each time point
1180 can be seen. The posterior predictive distributions for the time points that fall in the observed data
1181 interval on which the posterior distribution is conditioned, are displayed in light red (50% CI) and dark-red
1182 (95% CI). The corresponding observations are marked as black dots. Additionally, the posterior predictive
1183 distributions for future data are presented in light blue (50% CI) and dark-blue (95% CI). The actual
1184 realisations of these dates are marked as black triangles (F).

1185 **Figure 7. Elements of reproducibility in the research workflow**

1186 The figure shows good research practices across the research workflow that can contribute to
1187 reproducibility and demonstrates where the Bayesian research cycle (see Figure 1) and the WAMBS
1188 checklist (see Box 4) fit in the wider context of transparency in research. Not all elements are applicable
1189 to all types of research, e.g. preregistration is typically used for hypothesis-driven research but the
1190 specification of the prior and likelihood may be pre-registered. There can be legitimate reasons why not
1191 all data can be shared openly, but all scripts for running the Bayesian models could be shared on a data
1192 repository. *Note that part of the figure is based on a figure originally used in the Utrecht University*
1193 *Summerschool on Open Science and Scholarship 2019²⁵² (licensed CC-BY).*

1194 Boxes

1195 Box 1 | Bayes' Theorem

1196 In Bayesian statistics, all observed and unobserved quantities in a system are given a joint probability
1197 distribution, and inference for unobserved quantities is based on their conditional distribution given the
1198 observed data. By construction, Bayesian inferences are optimal when averaged over this joint
1199 distribution; in Bayesian terminology, the prior and data distributions. Rényi's axiom of probability²⁵³ lends
1200 itself to examining conditional probabilities, where the probabilities of Event A and Event B occurring are
1201 dependent, or conditional. The basic conditional probability can be written as:

$$p(B|A) = \frac{p(B \cap A)}{p(A)}, \quad (1)$$

1202 where the probability of Event B occurring is conditional on Event A. Equation 1 sets the foundation for
1203 Bayes' rule, which is a mathematical expression of Bayes' theorem that recognizes $p(B|A) \neq p(A|B)$ but
1204 $p(B \cap A) = p(A \cap B)$. Bayes' rule can be written as:

$$p(A|B) = \frac{p(A \cap B)}{p(B)}, \quad (2)$$

1205 which, based on Equation 1, can be reworked as:

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}. \quad (3 - \text{Bayes' rule})$$

1206 These principles can be extended to the situation of data and model parameters. With dataset \mathbf{y} and
1207 model parameters $\boldsymbol{\theta}$, Equation 3 (Bayes' rule) can be written as follows:

$$p(\boldsymbol{\theta}|\mathbf{y}) = \frac{p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y})}, \quad (4)$$

1208 which is often simplified to:

$$p(\boldsymbol{\theta}|\mathbf{y}) \propto p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta}). \quad (5)$$

1209 The term $p(\boldsymbol{\theta}|\mathbf{y})$ represents a conditional probability, where the probability of the model parameters ($\boldsymbol{\theta}$)
1210 is computed conditional upon the data (\mathbf{y}), and this term is also known as the *posterior*. The term $p(\mathbf{y}|\boldsymbol{\theta})$
1211 represents the conditional probability of the data given the model parameters, and this term represents
1212 the *data likelihood*. Finally, the term $p(\boldsymbol{\theta})$ represents the probability of particular model parameter values
1213 existing in the population. This term is called a *prior*. The term $p(\mathbf{y})$ is often viewed as a normalizing factor
1214 across all outcomes \mathbf{y} , which can be removed from the equation because $\boldsymbol{\theta}$ does not depend on \mathbf{y} or
1215 $p(\mathbf{y})$. Given that $p(\mathbf{y})$ is not needed for the posterior, it can be removed, and we say that the posterior is

1216 *proportional to* (\propto) the likelihood times the prior.. Figure 2 illustrates the relationship between the
1217 likelihood, prior, and posterior.
1218

1219 **Box 2 | The likelihood function for a coin experiment**

1220

1221 Consider the following textbook example: we are given a coin and want to know what the probability of
1222 obtaining “heads” (θ) is. To examine this, we toss the coin a number of times and count the number of
1223 heads. Let the outcome of the i th flip be denoted by $h_i = 1$ for heads and $h_i = 0$ for tails. The total
1224 experiment yields a sample of n independent binary observations $\{h_1, \dots, h_n\} = \mathbf{h}$ with y as the total
1225 number of heads; $y = \sum_{i=1}^n h_i$. We can assume that the probability to obtain heads remains constant over
1226 the experiment, i.e. $p(h_i) = \theta, (i = 1, \dots, n)$. Therefore the probability of the observed number of heads
1227 is expressed by the binomial distribution, given by

$$P(y|\theta) = \binom{n}{h} \theta^h (1 - \theta)^{n-h}, 0 \leq \theta \leq 1 \quad (1)$$

1228 ²⁵⁴.

1229 When y is kept fixed and θ is varying, $P(y|\theta)$ becomes a continuous function of θ , called the binomial
1230 likelihood function²⁵⁴.

1231 Suppose we flipped the coin 10 times and observed 4 heads, the likelihood function of θ is defined by

$$f(y|\theta) = \binom{10}{4} \theta^4 (1 - \theta)^6, 0 \leq \theta \leq 1. \quad (2)$$

1232 .

1233

1234

1235 **Box 3 | Bayes Factors**

1236 Hypothesis testing consists of using data to evaluate the evidence for competing claims or hypotheses.
1237 In the Bayesian framework, this can be accomplished using the Bayes factor, which corresponds to the
1238 ratio of the posterior odds to the prior odds of distinct hypotheses^{38,62}. For two hypotheses, H_0 and H_1 ,
1239 and observed data \mathbf{y} , the Bayes factor in favor of H_1 is given by

$$BF_{10} = \frac{p(H_1|\mathbf{y})/p(H_0|\mathbf{y})}{p(H_1)/p(H_0)}, \quad (6)$$

1240 where $p(H_0)$ and $p(H_1) = 1 - p(H_0)$ are the prior probabilities. A larger value of BF_{10} provides
1241 stronger evidence against H_0 ⁶². The posterior probability $p(H_j|\mathbf{y})$ is obtained using Bayes theorem

$$p(H_j|\mathbf{y}) = \frac{f(\mathbf{y}|H_j)p(H_j)}{f(\mathbf{y})}, j = \quad (7)$$

0,1.

1242 Thus, the Bayes factor can equivalently be written as the ratio of the marginal likelihoods of the
1243 observed data under the two hypotheses

$$BF_{10} = \frac{f(\mathbf{y}|H_1)}{f(\mathbf{y}|H_0)}. \quad (8)$$

1244 The competing hypotheses can take various forms and could be, for example, two non-nested
1245 regression models (see Variable Selection subsection). If H_0 and H_1 are simple hypotheses in which the
1246 parameters are fixed (e.g., $H_0: \mu = \mu_0$ versus $H_1: \mu = \mu_1$), the Bayes factor is identical to the likelihood
1247 ratio test. When either or both hypotheses are composite (i.e., not simple) or there are additional
1248 unknown parameters, the marginal likelihood $f(\mathbf{y}|H_j)$ is obtained by integrating over the parameters $\boldsymbol{\theta}_j$
1249 with prior densities $p(\boldsymbol{\theta}_j|H_j)$

$$f(\mathbf{y}|H_j) \quad (9)$$
$$= \int f(\mathbf{y}|\boldsymbol{\theta}_j, H_j) p(\boldsymbol{\theta}_j|H_j) d\boldsymbol{\theta}_j.$$

1250 This integral is often intractable and must be computed by numerical methods. If $p(\boldsymbol{\theta}_j|H_j)$ is improper
1251 (i.e., $\int p(\boldsymbol{\theta}_j|H_j) d\boldsymbol{\theta}_j = \infty$) then $f(\mathbf{y}|H_j)$ will be improper and the Bayes factor will not be uniquely
1252 defined. Overly diffuse priors should also be avoided, as they result in a Bayes factor that favors H_0
1253 regardless of the information in the data¹⁰⁴. As a simple illustrative example, suppose one collects n
1254 random samples from a normally distributed population with an unknown mean μ and a known variance

1255 σ^2 , and wishes to test $H_0: \mu = \mu_0$ versus $H_1: \mu \neq \mu_0$. Let \bar{y} be the sample mean. H_0 is a simple
1256 hypothesis with a point mass at μ_0 , so $\bar{y}|H_0 \sim N(\mu_0, \sigma^2/n)$. Under H_1 , $\bar{y}|\mu, H_1 \sim N(\mu, \sigma^2/n)$ and
1257 assuming $\mu|H_1 \sim N(\mu_0, \tau^2)$ with τ^2 fixed, then $f(\bar{y}|H_1) = \int f(\bar{y}|\mu, H_1) p(\mu|H_1) d\mu$ reduces to
1258 $\bar{y}|H_1 \sim N(\mu_0, \tau^2 + \sigma^2/n)$. Thus, the Bayes factor in favor of H_1 is

1259

$$BF_{10} = \frac{f(\bar{y}|H_1)}{f(\bar{y}|H_0)} = \frac{(\tau^2 + \sigma^2/n)^{-1/2} \exp\left\{-\frac{(\bar{y} - \mu_0)^2}{2(\tau^2 + \sigma^2/n)}\right\}}{(\sigma^2/n)^{-1/2} \exp\left\{-\frac{(\bar{y} - \mu_0)^2}{2(\sigma^2/n)}\right\}} \quad (10)$$

1260 For example, for $n = 20$, $\bar{y} = 5.8$, $\mu_0 = 5$, $\sigma^2 = 1$ and $\tau^2 = 1$, the Bayes factor is $BF_{10} = 96.83$,
1261 which provides strong evidence that the mean μ is not 5.

1262

1263

1264 Box 4 | WAMBS-Checklist

1265 **[bH1]** The 10 checklist points of WAMBS-v2

1266 [b1] Ensure the prior distributions and the model (or likelihood) are well understood and described in
1267 detail in the text, including the hyperparameter settings and all details surrounding the model. In
1268 addition, prior-predictive checking can help identify any prior-data conflict.

1269 [b1] Assess each parameter for convergence. Use multiple convergence diagnostics if possible. This may
1270 involve examining trace-plots or ensuring diagnostics (e.g., \hat{R} or effective sample size) are being met for
1271 each parameter. For example, \hat{R} values smaller than 1.05 are typically recommended. Likewise, effective
1272 sample sizes of 10,000 or more are recommended as a general rule of thumb.

1273 [b1] Sometimes convergence diagnostics can fail at detecting non-convergence within the chain.
1274 Subsequent measures, such as the split- \hat{R} can be used to identify such situations. The split- \hat{R} can detect
1275 trends that are missed if the chains have similar marginal distributions (the \hat{R} may miss these trends).

1276 [b1] Ensure that there were sufficient chain iterations to construct a meaningful posterior distribution.
1277 The posterior distribution should consist of enough samples to visually examine the shape, scale, and
1278 central tendency of the distribution. Without enough samples, there is an incomplete picture of the full
1279 distribution.

1280 [b1] Check all parameters for strong degrees of autocorrelation (e.g., through examining the effective
1281 sample size for parameters), which may be a sign of model or prior misspecification.

1282 [b1] Visually examine the marginal posteriors distribution for each model parameter to ensure that they
1283 make substantive sense. Posterior predictive distributions can be used to aid in examining the
1284 posteriors.

1285 [b1] Fully examine multivariate priors through a sensitivity analysis. These priors can be particularly
1286 influential on the posterior, even with slight modifications to the hyperparameters.

1287 [b1] To fully understand the impact of subjective priors, compare the posterior results to an analysis
1288 using diffuse (or objective) priors. This comparison can facilitate a deeper understanding of the impact
1289 the subjective priors (i.e., the theory being implemented) are having on findings. Next, conduct a full

1290 sensitivity analysis of all priors to gain a clearer understanding of the robustness of the results to
1291 different prior settings.

1292 [b1] Given the subjectivity of the model, it is also important to conduct a sensitivity analysis of the
1293 model (or likelihood) to help uncover how robust results are to deviations in the model.

1294 [b1] Report findings by including Bayesian interpretations. Take advantage of explaining and capturing
1295 the entire posterior rather than simply a point estimate. For example, it may be helpful to examine the
1296 density at different quantiles to fully capture and understand the posterior distribution.

1297

1298 Glossary Terms

1299

1300 **Prior distribution:** Beliefs held by researchers about the parameters in a statistical model BEFORE seeing
1301 the data.

1302 **Hyperparameters:** Hyperparameters are the parameters that define the prior distribution. For
1303 example, the normal distribution is defined through a mean and variance, and these are
1304 referred to as the hyperparameters.

1305 **Informative prior:** Informative priors reflect a high degree of certainty or knowledge surrounding
1306 the population parameters and the hyperparameters are specified to express particular
1307 information reflecting a greater degree of certainty about the model parameters being
1308 estimated

1309 **Weakly informative prior:** The weakly informative prior incorporates some information about
1310 the population parameter but are not as restrictive as an informative prior.; some researchers
1311 find this to be a nice middle ground regarding the informativeness of the prior

1312 **Diffuse priors:** Diffuse priors reflect complete uncertainty about population parameters.

1313 **Shrinkage priors** A specific prior that shrinks the posterior estimate towards a particular value.

1314 **Spike-and-slab prior** A specific shrinkage prior distribution used for variable selection that
1315 corresponds to a mixture of two distributions, one spiked around 0 and the other with a large
1316 variance corresponding to the slab component.

1317 **Horseshoe prior** A prior for variable selection that uses a half-Cauchy scale mixture of normal
1318 distribution. This prior is characterized by a high concentration around zero to shrink small
1319 coefficients and heavy tails to avoid excessive shrinkage of large parameters.

1320 **Prior predictive distribution** All possible samples that could occur if the model is true based on
1321 the priors. In theory, a “correct” prior provides a prior predictive distribution similar to the true
1322 data generating distribution

1323 **Prior predictive p-value** An estimate to indicate how unlikely the observed data is to be
1324 generated by the model based on the prior predictive distribution

1325

1326 **likelihood function** The conditional probability distribution $p(y|\theta)$ of the data y given parameters θ .

1327

1328 **posterior distribution** The posterior distribution reflects one’s updated knowledge, balancing prior
1329 knowledge with observed data.

1330

1331

1332 **Markov chain Monte Carlo (MCMC)** = A method to indirectly obtain inference on the posterior
1333 distribution via simulation which combines two concepts: (i) obtain a set of parameter values from the
1334 posterior distribution (using the Markov chain, or the first “MC”); and (ii) given sampled parameter
1335 values obtain a distributional estimate of the posterior and associated posterior statistics of interest
1336 (using Monte Carlo, or the second “MC”).

1337

1338 **Trace plots** A plot describing the posterior parameter value at each iteration of the Markov chain (on the
1339 y-axis) against iteration number (on the x-axis)

1340

1341 **\hat{R} statistic** \hat{R} is defined to be the ratio of the within- and between-chain variability. Values close to 1 for
1342 all parameters and quantities of interest suggest the chain has sufficiently converged to the stationary
1343 distribution

1344

1345 **Bayes factor** Bayes factors (Box 3) can be used to compare and choose between candidate models,
1346 where each candidate model would correspond to a hypothesis

1347

1348 **kernel density estimation** A kernel density estimation is a non-parametric approach used to estimate a
1349 probability density function for the observed data.

1350 **transition kernel** determines the performance of the MCMC algorithm in terms of how long the Markov
1351 chain needs to be run to obtain reliable inference on the posterior distribution of interest.

1352 **auxiliary variables** additional variables entered in the model to improve the missing data model.

1353 **sparsity**: indicates that most parameter values are zero and only a few are non-zero.

1354

1355 **Stochastic Gradient Descent (SGD) algorithm.** SGD algorithms use a randomly chosen subset of data
1356 points to estimate the gradient of a loss function with respect to parameters. This can provide radical
1357 computational savings in optimisation problems involving many data points.

1358

1359 **Variational Inference (VI).** Variational methods refers to a class of approximate inference techniques in
1360 which deterministic posterior approximations are constructed from a family of predefined distributions.
1361 These approximations contain variational parameters which are optimised to match the approximating
1362 distribution as closely as possible to the true posterior. They are popular methods for achieving scalable
1363 but approximate Bayesian inference in large data scenarios where MCMC sampling-based inference
1364 would be prohibitive.

1365

1366

1367

Highlighted Refences

1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404

1. O'Hagan, A., Buck, C. E., Daneshkhan, A., Eiser, J. R., Garthwaite, P. H., Jenkinson, D. J.,..... & Rakow, T. (2006). *Uncertain judgements: eliciting experts' probabilities*. John Wiley & Sons.
This book is a great collection of information with respect to prior elicitation. It includes elicitation techniques, summarizes potential pitfalls, and describes examples across a wide variety of disciplines.

2. E.J. George and R.E. McCulloch (1993). Variable selection via Gibbs sampling. *Journal of the American Statistical Association*, 88: 881-889.
This is the paper that popularized the use of spike-and-slab priors for Bayesian variable selection and introduced MCMC techniques to explore the model space.

3. N.G. Polson and J.G. Scott (2010). Shrink globally, act locally: Sparse Bayesian regularization and prediction. *Bayesian Statistics 9*, 9: 501-538.
This paper provides a unified framework for continuous shrinkage priors, which allow global sparsity while controlling the amount of regularization for each regression coefficient.

4. M.G. Tadesse and M. Vannucci (2020). Handbook of Mixture Analysis. CRC Press, Chapman & Hall/CRC Handbooks of Modern Statistical Methods, in preparation.
This is a forthcoming edited book that presents a comprehensive review of Bayesian variable selection methods and highlights recent developments.

5. Gelfand, A. E. & Smith, A. F. M. Sampling-Based Approaches to Calculating Marginal Densities. *J. Am. Stat. Assoc.* 85, 398-409, doi:10.1080/01621459.1990.10476213 (1990).
Seminal paper that identified Markov chain Monte Carlo as a practical approach for Bayesian inference.

6. Lunn, D.J., Thomas, A., Best, N., and Spiegelhalter, D. (2000) WinBUGS — a Bayesian modelling framework: concepts, structure, and extensibility. *Statistics and Computing*, 10, 325–337
Provided an early user-friendly and freely-available black-box MCMC sampler opening up Bayesian inference to the wider scientific community.

7. Brooks, S. P., Gelman, A., Jones, G., Meng, X. (Eds) (2011) *Handbook of Markov chain Monte Carlo*. CRC Press.
Comprehensive review of Markov chain Monte Carlo and its use in many different applications.

1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441

- 8. Depaoli, S., and van de Schoot, R. (2017). Improving transparency and replication in Bayesian statistics: The WAMBS-checklist. *Psychological Methods*, 22, 240-261.**

This paper goes through, in a step-by-step manner, the various points that need to be checked when estimating a model via Bayesian statistics. It can be used as a guide for implementing Bayesian methods.

- 9. Kass, R.E., Raftery, A.E. (1995). Bayes factors. *Journal of the American Statistical Association*, 90: 773-795.**

This paper provides an extensive discussion of Bayes factors with several examples.

- 10. Blei, D. M., et al. (2017). "Variational inference: A review for statisticians." *Journal of the American statistical Association* 112(518): 859–877.**

Recent review of variational inference methods, including stochastic variants, which underpin popular approximate Bayesian inference methods for large data or complex modelling problems where computation using MCMC stochastic simulation would be prohibitively costly.

- 11. Kingma, D. P. and M. Welling (2019). An Introduction to Variational Autoencoders.**

Recent review of variational autoencoders, encompassing deep generative models, the reparameterisation trick and current inference methods. These are an important class of models in modern Bayesian machine learning that combines the use of Bayesian modelling with deep neural networks for flexible function parameterisation.

- 12. Neal, R. M. (1996). Priors for Infinite Networks. *Bayesian Learning for Neural Networks*. R. M. Neal. New York, NY, Springer New York: 29-53.**

A classic text highlighting the connection between neural networks and Gaussian processes and the application of Bayesian approaches for fitting neural networks.

- 13. Berger, J. (2006). The case for objective Bayesian analysis. *Bayesian Analysis*, 1(3), 385-402.**

"A discussion of objective Bayesian analysis, including criticisms of the approach and a personal perspective on the debate on the value of objective Bayesian versus subjective Bayesian analysis."

- 14. Tanner, M. A., & Wong, W. H. (1987). The calculation of posterior distributions by data augmentation. *1098 Journal of the American statistical Association*, 82(398), 528-540.**

"In this article the authors explain how to use data augmentation when direct computation of the posterior density of the parameters of interest is not possible."

References

- 1444 1 Bayes, M. & Price, M. LII. An essay towards solving a problem in the doctrine of chances. By the
 1445 late Rev. Mr. Bayes, F. R. S. communicated by Mr. Price, in a letter to John Canton, A. M. F. R. S.
 1446 *Philosophical Transactions of the Royal Society of London* **53**, 370-418, doi:10.1098/rstl.1763.0053
 1447 (1997).
- 1448 2 Laplace, P. S. *Essai Philosophique sur les Probabilités*. (Courcier, 1814).
- 1449 3 König, C. & van de Schoot, R. Bayesian statistics in educational research: a look at the current state
 1450 of affairs. *Educational Review*, 1-24 (2017).
- 1451 4 van de Schoot, R., Winter, S., Zondervan-Zwijnenburg, M., Ryan, O. & Depaoli, S. A systematic
 1452 review of Bayesian applications in psychology: The last 25 years. *Psychological Methods* **22**, 217-
 1453 239 (2017).
- 1454 5 Ashby, D. Bayesian statistics in medicine: a 25 year review. *Stat Med* **25**, 3589-3631,
 1455 doi:10.1002/sim.2672 (2006).
- 1456 6 Rietbergen, C., Debray, T. P. A., Klugkist, I., Janssen, K. J. M. & Moons, K. G. M. Reporting of
 1457 Bayesian analysis in epidemiologic research should become more transparent. *J Clin Epidemiol* **86**,
 1458 51-58 e52, doi:10.1016/j.jclinepi.2017.04.008 (2017).
- 1459 7 Spiegelhalter, D. J., Myles, J. P., Jones, D. R. & Abrams, K. R. Bayesian methods in health technology
 1460 assessment: a review. *Health Technol Assess* **4**, 1-130, doi:10.3310/hta4380 (2000).
- 1461 8 Kruschke, J. K., Aguinis, H. & Joo, H. The Time Has Come Bayesian Methods for Data Analysis in
 1462 the Organizational Sciences. *Organizational Research Methods* **15**, 722-752 (2012).
- 1463 9 Smid, S. C., McNeish, D., Miočević, M. & van de Schoot, R. Bayesian Versus Frequentist Estimation
 1464 for Structural Equation Models in Small Sample Contexts: A Systematic Review. *Structural*
 1465 *Equation Modeling: A Multidisciplinary Journal* **27**, 131-161,
 1466 doi:10.1080/10705511.2019.1577140 (2019).
- 1467 10 Rupp, A. A., Dey, D. K. & Zumbo, B. D. To bayes or not to bayes, from whether to when:
 1468 Applications of Bayesian methodology to modeling. *Structural Equation Modeling* **11**, 424-451
 1469 (2004).
- 1470 11 Depaoli, S. & Clifton, J. P. A Bayesian approach to multilevel structural equation modeling with
 1471 continuous and dichotomous outcomes. *Structural Equation Modeling: A Multidisciplinary Journal*
 1472 **22**, 327-351 (2015).
- 1473 12 Kim, S.-Y., Suh, Y., Kim, J.-S., Albanese, M. A. & Langer, M. M. Single and multiple ability estimation
 1474 in the SEM framework: A noninformative Bayesian estimation approach. *Multivariate behavioral*
 1475 *research* **48**, 563-591 (2013).
- 1476 13 Depaoli, S. Mixture class recovery in GMM under varying degrees of class separation: frequentist
 1477 versus Bayesian estimation. *Psychol Methods* **18**, 186-219, doi:10.1037/a0031609 (2013).
- 1478 14 Blaxter, L. *How to research*. (McGraw-Hill Education (UK), 2010).
- 1479 15 Neuman, W. L. *Understanding research*. (Pearson, 2016).
- 1480 16 Heo, I. & Van de Schoot, R. Tutorial: Advanced Bayesian regression in JASP. (2020).
- 1481 17 Van de Schoot, R., Yerkes, M. A., Mouw, J. M. & Sonneveld, H. What took them so long? Explaining
 1482 PhD delays among doctoral candidates. *PloS one* **8**, e68839 (2013).
- 1483 18 Muthen, B. & Asparouhov, T. Bayesian structural equation modeling: a more flexible
 1484 representation of substantive theory. *Psychol Methods* **17**, 313-335, doi:10.1037/a0026802
 1485 (2012).
- 1486 19 van de Schoot, R. *et al.* Facing off with Scylla and Charybdis: a comparison of scalar, partial, and
 1487 the novel possibility of approximate measurement invariance. *Front Psychol* **4**, 770,
 1488 doi:10.3389/fpsyg.2013.00770 (2013).

- 1489 20 O'Hagan, A. *et al.* *Uncertain judgements: eliciting experts' probabilities*. (John Wiley & Sons,
1490 2006).
- 1491 21 Howard, G. S., Maxwell, S. E. & Fleming, K. J. The proof of the pudding: an illustration of the
1492 relative strengths of null hypothesis, meta-analysis, and Bayesian analysis. *Psychol Methods* **5**,
1493 315-332, doi:10.1037/1082-989x.5.3.315 (2000).
- 1494 22 Veen, D., Stoel, D., Zondervan-Zwijnenburg, M. & van de Schoot, R. Proposal for a Five-Step
1495 Method to Elicit Expert Judgement. *Frontiers in psychology* **8**, 2110 (2017).
- 1496 23 Johnson, S. R., Tomlinson, G. A., Hawker, G. A., Granton, J. T. & Feldman, B. M. Methods to elicit
1497 beliefs for Bayesian priors: a systematic review. *J Clin Epidemiol* **63**, 355-369,
1498 doi:10.1016/j.jclinepi.2009.06.003 (2010).
- 1499 24 Ibrahim, J. G., Chen, M. H., Gwon, Y. & Chen, F. The power prior: theory and applications. *Stat*
1500 *Med* **34**, 3724-3749, doi:10.1002/sim.6728 (2015).
- 1501 25 Rietbergen, C., Klugkist, I., Janssen, K. J., Moons, K. G. & Hoijtink, H. J. Incorporation of historical
1502 data in the analysis of randomized therapeutic trials. *Contemp Clin Trials* **32**, 848-855,
1503 doi:10.1016/j.cct.2011.06.002 (2011).
- 1504 26 van de Schoot, R. *et al.* Bayesian PTSD-Trajectory Analysis with Informed Priors Based on a
1505 Systematic Literature Search and Expert Elicitation. *Multivariate Behav Res* **53**, 267-291,
1506 doi:10.1080/00273171.2017.1412293 (2018).
- 1507 27 Smeets, L. & van de Schoot, R. Code for the ShinyApp to Determine the Plausible Parameter Space
1508 for the PhD-delay Data (Version v1.0). . (2020).
- 1509 28 Berger, J. The case for objective Bayesian analysis. *Bayesian analysis* **1**, 385-402 (2006).
- 1510 29 Brown, L. D. In-season prediction of batting averages: A field test of empirical Bayes and Bayes
1511 methodologies. *The Annals of Applied Statistics*, 113-152 (2008).
- 1512 30 Candel, M. J. & Winkens, B. Performance of empirical Bayes estimators of level-2 random
1513 parameters in multilevel analysis: A Monte Carlo study for longitudinal designs. *Journal of*
1514 *Educational and Behavioral Statistics* **28**, 169-194 (2003).
- 1515 31 van der Linden, W. J. Using response times for item selection in adaptive testing. *Journal of*
1516 *Educational and Behavioral Statistics* **33**, 5-20 (2008).
- 1517 32 Darnieder, W. F. *Bayesian methods for data-dependent priors*, The Ohio State University, (2011).
- 1518 33 Richardson, S. & Green, P. J. On Bayesian Analysis of Mixtures with an Unknown Number of
1519 Components (with discussion). *Journal of the Royal Statistical Society: Series B (Statistical*
1520 *Methodology)* **59**, 731-792, doi:10.1111/1467-9868.00095 (1997).
- 1521 34 Wasserman, L. Asymptotic inference for mixture models by using data-dependent priors. *Journal*
1522 *of the Royal Statistical Society: Series B (Statistical Methodology)* **62**, 159-180, doi:10.1111/1467-
1523 9868.00226 (2000).
- 1524 35 Chung, Y., Gelman, A., Rabe-Hesketh, S., Liu, J. & Dorie, V. Weakly informative prior for point
1525 estimation of covariance matrices in hierarchical models. *Journal of Educational and Behavioral*
1526 *Statistics* **40**, 136-157 (2015).
- 1527 36 Gelman, A., Jakulin, A., Pittau, M. G. & Su, Y.-S. A weakly informative default prior distribution for
1528 logistic and other regression models. *The annals of applied statistics* **2**, 1360-1383 (2008).
- 1529 37 Gelman, A., Carlin, J. B., Stern, H. S. & Rubin, D. B. *Bayesian data analysis*. Vol. 2
1530 (Chapman&HallCRC, 2004).
- 1531 38 Jeffreys, H. *Theory of probability*. Vol. 3 (Clarendon Press, 1961).
- 1532 39 Seaman III, J. W., Seaman Jr, J. W. & Stamey, J. D. Hidden dangers of specifying noninformative
1533 priors. *The American Statistician* **66**, 77-84, doi:<https://doi.org/10.1080/00031305.2012.695938>
1534 (2012).
- 1535 40 Gelman, A. Prior distributions for variance parameters in hierarchical models (comment on article
1536 by Browne and Draper). *Bayesian analysis* **1**, 515-534 (2006).

1537 41 Lambert, P. C., Sutton, A. J., Burton, P. R., Abrams, K. R. & Jones, D. R. How vague is vague? A
1538 simulation study of the impact of the use of vague prior distributions in MCMC using WinBUGS.
1539 *Stat Med* **24**, 2401-2428, doi:10.1002/sim.2112 (2005).

1540 42 McNeish, D. On using Bayesian methods to address small sample problems. *Structural Equation*
1541 *Modeling: A Multidisciplinary Journal* **23**, 750-773 (2016).

1542 43 van de Schoot, R. & Miocević, M. *Small sample size solutions: A guide for applied researchers and*
1543 *practitioners*. (Taylor & Francis, 2020).

1544 44 Schuurman, N. K., Grasman, R. P. & Hamaker, E. L. A Comparison of Inverse-Wishart Prior
1545 Specifications for Covariance Matrices in Multilevel Autoregressive Models. *Multivariate Behav*
1546 *Res* **51**, 185-206, doi:10.1080/00273171.2015.1065398 (2016).

1547 45 Liu, H., Zhang, Z. & Grimm, K. J. Comparison of inverse Wishart and separation-strategy priors for
1548 Bayesian estimation of covariance parameter matrix in growth curve analysis. *Structural Equation*
1549 *Modeling: A Multidisciplinary Journal* **23**, 354-367 (2016).

1550 46 Ranganath, R. & Blei, D. M. Population predictive checks. *arXiv preprint arXiv:1908.00882* (2019).

1551 47 Daimon, T. Predictive checking for Bayesian interim analyses in clinical trials. *Contemp Clin Trials*
1552 **29**, 740-750, doi:10.1016/j.cct.2008.05.005 (2008).

1553 48 Morris, D. E., Oakley, J. E. & Crowe, J. A. A web-based tool for eliciting probability distributions
1554 from experts. *Environmental Modelling & Software* **52**, 1-4 (2014).

1555 49 Garthwaite, P. H., Al-Awadhi, S. A., Elfadaly, F. G. & Jenkinson, D. J. Prior distribution elicitation
1556 for generalized linear and piecewise-linear models. *Journal of Applied Statistics* **40**, 59-75 (2013).

1557 50 Elfadaly, F. G. & Garthwaite, P. H. Eliciting Dirichlet and Gaussian copula prior distributions for
1558 multinomial models. *Statistics and Computing* **27**, 449-467 (2017).

1559 51 Veen, D., Egberts, M. R., van Loey, N. E. E. & van de Schoot, R. Expert Elicitation for Latent Growth
1560 Curve Models: The Case of Posttraumatic Stress Symptoms Development in Children With Burn
1561 Injuries. *Front Psychol* **11**, 1197, doi:10.3389/fpsyg.2020.01197 (2020).

1562 52 Runge, A. K., Scherbaum, F., Curtis, A. & Riggelsen, C. An interactive tool for the elicitation of
1563 subjective probabilities in probabilistic seismic-hazard analysis. *Bulletin of the Seismological*
1564 *Society of America* **103**, 2862-2874 (2013).

1565 53 Zondervan-Zwijnenburg, M., van de Schoot-Hubeek, W., Lek, K., Hoijtink, H. & van de Schoot, R.
1566 Application and Evaluation of an Expert Judgment Elicitation Procedure for Correlations. *Front*
1567 *Psychol* **8**, 90, doi:10.3389/fpsyg.2017.00090 (2017).

1568 54 Cooke, R. M. & Goossens, L. H. J. TU Delft expert judgment data base. *Reliability Engineering &*
1569 *System Safety* **93**, 657-674 (2008).

1570 55 Hanea, A. M., Nane, G. F., Bedford, T. & French, S. *Expert Judgment in Risk and Decision Analysis*.
1571 (Springer, 2020).

1572 56 Dias, L. C., Morton, A. & Quigley, J. *Elicitation*. (Springer, 2018).

1573 57 Box, G. E. Sampling and Bayes' inference in scientific modelling and robustness. *Journal of the*
1574 *Royal Statistical Society: Series A (General)* **143**, 383-404 (1980).

1575 58 Nott, D. J., Drovandi, C. C., Mengersen, K. & Evans, M. Approximation of Bayesian Predictive p-
1576 Values with Regression ABC. *Bayesian Analysis* **13**, 59-83 (2018).

1577 59 Evans, M. & Moshonov, H. Checking for prior-data conflict with hierarchically specified priors.
1578 *Bayesian statistics and its applications*, 145-159 (2007).

1579 60 Evans, M. & Jang, G. H. A limit result for the prior predictive applied to checking for prior-data
1580 conflict. *Statistics & Probability Letters* **81**, 1034-1038, doi:10.1016/j.spl.2011.02.025 (2011).

1581 61 Young, K. & Pettit, L. Measuring discordancy between prior and data. *Journal of the Royal*
1582 *Statistical Society: Series B (Methodological)* **58**, 679-689 (1996).

1583 62 Kass, R. E. & Raftery, A. E. Bayes factors. *Journal of the american statistical association* **90**, 773-
1584 795 (1995).

1585 63 Bousquet, N. Diagnostics of prior-data agreement in applied Bayesian analysis. *Journal of Applied*
1586 *Statistics* **35**, 1011-1029, doi:<https://doi.org/10.1080/02664760802192981> (2008).

1587 64 Veen, D., Stoel, D., Schalken, N., Mulder, K. & van de Schoot, R. Using the Data Agreement
1588 Criterion to Rank Experts' Beliefs. *Entropy* **20**, 592, doi:10.3390/e20080592 (2018).

1589 65 Nott, D. J., Xueou, W., Evans, M. & Englert, B. Checking for prior-data conflict using prior to
1590 posterior divergences. *arXiv preprint arXiv:1611.00113* (2016).

1591 66 Lek & Van De, S. How the Choice of Distance Measure Influences the Detection of Prior-Data
1592 Conflict. *Entropy* **21**, 446, doi:10.3390/e21050446 (2019).

1593 67 Gabry, J., Simpson, D., Vehtari, A., Betancourt, M. & Gelman, A. Visualization in Bayesian
1594 workflow. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* **182**, 389-402
1595 (2019).

1596 68 O'Hagan, A. Bayesian statistics: principles and benefits. *Frontis*, 31-45 (2004).

1597 69 Etz, A. Introduction to the Concept of Likelihood and Its Applications. *Advances in Methods and*
1598 *Practices in Psychological Science* **1**, 60-69, doi:10.1177/2515245917744314 (2018).

1599 70 Pawitan, Y. *In all likelihood: statistical modelling and inference using likelihood*. (Oxford University
1600 Press, 2001).

1601 71 Gelman, A., Simpson, D. & Betancourt, M. The prior can often only be understood in the context
1602 of the likelihood. *Entropy* **19**, 555, doi:<https://doi.org/10.3390/e19100555s> (2017).

1603 72 Aczel, B. *et al.* Discussion points for Bayesian inference. *Nat Hum Behav* **4**, 561-563,
1604 doi:10.1038/s41562-019-0807-z (2020).

1605 73 Gelman, A. *et al.* *Bayesian data analysis*. (CRC press, 2013).

1606 74 Greco, L., Racugno, W. & Ventura, L. Robust likelihood functions in Bayesian inference. *Journal of*
1607 *Statistical Planning and Inference* **138**, 1258-1270, doi:10.1016/j.jspi.2007.05.001 (2008).

1608 75 Shyamalkumar, N. D. in *Robust Bayesian Analysis Lecture Notes in Statistics* Ch. Chapter 7, 127-
1609 143 (Springer, 2000).

1610 76 Agostinelli, C. & Greco, L. A weighted strategy to handle likelihood uncertainty in Bayesian
1611 inference. *Computational Statistics* **28**, 319-339 (2013).

1612 77 Rubin, D. B. Bayesianly justifiable and relevant frequency calculations for the applies statistician.
1613 *The Annals of Statistics*, 1151-1172 (1984).

1614 78 Gelfand, A. E. & Smith, A. F. M. Sampling-Based Approaches to Calculating Marginal Densities.
1615 *Journal of the American Statistical Association* **85**, 398-409,
1616 doi:10.1080/01621459.1990.10476213 (1990).

1617 79 Geyer, C. J. Markov chain Monte Carlo maximum likelihood. 156-163 (1991).

1618 80 Van de Schoot, R., Veen, D., Smeets, L., Winter, S. D. & Depaoli, S. A tutorial on using the WAMBS
1619 checklist to avoid the misuse of Bayesian statistics. *Small Sample Size Solutions: A Guide for*
1620 *Applied Researchers and Practitioners; van de Schoot, R., Miocevic, M., Eds*, 30-49 (2020).

1621 81 Veen, D. & Egberts, M. The Importance of Collaboration in Bayesian Analyses with Small Samples.
1622 *SMALL SAMPLE SIZE SOLUTIONS*, 50 (2020).

1623 82 Robert, C. & Casella, G. *Monte Carlo statistical methods*. (Springer Science & Business Media,
1624 2013).

1625 83 Silverman, B. W. *Density estimation for statistics and data analysis*. Vol. 26 (CRC press, 1986).

1626 84 Geman, S. & Geman, D. Stochastic relaxation, gibbs distributions, and the bayesian restoration of
1627 images. *IEEE Trans Pattern Anal Mach Intell* **6**, 721-741, doi:10.1109/tpami.1984.4767596 (1984).

1628 85 Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H. & Teller, E. Equation of state
1629 calculations by fast computing machines. *The journal of chemical physics* **21**, 1087-1092 (1953).

1630 86 Hastings, W. K. Monte Carlo sampling methods using Markov chains and their applications.
1631 *Biometrika* **57**, 97-109 (1970).

1632 87 Duane, S., Kennedy, A. D., Pendleton, B. J. & Roweth, D. Hybrid monte carlo. *Physics letters B* **195**,
1633 216-222 (1987).

1634 88 Tanner, M. A. & Wong, W. H. The calculation of posterior distributions by data augmentation.
1635 *Journal of the American statistical Association* **82**, 528-540 (1987).

1636 89 Gelman, A. *Burn-in for MCMC, why we prefer the term warm-up*,
1637 <[https://statmodeling.stat.columbia.edu/2017/12/15/burn-vs-warm-iterative-simulation-
1638 algorithms/](https://statmodeling.stat.columbia.edu/2017/12/15/burn-vs-warm-iterative-simulation-algorithms/)> (2017).

1639 90 Gelman, A. & Rubin, D. B. Inference from iterative simulation using multiple sequences. *Statistical
1640 Science* **7**, 457-511 (1992).

1641 91 Brooks, S. P. & Gelman, A. General methods for monitoring convergence of iterative simulations.
1642 *Journal of Computational and Graphical Statistics* **7**, 434-455 (1998).

1643 92 Roberts, G. O. Markov chain concepts related to sampling algorithms. *Markov chain Monte Carlo
1644 in practice* **57**, 45-58 (1996).

1645 93 Vehtari, A., Gelman, A., Simpson, D., Carpenter, B. & Bürkner, P. (2020).

1646 94 Gamerman, D. & Lopes, H. F. *Markov chain Monte Carlo: stochastic simulation for Bayesian
1647 inference*. (CRC Press, 2006).

1648 95 Brooks, S. P., Gelman, A., Jones, G. & Meng, X.-L. *Handbook of markov chain monte carlo*. (CRC
1649 press, 2011).

1650 96 Bürkner, P.-C. Advanced Bayesian multilevel modeling with the R package brms. *arXiv preprint
1651 arXiv:1705.11123* (2017).

1652 97 Merkle, E. C. & Rosseel, Y. blavaan: Bayesian structural equation models via parameter expansion.
1653 *arXiv preprint arXiv:1511.05604* (2015).

1654 98 Carpenter, B. *et al.* Stan: A Probabilistic Programming Language. *Journal of Statistical Software*
1655 **76**, doi:10.18637/jss.v076.i01 (2017).

1656 99 Blei, D. M., Kucukelbir, A. & McAuliffe, J. D. Variational inference: A review for statisticians. *Journal
1657 of the American statistical Association* **112**, 859–877 (2017).

1658 100 Minka, T. P. 362–369.

1659 101 Hoffman, M. D., Blei, D. M., Wang, C. & Paisley, J. Stochastic variational inference. *The Journal of
1660 Machine Learning Research* **14**, 1303–1347 (2013).

1661 102 Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980
1662* (2014).

1663 103 Li, Y., Hernández-Lobato, J. M. & Turner, R. E. 2323–2331.

1664 104 Liang, F., Paulo, R., Molina, G., Clyde, M. A. & Berger, J. O. Mixtures of g priors for Bayesian variable
1665 selection. *Journal of the American Statistical Association* **103**, 410-423 (2008).

1666 105 Forte, A., Garcia-Donato, G. & Steel, M. Methods and tools for Bayesian variable selection and
1667 model averaging in normal linear regression. *International Statistical Review* **86**, 237-258 (2018).

1668 106 Mitchell, T. J. & Beauchamp, J. J. Bayesian variable selection in linear regression. *Journal of the
1669 American Statistical association* **83**, 1023-1032 (1988).

1670 107 George, E. J. & McCulloch, R. E. Variable selection via Gibbs sampling. *Journal of the American
1671 Statistical Association* **88**, 881-889-889 (1993).

1672 108 Ishwaran, H. & Rao, J. S. Spike and slab variable selection: frequentist and Bayesian strategies.
1673 *Annals of Statistics* **33**, 730-773 (2005).

1674 109 Bottolo, L. & Richardson, S. Evolutionary stochastic search. *Bayesian Analysis* **5**, 583-618 (2010).

1675 110 Ročková, V. & George, E. I. EMVS: the EM approach to Bayesian variable selection. *Journal of the
1676 American Statistical Association* **109**, 828-846 (2014).

1677 111 Park, T. & Casella, G. The Bayesian lasso. *Journal of the American Statistical Association* **103**, 681-
1678 686 (2008).

1679 112 Carvalho, C. M., Polson, N. G. & Scott, J. G. The horseshoe estimator for sparse signals. *Biometrika*
1680 **97**, 465-480 (2010).

1681 113 Polson, N. G. & Scott, J. G. Shrink globally, act locally: Sparse Bayesian regularization and
1682 prediction. *Bayesian statistics* **9**, 105 (2010).

1683 114 Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical*
1684 *Society: Series B (Methodological)* **58**, 267-288 (1996).

1685 115 Van Erp, S., Oberski, D. L. & Mulder, J. Shrinkage priors for Bayesian penalized regression. *Journal*
1686 *of Mathematical Psychology* **89**, 31-50 (2019).

1687 116 Brown, P. J., Vannucci, M. & Fearn, T. Multivariate Bayesian variable selection and prediction.
1688 *Journal of the Royal Statistical Society, Series B* **60**, 627-641 (1998).

1689 117 Lee, K. H., Tadesse, M. G., Baccarelli, A. A., Schwartz, J. & Coull, B. A. Multivariate Bayesian variable
1690 selection exploiting dependence structure among outcomes: Application to air pollution effects
1691 on DNA methylation. *Biometrics* **73**, 232-241 (2017).

1692 118 Frühwirth-Schnatter, S. & Wagner, H. Stochastic model specification search for Gaussian and
1693 partially non-Gaussian state space models. *Journal of Econometrics* **154**, 85-100 (2010).

1694 119 Scheipl, F., Fahrmeir, L. & Kneib, T. Spike-and-slab priors for function selection in structured
1695 additive regression models. *Journal of the American Statistical Association* **107**, 1518-1532 (2012).

1696 120 Tadesse, M. G., Sha, N. & Vannucci, M. Bayesian variable selection in clustering high dimensional
1697 data. *Journal of the American Statistical Association*, 602-617 (2005).

1698 121 Wang, H. Scaling it up: Stochastic search structure learning in graphical models. *Bayesian Analysis*
1699 **10**, 351-377 (2015).

1700 122 Peterson, C. B., Stingo, F. C. & Vannucci, M. Bayesian Inference of Multiple Gaussian Graphical
1701 Models. *J Am Stat Assoc* **110**, 159-174, doi:10.1080/01621459.2014.896806 (2015).

1702 123 Tadesse, M. G. & Vannucci, M. *Handbook of Bayesian variable selection*. (CRC Press, 2020).

1703 124 Li, F. & Zhang, N. R. Bayesian variable selection in structured high-dimensional covariate spaces
1704 with applications in genomics. *Journal of the American Statistical association* **105**, 1978-2002
1705 (2010).

1706 125 Stingo, F., Chen, Y., Tadesse, M. G. & Vannucci, M. Incorporating biological information into linear
1707 models: A Bayesian approach to the selection of pathways and genes. *Annals of Applied Statistics*
1708 **5**, 1202-1214 (2011).

1709 126 Guan, Y. & Stephens, M. Bayesian variable selection regression for genome-wide association
1710 studies and other large-scale problems. *Annals of Applied Statistics* **5**, 1780-1815 (2011).

1711 127 Bottolo, L. *et al.* GUESS-ing polygenic associations with multiple phenotypes using a GPU-based
1712 evolutionary stochastic search algorithm. *PLoS Genetics* **9(8)**, e1003657-e1003657 (2013).

1713 128 Banerjee, S., Carlin, B. P. & Gelfand, A. E. *Hierarchical modeling and analysis for spatial data*. (CRC
1714 press, 2014).

1715 129 Vock, L. F. B., Reich, B. J., Fuentes, M. & Dominici, F. Spatial variable selection methods for
1716 investigating acute health effects of fine particulate matter components. *Biometrics* **71**, 167-177
1717 (2015).

1718 130 Penny, W. D., Trujillo-Barreto, N. J. & Friston, K. J. Bayesian fMRI time series analysis with spatial
1719 priors. *Neuroimage* **24**, 350-362, doi:10.1016/j.neuroimage.2004.08.034 (2005).

1720 131 Smith, M., Pütz, B., Auer, D. & Fahrmeir, L. Assessing brain activity through spatial Bayesian
1721 variable selection. *Neuroimage* **20**, 802-815 (2003).

1722 132 Zhang, L., Guindani, M., Versace, F. & Vannucci, M. A spatio-temporal nonparametric Bayesian
1723 variable selection model of fMRI data for clustering correlated time courses. *Neuroimage* **95**, 162-
1724 175, doi:10.1016/j.neuroimage.2014.03.024 (2014).

1725 133 Gorrostieta, C., Fiecas, M., Ombao, H., Burke, E. & Cramer, S. Hierarchical vector auto-regressive
1726 models and their applications to multi-subject effective connectivity. *Frontiers on Computational*
1727 *Neurosciences* **7**, 159-159 (2013).

1728 134 Chiang, S. *et al.* Bayesian vector autoregressive model for multi-subject effective connectivity
1729 inference using multi-modal neuroimaging data. *Human Brain Mapping* **38**, 1311-1332 (2017).

1730 135 Schad, D. J., Betancourt, M. & Vasishth, S. Toward a principled Bayesian workflow in cognitive
1731 science. *arXiv preprint arXiv:1904.12765* (2019).

1732 136 Gelman, A., Meng, X.-L. & Stern, H. Posterior predictive assessment of model fitness via realized
1733 discrepancies. *Statistica sinica*, 733-760 (1996).

1734 137 Meng, X.-L. Posterior predictive p-values. *The annals of statistics* **22**, 1142-1160 (1994).

1735 138 Asparouhov, T., Hamaker, E. L. & Muthén, B. Dynamic structural equation models. *Structural*
1736 *Equation Modeling: A Multidisciplinary Journal* **25**, 359-388 (2018).

1737 139 Zhang, Z., Hamaker, E. L. & Nesselroade, J. R. Comparisons of four methods for estimating a
1738 dynamic factor model. *Structural Equation Modeling: A Multidisciplinary Journal* **15**, 377-402
1739 (2008).

1740 140 Hamaker, E., Ceulemans, E., Grasman, R. & Tuerlinckx, F. Modeling affect dynamics: State of the
1741 art and future challenges. *Emotion Review* **7**, 316-322 (2015).

1742 141 Meissner, P. *wikipediatrend: Public Subject Attention via Wikipedia Page View Statistics*. (2019).

1743 142 Harvey, A. C. & Peters, S. Estimation procedures for structural time series models. *Journal of*
1744 *Forecasting* **9**, 89-108 (1990).

1745 143 Taylor, S. J. & Letham, B. Forecasting at scale. *The American Statistician* **72**, 37-45 (2018).

1746 144 Gopnik, A. & Bonawitz, E. Bayesian models of child development. *Wiley Interdiscip Rev Cogn Sci*
1747 **6**, 75-86, doi:10.1002/wcs.1330 (2015).

1748 145 Gigerenzer, G. & Hoffrage, U. How to improve Bayesian reasoning without instruction: frequency
1749 formats. *Psychological review* **102**, 684 (1995).

1750 146 Slovic, P. & Lichtenstein, S. Comparison of Bayesian and regression approaches to the study of
1751 information processing in judgment. *Organizational behavior and human performance* **6**, 649-744
1752 (1971).

1753 147 Hoiijtink, H., Beland, S. & Vermeulen, J. A. Cognitive diagnostic assessment via Bayesian evaluation
1754 of informative diagnostic hypotheses. *Psychol Methods* **19**, 21-38, doi:10.1037/a0034176 (2014).

1755 148 Lee, M. D. How cognitive modeling can benefit from hierarchical Bayesian models. *Journal of*
1756 *Mathematical Psychology* **55**, 1-7 (2011).

1757 149 Depaoli, S., Rus, H. M., Clifton, J. P., van de Schoot, R. & Tiemensma, J. An introduction to Bayesian
1758 statistics in health psychology. *Health Psychol Rev* **11**, 248-264,
1759 doi:10.1080/17437199.2017.1343676 (2017).

1760 150 Kruschke, J. K. Bayesian estimation supersedes the t test. *J Exp Psychol Gen* **142**, 573-603,
1761 doi:10.1037/a0029146 (2013).

1762 151 Bolt, D. M., Piper, M. E., Theobald, W. E. & Baker, T. B. Why two smoking cessation agents work
1763 better than one: Role of craving suppression. *Journal of Consulting and Clinical Psychology* **80**, 54-
1764 65 (2012).

1765 152 Billari, F. C., Graziani, R. & Melilli, E. Stochastic population forecasting based on combinations of
1766 expert evaluations within the Bayesian paradigm. *Demography* **51**, 1933-1954 (2014).

1767 153 Fallesen, P. & Breen, R. Temporary Life Changes and the Timing of Divorce. *Demography* **53**, 1377-
1768 1398, doi:10.1007/s13524-016-0498-2 (2016).

1769 154 Hansford, T. G., Depaoli, S. & Canelo, K. S. Locating U.S. Solicitors General in the Supreme Court_s
1770 policy space. *Presidential Studies Quarterly* **49**, 855-869 (2019).

1771 155 Phipps, D. J., Hagger, M. S. & Hamilton, K. Predicting limiting_free sugar_consumption using an
1772 integrated model of health behavior. *Appetite* (2020).

- 1773 156 Royle, J. & Dorazio, R. Hierarchical Modeling and Inference in Ecology.,(Academic Press:
1774 Amsterdam.). (2008).
- 1775 157 Gimenez, O. *et al.* in *Modeling demographic processes in marked populations* 883-915 (Springer,
1776 2009).
- 1777 158 King, R., Morgan, B., Gimenez, O. & Brooks, S. P. *Bayesian analysis for population ecology*. (CRC
1778 press, 2009).
- 1779 159 Kéry, M. & Schaub, M. *Bayesian population analysis using WinBUGS: a hierarchical perspective*.
1780 (Academic Press, 2011).
- 1781 160 McCarthy, M. (New York, New York: Cambridge University Press, 2012).
- 1782 161 Korner-Nievergelt, F. *et al.* *Bayesian data analysis in ecology using linear models with R, BUGS,
1783 and Stan*. (Academic Press, 2015).
- 1784 162 Monnahan, C. C., Thorson, J. T. & Branch, T. A. Faster estimation of Bayesian models in ecology
1785 using Hamiltonian Monte Carlo. *Methods in Ecology and Evolution* **8**, 339-348 (2017).
- 1786 163 Ellison, A. M. Bayesian inference in ecology. *Ecology letters* **7**, 509-520 (2004).
- 1787 164 Choy, S. L., O'Leary, R. & Mengersen, K. Elicitation by design in ecology: using expert opinion to
1788 inform priors for Bayesian statistical models. *Ecology* **90**, 265-277 (2009).
- 1789 165 Kuhnert, P. M., Martin, T. G. & Griffiths, S. P. A guide to eliciting and using expert knowledge in
1790 Bayesian ecological models. *Ecology letters* **13**, 900-914 (2010).
- 1791 166 King, R., Brooks, S. P., Mazzetta, C., Freeman, S. N. & Morgan, B. J. Identifying and diagnosing
1792 population declines: a Bayesian assessment of lapwings in the UK. *Journal of the Royal Statistical
1793 Society: Series C (Applied Statistics)* **57**, 609-632 (2008).
- 1794 167 Newman, K. *et al.* *Modelling population dynamics*. (Springer, 2014).
- 1795 168 Bachl, F. E., Lindgren, F., Borchers, D. L. & Illian, J. B. inlabru: an R package for Bayesian spatial
1796 modelling from ecological survey data. *Methods in Ecology and Evolution* **10**, 760-766 (2019).
- 1797 169 King, R. & Brooks, S. P. On the Bayesian estimation of a closed population size in the presence of
1798 heterogeneity and model uncertainty. *Biometrics* **64**, 816-824, doi:10.1111/j.1541-
1799 0420.2007.00938.x (2008).
- 1800 170 Saunders, S. P., Cuthbert, F. J. & Zipkin, E. F. Evaluating population viability and efficacy of
1801 conservation management using integrated population models. *Journal of Applied Ecology* **55**,
1802 1380-1392 (2018).
- 1803 171 McClintock, B. T. *et al.* A general discrete-time modeling framework for animal movement using
1804 multistate random walks. *Ecological Monographs* **82**, 335-349 (2012).
- 1805 172 Dennis, B., Ponciano, J. M., Lele, S. R., Taper, M. L. & Staples, D. F. Estimating density dependence,
1806 process noise, and observation error. *Ecological Monographs* **76**, 323-341 (2006).
- 1807 173 Aeberhard, W. H., Mills Flemming, J. & Nielsen, A. Review of state-space models for fisheries
1808 science. *Annual Review of Statistics and Its Application* **5**, 215-235 (2018).
- 1809 174 Isaac, N. J. B. *et al.* Data Integration for Large-Scale Models of Species Distributions. *Trends Ecol
1810 Evol* **35**, 56-67, doi:10.1016/j.tree.2019.08.006 (2020).
- 1811 175 McClintock, B. T. *et al.* Uncovering ecological state dynamics with hidden Markov models. *arXiv
1812 preprint arXiv:2002.10497* (2020).
- 1813 176 King, R. Statistical ecology. *Annual Review of Statistics and its Application* **1**, 401-426 (2014).
- 1814 177 Fearnhead, P. MCMC for state-space models. (2011).
- 1815 178 Andrieu, C., Doucet, A. & Holenstein, R. Particle markov chain monte carlo methods. *Journal of
1816 the Royal Statistical Society: Series B (Statistical Methodology)* **72**, 269-342 (2010).
- 1817 179 Knappe, J. & de Valpine, P. Fitting complex population models by combining particle filters with
1818 Markov chain Monte Carlo. *Ecology* **93**, 256-263, doi:10.1890/11-0797.1 (2012).

1819 180 Finke, A., King, R., Beskos, A. & Dellaportas, P. Efficient sequential Monte Carlo algorithms for
1820 integrated population models. *Journal of Agricultural, Biological and Environmental Statistics* **24**,
1821 204-224 (2019).

1822 181 Stephens, M. & Balding, D. J. Bayesian statistical methods for genetic association studies. *Nat Rev*
1823 *Genet* **10**, 681-690, doi:10.1038/nrg2615 (2009).

1824 182 Mimno, D., Blei, D. M. & Engelhardt, B. E. Posterior predictive checks to quantify lack-of-fit in
1825 admixture models of latent population structure. *Proc Natl Acad Sci U S A* **112**, E3441-3450,
1826 doi:10.1073/pnas.1412301112 (2015).

1827 183 Schaid, D. J., Chen, W. & Larson, N. B. From genome-wide associations to candidate causal variants
1828 by statistical fine-mapping. *Nat Rev Genet* **19**, 491-504, doi:10.1038/s41576-018-0016-z (2018).

1829 184 Marchini, J. & Howie, B. Genotype imputation for genome-wide association studies. *Nat Rev*
1830 *Genet* **11**, 499-511, doi:10.1038/nrg2796 (2010).

1831 185 Allen, N. E., Sudlow, C., Peakman, T., Collins, R. & Biobank, U. K. UK biobank data: come and get
1832 it. *Sci Transl Med* **6**, 224ed224, doi:10.1126/scitranslmed.3008601 (2014).

1833 186 Cortes, A. *et al.* Bayesian analysis of genetic association across tree-structured routine healthcare
1834 data in the UK Biobank. *Nat Genet* **49**, 1311-1318, doi:10.1038/ng.3926 (2017).

1835 187 Argelaguet, R. *et al.* Multi-Omics Factor Analysis-a framework for unsupervised integration of
1836 multi-omics data sets. *Mol Syst Biol* **14**, e8124, doi:10.15252/msb.20178124 (2018).

1837 188 Stuart, T. & Satija, R. Integrative single-cell analysis. *Nat Rev Genet* **20**, 257-272,
1838 doi:10.1038/s41576-019-0093-7 (2019).

1839 189 Yau, C. & Campbell, K. Bayesian statistical learning for big data biology. *Biophys Rev* **11**, 95-102,
1840 doi:10.1007/s12551-019-00499-1 (2019).

1841 190 Vallejos, C. A., Marioni, J. C. & Richardson, S. BASiCS: Bayesian Analysis of Single-Cell Sequencing
1842 Data. *PLoS Comput Biol* **11**, e1004333, doi:10.1371/journal.pcbi.1004333 (2015).

1843 191 Wang, J. *et al.* Data denoising with transfer learning in single-cell transcriptomics. *Nat Methods*
1844 **16**, 875-878, doi:10.1038/s41592-019-0537-1 (2019).

1845 192 Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for single-cell
1846 transcriptomics. *Nat Methods* **15**, 1053-1058, doi:10.1038/s41592-018-0229-2 (2018).

1847 193 Institute, N. C. & National Cancer, I. The Cancer Genome Atlas. *Definitions*, doi:10.32388/e1plqh
1848 (2020).

1849 194 Kuipers, J. *et al.* Mutational interactions define novel cancer subgroups. *Nat Commun* **9**, 4353,
1850 doi:10.1038/s41467-018-06867-x (2018).

1851 195 Schwartz, R. & Schaffer, A. A. The evolution of tumour phylogenetics: principles and practice. *Nat*
1852 *Rev Genet* **18**, 213-229, doi:10.1038/nrg.2016.170 (2017).

1853 196 Munafò, M. R. *et al.* A manifesto for reproducible science. *Nature Human Behaviour* **1**,
1854 doi:10.1038/s41562-016-0021 (2017).

1855 197 van Erp, S., Mulder, J. & Oberski, D. L. Prior sensitivity analysis in default Bayesian structural
1856 equation modeling. *Psychol Methods* **23**, 363-388, doi:10.1037/met0000162 (2018).

1857 198 Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and
1858 stewardship. *Sci Data* **3**, 160018, doi:10.1038/sdata.2016.18 (2016).

1859 199 Lamprecht, A.-L. *et al.* Towards FAIR principles for research software. *Data Science* **3**, 37-59,
1860 doi:10.3233/ds-190026 (2020).

1861 200 *re3data.org - Registry of Research Data Repositories.*

1862 201 Smith, A. M., Katz, D. S. & Niemeyer, K. E. Software citation principles. *PeerJ Computer Science* **2**,
1863 e86, doi:10.7717/peerj-cs.86 (2016).

1864 202 Clyburne-Sherin, A., Fei, X. & Green, S. A. Computational Reproducibility via Containers in
1865 Psychology. *Meta-Psychology* **3** (2019).

1866 203 Nosek, B. A. *et al.* SCIENTIFIC STANDARDS. Promoting an open research culture. *Science* **348**, 1422-
1867 1425, doi:10.1126/science.aab2374 (2015).

1868 204 Vehtari, A. & Ojanen, J. A survey of Bayesian predictive methods for model assessment, selection
1869 and comparison. *Statistics Surveys* **6**, 142-228 (2012).

1870 205 Abadi, M. *et al.* in *USENIX symposium on operating systems design and implementation*
1871 *(OSDI'16)*. 12 edn 265-283.

1872 206 Paszke, A. *et al.* in *Advances in neural information processing systems*. 8026-8037.

1873 207 Kingma, D. P. & Welling, M. *An Introduction to Variational Autoencoders*. (2019).

1874 208 Higgins, I. *et al.* beta-VAE: Learning Basic Visual Concepts with a Constrained Variational
1875 Framework. *Iclr* **2**, 6 (2017).

1876 209 Märtens, K. & Yau, C. BasisVAE: Translation-invariant feature-level clustering with Variational
1877 Autoencoders. *arXiv [stat.ML]* (2020).

1878 210 Liu, Q., Allamanis, M., Brockschmidt, M. & Gaunt, A. in *Advances in Neural Information Processing*
1879 *Systems 31* (eds S. Bengio *et al.*) 7795-7804 (Curran Associates, Inc., 2018).

1880 211 Louizos, C., Shi, X., Schutte, K. & Welling, M. in *Advances in Neural Information Processing Systems*
1881 8743-8754 (2019).

1882 212 Garnelo, M. *et al.* in *Proceedings of the 35th International Conference on Machine Learning Vol.*
1883 *80* (eds Jennifer Dy & Andreas Krause) 1704-1713 (PMLR, 2018).

1884 213 Kim, H. *et al.* Attentive Neural Processes. *arXiv [cs.LG]* (2019).

1885 214 Rezende, D. & Mohamed, S. in *Proceedings of the 32nd International Conference on Machine*
1886 *Learning Vol. 37* (eds Francis Bach & David Blei) 1530-1538 (PMLR, 2015).

1887 215 Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S. & Lakshminarayanan, B. Normalizing
1888 Flows for Probabilistic Modeling and Inference. *arXiv [stat.ML]* (2019).

1889 216 Korshunova, I. *et al.* in *Advances in Neural Information Processing Systems 31* (eds S. Bengio *et*
1890 *al.*) 7190-7198 (Curran Associates, Inc., 2018).

1891 217 Zhang, R., Li, C., Zhang, J., Chen, C. & Wilson, A. G. Cyclical stochastic gradient mcmc for bayesian
1892 deep learning. *arXiv preprint arXiv:1902. 03932* (2019).

1893 218 Neal, R. M. *Bayesian Learning for Neural Networks*. (Springer Science & Business Media, 2012).

1894 219 Neal, R. M. in *Bayesian Learning for Neural Networks Lecture Notes in Statistics* (ed Radford M.
1895 Neal) Ch. Chapter 2, 29-53 (Springer New York, 1996).

1896 220 Williams, C. K. I. in *Advances in neural information processing systems* 295-301 (1997).

1897 221 MacKay David, J. C. A practical bayesian framework for backprop networks. *Neural Comput.*
1898 (1992).

1899 222 Sun, S., Zhang, G., Shi, J. & Grosse, R. in *International Conference on Learning Representations*
1900 (2019).

1901 223 Lakshminarayanan, B., Pritzel, A. & Blundell, C. in *Advances in neural information processing*
1902 *systems* 6402-6413 (2017).

1903 224 Wilson, A. G. The Case for Bayesian Deep Learning. *arXiv [cs.LG]* (2020).

1904 225 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: a simple way
1905 to prevent neural networks from overfitting. *The journal of machine learning research* **15**, 1929-
1906 1958 (2014).

1907 226 Gal, Y. & Ghahramani, Z. 1050-1059.

1908 227 Green, P. J. Reversible jump Markov chain Monte Carlo computation and Bayesian model
1909 determination. *Biometrika* **82**, 711-732 (1995).

1910 228 Hoffman, M. D. & Gelman, A. The No-U-Turn sampler: adaptively setting path lengths in
1911 Hamiltonian Monte Carlo. *J. Mach. Learn. Res.* **15**, 1593-1623 (2014).

1912 229 Liang, F. & Wong, W. H. Evolutionary Monte Carlo: applications to C p model sampling and change
1913 point problem. *Statistica sinica*, 317-342 (2000).

- 1914 230 Liu, J. S. & Chen, R. Sequential Monte Carlo methods for dynamic systems. *Journal of the American*
1915 *statistical association* **93**, 1032-1044 (1998).
- 1916 231 Sisson, S., Fan, Y. & Beaumont, M. *Handbook of approximate Bayesian computation*. (Chapman
1917 and Hall/CRC 2018).
- 1918 232 Rue, H., Martino, S. & Chopin, N. Approximate Bayesian inference for latent Gaussian models by
1919 using integrated nested Laplace approximations. *Journal of the royal statistical society: Series B*
1920 *(Statistical Methodology)* **71**, 319-392 (2009).
- 1921 233 WinBUGS user manual (Citeseer, 2003).
- 1922 234 Ntzoufras, I. *Bayesian modeling using WinBUGS*. Vol. 698 (John Wiley & Sons, 2011).
- 1923 235 OpenBUGS user manual, version 3.0. 2 (2007).
- 1924 236 Plummer, M. in *Proceedings of the 3rd international workshop on distributed statistical*
1925 *computing*. 1-10.
- 1926 237 Goudie, R. J., Turner, R. M., De Angelis, D. & Thomas, A. MultiBUGS: A parallel implementation of
1927 the BUGS modelling framework for faster Bayesian inference. *arXiv preprint arXiv:1704.03216*
1928 (2017).
- 1929 238 Sturtz, S., Ligges, U. & Gelman, A. E. R2WinBUGS: a package for running WinBUGS from R. *Journal*
1930 *of Statistical Software* **12**, 1-16 (2005).
- 1931 239 Thomas, A., O'Hara, B., Ligges, U. & Sturtz, S. Making BUGS Open. *R News* **6** (2006).
- 1932 240 Salvatier, J., Wiecki, T. V. & Fonnesbeck, C. Probabilistic programming in Python using PyMC3.
1933 *PeerJ Computer Science* **2**, e55, doi:10.7717/peerj-cs.55 (2016).
- 1934 241 de Valpine, P. *et al.* Programming with models: writing statistical algorithms for general model
1935 structures with NIMBLE. *Journal of Computational and Graphical Statistics* **26**, 403-413 (2017).
- 1936 242 Dillon, J. V. *et al.* Tensorflow distributions. *arXiv preprint arXiv:1711.10604* (2017).
- 1937 243 Keydana, S. *tfprobability: R interface to TensorFlow Probability*,
1938 <https://rstudio.github.io/tfprobability/index.html> (2020).
- 1939 244 Bingham, E. *et al.* Pyro: Deep universal probabilistic programming. *The Journal of Machine*
1940 *Learning Research* **20**, 973-978 (2019).
- 1941 245 Bezanson, J., Karpinski, S., Shah, V. B. & Edelman, A. Julia: A fast dynamic language for technical
1942 computing. *arXiv preprint arXiv:1209.5145* (2012).
- 1943 246 Ge, H., Xu, K., Scibior, A. & Ghahramani, Z. in *Artificial Intelligence and Statistics*.
- 1944 247 Smith, B. J. *et al.* brian-j-smith/Mamba.jl: v0.12.4. *Zenodo*,
1945 doi:<http://doi.org/10.5281/zenodo.3740216> (2020).
- 1946 248 JASP Team. JASP (Version 0.13.1)[Computer software]. (2020).
- 1947 249 Lindgren, F. & Rue, H. Bayesian spatial modelling with R-INLA. *Journal of Statistical Software* **63**,
1948 1-25 (2015).
- 1949 250 Vanhatalo, J. *et al.* GPstuff: Bayesian modeling with Gaussian processes. *Journal of Machine*
1950 *Learning Research* **14**, 1175-1179 (2013).
- 1951 251 Betancourt, M. *Towards A Principled Bayesian Workflow*,
1952 https://betanalpha.github.io/assets/case_studies/principled_bayesian_workflow.html (April
1953 2020).
- 1954 252 Kramer, B. & Bosman, J. Summerschool Open Science and Scholarship 2019 - Utrecht University.
1955 doi:10.5281/ZENODO.3925004 (2020).
- 1956 253 Rényi, A. On a new axiomatic theory of probability. *Acta Mathematica Hungarica* **6**, 285-335
1957 (1955).
- 1958 254 Lesaffre, E. & Lawson, A. B. *Bayesian biostatistics*. (John Wiley & Sons, 2012).

1959