

Applying Bayesian Analysis Guidelines to Empirical Software Engineering Data

The Case of Programming Languages and Code Quality

CARLO A. FURIA, Software Institute, USI Università della Svizzera italiana, Switzerland

RICHARD TORKAR, University of Gothenburg, Sweden and Stellenbosch Institute for Advanced Study (STIAS), South Africa

ROBERT FELDT, Chalmers and University of Gothenburg, Sweden

Statistical analysis is the tool of choice to turn data into information, and then information into empirical knowledge. The process that goes from data to knowledge is, however, long, uncertain, and riddled with pitfalls. To be valid, it should be supported by detailed, rigorous guidelines, which help ferret out issues with the data or model, and lead to qualified results that strike a reasonable balance between generality and practical relevance. Such guidelines are being developed by statisticians to support the latest techniques for *Bayesian* data analysis. In this article, we frame these guidelines in a way that is apt to empirical research in software engineering.

To demonstrate the guidelines in practice, we apply them to reanalyze a GitHub dataset about code quality in different programming languages. The dataset's original analysis (Ray et al., 2014) and a critical reanalysis (Berger et al., 2019) have attracted considerable attention—in no small part because they target a topic (the impact of different programming languages) on which strong opinions abound. The goals of our reanalysis are largely orthogonal to this previous work, as we are concerned with demonstrating, on data in an interesting domain, how to build a principled Bayesian data analysis and to showcase its benefits. In the process, we will also shed light on some critical aspects of the analyzed data and of the relationship between programming languages and code quality—such as the impact of project-specific characteristics other than the used programming language.

The high-level conclusions of our exercise will be that Bayesian statistical techniques can be applied to analyze software engineering data in a way that is principled, flexible, and leads to convincing results that inform the state of the art while highlighting the boundaries of its validity. The guidelines can support building solid statistical analyses and connecting their results, and hence help buttress continued progress in empirical software engineering research.

CCS Concepts: • **Mathematics of computing** → **Bayesian computation**; • **Software and its engineering** → **Empirical software validation**.

Additional Key Words and Phrases: Bayesian data analysis, statistical analysis, guidelines, empirical software engineering, programming languages

ACM Reference Format:

Carlo A. Furia, Richard Torkar, and Robert Feldt. 2022. Applying Bayesian Analysis Guidelines to Empirical Software Engineering Data: The Case of Programming Languages and Code Quality. *ACM Trans. Softw. Eng. Methodol.* 1, 1 (January 2022), 40 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Empirical disciplines, including a substantial part of software engineering research, mine data for information, and then use the information as evidence to build, extend, and refine empirical knowledge. Statistical analysis is key to implementing this process; but statistical techniques are just tools, which need detailed *guidelines* to be applied properly

© 2022

1049-331X/2022/1-ART

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

50 and consistently. It is only through the combination of powerful statistical techniques and
51 rigorous guidelines to apply them that we can distill empirical knowledge following a
52 process that is consistent, rests on solid principles, and ultimately is more likely to lead to
53 valid results with a higher degree of confidence.

54 Whereas frequentist statistical techniques have been commonplace in science for over a
55 century—since the influential work of the likes of Pearson [68] and Fisher [20]—the state
56 of the art in applied statistics is moving towards using *Bayesian* analysis techniques. As
57 we discussed in previous work [21], recent developments in Bayesian analysis techniques
58 (such as using Hamiltonian Monte Carlo fitting algorithms [11]) coupled with an increasing
59 availability of the computing power needed to run them on large datasets have convincingly
60 demonstrated the advantages of using Bayesian statistics and the flexibility and rigor of
61 the analysis they support. More recently, applied statisticians have also been working
62 out practical *guidelines* that can boost usability and impact of Bayesian statistical data
63 analysis [1, 23, 31, 57]. In this paper, we present some of these guidelines and frame them
64 in a way that is suitable for empirical research in the software engineering domain—with
65 the goal of demonstrating how they can support a principled way of building statistical
66 analyses of software engineering data.

67 To demonstrate the guidelines in practice, we follow them to analyze a large dataset
68 about the code quality of projects written in disparate programming languages and hosted
69 on GitHub [55]. The empirical study that curated this dataset and performed the original
70 analysis [55] was followed by a critical reanalysis by a different group of researchers [6]; as
71 we recall in Section 1.1, the topic has received much attention and stirred some controversy.
72 This visibility makes the dataset an attractive target for our own purposes.

73 In the paper, we go through various aspects of the data analysis performed in the previous
74 studies [6, 55], illustrating the versatile features of Bayesian statistical models in practice.
75 We demonstrate how the guidelines support an incremental and iterative analysis process,
76 where several key features of a statistical model can be validated; this, in turn, encourages
77 trying out different models and comparing them in a rigorous way—as opposed to blindly
78 relying on one-size-fits-all rules of thumb. Following this process, we demonstrate that
79 some issues of the original analysis [55] or criticized by the follow-up reanalysis [6] could
80 have been identified more easily. Furthermore, the limitations and actual impact of previous
81 studies could have been framed more straightforwardly and more transparently. The
82 conclusion of our exercise will be that flexible statistical techniques coupled with principled
83 and structured guidelines can help address empirical research questions directly and
84 transparently. This can lead to explanations that are nuanced and detailed, and hence,
85 ultimately, that can become convincing foundations for building shared knowledge.

86

87 1.1 Dataset and previous studies

88 In this paper, we reuse the dataset collected and analyzed by Ray et al. in a paper published
89 at the FSE¹ conference [55]. A critical reproduction [6] of the original study, written by
90 Berger et al. and published in the TOPLAS² journal, triggered a prolonged controversy that
91 reverberated on social media.

92

93

94

95

96

97

98

¹The ACM SIGSOFT International Symposium on Foundations of Software Engineering.

²The ACM Transactions on Programming Languages and Systems.

99 Here is the story so far, in the shortest terms possible:³ based on their analysis of projects
100 hosted by GitHub, the original study [55] (henceforth, “FSE”) claimed to have found an
101 association between certain programming languages and the bug proneness of code written
102 in them. The reproduction study [6] (henceforth, “TOPLAS”) criticized several aspects of
103 FSE—most prominently, its data collection and classification practices—and questioned the
104 soundness of some of its results. In a rebuttal, the authors of FSE defended their results;⁴
105 and in a rebuttal of the rebuttal the authors of TOPLAS maintained their criticism.⁵

106 Our paper is *emphatically not* our attempt to jump into the fray: we do not have much to
107 add to the subject matter of the controversy. However, we appreciate the interest that the
108 controversial topic received, and see it as an opportunity to present our views on a different,
109 but related, aspect: practices in statistical analysis. Both FSE and TOPLAS primarily use
110 *frequentist* statistical techniques. Even in the best conditions, these techniques’ flexibility is
111 limited in comparison to the *Bayesian* statistical techniques we have been advocating [21, 65].

112 Overall, our contributions fall largely outside the focus of FSE and TOPLAS—except to
113 the extent that they target the same domain and the same data. The core of both papers
114 revolves around GitHub data, how it was collected and processed, and how the variables of
115 interest have been operationalized. TOPLAS’s main goal was to attempt to reproduce FSE’s
116 results, and hence it deliberately makes mostly limited changes to the statistical models.
117 Our analysis takes the data as it was collected and made available by FSE’s original study,
118 and tries to make the most out of it following rigorous guidelines to apply flexible statistical
119 practices—Bayesian statistics, that is.

120

121 1.2 Overview

122 The overall goal of this paper is demonstrating how Bayesian statistical techniques can
123 be applied in a principled way to build suitable statistical models. These models can then
124 be used to answer research questions in a flexible⁶ way, and to quantify limitations and
125 uncertainties about what one can reliably infer from the models.

126 This complements our earlier work on Bayesian data analysis for empirical software
127 engineering [21, 65], which:

- 128 • Argued for using Bayesian over frequentist statistics and showcased the former’s
129 flexibility on software engineering data [21].
- 130 • Suggested to analyze *practical* significance using a combination of Bayesian statistics
131 and cumulative prospect theory, which helps stakeholders evaluate the impact of a
132 technique or a practice in a way that takes into account their constraints, available
133 resources, and intuitive reasoning [65].

134

135 *1.2.1 Key benefits of Bayesian statistics.* Before we go into the novel contributions of the
136 present paper, let us briefly summarize the benefits of Bayesian statistical techniques—which
137 we presented in detail in our previous work [21, 65]. Section 4 will further demonstrate
138 several of these benefits on the programming language case study.

139 There is a growing awareness in several empirical scientific disciplines that “classical
140 [frequentist] statistical tools are not diverse enough to handle many common research

141

142 ³Hillel Wayne provides a much more detailed account <https://www.hillelwayne.com/post/this-is-how-science-happens/>.

143 ⁴<https://arxiv.org/abs/1911.07393>

144 ⁵<http://janvitek.org/var/rebuttal-rebuttal.pdf>

145 ⁶Here, “flexible” means that it can be adapted to different kinds of scientific questions and analysis domains
146 while remaining effective.

147

148 questions” [42]. Bayesian statistics, in contrast, are much more flexible, as they provide
149 general methods to connect data, models, and research questions [27]. Bayesian statistics
150 are more flexible because they are centered around *modeling*: how we fit a Bayesian model
151 is largely independent of the details of how the model was built. In contrast, different
152 frequentist models often require widely different analysis procedures: if we need to tweak
153 the model or its underlying assumptions even slightly, the frequentist analysis results may
154 become unreliable.

155 The main output of a Bayesian data analysis is a *distribution* of model parameters fitted
156 on the data. This includes rich quantitative information, in contrast to the point estimates
157 that are the usual outcome of applying frequentist statistics. Providing distributional infor-
158 mation is a key strength of Bayesian statistics. First, it supports quantitative and nuanced
159 analyses instead of a purely dichotomous (yes/no) view—which is prevalent with sta-
160 tistical hypothesis testing (a core technique of frequentist statistics that has been under
161 intense scrutiny [2, 70]). Second, distributional information is easier to understand, since
162 it measures quantities of interest in the specific domain. Contrast this to purely statistical
163 metrics such as p -values or confidence intervals, which are notoriously hard to interpret
164 correctly [34, 35]. Third, the quantitative distributional information that is provided by a
165 fitted Bayesian model supports simulating the *derived* distributions of a variety of quantities
166 of interest, such as outcomes in a specific scenario—which is especially useful to analyze
167 practical significance [65].

168 The current paper focuses on how to apply Bayesian data analysis in a principled way:
169 following detailed guidelines and a structured workflow. Demonstrating the guidelines on
170 FSE’s dataset, our contributions address four aspects that are relevant to every study that
171 involves statistical data analysis.

172 *1.2.2 How to design a statistical model?* Modeling requires to exercise *judgement*—something
173 that can be based on practices, customs, and heuristics, but is not completely reducible to a
174 fixed set of rigid rules. Bayesian statistics emphasizes the *modeling* aspect of data analysis,
175 and provides quantitative techniques to help ground heuristics and practices onto a robust
176 and sound statistical framework.

177 Section 2 presents guidelines to build a Bayesian statistical model incrementally (adding
178 features as needed), iteratively (improving a model based on the shortcomings of the
179 previous ones), and rigorously (with quantitative criteria to assess a model’s suitability).
180 Our guidelines customize general guidelines developed by the Bayesian data analysis
181 community to the scenarios that are common in empirical software engineering. Section 3
182 demonstrates, on the FSE dataset, that our guidelines provide principled ways of assessing
183 the strengths and weaknesses of any statistical model for the analysis at hand.

184 *1.2.3 How to spot data problems?* TOPLAS’s criticism of FSE’s analysis questions the accuracy
185 of some of the *data* that was collected and how it was processed. For example, it says that
186 “project size, computed in the FSE paper as the sum of inserted lines, is not accurate—
187 as it does not take deletions into account” [6, §3.2]. Can Bayesian statistical techniques
188 help discover problems with the data—such as inconsistencies, sparseness, and lack of
189 homogeneity—that limit the validity and generalizability of the statistical analysis’s results?

190 Naturally, no statistical technique (no matter how powerful) can supersede a careful
191 analysis of construct validity [19, 53], which should precede the statistical analysis and lay
192 the foundations for it. Still, applying the Bayesian guidelines that we present can ferret out
193 issues with the data and highlight where uncertainty is more or less pronounced, so that
194 we can heed any limitations when drawing conclusions. For example, Section 4.2 finds that
195

196

197 the number of inserted lines performs poorly as a predictor, echoing TOPLAS’s observation
198 that it may not be a suitable measure of size.

199
200 *1.2.4 How to assess significant results?* In previous work [21], we demonstrated that Bayesian
201 statistical techniques can help move away from a dichotomous (significant/not significant)
202 framing of research questions—which comes typically with frequentist null hypothesis
203 testing and is often artificially restrictive—and instead focus on *practical* significance [65].

204 The gap between statistical significance and practical significance is more likely to be
205 wide when studying complex domains with plenty of confounding factors. The analysis of
206 programming language data is a clear example of such complex domains. In this paper, we
207 show how the Bayesian statistics guidelines support a nuanced analysis of complex models,
208 and help keep the focus on concrete scenarios and practically relevant measures. Concretely,
209 we show that Bayesian data analysis provides a flexible model of data distributions, which
210 can be used to *predict* outcomes in different scenarios directly in terms of statistics that are
211 based on variables in the problem domain.

212 Our analysis’s conclusion will be that the key question “which programming languages
213 are more fault prone” does not admit a simple straightforward answer—not with the ana-
214 lyzed data at least. Nevertheless, as we argued in [21] and now demonstrate in Section 4.3,
215 practitioners can ask specific questions and answer them by running simulations on the
216 Bayesian model, rather than having to rely on general results that may not be meaningful
217 in their context.

218 *1.2.5 How to build knowledge incrementally?* Every empirical study has limitations; lifting
219 them requires to perform new experiments. Another advantage of Bayesian statistical
220 models built using an incremental process is that they can be refined as we collect more data.
221 This way, our models become better over time since they accurately reflect the evolving
222 scientific knowledge in a certain area.

223 Section 4.4 discusses how applying Bayesian analysis guidelines helps plan for additional
224 data collection based on the limitations of the analyzed data. Different experiments are no
225 longer merely a loose collection around the same themes, but can be planned back-to-back
226 in a way that progressively reduces the uncertainty in knowledge.
227

228 1.3 Contributions

229 This paper makes the following contributions:
230

- 231 • It presents guidelines to apply Bayesian statistics following a systematic process that
232 goes from building and validating the model to fitting and analyzing it.
- 233 • It demonstrates the guidelines by showing how to incrementally build a suitable
234 statistical model to capture FSE’s language quality data.
- 235 • It analyzes the fitted model to investigate the original questions of the effect of
236 programming languages on fault proneness with a focus on practical scenarios.
- 237 • For reproducibility, all analysis scripts are available online together with additional
238 results and detailed data visualization:

239 REPLICATION PACKAGE: <https://doi.org/10.5281/zenodo.4472963> [22].
240

241 *1.3.1 Scope.* To a large degree, the guidelines we present are *not specific* to certain classes
242 of statistical models or analysis domains. The case study we detail in this paper is about
243 programming language quality, which we model using several generalized linear models
244 of different complexity—a broad class of statistical models widely used for their flexibility.
245

This does not mean that the guidelines are only applicable to programming language data, nor that they only work for generalized linear models.

Since the guidelines are largely independent of the specific features of the chosen statistical model, the domain-specific details of how to operationalize a certain data analysis problem and how to build a valid “construct” (a statistical model) are largely outside the scope of the present paper. As we remarked above, we selected the programming language data as case study because it has been already thoroughly analyzed and scrutinized (albeit in a frequentist setting); thus, we can build on FSE’s and TOPLAS’s work to demonstrate the additional steps to be taken to bolster the validity of a statistical data analysis. This leaves room for *different* analyses of the same research questions but using different data collection processes or different statistical models. Our guidelines remain valid as a safeguard against modeling mistakes or shortcomings; since they promote an *iterative* approach, they can also suggest what to change when they fail to validate a candidate model.

1.3.2 Organization. The rest of the paper is organized as follows. Section 2 illustrates Bayesian data analysis guidelines with an angle that is relevant for empirical software engineering. Section 3 follows the guidelines to incrementally build a model that is suitable to capture FSE’s programming language data. Various models are rigorously evaluated and compared, so that the final model is arguably the “best” among them according to certain quantitative criteria. Section 4 analyzes the fitted model to study the original questions of which programming languages are associated with more or fewer faults. The results look at different scenarios and outline how further custom analyses could be built atop the same model. Finally, Section 5 discusses related work and Section 6 concludes with a brief summary and closing discussion.

2 BAYESIAN DATA ANALYSIS GUIDELINES

In the last decade, powerful Bayesian statistical analysis tools and languages have become widely available together with computational resources adequate to run them [13, 25, 50]. More recently, statisticians have also been introducing and refining guidelines on how to use these tools in a systematic way to perform principled Bayesian data modeling [1, 23, 31, 57]. In this section, we summarize these state-of-the-art guidelines while recasting them in a form suitable for empirical software engineering research.

A Bayesian model defines a statistical data-generating process in terms of a *prior* distribution of parameters θ and a *likelihood* that certain data is observed for each value of the parameters. *Fitting* such a model on some empirical data D then gives a *posterior* distribution of the same parameters that follows Bayes’ theorem:

$$\underbrace{P(\theta | D)}_{\text{posterior}} \propto \underbrace{P(D | \theta)}_{\text{likelihood}} \times \underbrace{P(\theta)}_{\text{prior}}. \quad (1)$$

The posterior can then be used to compute the probability of other observations of interest in a predictive fashion. Our previous work [21] presented more details about Bayes’ theorem and the roles of prior, likelihood, and posterior. In this paper, we focus on how to build a Bayesian model in practice: Bayesian modeling involves choosing components in a way that is sound and principled, and that works for the data and domain that we are targeting.

Figure 1 illustrates a key idea of the *guidelines* for Bayesian analysis presented here: developing a statistical model is a process of *iterative refinement*, which starts from a very simple (possibly simplistic) *initial model* that is gradually refined. Each iteration goes through a series of *steps* that assess the model’s suitability in terms of the following characteristics:

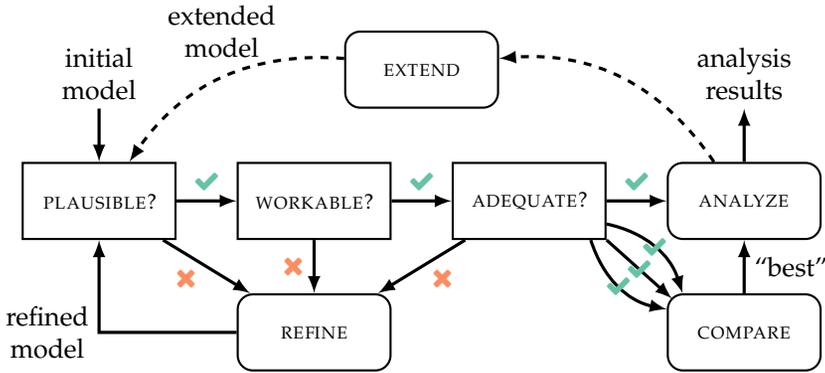


Fig. 1. Process for Bayesian data analysis: starting from an initial model, assess whether it is *plausible*, *workable*, and *adequate*. If it lacks any of these characteristics, *refine* the model by adding detail and features. Models that pass all checks can be fitted and used to answer the analysis’s specific questions. Different models that pass all checks can be rigorously *compared* to select those that perform “best” according to suitable criteria. The outer loop (dashed arrows) indicates that an analysis’s results may also suggest to *extend* an adequate model so that it can answer more precise, or just different, questions; this outer loop is another source of multiple models that can be compared.

- Plausibility:** Is the model consistent with (expert) knowledge about the data domain?
- Workability:** Can the model effectively and accurately be fitted using the available numerical algorithms?
- Adequacy:** Can the model capture the characteristics of the empirical data?

These steps help assess the *utility* (or *suitability*) of a model and its trade-offs. As we illustrate in Tables 1 and 2, each step puts additional requirements on a model, and checks whether the model is well-equipped to faithfully capture the observed data and to analyze it. Table 1 shows how each step broadens the scope of what model components are checked; in particular, the actual empirical data is only used in the *adequate* step, whereas the previous steps generate simulated data using priors and likelihood. Table 2 details how each step has a possible *outcome* (what it establishes about the model), which is supported by analysis *artifacts* that document the step. The following sections describe the steps, artifacts, and outcomes in some detail. Section 3 will apply the steps on the main case study of programming language data.

2.1 Modeling

To make the description concrete, and thus easier to follow, we illustrate what the steps compute on a toy problem: predicting an adult person’s height h in centimeters. To this end, we build this statistical model:

$$h \sim \text{Normal}(\mu, \sigma) \tag{2}$$

$$\mu \sim \text{Normal}(170, 50) \tag{3}$$

$$\sigma \sim \text{HalfCauchy}(0, 1) \tag{4}$$

The following paragraphs introduce its components one by one.

Parameters. A Bayesian statistical model consists of three components: parameters to estimate, likelihood, and priors. In our example, the model’s *parameters* θ are the mean μ and standard deviation σ of a person’s height. The data D records the value of outcome

344 variable h for several persons, which we can use to estimate μ and σ . There are no predictor
345 variables in this simplistic model, but otherwise these would also be recorded in the data
346 for every person.

347 **Likelihood.** The *likelihood* is a probability distribution of the data h given parameters μ
348 and σ . The simplest (yet extremely common) choice is a *normal* distribution, which encodes
349 no additional information about the data other than that it has a mean μ and a standard
350 deviation σ . This likelihood is defined by (2), which is in fact a probability distribution of h
351 given μ and σ .

352 **Priors.** Finally, we need *priors* for μ and σ . Specifying a prior means defining an initial
353 probability distribution for a parameter of the model—a probability distribution without
354 any dependency on the data.

355 The least informative priors are completely *flat* distributions, which assign the same
356 infinitesimal probability to any value of the parameter; this is the default behavior in
357 frequentist statistics. A flat prior for μ , for example, would be a uniform distribution with
358 support from $-\infty$ to $+\infty$. Flat priors are usually a poor choice: first, since they stretch a
359 probability distribution over an infinitely large support, they tend to generate infinitesimal
360 probabilities that may cause *numerical* rounding errors; second, a prior with no information
361 whatsoever about the realistic parameter domain is prone to *overfitting* the data. We can
362 see it clearly even in the simple example of estimating heights: a flat prior would give the
363 same a priori probability to height values -10 , 170 , and 10^9 , but only the second value is a
364 plausible human height!

365 A better choice are *weakly informative* priors, which still carry very little specific informa-
366 tion but perform much better than flat priors computationally and protect against overfitting
367 the data. As prior for μ , we select the normal distribution (3), with mean 170 and standard
368 deviation 50 ; this means that we expect most heights to be between $20 = 170 - 3 \cdot 50$ and
369 $320 = 170 + 3 \cdot 50$ centimeters. This is still an extremely broad range of values, but it favors
370 values that are in the ballpark of realistic human heights. By the way, there is nothing special
371 about the values 170 and 50 : the priors are only a starting point, which should just identify
372 a plausible range of heights without being unnecessarily constraining. Different, reasonable
373 choices for the priors would still lead to very similar outcomes.

374 A prior for σ should rule out negative values (that is, assign zero probability to them),
375 since a standard deviation must be a nonnegative number. A common choice for priors of
376 standard deviations is a so-called half-Cauchy distribution, which is a truncated Cauchy.
377 Precisely, we set the first (location) parameter of the prior (4) for σ to zero, so that the
378 distribution's support is restricted to the nonnegative reals. We set the second (scale)
379 parameter to one, which spreads out the probabilities smoothly while still preferring
380 moderate values of σ .

381 Bayesian data analysis tools can often suggest *default* weakly informative priors that may
382 work well in many cases. In the analysis of Section 3, we will define our priors—following
383 standard recommendations [42]—but very often using default priors would have lead to
384 overall similar results. In any case, the plausibility checks described next will validate our
385 choice of priors.

386

387

388

389

390

391

392

2.2 Plausible model

Once we have chosen parameters, likelihood, and priors our model definition includes all required parts. Then, we can “run” the model—that is, sample from it—and analyze how likelihood, data, and priors constrain the model parameters of interest [30]. The first

STEP	PRIOR	LIKELIHOOD	ALTERNATIVE MODELS	EMPIRICAL DATA	NEW DATA FOR PREDICTION
plausible?	✓				
workable?	✓	✓			
adequate?	✓	✓		✓	
compare			✓		
analyze	✓	✓		✓	✓

Table 1. For each step of a Bayesian statistical analysis (checks of plausibility, workability, adequacy, model comparison, and analysis), the model components (PRIOR and LIKELIHOOD), competing ALTERNATIVE models, and kinds of DATA (EMPIRICAL and NEW FOR PREDICTION) that the step primarily tests.

STEP	ARTIFACTS	OUTCOME
plausible?	(1) Prior predictive simulation plots; (2) Justification for priors if they disallow certain values.	The priors allow a broad range of possible values and give low probability to values that are unlikely to occur in the domain.
workable?	(1) Simulation-based calibration of z score and shrinkage; (2) Fitting diagnostic metrics.	Fitting the model works computationally and does not exhibit pathological behavior.
adequate?	Posterior predictive checks plots.	The model can generate data similar to the empirical observations.
compare	Information-criteria ranking and scores of competing alternative models.	The chosen model achieves a bias-variance trade-off better than the alternative models.
analyze	(1) Posterior plots based on the empirical data; (2) Distribution plots and summary statistics of any domain-specific variables of interest.	Quantitative answers to the analysis's specific questions.

Table 2. The ARTIFACTS that are typically produced, and the OUTCOME that follows from each step of a Bayesian statistical analysis (when the step succeeds).

step of this analysis focuses on the priors, which should be neither too constraining nor unreasonably permissive. This step is typically called *prior predictive simulations* or *prior predictive checks* [42, 57] and works as follows: sample the priors broadly; using the sampled distribution (ignoring the actual empirical data, which are not used in this step), run the model to get a distribution of the model variables of interests (typically, the outcome variable); check that this distribution is plausible for the variables that it measures.

The key principle is that, if the priors are properly chosen, this process should determine a distribution that allows all plausible values for the variables but gives vanishing small probability to values that are practically impossible or contradict established scientific knowledge. In summary, prior predictive checks answer the question: *Does sampling from the priors lead to a plausible range of parameter values?*

Here is how prior predictive simulation would work on our toy example. First, we sample random values for parameters μ and σ from their prior distributions (3),(4). We then plug each sampled pair of values $\bar{\mu}, \bar{\sigma}$ into the likelihood (2) and sample values of h from $\text{Normal}(\bar{\mu}, \bar{\sigma})$. Since the outcome variable h measures an adult person's height, the priors should be such that this sampled distribution of h freely allows heights between, say, 0

442 and 300 cm, whereas it disallows negative heights and assigns very small probabilities to
443 heights above 300 cm—since no human on record has ever been that tall.

444 Prior predictive simulation is not cheating [42]. As long as we do not set priors based
445 on the actual empirical data that we are going to analyze, but only through what we
446 know about the data *domain* independent of how we measured it, it is sensible to use our
447 existing knowledge to rule out priors that would lead to impossible or clearly implausible
448 results. Seen in this light, the possibility of choosing priors is a big advantage of Bayesian
449 analysis that is highly valuable for any empirical science. Using existing knowledge to
450 guide new analyses, we can develop sequences of studies that, taken together, progressively
451 sharpen knowledge in a specific area. The alternative is that every software engineering
452 research contribution remains an “island unto itself” without clear connections to the related
453 literature and the field as a whole [21].

454 The reasons for choosing certain priors that make the model *plausible* should be explicitly
455 justified. In practice, and to the extent that it is possible, empirical software engineering
456 studies should explicitly state which published results, common sense, or “folk knowledge”
457 justify the choice of priors and their plausibility behavior. Section 3.3 demonstrates how to
458 do that for the paper’s case study.

459 Selecting informative priors gives Bayesian statistics more flexibility, but does not limit
460 its applicability. When very little is known about the problem domain—for example, in
461 an exploratory first study about a certain practice—one can always fall back to using
462 completely uninformative priors, which require no specific knowledge, and hence are
463 vacuously plausible. When prior knowledge exists, however, defining more selective priors
464 can help sharpen the model and specialize it to the characteristics of the analysis domain.

466 2.3 Workable model

467 Once we have ascertained that the chosen priors are consistent with plausible parameter
468 values, the second step checks whether our model *works* computationally—that is, *fitting*
469 the model does not incur divergence or other numerical problems, and the fitting process
470 eventually reaches a stationary state that properly identifies a posterior distribution.

471 An emerging technique to do so is *simulation-based calibration* [57, 62], which relies on a
472 consistency property of Bayesian models: first, simulate parameter and data values from the
473 priors and likelihood as done in prior predictive simulations; then, using Bayes’ theorem,
474 combine the simulated parameter and data samples to get a *posterior* distribution of the
475 model’s parameter; if the model is consistent, the posterior distribution obtained in this
476 way should resemble the prior distribution. To perform simulation-based calibration on
477 our toy example, we would (i) sample parameter values from the priors; (ii) use those to
478 build samples of the outcome variable h ; (iii) use these outcome samples as data (instead of
479 the actual empirical data) and combine them again with priors and likelihood using Bayes’
480 theorem (1). These steps give a new sampled distribution of the parameters μ and σ , which
481 we compare with that obtained by sampling the priors directly in the first step.

482 While promising, simulation-based calibration is a cutting-edge technique that is still
483 undergoing major developments; none of the statistical analysis tools that are more widely
484 used for Bayesian analysis support it out-of-the-box. Instead, these tools offer other metrics
485 to assess workability that are specific to the fitting algorithms based on dynamic Hamilton-
486 ian Monte Carlo which they implement. Here are the metrics that are usually available, and
487 how they help us assess workability:

- A *divergent transition* in the sequences of samples indicates a possible numerical error; workable models should have few divergent transitions—ideally none.
- The sampling process is repeated a few (usually 2–4) times independently; each sequence of sampling is called a *chain*. In a workable model, different chains should be statistically similar: the ratio \hat{R} of within-to-between chain variance should converge to 1 as the number of samples grows. A common rule of thumb for finite sampling is that $\hat{R} < 1.01$, which indicates a stationary posterior distribution.
- The *effective sample size* is the fraction of all samples that are independent, that is not autocorrelated. We typically want it to be at least 10% for each parameter we estimate (and that the absolute number of independent samples be a few hundreds); lower values may indicate that sampling is ineffective.
- Finally, we can also visually inspect the plots that trace the samples in every chain. When the different lines look mixed up (like a “hairy caterpillar” [42]), it is one more sign that the fitting process works well.

Section 3.4 uses these metrics to analyze the workability of our models.

When a workability check fails, it suggests that there is a mismatch between the model and the algorithm used to fit it. Sometimes, this is due to the data—for instance it is too sparse to effectively sample from it. More commonly, it indicates that the model itself is unsuitable for the analysis at hand. A clear example is the problem of *multicollinearity*: when two variables are strongly correlated, their exact contribution to the outcome is undetermined; thus, the model may not be workable because it cannot be used to discover a definite value for each variable independent of the other.

2.4 Adequate model

If the previous analysis steps were successful, we determined that the priors are sensible (*plausibility*) and that fitting the model is a converging process (*workability*); it remains to check whether the model adequately captures reality. The third step thus fits the model using the actual empirical data (which was not used in the previous two steps) and performs *posterior predictive checks*: using the posterior distribution of parameters fitted on the actual empirical data, simulate new observations and compare them to the data. If the two are consistent, it means that the model can generate data similar to the observed data, and hence it captures the empirical observations *adequately*.

Here is how posterior predictive checks would work on our toy example. Similarly as in simulation-based calibration, we combine data and prior samples using Bayes’ theorem (1); the key difference is that we now use the actual observed data (the height of real people) instead of simulated data. This gives a posterior predictive distribution of parameters μ and σ , which, in turn, we sample; then, we plug the sampled parameter values into the likelihood to get a distribution of h —the so-called *posterior predictive distribution*, since it expresses the information about the posterior indirectly in terms of prediction of model (outcome) variables. In an adequate model, the posterior predictive distribution generates data somewhat similar to the actual observed data.

Section 3.5 discusses the results of posterior predictive checks on the programming language case study.

2.5 Model comparison

Information criteria such as WAIC (Widely Applicable Information Criterion, also known as Watanabe-Akaike Information Criterion) [71] and PSIS-LOO (Pareto-Smoothed Importance

540 Sampling Leave One Out validation) [67] assess a kind of relative adequacy by measuring
541 deviance or other information-theoretic metrics between a model's predictions and the data.
542 In a nutshell, these metrics assess how well each model performs out-of-sample predictions
543 compared to other competing models. Thus, information criteria measures are *relative*:
544 they are useful to compare the adequacy of a model relative to another but cannot gauge
545 a model's adequacy in absolute terms. Section 3.6 uses information criteria to compare
546 different models for the programming language data analysis.

547

548 2.6 Iterative refinement

549 After a candidate model goes through the steps described above, we have a clear under-
550 standing of its strengths and weaknesses. When the model fails specific steps, we also learn
551 what aspects we have to change to refine it: the priors of an implausible model need chang-
552 ing; an unworkable model needs to be refactored in a way that works computationally;
553 an inadequate model may require more information (typically in the form of additional
554 variables or parameters) for it to be consistent with the data (for example, to properly
555 capture inter-group variability).

556 Model design is an *iterative* process which gradually refines an initial model to improve
557 it. Usually, we start from a deliberately very simple model and make it more complex as
558 needed [57]. However, we can also do the opposite: start from a so-called *maximal* model,
559 and then simplify it as long as it retains the characteristics of plausibility, workability, and
560 adequacy [49]. In practice, we may even alternate simplification and refinement (detail-
561 adding) steps starting from a canonical model [47, pp. 103–104] until we are satisfied with
562 the results.

563 The presentation of the results of a Bayesian data analysis need not discuss the models in
564 the same order in which they were designed and evaluated; it does not even need to present
565 all models, but can simply present the final model as long as its choice can be soundly
566 justified a posteriori (and, preferably, a reproducibility package exists). Regardless of how
567 we choose to present the overall outcome of an analysis, considering different models
568 expands the flexibility of the modeling process, supports making informed choices about
569 each aspect of a model, and helps focus on and quantify the relative benefits of each model
570 in terms of the trade-offs that matter for the ongoing analysis.

571

572 *2.6.1 Uniqueness and optimality of models.* When should we stop refining our model? Para-
573 phrasing George Box's famous aphorism [9, 10], we could say that the goal of statistical
574 modeling is building a useful model, not a correct one. In other words, we cannot expect
575 that following our guidelines leads to designing a unique or optimal model.

576 Model comparison can identify which models perform better predictions than other
577 models, but it cannot assess a model's absolute predictive capabilities. The steps in Figure 1
578 make up a *validation* process, which can identify a model's shortcomings or confirm that it
579 is of suitable quality; they cannot say anything about the infinitely many *other* models that
580 were not considered. Building useful models still requires human intuition, knowledge,
581 and ingenuity—skills that no supporting process can completely replace.

582

583 2.7 Tools for Bayesian data analysis

584 Let us briefly mention which tools are available to support the kind of Bayesian data
585 analysis process that we discuss in this paper. Stan [13] and JAGS [51] are state-of-the-art
586 frameworks that offer a probabilistic language to express Bayesian models and implement
587 very efficient algorithms to fit such models on data. Commonly, one uses these frameworks

588

589 through a front-end library in a high-level programming language suitable for data analysis.
590 Libraries such as `brms` for R [12], `Stan.jl` for Julia [39], and `PyStan` for Python [52] provide
591 a rich interface to Stan, including support for the main steps of our guidelines (for example,
592 prior predictive simulations). `Turing.jl` for Julia [25] also provides a high-level interface to
593 perform Bayesian data analysis, but includes its own implementation of Bayesian sampling
594 instead of relying on Stan's (which may offer some advantages in terms of flexibility and
595 generality for the most advanced applications).

596 3 BAYESIAN DATA ANALYSIS OF PROGRAMMING LANGUAGE DATA 597

598 Equipped with a high-level understanding of the modeling guidelines that we outlined in
599 Section 2, we apply them to perform the analysis of the FSE data. The overall outcome of
600 the work described in this section will be a carefully designed, suitable statistical model of
601 this data. In Section 4, we will analyze this model to understand what it tells us about the
602 original questions on programming languages and code quality.

603 To mitigate the risk of mono-operational bias, the first author prepared the data for
604 analysis in R, and the second author developed the first complete analysis, which then the
605 first and third author revised. Finally, all three authors validated the final revised analysis,
606 which is presented here. In addition, the second author did not read the publication that
607 originated the dataset [55] or its reanalysis [6] until after completing the first complete
608 analysis. This reduced the chance that the others' design decisions, or some characteristics
609 of the data they highlighted, biased our application of the modeling guidelines.⁷
610

611 3.1 Data

612 FSE's authors released the original dataset—obtained by mining information from GitHub
613 repositories—upon request from TOPLAS's authors. TOPLAS performed first a repetition
614 of FSE's analysis on the same dataset, and then a reanalysis on a revised dataset obtained
615 by "alternative data processing and statistical analysis to address what [they] identified as
616 methodological weaknesses of the original work" [6, Sec. 4]. The main difference between
617 FSE's original dataset and TOPLAS's revised dataset is that the latter removes some dupli-
618 cated data, TypeScript projects (which often do not include much actual TypeScript code),
619 and the V8 project (whose JavaScript code in the dataset is mostly tests). Finally, TOPLAS's
620 replication package includes FSE's original dataset alongside TOPLAS's revised dataset.

621 In our analysis, we focus on the original FSE dataset,⁸ because we would like to see
622 whether a Bayesian data analysis can help spot issues and inconsistencies in the data that
623 may hinder replication attempts—and, conversely, that may make replication run-of-the-
624 mill if addressed early on. Our replication package includes all analysis details, including
625 the results of fitting the same models on TOPLAS's revised dataset (which we do not discuss
626 here for brevity).

627 FSE's dataset includes information about 1 578 165 commits, which we group by project
628 and language giving 1 127 datapoints. The attributes that are relevant for our analysis are:

629 *project*: the project's name
630 *language*: the used programming language
631 *commits*: the total number of commits in the project
632 *insertions*: the total number of inserted lines in all commits
633

634 ⁷We used Stan through its `brms` R front-end to perform the analysis described in the rest of the paper.

635 ⁸Which we obtained from TOPLAS's public replication package (available at [https://github.com/PRL-PRG/](https://github.com/PRL-PRG/TOPLAS19_Artifact)
636 `TOPLAS19_Artifact`).

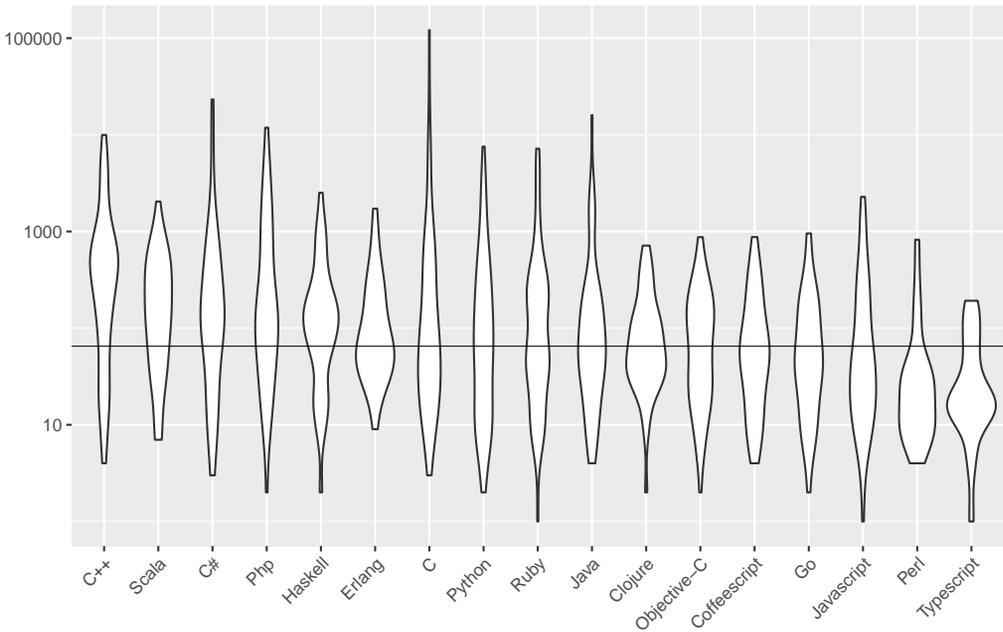


Fig. 2. Violin plots of the distributions of number of bugs per project for each programming language in the original FSE dataset. Languages are sorted, left-to-right, by decreasing values of the distributions' medians. The vertical axis's scale is logarithmic in base 10. An horizontal line marks the median number of bugs per project across all languages.

age: the time passed since the oldest recorded commit in the project
devs: the total number of users committing code to the project
bugs: the number of commits classified as "bugs"

The values of attributes *commits*, *insertions*, *age*, and *devs* vary greatly between projects. When this happens, it is customary to transform the data using a logarithmic function, so that the variability is over a smaller range whose unit corresponds to an order of magnitude. Both FSE's and TOPLAS's analyses log-transformed these attributes; we do the same: henceforth, *commits*, *insertions*, *age*, and *devs* represent the natural *logarithm* of the total number of commits, inserted lines, and so on.

Figure 2 provides an overview of the FSE dataset, showing the distribution of bugs per project grouped by programming language. Visualizing the raw data can be useful to get a broad idea of what is in the dataset; however, the information that such visualizations provide is mostly qualitative and we should be aware of its limitations. If one includes all data, any outliers may skew the picture at extreme values; conversely, if one excludes some data, deciding which data to exclude is itself a source of possible bias, and discards potentially useful information thus increasing uncertainty. In this dataset specifically, projects vary broadly in terms of size and other characteristics. Figure 2 conflates these differences, and hence a comparison of different languages based on it may be misleading. We could display a subset of the data that only includes projects with homogeneous characteristics; however, doing so would drop significant amounts of information, introduce a somewhat

STEP	\mathcal{M}_1	\mathcal{M}_2	\mathcal{M}_3
plausible?	✓	✓	✓
workable?	✓	✓	✓
adequate?	✗	✓	✓
compare	–	✗	✓

Table 3. Plausibility, workability, adequacy, and comparison of models \mathcal{M}_1 , \mathcal{M}_2 , \mathcal{M}_3 . All three models are *plausible* and *workable*, but model \mathcal{M}_1 is *not adequate* because it cannot accurately capture the regular features of the dataset. Model *comparison* between \mathcal{M}_2 and \mathcal{M}_3 shows that the latter performs much better concerning out-of-sample predictions, and hence we will use \mathcal{M}_3 for the rest of the analysis.

arbitrary partitioning (what projects are “similar?”), and increase the risk of overfitting other accidental characteristics of the data. By abstracting the information in the raw data and combining it with expert knowledge, a suitable statistical model can lessen several of these problems, thus supporting more robust and general inferences about the impact of programming languages.

3.2 Modeling

We build three models— \mathcal{M}_1 , \mathcal{M}_2 , and \mathcal{M}_3 —of increasing complexity.⁹ Table 3 summarizes the outcome of the steps in Figure 1 for the three models. Mirroring Table 2’s structure, Table 4 outlines the artifacts produced in each step, and the conclusions that the analysis draws about each model’s suitability. The rest of this section details the models and the outcome of the guidelines’ suitability checks presented in Section 2. This section presents the models and the outcome of their suitability analysis in detail; later, Section 4.1 will discuss, at a higher level, what this analysis reveals about the relations between model features and data.

As we remarked in Section 1.3.1, our guidelines can be used to validate different kinds of models. We consider these three models because they belong to a widely used family of statistical models, and for their similarity with the models of FSE and TOPLAS. An analysis with different goals or done by analysts with different expertise could end up building very different kinds of models—but they should still undergo the same validation steps.

3.2.1 Likelihood (and parameters). Generalized linear models are a broad category of statistical models that are so flexible that they can be “applied to just about any problem” [28] when modeling empirical data. The likelihood of a generalized linear model is a probability distribution over certain parameters, which are generalized linear functions of the variables chosen as predictors. The values drawn from the distribution correspond to the outcome that we are modeling.

Distribution family. In our case, the *outcome* variable is *bugs*, which always is a nonnegative integer. Therefore, we should select a likelihood distribution suitable for “counting”—that is, one in the Poisson family. The single-parameter Poisson is the distribution in this family with the highest information entropy [37], and hence it should be the customary initial choice.

Nevertheless, building a model using the single-parameter Poisson quickly reveals that it cannot account for the fact that the distribution of *bugs* in the data is *overdispersed*: its mean

⁹As discussed in Section 2.6, for clarity we present the three models at once but we actually designed them over several iterated applications of the guidelines.

STEP	ARTIFACTS	OUTCOME
plausible?	(1) Prior predictive simulation plots in Figure 5; (2) Justification for priors: typical relations between project size and number of known bugs [58].	The priors of all three models allow a very broad range of possible values for the number of bugs that may exist; extremely high numbers are still possible but with low probability.
workable?	Fitting diagnostic metrics reported in Section 3.4: \hat{R} , effective sample size, no divergent transitions, trace plots (details in the replication package).	Fitting all three models works computationally and reaches convergence.
adequate?	Posterior predictive checks plots in Figure 7.	Model \mathcal{M}_1 cannot generate a distribution similar to that observed in the data (top plot in Figure 7), whereas models \mathcal{M}_2 and \mathcal{M}_3 can.
compare	Information-criteria scores of competing models in Table 5.	Model \mathcal{M}_3 clearly outperforms \mathcal{M}_2 in how it can predict data out of the sample used for fitting.
analyze	(1) Posterior plots based on the empirical data in Figures 11 and 13; (2) Distribution plots and summary statistics of domain-specific variables of interest in Figures 8 and 12.	Quantitative answers to the analysis's specific questions in Section 4.

Table 4. A summary of the ARTIFACTS produced by the analysis of the three models, and the OUTCOME of each step of the analysis in terms model suitability. This summary instantiates Table 2 for the programming language data analysis.

$\mu_{bugs} = 501$ is much smaller than its variance $\sigma_{bugs}^2 = 15\,031\,006$. This justifies selecting the slightly more complex *negative binomial* distribution $\text{NegativeBinomial}(\lambda, \phi)$. The two parameters λ and ϕ represent¹⁰ rates that together determine the distribution's mean λ and variance $\lambda + \lambda^2/\phi$, which can take different values to accurately capture overdispersion. This is in contrast to the Poisson distribution whose mean and variance coincide. The negative binomial distribution is also the same distribution selected, for the same reason, by both FSE's original analysis and TOPLAS's reanalysis.

Model \mathcal{M}_1 . The first model we consider, called \mathcal{M}_1 , is very simple: it assumes that the rate λ is a function of two terms only. The first term Π is a constant intercept α ; the symbol Π highlights that it is a population-level term. The second term L is an additional intercept $\alpha_{language}$ that depends only on the *language* used in each observation; the symbol L highlights that it is a language-level term. Figure 3a shows \mathcal{M}_1 's overall likelihood, where the logarithm function *links*¹¹ the linear function of the parameters and λ so that the latter is always a nonnegative number—as it should be in a “counting” distribution.

Model \mathcal{M}_1 is obviously too simple to capture the variability in the data with high accuracy. Nonetheless, it is a useful starting point to understand the key relations between variables and to bootstrap the process that leads to incrementally more refined and precise models. In its simplicity, it highlights that the key predictor (the “treatment”) is the programming language used in each project, whose relation with the number of bugs we would like to capture. Finally, even a simplistic model serves as a useful *baseline* to compare to more

¹⁰<https://mc-stan.org/docs/2.20/functions-reference/nbalt.html>

¹¹In other words, the link function converts measures from the probability space to the outcome space.

$$\begin{array}{ll}
785 & \text{bugs}_i \sim \text{NegativeBinomial}(\lambda_i, \phi) \\
786 & \text{bugs}_i \sim \text{NegativeBinomial}(\lambda_i, \phi) & \log(\lambda_i) = \Pi_i + L_i \\
787 & \log(\lambda_i) = \Pi_i + L_i & \Pi_i = \alpha + \beta^C \cdot \text{commits}_i + \beta^I \cdot \text{insertions}_i \\
788 & \Pi_i = \alpha & \quad + \beta^A \cdot \text{age}_i + \beta^D \cdot \text{devs}_i \\
789 & L_i = \alpha_{\text{language}_i} & L_i = \alpha_{\text{language}_i} \\
790 & & \\
791 & & \\
792 & \text{(a) Model } \mathcal{M}_1 & \text{(b) Model } \mathcal{M}_2 \\
793 & & \\
794 & \text{bugs}_i \sim \text{NegativeBinomial}(\lambda_i, \phi) \\
795 & \log(\lambda_i) = \Pi_i + L_i + P_i \\
796 & \Pi_i = \alpha + \beta^C \cdot \text{commits}_i + \beta^I \cdot \text{insertions}_i \\
797 & \quad + \beta^A \cdot \text{age}_i + \beta^D \cdot \text{devs}_i \\
798 & L_i = \alpha_{\text{language}_i} + \beta_{\text{language}_i}^C \cdot \text{commits}_i \\
799 & \quad + \beta_{\text{language}_i}^I \cdot \text{insertions}_i + \beta_{\text{language}_i}^A \cdot \text{age}_i \\
800 & \quad + \beta_{\text{language}_i}^D \cdot \text{devs}_i \\
801 & P_i = \alpha_{\text{project}_i} \\
802 & \\
803 & \text{(c) Model } \mathcal{M}_3 \\
804 & \\
805 & \\
806 & \\
807 & \\
808 & \\
809 & \\
810 & \\
811 & \\
812 & \\
813 & \\
814 & \\
815 & \\
816 & \\
817 & \\
818 & \\
819 & \\
820 & \\
821 & \\
822 & \\
823 & \\
824 & \\
825 & \\
826 & \\
827 & \\
828 & \\
829 & \\
830 & \\
831 & \\
832 & \\
833 &
\end{array}$$

Fig. 3. The likelihoods of statistical models \mathcal{M}_1 , \mathcal{M}_2 , and \mathcal{M}_3 . Colors highlight the terms that are added to each model compared to the previous ones.

complex models—as a sanity check that the additional complexity that we are going to add to the models brings measurable improvements over the baseline.

Model \mathcal{M}_2 . The second model we consider, called \mathcal{M}_2 , is a standard linear-regressive model with negative binomial likelihood. Model \mathcal{M}_2 's population-level term Π is a linear function with intercept α and a slope β for each predictor variable *commits*, *insertions*, *age*, and *devs*. In addition, like model \mathcal{M}_1 , \mathcal{M}_2 includes a language-level term L that consists of an intercept that depends on the language used in each observation (a so-called “varying intercept” model [29]). Figure 3b shows \mathcal{M}_2 's overall likelihood.

Model \mathcal{M}_2 is the closest to the regressive models used in FSE and TOPLAS. The only difference is how each model accounts for the dependence on the programming *language*: FSE and TOPLAS use different kinds of *contrasts* [6, 55], whereas we simply add an intercept language-level term—thus making our models *multilevel* [28]. Multilevel modeling comes natural with Bayesian statistics, both because we do not have to worry too much about adding layers to the model (unlike with frequentist techniques, changing such characteristics of the model does not require changing the fitting algorithm) and because we can just model the quantities of interest directly and *compute* any derived quantity *after* we fit the model's *posterior* distribution (unlike with frequentist techniques, which mostly provide only point estimates without distributional information).

Model \mathcal{M}_3 . The third model we consider, called \mathcal{M}_3 , is a multilevel model that tries to capture the effect of the programming language with greater detail. Model \mathcal{M}_3 's population-level term Π is identical to \mathcal{M}_2 's. Its language-level term L is considerably more complex, since it introduces a linear model with different intercepts α_{language} and slopes β_{language} for each programming *language* (a so-called “varying intercepts and varying slopes” model,

834 also commonly known as “varying effects” model [29]). Unlike the population-level term Π ,
835 the language-level term L pools the information about each data cluster—where clusters are
836 identified by the used programming language. Since it clusters by programming language,
837 this partial pooling may help capture more accurately the effects of choosing a programming
838 language instead of another; at the same time, it also shares information among clusters
839 so that some information from larger clusters (languages with many projects) can sharpen
840 the information from smaller clusters (languages with fewer projects). This also means that
841 partial pooling helps protect from *overfitting*, as learning takes place first separately on each
842 cluster, and then is “regularized” by sharing its results among different clusters.

843 The three models’ focus on the programming language reflects our intuitive expectation
844 that the relation between programming languages and proneness to bugs is an important
845 one—regardless of whether it turns out to be significant or negligible in the end. At the same
846 time, adding predictors other than the programming language *accounts for* confounding
847 factors that may have a stronger correlation with the number of bugs. But what if the
848 intrinsic differences between *projects* turn out to dominate the discrepancies in code quality?
849 For instance, different projects may have wildly different protocols to report, triage, and fix
850 bugs, which might have an effect on the observed number of bugs.

851 In order to hedge against this possible confounding factor, model \mathcal{M}_3 also includes a
852 term P : an additional intercept $\alpha_{project}$ that depends only on each observation’s *project*;
853 the symbol P highlights that it is a project-level term, which will help in quantifying the
854 intrinsic variability across projects. Figure 3c shows \mathcal{M}_3 ’s overall likelihood.

855
856 **3.2.2 Priors.** As we demonstrated in the previous section, choosing the likelihood typically
857 requires making justified modeling choices, which depend on the kind of analysis we would
858 like to carry out.

859 When choosing the priors, in contrast, we can often rely on standard recommendations
860 that primarily depend on the domain of each variable. This does not mean that priors (or
861 likelihoods, for that matter) can be always chosen blindly using a fixed table of recommen-
862 dations. In the following sections, as we go through the various steps of the Bayesian data
863 analysis workflow, we will validate our choices of priors and likelihood. If validation fails,
864 we have to go back and revise the model: priors, likelihoods, or both.

865 **Model \mathcal{M}_1 .** As shown in Figure 4a, we use weakly informative priors for model \mathcal{M}_1
866 that are based on the normal distribution. As we will see during the plausibility analysis
867 (Section 3.3), model \mathcal{M}_1 is so simplistic that its performance is not affected much by the
868 choice of priors; nonetheless, we discuss its priors in some detail because we will build on
869 them to choose priors for the more complex models.

870 The intercept α ’s prior has mean 0 (that is, we do not know a priori whether the intercept
871 is positive or negative) and standard deviation 5. Remember that the estimated parameter
872 λ is log-transformed (see Figure 3b); therefore, we can appreciate how weakly constraining
873 this prior is: two standard deviations on each side of zero span the interval from $e^{-10} \simeq 0$ to
874 $e^{10} \simeq 22\,000$ on the bug counting scale; that is, the prior only assumes that a project’s bugs
875 are up to 22 000 with 95% probability—which is not a strong assumption at all [58]. Besides,
876 a normal distribution has *infinite* support, and hence it does not rule out any count of bugs
877 if the data provides evidence for it. The prior for the language-level intercept $\alpha_{language}$ is
878 also a normal distribution with mean 0; however, choosing the same standard deviation σ_α
879 for every language would defeat the purpose of having language-level intercepts. Instead,
880 we let σ_α be a random variable, and assign a prior to it. Distributions with support limited
881 to positive values are suitable priors for standard deviations—which must be nonnegative
882

883		$\alpha \sim \text{Normal}(0, 5)$
884		$\beta \sim \text{Normal}(0, 0.5)$
885	$\alpha \sim \text{Normal}(0, 5)$	
886	$\alpha_{\text{language}} \sim \text{Normal}(0, \sigma_\alpha)$	$\alpha_{\text{language}} \sim \text{Normal}(0, \sigma_\alpha)$
887	$\sigma_\alpha \sim \text{Weibull}(2, 1)$	$\sigma_\alpha \sim \text{Weibull}(2, 1)$
888	$\phi \sim \text{gamma}(0.01, 0.01)$	$\phi \sim \text{gamma}(0.01, 0.01)$
889		
890	(a) Priors of model \mathcal{M}_1	(b) Priors of model \mathcal{M}_2
891		
892		$\alpha \sim \text{Normal}(0, 5)$
893		$\beta \sim \text{Normal}(0, 0.5)$
894		$\alpha_{\text{language}} \sim \text{Normal}(0, \sigma_\alpha)$
895		$\beta_{\text{language}} \sim \text{Normal}(0, \sigma_\beta)$
896		$\alpha_{\text{project}} \sim \text{Normal}(0, \sigma_\gamma)$
897		$\sigma_\alpha, \sigma_\beta, \sigma_\gamma \sim \text{Weibull}(2, 1)$
898		$\mathcal{L} \sim \text{LKJ}(2)$
899		$\phi \sim \text{gamma}(0.01, 0.01)$
900		
901		
902		(c) Priors of model \mathcal{M}_3
903		

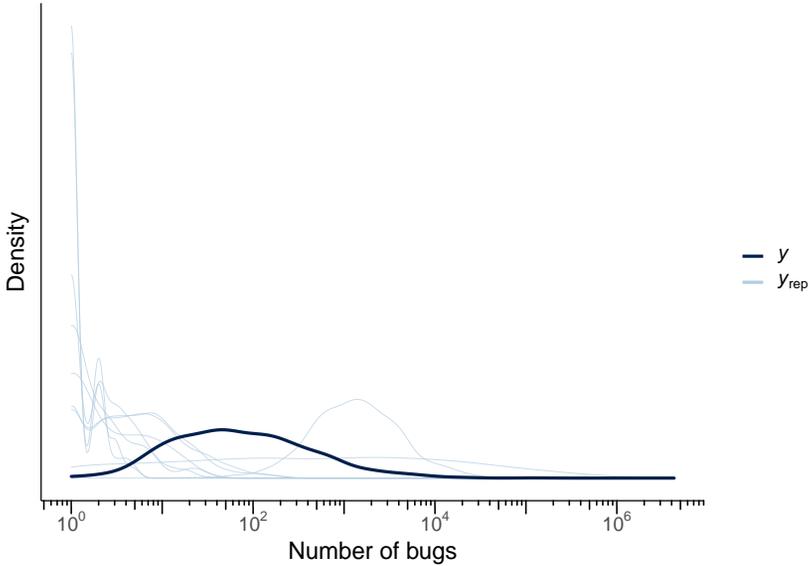
Fig. 4. The priors of statistical models \mathcal{M}_1 , \mathcal{M}_2 , and \mathcal{M}_3 . Colors highlight the terms that are added to each model compared to the previous ones.

values. In this case, we use a Weibull for σ_α and the default Gamma for the dispersion parameter ϕ of the negative binomial.

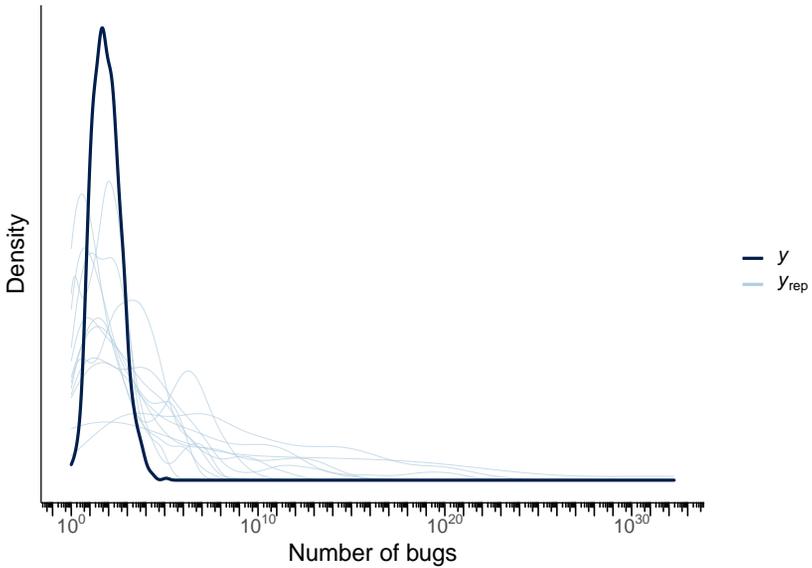
Model \mathcal{M}_2 . In addition to \mathcal{M}_1 's weakly informative priors for α , α_{language} , and ϕ , \mathcal{M}_2 needs a prior for the slope parameter vector β . Here too we use a simple normal distribution with mean 0 (so that there is no bias in the possible direction of each predictor's effect) and standard deviation 0.5 (which still allows for a broad variability on the logarithmic scale). The β 's prior standard deviations are smaller than the α 's because α determines the population average, and then β moves this average according to each predictor's effect (which needs only introduce a smaller variation relative to the average). Figure 4b shows the overall priors for \mathcal{M}_2 .

Model \mathcal{M}_3 . We choose the prior for the new part of \mathcal{M}_3 —the language-level slopes β_{language} —similarly to how we chose the language-level intercepts α_{language} : a normal with mean 0 and a random variable σ_β for standard deviation. However, there is an additional technicality that we need to handle: vectors α_{language} and β_{language} are not independent but are components of a single *multivariate* normal distribution with a variance *matrix* S . Variance matrix S combines diagonal matrices with the components of α_{language} and β_{language} and a covariance matrix \mathcal{L} . The customary prior for covariance matrices is a multivariate Lewandowski-Kurowicka-Joe distribution; $\text{LKJ}(2)$ is a weakly informative prior using this distribution, which assigns low probabilities to extreme correlations. Finally, the project-level intercept α_{project} 's prior is also a normal with mean 0 and a random variable σ_γ for standard deviation. Just like for the other standard deviations, we choose a Weibull as prior distribution of σ_γ —a weakly informative distribution that constrains α_{project} 's standard deviation to be a non-negative value. Figure 4c shows the overall priors for \mathcal{M}_3 .

932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980



(a) Prior predictive simulation for \mathcal{M}_2 .



(b) Prior predictive simulation for \mathcal{M}_3 .

Fig. 5. Prior predictive simulation plots for models \mathcal{M}_2 and \mathcal{M}_3 : each thin light blue line pictures one simulated distribution of the number of bugs in a project drawn from the priors. For comparison, the thick dark blue line pictures the distribution of the number of bugs in the measured data. The horizontal scale is logarithmic in base 10.

3.3 Plausibility

The prior predictive checks are straightforward for all three models, confirming that our choice of priors—based on standard recommendations for these kinds of models—leads to plausible outcomes. As an example, Figure 5a shows several distributions of the outcome variable *bugs* obtained with prior predictive simulations of \mathcal{M}_2 . These span a very wide support that goes from zero¹² up to over a million bugs per project. While there are no theoretical limits on the number of bugs in a project independent of its size, it is realistic that most projects have less than one million *known* bugs, and the majority of projects have less than a few thousands—simply because not many projects have more than one known bug for each line of code [58], and hence a project’s size in lines of code is a workable upper bound on the number of distinct bugs. Anyway, the priors still allow even larger bug counts, but assign to them increasingly smaller probabilities. Figure 5a also displays the empirical distribution of bug counts in FSE’s dataset (thick dark blue line); this visually confirms that the priors are not too restrictive and reflect reasonable expectations. The prior predictive checks of model \mathcal{M}_3 is shown in Figure 5b, and leads to qualitatively similar conclusions—in fact even stronger, given that the priors stretch past an astronomical number of bugs—about the model’s plausibility. So do the prior predictive checks of model \mathcal{M}_1 , which we do not show for brevity.

3.4 Workability

Section 2.3 outlined *simulation-based calibration* and Hamiltonian Monte Carlo validation metrics to assess a model’s workability. We use the latter, which are extensively supported by Stan, to determine whether the sampling process for each of the three models reached a stable state.

\hat{R} —the ratio of within-to-between chain variance—is < 1.01 for all three models; this indicates that the chains have converged towards a stationary posterior probability distribution. The *effective sample size* is at least 0.11, 0.17, 0.13 for all parameters in each of the three models, and confirms that sampling effectively converged. Fitting all three models does not run into any *divergent transitions*, and the trace plots of the models (included in the replication package) look well-mixed. In summary, all three models work well computationally.

3.5 Adequacy

Let us first study the *adequacy* of our models with posterior predictive simulations: we visually compare the distribution of number of bugs per project in our dataset to several simulated distributions using the fitted models. The top plot in Figure 7 indicates that \mathcal{M}_1 is *not* adequate: it is too simplistic to capture the data’s features; in particular, the means of the simulated distributions are more than ten times larger than the mean of the data (thick dark blue line). In contrast, \mathcal{M}_2 passes this adequacy test: the middle plot in Figure 7 shows that the model’s predictions look similar to the data. Model \mathcal{M}_3 also passes the visual adequacy test based on the posterior predictive simulations, as shown by the bottom plot in Figure 7. In Section 4.1, we will discuss what these adequacy results tell us about the interplay between model features and data.

3.6 Model comparison

It is now clear that \mathcal{M}_1 is too simplistic, but how to choose between \mathcal{M}_2 and \mathcal{M}_3 ? Information criteria, which measure the relative adequacy of different models fitted on the

¹²The logarithmic link function guarantees a lower bound of zero.

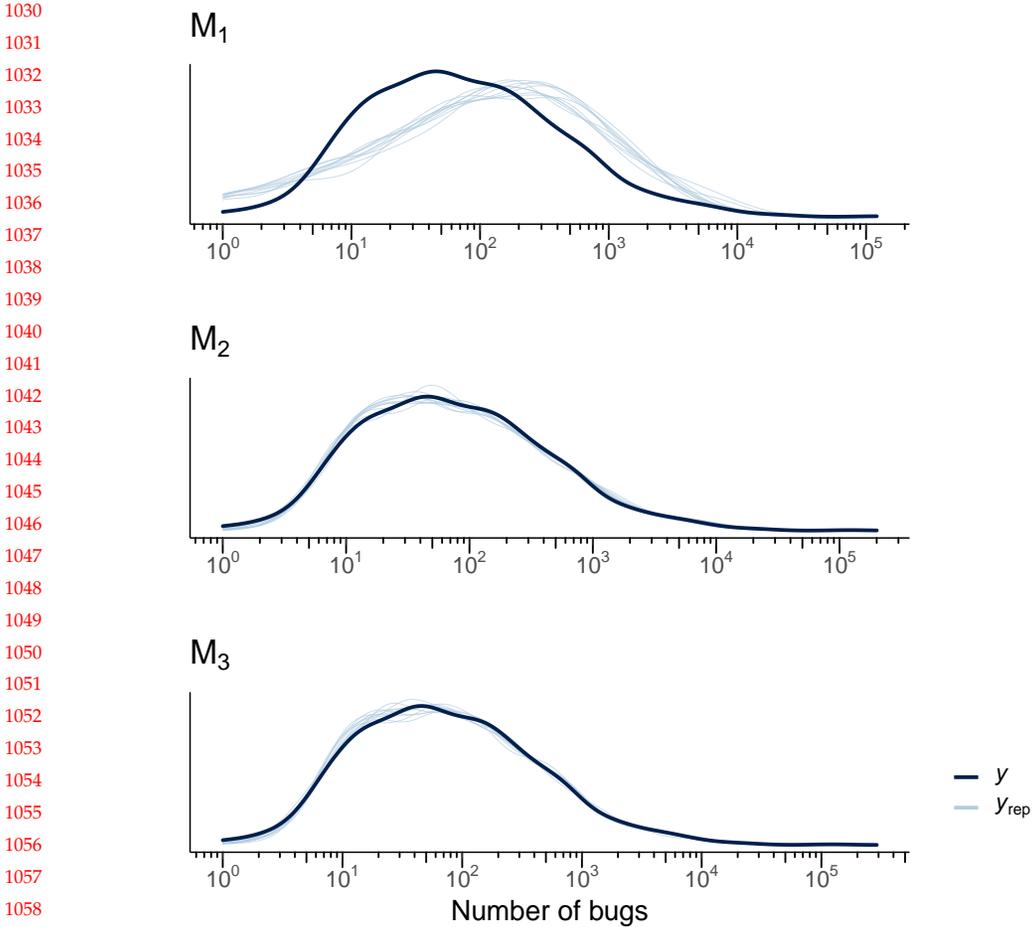


Fig. 6. Posterior predictive checks for \mathcal{M}_2 .

Fig. 7. Posterior predictive checks plots for models \mathcal{M}_1 , \mathcal{M}_2 , and \mathcal{M}_3 : each thin light blue line pictures one simulated distribution of the number of bugs in a project drawn from the posterior. For comparison, the thick dark blue line pictures the distribution of the number of bugs in the measured data's; in contrast, \mathcal{M}_1 fails the check because the simulated distributions deviate substantially from the data's; in contrast, \mathcal{M}_2 and \mathcal{M}_3 pass the check because the simulated distributions are similar to the data's. The horizontal scale is logarithmic in base 10.

same data, can help answer this question. We use the increasingly popular PSIS-LOO information criterion [67], which works well with models fitted using dynamic Hamiltonian Monte Carlo.¹³ In a nutshell, the criterion ranks the three models according to their relative adequacy. It also gives a *difference* score that measures how well each model performs out-of-sample predictions relative to the next one in the ranking, and a *standard error* of

¹³Compared to more traditional information criteria—such as AIC, BIC, and WAIC—PSIS-LOO can handle non-Gaussian likelihoods (as can WAIC) and also provides diagnostics useful for further analyzing whether a model's posterior behaves well numerically.

MODEL	RANK	DIFFERENCE	STANDARD ERROR
\mathcal{M}_3	1	–	–
\mathcal{M}_2	2	–172.0	23.6
\mathcal{M}_1	3	–1963.0	57.0

Table 5. Ranking of models (from better to worse) according to the PSIS-LOO information criterion. Each MODEL is RANKED (from better to worse) according to the PSIS-LOO information criterion. The score DIFFERENCE between each model and the immediately better one in the ranking, as well as the STANDARD ERROR of such difference, quantify the difference in adequacy between models.

the difference, which quantifies how much more adequate a model is compared to another. Table 5 displays the scores for the three models. Model \mathcal{M}_3 is ranked first; \mathcal{M}_2 comes second but its score is a whopping 7 standard errors worse than \mathcal{M}_3 's; and \mathcal{M}_1 is, unsurprisingly, a distant last.

We conclude that \mathcal{M}_3 is the “best” model among the three according to a variety of criteria. Our analysis will thus use \mathcal{M}_3 —starting with a discussion of its features from the point of view of the analysis’s goals in Section 4.1.

At some point, one has to stop adding model features and finalize a model for the current analysis. Nevertheless, as we discussed in Section 2.6.1, statistical modeling is never really done: as more insights, more data, or new techniques become available, we could go back to the drawing board and refine the latest model to better capture all available information.

4 BAYESIAN STATISTICAL ANALYSIS: RESULTS

When applied following a structured process—like the one we described in Section 2—Bayesian statistics does not simply produce dichotomous answers to research questions. The outcome of a Bayesian data analysis is a posterior probability distribution, which we can probe from different angles to get nuanced answers that apply to specific scenarios. In this section we are going to do this for our case study.

Section 4.1 discusses how the Bayesian data analysis process guided our choice of *models*: modeling is a looping process, whose feedback also informs us about key characteristics of the data we are analyzing. Not all data features have the same influence on the statistics. Section 4.2 illustrates how Bayesian analysis can point to variables with brittle or negligible predictive power that may indicate problems in how certain measures were operationalized. After understanding the features and limitations of the fitted model, Section 4.3 addresses the original study’s research questions in practical settings; and Section 4.4 outlines follow-up studies that could address some of the outstanding limitations.

4.1 Modeling

When following the Bayesian data analysis process in Figure 1, statistical modeling is based on principles and on checks that the models are suitable. This is in contrast to most frequentist statistical practices, which are primarily based on rules of thumb, conventions (“recipes”), and generic results, but may lack operational model-checking processes that can assess how much confidence we can put in a certain modeling choice.

Section 3 described such a principled Bayesian modeling process applied to the programming language data. The first outcome was ruling out \mathcal{M}_1 as inadequate, which was unsurprising given that \mathcal{M}_1 ignores most of the information that could explain the dataset’s variability. Still, even checks with predictable outcomes are useful: if they fail, they confirm

1128 that a more realistic model is needed and provide a minimal effectiveness yardstick; if they
1129 succeed, they avoid an overly complicated model. This can be especially important in the
1130 context of the software engineering industry, where a complex model can be more costly to
1131 understand, collect data for, and maintain.

1132 The second outcome of Section 3's analysis was indicating that, while both \mathcal{M}_2 and
1133 \mathcal{M}_3 are adequate, \mathcal{M}_3 clearly outperforms \mathcal{M}_2 in out-of-sample predictive capabilities.
1134 Informally, this means that \mathcal{M}_3 fits the data well, while still avoiding overfitting. In other
1135 words, \mathcal{M}_3 's additional complexity over \mathcal{M}_2 is justified by its much better effectiveness.
1136 This outcome is specific to the data that we are analyzing, and is not something that can be
1137 determined a priori for all models. Iteratively creating multiple models and then comparing
1138 them is thus an important part of any analysis.

1139 If we compare the definitions of \mathcal{M}_2 and \mathcal{M}_3 —in particular, their model specifications in
1140 Figure 3b and Figure 3c—we can attribute \mathcal{M}_3 's superior performance to its unique features.
1141 Unlike \mathcal{M}_2 , which only includes population-level effects that control for project character-
1142 istics other than the programming language, \mathcal{M}_3 includes a project-specific intercept
1143 and controls for the same characteristics with language-specific slopes. Thus, we see that
1144 clustering per project *and* per language captures the dataset's characteristics much better: if
1145 we do not do that, we may lose some of the "signal" in the data, or conflate different effects
1146 and associate them with a single generic predictor.

1147 An indirect advantage of Bayesian data analysis comes from the techniques that are
1148 commonly used to fit Bayesian models: flexible algorithmic techniques such as dynamic
1149 Hamiltonian Monte Carlo that can fit, in principle, models of arbitrary complexity—in
1150 contrast to *ad hoc* frequentist techniques that only work for specific, and often limited,
1151 distributional families. On the other hand, Bayesian models are often more effective not
1152 simply because they can be more complex. More complex models invariably fit better, but
1153 unwarranted complexity leads to *overfitting*: a model fits the data perfectly but fails to
1154 generalize. Bayesian analysis techniques include several features that specifically limit the
1155 risk of overfitting when exploring more expressive models:

- 1156
- 1157
- 1158 • Multi-level models, such as \mathcal{M}_3 , introduce *partial pooling*, which smoothens dif-
1159 ferences between groups of different size. In our case study, the data about some
1160 programming languages is more scarce than the data about others. For example,
1161 only 25 projects use Perl, whereas more than 200 use JavaScript; thus, overfitting
1162 Perl's data is a more serious risk than overfitting JavaScript's. Partial pooling works
1163 by transferring some of the information learned by fitting the larger groups to tune
1164 the fitting of the smaller groups, thus reducing the risk of overfitting the latter (and
1165 the whole dataset as a result).
- 1166 • Prior predictive simulations—discussed in Section 3.5—check that the priors we
1167 have chosen are *regularizing*: they are not so constraining that they prevent learning
1168 from the data, but they are also not so weak that they cannot prevent overfitting
1169 the data. Being able to choose priors, to select different priors, and to quantitatively
1170 compare their effectiveness is a distinct advantage of Bayesian statistics. Frequentist
1171 statistics usually have flat priors, which are the most prone to overfitting.
- 1172 • The information criteria that we used to select \mathcal{M}_3 measure the out-of-sample predic-
1173 tion performance of one model relative to the others. Models that are unnecessarily
1174 complex will overfit the data, and hence perform worse predictions for *new* data
1175 (different from the sample that has been used for fitting).
- 1176

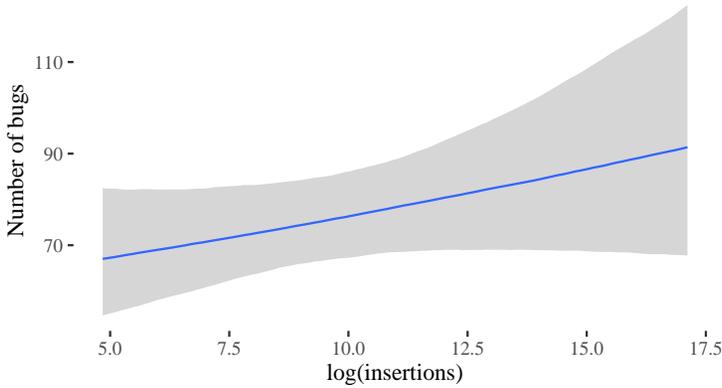


Fig. 8. Conditional effects of variable *insertions* (the logarithm of the total number of lines added to a project) on the outcome variable *bugs* (the number of bugs in a project), corresponding to the marginal distribution derived from the posterior of model \mathcal{M}_3 .

4.2 Spotting data problems

The rich information provided by a Bayesian data analysis may also highlight issues with the quality of (parts of) the data that is analyzed, and suggest which measures need to be cleaned up or improved.

Measuring size. Code size is a basic yet essential measure of complexity, which correlates with lots of other useful metrics of quality [32]. Therefore, controlling for project size is essential when analyzing heterogeneous projects. To this effect, the FSE study included a variable *size* in their regressive model, which measures the total number of inserted lines in all project commits—and which we called *insertions* in our models to make its actual meaning more transparent. The TOPLAS analysis criticized this choice of size metric—which does not take deletions and merges into account—and reported discrepancies between the raw commit data and the totals in FSE’s dataset. Does our Bayesian analysis offer any hints about the reliability of *insertions* as a measure of size?

A few results actually single out *insertions* as a poor predictor compared to the others:

- The 95% probability estimate of its population-level effect includes zero (namely, the (credibility) interval is $[-0.01, 0.06]$), which indicates some uncertainty about whether more inserted lines are associated with more or fewer bugs on average. Variable *insertions*’s mean estimated effect is still positive, but other predictors have more clearly defined effects.
- The plot of *insertions*’s conditional effect on the number of *bugs* in Figure 8 visually confirms a large uncertainty (again, compared with the other predictors’), which also increases with larger values of *insertions*.
- The *varying* effect of *insertions* tend to have larger variance than other predictors—for every language.
- If we remove *insertions* from \mathcal{M}_3 , the resulting model’s predictive performance is practically indistinguishable from \mathcal{M}_3 ’s.¹⁴

¹⁴Variable selection [18]—an analysis technique that we do not describe in the paper for brevity—also suggests to drop variable *insertions*. See the paper’s replication package for details about this additional analysis of suitability.

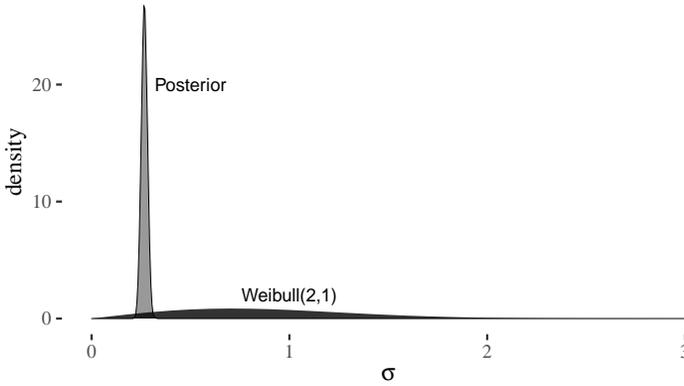


Fig. 9. Comparison of prior Weibull(2, 1) and posterior for project-level intercept $\alpha_{project}$'s standard deviation σ_γ in model \mathcal{M}_3 . The drastic restriction in uncertainty indicates that the data *swamps* the priors.

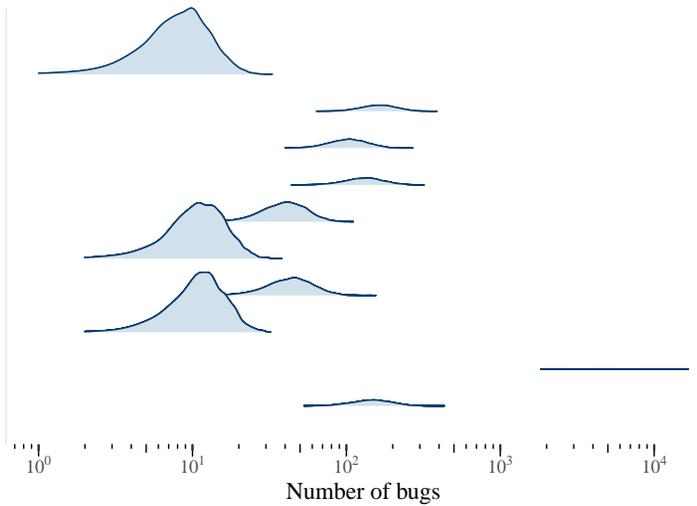


Fig. 10. Posterior predictions of the bug distributions of 10 projects drawn randomly from the posterior of \mathcal{M}_3 . The horizontal axis is logarithmic in base 10. The marked differences in shape and location among the distributions indicate that projects are heterogeneous.

All in all, our analysis indicates that *insertions* does not appear to be a particularly useful predictor, and hence it may not be a reliable measure of code size.

Inter-project variability. Section 4.1 showed that a project-specific intercept—which we introduced in \mathcal{M}_3 —provides better out-of-sample prediction capabilities. The flip side is that several features of the data vary considerably from project to project.

The 1 127 data rows are somewhat sparse among the 729 projects: 64% of all projects appear in a single row; another 26% in two rows. Despite these characteristics, the data *swamps* the priors: it determines a very precise (that is, narrow) posterior distribution of the project-specific intercept $\alpha_{project}$'s standard deviation σ_γ —shown in Figure 9. In other words,

1275 the uncertainty about the contribution of each project to the overall number of bugs is quite
1276 limited: the data characterizes each project's contribution precisely. This does not mean that
1277 the projects' number of bugs is similar; on the contrary, the bug distributions of randomly
1278 drawn projects that we get by simulating from \mathcal{M}_3 's fitted posterior differ considerably
1279 in shape, support, and mean (see Figure 10). In all, project-specific characteristics are an
1280 important and well-defined source of information in the data, which other control variables
1281 cannot fully capture.

1282 4.3 Practical significance

1284 Let us now address the original study's research questions. For brevity, our analysis won't
1285 consider criteria to classify languages ("language classes" such as procedural, functional,
1286 scripting, and so on), projects ("application domains" such as application, database, frame-
1287 work, and so on), or bugs ("bug types" such as algorithm, concurrency, performance, and
1288 so on). These are largely orthogonal to the main focus of the present paper. Instead, we
1289 focus on the key first research question:

1290 **RQ.** *Are some languages more defect-prone than others?*

1291 Our analysis's information is condensed in the fitted model \mathcal{M}_3 , which we can use to
1292 generate a distribution of bugs for every language. In order to do this, we have to pick the
1293 other inputs of the model: the number of commits, insertions, age, and developers of the
1294 hypothetical projects whose number of bugs we are estimating. While, in principle, these
1295 inputs could be any situation or scenario that we want to investigate, it is sensible to start
1296 exploring values that are close to those observed in the data used to fit the model (following
1297 the usual assumption that the sample is representative of the entire population).

1299 **4.3.1 Ranking all languages.** Figure 11 displays the distributions as violin plots for five
1300 combinations of input values: the dataset's *minimum*, *25th percentile*, *median* (50th percentile),
1301 *75th percentile*, and *maximum* number of commits, insertions, age, and developers.¹⁵ Each
1302 plot lists the languages in decreasing order of median predicted number of bugs per project:
1303 from most error prone (left) to least (right).

1304 The plots indicate that the relative ordering of languages can change conspicuously
1305 according to the conditions. For example, C#'s defect proneness is average for projects
1306 with large or median size and age; but it becomes better than average for smaller, younger
1307 projects. In contrast, C++ is less defect prone only in the largest projects, whereas it is the
1308 most or second most defect prone languages for projects of non-maximal size. Similarly, the
1309 relative rank of some language pairs varies considerably: for example, Erlang is less error
1310 prone than Go in the largest projects; the opposite is true in projects of smaller size.

1311 A few languages' ranks fluctuate wildly: Objective-C is among the most defect prone
1312 languages except in small projects, when it is among the least; TypeScript even goes from
1313 least defect prone on large and median projects to most defect prone on the smallest projects.
1314 These jumps are so extreme that they may indicate that the data about these languages is
1315 somewhat inconsistent or at least patchwork. Indeed, TOPLAS's reanalysis reported that
1316 only about a third of the commits classified as TypeScript in FSE's data actually included
1317 TypeScript code; and Ray et al.'s extended version [54] of their original FSE study dropped
1318 several projects classified as TypeScript. We did not further look into Objective-C's data,
1319 despite its high rank fluctuations, because we wanted to use the original data without
1320

1321 ¹⁵Unlike Figure 2, which plots the raw data, Figure 11's simulated projects are directly comparable in terms of
1322 defect proneness, as they only differ in the used programming language.

1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372

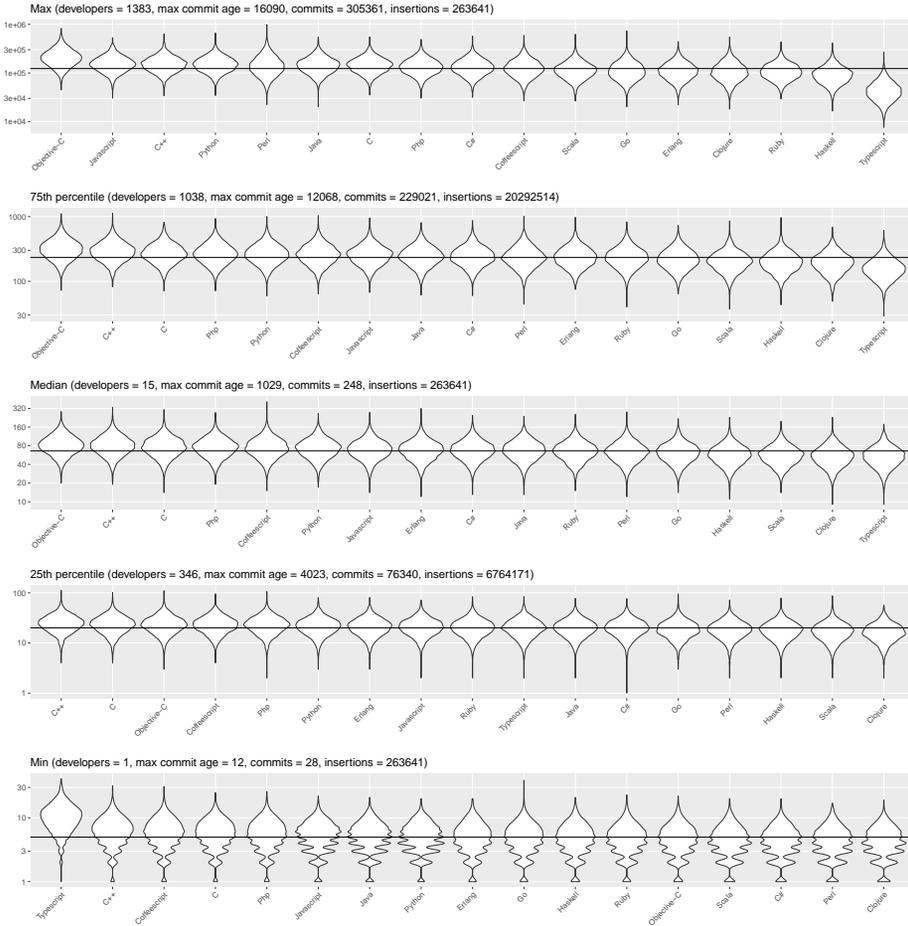


Fig. 11. Violin plots of the distributions of number of bugs per project per language, obtained from the posterior of \mathcal{M}_3 for five simulated scenarios. The plot in each row corresponds to a different scenario: from top to bottom plot, the input variables other than *language* are set to the empirical dataset's *maximum*, *75h percentile* (3rd quartile), *median*, *25th quartile* (1st quartile), and *minimum* values. Languages are sorted, left-to-right in each plot, by decreasing values of the distributions' medians. The vertical axes' scales are logarithmic in base 10. The horizontal line in each plot marks the median number of bugs per project across all languages.

changes. Nevertheless, this is one clear example of how Bayesian analysis can help spot data problems, and hence bolster better substantiated analyses. This observation about the fickle influence of some languages also corroborates the evidence that other project-specific characteristics might weigh comparatively more than the used programming language.

Figure 11 also shows that the bug distributions per language are spread out widely—especially for some languages and especially for projects that are large and long-running—and their ranges extensively overlap. The heterogeneity of project-specific characteristics may also contribute to these features; for example, if projects written in language *X* tend to be on the large side compared to projects written in language *Y*, the uncertainty in *Y*'s

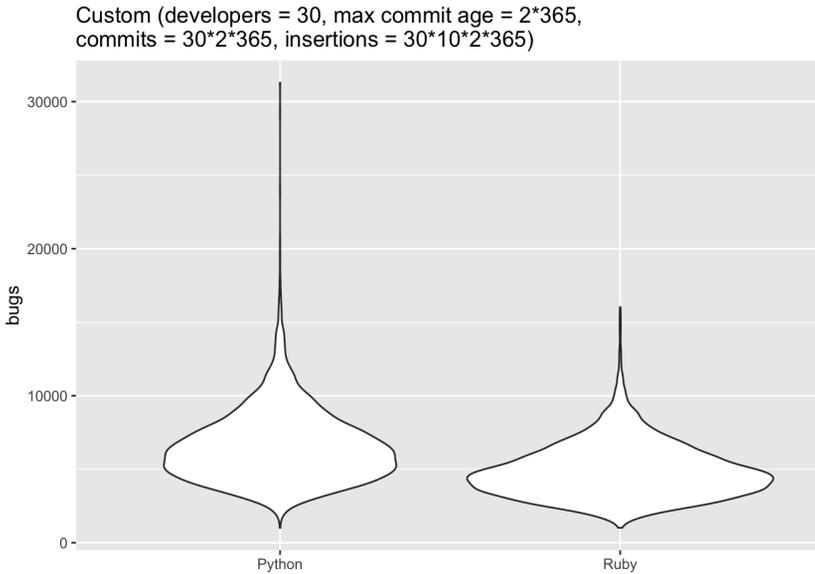


Fig. 12. Violin plots of the bug distributions of Python and Ruby obtained from the posterior of \mathcal{M}_3 . The input variables other than *language* capture a hypothetical project with 30 developers, an age of 2 years, 1 commit per developer every day, and 10 lines added by each developer every day.

error proneness when used for large projects would dominate the comparison with X. This suggests that the data we analyzed does not warrant summarizing the language differences using a single ranking of defect-proneness.

4.3.2 Custom scenarios. While a single ranking of languages according to their absolute defect-proneness would have little practical meaning, we can still zoom in on specific conditions that are relevant *in practice* for a specific project and see what the fitted model can tell us concerning those conditions.

Imagine, for example, we are planning a project that involves around 30 developers who can code in Python or Ruby; we estimate the project will run over 2 years, generating an average of 1 commit and 10 lines inserted per programmer per day. Plugging these numbers (*developers* = 30, *age* = 2×365 , *commits* = $30 \times 1 \times 2 \times 365$, and *insertions* = $30 \times 10 \times 2 \times 365$) into the fitted model \mathcal{M}_3 , we get the estimated bug distributions for Python and Ruby shown in Figure 12. In this scenario, Python tends to be worse (more bugs) than Ruby, since the latter's distribution has a lower mean, a shorter tail towards high number of bugs, and more mass around lower values.

Whether this evidence is sufficiently strong to decide to choose one language over another depends on myriad other factors that are incidental, such as the availability of programmers familiar with one language, the cost of training new ones, the usability of the programming language for the project at hand, and so on. Whatever the practical constraints and requirements may be, the fitted model can help us meet them by providing estimates complete with a quantification of their uncertainty. Realistically, any estimate about the size and development time of a project is also likely to be somewhat uncertain; therefore, we

would run *multiple* simulations and weigh the evidence summarized by each one against the confidence we have in the corresponding scenario occurring.

More generally, the results of a principled Bayesian analysis facilitate a quantitatively accurate transfer of knowledge to practitioners and other researchers. Rather than relying only on overly broad conclusions about the impact of different programming languages, using simulations of custom scenarios drives follow-up work more precisely: a practitioner can judge whether the uncertainty in a specific comparison is too large to base a decision on it; a researcher can decide whether more data is needed to claim more general conclusions.

4.3.3 Statistical significance. Our analysis so far has focused on concrete scenarios defined in terms of tangible measures in the data domain—such as number of bugs and project age. In contrast, widespread statistical practices (mostly of a frequentist flavor) try to answer research questions by analyzing *statistical significance*, which measures generic characteristics of a statistical model.

In a standard regression analysis, one usually assesses the *statistical significance* of each coefficient in the model (also called “effect”) by checking whether it differs from zero with a certain probability. For example, we could compute the distribution of the estimate of coefficient α_{language} for every language. If α_X is negative with, say, 95% probability, we would conclude that language X is associated with fewer bugs than average with that probability; in other words, X is “statistically significantly” less error prone than other languages.

FSE and TOPLAS both proceed in such a way, but using frequentist coefficient estimates instead of a posterior probability distribution on their models—which are similar to our \mathcal{M}_2 —leading to their findings about which languages are more error prone than others.¹⁶ What about model \mathcal{M}_3 fitted on the same data? The 95% probability intervals of α_{language} include the origin for *every language*, except TypeScript whose α_{language} is strictly positive—but, as we have commented above, the uncertainty about TypeScript’s data puts any results about this language on shaky grounds. Overall, the canonical analysis of statistical significance is just inconclusive on our model.

To some extent, this outcome is a side effect of \mathcal{M}_3 ’s greater complexity over simpler models. There is a trade-off between the complexity of a model (which brings greater expressiveness and better predictive performance) and its interpretability. The criteria we used to choose \mathcal{M}_3 over the simpler \mathcal{M}_2 ensure that the former’s additional complexity is justified by its much better effectiveness. However, a simple interpretation is no longer feasible: \mathcal{M}_3 includes slope coefficients that also vary with each language, as well as a project-level contribution; how each language-specific term interacts with the others is not something that can be simply estimated with a single coefficient independent of the predictors’ values.

We should appreciate that this is more a feature than it is a limitation. While mathematically simple models are nice to have, not all data analysis problems can be addressed with a basic model. Bayesian analysis techniques do not just support fitting complex models but provide the means to handle their complexity and to perform a convincing analysis without resorting to formulaic measures of “significance”. The individual model characteristics are not easy to interpret in isolation, so that we are forced to interpret the model by providing concrete conditions—the number of commits, age, and so on—which ground our generic research question onto scenarios that are realistic and meaningful for our purposes. In other

¹⁶This is explained in detail in TOPLAS’s repetition [6, § 2.2.3], which uses FSE’s model; TOPLAS’s reanalysis [6, § 4.2.1] suggests a different statistical measure.

1471 words, it may be cumbersome to reason about statistical significance in a Bayesian model
 1472 but it is always natural to reason about *practical significance*—which is what matters most in
 1473 the end to answer our research questions.¹⁷

1474 **Model interpretability.** The focus on practical significance follows from the specific
 1475 research problem we considered: answering the question of whether some programming
 1476 languages are more prone to defects requires precise *predictions* about the defect-proneness
 1477 of projects written in different programming languages. On the other hand, the inter-
 1478 pretability of a statistical model may become crucial when targeting other kinds of research
 1479 questions. In these scenarios, Bayesian models can still be practically effective. A relatively
 1480 complex model like \mathcal{M}_3 is not easy to interpret *directly*; that is, we cannot easily and unam-
 1481 biguously assign a direct interpretation to each individual fitted parameter. However, it is
 1482 amenable to interpret *indirectly*: we formulate some scenarios in terms of model variables,
 1483 and then we simulate those scenarios on the fitted model. Interpreting the simulation’s
 1484 results is how we indirectly interpret the model’s characteristics. Ultimately, a key feature of
 1485 Bayesian data analysis techniques is what makes these analyses so flexible: we can combine
 1486 (indirect) interpretability and predictive capabilities because the outcome of an analysis is a
 1487 (sampled) posterior probability distribution—a rich and actionable source of information.

1488 **4.3.4 Effect sizes.** Section 4.3.1 demonstrated that the fault proneness of a language over
 1489 another strongly depends on the conditions in which the languages are to be used. If we
 1490 have specific scenarios in mind, we can just simulate those as discussed in Section 4.3.2.

1491 Another approach is to compare languages pairwise by simulating their performance on
 1492 a population that resembles the observed data. Since the comparisons are quantitative—in
 1493 the form of derived distributions—they can be seen as an effect size, but relative to each
 1494 language pair instead of absolute for all languages at once.

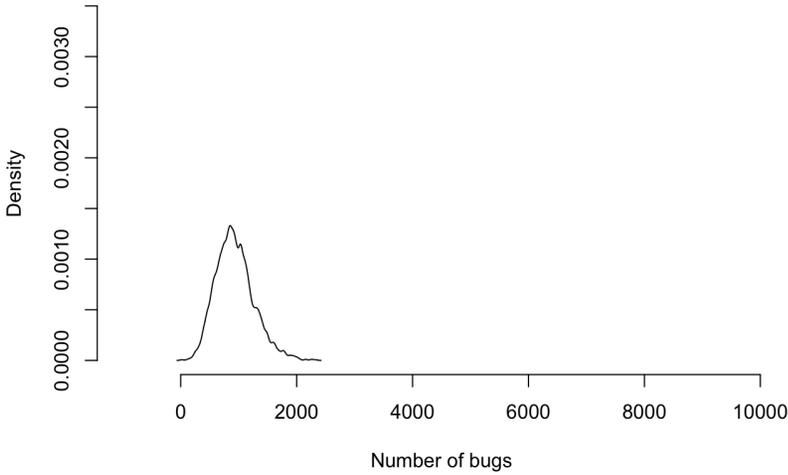
1495 As usual, simulations are derived from the posterior, which entails that there is no
 1496 multiple comparisons problem [45]: all information is encoded jointly by the posterior; the
 1497 pairwise comparisons are just projections of some of that information. For the same reason,
 1498 we do not have to commit to a certain way of comparing languages when we build the
 1499 model (for example, by choosing how to encode contrasts): we just select the “best” model
 1500 according to its performance, and then derive all the information we are interested in from
 1501 the model fitted on the data.

1502 Concretely, take two languages ℓ_1 and ℓ_2 that we want to compare for bug proneness. For
 1503 every data point d in the empirical data, we set, in the posterior, all predictors except the
 1504 language to their values in d . Then, we simulate the distribution of the expected difference
 1505 $bugs_{\ell_1} - bugs_{\ell_2}$ in bugs produced when using one language over the other.¹⁸

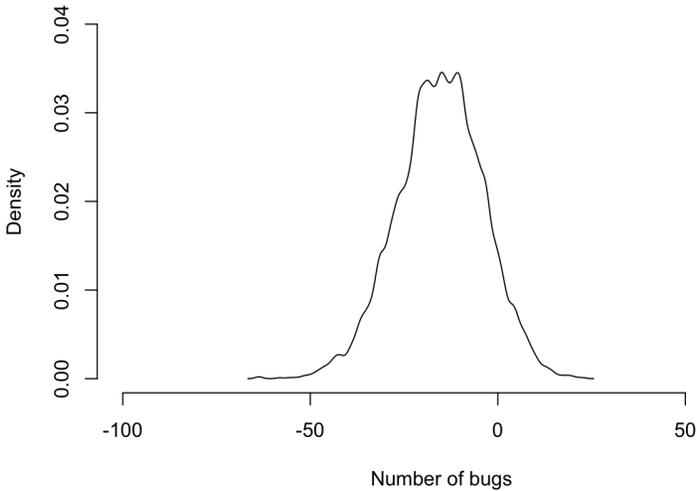
1506 Figure 13 plots the distributions comparing two pairs of languages, which we selected
 1507 to demonstrate qualitatively different outcomes of the pairwise comparisons. The distri-
 1508 bution of $bugs_{C\#} - bugs_C$ in Figure 13a covers only nonnegative values, which means C#
 1509 was consistently more fault prone than C. The distribution of $bugs_{CoffeeScript} - bugs_{Go}$ in
 1510 Figure 13b covers negative values more often than positive ones, denoting that Go tended
 1511 to be more fault prone. Precisely, we can compute that CoffeeScript was more error prone
 1512 than Go only around 8.5% of the times.

1514
 1515 ¹⁷In related work, we discuss in greater detail methods to analyze practical significance based on Bayesian data
 1516 analysis [65].

1517 ¹⁸We can compute the absolute difference in number of bugs because the difference is between samples where
 1518 the language is the only project characteristic that changes.

C# vs. C

(a) Posterior distribution of the difference $bugs_{C\#} - bugs_C$ when predictors are set to the same values as in the empirical data. This indicates that C# is consistently more fault-prone than C, since it leads to more bugs.

Coffeescript vs. Go

(b) Posterior distribution of the difference $bugs_{CoffeeScript} - bugs_{Go}$ when predictors are set to the same values as in the empirical data. This indicates that Coffeescript is somewhat less fault-prone than Go.

Fig. 13. Probability distributions of the difference in bug proneness between pairs of languages according to the posterior distribution with population data.

1569 In the end, we did not really answer the original research question—not with a definitive,
1570 straightforward answer at least. Instead, our analysis identified sources of uncertainty in the
1571 data, provided means of simulating custom scenarios, and compared pairs of languages in
1572 conditions similar to the collected data's. This is a solid basis to understand what questions
1573 can and cannot be answered by the data, and to plan follow-up data collections and analyses
1574 that zero in on understanding specific outcomes.

1575

1576

1577

4.4 Planning the next study

1578

1579

1580

1581

1582

1583

1584

1585

1586

1587

1588

1589

1590

1591

1592

1593

1594

1595

1596

1597

1598

1599

1600

1601

1602

1603

1604

1605

1606

1607

1608

1609

1610

1611

5 RELATED WORK

1612

1613

1614

1615

1616

1617

We discuss related work in three areas, broadly connected to the paper's contributions: Section 5.1 briefly reviews some widely-used statistical models other than those we deployed in this paper; Section 5.2 summarizes other work about statistical analysis guidelines (for frequentist and for Bayesian techniques); and Section 5.3 outlines the state of replication studies in empirical software engineering research.

1618 5.1 Statistical data modeling and analysis

1619 As we remarked in Section 1.3.1, this paper’s guidelines are largely independent of the
1620 specific features of the chosen statistical models. All our examples used (generalized/hier-
1621 archical) linear *regressive* models—the workhorse of statistical analysis [29]. Besides their
1622 great flexibility, another practical advantage is that they have been widely used also with
1623 frequentist statistics; therefore, they offer a convenient bridge for a gradual transition to
1624 Bayesian statistics.

1625 Nevertheless, other classes of statistical models can be used for similar analyses. One
1626 alternative, broad class of statistical models are so-called *graphical* models [38], which use
1627 graphs to encode the probabilistic relations between variables. Bayesian networks [59] are
1628 arguably the best-known kind of graphical models, which can be used both directly as
1629 probabilistic classifiers [17, 40] but have also become the basis to encode *causal* relations
1630 that go beyond mere correlations [48]. The expressive power of Bayesian networks and
1631 hierarchical regressive models significantly overlap: among other things, one can encode
1632 a regressive model as a Bayesian network, and then use network’s fitting techniques to
1633 analyze it [36]; conversely, one can encode a Bayesian network as a hierarchical regressive
1634 model [59, Ch. 5], and then apply similar analysis techniques as those we demonstrated in
1635 the paper.

1636 Machine learning toolkits such as Weka [72] provide a convenient way of experimenting
1637 with a wide variety of classical statistical analysis models, which have been frequently used
1638 in the analysis of software engineering empirical data [43, 73] and provide additional ser-
1639 viceable classes of statistical models. While our paper’s guidelines would remain applicable,
1640 at least at a high level, to compare other, widely different statistical models, doing so in
1641 practice may require developing new analysis techniques or extending existing ones. In
1642 particular, some information criteria—which are used for Bayesian model comparison as
1643 discussed in Section 2.5—are only applicable to statistical models that can provide multiple
1644 samples from a posterior when the fitted model is used for prediction [67]. Obviously,
1645 generalizing the model comparison criteria (and the other techniques for model analysis)
1646 so that they are applicable to all inductive machine learning approaches falls outside this
1647 paper’s scope.

1648 It is interesting that several modern machine learning algorithms, such as deep neural
1649 networks and active learning, are applicable both in a frequentist [74] and in a Bayesian [24,
1650 46, 69] context. Our paper’s guidelines could remain broadly useful for frequentist models,
1651 but they do not cover *online* approaches (for example, active learning), where each iteration
1652 of a statistical analysis influences which additional data is collected. Extending some of our
1653 guidelines to online approaches is an interesting direction for future work. At the same time,
1654 the increasingly recognized value of practices such as pre-registered studies [14] suggests
1655 that the *offline* analysis of fixed, previously collected, datasets will remain an important and
1656 common approach in software engineering empirical research.

1658 5.2 Guidelines about statistical analysis

1659 **Bayesian data analysis guidelines.** The last decade’s progress in algorithms and tools for
1660 Bayesian data analysis has been impressive [13, 25, 41, 50] but, without practical support,
1661 it is not sufficient to promote widespread usage in the empirical sciences. Recent work
1662 about developing guidelines to apply Bayesian data analysis techniques [23, 31, 57], which
1663 Section 2 summarized in a form amenable to software engineering empirical research, has
1664 been trying to close this gap. While these proposals differ in their intended audience and
1665
1666

1667 level of detail, they all build on the basic view [26] that an analysis should go through
1668 multiple models, refine them in several iterations, and compare them. Then, Gabry et al. [23]
1669 focus on visualization and how to use it throughout a workflow; Schad et al. [57], instead,
1670 introduce quantitative checks and illustrate them for a specific scientific area (the cognitive
1671 sciences). Our guidelines combine elements from both [23, 57] but illustrate them in a way
1672 that is amenable to empirical software engineering research practices. Very recently, Gelman
1673 et al. [31] presented an early draft of a book that will further refine some of these Bayesian
1674 analysis workflows and guidelines. Also very recently, van de Schoot et al. [66] published
1675 an accessible primer on Bayesian statistics and modeling for scientists.

1676 **Guidelines on using statistics in empirical software engineering.** Over the years, several
1677 guidelines for using statistics in empirical software engineering have been proposed—all
1678 of them focusing on frequentist statistics, which remain the norm in empirical software
1679 engineering [16]. Arcuri and Briand [3] focus on analyzing experiments with randomized
1680 algorithms, and highlight the importance of checking the assumptions of each statistical
1681 significance test. They also advocate for extensively using non-parametric statistical tests
1682 and effect size measures. Menzies and Shepperd [44] catalog “bad smells” in data analytics
1683 studies and discuss remedies to excise them. Among the techniques they recommend are
1684 up-front power analysis, reporting effect sizes and confidence limits, and using robust
1685 statistics and sensitivity analysis. A recent literature review of ours [16] found evidence
1686 of a positive impact of such empirical guidelines on the maturity of statistical practice in
1687 empirical software engineering research: statistical testing, non-parametric tests, and effect
1688 sizes have all been increasingly used in the field over the last 5–10 years.

1689 5.3 Replication in software engineering research

1691 Recent years have finally seen replication studies become more popular in software en-
1692 gineering research. Nevertheless, Da Silva et al.’s systematic literature review found that
1693 internal replications (done by the same authors as the original study) are still much more
1694 common than external replications (done by an independent group of authors) [15]. Un-
1695 surprisingly, Bezerra et al.’s related literature review found that internal replications are
1696 much more likely to confirm the results of the replicated study than external replications [7],
1697 and used this result to question the value of replications compared to meta-analyses. Both
1698 literature reviews found hardly any examples of *reanalyses* (replications limited to data
1699 analysis); similarly, a taxonomy for replications in software engineering does not explicitly
1700 mention reanalysis [4].

1701 In fact, we tried searching for “*reanalysis + software engineering*” in publication databases
1702 and found very few relevant hits—mostly papers revisiting qualitative data such as in-
1703 terview transcripts, and reanalyzing them to address new questions or theories. As one
1704 example, Bjarnason et al. [8] developed a new theory by reanalyzing interview transcripts
1705 from an earlier study of theirs. In contrast, Tantithamthavorn et al. [63] revised a meta-
1706 analysis of machine learning in software defect prediction [60] and found that several
1707 predictor variables of the original study were co-linear. Based on a reanalysis of a subset
1708 of the same data, they also questioned some of the original results and implications. This
1709 criticism was later disputed, on statistical grounds, by the original study’s authors [61].
1710 Our previous work about using Bayesian analysis in empirical software engineering also
1711 performed reanalyses of previous studies using Bayesian techniques [21, 65].¹⁹ Another
1712

1713 ¹⁹Our previous work targets various applications of Bayesian statistics such as analyzing practical signifi-
1714 cance [65] and dealing with missing data [64]; the case studied developed there also follow some of the guidelines
1715

1716 noticeable external reanalysis is of course Berger et al. [6]’s of Ray et al. [55], which we
1717 summarized in Section 1.

1718 Overall, reanalyses of software engineering data remain uncommon—especially com-
1719 pared to other scientific areas where they are widespread forms of publication, including
1720 those using Bayesian statistics (for example, in astronomy [33] and medicine [5]).

1721

1722 6 CONCLUSIONS

1723 Reaping the benefits of Bayesian statistics requires more than powerful analysis techniques
1724 and tools. In this paper, we presented practical guidelines to build, check, and analyze a
1725 Bayesian statistical model that summarize recently developed suggestions brought forward
1726 by prominent statisticians and cast them in a format that is amenable to empirical software
1727 engineering research.

1728 We then applied the guidelines to analyze a large dataset of GitHub projects that was
1729 previously used to study the impact of programming languages on code quality [55].
1730 This study was later criticized by a reproduction attempt that failed to confirm some
1731 of the originally claimed results [6]. Our reanalysis using Bayesian statistics identified
1732 some shortcomings of the data that also emerged in the reproduction attempt (such as the
1733 large uncertainty associated with data for programming languages such as TypeScript)
1734 and pointed to other possible effects that were not fully accounted for by the frequentist
1735 models of the previous studies [6, 55] (such as the disproportionate differences that are
1736 project-specific rather than language-specific). Moving on to the previous studies’ main
1737 research question (“Are some languages more defect-prone than others?”), our Bayesian
1738 model lent itself to evaluating the effect of programming languages in different concrete
1739 scenarios rather than in terms of generic “statistical significance”. We found that the impact
1740 of programming languages can vary considerably with other contextual conditions, and
1741 hence the original research question does not admit a simple, generally valid answer—at
1742 least not with the analyzed data.

1743 Throughout our reanalysis, a key advantage of Bayesian techniques was that they can be
1744 used to *quantify* any derived measures of interest, as well as the *uncertainty* that comes with
1745 each measure. Such capabilities are useful not only to infer results in each study, but also to
1746 present and share them in a robust way with other researchers and practitioners. A Bayesian
1747 quantitative framework focused on practical significance can also help plan the next studies
1748 in a research area—thus steadying the long-term progress of software engineering empirical
1749 research and enhancing its broader impact.

1750

1751 7 ACKNOWLEDGEMENTS

1752 The computations were enabled by resources provided by the Swedish National Infrastruc-
1753 ture for Computing (SNIC), partially funded by the Swedish Research Council through
1754 grant agreement no. 2018–05973. We thank Jonah Gabry for reading an earlier draft of this
1755 paper and providing helpful comments.

1756

1757 REFERENCES

1758 [1] Balazs Aczel, Rink Hoekstra, Andrew Gelman, Eric-Jan Wagenmakers, Irene G. Klugkist, Jeffrey N. Rouder,
1759 Joachim Vandekerckhove, Michael D. Lee, Richard D. Morey, Wolf Vanpaemel, Zoltan Dienes, and Don van

1760

1761 that we explicitly and specifically present in the present paper. Therefore, they provide further examples of
1762 applications of the guidelines to analyze software engineering empirical data—especially in their replication
1763 packages, since the papers’ presentations have a different focus.

1764

- 1765 Ravenzwaaij. 2020. Discussion points for Bayesian inference. *Nature Human Behaviour* 4, 6 (2020), 561–563.
1766 <https://doi.org/10.1038/s41562-019-0807-z>
- 1767 [2] Valentin Amrhein, Sander Greenland, and Blake McShane. 2019. Scientists rise up against statistical signifi-
1768 cance. *Nature* 567 (2019), 305–307.
- 1769 [3] Andrea Arcuri and Lionel Briand. 2011. A practical guide for using statistical tests to assess randomized
1770 algorithms in software engineering. In *2011 33rd International Conference on Software Engineering (ICSE)*. IEEE
1771 Computer Society, Hawaii, USA, 1–10. <https://doi.org/10.1145/1985793.1985795>
- 1772 [4] Maria Teresa Baldassarre, Jeffrey C. Carver, Oscar Dieste, and Natalia Juristo Juzgado. 2014. Replication
1773 types: Towards a shared taxonomy. In *18th International Conference on Evaluation and Assessment in Software
1774 Engineering, EASE '14*, Martin J. Shepperd, Tracy Hall, and Ingunn Myrvtveit (Eds.). ACM, London, UK,
1775 18:1–18:4. <https://doi.org/10.1145/2601248.2601299>
- 1776 [5] Philip M. W. Bath. 2007. Can we improve the statistical analysis of stroke trials? Statistical re-analysis of
1777 functional outcomes in stroke trials. *Stroke* 38, 6 (2007), 1911–1915. [https://doi.org/10.1161/STROKEAHA.106.
1778 474080](https://doi.org/10.1161/STROKEAHA.106.474080)
- 1779 [6] Emery D. Berger, Celeste Hollenbeck, Petr Maj, Olga Vitek, and Jan Vitek. 2019. On the Impact of Programming
1780 Languages on Code Quality: A Reproduction Study. *ACM Transactions on Programming Languages and Systems*
1781 41, 4 (2019), 21:1–21:24. <https://doi.org/10.1145/3340571>
- 1782 [7] Roberta M. M. Bezerra, Fabio Q. B. da Silva, Anderson M. Santana, Cleyton V. C. de Magalhães, and
1783 Ronnie E. S. Santos. 2015. Replication of Empirical Studies in Software Engineering: An Update of a
1784 Systematic Mapping Study. In *2015 ACM/IEEE International Symposium on Empirical Software Engineering and
1785 Measurement, ESEM 2015, Beijing, China, October 22-23, 2015*. IEEE Computer Society, Beijing, China, 132–135.
1786 <https://doi.org/10.1109/ESEM.2015.7321213>
- 1787 [8] Elizabeth Bjarnason, Kari Smolander, Emelie Engström, and Per Runeson. 2016. A theory of distances in
1788 software engineering. *Information and Software Technology* 70 (2016), 204–219. [https://doi.org/10.1016/j.
1789 infsof.2015.05.004](https://doi.org/10.1016/j.infsof.2015.05.004)
- 1790 [9] George E. P. Box. 1976. Science and Statistics. *J. Amer. Statist. Assoc.* 71, 356 (1976), 791–799. <https://doi.org/10.1080/01621459.1976.10480949>
- 1791 [10] George E. P. Box. 1979. Robustness in the Strategy of Scientific Model Building. In *Robustness in Statistics*.
1792 Academic Press, 201–236. <https://doi.org/10.1016/B978-0-12-438150-6.50018-2>
- 1793 [11] Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. 2011. *Handbook of Markov Chain Monte Carlo*.
1794 CRC press, Florida, USA.
- 1795 [12] P. C. Bürkner. 2017. brms: An R Package for Bayesian Multilevel Models using Stan. *Journal of Statistical
1796 Software* 80, 1 (2017), 1–28. [https://doi.org/doi.org/10.18637/jss.v080.i01](https://doi.org/10.18637/jss.v080.i01)
- 1797 [13] Bob Carpenter, Andrew Gelman, Matthew D Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt,
1798 Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. 2017. Stan: A probabilistic programming language.
1799 *Journal of statistical software* 76, 1 (2017), 1–32.
- 1800 [14] Christopher D Chambers, Zoltan Dienes, Robert D McIntosh, Pia Rotshtein, and Klaus Willmes. 2015.
1801 Registered reports: realigning incentives in scientific publishing. *Cortex* 66 (2015), A1–A2.
- 1802 [15] Fabio QB Da Silva, Marcos Suassuna, A César C França, Alicia M Grubb, Tatiana B Gouveia, Cleviton VF
1803 Monteiro, and Igor Ebrahim dos Santos. 2014. Replication of empirical studies in software engineering
1804 research: A systematic mapping study. *Empirical Software Engineering* 19, 3 (2014), 501–557. [https://doi.org/
1805 10.1007/s10664-012-9227-7](https://doi.org/10.1007/s10664-012-9227-7)
- 1806 [16] Francisco Gomes de Oliveira Neto, Richard Torkar, Robert Feldt, Lucas Gren, Carlo A. Furia, and Ziwei Huang.
1807 2019. Evolution of statistical analysis in empirical software engineering research: Current state and steps
1808 forward. *Journal of Systems and Software* 156 (2019), 246–267. <https://doi.org/10.1016/j.jss.2019.07.002>
- 1809 [17] Tapajit Dey and Audris Mockus. 2020. Deriving a usage-independent software quality metric. *Empir. Software
1810 Eng.* 25 (2020), 1596–1641. <https://doi.org/10.1007/s10664-019-09791-w>
- 1811 [18] Michaela Dvorzak and Helga Wagner. 2016. Sparse Bayesian modelling of underreported count data.
1812 *Statistical Modelling* 16, 1 (2016), 24–46. <https://doi.org/10.1177/1471082X15588398>
- 1813 [19] Robert Feldt and Ana Magazinius. 2010. Validity threats in empirical software engineering research—an initial
1814 survey. In *International Conference on Software Engineering and Knowledge Engineering*. 374–379.
- [20] R. A. Fisher. 1925. *Statistical Methods for Research Workers* (1 ed.). Oliver and Boyd.
- [21] Carlo A. Furia, Robert Feldt, and Richard Torkar. 2019. Bayesian data analysis in empirical software
engineering research. *IEEE Transactions on Software Engineering* -, - (2019), 1–1. [https://doi.org/10.1109/
TSE.2019.2935974](https://doi.org/10.1109/TSE.2019.2935974)
- [22] Carlo A. Furia, Richard Torkar, and Robert Feldt. 2021. *Replication Package*. [https://doi.org/10.5281/zenodo.
4472963](https://doi.org/10.5281/zenodo.4472963)

- 1814 [23] Jonah Gabry, Daniel Simpson, Aki Vehtari, Michael Betancourt, and Andrew Gelman. 2019. Visualization in
 1815 Bayesian workflow. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 182, 2 (2019), 389–402.
 1816 <https://doi.org/10.1111/rssa.12378>
- 1817 [24] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. 2017. Deep Bayesian Active Learning with Image Data. In
 1818 *Proceedings of the 34th International Conference on Machine Learning (ICML)*. JMLR.org, 1183–1192.
- 1819 [25] Hong Ge, Kai Xu, and Zoubin Ghahramani. 2018. Turing: a language for flexible probabilistic inference. In
 1820 *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 1682–1690. <http://proceedings.mlr.press/v84/ge18b.html>
- 1821 [26] Andrew Gelman. 2004. Exploratory data analysis for complex models. *Journal of Computational and Graphical
 1822 Statistics* 13, 4 (2004), 755–779.
- 1823 [27] Andrew Gelman. 2016. Bayesian statistics: What’s it all about? [https://statmodeling.stat.columbia.edu/
 1824 2016/12/13/bayesian-statistics-whats/](https://statmodeling.stat.columbia.edu/2016/12/13/bayesian-statistics-whats/).
- 1825 [28] Andrew Gelman and Jennifer Hill. 2007. *Data analysis using regression and multilevel/hierarchical models*.
 1826 Vol. Analytical methods for social research. Cambridge University Press, Cambridge, UK. xxii, 625 p pages.
- 1827 [29] Andrew Gelman, Jennifer Hill, and Aki Vehtari. 2020. *Regression and other stories*. Cambridge University
 1828 Press, Cambridge, UK. <https://books.google.se/books?id=SZFkzQEACAAJ>
- 1829 [30] Andrew Gelman, Daniel Simpson, and Michael Betancourt. 2017. The prior can often only be understood in
 1830 the context of the likelihood. *Entropy* 19, 10 (Oct 2017), 555. <https://doi.org/10.3390/e19100555>
- 1831 [31] Andrew Gelman, Aki Vehtari, Daniel Simpson, Charles C. Margossian, Bob Carpenter, Yuling Yao, Lau-
 1832 ren Kennedy, Jonah Gabry, Paul-Christian Bürkner, and Martin Modrák. 2020. Bayesian workflow.
 1833 arXiv:2011.01808 [stat.ME]
- 1834 [32] Yossi Gil and Gal Lalouche. 2017. On the correlation between size and metric validity. *Empirical Software
 1835 Engineering* 22, 5 (Oct. 2017), 2585–2611. <https://doi.org/10.1007/s10664-017-9513-5>
- 1836 [33] Philip C Gregory. 2011. Bayesian re-analysis of the Gliese 581 exoplanet system. *Monthly Notices of the Royal
 1837 Astronomical Society* 415, 3 (2011), 2523–2545.
- 1838 [34] Heiko Haller and Stefan Kraus. 2002. Misinterpretations of significance: A problem students share with their
 1839 teachers? *Methods of Psychological Research* 7, 1 (2002), 1–20.
- 1840 [35] Rink Hoekstra, Richard D. Morey, Jeffrey N. Rouder, and Eric-Jan Wagenmakers. 2014. Robust misinterpretation
 1841 of confidence intervals. *Psychon. Bull. Rev.* (2014), 1157–1164.
- 1842 [36] C. H. Jackson, N. G. Best, and S. Richardson. 2008. Bayesian graphical models for regression on multiple data
 1843 sets with different variables. *Biostatistics* 10, 2 (11 2008), 335–351. [https://doi.org/10.1093/biostatistics/
 1844 kxn041](https://doi.org/10.1093/biostatistics/kxn041)
- 1845 [37] Edwin T. Jaynes. 2003. *Probability theory: The logic of science*. Cambridge University Press, Cambridge.
- 1846 [38] Michael I. Jordan. 2004. Graphical Models. *Statist. Sci.* 19, 1 (2004), 140–155. [https://doi.org/10.1214/
 1847 088342304000000026](https://doi.org/10.1214/088342304000000026)
- 1848 [39] JuliaStan [n.d.]. Stan.jl. <https://mc-stan.org/users/interfaces/julia-stan>.
- 1849 [40] Andrey Krutauz, Tapajit Dey, Peter C. Rigby, and Audris Mockus. 2020. Do code review measures explain the
 1850 incidence of post-release defects? *Empir. Softw. Eng.* 25, 5 (2020), 3323–3356. [https://doi.org/10.1007/s10664-
 1851 020-09837-4](https://doi.org/10.1007/s10664-020-09837-4)
- 1852 [41] David Lunn, David Spiegelhalter, Andrew Thomas, and Nicky Best. 2009. The BUGS project: Evolution,
 1853 critique and future directions. *Statistics in medicine* 28, 25 (2009), 3049–3067.
- 1854 [42] Richard McElreath. 2020. *Statistical rethinking: A Bayesian course with examples in R and Stan* (2 ed.). CRC press,
 1855 Florida, USA.
- 1856 [43] Tim Menzies, Jeremy Greenwald, and Art Frank. 2007. Data Mining Static Code Attributes to Learn Defect
 1857 Predictors. *IEEE Trans. Software Eng.* 33, 1 (2007), 2–13. <https://doi.org/10.1109/TSE.2007.256941>
- 1858 [44] Tim Menzies and Martin Shepperd. 2019. “Bad smells” in software analytics papers. *Information and Software
 1859 Technology* 112 (2019), 35–47. <https://doi.org/10.1016/j.infsof.2019.04.005>
- 1860 [45] Rupert G. Miller. 1981. *Simultaneous statistical inference* (2nd ed.). Springer-Verlag, Berlin, Heidelberg.
- 1861 [46] Salman Mohamadi and Hamidreza Amindavar. 2020. Deep Bayesian Active Learning, A Brief Survey on
 1862 Recent Advances. arXiv:2012.08044 [cs.LG]
- [47] Radford M. Neal. 1996. *Bayesian Learning for Neural Networks*. Springer-Verlag, Berlin, Heidelberg.
- [48] Judea Pearl. 2009. *Causality: Models, Reasoning and Inference* (2nd ed.). Cambridge University Press, USA.
- [49] Juho Piironen, Markus Paasiniemi, and Aki Vehtari. 2020. Projective inference in high-dimensional problems: Prediction and feature selection. *Electronic Journal of Statistics* 14, 1 (2020), 2155–2197. <https://doi.org/10.1214/20-EJS1711>
- [50] Martin Plummer. 2003. JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling. In *Proceedings of the 3rd International Workshop on Distributed Statistical Computing: Workshop on Distributed*

- 1863 *Statistical Computing*. Achim Zeileis, Vienna, Austria, 10 pages.
- 1864 [51] Martyn Plummer. 2003. JAGS: A Program for Analysis of Bayesian Graphical Models Using Gibbs Sampling. In *Proceedings of the 3rd International Workshop on Distributed Statistical Computing (DSC)*.
- 1865 [52] PyStan [n.d.]. Stan.jl. <https://pystan.readthedocs.io/en/latest/>.
- 1866 [53] Paul Ralph and Ewan Tempero. 2018. Construct Validity in Software Engineering Research and Software Metrics. In *Proceedings of the 22nd International Conference on Evaluation and Assessment in Software Engineering 2018 (Christchurch, New Zealand) (EASE'18)*. Association for Computing Machinery, New York, USA, 13–23. <https://doi.org/10.1145/3210459.3210461>
- 1867 [54] Baishakhi Ray, Daryl Posnett, Premkumar Devanbu, and Vladimir Filkov. 2017. A large-scale study of programming languages and code quality in GitHub. *Commun. ACM* 60, 10 (Sept. 2017), 91–100. <https://doi.org/10.1145/3126905>
- 1870 [55] Baishakhi Ray, Daryl Posnett, Vladimir Filkov, and Premkumar Devanbu. 2014. A large scale study of programming languages and code quality in Github. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering (Hong Kong, China) (FSE 2014)*. Association for Computing Machinery, New York, NY, USA, 155–165. <https://doi.org/10.1145/2635868.2635922>
- 1874 [56] Bin Xin Ru, Mark McLeod, Diego Granzio, and Michael A. Osborne. 2018. Fast information-theoretic Bayesian optimisation. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018 (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer G. Dy and Andreas Krause (Eds.). PMLR, Stockholm, Sweden, 4381–4389. <http://proceedings.mlr.press/v80/ru18a.html>
- 1875 [57] Daniel J. Schad, Michael Betancourt, and Shravan Vasishth. 2020. Toward a principled Bayesian workflow in cognitive science. *Psychological methods* -, - (June 2020), -. <https://doi.org/10.1037/met0000275>
- 1876 [58] Maximilian Scholz and Richard Torkar. 2020. An empirical study of Linespots: A novel past-fault algorithm. *arXiv e-prints* -, -, Article arXiv:2007.09394 (July 2020), 19 pages. arXiv:2007.09394 [cs.SE]
- 1877 [59] Marco Scutari and Jean-Baptiste Denis. 2021. *Bayesian Networks (with Examples in R)* (2 ed.). Chapman and Hall/CRC.
- 1878 [60] Martin Shepperd, David Bowes, and Tracy Hall. 2014. Researcher bias: The use of machine learning in software defect prediction. *IEEE Transactions on Software Engineering* 40, 6 (2014), 603–616. <https://doi.org/10.1109/TSE.2014.2322358>
- 1880 [61] Martin Shepperd, Tracy Hall, and David Bowes. 2017. Authors' Reply to "Comments on 'Researcher Bias: The Use of Machine Learning in Software Defect Prediction'". *IEEE Transactions on Software Engineering* 44, 11 (2017), 1129–1131. <https://doi.org/10.1109/TSE.2017.2731308>
- 1881 [62] Sean Talts, Michael Betancourt, Daniel Simpson, Aki Vehtari, and Andrew Gelman. 2018. Validating Bayesian inference algorithms with simulation-based calibration. *arXiv preprint arXiv:1804.06788* (2018).
- 1882 [63] Chakkrit Tantithamthavorn, Shane McIntosh, Ahmed E Hassan, and Kenichi Matsumoto. 2016. Comments on "Researcher bias: the use of machine learning in software defect prediction". *IEEE Transactions on Software Engineering* 42, 11 (2016), 1092–1094. <https://doi.org/10.1109/TSE.2016.2553030>
- 1883 [64] Richard Torkar, Robert Feldt, and Carlo A. Furia. 2020. Bayesian data analysis in empirical software engineering—The case of missing data. In *Contemporary Empirical Methods in Software Engineering*, Michael Felderer and Guilherme Horta Travassos (Eds.). Springer, Chapter 11.
- 1884 [65] Richard Torkar, Carlo A. Furia, Robert Feldt, Francisco Gomes de Oliveira Neto, Lucas Gren, Per Lenberg, and Neil A. Ernst. 2021. A method to assess and argue for practical significance in software engineering. *Transactions on Software Engineering* -, - (2021), 13 pages. <https://doi.org/10.1109/TSE.2020.3048991>
- 1885 [66] Rens van de Schoot, Sarah Depaoli, Ruth King, Bianca Kramer, Kaspar Märtens, Mahlet G. Tadesse, Marina Vannucci, Andrew Gelman, Duco Veen, Joukje Willemsen, and Christopher Yau. 2021. Bayesian statistics and modelling. *Nature Reviews Methods Primers* 1, 1 (14 Jan 2021), 1. <https://doi.org/10.1038/s43586-020-00001-2>
- 1886 [67] Aki Vehtari, Andrew Gelman, and Jonah Gabry. 2017. Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Statistics and Computing* 27, 5 (2017), 1413–1432. <https://doi.org/10.1007/s11222-016-9696-4>
- 1887 [68] Helen M. Walker. 1958. The Contributions of Karl Pearson. *J. Amer. Statist. Assoc.* 53, 281 (1958), 11–22.
- 1888 [69] Hao Wang and Dit-Yan Yeung. 2016. Towards Bayesian deep learning: A framework and some existing methods. *IEEE Transactions on Knowledge and Data Engineering* 28, 12 (2016), 3395–3408.
- 1889 [70] Ronald L. Wasserstein and Nicole A. Lazar. 2016. The ASA Statement on *p*-Values: Context, Process, and Purpose. *The American Statistician* 70, 2 (2016), 129–133. <https://www.amstat.org/asa/files/pdfs/P-ValueStatement.pdf>.
- 1890 [71] Sumio Watanabe. 2010. Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *Journal of Machine Learning Research* 11 (Dec. 2010), 3571–3594.
- 1891

- 1912 [72] Ian H. Witten, Eibe Frank, Mark A. Hall, and Christopher J. Pal. 2016. *Data Mining: Practical Machine Learning*
1913 *Tools and Techniques* (4 ed.). Morgan Kaufmann.
- 1914 [73] Zhou Xu, Li Li, Meng Yan, Jin Liu, Xiapu Luo, John Grundy, Yifeng Zhang, and Xiaohong Zhang. 2021. A
1915 comprehensive comparative study of clustering-based unsupervised defect prediction models. *J. Syst. Softw.*
1916 172 (2021), 110862. <https://doi.org/10.1016/j.jss.2020.110862>
- 1917 [74] Zhe Yu, Fahmid Morshed Fahid, Huy Tu, and Tim Menzies. 2021. Identifying Self-Admitted Technical
1918 Debts with Jitterbug: A Two-Step Approach. *IEEE Transactions on Software Engineering* (2021). <https://doi.org/10.1109/TSE.2020.3031401> Preprint: <https://arxiv.org/abs/2002.11049>.

1919

1920

1921

1922

1923

1924

1925

1926

1927

1928

1929

1930

1931

1932

1933

1934

1935

1936

1937

1938

1939

1940

1941

1942

1943

1944

1945

1946

1947

1948

1949

1950

1951

1952

1953

1954

1955

1956

1957

1958

1959

1960