








JAX à la Stan

Bayesian modeling directly in Python

ISSN 2824-7795

Brian Ward  Center for Computational Mathematics, Flatiron Institute
Matthijs Vákár  Information and Computing Sciences, Utrecht University
Mitzi Morris  Independent Contractor
Andrew Gelman  Department of Statistics, Columbia University
Bob Carpenter ¹ Center for Computational Mathematics, Flatiron Institute

Date published: 2026-03-29 Last modified: 2026-03-29

Abstract

We introduce a methodology for coding Bayesian statistical models in Python with JAX that follows the design pattern of the Stan probabilistic programming language. This allows a direct, line-by-line translation into JAX of all of the courses, textbooks, and case studies for Stan across the physical, biological, and social sciences, engineering, business, health, education, policy, economics, and sports. It also provides a transparent framework for further model development. Coupled with modern hardware (e.g., multi-core, graphics processing units, and tensor processing units), compiled JAX far exceeds the efficiency and scalability of Stan for computing the log densities and gradients needed by state-of-the-art inference algorithms. JAX's implementation of NumPy and SciPy, along with the packages TensorFlow (including TensorFlow Probability) and Distrax, provide a much wider range of special function support than Stan, including partial and stochastic differential equations and neural networks. The package ArviZ provides the same posterior analysis tools as Stan, Blackjax provides a wider range of inference algorithms, and TensorFlow Probability provides a wider range of variable transforms. Together, these tools provide an environment to code models in the style of Stan targeting modern hardware without leaving an integrated Python programming environment.

Keywords: JAX, Stan, probabilistic programming, differentiable programming

1 Contents

2	1 Motivation	2
3	1.1 Bayesian workflow	2
4	1.2 Why not just use Stan?	3
5	1.3 Why not just use PyMC or NumPyro?	4
6	1.4 Special function support in JAX	5
7	1.5 Constrained parameter support in JAX	7
8	1.6 Modularity in JAX	7

¹Corresponding author: bcarpenter@flatironinstitute.org

9	2 Bayesian inference	7
10	2.1 Bayesian models	7
11	2.2 Generative modeling	8
12	2.3 Unmodeled data	8
13	2.4 Bayesian inference with posterior expectations and simulations	8
14	2.5 Computation with Markov chain Monte Carlo methods	10
15	2.6 Unconstrained parameterizations	10
16	3 Coding models in Stan	11
17	3.1 A multivariate linear regression model	11
18	3.2 Coding a linear regression in Stan	12
19	3.3 Simulating data	12
20	4 Sampling with Stan	13
21	4.1 The CmdStanPy interface	13
22	4.2 Transpilation and compilation	14
23	4.3 Sampling	14
24	4.4 Posterior summaries with CmdStanPy and ArviZ	14
25	5 Stan to C++ transpilation	16
26	5.1 Stan’s block structure	16
27	5.2 Stan’s transpiled C++ model class	17
28	6 Coding models directly in JAX	19
29	6.1 Linear regression directly in JAX	19
30	6.2 The untransformed log posterior in JAX	19
31	6.3 Transforms and their inverses in JAX	19
32	6.4 The transformed log density in JAX	20
33	6.5 Random initialization in JAX	20
34	6.6 Sampling in Blackjax	20
35	6.7 Posterior analysis in JAX	21
36	6.8 Posterior predictive quantities in JAX	21
37	6.9 Key differences between the Stan and JAX implementations	21
38	7 Linear regression in JAX with dense_jax	22
39	7.1 Abstracting parameter transformations	22
40	7.2 Simplifying inference	23
41	References	24

1 Motivation

1.1 Bayesian workflow

Gelman et al. (2013) begin their foundational textbook by factoring Bayesian data analysis into three steps.

1. Design a joint probability distribution for observable data and unobservable parameters.
2. Perform inference to generate a posterior sample over parameters and unobserved data conditioned on observed data.
3. Evaluate the model fit and what it tells us about our quantities of interest.

50 The authors further suggest that if the evaluation in step (3) is not sufficient, then one should go back
51 to step (1) and try to come up with a better model. More recently, Gelman et al. (2020) outlined a
52 workflow for Bayesian analysis that puts more emphasis on fitting multiple models and transparently
53 reporting their exploration and comparison.

54 A probabilistic programming language (PPL) primarily provides support for statisticians who wish
55 to code a statistical model defined in step (1) so that it can be used for inference in step (2). The
56 PPL syntax allows for a straightforward encoding of the joint probability distribution. PPLs are
57 typically coupled with samplers that allow step (2) to be performed with Monte Carlo methods and
58 are coupled with posterior analysis and model comparison tools for task (3). In addition, once the log
59 unnormalized density is set up, PPLs can compute other operations such as mode finding, Laplace
60 approximation, and variational inference.

61 In this paper, we are going to focus on specifying models and performing posterior sampling as
62 is done by a PPL, but without actually using a PPL. We will also provide guidance on integration
63 with tools for inference, model evaluation/criticism, and model comparison. Our goal is to develop a
64 methodology for implementing efficient and scalable differentiable Bayesian models in Python in a
65 way that is both easy to code and easy to read.

66 For step (2), we recommend the Python package Blackjax (Cabezas et al. 2024). Blackjax includes
67 implementations of Stan’s primary inference methods, the no-U-turn sampler (NUTS) (Hoffman
68 and Gelman 2014), automatic differentiation variational inference (ADVI) (Kucukelbir et al. 2017),
69 and Pathfinder variational inference (Zhang et al. 2022). Blackjax is being actively maintained
70 and extended, and already includes several other useful algorithms, including microcanonical (aka
71 isokinetic) sampling (Robnik et al. 2025), sequential Monte Carlo (SMC) (Doucet et al. 2001), elliptical
72 slice sampling (Murray et al. 2010), generalized HMC (GHMC) (Horowitz 1991), and even random-
73 walk Metropolis (RWM) (Hastings 1970).

74 For step (3), we recommend the Python package Arviz (Kumar et al. 2019). Arviz provides state-of-
75 the-art convergence monitoring using split, ranked \hat{R} , estimation of bulk and tail effective sample
76 sizes, standard errors, posterior means, standard deviations, and quantiles, as well as approximate
77 leave-one-out (LOO) cross-validation (Vehtari et al. 2021, 2017).

78 1.2 Why not just use Stan?

79 Stan (Carpenter et al. 2017) is a domain-specific language for expressing differentiable probability
80 densities and posterior predictive quantities. Stan is a probabilistic programming language in the
81 sense that its variables can be interpreted as random variables. Stan has been used in almost every
82 area where statistics is applied, and as such, has accumulated an unmatched depth and breadth
83 of training materials around different classes of probabilistic models. There are textbooks, college
84 classes, and reproducible case studies in all of these areas. There’s a vibrant community with a high
85 volume [discussion forum](#). The language is still being expanded and so is its math library.

86 The Stan project introduced several state of the art gradient-based inference algorithms including
87 the no-U-turn sampler (NUTS) (Hoffman and Gelman 2014), automatic differentiation variational
88 inference (ADVI) (Kucukelbir et al. 2017), Pathfinder variational inference (Zhang et al. 2022), and
89 black-box nested Laplace approximations (Margossian 2023) as well as posterior analysis tools such
90 as split- and ranked- \hat{R} and corresponding bulk and tail effective sample size (Vehtari et al. 2021),
91 leave-one-out cross-validation (Vehtari et al. 2017), refined simulation-based calibration checks
92 (Cook et al. 2006), and prior predictive checks (Gabry et al. 2019). It is often used as the basis of
93 methodological developments such as Bayesian workflow (Gelman et al. 2020).

94 So why not just use Stan? The first reason is that Stan is largely CPU-bound. All of its analysis

95 tools and algorithms run on the CPU. Although there are ways to call individual functions in a Stan
96 program on the GPU (e.g., Cholesky decomposition) and ways to apply map-reduce across multiple
97 cores, this is not enough. Stan lacks a way to keep computation in-kernel on the GPU or organize
98 inference to enable single-instruction multiple-data (SIMD) parallelism. Although Stan has been
99 faster than JAX on CPU, it is not competitive on modern hardware (Sountsov et al. 2024; Maskell
100 2024).

101 The second obstacle to using Stan is the need to learn a second language. While Stan is not particularly
102 complicated, it does present several additional difficulties beyond unfamiliar syntax and semantics.

- 103 • Stan is indexed from 1, like much of mathematics and in particular, linear algebra, whereas
104 Python is indexed from 0, like most programming languages. Translating between 0-based
105 and 1-based indexing is tedious, error-prone, and obfuscates code.
- 106 • Stan is strongly typed and statically compiled, which has the benefit of being type-safe at run
107 time and leads to fast C++ computation. The downside is that this kind of static typing is
108 unfamiliar to most of Stan’s intended users. Stan’s typing can be an annoyance and performance
109 bottleneck even for experienced programmers due to its poor representational choice for
110 containers that mixes C++ standard vectors (Josuttis 2012) for arrays and Eigen matrices
111 (Guennebaud and Jacob 2010) for linear algebra (this may sound harsh, but it was our fault as
112 the original developers of Stan).
- 113 • Stan requires a scripting language like R, Python, and Julia, to combine modeling with data
114 preparation and posterior analysis, rather than providing a seamless single-language experience.
115 The interface between these languages and Stan is minimal—Stan is just being called as a black
116 box to return samples.
- 117 • There is relatively little tooling to aid with Stan development. Currently, it’s just autocomplete,
118 syntax highlighting, and debug-by-print. Python, in contrast, has several well supported
119 integrated development environments with tooling for debugging, generating notebooks,
120 integration with chatbots, integration with documentation, automatic refactoring tools, etc.
- 121 • While there is a great deal of tutorial and onboarding material for Stan, it has to split its
122 attention among several interfaces (two interfaces in R and Python and one in Julia). Taken
123 together, even Stan’s extensive documentation is dwarfed by the pedagogical material around
124 scientific, differentiable, and probabilistic computing in Python.

125 1.3 Why not just use PyMC or NumPyro?

126 The very first probabilistic programming language of which we are aware is BUGS (Bayesian inference
127 using Gibbs Sampling) (Lunn et al. 2009, 2012), which was released way back in 1991. In BUGS,
128 Bayesian models are specified with deterministic and stochastic nodes arranged in a directed acyclic
129 graph. Each node was either input as data, or defined as a (possibly stochastic) function of its
130 direct ancestors in the graph. This enabled a BUGS model to be used to infer any of the stochastic
131 variables in the model given values for the data nodes and all other stochastic nodes. This provides a
132 clean way to perform analyses such as prior predictive inference and posterior predictive inference
133 automatically through the graphical structure of the model. BUGS samples using generalized Gibbs
134 sampling, which does not scale well in dimension.

135 PyMC (Salvatier et al. 2016) and NumPyro (Phan et al. 2019) are Python packages that take a directed
136 graphical modeling approach to specifying Bayesian models and are capable of generating JAX code
137 as output. There are similar packages in other languages, but they do not generate JAX code. JAGS
138 (Plummer 2003) is a standalone language that reimplements and extends BUGS and is typically used
139 through R, NIMBLE (Valpine et al. 2017) is coded in R, and Turing.jl (Ge et al. 2018) is coded in Julia.

140 Like Stan, all of the modern PPLs efficiently scale in dimension by using gradient-based inference

141 methods. Like BUGS, they are able to exploit the graphical model structure directly to automate
142 a number of functions that are painful to code in Stan and will largely remain painful to code in
143 what we are proposing here, such as prior and posterior predictive checks and simulation-based
144 calibration, and at least in the case of PyMC, general patterns of missing data.

145 When models get more complicated in terms of novel parameter constraints, densities, conditional
146 structures, etc., both PyMC and NumPyro provide escape hatches to let you define transformations
147 and log densities directly in the same way as Stan. This feels dirty in the same way as using the
148 “ones-trick” in BUGS (Lunn et al. 2012). And while it is great for allowing general models to be
149 defined, it defeats all the benefits of having a clean generative graphical model in the first place. At
150 the point models start getting more complicated, we believe it is more straightforward to code the
151 models directly in JAX rather than working around the graphical modeling paradigm of PyMC or
152 NumPyro.

153 A second reason to prefer the approach we are presenting here is that it is much more direct. By that,
154 we mean that like Stan, the resulting code is implemented transparently in an imperative fashion
155 rather than indirectly through the structure of the directed acyclic graph.

156 1.4 Special function support in JAX

157 Stan has an extensive library of special mathematical and statistical functions, as well as restructuring
158 functions for arrays and matrices. Many of these are needed to differentiate cumulative distribution
159 functions and to define custom densities. Here’s a brief overview of the coverage available in JAX
160 compared to Stan. The bottom line is that support for special functions is deeper in JAX.

- 161 • *Matrix library*: This is JAX’s main focus and it far exceeds Stan’s collection of familiar ma-
162 trix functions and reshaping tools by punning NumPy ([jax.numpy](#)) and SciPy ([jax.scipy](#)).
163 There is even limited (and experimental) support for sparse matrices and solvers natively
164 ([jax.experimental.sparse](#)).
- 165 • *Special functions*: These are available all over the JAX ecosystem, including in JAX’s NumPy
166 and SciPy modules. The differentiable SciPy module ([jax.scipy.special](#)) is not complete
167 compared to SciPy. The deficit is more than made up by the special function library provided
168 by TensorFlow Probability (TFP) ([tfp.math](#)) and maintained by Google. For example, the
169 Lambert W function available in Stan has not been ported from SciPy but is available through
170 TFP. The bottom line is that JAX provides a *better* selection of well supported special functions.
- 171 • *Probability distributions*: Stan implements dozens of probability distributions, including
172 most (but not all) of the ones in common use for statistical models. While the basics
173 are available through JAX’s NumPy and SciPy modules (often redundantly), the go-to
174 library is TensorFlow Probability (TFP) (e.g., [tfp.distributions](#); for JAX specifically,
175 [tfp.substrates.jax.distributions](#)) for probability-related functions (e.g., probability
176 density functions, probability mass functions, and cumulative distribution functions). In some
177 cases, there are also quantile functions, which are poorly supported in Stan. The bottom line is
178 that the native Google-supported JAX ecosystem provides a *better* selection of well supported
179 probability functions and random number generators. There is even wider support beyond
180 native JAX and TensorFlow, including the probabilistic programming language NumPyro
181 (Phan et al. 2019) and Google DeepMind’s library Distrax (DeepMind et al. 2020).
- 182 • *Complex-valued functions*: JAX and Stan both have core library support for complex-valued
183 functions built-in, including fast Fourier transforms and complex matrix operations.
- 184 • *Neural networks*: JAX is integrated tightly with a range of neural network constructions through
185 the Flax package (Heek et al. 2024), which is maintained by a team at Google DeepMind, but

186 is not an official Google product like TensorFlow or JAX. It supports everything from simple
187 multilayer perceptrons and convolutional neural networks to autoencoders and multi-head
188 attention.

- 189 • *Simulation-based inference (SBI)*: There is as of yet, no mature and commonly used SBI package
190 in JAX, nor is there any support in Stan.
- 191 • *Implicit Solvers*: Applied statistics often requires equations to be solved and differentiated and
192 Stan provides a fairly extensive library.
 - 193 – *Ordinary differential equation (ODE) solvers*: There is no built-in support in JAX for ODE
194 solvers, but the Difffrax package (Kidger 2021) is widely used and provides the same
195 kind of adjoint and analytic methods as Stan that provide sensitivity analysis without
196 automatically differentiating through the algorithm.
 - 197 – *Root finders*: There is no built-in support in JAX for root finders, but the JAXOpt package
198 (Blondel et al. 2021) is maintained by Google and provides a range of solvers including
199 the Newton method used by Stan.
 - 200 – *1D Integration*: These functions are useful for defining cumulative distribution functions
201 for novel densities. JAX does not ship an adaptive 1D quadrature routine in core, but ex-
202 ternal libraries such as [quadax](#) provide `jit/vmap`-able, differentiable adaptive quadrature
203 (e.g., Gauss–Kronrod-style `quadgk`).
 - 204 – *Hidden Markov models*: TFP (through `tfp.distributions.HiddenMarkovModel`)
205 supports likelihood evaluation (forward algorithm) and computation of latent-state
206 marginals (a form of the forward–backward algorithm, e.g. `posterior_marginals`).
207 Stan provides analogous built-ins (`hmm_marginal` for the marginal likelihood and
208 `hmm_hidden_state_prob` for marginal latent state probabilities).
 - 209 – *Kalman filters*: Stan has direct support for Gaussian dynamic linear model likelihoods com-
210 puted via the Kalman filter (e.g., `gaussian_dlm_obs_lpdf` and related functions). There is
211 also experimental support in TFP (in `tfp.experimental.parallel_filter.kalman_filter`).
 - 212 – *Partial differential equation (PDE) solvers*: Stan has no support for partial differential
213 equations. The JAX-CFD package (Kochkov et al. 2021) from Google, mostly focused
214 on fluid dynamics, includes a number of general-purpose PDE solvers, including finite
215 volume/difference methods, pseudospectral methods, and machine-learning methods. It
216 supports Navier-Stokes, advection-diffusion, period and wall boundary conditions, and
217 runs on GPU and TPU with gradients.
 - 218 – *Stochastic differential equation (SDE) solvers*: Stan has no support for stochastic differ-
219 ential equations. The Difffrax package (Kidger 2021), in addition to providing stiff ODE
220 solvers, also supports traditional SDE and SPDE solvers (Euler-Maruyama, Milstein,
221 Stratonovich/Itô) through spatial discretization.

222 With the caveat that functionality is spread across multiple libraries (core JAX, TFP, and smaller
223 add-ons such as `quadax` for adaptive quadrature), the special function, probability distribution, and
224 transform coverage available in the JAX ecosystem is broadly comparable to Stan’s, and in some
225 areas deeper. When it comes to even more complicated functions like PDE and SDE solvers and
226 neural networks, Stan isn’t even in the game. The JAX-based PPL NumPyro can also make use of
227 these JAX packages in many cases, especially if users are willing to forego its underlying graphical
228 model abstraction.

229 1.5 Constrained parameter support in JAX

230 Stan provides built-in transforms that map an unconstrained vector in \mathbb{R}^d into a desired constrained
231 space (e.g., a $(d + 1)$ -simplex), along with the corresponding change-of-variables adjustments to the
232 density. These are custom implementations with analytic Jacobian-adjoint product gradients and
233 vectorized application to containers.

234 The Oryx transform library in TensorFlow is built on top of the TFP bijector library ([tfp.bijectors](#))
235 (Dillon et al. 2017), which can also be used directly. TFP bijectors provide all of the transforms
236 provided by Stan and many more including softplus, various cdfs and hyperbolic tangent as replace-
237 ments for inverse-logit, more multivariate transforms such as cumulative sums and Householder
238 factorizations, as well as trained transforms like RealNVP normalizing flows (Dinh et al. 2016). Oryx
239 additionally allows transforms to be written down directly in such a way that Oryx can automatically
240 calculate inverse transforms and Jacobian determinants of inverse transforms.

241 1.6 Modularity in JAX

242 SlicStan (Gorinova et al. 2019) reconceived Stan without blocks—the sorting into data, parameters,
243 and generated quantities was carried out by data flow analysis. The primary motivation was to make
244 it possible to modularly express concepts like a hierarchical prior. With Stan itself, this is impossible
245 unless the modularity is in the form of a simple function. With SlicStan, the parameters, priors,
246 etc., could all be constructed modularly and reused. By allowing models to be expressed directly in
247 Python code, NumPyro and PyMC already support modular code reuse. Although it is rare to see
248 this feature used in example code, it is widely used in production.

249 By coding models directly in Python, we gain the same benefits of NumPyro and PyMC. We can
250 write general programs returning arbitrary components of a probabilistic program and combine
251 them at will. We will provide examples below.

252 2 Bayesian inference

253 In this section, we will lay out the precise definition of the Bayesian inference problem we are trying
254 to compute, show how it can be computed asymptotically exactly using Monte Carlo methods, and
255 apply some simple calculus to transform parameterizations to be unconstrained (i.e., have support
256 over all of \mathbb{R}^D).

257 2.1 Bayesian models

258 A Bayesian model is a joint distribution over parameters θ and data y .

259 In the simplest version of this formulation, Bayesian inference uses the posterior distribution of
260 the parameters given the data: $p(\theta | y) \propto p(\theta, y)$. Often our density is factored into a prior times
261 a likelihood, $p(\theta | y) = p(\theta) \cdot p(y | \theta)$, but nothing in Stan or what we are proposing for JAX
262 presupposes a clean factorization. In fact, our methods can be used to sample from any unnormalized
263 density with finite integral—it doesn't have to arise as a Bayesian posterior.

264 More generally, if “data” are defined as observable quantities, there can be observed data y^{obs} and
265 unobserved data y^{mis} (short for “missing data” but this can also include latent (unobserved but
266 potentially observable) data, future data, unrecorded past data, etc.), and Bayesian inference uses the
267 posterior distribution of all unobserved quantities conditional on observed data, thus $p(\theta, y^{\text{mis}} | y^{\text{obs}})$.
268 There is no strict rule for what aspects of a model count as “parameters” and what count as “data.”
269 But Bayesian inference doesn't really care if an unknown quantity is put in the “ θ ” bin or the “ y^{mis} ”
270 bin; all that matters is what is being conditioned on.

271 In that case, why not just label everything unknown in a problem as “parameters” and everything
272 observed as “data”? You can do this; indeed, that’s how variables are coded in Stan. The trouble
273 with using this as a general dividing line is that, the knowledge of what is observed and what is not
274 observed can be logically independent of the model. For a simple example, start with data y_1, \dots, y_N
275 from a first-order autoregressive model with parameters a, b, σ . Now suppose that you want to make
276 predictions for future data, y_{N+1}, \dots, y_{N+10} . It would be awkward to label these new data points as
277 parameters, even though they are unknown.

278 2.2 Generative modeling

279 In Bayesian terminology, the joint distribution, $p(\theta, y) = p(\theta)p(y | \theta)$, is a *generative model*, because a
280 random θ can be generated from the prior and a random y can be generated from the data distribution.
281 For complicated models the generation process can be further divided; for example, in a hierarchical
282 model with local parameters α and hyperparameters ϕ , so that $\theta = (\phi, \alpha)$, the joint distribution can
283 be factored as $p(\phi)p(\alpha | \phi)p(\theta | \phi, \alpha)$ which corresponds to the generated process in which ϕ, α , and
284 y are simulated in order. As models become even more complicated, it can be helpful to express the
285 factorization as a directed acyclic graph, and the joint distribution is called a *graphical model*.

286 2.3 Unmodeled data

287 So far we have discussed parameters θ , observed data y^{obs} , and latent or unobserve data y^{mis} .
288 Bayesian models typically also include *unmodeled data*, for example the sample size and predictors in
289 a regression model. If we write all the unmodeled data as x , then the posterior is $p(\theta, y^{\text{mis}} | y^{\text{obs}}, x)$.
290 This is not a full generative model because there is no distribution for x .

291 Different probabilistic programming languages handle unmodeled data in different ways.

292 A Stan program computes the target function—the unnormalized log posterior density function, and
293 there is no logical distinction between y^{obs} and x : both are input to the program as data. Indeed,
294 there is no logical distinction between θ and y^{mis} , although for computational efficiency it is best to
295 simulate y^{mis} in the generated quantities block if possible.

296 In contrast, a PyMC program specifies a graphical model for $(\theta, y | x)$, and it can be run to simulate
297 any subset of (θ, y) , conditional on x and whatever information is passed in as data: for example,
298 specify nothing but x and it will simulate from the joint distribution; specify θ and x and it will
299 simulate y from the data model; specify y and x and it will simulate θ from the posterior.

300 In any probabilistic programming language, you need to specify the unmodeled data. But it is also
301 possible, and often recommended, to expand to include a model for x .

302 2.4 Bayesian inference with posterior expectations and simulations

303 Bayesian inference involves visualizing and summarizing quantities of interest in the posterior and
304 posterior predictive distributions. For summaries, we are typically interested in posterior averages or
305 posterior quantiles of parameters, predictions, or some functions of parameters and data. Quantiles
306 are typically used to give us medians and uncertainty intervals, while averages give us everything
307 else. By averaging estimates of quantities of interest derived from parameters weighted by posterior
308 density, we can account for estimation uncertainty into account in our inferences.

309 Averages over the posterior density are most naturally expressed using expectation notation. Given a
310 random variable $\Theta \in \mathbb{R}^D$ representing knowledge about the parameters, the expectation of a function
311 $f : \mathbb{R}^D \rightarrow \mathbb{R}$ is given by

$$E[f(\Theta) | y] = \int_{\mathbb{R}^D} f(\theta) \cdot p(\theta | y) d\theta.$$

312 The rest of this section provides examples of common functions whose expectations are evaluated
 313 for Bayesian inference.

314 The parameter estimate that minimizes expected square error, assuming the model is correct, is the
 315 posterior mean,

$$\hat{\theta} = \mathbb{E}[\Theta | y].$$

316 This is just taking the function f to be the identity. The posterior covariance is

$$\text{var}[\Theta | y] = \mathbb{E}[(\Theta - \hat{\theta}) \cdot (\Theta - \hat{\theta})^\top].$$

317 Here the function is $f(\theta) = (\theta - \hat{\theta}) \cdot (\theta - \hat{\theta})^\top$.

318 If we have an event $A \subseteq \mathbb{R}^D$, then its probability is given by the expectation of its indicator function,

$$\Pr[\theta \in A | y] = \mathbb{E}[\mathbb{1}_A(\theta \in A) | y].$$

319 If we want to evaluate posterior predictive densities $p(\tilde{y} | y)$ for new data \tilde{y} , this can be expressed as
 320 an expectation of the likelihood $p(\tilde{y} | \theta)$,

$$p(\tilde{y} | y) = \mathbb{E}[p(\tilde{y} | \theta) | y].$$

321 We can reduce the high-dimensional integration problem to one of sampling from the posterior,

$$\theta^{(m)} \sim p(\theta | y).$$

322 Given independent, identically distributed (i.i.d.) Monte Carlo draws, the Monte Carlo estimator

$$\hat{I}_M = \frac{1}{M} \sum_{m=1}^M f(\theta^{(m)})$$

323 converges almost surely to the true posterior expectation $\mathbb{E}[f(\Theta) | y]$ by the strong law of large
 324 numbers. Moreover, if $\text{var}[f(\Theta) | y] < \infty$, then a central limit theorem implies the asymptotic normal
 325 approximation

$$\sqrt{M}(\hat{I}_M - \mathbb{E}[f(\Theta) | y]) \xrightarrow{d} \mathcal{N}(0, \text{var}[f(\Theta) | y]),$$

326 which justifies the familiar $\mathcal{O}(1/\sqrt{M})$ Monte Carlo error rate and standard-error estimates for averages.
 327 Finite-sample concentration bounds are also hold under additional regularity assumptions (e.g.,
 328 boundedness or sub-exponential tails of $f(\Theta)$), giving us explicit deviation probabilities for \hat{I}_M
 329 around $\mathbb{E}[f(\Theta) | y]$ (Gelman et al. 2013; Paulin 2015).

330 In practice, we often summarize Bayesian inference using posterior simulations, without deciding
 331 ahead of time what expectations or other posterior summaries will be needed. So it is convenient
 332 that, with a large number of independent draws from the posterior distribution of all parameters and
 333 unobserved data of potential interest, we can use these to estimate arbitrary expectations, as well as
 334 obtain an estimate of their Monte Carlo uncertainty.

335 In addition, posterior simulations can be used for multiple imputation of missing data (Rubin 1996)
 336 and to understand the fitted model. For example, in a logistic regression model, $\Pr(y = 1|x) =$
 337 $\text{logit}^{-1}(a + bx)$, one can make a scatterplot of posterior simulations of (a, b) , along with a graph of
 338 data $(x, y)_n, n = 1, \dots, N$ along with curves, $\text{logit}^{-1}(a + bx)$ corresponding to those draws of (a, b) .

2.5 Computation with Markov chain Monte Carlo methods

In practice, exact independent draws from $p(\theta | y)$ are rarely available. Fortunately, many weaker sampling schemes still yield consistent estimators of posterior expectations. These include rejection sampling and importance sampling (which produce weighted estimators), sequential Monte Carlo (which propagates and reweights a particle population), and Markov chain Monte Carlo (which constructs a dependent sequence whose stationary distribution is the posterior) (Gelman et al. 2013; Doucet et al. 2001; Roberts and Rosenthal 2004).

The idea of Markov chain Monte Carlo (MCMC) is that *correlated* draws from an ergodic Markov chain with invariant distribution $p(\theta | y)$ can still be used in the same sample-average estimator \hat{I}_M . Under standard conditions (e.g., Harris ergodicity) that are easily satisfied in practical applications, the Markov chain ergodic theorem guarantees $\hat{I}_M \rightarrow \mathbb{E}[f(\Theta) | y]$ even though the draws are dependent (Roberts and Rosenthal 2004).

When the chain is particularly well behaved (e.g., geometrically ergodic with suitable moment conditions) a Markov chain central limit theorem holds:

$$\sqrt{M}(\hat{I}_M - \mathbb{E}[f(\Theta) | y]) \xrightarrow{d} \mathcal{N}(0, \sigma_f^2),$$

with asymptotic variance

$$\sigma_f^2 = \text{var}_{\pi}(f(\Theta))\tau_{\text{int}},$$

where $\pi(\cdot) = p(\cdot | y)$, $\tau_{\text{int}} = 1 + 2 \sum_{k \geq 1} \text{corr}_{\pi}(f(\Theta_0), f(\Theta_k))$ and $\{\Theta_k\}$ denotes the stationary chain. This formula for the variance motivates the usual definition of the effective sample size $M_{\text{eff}} = M/\tau_{\text{int}}$ and the heuristic $\mathcal{O}(1/\sqrt{M_{\text{eff}}})$ error rate (Roberts and Rosenthal 2004; Vehtari et al. 2021). Under similarly strong assumptions (e.g., uniform or geometric ergodicity, or a spectral-gap condition), one can also obtain non-asymptotic deviation bounds and Bernstein/Hoeffding-type concentration inequalities for Markov chain averages (Paulin 2015).

Unfortunately, these strong conditions needed for MCMC central limit theorems and concentrations bounds are typically difficult to verify for the complex MCMC methods and posteriors encountered in modern Bayesian inference over continuous parameter spaces (Roberts and Rosenthal 2004). Even so, MCMC, especially modern gradient-based variants, has shown remarkable empirical success in applied Bayesian inference, with practical reliability assessed via convergence diagnostics, effective sample sizes, and predictive checks rather than direct verification of theoretical assumptions (Gelman et al. 2020; Vehtari et al. 2021).

2.6 Unconstrained parameterizations

It is much easier to define a sampling algorithm for situations where the posterior has support over all of \mathbb{R}^D , i.e., for all $\theta \in \mathbb{R}^D$, $p(\theta | y) > 0$. Stan and other PPLs transform any constrained parameters to be unconstrained. Then in practice, the unconstrained parameters are inverse transformed to satisfy their constraints. This requires an adjustment for the change of variables, which happens implicitly in PPLs. Stan requires all parameter constraints to be declared; the graphical modeling sublanguage of PyMC and NumPyro infer these constraints from the distributions in which variables participate (e.g., if a variable is given a Wishart distribution, it must be a symmetric and positive definite matrix).

Mathematically, given a constrained random variable $\Theta \in C \subseteq \mathbb{R}^D$, with a density $p_{\Theta}(\theta)$, and a smooth bijection $f : C \rightarrow \mathbb{R}^N$, we can derive the density of $\Phi = f(\Theta)$ as

$$p_{\Phi}(\phi) = p_{\Theta}(f^{-1}(\phi)) \cdot |\nabla f^{-1}(\phi)|,$$

378 where $|\cdot|$ denotes the absolute determinant operator and f^{-1} is the inverse of f . In the univariate
379 case, $\nabla f^{-1}(y)$ reduces to the derivative of the inverse transform at y (i.e., $\nabla f^{-1} = (f^{-1})'$).

380 If there is a sequence of variables being transformed one at a time, the overall Jacobian will be
381 block diagonal, with an absolute Jacobian determinant equal to the product of the absolute Jacobian
382 determinants of the blocks. Unconstrained parameters are transformed by the identity, which has a
383 unit Jacobian determinant. This makes it particularly simple to work on the unconstrained scale—we
384 just map unconstrained parameters back to the constrained space using the inverse transforms
385 and add the log absolute determinants of their Jacobians. For maximum likelihood estimation, the
386 Jacobians can be dropped from the target density with a flag.

387 Stan supplies constraints for variables that are lower bounded (for scales), upper bounded (for log
388 probabilities), range bounded (for probabilities), affine transforms (for non-centered parameteriza-
389 tions), ordered vectors (for cutpoints in ordinal regressions or identifying mixtures), unit vectors (for
390 points on a hypersphere), simplexes (for categorical probability distributions), sum-to-zero vectors
391 (for identifying varying effects), positive-definite symmetric matrices and their Cholesky factors (for
392 covariance or precision matrices), and for unit-diagonal positive-definite matrices and their Cholesky
393 factors (for correlation matrices). The Oryx package, which is native to JAX and built on top of the
394 TensorFlow Probability bijectors package (Dillon et al. 2017), provides an even wider range of useful
395 transforms than Stan (e.g., softplus, alternative sigmoid cdfs, tanh, autoregressions, and flows, and
396 many many more).

397 3 Coding models in Stan

398 To ground the discussion, let's consider the concrete example of coding a linear regression and using
399 it to predict new observations.

400 3.1 A multivariate linear regression model

401 We will assume a very simple multivariate regression formulation with an intercept, $P \in \mathbb{N}$ covariates,
402 and $N \in \mathbb{N}$ observations. Our data is made up of observations $y_n \in \mathbb{R}$ paired with covariates
403 $x_n \in \mathbb{R}^{N \times P}$. We will assume the usual parameters consisting of a slope $\alpha \in \mathbb{R}$, regression coefficients
404 $\beta \in \mathbb{R}^P$, and an error scale $\sigma \in (0, \infty)$. We will assume the data and covariates are approximately on
405 unit scale so that we can consider the following prior as weakly informative:

$$\alpha \sim \text{normal}(0, 5) \quad \beta_p \sim \text{normal}(0, 2.5) \quad \sigma \sim \text{exponential}(0.5).$$

406 We then add the conventional data generating process with independent normal errors,

$$y_n \sim \text{normal}(\alpha + x_n \cdot \beta, \sigma).$$

407 The joint density defining our Bayesian model (with data x taken as an unmodeled constant) is thus

$$p(y, \alpha, \beta, \sigma | x) = \text{exponential}(\sigma | 0.5) \cdot \text{normal}(\alpha | 0, 5) \\ \cdot \left(\prod_{p=1}^P \text{normal}(\beta_p | 0, 2.5) \right) \cdot \left(\prod_{n=1}^N \text{normal}(y_n | \alpha + x_n \cdot \beta, \sigma) \right).$$

408 Bayes's rule allows us to use the joint density as an unnormalized posterior,

$$p(\alpha, \beta, \sigma | y, x) \propto p(y, \alpha, \beta, \sigma | x).$$

409 The only constrained parameter is $\sigma > 0$. We transform positive-constrained parameters using the
410 log transform, i.e., $\sigma^{\text{unc}} = \log \sigma$. The inverse transform is the exponential. Applying the change-
411 of-variables formula and using the fact that $|\nabla \exp(u)| = |\exp'(u)| = \exp(u)$, the corresponding

412 unconstrained density is

$$p^{\text{unc}}(\alpha, \beta, \sigma^{\text{unc}} | x, y) = p(\alpha, \beta, \exp(\sigma^{\text{unc}}) | x, y) \cdot \exp(\sigma^{\text{unc}}).$$

413 On the log scale, where we operate to prevent underflow and maintain precision, we have

$$\log p(\alpha, \beta, \sigma^{\text{unc}} | x, y) = \log p(\alpha, \beta, \exp(\sigma^{\text{unc}}) | x, y) + \sigma^{\text{unc}}.$$

414 3.2 Coding a linear regression in Stan

415 Here's an example Stan program defining a linear regression, which we have placed into the file
416 `linear-regression.stan`.

```
data {  
  int<lower=0> N, N_new, P;  
  matrix[N, P] x;  
  vector[N] y;  
  matrix[N_new, P] x_new;  
}  
parameters {  
  real alpha;  
  vector[P] beta;  
  real<lower=0> sigma;  
}  
model {  
  alpha ~ normal(0, 5);  
  beta ~ normal(0, 2.5);  
  sigma ~ exponential(0.5);  
  y ~ normal(alpha + x * beta, sigma);  
}  
generated quantities {  
  array[N_new] real y_new = normal_rng(alpha + x_new * beta, sigma);  
}
```

417 Stan is very flexible. For example, if y were binary data and we wanted to fit a logistic regression,
418 y would be declared as `array[N] int<lower=0, upper=1>` and `normal()` would be replaced with
419 `bernoulli_logit()` (and similarly for the generated quantities).

420 To ease the transition to JAX, note that the distribution statements in the model block are just
421 syntactic sugar for incrementing the log target density (Carpenter et al. 2017). The model block
422 could have been coded as follows to generate the same C++ code.

```
target += normal_lupdf(alpha | 0, 5);  
target += normal_lupdf(beta | 0, 2.5);  
target += lognormal_lupdf(sigma | 0, 1);  
target += normal_lupdf(alpha + x * beta, sigma);
```

423 Here, the `_lupdf` indicates a log (l), unnormalized (u), probability density function (pdf). For
424 probability mass functions, replace `pdf` with `pmf`; to preserve normalizing constants, drop the `u`.

425 3.3 Simulating data

426 The sizes and covariates x cannot be simulated from the model as they are not modeled. We'll have to
427 invent a distribution for those in order to simulate data, and we'll take $p = 2$ and generate standard

428 normal variates $x_{0:N,0} \sim \text{normal}(0, 1)$, and take $x_{0:N,1} \sim \text{chiSq}(1)$ by squaring a standard normal
429 variate.

430 The simulated parameter values are as follows.

431 `alpha: -9.15; beta[0]: -4.82; beta[1]: 1.15; sigma: 0.60`

432 In Figure 1, we provide a heatmap of the expected values of $y \sim \text{normal}(\alpha + x \cdot \text{beta}, \sigma)$, given values
433 of x . The conditional expectation is $E[y | x] = \alpha + x \cdot \beta$. Here, x is being read as a row vector like in
434 the Stan program.

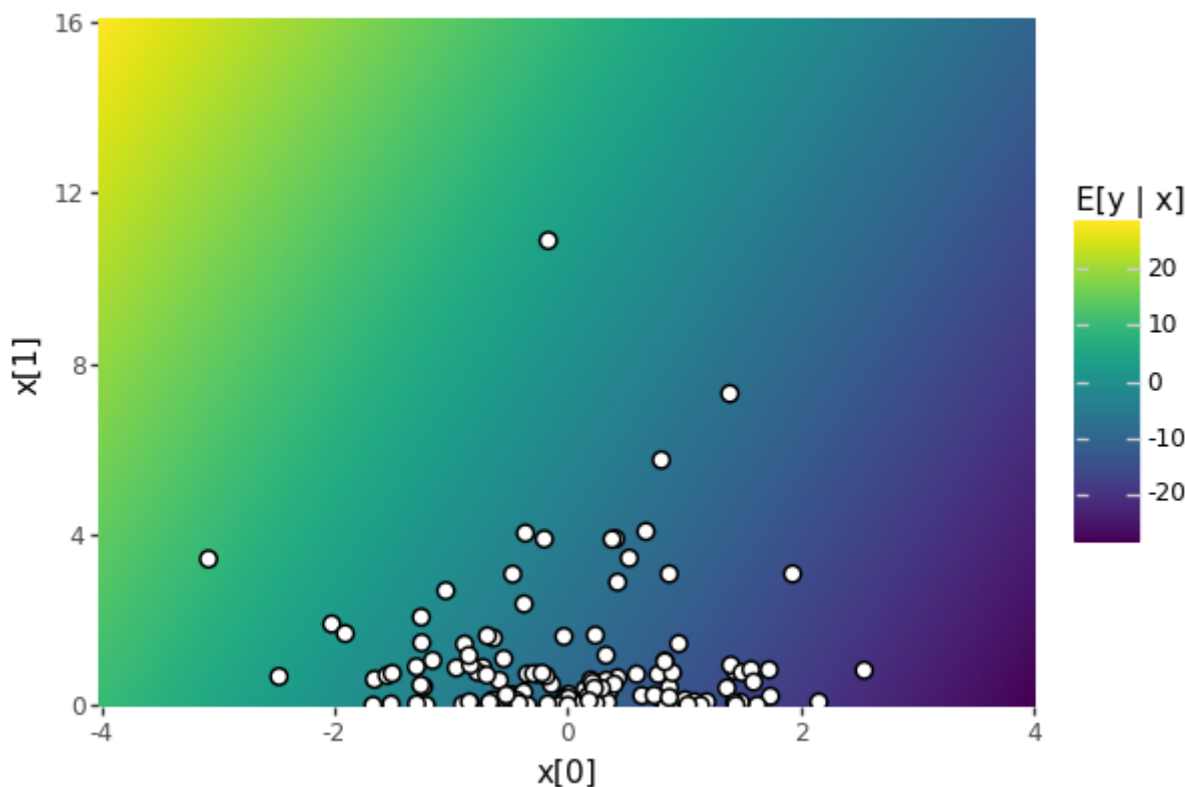


Figure 1: Heatmap of the expected value surface $E[y | x] = \alpha + x @ \text{beta}$ over 4 marginal standard deviations in x , where $x[0] \sim \text{normal}(0, 1)$ and $x[1] \sim \text{chi_sq}(1)$. The observed covariates overlaid as white disks.

435 4 Sampling with Stan

436 Now that we have defined a Stan program and simulated data, we can perform inference from the
437 simulated data (pretending that we did not already know the parameters). We will use Stan's defaults.

438 4.1 The CmdStanPy interface

439 We will use the Python interface CmdStanPy (Stan Development Team 2024), which makes an
440 external system call to the reference command-line interface, CmdStan. Before that, we import the
441 cmdstanpy package and make sure we have a compiled CmdStan installation and turn the logging
442 level down to ERROR so that our output isn't cluttered with progress messages.

443 `CmdStan install directory: /home/runner/.cmdstan`

444 `CmdStan version 2.38.0 already installed`

445 Test model compilation

446 True

447 4.2 Transpilation and compilation

448 Once CmdStan is installed, transpiling the model to C++ and compiling the C++ is a one-liner.

449 4.3 Sampling

450 Given the model `m`, sampling with Stan's default parameters is a one liner. We call the sampler with
451 the data variable simulated above, and we turn off progress bars to render in a paper.

452 The variable data was defined above in the simulation. The call to `sample()` uses CmdStanPy's
453 default configuration parameters. By default, Stan uses the multinomial, biased-progressive no-
454 U-turn sampler (Betancourt 2017), which invokes three stages of warmup (finding the bulk of the
455 probability mass, estimating the inverse mass matrix, and estimating the step size). By default, Stan
456 runs 1000 warmup and 1000 sampling iterations in 4 independent chains—this is often more than are
457 necessary for model development and many applications.

458 4.4 Posterior summaries with CmdStanPy and ArviZ

459 Now that we have the fit, we can print a summary of the posterior from CmdStan, reducing its default
460 number of digits printed.

	Mean	MCSE	StdDev	MAD	5%	50%	95%	ESS_bulk	\
461 lp__	10.600	0.034200	1.5100	1.3200	7.830	11.000	12.400	2070.0	
462 alpha	-9.170	0.000943	0.0587	0.0582	-9.270	-9.170	-9.070	3880.0	
463 beta[1]	-4.810	0.000828	0.0490	0.0479	-4.890	-4.810	-4.730	3580.0	
464 beta[2]	1.150	0.000557	0.0326	0.0326	1.090	1.150	1.200	3480.0	
465 sigma	0.542	0.000588	0.0354	0.0350	0.488	0.541	0.603	3700.0	
466 y_new[1]	-10.600	0.008640	0.5480	0.5400	-11.400	-10.600	-9.640	4030.0	
467 y_new[2]	-16.300	0.008700	0.5560	0.5520	-17.200	-16.300	-15.400	4100.0	
468 y_new[3]	-17.800	0.008700	0.5490	0.5510	-18.700	-17.800	-16.900	4020.0	
469 y_new[4]	-8.530	0.008720	0.5590	0.5440	-9.450	-8.530	-7.620	4110.0	
470									
471									
472	ESS_tail	ESS_bulk/s	R_hat						
473 lp__	2510.0	12000.0	1.0						
474 alpha	2630.0	22600.0	1.0						
475 beta[1]	2570.0	20800.0	1.0						
476 beta[2]	2750.0	20300.0	1.0						
477 sigma	2930.0	21500.0	1.0						
478 y_new[1]	3930.0	23500.0	1.0						
479 y_new[2]	3750.0	23800.0	1.0						
480 y_new[3]	3710.0	23400.0	1.0						
481 y_new[4]	4080.0	23900.0	1.0						

482 The `summary_csp` variable provides programmatic access to the data this printed by default along
483 with much more debugging and configuration information. The variable `lp__` reported in the first
484 row is the unnormalized log posterior density, the variables `alpha`, `beta`, and `sigma` report model
485 parameters, and `y_new` reports the models posterior predictions. The remaining columns are defined
486 per variable and include the posterior mean, Monte Carlo standard error, posterior standard deviation,
487 mean absolute deviation (like standard deviation for medians), three quantiles (5%, 50%, and 95%), as

488 well as some convergence diagnostics. The rank-normalized \hat{R} statistic is reported in the last column,
 489 and the effective sample size in the bulk and tail of the density, as well as bulk effective sample size
 490 per second (where we see this all ran quickly); definitions of all of these are provided by Vehtari et al.
 491 (2021).

492 The summary shows that the model fits very well. While it doesn't report timing, the effective sample
 493 size per second statistic lets you infer that this all ran in less than 0.1s. The \hat{R} statistics are all below
 494 1.005 and the effective sample sizes are all above 2000 with 4000 total sampling draws across four
 495 chains.

496 A very similar summary is available from the Python package ArviZ, which shares many developers
 497 with Stan. Using the method `draws_xr()`, `CmdStanPy` produces an `Xarray` (defined in package
 498 `xarray`) that can be directly consumed by ArviZ's `summary()` function.

	mean	sd	hdi_3%	hdi_97%	mcse_mean	mcse_sd	ess_bulk	\
499 alpha	-9.171	0.059	-9.280	-9.056	0.001	0.001	3879.0	
500 beta[0]	-4.808	0.049	-4.902	-4.717	0.001	0.001	3577.0	
501 beta[1]	1.146	0.033	1.081	1.203	0.001	0.001	3484.0	
502 sigma	0.542	0.035	0.474	0.606	0.001	0.001	3702.0	
503 y_new[0]	-10.555	0.548	-11.628	-9.554	0.009	0.006	4035.0	
504 y_new[1]	-16.286	0.556	-17.287	-15.218	0.009	0.006	4101.0	
505 y_new[2]	-17.792	0.549	-18.794	-16.766	0.009	0.006	4021.0	
506 y_new[3]	-8.530	0.559	-9.536	-7.417	0.009	0.006	4113.0	
507								
508								
	ess_tail	r_hat						
509 alpha	2625.0	1.0						
510 beta[0]	2568.0	1.0						
511 beta[1]	2748.0	1.0						
512 sigma	2933.0	1.0						
513 y_new[0]	3933.0	1.0						
514 y_new[1]	3748.0	1.0						
515 y_new[2]	3705.0	1.0						
516 y_new[3]	4082.0	1.0						
517								

518 The main difference is that ArviZ defaults to reporting 94% highest density intervals rather than
 519 traditional central intervals based on quantiles. ArviZ can be configured to include Stan's quantile
 520 output, if desired and the highest-density interval changed to 90% coverage to match Stan's output
 521 using NumPy's quantile functions. These packages typically report 90% rather than 95% intervals
 522 because the standard error of quantile estimates increases dramatically in the tails (Bahadur 1966;
 523 Van der Vaart 1998).

	mean	sd	hdi_5%	hdi_95%	mcse_mean	mcse_sd	ess_bulk	\
524 alpha	-9.171	0.059	-9.268	-9.073	0.001	0.001	3879.0	
525 beta[0]	-4.808	0.049	-4.883	-4.722	0.001	0.001	3577.0	
526 beta[1]	1.146	0.033	1.093	1.200	0.001	0.001	3484.0	
527 sigma	0.542	0.035	0.485	0.598	0.001	0.001	3702.0	
528 y_new[0]	-10.555	0.548	-11.428	-9.638	0.009	0.006	4035.0	
529 y_new[1]	-16.286	0.556	-17.199	-15.401	0.009	0.006	4101.0	
530 y_new[2]	-17.792	0.549	-18.680	-16.881	0.009	0.006	4021.0	
531 y_new[3]	-8.530	0.559	-9.396	-7.574	0.009	0.006	4113.0	
532								
533								
	ess_tail	r_hat	q5	q50	q95			
534								

535	alpha	2625.0	1.0	-9.269	-9.171	-9.074
536	beta[0]	2568.0	1.0	-4.889	-4.807	-4.727
537	beta[1]	2748.0	1.0	1.092	1.146	1.199
538	sigma	2933.0	1.0	0.488	0.541	0.603
539	y_new[0]	3933.0	1.0	-11.436	-10.553	-9.643
540	y_new[1]	3748.0	1.0	-17.183	-16.291	-15.382
541	y_new[2]	3705.0	1.0	-18.706	-17.793	-16.898
542	y_new[3]	4082.0	1.0	-9.449	-8.529	-7.617

543 5 Stan to C++ transpilation

544 A Stan program translates almost line-for-line to a C++ class that implements all the interfaces
 545 around a model required for log density evaluation, posterior predictive quantity generation, and
 546 variable transforms (Stan Development Team 2025).

547 5.1 Stan's block structure

- 548 • *Functions block*: Each function in the functions block is translated to a C++ function that is
 549 templated flexibly enough to allow automatic differentiation. The main limitation to Stan's
 550 functions is that they cannot cross blocks and they cannot introduce block-level variables such
 551 as data or parameters; a secondary limitation is that arguments must be real or integer based.
 552 Thus it is not possible to define functions like the built-in ODE solvers that take functions as
 553 arguments.
- 554 • *Data and transformed data blocks*: Each data declaration in the data block is translated to
 555 an instance variable of the generated class. The data block can only contain declarations,
 556 and thus behaves like the signature for the data ingestion function. The transformed data
 557 block is executed after the data is read in as the model object is being constructed. After the
 558 model object is constructed with the data, it remains immutable to allow safe multi-threaded
 559 application.
- 560 • *Parameter and transformed parameter blocks*: The log density function is defined for uncon-
 561 strained inputs gathered into a vector. The first part of the function unpacks the entries of that
 562 vector, applies the constraining transform to define constrained variables locally, then adds
 563 the log Jacobian determinant to the Jacobian accumulator to adjust for the change-of-variables
 564 making up the constraint. Transformed parameters also define local variables. Everything
 565 executes serially as coded in the Stan program.
- 566 • *Model block*: The rest of the body of the log density function is the line-by-line translation
 567 of Stan's model block. Each distribution and target increment statement increments the log
 568 density accumulator, which is returned as the value of the function. Everything is templated
 569 generally enough and coded in the underlying math library to allow automatic differentiation;
 570 the automatic differentiation code dominates the project in terms of size and presents the most
 571 challenges to developer recruitment because of its modern C++ architecture (Carpenter et al.
 572 2015).
- 573 • *Generated quantities block*: This block is motivated by the goal of providing efficient forward-
 574 simulated predictive inferences. The generated quantities block translates to a function that
 575 maps the parameters and a random number generator to the variables of interest. These are
 576 gathered by this function along with the parameters and transformed parameters to form
 577 vector output. The names of all the parameters define the columns of output and output is on
 578 the constrained scale where the user and Stan program model block operate.

579 5.2 Stan’s transpiled C++ model class

580 Stan programs are transpiled to C++, and then the C++ is compiled down to machine instructions.
581 The C++ class generated for the model has the following signatures, which have been simplified
582 to remove debugging traces, template traits restrictions, and some fine-grained control parameters.
583 This is the same interface that is exposed by BridgeStan in Python, R, Julia, C, and Rust (Roualdes et
584 al. 2023).

```
namespace lr_model_namespace {  
  
class lr_model : model_base<lr_model> {  
  // data, transformed data  
  int N, N_new, P;    MatrixXd x;  
  VectorXd y;        MatrixXd x_new;  
  
  // read data from c, compute transformed data  
  linear_regression_model(var_context& c, int seed,  
                          ostream* msgs);  
  
  // log density type T, propto drops constants, jacobian adjusts  
  template <bool propto, bool jacobian, typename T>  
  Vector<T> log_prob(Vector<T>& params, ostream* msgs) const;  
  
  // evaluate generated quantities and write csv  
  template <typename RNG>  
  void write(RNG& rng, const VectorXd& params,  
             ostream* csv_stream) const;  
  
  string model_name() const noexcept;  
  
  void unconstrain(const VectorXd& constrained_params,  
                  VectorXd& unconstrained_params) const;  
  void constrain(var_context& vars, VectorXd& constrained_params,  
                ostream* msgs) const;  
  
  void param_names(vector<string>& names) const;  
  void constrained_param_names(vector<string>& names) const;  
};  
} // lr_model_namespace  
  
// global namespace factory for lr_model  
model_base& new_model(var_context& c, int seed, ostream* msgs);
```

585 The data variables are specified in the data block in the Stan program. Here, we have the sizes, the
586 covariate matrices (x plus x_new for posterior prediction), and the outcomes (y). The constructed
587 C++ class is immutable and provides several methods, the most central of which is an unconstrained
588 (in the sense of having support over all of R^D) log density function that is templated in order to
589 support automatic differentiation (Carpenter et al. 2015). In math, the constraining transform maps
590 $(\alpha, \beta, \sigma^u)$ to $(\alpha, \beta, \exp(\sigma^u))$. The transforms are independent and the first two are the identity, so the
591 Jacobian determinant works out to $\exp(\sigma^u)$. Thus the additive change-of-variables adjustment on

592 the log scale is just $\log \exp(\sigma^u) = \sigma^u$.

593 The model block defines a density function p over constrained parameters. Although Stan can be
594 used to sample from any density, it is typically used to code the unnormalized posterior log density
595 of a Bayesian mode. Densities are read elementwise, with scalar arguments being broadcast where
596 necessary. With this translation, the unnormalized log posterior over the constrained variables
597 defined by the Stan program's model block is

$$\log p(\alpha, \beta, \sigma \mid x, y) = \log \text{normal}(\alpha \mid 0, 5) + \sum_{p=1}^P \log \text{normal}(\beta_p \mid 0, 2.5) + \log \text{gamma}(\sigma \mid 0.5) + \sum_{n=1}^N \log \text{normal}(y_n \mid \alpha + \beta \cdot x_n, \sigma).$$

598 The corresponding unconstrained log density over which inference is performed, adds the log
599 Jacobian adjustment for the change of variables,

$$\log p^u(\alpha, \beta, \sigma^u \mid x, y) = \log p(\alpha, \beta, \exp(\sigma^u) \mid x, y) + \sigma^u.$$

600 The transpiled C++ code implementing the log density function performs the following sequence of
601 operations, which directly follows the Stan program. The typing, logging, debugging traces, message
602 I/O, and name mangling have all been simplified.

```
template <bool propto, bool jacobian, typename T>
inline auto log_prob(Vector<T>& params) const {
    Accumulator<T> accum;

    // transpiled parameters block
    T log_jacobian;
    Deserializer<T> in(params)
    auto alpha = in.template read<T>();
    auto beta = in.template read<Vector<T>>(P);
    auto sigma = in_.template read_constrain_lb<T, jacobian>(0, log_jacobian);
    accum.add(log_jacobian);

    // transpiled model block
    accum.add(normal_lpdf<propto>(alpha, 0, 5));
    accum.add(normal_lpdf<propto>(beta, 0, 2.5));
    accum.add(exponential_lpdf<propto>(sigma, 0.5));
    accum.add(normal_lpdf<propto>(y, add(alpha, multiply(beta, x)), sigma));

    return accum.sum();
}
```

603 The read methods of the deserializer read the next variable from the sequence of parameters, perform
604 any necessary constraining transform and add any Jacobian adjustment to $\log p$, controlled by the
605 flag `jacobian`. Arithmetic in the Stan program is replaced with overloaded internal library calls
606 `add` and `multiply` which can handle vectors and scalars. There is implicit broadcasting and internal
607 summation in the `normal_lpdf` log density function. When applied to a vector like `beta` or `y`, it
608 sums the log densities of the components, broadcasting any scalar arguments as necessary. The
609 accumulator collects terms and then its method `sum` returns a flattened autodiff tree with a single
610 root pointing to all the terms with implicit unit derivatives.

611 To support the changes of variables for reporting constrained output, the compiled C++ code exposes
612 the constraining transform and its Jacobians. For initialization from constrained parameters, there is

613 a matching unconstraining transform. The model class also supplies methods for determining the
614 shape and names of the constrained and unconstrained parameters.

615 The final component of a compiled Stan model is a function to perform predictive inference as defined
616 by the generated quantities block (this can also be done post-hoc with a new program and generated
617 quantities block). In particular, the model compiles a generated quantities function `write()` that
618 takes a random number generator and produces the output defined in the generated quantities block
619 purely by forward sampling without the need for any automatic differentiation.

620 **6 Coding models directly in JAX**

621 In this section, we will translate the functionality provided by Stan’s C++ class directly into a more
622 functional JAX style. It’s much simpler than Stan’s transpiled C++ code because of the built-in
623 serialization and optimization of PyTree and the binding of data. In the next section, we will wrap
624 up some utility functions to make it much easier to achieve the same functionality.

625 **6.1 Linear regression directly in JAX**

626 This section provides one way to code our linear regression model directly in JAX. We will start
627 by importing the libraries we need, which include JAX’s version of NumPy, JAX’s random number
628 generator, and the `stats` library from JAX’s version of SciPy. Finally, we import the Blackjax library
629 for sampling. We first import all of JAX itself as we will need it for just-in-time compilation.

630 Before we begin, we will also convert the Python and NumPy data types to JAX.

631 **6.2 The untransformed log posterior in JAX**

632 The log density function from the Stan model block can then be translated almost line by line into
633 JAX; we only break out the definition of `mu` for readability—it could have been a nested expression.

634 The `log_posterior` function binds the covariate matrix `x` and the outcome vector `y` from the top-level
635 environment. An alternative would have been to bind a JAX conversion of top-level container data,
636 but we never need the encapsulated data, so just define the variables directly. The log probability
637 density functions are implemented in JAX’s version of SciPy’s `stats` library. Unlike in Stan, where
638 density functions automatically reduce by summation, the functions in `stats.norm` return JAX
639 arrays, which are here summed by JAX’s version of NumPy before being added to the target log
640 density accumulator `lp`.

641 The parameters themselves are read out of the argument dictionary. While this may appear inefficient,
642 JAX’s just-in-time (JIT) compiler is smart enough to reduce the dictionary lookups to direct accesses
643 at run time.

644 **6.3 Transforms and their inverses in JAX**

645 Before we can use this log posterior for sampling, we need a few support functions. To mirror Stan,
646 the sampler needs to work on the transformed scale. Here, as in most programs, the transform
647 maps constrained variables to unconstrained ones. In our regression example, the scale, which is
648 constrained to $(0, \infty)$, is log-transformed so that the transformed scale ranges over the whole real
649 number line, $(-\infty, \infty)$.

650 We will start with the transform, which is only needed in practice for initializations from untrans-
651 formed parameter values. In the linear regression example, the transform maps the regression
652 coefficients to themselves and the scale to the log scale.

653 The `transform()` function could have been implemented in one line. Instead, we chose to follow a
654 variable-by-variable pattern that will make it easier to automate later.

655 The inverse transform is needed at runtime to map the transformed variables back to their untrans-
656 formed state so that the log posterior may be applied. For the inverse transform, any adjustment to
657 ensure the implied distribution on the untransformed parameter space is uniform is returned along
658 with the transformed parameters. Typically, this is a Jacobian determinant to adjust for a bijective
659 change-of-variables. We use the same variable names for both transformed and untransformed
660 parameters. The adjustment for the log change of variables here is just the log scale itself, as we
661 derived above. To keep the code structure regular and to ease the transition to automation, we use a
662 log adjustment variable that is initialized to zero and updated after each inverse transform.

663 6.4 The transformed log density in JAX

664 The transformed log posterior function we will use for sampling can be defined to apply the inverse
665 transform and add the adjustment.

666 6.5 Random initialization in JAX

667 Following Stan, we need a way to generate a transformed random initialization. Here we use JAX's
668 version of SciPy's `stats.random` library (imported above as `jrd`).

669 The random number generators are supplied with shapes for variables that are not scalars. Let's
670 generate an initialization.

```
671 t_params_init={'alpha': Array(-0.48995897, dtype=float32), 'beta': Array([0.19053426, 1.4290651 ], dtype=float32)}
```

672 As a consistency check, we verify that the parameters we used to generate data are round-trippable
673 through the transform and inverse transform.

```
674 params_init={'alpha': Array(-0.48995897, dtype=float32), 'beta': Array([0.19053426, 1.4290651 ], dtype=float32)}  
675 log_adjust=Array(0.23013926, dtype=float32)  
676 t_params_init_round_trip={'alpha': Array(-0.48995897, dtype=float32), 'beta': Array([0.19053426, 1.4290651 ], dtype=float32)}
```

677 The round trip arithmetic is not exact for σ , because JAX works with 32-bit float types (`float32`)
678 by default.

679 6.6 Sampling in Blackjax

680 Before sampling, we will define a top-level Markov chain sampler and then we will instantiate it
681 with the NUTS transition kernel configured to run in a way that matches Stan's defaults.

682 Working outward, the nested function `one_step` takes one step of a Markov chain from a given state
683 and returns the next state. It uses a supplied transition kernel and random key. It starts from the
684 specified initial state and produces the specified number of draws. It uses JAX's `lax.scan` to execute
685 the loop, repeatedly applying the one-step transition function. The NUTS sample function is the top
686 level caller. it takes a PRNG key, the target log density, an initial position, and a number of draws
687 that specifies the number of warmup draws and the number of sampling draws. It uses Blackjax's
688 windowed warmup, as defined by Stan, then the NUTS sampler. The algorithm was defined to mirror
689 Stan's reference implementation.

690 We've already generated a transformed initialization `t_params_init`, so all that remains is to call
691 the sample function.

6.7 Posterior analysis in JAX

We will do just a bit of posterior analysis manually to make sure we're on the right track. Before that, we need to inverse transform our draws. We'll do that with a vectorized function and a single call.

Now we can perform the posterior analysis and see that it matches the results from Stan to within the tolerances to be expected from the standard errors.

```
posterior_means={'alpha': Array(-9.174101, dtype=float32), 'beta': Array([-4.808164 ,  1.1481543]),
posterior_stds={'alpha': Array(0.05411908, dtype=float32), 'beta': Array([0.04893167, 0.03037017]),
```

6.8 Posterior predictive quantities in JAX

In Stan, posterior predictive quantities are typically defined in the generated quantities block. In JAX, we can lazily define posterior predictive functions after sampling and efficiently map them over the posterior draws. Here's a function that matches the generated quantities block defined in the Stan implementation of linear regression. We have already converted the covariates `x_new` to JAX above.

Unlike the `write()` function in Stan's transpiled C++ class, there is no need to inverse transform—we already have the parameters on the natural scale at this point and can just pass them to the predictive function.

The `vmap` is the way that JAX maps a function in parallel. The call above is equivalent to

```
results = []
for i in range(S):
    results.append(posterior_predictive(keys[i], params_draws[i]))
return results
```

The specification `in_axes=(0, 0)` tells the application of `posterior_predictive` to bind to the first argument (i.e., to index position 0).

We can summarize the results as before and verify they line up with the output from Stan's generated quantities block.

```
posterior_pred_means={'y_new': Array([-10.565608, -16.299402, -17.768396, -
8.535892], dtype=float32)}
posterior_pred_stds={'y_new': Array([0.5454136 , 0.5331531 , 0.5449212 , 0.55034816], dtype=float32)}
```

6.9 Key differences between the Stan and JAX implementations

Some of the bigger differences between Stan and its JAX implementation are as follows.

1. There are no data declarations. The functions simply bind (i.e., close over) the data variables `x` and `y` (i.e., it reads them from the environment where it was defined); we will shortly show how to abstract this into a function.
2. JAX does not require parameters to be serialized. With the wizardry of PyTree, Blackjax and other JAX-based packages can work directly with the `log_prob` function without any need for serialization. If serialized log density functions *are* needed elsewhere, they are convertible with a single function call as we show below.
3. There is no need to work with accumulators—JAX's just-in-time compilation is sufficient to optimize sequences of operations at the XLA substrate to which JAX is compiled.
4. The generated quantities are not part of the basic model specification, but rather flexibly called later. While we could have done the same with Stan, the Stan functions are not available in either R or Python, making this another two-language problem. Furthermore, JAX's ability to efficiently scan accelerates posterior predictive inference in JAX.

730 5. There are no function blocks because plain old Python functions can be applied to JAX objects.
731 The only caveat is that these functions will be traced, so cannot contain any runtime branching
732 on parameter values.

733 Aside from the verbosity of all the namespace qualifiers and the need for all the intermediate calls to
734 `sum`, the biggest obstacle to writing code this way is having to manually deal with the transforms
735 and inverse transforms. Our goal is to make it as simple as Stan in the next section.

736 7 Linear regression in JAX with `densejax`

737 The Stan linear regression can be translated almost line for line into Python using `densejax`.

```
from densejax import (
    real, positive, normal, exponential, normal_rng, model
)

def linear_regression(x, y, x_new):
    N, P = x.shape

    parameters = {
        'a': real(),
        'b': real(size=P),
        's': positive()
    }

    def log_density(a, b, s):
        lp = 0
        lp += normal(a, 0, 2)
        lp += normal(b, 0, 1)
        lp += exponential(s, 0.5)
        lp += normal(y, a + x @ b, s)
        return lp

    def generate(rng, a, b, s):
        y_new = normal_rng(rng, a + x_new @ b, s)
        return { 'y_new': y_new }

    return model(parameters, log_density, generate)
```

738 After calling this function with data, the resulting `model` object is feature equivalent to all the manual
739 code written above for transforming parameters and performing initialization. It mirrors the C++
740 object produced by Stan (also the model object produced by BridgeStan (Roualdes et al. 2023)).
741 Because these simple functions generate fully JAX-embedded code, the log density function can be
742 automatically differentiated and all of the functions can be just-in-time compiled.

743 7.1 Abstracting parameter transformations

744 In order to automate the transformation functionality that were coded manually above, we introduce
745 the idea of a “parameter specification” (the `parameters` dictionary in the code) which provides all
746 the information about how the parameters should be transformed. The leaves of this parameter
747 specification are simple classes like `positive`, which is reproduced here in full:

```

class positive(ParameterConstraint):
    def __init__(self, shape=(), dtype=jnp.float32):
        self.shape = shape
        self.dtype = dtype

    def inverse_transform(self, x):
        return jnp.exp(x)

    def transform(self, y):
        return jnp.log(y)

    def jacobian(self, x):
        return jnp.sum(x)

```

748 The densejax package provides a menu of existing transforms, but the above class is designed to be
749 simple enough to be easily extensible.

750 The real class is simply a transform that has the identity for both sides of the transformation; we
751 rely on JAX's jit compiler to remove any overhead from this implementation style.

752 Turning this parameter spec into functions that perform the inverse transforms, transforms, and
753 (possibly random) initialization is trivial. All of these reduce to simple map-like operations over the
754 specification and the provided parameter values where required.

755 The majority of the code in densejax consists of these transformation classes, utility functions
756 consuming them, and the small wrappers provided for common distributions like normal as seen
757 above.

758 7.2 Simplifying inference

759 The model object remains general enough for use in any gradient based sampler implemented in
760 JAX; the one-liner

```
log_posterior_transformed = jax.jit(model.log_density)
```

761 would be enough to use blackjax in the same manner as above.

762 However, bundling more of the desired functionality into one object allows for more turn-key helpers.
763 Consider the following new definition of nuts_sample:

```

def nuts_sample(key, model, num_draws, init=None):
    log_density = jax.jit(model.log_density)

    init_key, warmup_key, sample_key = jrd.split(key, 3)

    if init is None:
        init_position = model.initialize_random(init_key)
    else:
        init_position = model.initialize(init)

    warmup = blackjax.window_adaptation(blackjax.nuts, log_density)
    (state, params), _ = warmup.run(warmup_key, init_position, num_steps=num_draws)

    kernel = blackjax.nuts(log_density, **params).step

```

```

constrain = jax.jit(model.inv_transform)
@jax.jit
def one_step(state, key):
    state, _ = kernel(key, state)
    return state, constrain(**state.position)

keys = jrd.split(sample_key, num_draws)
_, draws = jax.lax.scan(one_step, state, keys)

return draws

```

764 This version of `nuts_sample` has the advantage over the previous that both the input initial values
765 *and the returned draws* are in the model space. The user on the outside does not need to worry about
766 transforming the initial values and inverse-transforming the results.

767 We could have written `nuts_sample` to do this for us before, if we had been willing to pass it additional
768 arguments; the `model` class makes it very natural to define the most useful version out of the box.

769 Given these snippets, one can recreate the complete sampling workflow from the prior section with
770 the following brief code:

```

seed = 441_582
key = jrd.key(seed)

model = linear_regression(x, y, x_new)
num_draws = 1_000
draws = nuts_sample(key, model, num_draws)

mean = functools.partial(jax.tree.map(functools.partial(jnp.mean, axis=0)))
std = functools.partial(jax.tree.map(functools.partial(jnp.mean, axis=0)))

print(f"{mean(draws)=}")
print(f"{std(draws)=}")

# posterior predictive
key, gq_keys = jrd.split(key, 1+num_draws)
pred_draws = jax.vmap(model.forward, in_axes=(0, 0))(keys, draws)

print(f"{mean(pred_draws)=}")
print(f"{std(draws)=}")

```

771 References

772 Bahadur, R Raj. 1966. "A Note on Quantiles in Large Samples." *The Annals of Mathematical Statistics*
773 37 (3): 577–80.

774 Betancourt, Michael. 2017. "A Conceptual Introduction to Hamiltonian Monte Carlo." *arXiv*
775 1701.02434.

776 Blondel, Mathieu, Quentin Berthet, Marco Cuturi, et al. 2021. "Efficient and Modular Implicit
777 Differentiation." *arXiv* 2105.15183.

- 778 Cabezas, Alberto, Adrien Corenflos, Junpeng Lao, et al. 2024. “BlackJAX: Composable Bayesian
779 Inference in JAX.” *arXiv* 2402.10797.
- 780 Carpenter, Bob, Andrew Gelman, Matthew D Hoffman, et al. 2017. “Stan: A Probabilistic Program-
781 ming Language.” *Journal of Statistical Software* 76: 1–32.
- 782 Carpenter, Bob, Matthew D Hoffman, Marcus Brubaker, Daniel Lee, Peter Li, and Michael Betan-
783 court. 2015. “The Stan Math Library: Reverse-Mode Automatic Differentiation in C++.” *arXiv*
784 1509.07164.
- 785 Cook, Samantha, Andrew Gelman, and Donald B. Rubin. 2006. “Validation of Software for Bayesian
786 Models Using Posterior Quantiles.” *Journal of Computational and Graphical Statistics* 15 (3):
787 675–92.
- 788 DeepMind, Igor Babuschkin, Kate Baumli, et al. 2020. *The DeepMind JAX Ecosystem*. Released.
789 <http://github.com/deepmind>.
- 790 Dillon, Joshua V, Ian Langmore, Dustin Tran, et al. 2017. “Tensorflow Distributions.” *arXiv* 1711.10604.
- 791 Dinh, Laurent, Jascha Sohl-Dickstein, and Samy Bengio. 2016. “Density Estimation Using Real NVP.”
792 *arXiv* 1605.08803.
- 793 Doucet, Arnaud, Nando De Freitas, and Neil Gordon. 2001. “An Introduction to Sequential Monte
794 Carlo Methods.” In *Sequential Monte Carlo Methods in Practice*. Springer.
- 795 Gabry, Jonah, Daniel Simpson, Aki Vehtari, Michael Betancourt, and Andrew Gelman. 2019. “Visual-
796 ization in Bayesian Workflow.” *Journal of the Royal Statistical Society Series A: Statistics in Society*
797 182 (2): 389–402.
- 798 Ge, Hong, Kai Xu, and Zoubin Ghahramani. 2018. “Turing: A Language for Flexible Probabilistic
799 Inference.” *International Conference on Artificial Intelligence and Statistics*, 1682–90.
- 800 Gelman, Andrew, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin.
801 2013. *Bayesian Data Analysis*. Third. Chapman; Hall/CRC.
- 802 Gelman, Andrew, Aki Vehtari, Daniel Simpson, et al. 2020. “Bayesian Workflow.” *arXiv* 2011.01808.
- 803 Gorinova, Maria I, Andrew D Gordon, and Charles Sutton. 2019. “Probabilistic Programming
804 with Densities in SlicStan: Efficient, Flexible, and Deterministic.” *Proceedings of the ACM on*
805 *Programming Languages* 3 (POPL): 1–30.
- 806 Guennebaud, Gaël, and Benoît Jacob. 2010. “Eigen: A C++ linear algebra library.” *Eurograph-
807 ics/CGLibs*.
- 808 Hastings, W. K. 1970. “Monte Carlo Sampling Methods Using Markov Chains and Their Applications.”
809 *Biometrika* 57: 97–109.
- 810 Heek, Jonathan, Anselm Levskaya, Avital Oliver, et al. 2024. *Flax: A Neural Network Library and*
811 *Ecosystem for JAX*. V. 0.12.3. Released. <http://github.com/google/flax>.

- 812 Hoffman, Matthew D, and Andrew Gelman. 2014. “The No-U-Turn Sampler: Adaptively Setting Path
813 Lengths in Hamiltonian Monte Carlo.” *Journal of Machine Learning Research* 15 (1): 1593–623.
- 814 Horowitz, Alan M. 1991. “A Generalized Guided Monte Carlo Algorithm.” *Physics Letters B* 268 (2):
815 247–52.
- 816 Josuttis, Nicolai M. 2012. *The C++ Standard Library: A Tutorial and Reference*. Addison-Wesley.
- 817 Kidger, Patrick. 2021. “On Neural Differential Equations.” PhD thesis, University of Oxford.
- 818 Kochkov, Dmitrii, Jamie A. Smith, Ayya Alieva, Qing Wang, Michael P. Brenner, and Stephan Hoyer.
819 2021. “Machine Learning–accelerated Computational Fluid Dynamics.” *Proceedings of the National*
820 *Academy of Sciences* 118 (21). <https://doi.org/10.1073/pnas.2101784118>.
- 821 Kucukelbir, Alp, Dustin Tran, Rajesh Ranganath, Andrew Gelman, and David M Blei. 2017. “Automatic
822 Differentiation Variational Inference.” *Journal of Machine Learning Research* 18 (14): 1–45.
- 823 Kumar, Ravin, Colin Carroll, Ari Hartikainen, and Osvaldo Martin. 2019. “ArviZ a Unified Library
824 for Exploratory Analysis of Bayesian Models in Python.” *Journal of Open Source Software* 4 (33):
825 1143.
- 826 Lunn, David, Chris Jackson, Nicky Best, Andrew Thomas, and David Spiegelhalter. 2012. *The BUGS*
827 *Book: A Practical Introduction to Bayesian Analysis*. CRC Press.
- 828 Lunn, David, David Spiegelhalter, Andrew Thomas, and Nicky Best. 2009. “The BUGS Project:
829 Evolution, Critique and Future Directions.” *Statistics in Medicine* 28 (25): 3049–67.
- 830 Margossian, Charles C. 2023. “General Adjoint-Differentiated Laplace Approximation.” *arXiv*
831 2306.14976.
- 832 Maskell, Simon. 2024. *Running Multiple Short MCMC Chains on a GPU Using jAX for Fast Inference*
833 *with Stan*. YouTube! video. <https://www.youtube.com/watch?v=KpLZEYX8MpY>.
- 834 Murray, Iain, Ryan Adams, and David MacKay. 2010. “Elliptical Slice Sampling.” *Proceedings of the*
835 *13th International Conference on Artificial Intelligence and Statistics*, 541–48.
- 836 Paulin, Daniel. 2015. “Concentration Inequalities for Markov Chains by Marton Couplings and
837 Spectral Methods.” *Electronic Journal of Probability* 20 (79): 1–32. [https://doi.org/10.1214/EJP.v20-](https://doi.org/10.1214/EJP.v20-4039)
838 [4039](https://doi.org/10.1214/EJP.v20-4039).
- 839 Phan, Du, Neeraj Pradhan, and Martin Jankowiak. 2019. “Composable Effects for Flexible and
840 Accelerated Probabilistic Programming in NumPyro.” *arXiv* 1912.11554.
- 841 Plummer, Martyn. 2003. “JAGS: A Program for Analysis of Bayesian Graphical Models Using Gibbs
842 Sampling.” *Proceedings of the 3rd International Workshop on Distributed Statistical Computing*
843 (Vienna, Austria) 124: 1–10.
- 844 Roberts, Gareth O, and Jeffrey S Rosenthal. 2004. “General State Space Markov Chains and MCMC
845 Algorithms.” *Probability Surveys* 1: 20–71.

- 846 Robnik, Jakob, Reuben Cohn-Gordon, and Uroš Seljak. 2025. “Metropolis Adjusted Microcanonical
847 Hamiltonian Monte Carlo.” *arXiv* 2503.01707.
- 848 Roualdes, Edward A, Brian Ward, Bob Carpenter, Adrian Seyboldt, and Seth D Axen. 2023. “BridgeS-
849 tan: Efficient in-Memory Access to the Methods of a Stan Model.” *Journal of Open Source Software*
850 8 (87): 5236.
- 851 Rubin, Donald B. 1996. “Multiple Imputation After 18+ Years.” *Journal of the American Statistical*
852 *Association* 91 (434): 473–89.
- 853 Salvatier, John, Thomas V Wiecki, and Christopher Fonnesbeck. 2016. “Probabilistic Programming
854 in Python Using PyMC3.” *PeerJ Computer Science* 2: e55.
- 855 Sountsov, Pavel, Colin Carroll, and Matthew D Hoffman. 2024. “Running Markov Chain Monte Carlo
856 on Modern Hardware and Software.” *arXiv* 2411.04260.
- 857 Stan Development Team. 2024. *CmdStanPy*. V. 1.3.0. Released. [https://doi.org/10.5281/zenodo.
858 5733022](https://doi.org/10.5281/zenodo.5733022).
- 859 Stan Development Team. 2025. *Stan Reference Manual*. Version 2.38. Stan Project. [https://mc-
860 stan.org/docs/reference-manual/](https://mc-stan.org/docs/reference-manual/).
- 861 Valpine, Perry de, Daniel Turek, Christopher J Paciorek, Clifford Anderson-Bergman, Duncan Temple
862 Lang, and Rastislav Bodik. 2017. “Programming with Models: Writing Statistical Algorithms for
863 General Model Structures with NIMBLE.” *Journal of Computational and Graphical Statistics* 26 (2):
864 403–13.
- 865 Van der Vaart, A W. 1998. *Asymptotic Statistics*. Cambridge Series on Statistical and Probabilistic
866 Mathematics 3. Cambridge University Press.
- 867 Vehtari, Aki, Andrew Gelman, and Jonah Gabry. 2017. “Practical Bayesian Model Evaluation Using
868 Leave-One-Out Cross-Validation and WAIC.” *Statistics and Computing* 27 (5): 1413–32.
- 869 Vehtari, Aki, Andrew Gelman, Daniel Simpson, Bob Carpenter, and Paul-Christian Bürkner. 2021.
870 “Rank-Normalization, Folding, and Localization: An Improved R-Hat for Assessing Convergence
871 of MCMC (with Discussion).” *Bayesian Analysis* 16 (2): 667–718.
- 872 Zhang, Lu, Bob Carpenter, Andrew Gelman, and Aki Vehtari. 2022. “Pathfinder: Parallel Quasi-
873 Newton Variational Inference.” *Journal of Machine Learning Research* 23 (306): 1–49.