

Fast and Accurate Approximation of the Full Conditional for Gamma Shape Parameters

Jeffrey W. Miller*

Harvard University, Department of Biostatistics

July 30, 2018

Abstract

The gamma distribution arises frequently in Bayesian models, but there is not an easy-to-use conjugate prior for the shape parameter of a gamma. This inconvenience is usually dealt with by using either Metropolis–Hastings moves, rejection sampling methods, or numerical integration. However, in models with a large number of shape parameters, these existing methods are slower or more complicated than one would like, making them burdensome in practice. It turns out that the full conditional distribution of the gamma shape parameter is well approximated by a gamma distribution, even for small sample sizes, when the prior on the shape parameter is also a gamma distribution. This article introduces a quick and easy algorithm for finding a gamma distribution that approximates the full conditional distribution of the shape parameter. We empirically demonstrate the speed and accuracy of the approximation across a wide range of conditions. If exactness is required, the approximation can be used as a proposal distribution for Metropolis–Hastings.

Keywords: Bayesian, Generalized Newton, Hierarchical models, Markov chain Monte Carlo, Sampling.

*Contact: jwmiller@hsph.harvard.edu. Department of Biostatistics, Harvard School of Public Health, 655 Huntington Ave 1-419, Boston MA, 02115.

1 Introduction

The lack of a nice conjugate prior for the shape parameter of the gamma distribution is a commonly occurring nuisance in many Bayesian models. Conjugate priors for the shape parameter do exist, but unfortunately, they are not analytically tractable (Damsleth, 1975; Miller, 1980). In Markov chain Monte Carlo (MCMC) algorithms, one can easily use Metropolis–Hastings (MH) moves to update the shape parameter, but mixing can be slow if the proposal distribution is not well-calibrated. Rejection sampling schemes can also be used for sampling the shape parameter (Son and Oh, 2006; Pradhan and Kundu, 2011), however, the more efficient schemes such as adaptive rejection sampling (Gilks and Wild, 1992) are complicated and require log-concavity, which does not always hold under some important choices of prior. Modern applications involve models with a large number of parameters, necessitating fast inference algorithms that work very generally with minimal tuning. Indeed, our motivation for developing the approach in this article is a gene expression model involving tens of thousands of gamma shape parameters.

In this article, we introduce a fast and simple algorithm for finding a gamma distribution that approximates the full conditional distribution of the gamma shape parameter, when the prior on the shape parameter is also a gamma distribution. This algorithm can be used to perform approximate Gibbs updates by sampling from the approximating gamma distribution. Alternatively, the approximation can be used as an MH proposal distribution to make a move that exactly preserves the full conditional and has high acceptance rate.

The basic idea of the algorithm is to approximate the full conditional density f by a gamma density g chosen such that the first and second derivatives of $\log g$ match those of $\log f$ at a point near the mean of f . Since the mean of f is not known in closed form, the approximation is iteratively refined by matching derivatives at the mean of the current g .

The article is organized as follows. Section 2 describes the algorithm, and Section 3 contains an empirical assessment of the accuracy and speed of the algorithm. Section 4 discusses previous work, and Section 5 provides mathematical details on the derivation of the algorithm. The supplementary material contains a characterization of the fixed points of the algorithm, additional derivations, and additional empirical results.

2 Algorithm

Consider the model $X_1, \dots, X_n | a, \mu \sim \text{Gamma}(\text{shape} = a, \text{rate} = a/\mu)$ i.i.d., where $a, \mu > 0$, and note that this makes $\mu = \mathbb{E}(X_i | a, \mu)$. This parametrization is convenient, since a controls the concentration and μ controls the mean. We assume a gamma prior on the shape, say, $a \sim \text{Gamma}(\text{shape} = a_0, \text{rate} = b_0)$ where, if desired, a_0 and b_0 can depend on μ . The following algorithm produces A and B such that

$$p(a | x_1, \dots, x_n, \mu, a_0, b_0) \approx \text{Gamma}(a | \text{shape} = A, \text{rate} = B).$$

We use $\psi(x)$ and $\psi'(x)$ to denote the digamma and trigamma functions, respectively, and $\log(x)$ denotes the natural log.

Algorithm 1: Approximating the full conditional of the shape parameter

input : data $x_1, \dots, x_n > 0$, parameters $\mu, a_0, b_0 > 0$, tolerance $\epsilon > 0$, and maximum number of iterations M .

output: A and B .

begin

$R \leftarrow \sum_{i=1}^n \log(x_i)$, $S \leftarrow \sum_{i=1}^n x_i$, and $T \leftarrow S/\mu - R + n \log(\mu) - n$

$A \leftarrow a_0 + n/2$ and $B \leftarrow b_0 + T$

for $j = 1, \dots, M$ **do**

$a \leftarrow A/B$

$A \leftarrow a_0 - na + na^2\psi'(a)$

$B \leftarrow b_0 + (A - a_0)/a - n \log(a) + n\psi(a) + T$

if $|a/(A/B) - 1| < \epsilon$ **then return** A, B

return A, B

We recommend setting $\epsilon = 10^{-8}$ and $M = 10$. Using $\epsilon = 10^{-8}$, the algorithm terminates in four iterations or less in all of the simulations we have done, so $M = 10$ is conservative. Most languages have routines for the digamma and trigamma functions, which are defined by $\psi(x) = \frac{\partial}{\partial x} \log \Gamma(x)$ and $\psi'(x) = \frac{\partial}{\partial x} \psi(x)$, where $\Gamma(x)$ is the gamma function. See Section 5 for the derivation of the algorithm. See the supplement for a fixed point analysis.

To perform an approximate Gibbs update to a in an MCMC sampler, run Algorithm 1 to obtain A, B using the current values of $x_{1:n}, \mu, a_0, b_0$, and then sample $a \sim \text{Gamma}(A, B)$.

(For brevity, in the rest of the paper, we write $x_{1:n}$ to denote (x_1, \dots, x_n) , and we write $\text{Gamma}(\alpha, \beta)$ to denote the gamma distribution with shape α and rate β .)

Alternatively, to perform a Metropolis–Hastings move to update a , one can run Algorithm 1 to obtain A, B using the current values of $x_{1:n}, \mu, a_0, b_0$, then sample a proposal $a' \sim \text{Gamma}(A, B)$, and accept with probability

$$\min \left\{ 1, \frac{p(a' | x_{1:n}, \mu, a_0, b_0) \text{Gamma}(a|A, B)}{p(a | x_{1:n}, \mu, a_0, b_0) \text{Gamma}(a'|A, B)} \right\},$$

where a is the current value of the shape parameter. Note that the normalization constant of $p(a | x_{1:n}, \mu, a_0, b_0)$ does not need to be known since it cancels.

To infer the mean μ , note that the inverse-gamma family is a conditionally conjugate prior for μ given a . Therefore, Gibbs updates to μ can easily be performed when using an inverse-gamma prior on μ , independently of a , given a_0 and b_0 .

3 Performance assessment

To assess the performance of Algorithm 1, we evaluated its accuracy in approximating the true full conditional distribution, and its speed in terms of the number of iterations until convergence, across a wide range of simulated conditions.

Specifically, for each combination of $n \in \{1, 10, 100\}$, $r \in \{0.5, 1, 2\}$, $a_{\text{true}} \in \{10^{-6}, 10^{-5}, \dots, 10^5, 10^6\}$, $\mu_{\text{true}} \in \{10^{-6}, 10^{-5}, \dots, 10^5, 10^6\}$, and $a_0 \in \{1, 0.1, 0.01\}$, we generated five independent data sets $x_1, \dots, x_n \sim \text{Gamma}(a_{\text{true}}, a_{\text{true}}/\mu_{\text{true}})$ i.i.d., ran Algorithm 1 with inputs $x_{1:n}, \mu = r\mu_{\text{true}}, a_0$, and $b_0 = a_0$, and then compared the resulting $\text{Gamma}(a|A, B)$ approximation to the true full conditional distribution $p(a | x_{1:n}, \mu, a_0, b_0)$. For the prior parameters, we set $b_0 = a_0$ in each case so that the prior mean is 1. In each case, we used $\epsilon = 10^{-8}$ for the convergence tolerance, and $M = 10$ for the maximum number of iterations.

3.1 Approximation accuracy

First, to get a visual sense of the closeness of the approximation, Figure 1 shows the cumulative distribution functions (CDFs) of the true and approximate full conditional

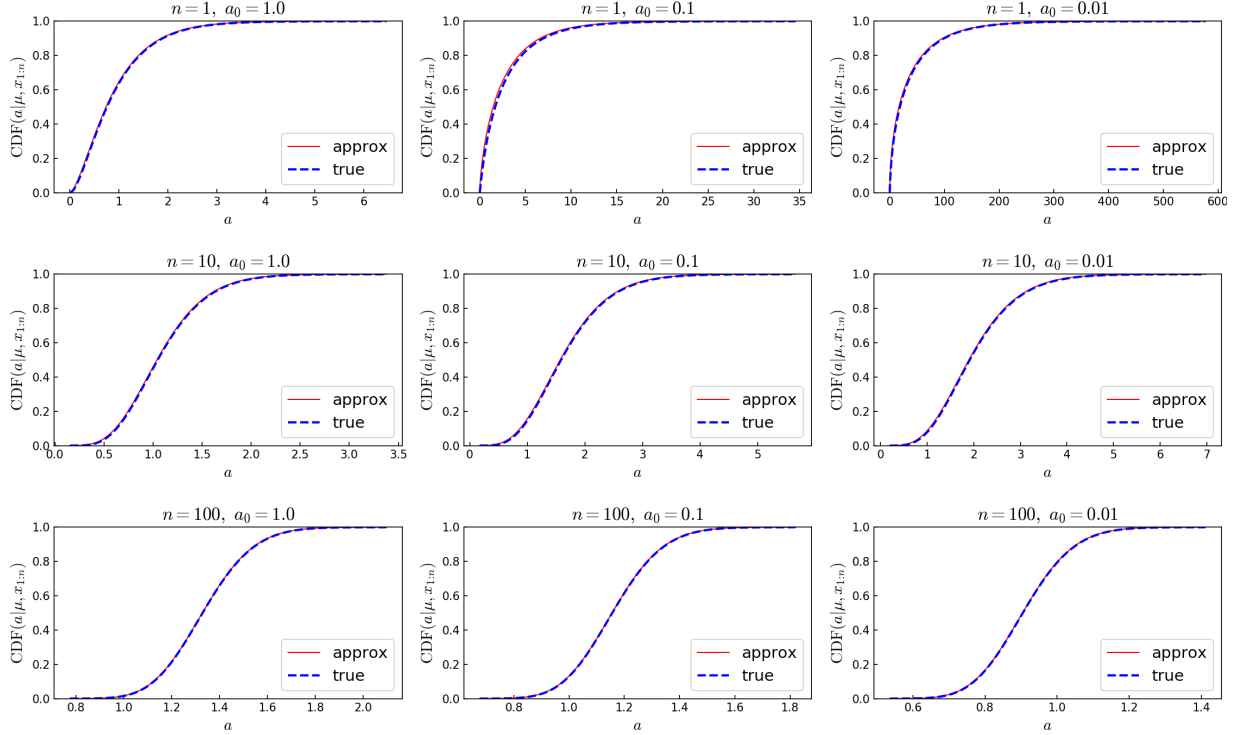


Figure 1: Cumulative distribution function (CDF) of the true full conditional $p(a | x_{1:n}, \mu, a_0, b_0)$ and approximate full conditional, on simulated data $x_1, \dots, x_n \sim \text{Gamma}(1, 1)$, when conditioning on $\mu = 1$ and the prior is $a \sim \text{Gamma}(a_0, a_0)$.

distributions for a single simulated data set for each case with $r = 1$, $a_{\text{true}} = 1$, and $\mu_{\text{true}} = 1$. In each case, the approximation is close to the true distribution, and in most cases the approximation is visually indistinguishable from the truth. These plots are representative of the full range of cases considered.

To quantify the discrepancy between the truth and the approximation, for each simulation run we computed:

1. the total variation distance, $d_{\text{TV}}(f, g) = \frac{1}{2} \int |f(a) - g(a)| da$,
2. the Kullback–Leibler divergence, $d_{\text{KL}}(f, g) = \int f(a) \log(f(a)/g(a)) da$, and
3. the reverse Kullback–Leibler divergence, $d_{\text{KL}}(g, f) = \int g(a) \log(g(a)/f(a)) da$,

where $f(a) = p(a | x_{1:n}, \mu, a_0, b_0)$ is the true density and $g(a) = \text{Gamma}(a|A, B)$ is the approximate density. Numerical integration was performed to compute these integrals in

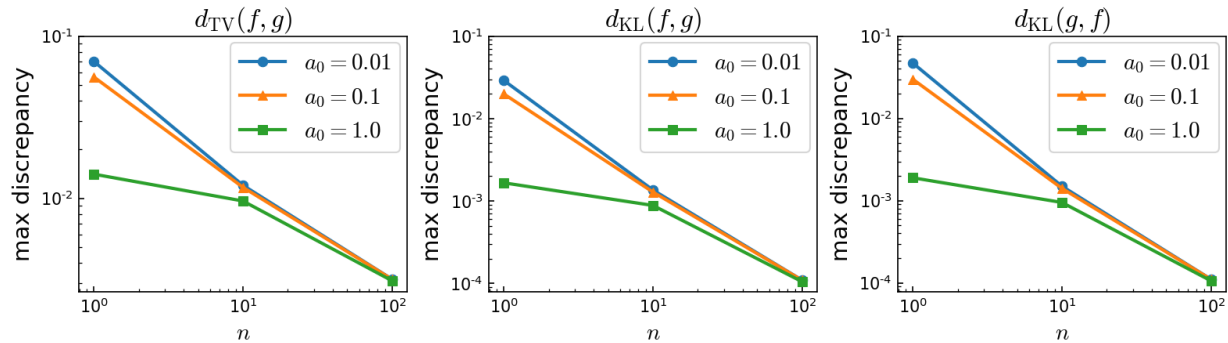


Figure 2: Largest observed discrepancy between the true full conditional f and the approximate full conditional g , for the total variation distance $d_{\text{TV}}(f, g)$ and the Kullback–Leibler divergences $d_{\text{KL}}(f, g)$ and $d_{\text{KL}}(g, f)$, for each a_0 and n .

order to compare with the truth; see the supplementary material for details.

Figure 2 shows the largest observed value (that is, the worst-case value) of $d_{\text{TV}}(f, g)$, $d_{\text{KL}}(f, g)$, and $d_{\text{KL}}(g, f)$ across all cases for each a_0 and n , using the average discrepancy over the five data sets for each case. This shows that the approximation is quite good across a very wide range of cases. The approximation accuracy improves as n increases, and the accuracy is decent even when n is small. The accuracy improves as a_0 increases, presumably because this makes the prior stronger and the prior is itself a gamma distribution, so it makes sense that the fit of a gamma approximation would improve.

For a more fine-grained breakdown of the discrepancies, case-by-case, see the supplementary material for heatmaps showing $d_{\text{TV}}(f, g)$, $d_{\text{KL}}(f, g)$, and $d_{\text{KL}}(g, f)$ across a range of cases. The accuracy tends to be worst in the cases with $a_{\text{true}} = 1$, except that sometimes when $n = 1$ it is worst when a_{true} is larger. The accuracy tends to level off as a_{true} grows larger than 100 or smaller than 0.01. Across all cases, the approximation is particularly good when a_{true} is 0.01 or smaller, or when n is large and $\mu = \mu_{\text{true}}$.

3.2 Speed of convergence

To evaluate the speed of Algorithm 1 in terms of the number of iterations required, Table 1 shows the number of runs in which the algorithm required k iterations to reach the termination condition, for each k . All runs terminated in four iterations or less. (For each a_0 ,

Table 1: Number of runs that required k iterations before terminating.

k	$a_0 = 1$	$a_0 = 0.1$	$a_0 = 0.01$
1 iteration	0	0	0
2 iterations	0	318	631
3 iterations	5751	4699	4308
4 iterations	1854	2588	2666
≥ 5 iterations	0	0	0

these results are over all of the 1521 combinations of n , r , a_{true} , μ_{true} , times five replicates each, as described above.) The algorithm terminates very rapidly in all cases considered.

4 Previous work

The original inspiration for Algorithm 1 was the generalized Newton algorithm of Minka (2002) for maximum likelihood estimation of gamma parameters. That said, our algorithm is significantly different since we are considering the full conditional distribution of the shape, rather than finding the unconditional MLE of the shape. Further, we iteratively set a to the mean rather than the mode of the approximating density. This is an important difference since the mode can be at 0 when $a_0 < 1$, but the derivatives blow up at 0, making it incoherent to match the derivatives at 0. Meanwhile, our algorithm handles $a_0 < 1$ without difficulties. Also, empirical evidence suggests that even when the mode is valid, using the mean tends to yield a more accurate approximation.

Damsleth (1975) found conjugate priors for the shape a in the case of known rate b , as well as the case of unknown b . However, it seems that the distributions in these conjugate families are not easy to work with, analytically. Consequently, Damsleth (1975) resorts to numerical integration over a , for each value of b , in order to obtain $p(b | x_{1:n}) = \int p(a, b | x_{1:n}) da$, the posterior density of b . Similarly, Miller (1980) numerically integrates over a to obtain the posterior of b (or of the mean μ), but instead of doing so for each value of b , he numerically computes the mean, variance, skewness, and kurtosis of $p(b | x_{1:n})$, and then constructs an approximation to $p(b | x_{1:n})$ by using the method of moments to fit a member

of the Pearson family to it. Miller (1980) uses improper priors as well as the conjugate priors of Damsleth (1975). These numerical integration approaches would be fine for a small number of shape parameters with fixed data $x_{1:n}$ and fixed hyperparameters a_0, b_0 . However, modern applications involve hierarchical models with many shape parameters, with latent x 's and unknown hyperparameters. In such cases, the numerical integration approaches of Damsleth (1975) and Miller (1980) would be very computationally intensive.

Son and Oh (2006) propose to use the adaptive rejection sampling (ARS) method of Gilks and Wild (1992) to sample from the full conditional for the shape a in a Gibbs sampler for a and b . There are a few disadvantages to this approach. ARS is quite a bit more complicated to implement than our procedure, and requires additional bookkeeping that is burdensome when dealing with a large number of shape parameters. Further, ARS requires the target distribution to be log-concave, but log-concavity of the full conditional of a is not always guaranteed. For instance, when using a $\text{Gamma}(a_0, b_0)$ prior on a , if $0 < a_0 < 1$ then the prior is not log-concave, which can cause the full conditional given $\mu, x_{1:n}$ to fail to be log-concave when a_0 is sufficiently small and $n = 1$. In practice, it is important to allow a_0 to be small in order to obtain a less informative prior while holding the prior mean a_0/b_0 fixed. Rejection sampling (RS) approaches have also been used by other authors, for instance, Pradhan and Kundu (2011) use the log-concave RS method of Devroye (1984), and Tsionas (2001) uses RS with a calibrated Gaussian proposal.

5 Derivation of the algorithm

Let $f(a) = p(a \mid x_{1:n}, \mu, a_0, b_0)$ denote the true full conditional density, and let $g(a) = \text{Gamma}(a \mid A, B)$ for some A, B to be determined. The basic idea of the algorithm is to choose A and B so that the first and second derivatives of $\log g$ match those of $\log f$ at a point a near the mean of f . To find a point near the mean of f , the algorithm is initialized with a choice of A and B based on an approximation of the gamma function (see the supplementary material), and then iteratively, A and B are refined by matching the first and second derivatives of $\log f$ at the mean A/B from the previous iteration.

It might seem like a better idea to use a point near the mode of f rather than the mean, by analogy with the Laplace approximation. However, on this problem, empirical

evidence suggests that using the mean tends to provide better accuracy than the mode. Using the mean also simplifies the algorithm since the mode of g is sometimes at 0, but the derivatives of $\log f$ and $\log g$ blow up at 0. To interpret the limiting behavior of the algorithm, we characterize its fixed points in the supplementary material.

5.1 Matching derivatives

First, to find the formulas for matching the first and second derivatives of $\log g$ to those of $\log f$ at some point $a > 0$, note that

$$\begin{aligned}\log g(a) &= A \log(B) - \log \Gamma(A) + (A - 1) \log(a) - Ba \\ \frac{\partial}{\partial a} \log g(a) &= \frac{A - 1}{a} - B \\ \frac{\partial^2}{\partial a^2} \log g(a) &= -\frac{A - 1}{a^2}.\end{aligned}$$

Meanwhile, $f(a) \propto p(x_{1:n} | a, \mu, a_0, b_0) p(a | \mu, a_0, b_0)$ and

$$\begin{aligned}p(x_{1:n} | a, \mu, a_0, b_0) &= \prod_{i=1}^n \frac{(a/\mu)^a}{\Gamma(a)} x_i^{a-1} \exp(- (a/\mu)x_i) \\ &= \frac{(a/\mu)^{na}}{\Gamma(a)^n} \exp((a - 1)R) \exp(- (a/\mu)S) \\ &= \frac{a^{na}}{\Gamma(a)^n} \exp(- R - (T + n)a)\end{aligned}\tag{5.1}$$

where $R = \sum_{i=1}^n \log(x_i)$, $S = \sum_{i=1}^n x_i$, and $T = S/\mu - R + n \log(\mu) - n$. Thus,

$$\begin{aligned}\log f(a) &= \text{const} + \log p(x_{1:n} | a, \mu, a_0, b_0) + \log p(a | \mu, a_0, b_0) \\ &= \text{const} + na \log(a) - n \log \Gamma(a) - (T + n)a + (a_0 - 1) \log(a) - b_0 a \\ \frac{\partial}{\partial a} \log f(a) &= n \log(a) + n - n\psi(a) - (T + n) + \frac{a_0 - 1}{a} - b_0 \\ \frac{\partial^2}{\partial a^2} \log f(a) &= n/a - n\psi'(a) - \frac{a_0 - 1}{a^2}\end{aligned}\tag{5.2}$$

since $\psi(a) = \frac{\partial}{\partial a} \log \Gamma(a)$ and $\psi'(a) = \frac{\partial^2}{\partial a^2} \log \Gamma(a)$. Setting $\frac{\partial^2}{\partial a^2} \log g(a) = \frac{\partial^2}{\partial a^2} \log f(a)$ yields

$$A = a_0 - na + na^2\psi'(a),\tag{5.3}$$

and setting $\frac{\partial}{\partial a} \log g(a) = \frac{\partial}{\partial a} \log f(a)$ yields

$$B = b_0 + \frac{A - a_0}{a} - n \log(a) + n\psi(a) + T.\tag{5.4}$$

The updates to A and B in Algorithm 1 are defined by Equations 5.3 and 5.4.

5.2 Updating the matching point a

Given A and B such that $\text{Gamma}(A, B)$ approximates f , we can choose a to approximate the mean of f by setting it to the mean of $\text{Gamma}(A, B)$, namely, $a = A/B$. By alternating between updating $a \leftarrow A/B$, and updating A, B via Equations 5.3 and 5.4, we iteratively refine the choice of a . This yields the loop in Algorithm 1. Convergence is assessed by checking if the relative change in a is smaller than the tolerance ϵ , that is, if $|a_{\text{old}} - a_{\text{new}}| < \epsilon a_{\text{new}}$, or equivalently, $|a_{\text{old}}/a_{\text{new}} - 1| < \epsilon$.

SUPPLEMENTARY MATERIAL

Gshape.jl: Source code implementing the algorithm described in the article, as well as the code used to generate the figures in the article. (file type: Julia source code)

Supplement.pdf: Additional analysis and empirical results. (file type: PDF)

References

- E. Damsleth. Conjugate classes for gamma distributions. *Scandinavian Journal of Statistics*, pages 80–84, 1975.
- L. Devroye. A simple algorithm for generating random variates with a log-concave density. *Computing*, 33(3-4):247–257, 1984.
- W. R. Gilks and P. Wild. Adaptive rejection sampling for Gibbs sampling. *Applied Statistics*, 41(2):337–348, 1992.
- R. B. Miller. Bayesian analysis of the two-parameter gamma distribution. *Technometrics*, 22(1):65–69, 1980.
- T. P. Minka. Estimating a gamma distribution. Technical report, Microsoft Research, 2002.
- B. Pradhan and D. Kundu. Bayes estimation and prediction of the two-parameter gamma distribution. *Journal of Statistical Computation and Simulation*, 81(9):1187–1198, 2011.
- Y. S. Son and M. Oh. Bayesian estimation of the two-parameter Gamma distribution. *Communications in Statistics—Simulation and Computation*, 35(2):285–293, 2006.
- E. G. Tsionas. Exact inference in four-parameter generalized gamma distributions. *Communications in Statistics—Theory and Methods*, 30(4):747–756, 2001.

Supplementary material for “Fast and Accurate Approximation of the Full Conditional for Gamma Shape Parameters”

Jeffrey W. Miller

S1 Fixed points of the algorithm

In this section, we analyze the limiting behavior of Algorithm 1 by characterizing its fixed points. This provides a direct interpretation of the final point a at which we match the first and second derivatives of $\log g$ to those of $\log f$ (in the notation of Section 5) in the last iteration of the algorithm.

In Theorem S1.1, we show that $a > 0$ is a fixed point of Algorithm 1 if and only if $\frac{\partial}{\partial a} \log f(a) + 1/a = 0$, and further, the algorithm is guaranteed to have a fixed point. It might seem more natural to instead match derivatives at a point where $\frac{\partial}{\partial a} \log f(a) = 0$, however, empirical evidence indicates that this tends to yield a worse approximation.

To interpret this result, note that if $f(a)$ were exactly equal to a gamma density, then the mean of f would be the unique point a such that $\frac{\partial}{\partial a} \log f(a) + 1/a = 0$. There is also a natural interpretation of $\frac{\partial}{\partial a} \log f(a) + 1/a = 0$ in terms of optimization with a logarithmic barrier. Barrier functions are frequently used in optimization to deal with inequality constraints. Adding the logarithmic barrier function $\log a$ to the objective function $\log f(a)$ allows one to smoothly enforce the constraint that $a > 0$. This leads to the penalized objective function $\log f(a) + \log a$, and it turns out that the fixed points of Algorithm 1 coincide with the points where the derivative of this penalized objective equals 0.

Theorem S1.1. *Let $x_1, \dots, x_n > 0$, where $n \geq 0$, and let $\mu, a_0, b_0 > 0$. For any $a > 0$, $\frac{\partial}{\partial a} \log f(a) + 1/a = 0$ if and only if a is a fixed point of Algorithm 1. Further, there exists $a > 0$ such that $\frac{\partial}{\partial a} \log f(a) + 1/a = 0$.*

Proof. First, we show that $\frac{\partial}{\partial a} \log f(a) + 1/a = 0$ characterizes the fixed points. Define $d_1(a) = \frac{\partial}{\partial a} \log f(a)$ and $d_2(a) = \frac{\partial^2}{\partial a^2} \log f(a)$. Recall from Section 5 that at each iteration,

the algorithm updates a by setting $a \leftarrow A/B$ where A and B are defined such that $d_1(a) = \frac{\partial}{\partial a} \log g(a) = (A-1)/a - B$ and $d_2(a) = \frac{\partial^2}{\partial a^2} \log g(a) = -(A-1)/a^2$ for the current value of a . Solving for A and B implies that $A = -a^2 d_2(a) + 1$ and $B = -a d_2(a) - d_1(a)$. Thus, a is a fixed point of the algorithm precisely when $a = A/B = (a^2 d_2(a) - 1)/(a d_2(a) + d_1(a))$, or equivalently, by rearranging and canceling terms, when $d_1(a) + 1/a = 0$.

To justify rearranging and cancelling, note that $d_1(a)$ and $d_2(a)$ are finite for all $a \in (0, \infty)$; see Equation 5.2. If a is a fixed point then $a d_2(a) + d_1(a) \neq 0$ since otherwise A/B is infinite or undefined. On the other hand, if $d_1(a) + 1/a = 0$ then $a d_2(a) + d_1(a) \neq 0$, since otherwise $1/a^2 = d_2(a) = n/a - n\psi'(a) - \frac{a_0-1}{a^2}$, which would imply that $n(\psi'(a) - 1/a) = -a_0/a^2 < 0$, but this is a contradiction since $\psi'(a) - 1/a > 0$ for all $a > 0$ (Qi and Guo, 2009, Eqn. 13).

Now, we show that there always exists $a > 0$ such that $d_1(a) + 1/a = 0$. By Equation 5.2,

$$\begin{aligned} d_1(a) + 1/a &= n \log(a) + n - n\psi(a) - (T + n) + a_0/a - b_0 \\ &= n(\log(a) - \psi(a)) + a_0/a - b_0 - T. \end{aligned}$$

Below, we show that $T \geq 0$. Assuming $T \geq 0$ for the moment, it follows from the Intermediate Value Theorem that $d_1(a) + 1/a = 0$ for some $a > 0$, since $\log(a) - \psi(a) \rightarrow \infty$ as $a \rightarrow 0$ and $\log(a) - \psi(a) \rightarrow 0$ as $a \rightarrow \infty$ (Alzer, 1997, Eqn. 2.2).

To verify that $T \geq 0$, recall from Equation 5.2 that $T = S/\mu - R + n \log(\mu) - n$, where $S = \sum_{i=1}^n x_i$ and $R = \sum_{i=1}^n \log(x_i)$. By rearranging terms, $T = \sum_{i=1}^n (x_i/\mu - \log(x_i/\mu) - 1)$, and by straightforward calculus, one can check that $x - \log(x) - 1 \geq 0$ for all $x > 0$. □

S2 Initializing via approximation of $\Gamma(a)$

To find good starting values of A and B , we consider approximations to the gamma function $\Gamma(a)$ in two regimes: (i) when a is large, and (ii) when a is small.

By Stirling's approximation to $\Gamma(a)$, we have that $\Gamma(a)/a^a \sim \sqrt{2\pi} a^{-1/2} e^{-a}$ as $a \rightarrow \infty$.

(Here, $h_1(a) \sim h_2(a)$ as $a \rightarrow a^*$ means that $h_1(a)/h_2(a) \rightarrow 1$ as $a \rightarrow a^*$.) Plugging this into Equation 5.1, we get

$$p(x_{1:n} \mid a, \mu, a_0, b_0) \approx (2\pi)^{-n/2} a^{n/2} e^{na} \exp(-R - (T+n)a) \propto a^{n/2} \exp(-Ta),$$

and therefore, $f(a)$ is approximately proportional to $\text{Gamma}(a \mid a_0 + n/2, b_0 + T)$ when a is large. This suggests initializing with $A = a_0 + n/2$ and $B = b_0 + T$, as in Algorithm 1.

As $a \rightarrow 0$, we have $\Gamma(a)/a^a \sim a^{-1}$, since $a^a \rightarrow 1$ and $\Gamma(a)/a^{-1} = a\Gamma(a) = \Gamma(a+1) \rightarrow \Gamma(1) = 1$ as $a \rightarrow 0$. Plugging this into Equation 5.1, we see that

$$p(x_{1:n} \mid a, \mu, a_0, b_0) \approx a^n \exp(-R - (T+n)a) \propto a^n \exp(-(T+n)a),$$

and therefore, $f(a)$ is approximately proportional to $\text{Gamma}(a \mid a_0 + n, b_0 + T + n)$ when a is small. This suggests initializing with $A = a_0 + n$ and $B = b_0 + T + n$. While similar to the large a version, empirically we find that this initialization tends to converge slightly slower, by 1 iteration or so.

S3 Numerical integration for truth comparison

Since the true full conditional distribution is not analytically tractable, it is necessary to approximate the integrals for the total variation distance and the Kullback–Leibler divergences between the true and approximate distributions. To do this, we use a simple adaptive numerical integration technique using an importance distribution to choose points (basically, a deterministic importance sampler).

Specifically, suppose we want to approximate the integral $\int h(a)g(a)da$ where $g(a)$ is a probability density on \mathbb{R} , and $h(a)$ is some function on \mathbb{R} . Let $G(a)$ be the CDF of g , and let $a_i = G^{-1}(u_i)$ where $u_i = (i - 0.5)/N$ for $i = 1, \dots, N$. (If the inverse of G does not exist, use the generalized inverse distribution function.) Then $\int h(a)g(a)da \approx \frac{1}{N} \sum_{i=1}^N h(a_i)$ when N is large. This is similar to Monte Carlo, but is much more accurate due to using equally-spaced points u_i rather than random draws $u_i \sim \text{Uniform}(0, 1)$.

We apply this approximation with $g(a) = \text{Gamma}(a|A, B)$, where $A, B > 0$ are the output from Algorithm 1. Let $f(a) = p(a | x_{1:n}, \mu, a_0, b_0)$ be the true full conditional. Since the normalizing constant of f is unknown, write $f(a) = \tilde{f}(a)/Z$ where $\tilde{f}(a)$ can be easily computed, and observe that

$$Z = \int_{-\infty}^{\infty} \tilde{f}(a) da = \int_0^{\infty} \frac{\tilde{f}(a)}{g(a)} g(a) da = \int_{-\infty}^{\infty} w(a) g(a) da$$

where $w(a) = \tilde{f}(a)/g(a)$ for $a > 0$ and $w(a) = 0$ for $a \leq 0$, since $g(a) > 0$ when $a > 0$, and $f(a) = g(a) = 0$ when $a \leq 0$. Thus, $Z \approx \frac{1}{N} \sum_{i=1}^N w(a_i)$.

Define $\hat{f}(a) = \tilde{f}(a)/(\frac{1}{N} \sum_{i=1}^N w(a_i))$, so that $\hat{f}(a) \approx f(a)$. The integral for total variation distance can be written in the form $\int h(a)g(a)da$ as follows:

$$d_{\text{TV}}(f, g) = \frac{1}{2} \int_{-\infty}^{\infty} |f(a) - g(a)| da = \frac{1}{2} \int_0^{\infty} \left| \frac{f(a)}{g(a)} - 1 \right| g(a) da = \int_{-\infty}^{\infty} h(a) g(a) da$$

where $h(a) = \frac{1}{2} |f(a)/g(a) - 1|$ for $a > 0$ and $h(a) = 0$ for $a \leq 0$. Thus, since $a_i > 0$,

$$d_{\text{TV}}(f, g) \approx \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \left| \frac{\hat{f}(a_i)}{g(a_i)} - 1 \right|.$$

Similarly, $d_{\text{KL}}(f, g) \approx \frac{1}{N} \sum_{i=1}^N \frac{\hat{f}(a_i)}{g(a_i)} \log \frac{\hat{f}(a_i)}{g(a_i)}$ and $d_{\text{KL}}(g, f) \approx \frac{1}{N} \sum_{i=1}^N \log \frac{g(a_i)}{\hat{f}(a_i)}$.

S4 Additional empirical results

Figures S1, S2, and S3 show the total variation distance for all the cases with $a_0 = 1$, $a_0 = 0.1$, and $a_0 = 0.01$, respectively. Figures S4 and S5 show the Kullback–Leibler divergences for all the cases with $a_0 = 1$.

References

- H. Alzer. On some inequalities for the gamma and psi functions. *Mathematics of Computation of the American Mathematical Society*, 66(217):373–389, 1997.
- F. Qi and B.-N. Guo. Refinements of lower bounds for polygamma functions. *arXiv:0903.1996*, 2009.

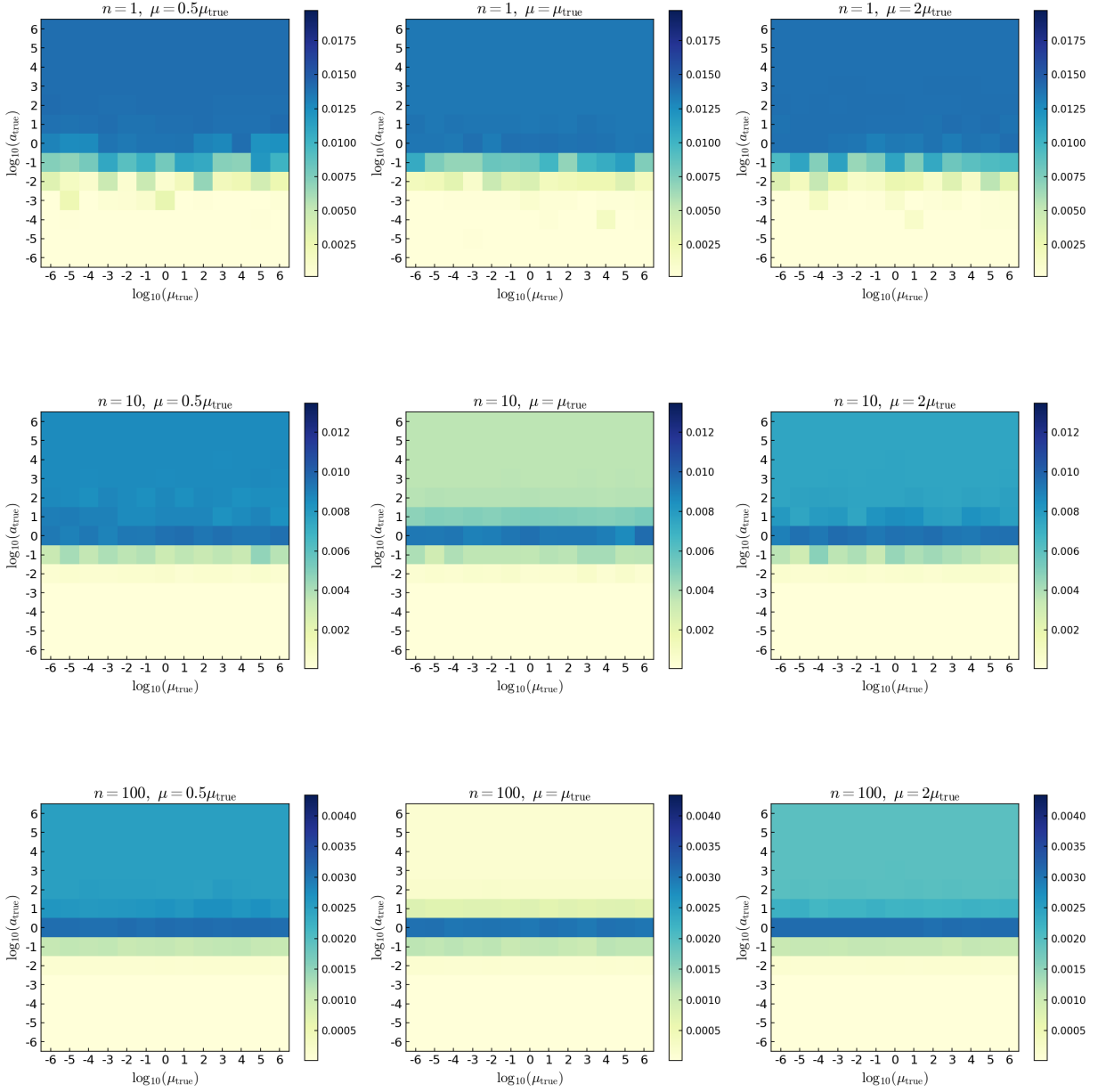


Figure S1: Total variation distance $d_{TV}(f, g)$ between the true full conditional f and the approximate full conditional g , for every case with $a_0 = 1$. The values shown are the averages over the five independent data sets for each case. Note the scale at the right of each plot.

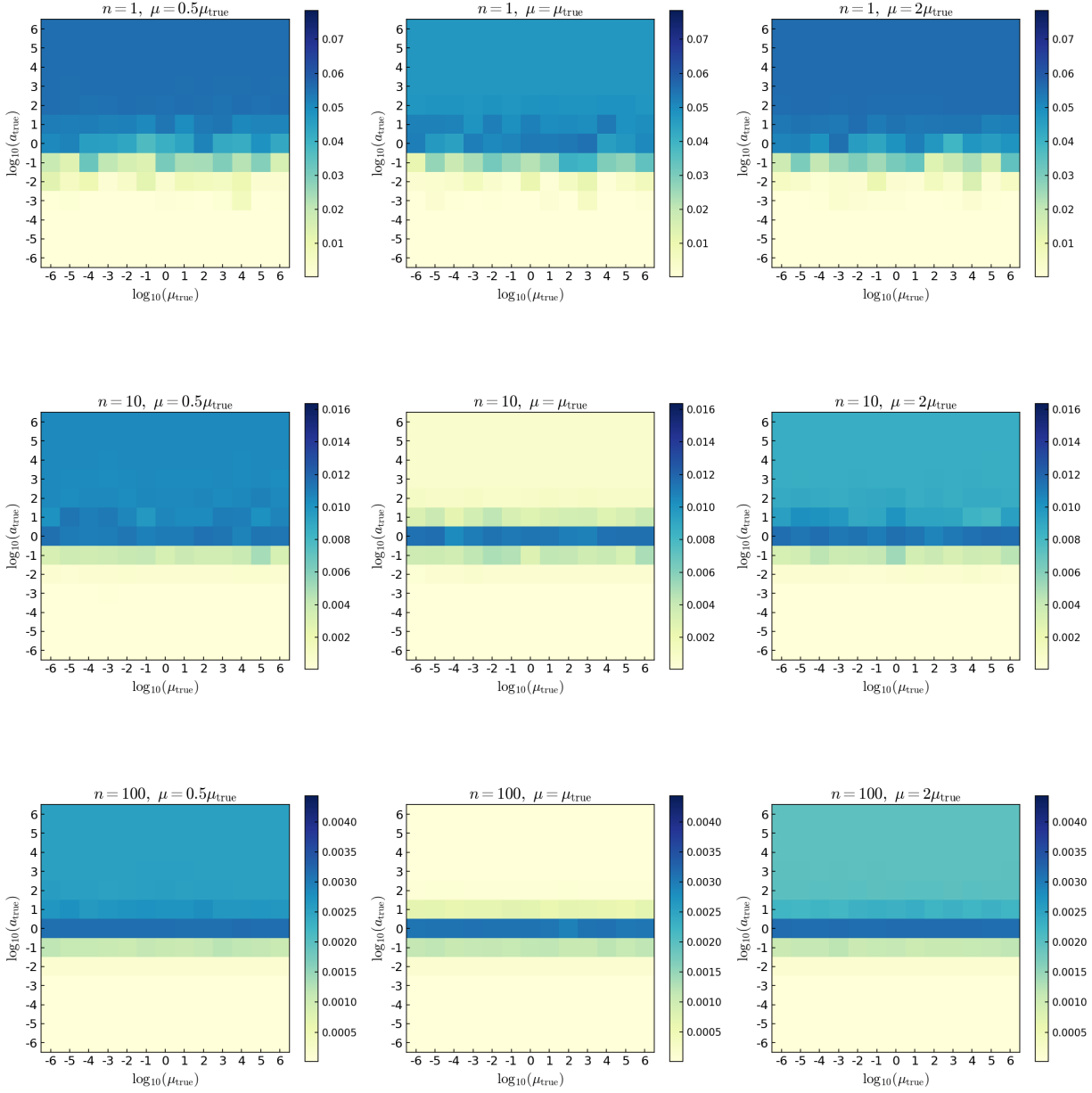


Figure S2: Total variation distance $d_{TV}(f, g)$ between the true full conditional f and the approximate full conditional g , for every case with $a_0 = 0.1$. The values shown are the averages over the five independent data sets for each case. Note the scale at the right of each plot.

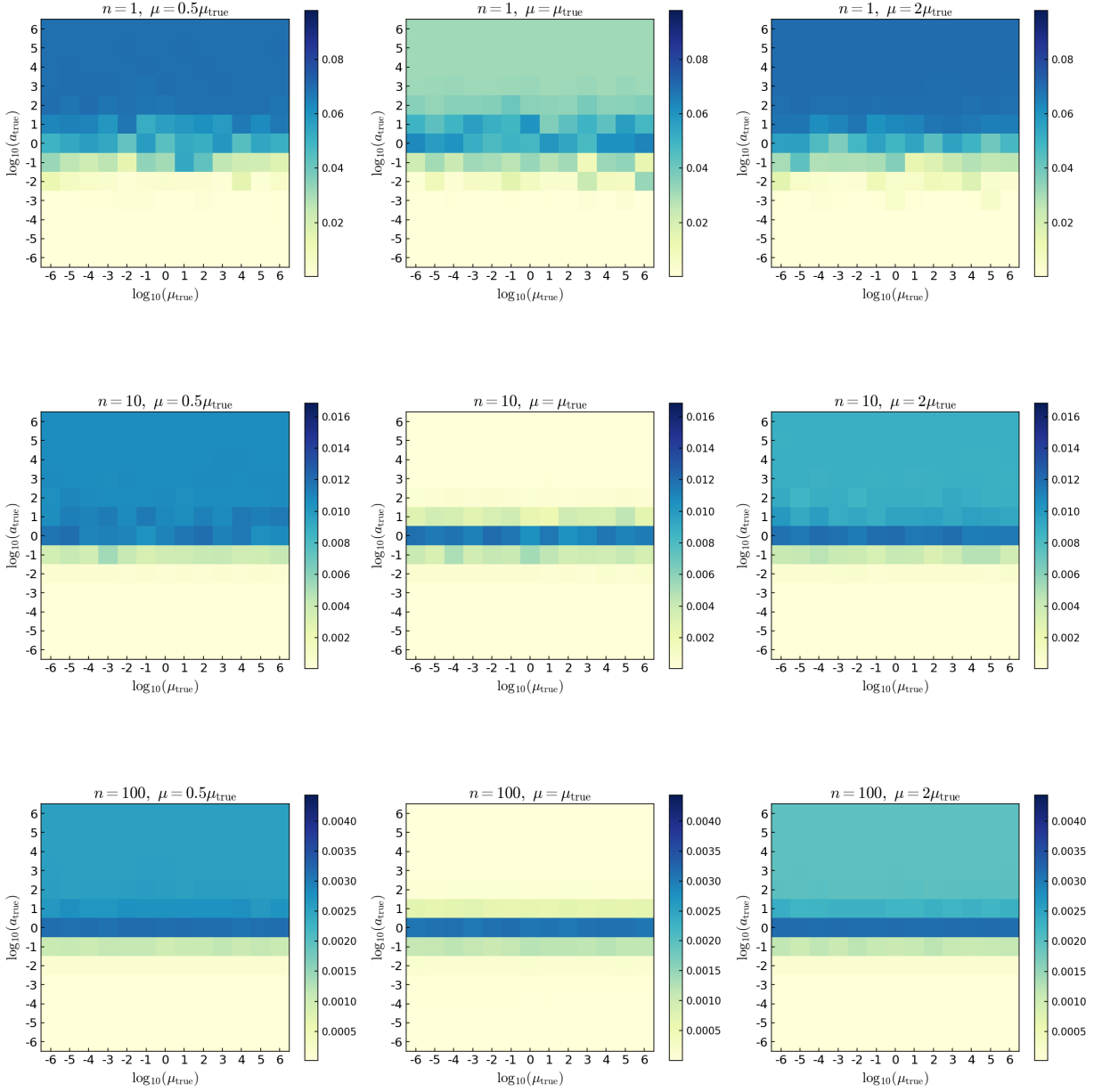


Figure S3: Total variation distance $d_{TV}(f, g)$ between the true full conditional f and the approximate full conditional g , for every case with $a_0 = 0.01$. The values shown are the averages over the five independent data sets for each case. Note the scale at the right of each plot.

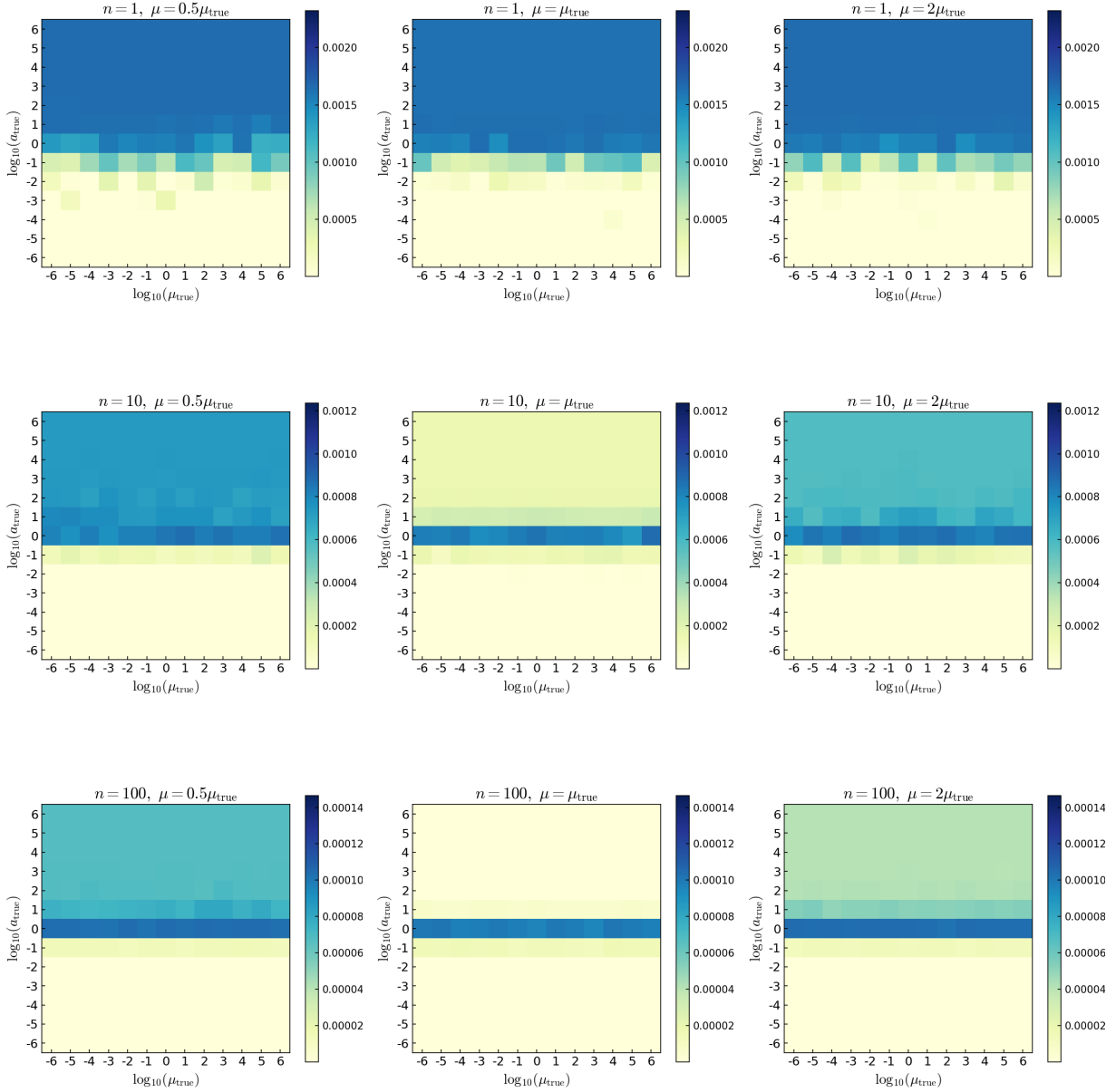


Figure S4: Kullback–Leibler divergence $d_{KL}(f, g)$ between the true full conditional f and the approximate full conditional g , for every case with $a_0 = 1$. The values shown are the averages over the five independent data sets for each case. Note the scale at the right of each plot.

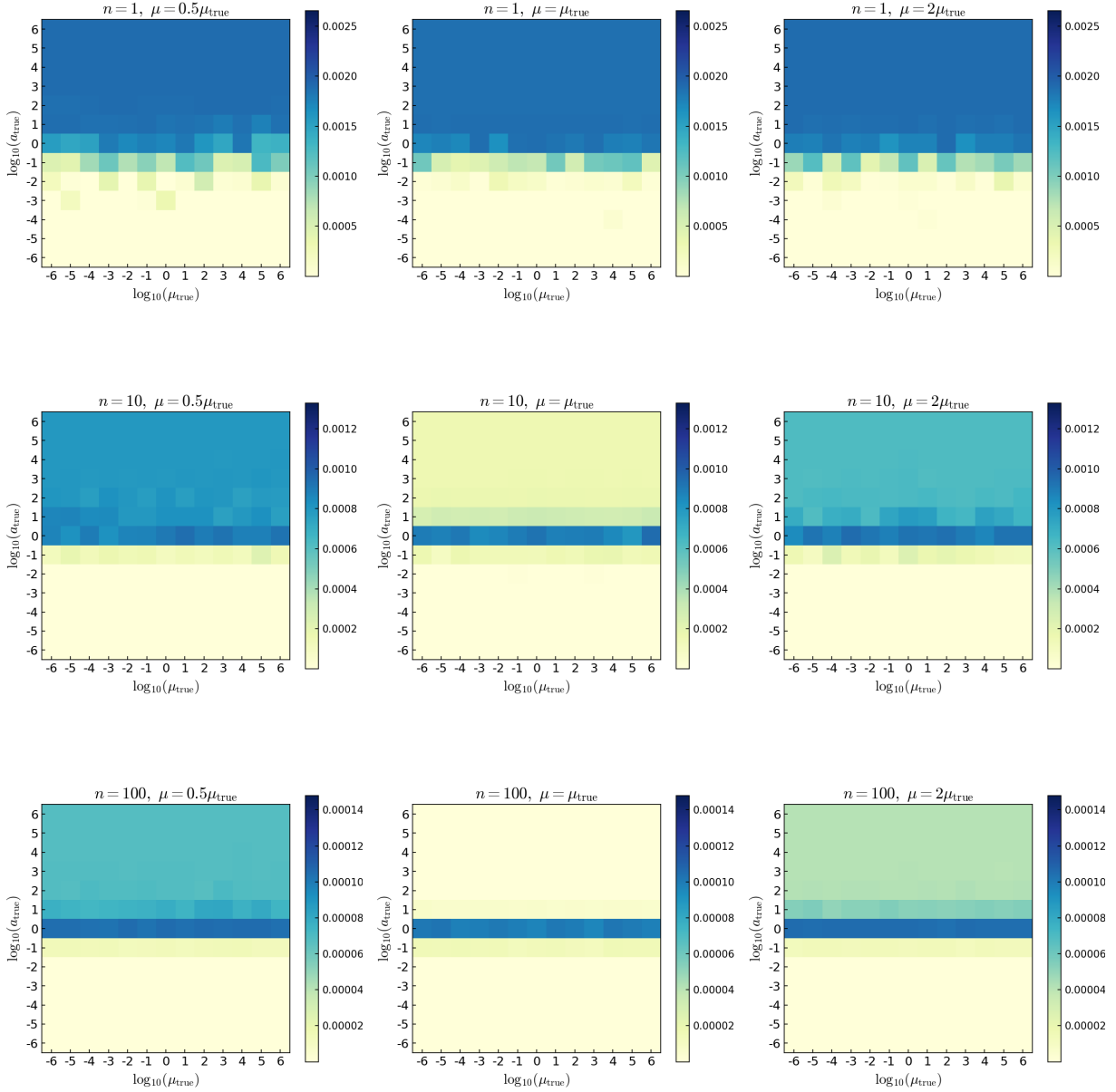


Figure S5: Kullback–Leibler divergence $d_{KL}(g, f)$ between the true full conditional f and the approximate full conditional g , for every case with $a_0 = 1$. The values shown are the averages over the five independent data sets for each case. Note the scale at the right of each plot.