

Author Accepted Manuscript

**Revisiting Scalable Targeted Marketing with Distributed Markov Chain Monte Carlo**

Journal:	<i>Journal of Marketing Research</i>
Manuscript ID	JMR-22-0154.R5
Manuscript Type:	Revised Submission
Topics and Methods:	Bayesian estimation < Theoretical Foundation, Measurement and inference < Theoretical Foundation, Target marketing, Big data

SCHOLARONE™
Manuscripts

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

Revisiting Scalable Targeted Marketing with
Distributed Markov Chain Monte Carlo

Michael Braun
Cox School of Business
Southern Methodist University
braunm@smu.edu

April 29, 2024

24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

Abstract

Bumbaca, Misra, and Rossi (2020) propose a parallelizable algorithm for estimating a large number of customer-level parameters in a Bayesian hierarchical model. However, the algorithm follows from a mathematical error in the derivation of the target posterior density, which calls into question the theoretical support for the algorithm sampling from the specified model. Adapting the algorithm to be consistent with the corrected math nullifies the claimed benefits in scalability and efficiency. Notwithstanding that error, unbiasedness requires the number of customers to be asymptotic per computational node, which is more restrictive than being asymptotic in the size of the dataset as a whole. The more the algorithm is parallelized, the greater the bias. Potential adopters should be aware that the algorithm does not sample from the exact posterior distribution, and that its ability to take advantage of distributed computing infrastructure is limited.

40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Keywords: parallel Bayesian estimation, Bayesian hierarchical models, target marketing, big data, hierarchical models

Author Accepted Manuscript

1
2
3 Tailoring marketing strategies to specific consumers can often improve their effectiveness (e.g.,
4
5
6 Ascarza 2018; Danaher 2023). Doing so requires estimating each customer's unobserved ten-
7
8 dencies and propensities to respond to marketing interventions, which can be methodologically
9
10 difficult. One approach is to conduct inference on customer-level parameters in a hierarchical
11
12 Bayesian model (Rossi et al. 1993). But standard Bayesian estimation algorithms like Markov
13
14 chain Monte Carlo (MCMC) become computationally intractable when the number of customers
15
16 in a database is very large.

17
18
19
20 Some researchers have proposed scalable MCMC variants that employ parallel computing, where
21
22 the model is estimated from separate “shards” (partitions) of the dataset on distributed computa-
23
24 tional nodes. MCMC draws are then post-processed to allow for approximate posterior inference
25
26 (e.g., Neiswanger et al. 2014; Scott et al. 2016; Vyner et al. 2023). The post-processing stage is
27
28 necessary because combining samples generated from separate nodes is not equivalent to sam-
29
30 pling from a posterior that conditions on the aggregate database. Potential adopters of these
31
32 algorithms must consider the tradeoffs among computational efficiency, the scope of applications
33
34 for which the various methods are valid, and the magnitude of the biases that these methods may
35
36 introduce in practice.
37
38
39
40

41
42
43 Bumbaca, Misra, and Rossi (2020, henceforth BMR) introduce an algorithm in that same class of
44
45 “divide and conquer” parallel MCMC methods, with a specific focus on estimating customer-level
46
47 parameters. However, their paper contains a material mathematical error where they treat a pro-
48
49 portional relationship as an equality, and make an impermissible substitution in their derivation of
50
51 the target posterior density. This calls into question BMR's claim of asymptotic exactness. This
52
53 paper identifies and corrects the error, and shows that the published version of the algorithm is
54
55
56
57
58
59
60

inconsistent with the corrected math. Adapting the algorithm to be consistent with the corrected math reduces its scalability and efficiency.

Furthermore, the absence of a post-processing step means that the distributed nature of the BMR algorithm induces “parallelization bias.” In practice, this is a problem because the algorithm’s asymptotic unbiasedness assumes a sufficiently large number of customers *per shard*, which is a stronger (and less typical) assumption than a large number of customers in the aggregate dataset. Because real-world databases contain data for a finite number of customers, the more a practitioner tries to parallelize the algorithm (which is the primary motivation of BMR’s paper), the more biased the results will be. While the BMR algorithm could produce a sample that may be close to the exact posterior, potential adopters of the BMR algorithm should be aware of these theoretical inconsistencies and practical limitations.

THE BMR ALGORITHM AND THE MATHEMATICAL ERROR

The objective of BMR’s algorithm is to estimate customer-level parameters β_i in a Bayesian hierarchical model when N , the number of customers in the dataset, is large. It does this by partitioning the dataset Y into S shards of size $N_s = N/S$ (denoted as $Y_{1:S} = \{Y_1, \dots, Y_S\}$) and parceling computation across S computational nodes. In Stage 1 of the algorithm, R samples of a population-level parameter θ are generated from each of $p(\theta | Y_1, \tau), \dots, p(\theta | Y_S, \tau)$ in parallel. These samples are gathered from the various nodes and combined into $\{\theta_s^r\}$, a collection of $R \times S$ draws. Identical copies of $\{\theta_s^r\}$ are redistributed to the nodes in preparation for Stage 2, which involves iteratively sampling β_i from $p(\beta_i | \theta_s^r) \propto p(y_i | \beta_i)p(\beta_i | \theta_s^r)$. BMR claim that samples from $p(\beta_i | \theta_s^r)$ constitute an asymptotically unbiased estimate of the target posterior

Author Accepted Manuscript

distribution $p(\beta_i | Y, \tau)$. Web Appendix A summarizes the BMR algorithm and the hierarchical model it is designed to estimate.

BMR's error is in their derivation of $p(\beta_i | Y, \tau)$. To begin, consider their Eqs. 5 and 7:

$$(BMR-5) \quad p(\beta_i | y_i, \theta) \propto p(\beta_i | \theta)p(y_i | \beta_i)$$

$$(BMR-7) \quad p(\beta_i | Y, \tau) = \int p(\beta_i | y_i, \theta)p(\theta | Y, \tau)d\theta.$$

BMR substitute the $p(\beta_i | y_i, \theta)$ term in the integrand of BMR-7 with BMR-5, resulting in

$$(BMR-8) \quad p(\beta_i | Y, \tau) \propto \int p(\beta_i | \theta)p(y_i | \beta_i)p(\theta | Y, \tau)d\theta$$

$$(BMR-13) \quad = \frac{p_{\theta|Y,\tau}(\beta_i)p(y_i | \beta_i)}{p(y_i)},$$

where $p_{\theta|Y,\tau}(\beta_i) = \int p(\beta_i | \theta)p(\theta | Y, \tau)d\theta$. But this step is incorrect because BMR-5 is a proportional relationship (\propto), not an equality ($=$). Instead, by Bayes' Theorem,

$$(1) \quad p(\beta_i | y_i, \theta) = \frac{p(y_i | \beta_i, \theta)p(\beta_i | \theta)}{p(y_i | \theta)}.$$

The $p(y_i | \theta)$ term in the denominator of Eq. 1 is consequential because it depends on θ , the variable of integration in BMR-7, and cannot be factored out of the integrand. Web Appendix B shows that substituting Eq. 1 into BMR-7 corrects the target posterior density in BMR-13.

$$(2) \quad p(\beta_i | Y, \tau) = \frac{p_{\theta|Y_{-i},\tau}(\beta_i)p(y_i | \beta_i)}{p(y_i | Y_{-i})}.$$

In Eq. 2, Y_{-i} is data for all customers excluding i , and $p_{\theta|Y_{-i},\tau}(\beta_i) = \int p(\beta_i | \theta)p(\theta | Y_{-i}, \tau)d\theta$.

Author Accepted Manuscript

As published, the BMR algorithm is inconsistent with the corrected math, and appears to be scalable only because it samples from the incorrect BMR-13. Using all S nodes to generate a single $\{\theta_s^r\}$ in Stage 1 was justified only because of how BMR derived $p_{\theta|Y,\tau}(\beta_i)$: integrating $p(\beta_i | \theta)$ over the same $p(\theta | Y, \tau)$ for all customers. BMR's reports of computational efficiency depend on all customers' draws from $p(\beta_i | \theta_s^r)$ being conditional on elements of a common $\{\theta_s^r\}$.

But the $p_{\theta|Y_{-i},\tau}(\beta_i)$ term in the corrected Eq. 2 is an integral over $p(\theta | Y_{-i}, \tau)$, which is different for each customer. Fixing the algorithm to be consistent with the corrected math would involve conditioning Stage 2 samples on a distinct $\{\theta_s^r\}_i$ for each customer. The efficiency of running only S parallel MCMC instances in Stage 1 would be lost, especially since N/S must be sufficiently large for the algorithm to be asymptotically unbiased (see below). Also, Stage 2 would incur additional communication overhead because each customer's θ_s^r draws would come from a different $\{\theta_s^r\}_i$. Gains in computational efficiency would be much more modest. See Web Appendix C.

Web Appendix D explains how BMR algorithm systematically biases the estimate of $p(\beta_i | Y, \tau)$, relative to the target posterior density in the model. The distinction between BMR-13 and Eq. 2 is in whether the prior on β_i conditions on Y , which includes on y_i , or on Y_{-i} , which does not. The practical effect of this error may be small, but remains to be proven. What is true is that even after conditioning on θ_s^r , the BMR algorithm does not sample from the model that is specified in their paper. The corrected derivation calls into question the theoretical support for BMR's claims that their algorithm samples from the exact posterior.

Further, BMR's salient theoretical contribution is that their approximations to $p_{\theta|Y,\tau}(\beta_i)$ and

Author Accepted Manuscript

$p(\beta_i | Y, \tau)$ converge in distribution to their true densities (Theorems 5 and 6). The error means that many of their theorems need to be either restated or reproven. See Web Appendix E.

PARALLELIZATION BIAS

The key to BMR's faster computational speed comes from dividing the data into shards so that computation can be performed in parallel. BMR's claimed exactness also depends on each shard including a very large number of customers (approaching infinity). One concern for potential adopters is that the algorithm induces additional "parallelization bias" into estimates of $p(\beta_i | Y, \theta)$. This is likely of more concern to potential adopters than the mathematical error.

The problem comes from how in Stage 1, $\{\theta_s^r\}$ is not actually sampled from $p(\theta | Y, \tau)$, but rather from a mixture of posterior distributions that are conditional on each shard:

$$(3) \quad p^*(\theta | Y_{1:S}, \tau) = \frac{1}{S} \sum_{s=1}^S p(\theta | Y_s, \tau)$$

The algorithm is asymptotically unbiased as $p^*(\theta | Y_{1:S}, \tau)$ converges to $p(\theta | Y, \tau)$, which requires the number of customers *per shard* (N/S) to be sufficiently large (see BMR's Theorem 4.1.2). This is more restrictive than the typical asymptotic assumption on N ($N/S \rightarrow \infty$ implies $N \rightarrow \infty$ trivially), and N/S gets smaller the more the practitioner parallelizes the algorithm (increasing S). Intuitively, BMR achieve asymptotic unbiasedness because when each shard has an infinite number of customers, each summation term in Eq. 3 conditions on the same infinitely-sized dataset. Web Appendix F shows that for *finite* N and $S > 1$, $p^*(\theta | Y_{1:S}, \tau)$ has higher variance than $p(\theta | Y, \tau)$, which attenuates Bayesian shrinkage. In BMR-13, $p_{\theta|Y, \tau}(\beta_i)$ acts as an infor-

Author Accepted Manuscript

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

mative prior on β_i (as does $p_{\theta|Y_{-i},\tau}(\beta_i)$ in Eq. 2). Parallelization bias weakens that information, allowing y_i to pull the estimate of $p(\beta_i | Y, \tau)$ away from the target posterior in the model. That is, the estimates of $\beta_i | Y$ potentially overfit customer i 's data when the shards are small.

Because the size of each shard decreases as the number of computing nodes increases, this bias will increase as the practitioner takes greater advantage of parallel computing. BMR address this problem with a heuristic that computes S_{\max} , which is an upper bound on S that is a function of a *practitioner-specified* bias tolerance, ϵ_{\max}^2 . This approach has some limitations. First, by constraining the number of shards that can be used, using the BMR algorithm as specified fails to fully take advantage of distributed computing infrastructure. BMR's heuristic asks practitioners to process larger shards on fewer nodes than are available, and to leave the remaining nodes idle. And if S_{\max} is very small, the time to run Stage 1 may still be prohibitively large.

Second, it is not clear how a practitioner should select or interpret the bias tolerance ϵ_{\max}^2 , which is a maximum allowable squared difference between densities. BMR's Theorem 9 requires that N'/S' in the *pilot run* be large enough for asymptotics to hold, which is even more restrictive than N/S being large. It is not clear how confident researchers can be that S_{\max} bounds the bias to be within ϵ_{\max}^2 . A tolerance based on moments would be more useful.

Finally, researchers should consider whether their data are a good fit for the BMR algorithm. In BMR's simulation study, customers are generated with either 5, 15 or 45 observations. These are situations in which the customer's data will overwhelm $p_{\theta|Y,\tau}^*(\beta_i)$ anyway, so attenuated shrinkage from parallelization bias is less likely to come into play. The algorithm should predict well for those customers. In BMR's empirical application, the algorithm's predictive ability is tested only among customers with at least 16 observations. Yet, BMR's Fig. 4 shows many

Author Accepted Manuscript

instances where the BMR posterior sample is quite different from the Gibbs sample (“a source of truth”). If the database includes customers with only a few observations, for whom the posterior depends more on the prior, or if both the data and prior on individual level parameters are highly informative, parallelization bias is likely to have a more detrimental effect on the results.

DISCUSSION

In their “What is Novel” section, BMR write: “The shard-splitting idea in each of the two stages is not new and common to both Scott et al. (2016) and Neiswanger et al. (2014). The novelties of the first stage are (1) constructing the posterior predictive density and (2) drawing from the posterior predictive density in parallel to reduce communication overhead between stages” (Bumbaca et al. 2020, p.1004). *These are the two aspects of their paper that are directly affected by their error.* Even if a practitioner were to look past that error, estimating population-level parameters conditional on shards of data induces bias in finite samples that practitioners cannot ignore, particularly when making maximal use of available parallel computing resources. While BMR do not promise unbiased distributions of population-level parameters as a “deliverable,” there is no avoiding the fact that Stage 1 generates biased population-level samples, and nothing in Stage 2 adjusts for that.

Despite the concerns about the BMR algorithm that are described in this paper, some practitioners may still find the trade-offs between real-time efficiency gains and biased estimation to be favorable. Potential adopters should consider whether these biases are likely to be small enough for their own practical applications.

REFERENCES

- Ascarza, Eva (2018). “Retention Futility: Targeting High-Risk Customers Might Be Ineffective.” *Journal of Marketing Research*, 60(1):80–98.
- Bumbaca, Federico, Sanjog Misra, and Peter E. Rossi (2020). “Scalable Target Marketing: Distributed Markov Chain Monte Carlo for Bayesian Hierarchical Models.” *Journal of Marketing Research*, 57(6):999–1018.
- Danaher, Peter J. (2023). “Optimal Microtargeting of Advertising.” *Journal of Marketing Research*, 60(3):564–584.
- Neiswanger, Willie, Chong Wang, and Eric Xing (2014). “Asymptotically Exact, Embarrassingly Parallel MCMC.” *Proceedings of the 30th International Conference on Uncertainty in Artificial Intelligence*. Quebec City. ARXIV:1311.4780v2.
- Rossi, Peter E. and Greg M. Allenby (1993). “A Bayesian Approach to Estimating Household Parameters.” *Journal of Marketing Research*, 30(2):171–182.
- Scott, Steven L., Alexander W. Blocker, Fernando V. Bonassi, Hugh A. Chipman, Edward I. George, and Robert E. McCulloch (2016). “Bayes and Big Data: The Consensus Monte Carlo Algorithm.” *International Journal of Management Science and Engineering Management*, 11(2):78–88.
- Vyner, Callum, Christopher Nemeth, and Chris Sherlock (2023). “SwISS: A Scalable Markov chain Monte Carlo Divide-and-Conquer Strategy.” *Stat*, 12(1):1–11.

Author Accepted Manuscript

Web Appendices

Revisiting Scalable Targeted Marketing with Distributed Markov Chain Monte Carlo

Michael Braun

Cox School of Business

Southern Methodist University

braunm@smu.edu

Contents

A	Context, background, and objectives	2
B	Deriving the Correct Posterior	4
C	Parallelizing a corrected algorithm	8
D	The BMR algorithm overfits customer data	9
E	Notes on BMR's Theoretical Proofs	11
F	Parallelization Bias	14

These materials have been supplied by the authors to aid in the understanding of their paper.

The AMA is sharing these materials at the request of the authors.

WEB APPENDIX A: CONTEXT, BACKGROUND, AND OBJECTIVES

To set the scene, Fig. A1 illustrates a hierarchical model with conditional independence across N customers (BMR's Eqs. 1 to 3). $Y = \{y_1, \dots, y_N\}$ represents the entire dataset, where each y_i is a (possibly multivariate) observed outcome for customer i . The customer-level latent parameters $\beta = \{\beta_1, \dots, \beta_N\}$ are associated only through their common dependence on a population parameter θ , which affects customer data y_i only through its corresponding β_i . $p(y_i | \beta_i)$ is the likelihood of each customer's data, $p(\beta_i | \theta)$ is the customer-level prior, and $p(\theta | \tau)$ is the population-level hyperprior.

BMR's managerial task involves sampling from the marginal posterior density $p(\beta_i | Y)$ for individual customers. BMR's proposed alternative is to sample customer-level parameters from an approximation to

$$(BMR-13) \quad p(\beta_i | Y) = \frac{p_{\theta|Y,\tau}(\beta_i)p(y_i | \beta_i)}{p(y_i)},$$

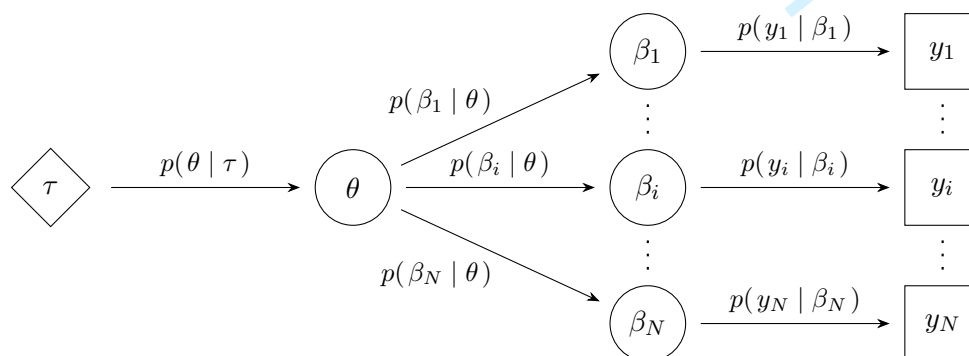
where

$$(A1) \quad p_{\theta|Y,\tau}(\beta_i) = \int p(\beta_i | \theta)p(\theta | Y, \tau)d\theta$$

Here, $p_{\theta|Y,\tau}(\beta_i)$ acts as an informative prior distribution on β_i that conditions on the entire dataset by marginalizing $p(\beta_i | \theta)$ over $p(\theta | Y, \tau)$.

To provide some intuition, $p(\beta_i | \theta)$ is not only a prior on a single customer's β_i , but it also describes the heterogeneous mixture of β_i across the population. By averaging $p(\beta_i | \theta)$ over $p(\theta | Y, \tau)$, $p_{\theta|Y,\tau}(\beta_i)$ takes two sources of uncertainty or variation into account: variation in β_i for a given value of θ , and the

Figure A1: A Bayesian Hierarchical Model with Conditional Independence



NOTE: Observed data: $Y = \{y_1, \dots, y_N\}$. Heterogeneous parameters: $\beta = \{\beta_1, \dots, \beta_N\}$. Population parameter: θ . Hyperprior parameter: τ .

Author Accepted Manuscript

informed uncertainty in the value of θ itself.

In Stage 1 of BMR's algorithm, customers are randomly assigned into S cohorts, with $N_S = N/S$ customers in each cohort, and with each cohort assigned to a computational node. Y_s is the "shard" of data for customers in cohort s in the partition $Y = \{Y_1, \dots, Y_S\}$. Using a standard MCMC algorithm, R samples are generated from each of $p(\theta | Y_1, \tau), \dots, p(\theta | Y_S, \tau)$ in parallel, across S separate nodes, conditioning on each node's cohort's respective shard of data. These samples are gathered from the various nodes and combined into a pool of $R \times S$ draws, denoted as $\{\theta_s^r\}$. Identical copies of $\{\theta_s^r\}$ are redistributed to the nodes in preparation for Stage 2.

Stage 2 flows from how Eq. A1 is a posterior expected value of $p(\beta_i | \theta)$, and therefore can be approximated as a sample mean across $\{\theta_s^r\}$.

$$(BMR-20) \quad p_{\theta|Y, \tau}(\beta_i) \approx \ddot{p}_{\theta|Y, \tau}(\beta_i) = \frac{1}{RS} \sum_{r,s} p(\beta_i | \theta_s^r).$$

The goal of Stage 2 is to sample each β_i from an approximation to $p(\beta_i | Y)$, denoted as $\ddot{p}(\beta_i | Y)$. BMR do this by sampling uniformly from $\{\theta_s^r\}$, and then sampling $\beta_i | \theta_s^r$ from

$$(A2) \quad p(\beta_i | \theta_s^r) \propto p(y_i | \beta_i) p(\beta_i | \theta_s^r)$$

Repeating these steps samples from the marginal distribution $\ddot{p}(\beta_i | Y)$ by numerically integrating

$$(BMR-21) \quad \ddot{p}(\beta_i | \{\theta_s^r\}, Y) \propto p(y_i | \beta_i) \ddot{p}_{\theta|Y, \tau}(\beta_i)$$

over the empirical posterior distribution of $\{\theta_s^r\}$, which is conditional on Y . In this stage the only data needed to sample $\beta_i | Y$ for a single customer are y_i and $\{\theta_s^r\}$. Because $\{\theta_s^r\}$ is the same for all customers, β_i can be sampled from Eq. A2 in parallel.

WEB APPENDIX B: DERIVING THE CORRECT POSTERIOR

Definitions and Axioms

Define $Y = \{y_1, \dots, y_N\}$ and $\beta = \{\beta_1, \dots, \beta_N\}$. Further, define Y_{-i} as all elements of Y excluding y_i and β_{-i} as all elements of β excluding β_i . Hence, $Y = \{y_i, Y_{-i}\}$ and $\beta = \{\beta_i, \beta_{-i}\}$. A consequence of these definitions is that

$$(B1) \quad p(Y | \beta, \cdot) = f(y_i, Y_{-i} | \beta_i, \beta_{-i}, \cdot)$$

$$(B2) \quad p(\beta | \cdot) = p(\beta_i, \beta_{-i} | \cdot)$$

The following expressions formalize BMR's definition of conditional independence.

$$(B3) \quad p(y_i | \beta_{-i}, \theta) = p(y_i | \theta)$$

$$(B4) \quad p(Y | \beta, \theta) = p(Y | \beta)$$

$$(B5) \quad p(Y | \beta) = \prod_{i=1}^N p(y_i | \beta_i)$$

$$(B6) \quad p(\beta | \theta) = \prod_{i=1}^N p(\beta_i | \theta)$$

It follows from Eqs. B3 to B6 that

$$(B7) \quad p(Y | \beta, \theta) = p(y_i | \beta_i, \theta) p(Y_{-i} | \beta_{-i}, \theta)$$

$$(B8) \quad p(Y | \beta) = p(y_i | \beta_i) p(Y_{-i} | \beta_{-i})$$

$$(B9) \quad p(\beta | \theta) = p(\beta_i | \theta) p(\beta_{-i} | \theta)$$

$$(B10) \quad p(y_i | Y_{-i}, \beta_i, \theta) = p(y_i | \beta_i, \theta)$$

Full Posterior

Bayes Theorem gives the following posterior distributions.

$$(B11) \quad p(\beta | Y, \theta) = \frac{p(Y | \beta, \theta) p(\beta | \theta)}{p(Y | \theta)}$$

$$(B12) \quad p(\theta | Y, \tau) = \frac{p(Y | \theta) p(\theta | \tau)}{p(Y)}$$

Author Accepted Manuscript

The full posterior distribution is

$$(B13) \quad p(\beta, \theta | Y) = p(\beta | Y, \theta)p(\theta | Y, \tau)$$

$$(B14) \quad = \frac{p(Y | \beta)p(\beta | \theta)p(\theta | \tau)}{p(Y)}$$

Correcting the BMR Error

Because conditional independence assumptions are defined at the customer level, we need establish independence of y_i and Y_{-i} after marginalizing over β_{-i} .

Lemma 1. $p(Y | \theta) = p(y_i | \theta)p(Y_{-i} | \theta)$.

Proof. By the Law of Total Probability and Eq. B4,

$$(B15) \quad p(Y | \theta) = \int p(Y | \beta, \theta)p(\beta | \theta)d\beta$$

$$(B16) \quad = \int p(Y | \beta)p(\beta | \theta)d\beta$$

Substitute Eqs. B8 and B9 into Eq. B16.

$$(B17) \quad p(Y | \theta) = \iint p(y_i | \beta_i)p(Y_{-i} | \beta_{-i})p(\beta_i | \theta)p(\beta_{-i} | \theta)d\beta_i d\beta_{-i}$$

Rearrange terms, factor the integral into two parts, and integrate over β_i and β_{-i} separately to get the result.

$$(B18) \quad p(Y | \theta) = \int p(y_i | \beta_i)p(\beta_i | \theta)d\beta_i \cdot \int p(Y_{-i} | \beta_{-i})p(\beta_{-i} | \theta)d\beta_{-i}$$

$$(B19) \quad = p(y_i | \theta)p(Y_{-i} | \theta)$$

■

Next, we show that conditional on θ , the only data that affects β_i is y_i .

Lemma 2. $p(\beta_i | Y, \theta) = p(\beta_i | y_i, \theta)$

Proof. By Bayes' Theorem,

$$(B20) \quad p(\beta_i | Y, \theta) = \frac{p(Y | \beta_i, \theta)}{p(Y | \theta)}p(\beta_i | \theta)$$

Factor the numerator in Eq. B20.

$$(B21) \quad p(Y | \beta_i, \theta) = f(y_i, Y_{-i} | \beta_i, \theta)$$

$$(B22) \quad = p(y_i | Y_{-i}, \beta_i, \theta) p(Y_{-i} | \beta_i, \theta)$$

A corollary to Eq. B3 is that $p(Y_{-i} | \beta_i, \theta) = p(Y_{-i} | \theta)$. Substitute that and Eq. B10 into Eq. B22.

$$(B23) \quad p(Y | \beta_i, \theta) = p(y_i | \beta_i, \theta) p(Y_{-i} | \theta)$$

Back to Eq. B20, substitute Eq. B23 in the numerator, and the result of Lemma 1 in the denominator.

$$(B24) \quad p(\beta_i | Y, \theta) = \frac{p(y_i | \beta_i, \theta) p(Y_{-i} | \theta)}{p(y_i | \theta) p(Y_{-i} | \theta)} p(\beta_i | \theta)$$

After cancelling terms in Eq. B24, applying Bayes Theorem to the RHS gives the result. ■

The next lemma states the recursive property of Bayesian updating.

Lemma 3.

$$(B25) \quad p(\theta | Y, \tau) = \frac{p(y_i | \theta)}{p(y_i | Y_{-i})} p(\theta | Y_{-i}, \tau)$$

Proof. By Bayes' Theorem,

$$(B26) \quad p(\theta | Y, \tau) = \frac{p(Y | \theta) p(\theta | \tau)}{p(Y)}$$

$$(B27) \quad p(\theta | Y_{-i}, \tau) = \frac{p(Y_{-i} | \theta) p(\theta | \tau)}{p(Y_{-i})}$$

In Eq. B26, factor $p(Y | \theta)$ using Lemma 1 and factor $p(Y)$ using the definition of joint probability.

$$(B28) \quad p(\theta | Y, \tau) = \frac{p(y_i | \theta)}{p(y_i | Y_{-i})} \frac{p(Y_{-i} | \theta) p(\theta | \tau)}{p(Y_{-i})}$$

Replace the second fraction in the RHS of Eq. B28 with the LHS of Eq. B27 to get the result. ■

We can now derive the correct $p(\beta_i | Y, \tau)$ to replace BMR-13.

Author Accepted Manuscript

Proposition B1.

$$(B29) \quad p(\beta_i | Y, \tau) = \frac{p(y_i | \beta_i) p_{\theta | Y_{-i}, \tau}(\beta_i)}{p(y_i | Y_{-i})}$$

Proof. We restate BMR-7 here. Note that $p(\beta_i | Y, \theta) = p(\beta_i | y_i, \theta)$ is established in Lemma 2.

$$(BMR-7) \quad p(\beta_i | Y, \tau) = \int p(\beta_i | y_i, \theta) p(\theta | Y, \tau) d\theta,$$

Replace $p(\beta_i | y_i, \theta)$ with Eq. 1, and replace $p(\theta | Y, \tau)$ with the result from Lemma 3.

$$(B30) \quad p(\beta_i | Y, \tau) = \int \frac{p(y_i | \beta_i) p(\beta_i | \theta)}{p(y_i | \theta)} \frac{p(y_i | \theta)}{p(y_i | Y_{-i})} p(\theta | Y_{-i}, \tau) d\theta$$

Cancel the $p(y_i | \theta)$ terms and factor out terms that do not depend on θ .

$$(B31) \quad p(\beta_i | Y, \tau) = \frac{p(y_i | \beta_i)}{p(y_i | Y_{-i})} \int p(\beta_i | \theta) p(\theta | Y_{-i}, \tau) d\theta$$

The integrand in Eq. B31 is defined in Eq. D2, giving the result.

$$(B32) \quad p(\beta_i | Y, \tau) = \frac{p(y_i | \beta_i)}{p(y_i | Y_{-i})} p_{\theta | Y_{-i}, \tau}(\beta_i)$$

■

Author Accepted Manuscript

WEB APPENDIX C: PARALLELIZING A CORRECTED ALGORITHM

A practitioner may wonder if there is any way to amend the BMR algorithm to be consistent with the corrected math. To maintain the same spirit of BMR, a modified algorithm should retain two parallelizable stages with minimal communication overhead between nodes. Otherwise, it would be a different algorithm altogether.

As a constructive measure for BMR and potential adopters, we describe one possible approach, with no claims of optimality. Stage 1 would have to generate a separate $\{\theta_s^r\}_i$ for each customer, where y_i is held out for each run. First, generate “full shard” samples from $p(\theta | Y_s, \tau)$ on each node. Then, sample from $p(\theta | Y_{s-i}, \tau)$ on each customer i 's cohort's node, where Y_{s-i} is the shard for cohort s after removing y_i . This would happen N_S times on each node, once for each customer in the cohort. The pool for customer i , $\{\theta_s^r\}_i$, would combine the full shard samples generated on all the other nodes with samples from $p(\theta | Y_{-i}, \tau)$. This approach would take $N_S + 1$ times as long as BMR's algorithm.

Author Accepted Manuscript

WEB APPENDIX D: THE BMR ALGORITHM OVERFITS CUSTOMER DATA

In this section we offer a Bayesian theoretical perspective on why BMR's derivation of the target posterior density could not have been correct. BMR express the posterior density of β_i as

$$(BMR-13) \quad p(\beta_i | Y) = \frac{p_{\theta|Y,\tau}(\beta_i)p(y_i | \beta_i)}{p(y_i)},$$

where

$$(D1) \quad p_{\theta|Y,\tau}(\beta_i) = \int p(\beta_i | \theta)p(\theta | Y, \tau)d\theta$$

In BMR-13, $p_{\theta|Y,\tau}(\beta_i)$ acts as an informative prior on β_i . But BMR describe $p_{\theta|Y,\tau}(\beta_i)$ as both “the posterior predictive density of β_i : the density of β_i before observing y_i , given $p(\theta | Y, \tau)$ ” and “a highly informative prior distribution for β_i , before observing y_i ” (p. 1003, col. 1). These statements are difficult to reconcile because y_i is already included in Y . BMR-13 is double-counting the focal customer's data, once as part of the information contained in $p(\theta | Y, \tau)$, and then again in the data likelihood $p(y_i | \beta_i)$. Therefore, $p_{\theta|Y,\tau}(\beta_i)$ cannot be a valid prior on β_i .

A more suitable prior on β_i would contain information about Y_{-i} : all customer data *except* y_i .

$$(D2) \quad p_{\theta|Y_{-i},\tau}(\beta_i) = \int p(\beta_i | \theta)p(\theta | Y_{-i}, \tau)d\theta$$

Because $p(\theta | Y_{-i}, \tau)$ does not condition on data from customer i , $p_{\theta|Y_{-i},\tau}(\beta_i)$ truly captures prior information about β_i before y_i is observed. Then, once y_i is observed,

$$(D3) \quad p(\beta_i | Y, \tau) \propto p(y_i | \beta_i)p_{\theta|Y_{-i},\tau}(\beta_i)$$

reflects the properly updated posterior distribution of β_i . This is precisely what we get when we correct BMR's error.

The BMR algorithm reflects their incorrect derivation, so it creates *systematic* bias in $p(\beta_i | Y, \tau)$. Note that y_i is not a *random* observation, but the data that corresponds to the specific β_i being inferred. Unlike $p_{\theta|Y_{-i},\tau}(\beta_i)$, $p_{\theta|Y,\tau}(\beta_i)$ is influenced by y_i , so BMR-estimated posterior distributions of $\beta_i | Y$ are biased

Author Accepted Manuscript

toward the observed data. If y_i were, say, an outlier data point, and the algorithm were deployed on a large number of nodes (large S), the magnitude of the bias could be substantial. Potential adopters of the algorithm should be aware that the error and associated bias exist.

Peer Review Version

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Author Accepted Manuscript

WEB APPENDIX E: NOTES ON BMR'S THEORETICAL PROOFS

In this section we explain how the mathematical error affects BMR's theorems in their Web Appendix. We do not discuss theorems whose subject matter falls outside the scope of this paper. We include this section to assist BMR and the reader, and we do not claim rigorous results.

The following expressions are defined in BMR and used in this section, but were not used in our main text.

$$(BMR-10) \quad p(\beta_i | Y) = E_{\theta|Y, \tau}(p(\beta_i | \theta))p(y_i | \beta_i)$$

$$(BMR-14) \quad \dot{p}_{\theta|Y, \tau}(\beta_i) = \frac{1}{R} \sum_r p(\beta_i | \theta^r)$$

$$(BMR-15) \quad \dot{p}(\beta_i | \{\theta^r\}, Y) \propto \dot{p}_{\theta|Y, \tau}(\beta_i)p(y_i | \beta_i)$$

Because $E_{\theta|Y, \tau}(p(\beta_i | \theta)) = p_{\theta|Y, \tau}(\beta_i)$, BMR-10 is equivalent to BMR-13.

Theorem 1

Using our notation, Gelman et al. (2004, Eq. 1.4) defines a posterior predictive density (PPD) as

$$(E1) \quad p(y_i | Y_{-i}) = \int p(y_i | \beta_i, \theta)g(\beta_i, \theta | Y_{-i})d\beta_i d\theta$$

That is, it is a predictive distribution for new data, conditional on previously observed data. BMR call $p(\beta_i | Y_{-i})$ a PPD, but it is better described as a "marginal posterior distribution." Note that Theorem 1 is not referenced in the BMR paper at all. In the main text, BMR refer to $p_{\theta|Y, \tau}(\beta_i)$ (conditional on Y) as a PPD (p. 1003, col. 1).

Theorem 2

Restate this theorem as

$$(E2) \quad E(\dot{p}_{\theta|Y_{-i}, \tau}(\beta_i)) = E_{\theta|Y_{-i}}(p(\beta_i)) = p_{\theta|Y_{-i}, \tau}(\beta_i)$$

Let θ^r be a sample from $p(\theta | Y_{-i}, \tau)$. Replacing BMR-14 with

$$(E3) \quad \dot{p}_{\theta|Y_{-i}, \tau}(\beta_i) = \frac{1}{R} \sum_r p(\beta_i | \theta^r)$$

1
2
3 makes the restated theorem correct.
4
5

6 **Theorem 3**
7

8 This theorem does not need to be restated. It would be correct after making the following changes.
9

- 10 1. Replace BMR-10 with $p(\beta_i | Y) = E_{\theta|Y_{-i}, \tau}(p(\beta_i | \theta))p(y_i | \beta_i)$.
11
12 2. Define $\{\theta^r\}$ as samples from $p(\theta | Y_{-i}, \tau)$ instead of $p(\theta | Y, \tau)$.
13
14 3. Replace BMR-15 with $\dot{p}(\beta_i | \{\theta^r\}, Y, \tau) \propto \dot{p}_{\theta|Y_{-i}, \tau}(\beta_i)p(y_i | \beta_i)$
15
16
17 4. Apply the restated Theorem 2.
18
19

20 **Theorem 4**
21

22 BMR acknowledge taking a frequentist perspective to their asymptotic analysis, where “Y is a random
23 sample from a distribution for some fixed, nonrandom, unknown parameter.” In that sense, $\sqrt{N}(\theta_N^r - \theta) \xrightarrow{P}$
24 $N(0, I_\theta^{-1})$ represents convergence in probability to the fixed parameter θ . But in the rest of the paper, the
25 treatment of θ is Bayesian, as a random variable. In the convergence proofs, the fixed θ should be noted as
26 something else, like θ^* .
27
28
29

30 We think Theorem 4.1.1 should be correct if $\dot{p}_{\theta|Y, \tau}(\beta_i)$ were replaced with $\dot{p}_{\theta|Y_{-i}, \tau}(\beta_i)$ and $\{\theta^r\}$ are
31 samples from $p(\theta | Y_{-i}, \tau)$ instead of $p(\theta | Y, \tau)$. However, partitioning Y_{-i} into shards, with $Y_{-i} =$
32 $\{Y_{1:S-1}, Y_{S-i}\}$ would leave one shard with $N_S - 1$ observations. It is possible that the theorem is still valid in
33 the limit, but that remains to be proven.
34
35
36
37
38
39

40 **Theorem 5**
41

42 Restate the theorem as $\lim_{N \rightarrow \infty} [\ddot{p}_{\theta|Y_{-i}, \tau}(\beta_i)] = \dot{p}_{\theta|Y_{-i}, \tau}(\beta_i)$. If we accept that Theorem 4 holds, then the
43 theorem can be corrected with the following changes.
44
45

- 46 1. Denote θ_N^r as a draw from $p(\theta | Y_{-i}, \tau)$ instead of $p(\theta | Y, \tau)$.
47
48 2. Replace all instances of $\dot{p}_{\theta|Y, \tau}(\beta_i)$ with $\dot{p}_{\theta|Y_{-i}, \tau}(\beta_i)$.
49
50

51 **Theorem 6**
52

53 This theorem does not need to be restated. It would be correct after making the following changes.
54
55

- 56 1. In BMR-21, replace $\ddot{p}_{\theta|Y, \tau}(\beta_i)$ with $\ddot{p}_{\theta|Y_{-i}, \tau}(\beta_i)$.
57
58
59
60

Author Accepted Manuscript

2. Define $\{\theta^r\}$ as samples from $p(\theta | Y_{-i}, \tau)$ instead of $p(\theta | Y, \tau)$.
3. Apply the restated Theorem 5 and Theorem 2.

Theorem 8

Because one observation needs to be held out of Y_S , it is not clear that the assumptions of Theorem 4 apply.

Theorem 11

By Eq. B29, BMR-13 should be

$$(E4) \quad p(\beta_i | Y) = \frac{p(y_i | \beta_i) p_{\theta|Y_{-i}, \tau}(\beta_i)}{p(y_i | Y_{-i})}$$

To correct this theorem, replace $p_{\theta|Y, \tau}(\beta_i)$, $\dot{p}_{\theta|Y, \tau}(\beta_i)$, and $\ddot{p}_{\theta|Y, \tau}(\beta_i)$ with $p_{\theta|Y_{-i}, \tau}(\beta_i)$, $\dot{p}_{\theta|Y_{-i}, \tau}(\beta_i)$, and $\ddot{p}_{\theta|Y_{-i}, \tau}(\beta_i)$, respectively. Also, replace the normalizing constant $p(y_i)$ with $p(y_i | Y_{-i})$.

However, it is not immediately clear if the conclusion of this theorem can be supported by Theorem 4.

WEB APPENDIX F: PARALLELIZATION BIAS

Proposition F1. For $S > 1$, the expected variance of $p^*(\theta | Y_{1:S}, \tau)$ is strictly greater than the variance of $p(\theta | Y, \tau)$.

Proof. Stage 1 of the published BMR algorithm samples θ from

$$(F1) \quad p^*(\theta | Y_{1:S}, \tau) = \frac{1}{S} \sum_{s=1}^S p(\theta | Y_s, \tau)$$

The first two moments of this distribution are

$$(F2) \quad E^*(\theta | Y_{1:S}, \tau) = \frac{1}{S} \sum_{s=1}^S E(\theta | Y_s, \tau)$$

$$(F3) \quad E^*(\theta^2 | Y_{1:S}, \tau) = \frac{1}{S} \sum_{s=1}^S E(\theta^2 | Y_s, \tau)$$

So the variance is

$$(F4) \quad \text{var}^*(\theta | Y_{1:S}, \tau) = E^*(\theta^2 | Y_{1:S}, \tau) - E^*(\theta | Y_{1:S}, \tau)^2$$

$$(F5) \quad = \frac{1}{S} \sum_{s=1}^S E(\theta^2 | Y_s, \tau) - \left[\frac{1}{S} \sum_{s=1}^S E(\theta | Y_s, \tau) \right]^2$$

By the Law of Total Variance,

$$(F6) \quad E(\theta^2 | Y_s, \tau) = \text{var}(\theta | Y_s, \tau) + E(\theta | Y_s, \tau)^2$$

Substitute Eq. F6 into the first summand in Eq. F5.

$$(F7) \quad \text{var}^*(\theta | Y_{1:S}, \tau) = \frac{1}{S} \sum_{s=1}^S \text{var}(\theta | Y_s, \tau) + \frac{1}{S} \sum_{s=1}^S E(\theta | Y_s, \tau)^2 - \left[\frac{1}{S} \sum_{s=1}^S E(\theta | Y_s, \tau) \right]^2$$

By Jensen's inequality, $\sum_{s=1}^S E(\theta | Y_s, \tau)^2 \geq \left[\sum_{s=1}^S E(\theta | Y_s, \tau) \right]^2$. Therefore,

$$(F8) \quad \text{var}^*(\theta | Y_{1:S}, \tau) \geq \frac{1}{S} \sum_{s=1}^S \text{var}(\theta | Y_s, \tau)$$

So we know that the variance of $p^*(\theta | Y_{1:S}, \tau)$ will be at least as great as the variance of $p(\theta | Y_{1:S}, \tau)$.

If $S > 1$, then $N_S < N$ and $E_Y(\text{var}(\theta | Y_s, \tau)) > \text{var}(\theta | Y, \tau)$ (where $E_Y(\cdot)$ is a expectation across

Author Accepted Manuscript

random partitions in a frequentist sense). Averaging across shards,

$$(F9) \quad \frac{1}{S} \sum_{s=1}^S E_Y(\text{var}(\theta | Y_s, \tau)) > \text{var}(\theta | Y, \tau)$$

Taking the expectation across all terms in Eq. F8, and combining with Eq. F9, gives us

$$(F10) \quad E_Y(\text{var}(\theta | Y_{1:S}, \tau)) \geq \frac{1}{S} \sum_{s=1}^S E_Y(\text{var}(\theta | Y_s, \tau)) > \text{var}(\theta | Y, \tau)$$

In Eq. F7, we see that $\text{var}(\theta | Y_{1:S}, \tau)$ represents two sources of variation in θ . The first term on the RHS is the average of the shard-level posterior variances. The second and third terms on the RHS define the variation in the shard-level means across shards.

Proposition F2. For $S > 1$, the expected variance of $p^*(\beta_i | Y_{1:S}, \tau)$ is strictly greater than the variance of $p(\beta_i | Y, \tau)$.

Proof. By the Law of Total Variance,

$$(F11) \quad \text{var}(\beta_i | Y, \tau) = E(\text{var}(\beta_i | \theta, \tau)) + \text{var}(E(\beta_i | \theta, \tau))$$

$$(F12) \quad \text{var}^*(\beta_i | Y, \tau) = E^*(\text{var}(\beta_i | \theta, \tau)) + \text{var}^*(E(\beta_i | \theta, \tau))$$

If $E(\theta | Y, \tau) = E^*(\theta | Y_{1:S}, \tau)$, then $E(\text{var}(\beta_i | \theta, \tau)) = E^*(\text{var}(\beta_i | \theta, \tau))$. Therefore, if $\text{var}^*(E(\beta_i | \theta, \tau)) > \text{var}(E(\beta_i | \theta, \tau))$, then $\text{var}^*(\beta_i | Y, \tau) > \text{var}(\beta_i | Y, \tau)$. Intuitively, this occurs because an increase in the posterior variance of $\theta | Y$ results in greater variance in the means of $\beta_i | \theta$, so $\text{var}(\beta_i | Y, \tau)$ increases as well.

WEB APPENDIX REFERENCES

Gelman, Andrew, John B. Carlin, Hal S. Stern, and Donald B. Rubin (2004). *Bayesian Data Analysis*. 2nd ed. Boca Raton, Fla.: Chapman and Hall.