
Massively Parallel Importance Weighting and Sampling

Sam Bowyer

School of Mathematics
University of Bristol
Bristol

Thomas Heap

Department of Computer Science
University of Bristol
Bristol

Laurence Aitchison

Department of Computer Science
University of Bristol
Bristol
`laurence.aitchison@bristol.ac.uk`

Abstract

Importance sampling is a popular technique in Bayesian inference: by reweighting samples drawn from a proposal distribution we are able to obtain samples and moment estimates from a Bayesian posterior over some n latent variables. Recent work, however, indicates that importance sampling scales poorly — in order to accurately approximate the true posterior, the required number of importance samples grows exponentially in the number of latent variables [Chatterjee and Diaconis, 2018]. We work around this issue by drawing K samples for each of the n latent variables and reasoning about all K^n combinations of latent samples. In principle, we reason efficiently over K^n combinations of samples by exploiting conditional independencies in the generative model. However, in practice this requires complex algorithms. Instead, we exploit a trick from physics to compute all the required quantities — posterior expectations, marginals and samples — by differentiating through a marginal likelihood estimator.

1 Introduction

Importance weighting allows us to reweight samples drawn from a proposal in order to compute expectations of a different distribution, such as a Bayesian posterior. However, importance weighting breaks down in larger models. Chatterjee and Diaconis [2018] showed that the number of samples required to accurately approximate the true posterior scales as $\exp(D_{\text{KL}}(P(z|x) || Q(z)))$, where $P(z|x)$ is the true posterior over latent variables, z , given data x , and $Q(z)$ is the proposal. Problematically, we expect the KL divergence to scale with n , the number of latent variables. Indeed, if $P(z|x)$ and $Q(z)$ are IID with n variables, then the KL-divergence is exactly proportional to n . Thus, we expect the required number of importance samples to be exponential in the number of latent variables, and hence we expect accurate importance sampling to be intractable in larger models.

To resolve this issue, we introduce a massively parallel importance sampling scheme that in effect uses an exponential number of samples to compute posterior expectations, marginals and samples. We do this by drawing K samples of each of the n latent variables from the proposal, then individually reweighting all K^n combinations of all samples of all latent variables. While reasoning about all K^n combinations of samples might seem intractable, we should in principle be able to perform efficient computations by exploiting conditional independencies in the underlying graphical model.

Of course, many computations that are possible in principle are extremely complex in practice, and that turns out to be the case here.

We then noticed that we could perform this reasoning over K^n latent variables using methods from the discrete graphical model literature. However, this turned out to be less helpful than expected because these algorithms are highly complex, and different for computing posterior expectations, marginals and samples. We therefore develop a much simpler approach to computing posterior expectations, marginals and samples, which entirely avoids the need to explicitly write backwards computations. In particular, we show that posterior expectations, marginals and samples can be obtained simply by differentiating through (a slightly modified) forward computation that produces an estimate of the marginal likelihood. The required gradients can be computed straightforwardly using modern autodiff, and the resulting implicit backward computations automatically inherit potentially complex optimizations from the forward pass.

2 Related work

There is a considerable body of work in the discrete graphical model setting that computes posterior expectations, marginals and samples [Dawid, 1992, Pfeffer, 2005, Bidyuk and Dechter, 2007, Geldenhuys et al., 2012, Claret et al., 2013, Sankaranarayanan et al., 2013, Goodman and Stuhlmüller, 2014, Gehr et al., 2016, Narayanan et al., 2016, Albarghouthi et al., 2017, Wang et al., 2018, Obermeyer et al., 2019, Holtzen et al., 2020]. Our work differs in two respects. First, our massively parallel methods are not restricted to discrete graphical models, but can operate with arbitrary continuous latent variables. Second, this work involves complex implementations that — in one sense or another — “proceed by recording an adjoint compute graph alongside the forward computation and then traversing the adjoint graph backwards starting from the final result of the forward computation” [Obermeyer et al., 2019]. The forward computation is reasonably straightforward: it is just a big tensor product that can be computed efficiently using pre-existing libraries such as `opt-einsum`, and results in (an estimate of) the marginal likelihood. However, the backward traversal is much more complex, if for no other reason than the need to implement separate traversals for each operation of interest (computing posterior expectations, marginals and samples). Additionally, these traversals need to correctly handle all special cases (including e.g. optimized implementations of plates and timeseries). Importantly, the optimized forward is often relatively straightforward to implement, while the backward traversal is far more complex. For instance, the forward computation for a timeseries involves a product of T matrices arranged in a chain. Naively computing this product on GPUs is very slow, as it requires T separate matrix multiplications. However, it is possible to massively optimize this forward computation, converting $\mathcal{O}(T)$ to $\mathcal{O}(\log(T))$ tensor operations by multiplying adjacent pairs of matrices in a single batched matrix multiplication operation. This optimization is straightforward in the forward computation. However, applying this optimization as part of the backward computation is far more complex (see Corenflos et al., 2022 for details). This complexity (along with similar complexity for other important optimizations such as plates) was prohibitive for our small team (which ultimately aims to implement a new massively parallel probabilistic programming language). In contrast, we provide a much simpler approach, where we only explicitly perform the forward computation; we compute the posterior expectations, marginals and samples by differentiating through that forward computation.

There is work on fitting importance weighted autoencoders [Burda et al., 2015, IWAE;] and reweighted wake-sleep [RWS Bornschein and Bengio, 2014, Le et al., 2020] in the massively parallel setting [Aitchison, 2019, Geffner and Domke, 2022, Anonymous, 2023] for general probabilistic models (i.e. with continuous latent variables). However, this work only provides methods for performing massively parallel updates to approximate posteriors (e.g. by optimizing a massively parallel ELBO). This work does not provide a method to individually reweight the samples to provide accurate posterior expectations, marginals and samples. Instead, this previous work simply takes the learned approximate posterior as an estimate of the true posterior, and does not attempt to correct for inevitable biases.

Massively parallel importance weighting for timeseries have some similarities to particle filtering/SMC methods [Gordon et al., 1993, Doucet et al., 2009, Andrieu et al., 2010, Maddison et al., 2017, Le et al., 2017, Lindsten et al., 2017, Naesseth et al., 2018, Lai et al., 2022]. However, our massively parallel methods allow a very general class of proposal distributions in a very general class of probabilistic models. Additionally, we develop a new approach to computing posterior ex-

pectations, marginals and samples by differentiating through a slightly modified marginal likelihood estimator.

3 Background

Bayesian inference. In Bayesian inference, we have a prior, $P(z')$ over latent variables, z' , and a likelihood, $P(x|z')$ connecting the latents to the data, x . Here, we use z' rather than z because we reserve z for future use as a collection of K samples (Eq. 3). Our goal is to compute the posterior distribution over latent variables conditioned on observed data,

$$P(z'|x) = \frac{P(x|z') P(z')}{\sum_{z''} P(x, z'')}, \quad (1)$$

We often seek to obtain samples from the posterior or to compute posterior expectations,

$$m_{\text{post}} = \sum_{z'} P(z'|x) m(z') \quad (2)$$

However, the true posterior moment is usually intractable, so instead we are forced to use an alternative method such as importance weighting.

Importance weighting. In importance weighting, we draw a collection of K samples from the full joint state space. This collection of K samples is denoted $z \in \mathcal{Z}$, with a single sample denoted z^k ,

$$z = (z^1, z^2, \dots, z^K) \in \mathcal{Z}^K. \quad (3)$$

The collection of K samples, z , is drawn by sampling K times from the proposal,

$$Q(z) = \prod_{k \in \mathcal{K}} Q(z^k), \quad (4)$$

where \mathcal{K} is the set of possible indices, $\mathcal{K} = \{1, \dots, K\}$. As the true posterior moment is usually intractable, one approach is to use a self-normalized importance sampling estimate, $m_{\text{global}}(z)$. We call this a “global” importance weighted estimate following terminology in [Geffner and Domke, 2022] and in contrast with the massively parallel methods that we define later,

$$m_{\text{global}}(z) = \frac{1}{K} \sum_{k \in \mathcal{K}} \frac{r_k(z)}{\mathcal{P}_{\text{global}}(z)} m(z^k) \quad (5)$$

where,

$$r_k(z) = \frac{P(x, z^k)}{Q(z^k)} \quad (6)$$

$$\mathcal{P}_{\text{global}}(z) = \frac{1}{K} \sum_{k \in \mathcal{K}} r_k(z) \quad (7)$$

with samples z drawn from the proposal (Eq. 4). Here, $r_k(z)$ is the ratio of the generative and proposal probabilities, and $\mathcal{P}_{\text{global}}(z)$ is an unbiased estimator of the marginal likelihood,

$$\begin{aligned} E_{Q(z)} [\mathcal{P}_{\text{global}}(z)] &= E_{Q(z^k)} \left[\frac{P(x, z^k)}{Q(z^k)} \right] \\ &= \sum_{z^k} P(x, z^k) = P(x) \end{aligned} \quad (8)$$

The first equality arises because $\mathcal{P}_{\text{global}}(z)$ is the average of K IID terms, $P(x, z^k)/Q(z^k)$, so is equal to the expectation of a single term, and the second equality arises if we write the expectation as a sum.

Source term trick. Here, we outline a standard trick from physics that can be used to compute expectations of arbitrary probability distribution by differentiating a modified log-normalizing constant. This trick is used frequently in Quantum Field Theory, for instance [Weinberg, 1995] (Chapter

16), and also turns up in the theory of neural networks [Zavatone-Veth et al., 2021]. But the trick is simple enough that we can give a self-contained introduction here.

In our context, Bayes theorem (Eq. 1) defines an unnormalized density, $P(z|x) \propto P(x, z)$, with normalizing constant, $\sum_{z'} P(x, z')$. Of course, the normalizing constant is usually intractable, but in our massively parallel context, we can often obtain a reasonable estimate of the normalizing constant, and that will be sufficient. It turns out that we can compute posterior expectations using a slightly modified normalizing constant,

$$Z_m(J) = \sum_{z'} P(x, z') e^{Jm(z')}. \quad (9)$$

where $e^{Jm(z')}$ is known as a source term. Note that setting J to zero recovers the usual normalizing constant,

$$Z_m(J=0) = \sum_{z'} P(x, z'). \quad (10)$$

Now, we can extract the posterior moment by evaluating the gradient of $\log Z_m(J)$ at $J=0$,

$$\left. \frac{\partial}{\partial J} \right|_{J=0} \log Z_m(J) = \left. \frac{\partial}{\partial J} \right|_{J=0} \log \sum_{z'} P(x, z') e^{Jm(z')}. \quad (11)$$

Differentiating the logarithm at $J=0$,

$$\left. \frac{\partial}{\partial J} \right|_{J=0} \log Z_m(J) = \frac{\sum_{z'} P(x, z') \left. \frac{\partial}{\partial J} \right|_{J=0} e^{Jm(z')}}{\sum_{z''} P(x, z'')}. \quad (12)$$

Differentiating the exponential at $J=0$,

$$\left. \frac{\partial}{\partial J} \right|_{J=0} \log Z_m(J) = \sum_{z'} \frac{P(x, z')}{\sum_{z''} P(x, z'')} m(z'). \quad (13)$$

and identifying the posterior using Bayes theorem (Eq. 1).

$$\left. \frac{\partial}{\partial J} \right|_{J=0} \log Z_m(J) = \sum_{z'} P(z'|x) m(z) = m_{\text{post}} \quad (14)$$

This is exactly the form for the posterior moment in Eq. (2).

Massively parallel marginal likelihood estimators To get an accurate marginal likelihood estimator, we introduce a massively parallel estimator which individually weights all K^n combinations of K samples on n latent variables. We write the k th sample of the i th latent variable as z_i^k , and we can then write the collection of all K samples of the i th latent variable as,

$$z_i = (z_i^1, z_i^2, \dots, z_i^K). \quad (15)$$

And z is the collection of all K samples of all n latent variables, (as in Eq. 3),

$$z = (z_1, z_2, \dots, z_n). \quad (16)$$

As we have multiple latent variables, our massively-parallel proposals have a graphical model structure,

$$Q_{\text{MP}}(z) = \prod_{i=1}^n Q_{\text{MP}}(z_i | z_j \text{ for } j \in \text{qa}(i)), \quad (17)$$

where $\text{qa}(i)$ is the set of indices of parents of z_i under that graphical model. Note that this form implies that z_i may depend on all samples of the parent latent variables, and that the K samples of z_i , namely z_i^1, \dots, z_i^K , may have some dependencies. Usually, we derive the massively-parallel proposal from an underlying, user-specified, single-sample proposal. The simplest approach is just to draw K samples independently from the full joint latent space. Other alternatives are available, including to use a uniform mixture / permutation over samples of parent latent variables (see [Anonymous, 2023] for further details). The key property of these proposals is that the single-sample

marginals, $Q_{\text{MP}}(z_i^k | z_j \text{ for } j \in \text{pa}(i))$ are easily computable. In all our experiments, we will use a permutation over samples of the parent latent variables.

For the generative model, we need to explicitly consider all K^n combinations of K samples on n latent variables. To help us write down these combinations, we define a vector of indices,

$$\mathbf{k} = (k_1, k_2, \dots, k_n) \in \mathcal{K}^n, \quad (18)$$

with one index, k_i for each latent variable, z_i . Thus, we can write the “indexed” latent variables as,

$$z^{\mathbf{k}} = (z_1^{k_1}, z_2^{k_2}, \dots, z_n^{k_n}) \in \mathcal{Z}, \quad (19)$$

which represents a single sample from the full joint latent space. The generative model also has graphical model structure, with the set of indices of parents of the i th latent variable under the generative model begin denoted $\text{pa}(i)$ (contrast this with $\text{qa}(i)$ which is the parents of the i th latent variable under the proposal).

Thus, the generative probability for a single combination of samples, denoted $z^{\mathbf{k}}$, can be written as,

$$\begin{aligned} P(x, z^{\mathbf{k}}) &= P\left(x \mid z_j^{k_j} \text{ for all } j \in \text{pa}(x)\right) \\ &\quad \prod_{i=1}^n P\left(z_i^{k_i} \mid z_j^{k_j} \text{ for all } j \in \text{pa}(i)\right). \end{aligned} \quad (20)$$

Thus, we can write a massively parallel marginal likelihood estimator as,

$$r_{\mathbf{k}}(z) = \frac{P(x, z^{\mathbf{k}})}{\prod_i Q_{\text{MP}}\left(z_i^{k_i} \mid z_j \text{ for } j \in \text{qa}(i)\right)} \quad (21)$$

$$\mathcal{P}_{\text{MP}}(z) = \frac{1}{K^n} \sum_{\mathbf{k} \in \mathcal{K}^n} r_{\mathbf{k}}(z). \quad (22)$$

While this looks intuitively reasonable, proving that Eq. (22) is a valid marginal likelihood estimator: the full proofs are given in [Anonymous, 2023] (their Appendix C.1.3).

The next challenge is to compute the sum in Eq. (22). The sum looks intractable as we have to sum over K^n settings of \mathbf{k} . However, it turns out that these sums usually are tractable. The reason is that that if we fix the samples, z , then $r_{\mathbf{k}}(z)$ can be understood as a product of low-rank tensors,

$$r_{\mathbf{k}}(z) = f_{\mathbf{k}_{\text{pa}(x)}}^x(z) \prod_i f_{k_i, \mathbf{k}_{\text{pa}(i)}}^i(z) \quad (23)$$

$$f_{\mathbf{k}_{\text{pa}(x)}}^x(z) = P\left(x \mid z_j^{k_j} \text{ for all } j \in \text{pa}(x)\right), \quad (24)$$

$$f_{k_i, \mathbf{k}_{\text{pa}(i)}}^i(z) = \frac{P\left(z_i^{k_i} \mid z_j^{k_j} \text{ for all } j \in \text{pa}(i)\right)}{Q_{\text{MP}}\left(z_i^{k_i} \mid z_j \text{ for all } j \in \text{qa}(i)\right)}. \quad (25)$$

Here, $f_{\mathbf{k}_{\text{pa}(x)}}^x(z)$ is a rank $|\text{pa}(x)|$ tensor, and $f_{k_i, \mathbf{k}_{\text{pa}(i)}}^i(z)$ are rank $1 + |\text{pa}(i)|$ tensors, where $|\text{pa}(i)|$ denotes the number of parents of the i th latent variable. Thus, the sum in Eq. (22) is really a large tensor product,

$$\mathcal{P}_{\text{MP}}(z) = \frac{1}{K^n} \sum_{\mathbf{k} \in \mathcal{K}^n} f_{\mathbf{k}_{\text{pa}(x)}}^x(z) \prod_i f_{k_i, \mathbf{k}_{\text{pa}(i)}}^i(z) \quad (26)$$

which can be directly and efficiently computed using an opt-einsum implementation. This is discussed in more depth in [Anonymous, 2023].

Of course, the contributions of this paper are not in computing the unbiased marginal likelihood estimator. They are in using the marginal likelihood estimator to compute the other key quantities of interest in Bayesian computations, namely posterior expectations, marginals and samples, and we consider each in turn in the following sections.

4 Methods

Now, we can define an importance sampling scheme that operates on all K^n combinations of samples,

$$m_{\text{MP}}(z) = \frac{1}{K^n} \sum_{\mathbf{k} \in \mathcal{K}^n} \frac{r_{\mathbf{k}}(z)}{\mathcal{P}_{\text{MP}}(z)} m(z^{\mathbf{k}}). \quad (27)$$

This looks very similar to the standard global importance sampling scheme in Eq. (5), except that Eq. (5) averages only over K samples, whereas this massively parallel moment estimator averages over all K^n combinations of samples. Proving that this is a valid importance-sampled moment estimator is not trivial (see Appendix A for details).

4.1 Interpreting massively parallel importance weighting as inference in a discrete graphical model

Now, $\frac{1}{K^n} r_{\mathbf{k}}(z) / \mathcal{P}_{\text{MP}}(z)$ in (Eq. 27) can be understood as a normalized probability distribution over \mathbf{k} . In particular, this quantity is always positive, and we can show that it normalizes to 1 by substituting the definition of $\mathcal{P}_{\text{MP}}(z)$ from Eq. (22),

$$\frac{1}{K^n} \sum_{\mathbf{k}} \frac{r_{\mathbf{k}}(z)}{\mathcal{P}_{\text{MP}}(z)} = \frac{\frac{1}{K^n} \sum_{\mathbf{k}} r_{\mathbf{k}}(z)}{\frac{1}{K^n} \sum_{\mathbf{k}'} r_{\mathbf{k}'}(z)} = 1 \quad (28)$$

As such, we can in principle use methods for discrete graphical models, treating \mathbf{k} as a random variable. However, as discussed in the Related work, this treating the problem as inference in a discrete graphical model is still prohibitively difficult.

4.2 Computing expectations by differentiating an estimate of the normalizing constant

Instead, inspired by Sec. 3, we modify our marginal likelihood estimator by introducing a source term,

$$\mathcal{P}_{\text{MP}}^{\text{exp}}(z, J) = \frac{1}{K^n} \sum_{\mathbf{k} \in \mathcal{K}^n} r_{\mathbf{k}}(z) e^{Jm(z^{\mathbf{k}})}. \quad (29)$$

Remember that $r_{\mathbf{k}}(z)$ is a product of low-rank tensors, indexed by subsets of \mathbf{k} (Eq. 23), so the sum can be computed efficiently using opt-einsum. Critically, the source term is just another factor with indices given by a subset of \mathbf{k} . For instance, most often m (the function whose expectation we want to compute) will depend on only a single latent variable $m(z^{\mathbf{k}}) = m(z_i^{k_i})$, in which case the source term can be understood as just another tensor in the tensor product, with one index, k_i ,

$$f_{k_i}^m = e^{Jm(z^{\mathbf{k}})}. \quad (30)$$

Again, we have that $\mathcal{P}_{\text{MP}}^{\text{exp}}(z, J = 0) = \mathcal{P}_{\text{MP}}(z)$. Inspired by Sec. 3, we now differentiate $\log \mathcal{P}_{\text{MP}}^{\text{exp}}(z, J)$ at $J = 0$. Specifically,

$$\left. \frac{\partial}{\partial J} \right|_{J=0} \log \mathcal{P}_{\text{MP}}^{\text{exp}}(z, J) = \frac{\left. \frac{\partial}{\partial J} \right|_{J=0} \mathcal{P}_{\text{MP}}^{\text{exp}}(z, J)}{\mathcal{P}_{\text{MP}}^{\text{exp}}(z, 0)} \quad (31)$$

Substituting for $\mathcal{P}_{\text{MP}}^{\text{exp}}(z, J)$ (Eq. 29), and remembering that $\mathcal{P}_{\text{MP}}^{\text{exp}}(z, 0) = \mathcal{P}_{\text{MP}}(z)$,

$$\left. \frac{\partial}{\partial J} \right|_{J=0} \log \mathcal{P}_{\text{MP}}^{\text{exp}}(z, J) = \frac{\frac{1}{K^n} \sum_{\mathbf{k}} r_{\mathbf{k}}(z) \left. \frac{\partial}{\partial J} \right|_{J=0} e^{Jm(z^{\mathbf{k}})}}{\mathcal{P}_{\text{MP}}(z)} \quad (32)$$

Computing the gradient of $e^{Jm(z^{\mathbf{k}})}$ at $J = 0$,

$$\begin{aligned} \left. \frac{\partial}{\partial J} \right|_{J=0} \log \mathcal{P}_{\text{MP}}^{\text{exp}}(z, J) &= \frac{\frac{1}{K^n} \sum_{\mathbf{k}} r_{\mathbf{k}}(z) m(z^{\mathbf{k}})}{\mathcal{P}_{\text{MP}}(z)} \\ &= m_{\text{MP}}(z) \end{aligned} \quad (33)$$

where the final equality comes from the definition of $m_{\text{MP}}(z)$ in Eq. (27). Note that this derivation is quite different from the standard “source-term trick” from Physics described in Sec. 3, which works with either the true normalizing constant, or with a low-order perturbation to that normalizing constant. In contrast, here we use a quite different massively parallel sample-based estimate of the marginal likelihood. Importantly, the subsequent two derivations are even more different from uses of the “source-term trick” in Physics, as these uses are typically around computing a moment/expectation, while the subsequent two derivations use the same trick to compute quite different quantities (namely, probability distributions over samples).

4.3 Computing marginal importance weights

Computing expectations directly is very powerful and almost certainly necessary for computing complex quantities that depend on multiple latent variables. However, if we are primarily interested in expectations of individual variables, then it is considerably more flexible to compute “marginal” importance weights. Once we have these marginal importance weights, we can easily compute arbitrary expectations for individual variables, and other quantities such as effective sample sizes. To define the marginal weights for the i th latent, note that a moment for the i th latent variable can be written as a sum over k_i ,

$$m_{\text{MP}}(z) = \sum_{\mathbf{k} \in \mathcal{K}^n} \frac{r_{\mathbf{k}}(z)}{\mathcal{P}_{\text{MP}}(z)} m(z_i^{k_i}) = \sum_{k_i} w_{k_i}^i m(z_i^{k_i}), \quad (34)$$

where $w_{k_i}^i$ are the marginal importance weights for the i th latent variable, which are defined by,

$$w_{k_i}^i = \frac{\frac{1}{K^n} \sum_{\mathbf{k}/k_i \in \mathcal{K}^{n-1}} r_{\mathbf{k}}(z)}{\mathcal{P}_{\text{MP}}(z)}, \quad (35)$$

where the sum is over all \mathbf{k} except k_i . Formally,

$$\mathbf{k}/k_i = (k_1, \dots, k_{i-1}, k_{i+1}, \dots, k_n) \in \mathcal{K}^{n-1}. \quad (36)$$

Again we can compute the marginal importance weights using gradients of a slightly different modified marginal likelihood estimator. Specifically, we now use a vector-valued $\mathbf{J} \in \mathbb{R}^K$ in a slightly different modified marginal likelihood estimator,

$$\mathcal{P}_{\text{MP}}^{\text{marg}}(z, \mathbf{J}) = \frac{1}{K^n} \sum_{\mathbf{k}} r_{\mathbf{k}}(z) e^{J_{k_i}}. \quad (37)$$

Again, $\mathcal{P}_{\text{MP}}^{\text{marg}}(z, \mathbf{0}) = \mathcal{P}_{\text{MP}}(z)$. As before, we differentiate $\log \mathcal{P}_{\text{MP}}^{\text{marg}}(z, \mathbf{J})$ at $\mathbf{J} = \mathbf{0}$,

$$\left. \frac{\partial}{\partial J_{k'_i}} \right|_{\mathbf{J}=\mathbf{0}} \log \mathcal{P}_{\text{MP}}^{\text{marg}}(z, \mathbf{J}) = \frac{\left. \frac{\partial}{\partial J_{k'_i}} \right|_{\mathbf{J}=\mathbf{0}} \mathcal{P}_{\text{MP}}^{\text{marg}}(z, \mathbf{J})}{\mathcal{P}_{\text{MP}}^{\text{marg}}(z, \mathbf{0})}. \quad (38)$$

Substituting for $\mathcal{P}_{\text{MP}}^{\text{marg}}(z, \mathbf{J})$ in the numerator,

$$\left. \frac{\partial}{\partial J_{k'_i}} \right|_{\mathbf{J}=\mathbf{0}} \log \mathcal{P}_{\text{MP}}^{\text{marg}}(z, \mathbf{J}) = \frac{\frac{1}{K^n} \sum_{\mathbf{k}} r_{\mathbf{k}}(z) \left. \frac{\partial}{\partial J_{k'_i}} \right|_{\mathbf{J}=\mathbf{0}} e^{J_{k_i}}}{\mathcal{P}_{\text{MP}}(z)}. \quad (39)$$

The gradient is 1 when $k'_i = k_i$ and zero otherwise which can be represented using a Kronecker delta,

$$\left. \frac{\partial}{\partial J_{k'_i}} \right|_{\mathbf{J}=\mathbf{0}} \log \mathcal{P}_{\text{MP}}^{\text{marg}}(z, \mathbf{J}) = \frac{\frac{1}{K^n} \sum_{\mathbf{k}} r_{\mathbf{k}}(z) \delta_{k'_i, k_i}}{\mathcal{P}_{\text{MP}}(z)}. \quad (40)$$

We can rewrite this as a sum over all \mathbf{k} except k_i ,

$$\begin{aligned} \left. \frac{\partial}{\partial J_{k'_i}} \right|_{\mathbf{J}=\mathbf{0}} \log \mathcal{P}_{\text{MP}}^{\text{marg}}(z, \mathbf{J}) &= \frac{\frac{1}{K^n} \sum_{\mathbf{k}/k_i \in \mathcal{K}^{n-1}} r_{\mathbf{k}}(z)}{\mathcal{P}_{\text{MP}}(z)} \\ &= w_{k_i}^i, \end{aligned} \quad (41)$$

which is exactly the definition of the marginal importance weights in Eq. (35).

4.4 Computing conditional distributions for importance sampling

A common alternative to importance weighting is importance sampling. In importance sampling, we rewrite the usual estimates of the expectations in terms of a distribution over indices, $P(\mathbf{k})$,

$$m_{\text{MP}}(z) = \sum_{\mathbf{k} \in \mathcal{K}^n} P(\mathbf{k}) m(z^{\mathbf{k}}) \quad (42)$$

$$P(\mathbf{k}) = \frac{1}{K^n} \frac{1}{\mathcal{P}_{\text{MP}}(z)} r_{\mathbf{k}}(z) \quad (43)$$

Using \mathbf{k} samples from this distribution, $z^{\mathbf{k}}$ are approximate samples from the true posterior. However, sampling from $P(\mathbf{k})$ is difficult in our context, as there are K^n possible settings of \mathbf{k} , so we cannot explicitly compute the full distribution. Instead, we need to factorise the distribution in some way, and iteratively sample (e.g. we sample k_1 from $P(k_1)$ then sample k_2 from $P(k_2|k_1)$ etc.) However, this raises a question: how should we factorise the distribution over \mathbf{k} ? This is a difficult problem in the probabilistic programming setting, because we need a factorisation that is always valid, and at the same time, has as few indices in each term as possible, to ensure that the computations remain efficient. By substituting Eq. (23) into Eq. (43),

$$P(\mathbf{k}) = \frac{1}{K^n} \frac{1}{\mathcal{P}_{\text{MP}}(z)} f_{\mathbf{k}_{\text{pa}(x)}}^x(z) \prod_i f_{k_i, \mathbf{k}_{\text{pa}(i)}}^i(z) \quad (44)$$

we can see that one valid factorisation follows the factorisation of the generative model. This is particularly useful, as it is guaranteed to be a valid factorisation, likely to be small (if not minimal) and it is easy for us to extract. Formally, we use,

$$P(\mathbf{k}) = \prod_i P(k_i | \mathbf{k}_{\text{pa}(i)}) \quad (45)$$

where, remember $\text{pa}(i)$ is the set of indices of parents of the i th latent variable under the generative model, so,

$$\mathbf{k}_{\text{pa}(i)} = (k_j \text{ for all } j \in \text{pa}(i)) \quad (46)$$

Now, we have the problem of computing the conditionals, $P(k_i | \mathbf{k}_{\text{pa}(i)})$. We can compute the conditionals from the marginals using Bayes theorem,

$$P(k_i | \mathbf{k}_{\text{pa}(i)}) = \frac{P(k_i, \mathbf{k}_{\text{pa}(i)})}{\sum_{k'_i} P(k'_i, \mathbf{k}_{\text{pa}(i)})} \quad (47)$$

where the marginals are given by,

$$P(k_i, \mathbf{k}_{\text{pa}(i)}) = \sum_{\mathbf{k} / (k_i, \mathbf{k}_{\text{pa}(i)})} P(\mathbf{k}) \quad (48)$$

Again, we can compute these marginals efficiently by differentiating a modified estimate of the marginal likelihood. This time, we take a tensor-valued $\mathbf{J} \in \mathbb{R}^{K^{1+|\text{pa}(i)|}}$, where remember $|\text{pa}(i)|$ is the number of parents of the i th latent variable under the generative model.

$$\mathcal{P}_{\text{MP}}^{\text{samp}}(z, \mathbf{J}) = \frac{1}{K^n} \sum_{\mathbf{k}} \frac{P(x, z^{\mathbf{k}})}{Q(z^{\mathbf{k}})} e^{J_{k_i, \mathbf{k}_{\text{pa}(i)}}} \quad (49)$$

As usual, we differentiate with respect to \mathbf{J} at $\mathbf{J} = \mathbf{0}$,

$$\left. \frac{\partial}{\partial J_{k'_i, \mathbf{k}'_{\text{pa}(i)}}} \right|_{\mathbf{J}=\mathbf{0}} \log \mathcal{P}_{\text{MP}}^{\text{samp}}(z, \mathbf{J}) = \frac{\frac{1}{K^n} \sum_{\mathbf{k}} r_{\mathbf{k}}(z) \delta_{(k_i, \mathbf{k}_{\text{pa}(i)}), (k'_i, \mathbf{k}'_{\text{pa}(i)})}}{\mathcal{P}_{\text{MP}}(z)} \quad (50)$$

Here, $\delta_{(k_i, \mathbf{k}_{\text{pa}(i)}), (k'_i, \mathbf{k}'_{\text{pa}(i)})}$ is a generalisation of the Kronecker delta. It is 1 when all the indices match (i.e. $k_i = k'_i$, and $\mathbf{k}_{\text{pa}(i)} = \mathbf{k}'_{\text{pa}(i)}$) and zero otherwise.

$$\begin{aligned} \left. \frac{\partial}{\partial J_{k_i, \mathbf{k}_{\text{pa}(i)}}} \right|_{\mathbf{J}=\mathbf{0}} \log \mathcal{P}_{\text{MP}}^{\text{samp}}(z, \mathbf{J}) &= \frac{\frac{1}{K^n} \sum_{\mathbf{k}/(k_i, \mathbf{k}_{\text{pa}(i)})} r_{\mathbf{k}}(z)}{\mathcal{P}_{\text{MP}}(z)} \\ &= \sum_{\mathbf{k}/(k_i, \mathbf{k}_{\text{pa}(i)})} P(\mathbf{k}) \\ &= P(k_i, \mathbf{k}_{\text{pa}(i)}) \end{aligned} \quad (51)$$

These are precisely the marginals in Eq. (48).

5 Experiments

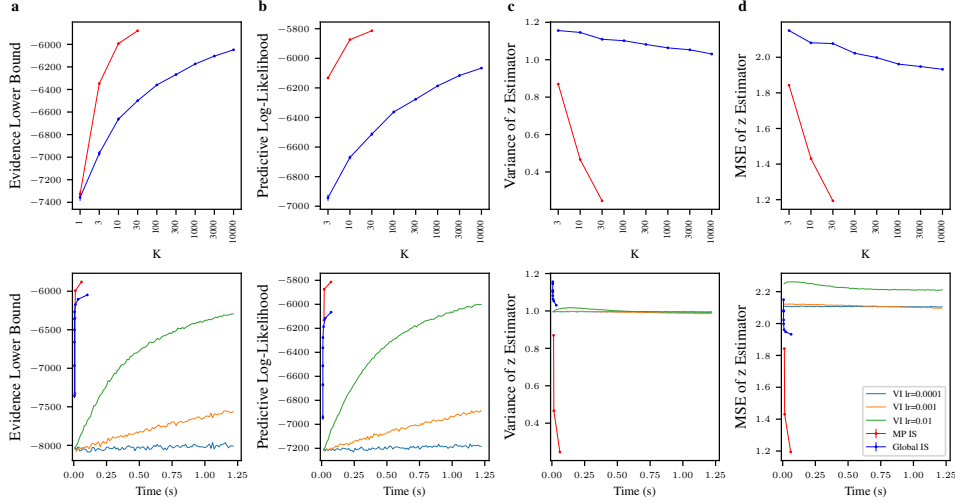


Figure 1: Results obtained in the MovieLens model. Columns **a–c** show the evidence lower bound, predictive log-likelihood and variance in the estimator of \mathbf{z}_m using the true MovieLens100K data. Column **d** shows the mean squared error in the estimator of \mathbf{z}_m when the data is sampled from the model and thus the true value of \mathbf{z}_m is known.

We provide empirical results comparing the global and massively parallel importance weighting/sampling methods. We consider four quantities in Fig. 1 and Fig. 2. First, we compute the ELBO that arises just from the forward pass. While the ELBO can be computed using methods in [Anonymous, 2023], and does not require the contributions in this paper, it is a useful sanity check. Second, we consider the predictive log-likelihood, computed using posterior samples of the latent variables. To obtain these posterior samples, we use the methods in Sec. 4.4. Finally, we consider the quality of the moment estimates (in particular an estimate of the posterior mean of one of the latent variables), where those moments are computed using the methods in Sec. 4.3. In particular, we use two measures of the quality of the estimates, based on real data and data sampled from the generative model. On real data, all we can do is to look at the variance of our estimator of the posterior mean (Fig. 1c, Fig. 2c). In contrast, if we sample from the generative model, we know the true value of the latent variable used to sample the observations. We can therefore compute the MSE between our estimate of the posterior mean and the true value of the latent variable used to generate the data. However, we must be careful when interpreting this quantity, as the true posterior mean is not necessarily equal to the value of the latent variable used to generate the data (they only become equal in the infinite data limit).

In order to compare the two methods in the simplest possible case, we use the prior as the proposal in all models. The values of interest (evidence lower bound, predictive log-likelihood and posterior expectation estimates both with real and generated data) were calculated 1000 times using different

importance weights, generated with both the global and the massively parallel scheme, and averaged to obtain the results presented in figures 1 and 2.

5.1 MovieLens Dataset

For our first experiment we use the MovieLens100K Harper and Konstan [2015] dataset, containing 100K ratings of $N = 1682$ films from among $M = 943$ users. In our experiments these ratings are binarised from $(0, 1, 2, 3)$ to 0 and from $(4, 5)$ to 1, as in Geffner and Domke [2022]. We use n to index films, and m to index users. Each film has a feature vector \mathbf{x}_n and we use the following hierarchical model, similar to that in [Anonymous, 2023]:

$$\begin{aligned}\boldsymbol{\mu} &\sim \mathcal{N}(\mathbf{0}_{18}, 0.25\mathbf{I}) \\ \boldsymbol{\psi} &\sim \mathcal{N}(\mathbf{0}_{18}, 0.25\mathbf{I}) \\ \mathbf{z}_m &\sim \mathcal{N}(\boldsymbol{\mu}, \exp(\boldsymbol{\psi})\mathbf{I}), m = 1, \dots, M \\ \text{Rating}_{mn} &\sim \text{Bernoulli}(\sigma(\mathbf{z}_m^\top \mathbf{x}_n)), n = 1, \dots, N\end{aligned}\tag{52}$$

After first sampling a global mean, $\boldsymbol{\mu}$, and variances, $\boldsymbol{\psi}$, the model samples a latent vector, \mathbf{z}_m , for each user, corresponding to the features of films \mathbf{x}_n they are likely to rate highly. We take the dot-product of the latent user-vector, \mathbf{z}_m , and the film’s feature vector, \mathbf{x}_n , to obtain the probability of a positive rating for film n from user m . A corresponding graphical model is given in Appendix B.1.

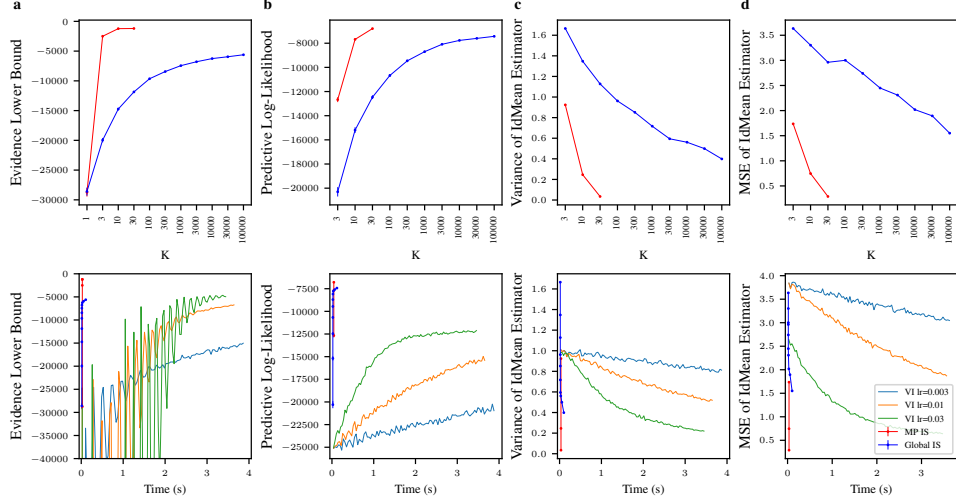


Figure 2: Results obtained in the NYC Bus Breakdown model. Columns **a–c** show the evidence lower bound, predictive log-likelihood and variance in the estimator of IdMean_{mj} using the true data. Column **d** shows the mean squared error in the estimator of IdMean_{mj} when the data is sampled from the model.

We use a random subset of $N = 20$ films and $M = 450$ users for our experiment, with an equally sized but disjoint subset held aside for calculation of the predictive log-likelihood. We ensure high levels of uncertainty in the per user mean \mathbf{z}_m by using far fewer films than users, and assume ratings of 0 for films which users have not previously rated. In order to obtain results with comparable computation times we test the massively parallel importance samples for $K \in \{3, 10, 30\}$ and the global importance samples for $K \in \{3, 10, 30, 100, 300, 1000, 3000, 10000\}$.

Our results (Fig. 1) show tighter evidence lower bounds and higher predictive log-likelihoods using massively parallel importance samples than using global importance samples for a given $K > 1$, with superior results obtained in less time once K is large enough in the massively parallel approach (in this case, for $K \geq 10$). We also see much lower variance in our importance weighted estimates for the expectation of per user means, \mathbf{z}_m , for every $K > 1$ and again see this reduced variance achieved in less time than the global method for large enough K (here, $K \geq 3$). Using data sampled from the model in order to obtain ground truth values of \mathbf{z}_m , we found the mean squared error of the

\mathbf{z}_m expectation estimates was far lower for all $K > 1$ when those estimates were calculated using massively parallel importance weights rather than global importance weights, with similar behaviour in the time required to reduce the MSE as with reducing the variance.

5.2 NYC Bus Breakdown Dataset

In our second experiment, we model the length of delay time of New York school bus journeys, working with a dataset supplied by the City of New York DOE [2023], based on the school year, borough, and the ID of the bus in question. These three levels motivate a hierarchical model similar to that in Anonymous [2023] — a mean and variance is sampled for each year, which are used to sample a borough mean which, along with a sampled borough variance, is used to sample an ID mean for each year and borough. It is these ID means, denoted by IdMean_{mj} for year m and borough j , over IDs $i = 1, \dots, I$, for which we estimate the posterior expectation. After this, the model samples a final variance that is used to sample two weight vectors which are multiplied with two vectors of covariates indicating the bus company and journey type. Finally, the predicted delay is sampled from a negative binomial distribution with logits given by the sum of these two multiplied terms and the ID mean. The full specification of the model can be found in Appendix B.2, with a corresponding graphical model provided in Appendix B.3.

The experiments were run using a subset of the data with $I = 30$ IDs from $J = 3$ boroughs in $M = 3$ years. As with the previous experiment, an equally sized but disjoint subset of the data is used to calculate predictive log-likelihoods. We test the massively parallel importance samples for $K \in \{3, 10, 30\}$ and the global importance samples for $K \in \{3, 10, 30, 100, 300, 1000, 3000, 10000, 30000, 100000\}$.

We run the same experiments as on the MovieLens model, but using posterior expectation estimates for IdMean_{mj} , as displayed in Fig. 2. We again observe improved performance from massively parallel importance sampling and weighting compared to global importance sampling and weighting.

6 Conclusion

We have shown how posterior moments, marginals and samples may be computed using massively parallel importance sampling/weighting methods by drawing K samples for n latent variables, and individually reasoning about all K^n combinations. We gave a new and far simpler method for computing these quantities based on differentiating a slightly modified marginal likelihood estimator. These approaches led to better importance samples (as measured by predictive log-likelihoods) and better importance weighted moment estimates.

References

- Laurence Aitchison. Tensor Monte Carlo: particle methods for the GPU era. *Advances in Neural Information Processing Systems*, 32, 2019.
- Aws Albarghouthi, Loris D’Antoni, Samuel Drews, and Aditya V Nori. Fairsquare: probabilistic verification of program fairness. *Proceedings of the ACM on Programming Languages*, 1 (OOPSLA):1–30, 2017.
- Christophe Andrieu, Arnaud Doucet, and Roman Holenstein. Particle markov chain monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3): 269–342, 2010.
- Anonymous. Massively parallel reweighted wake-sleep. *Submitted to UAI 2023, and included in the Appendix*, 2023.
- Bozhena Bidyuk and Rina Dechter. Cutset sampling for bayesian networks. *Journal of Artificial Intelligence Research*, 28:1–48, 2007.
- Jörg Bornschein and Yoshua Bengio. Reweighted wake-sleep. *arXiv preprint arXiv:1406.2751*, 2014.
- Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*, 2015.

- Sourav Chatterjee and Persi Diaconis. The sample size required in importance sampling. *The Annals of Applied Probability*, 28(2):1099–1135, 2018.
- Guillaume Claret, Sriram K Rajamani, Aditya V Nori, Andrew D Gordon, and Johannes Borgström. Bayesian inference using data flow analysis. In *Proceedings of the 2013 9th joint meeting on foundations of software engineering*, pages 92–102, 2013.
- Adrien Corenflos, Nicolas Chopin, and Simo Särkkä. De-sequentialized monte carlo: a parallel-in-time particle smoother. *arXiv preprint arXiv:2202.02264*, 2022.
- A Philip Dawid. Applications of a general propagation algorithm for probabilistic expert systems. *Statistics and computing*, 2(1):25–36, 1992.
- DOE. Bus breakdown and delays, 2023. url <https://data.cityofnewyork.us/Transportation/Bus-Breakdown-and-Delays/ez4e-fazm> Accessed on: 05.05.2023.
- Arnaud Doucet, Adam M Johansen, et al. A tutorial on particle filtering and smoothing: Fifteen years later. *Handbook of nonlinear filtering*, 12(656-704):3, 2009.
- Tomas Geffner and Justin Domke. Variational inference with locally enhanced bounds for hierarchical models. *arXiv preprint arXiv:2203.04432*, 2022.
- Timon Gehr, Sasa Misailovic, and Martin Vechev. Psi: Exact symbolic inference for probabilistic programs. In *International Conference on Computer Aided Verification*, pages 62–83. Springer, 2016.
- Jaco Geldenhuys, Matthew B Dwyer, and Willem Visser. Probabilistic symbolic execution. In *Proceedings of the 2012 International Symposium on Software Testing and Analysis*, pages 166–176, 2012.
- Noah D Goodman and Andreas Stuhlmüller. The design and implementation of probabilistic programming languages, 2014.
- Neil J Gordon, David J Salmond, and Adrian FM Smith. Novel approach to nonlinear/non-gaussian bayesian state estimation. In *IEE proceedings F (radar and signal processing)*, volume 140, pages 107–113. IET, 1993.
- F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)*, 5(4):1–19, 2015.
- Steven Holtzen, Guy Van den Broeck, and Todd Millstein. Scaling exact inference for discrete probabilistic programs. *Proceedings of the ACM on Programming Languages*, 4(OOPSLA):1–31, 2020.
- Jinlin Lai, Justin Domke, and Daniel Sheldon. Variational marginal particle filters. In *International Conference on Artificial Intelligence and Statistics*, pages 875–895. PMLR, 2022.
- Tuan Anh Le, Maximilian Igl, Tom Rainforth, Tom Jin, and Frank Wood. Auto-encoding sequential monte carlo. *arXiv preprint arXiv:1705.10306*, 2017.
- Tuan Anh Le, Adam R Kosiorek, N Siddharth, Yee Whye Teh, and Frank Wood. Revisiting re-weighted wake-sleep for models with stochastic control flow. In *Uncertainty in Artificial Intelligence*, pages 1039–1049. PMLR, 2020.
- Fredrik Lindsten, Adam M Johansen, Christian A Naesseth, Bonnie Kirkpatrick, Thomas B Schön, JAD Aston, and Alexandre Bouchard-Côté. Divide-and-conquer with sequential monte carlo. *Journal of Computational and Graphical Statistics*, 26(2):445–458, 2017.
- Chris J Maddison, John Lawson, George Tucker, Nicolas Heess, Mohammad Norouzi, Andriy Mnih, Arnaud Doucet, and Yee Teh. Filtering variational objectives. *Advances in Neural Information Processing Systems*, 30, 2017.
- Christian Naesseth, Scott Linderman, Rajesh Ranganath, and David Blei. Variational sequential monte carlo. In *International conference on artificial intelligence and statistics*, pages 968–977. PMLR, 2018.

- Praveen Narayanan, Jacques Carette, Wren Romano, Chung-chieh Shan, and Robert Zinkov. Probabilistic inference by program transformation in hakaru (system description). In *International Symposium on Functional and Logic Programming*, pages 62–79. Springer, 2016.
- Fritz Obermeyer, Eli Bingham, Martin Jankowiak, Neeraj Pradhan, Justin Chiu, Alexander Rush, and Noah Goodman. Tensor variable elimination for plated factor graphs. In *International Conference on Machine Learning*, pages 4871–4880. PMLR, 2019.
- Avi Pfeffer. The design and implementation of ibal: A generalpurpose probabilistic programming language. In *Harvard Univesity*. Citeseer, 2005.
- Sriram Sankaranarayanan, Aleksandar Chakarov, and Sumit Gulwani. Static analysis for probabilistic programs: inferring whole program properties from finitely many paths. In *Proceedings of the 34th ACM SIGPLAN conference on Programming language design and implementation*, pages 447–458, 2013.
- Di Wang, Jan Hoffmann, and Thomas Reps. Pmaf: an algebraic framework for static analysis of probabilistic programs. *ACM SIGPLAN Notices*, 53(4):513–528, 2018.
- Steven Weinberg. *The quantum theory of fields*, volume 2. Cambridge university press, 1995.
- Jacob Zavatore-Veth, Abdulkadir Canatar, Ben Ruben, and Cengiz Pehlevan. Asymptotics of representation learning in finite bayesian neural networks. *Advances in neural information processing systems*, 34:24765–24777, 2021.

A Derivations

A.1 Global Importance Sampling

Here, we give the derivation for standard global importance sampling. Ideally we would compute moments using the true posterior, $P(z|x)$,

$$m_{\text{post}} = E_{P(z^k|x)} [m(z^k)]. \quad (53)$$

However, the true posterior is not known. Instead, we write down the moment under the true posterior as an integral,

$$m_{\text{post}} = \int dz^k P(z^k|x) m(z^k). \quad (54)$$

Next, we multiply the integrand by $1 = Q(z^k)/Q(z^k)$,

$$m_{\text{post}} = \int dz^k Q(z^k) \frac{P(z^k|x)}{Q(z^k)} m(z^k). \quad (55)$$

Next, the integral can be written as an expectation,

$$m_{\text{post}} = E_{Q(z^k)} \left[\frac{P(z^k|x)}{Q(z^k)} m(z^k) \right]. \quad (56)$$

It looks like we should be able to estimate m_{post} by sampling from our approximate posterior, $Q(z^k)$. However, this is not yet possible, as we are not able to compute the true posterior, $P(z^k|x)$. We might consider using Bayes theorem,

$$P(z^k|x) = \frac{P(z^k, x)}{P(x)} \quad (57)$$

But this requires computing an intractable normalizing constant,

$$P(x) = \int dz^k P(z^k, x). \quad (58)$$

Instead, we use an unbiased, importance-sampled estimate of the normalizing constant, $\mathcal{P}_{\text{global}}(z)$ (Eq. 7). Additionally, Burda et al. [2015] showed that in the limit as $K \rightarrow \infty$, $\mathcal{P}_{\text{global}}(z)$ approaches $P(x)$. Using this estimate of the marginal likelihood our moment estimate becomes,

$$m_{\text{global}} = E_{Q(z^k)} \left[\frac{\frac{P(z^k, x)}{Q(z^k)}}{\mathcal{P}_{\text{global}}(z)} m(z^k) \right]. \quad (59)$$

Using $r_k(z)$ (Eq. 6), we can write this expression as,

$$m_{\text{global}} = E_{Q_\phi(z)} \left[\frac{r_k(z)}{\mathcal{P}_{\text{global}}(z)} m(z^k) \right]. \quad (60)$$

The approximate posterior and generative probabilities are the same for different values of k , so we can average over k , which gives Eq. (5) in the main text.

A.2 Massively Parallel Importance Sampling

Inspired by the global importance sampling derivation, we consider massively parallel importance sampling. In the global importance sampling derivation, the key idea was to show that the estimator was unbiased for each of the K samples, z^k , in which case the average over all K samples is also unbiased. In massively parallel importance sampling, we use the same idea, except that we now have K^n samples, denoted $z^{\mathbf{k}}$. As before,

$$m_{\text{post}} = E_{P(z^{\mathbf{k}}|x)} [m(z^{\mathbf{k}})] = \int dz^{\mathbf{k}} P(z^{\mathbf{k}}|x) m(z^{\mathbf{k}}). \quad (61)$$

Again, we multiply by $1 = \prod_i Q(z_i^{k_i} | z_{\text{qa}(i)}) / \prod_i Q(z_i^{k_i} | z_{\text{qa}(i)})$,

$$m_{\text{post}} = \int dz^{\mathbf{k}} \left(\prod_i Q(z_i^{k_i} | z_{\text{qa}(i)}) \right) \frac{P(z^{\mathbf{k}} | x)}{\prod_i Q(z_i^{k_i} | z_{\text{qa}(i)})} m(z^{\mathbf{k}}). \quad (62)$$

Overall, our goal is to convert the integral over the indexed latent variables in Eq. (62) into an integral over the full latent space, z , so that it can be written as an expectation over the proposal, $Q(z)$. To do that, we need to introduce the concept of non-indexed latent variables. These are all samples of the latent variables, except for the “indexed”, or k th sample. For the i th latent variable, the non-indexed samples are,

$$z_i^{/k_i} = (z_i^1, \dots, z_i^{k_i-1}, z_i^{k_i+1}, \dots, z_i^K) \in \mathcal{Z}_i^{K-1}. \quad (63)$$

We can also succinctly write the non-indexed samples of all latent variables as,

$$z^{/k} = (z_1^{/k_1}, z_2^{/k_2}, \dots, z_n^{/k_n}) \in \mathcal{Z}^{K-1}. \quad (64)$$

The joint distribution over the non-indexed latent variables, conditioned on the indexed latent variables integrates to 1,

$$1 = \int dz^{/k} Q(z^{/k} | z^{\mathbf{k}}) = \int dz^{/k} \prod_i Q(z_i^{/k_i} | z_i^{k_i}, z_{\text{qa}(i)}), \quad (65)$$

We use this to multiply the integrand in Eq. (62),

$$m_{\text{post}} = \int dz^{\mathbf{k}} \left(\prod_i Q(z_i^{k_i} | z_{\text{qa}(i)}) \right) \frac{P(z^{\mathbf{k}} | x)}{\prod_i Q(z_i^{k_i} | z_{\text{qa}(i)})} m(z^{\mathbf{k}}) \int dz^{/k} \prod_i Q(z_i^{/k_i} | z_i^{k_i}, z_{\text{qa}(i)}). \quad (66)$$

Next, we merge the integrals over $z^{\mathbf{k}}$ and $z^{/k}$ to form one integral over z ,

$$m_{\text{post}} = \int dz Q(z) \frac{P(z^{\mathbf{k}} | x)}{\prod_i Q(z_i^{k_i} | z_{\text{qa}(i)})} m(z^{\mathbf{k}}). \quad (67)$$

This integral can be written as an expectation,

$$m_{\text{post}} = E_{Q(z)} \left[\frac{P(z^{\mathbf{k}} | x)}{\prod_i Q(z_i^{k_i} | z_{\text{pa}(i)})} m(z^{\mathbf{k}}) \right]. \quad (68)$$

As in the derivation for global importance sampling, it looks like we might be able to estimate this by sampling from $Q(z|x)$, but this does not yet work as we do not yet have a form for the posterior. Again, we could compute the posterior using Bayes theorem,

$$P(z^{\mathbf{k}} | x) = \frac{P(z^{\mathbf{k}}, x)}{P(x)}, \quad (69)$$

but we cannot compute the model evidence,

$$P_{\theta}(x) = \int dz^{\mathbf{k}} P_{\theta}(z^{\mathbf{k}}, x). \quad (70)$$

As in the global importance sampling section, we instead use an estimate of the marginal likelihood. Here, we use a massively parallel estimate, $\mathcal{P}_{\text{MP}}(z)$,

$$m_{\text{MP}} = E_{Q(z|x)} \left[\frac{\frac{P(z^{\mathbf{k}}, x)}{\prod_i Q(z_i^{k_i} | z_{\text{pa}(i)})}}{\mathcal{P}_{\text{MP}}(z)} m(z^{\mathbf{k}}) \right]. \quad (71)$$

Again, we use $r_{\mathbf{k}}(z)$ (Eq. 21),

$$m_{\text{MP}} = E_{Q(z)} \left[\frac{r_{\mathbf{k}}(z)}{\mathcal{P}_{\text{MP}}(z)} m(z^{\mathbf{k}}) \right]. \quad (72)$$

So the value for a single set of latent variables, $z^{\mathbf{k}}$, has the right expectation. Thus, averaging over all K^n settings of \mathbf{k} , we get the unbiased estimator in the main text, (Eq. 27).

B Experiment Models

B.1 MovieLens Graphical Model

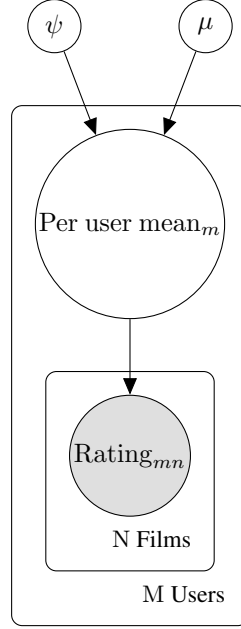


Figure 3: Graphical model for the MovieLens dataset

B.2 Bus Delay Model Specification

$$\begin{aligned}
 \text{YearVariance} &\sim \mathcal{N}(0, 10^{-4}) \\
 \text{YearMean} &\sim \mathcal{N}(0, 10^{-4}) \\
 \text{BoroughMean}_m &\sim \mathcal{N}(\text{YearMean}, \exp(\text{YearVariance})), \quad m = 1, \dots, M \\
 \text{BoroughVariance}_j &\sim \mathcal{N}(0, 0.25), \quad j = 1, \dots, J \\
 \text{IdMean}_{mj} &\sim \mathcal{N}(\text{BoroughMean}_m, \text{BoroughVariance}_j), \quad j = 1, \dots, J \\
 \text{WeightVariance}_i &\sim \mathcal{N}(0, 10^{-4}), \quad i = 1, \dots, I \\
 \mathbf{C}_i &\sim \mathcal{N}(\mathbf{0}_{\# \text{BusCo.s}}, \exp(\text{WeightVariance}_i)), \quad i = 1, \dots, I \\
 \mathbf{J}_i &\sim \mathcal{N}(\mathbf{0}_{\# \text{JourneyTypes}}, \exp(\text{WeightVariance}_i)), \quad i = 1, \dots, I \\
 \text{logits}_{mji} &= \text{IdMean}_{mj} + \mathbf{C}_i * \text{Bus company name}_{mji} + \mathbf{J}_i * \text{Journey type}_{mji} \\
 \text{Delay}_{mji} &\sim \text{NegativeBinomial}(\text{total count} = 130, \text{logits}_{mji}), \quad i = 1, \dots, I
 \end{aligned} \tag{73}$$

Here Bus company name_{mji} and Journey type_{mji} are one-hot encoded indicator variables indicating the bus company and the type of bus journey respectively. The dataset's largest recorded delay is 130, hence we use total count = 130.

B.3 NYC Bus Breakdown Graphical Model

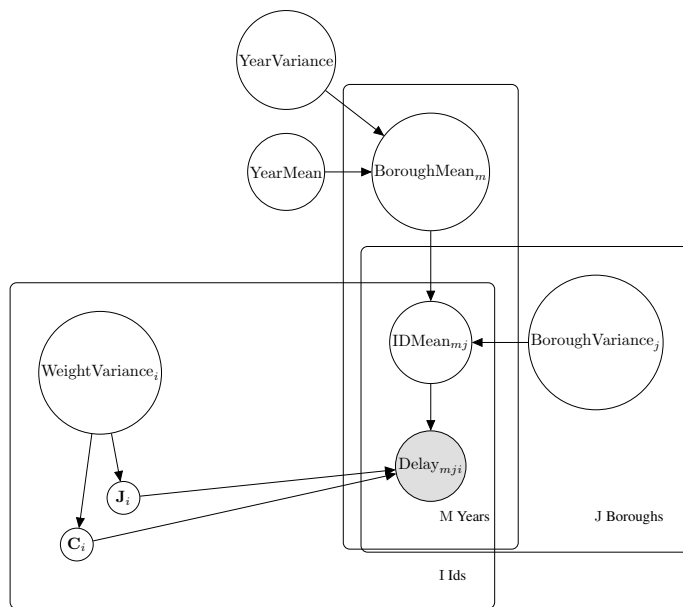


Figure 4: Graphical model for the NYC Bus Breakdown dataset