

# Tensorised probabilistic programming

laurence.aitchison

June 2020

## Abstract

Tensor Monte Carlo introduced the possibility of dramatically improving the tightness of the evidence lower bound objective by sampling each latent variable  $n$  times and considering all  $K^n$  combinations as importance samples. However, their particular approach to augmenting the state-space was inspired by particle filters, and was poorly suited to the full probabilistic programming setting. In particular, In particular ...

Examples:

- Can we compare speed to Justin’s differentiable particle filter thingy?
- Simple Heirarchical model.
- MNIST?

## 1 Methods

We have an underlying model with data,  $x$ , and  $n$  latent variables denoted  $(z'_1, \dots, z'_n)$ . Each latent variable,  $z'_i$  lives in a set,  $z'_i \in \mathcal{Z}_i$ . The full collection of latent variables is denoted  $z'$ ,

$$z' = (z'_1, \dots, z'_n) \in \mathcal{Z}_1 \times \dots \times \mathcal{Z}_n = \mathcal{Z}. \quad (1)$$

The full probability density includes one factor for the probability density of the data, conditioned on the latents,  $\pi_x(x; z'_{\text{pa}(x)})$  and one factor for each latent variable,  $\pi_i(z'_i; z'_{\text{pa}(i)})$ ,

$$\pi(x, z') = \pi_x(x; z'_{\text{pa}(x)}) \prod_{i=1}^n \pi_i(z'_i; z'_{\text{pa}(i)}). \quad (2)$$

where  $\text{pa}(i)$  is the indicies of all parents of  $z_i$ , and  $z'_{\text{pa}(i)}$  is the collection of all parent latent variables.

We embed this model in an extended state space. The full extended model consists of data,  $x$ , indices,  $k$ , and latent variables,  $z$ , and factorises as,

$$P(x, k, z) = P(x, z|k) P(k). \quad (3)$$

There are  $n$  latent variables, as in the underlying model, and each latent variable is copied  $K$  times to form an extended state space. We denote a single copy of a single latent variable as  $z_i^\kappa$ , where  $i \in \{1, \dots, n\}$  indexes the latent variable and  $\kappa \in \{1, \dots, K\}$  indexes the copy. That single latent variable,  $z_i^\kappa$ , is a random variable with support  $\mathcal{Z}_i$  (analogous to  $z_i'$ ). The tuple of all  $K$  copies of a single latent variable is denoted  $z_i$ ,

$$z_i = (z_i^1, \dots, z_i^K) \in \mathcal{Z}_i^K \quad (4)$$

and the tuple of all copies of all latent variables is denoted  $z$ ,

$$z = (z_1, \dots, z_n) \in \mathcal{Z}_1^K \times \dots \times \mathcal{Z}_n^K = \mathcal{Z}^K. \quad (5)$$

We also augment the state-space with a set of uniformly distributed random indices,  $k_i \in \mathcal{K} = \{1, \dots, K\}$ , one for each latent variable,

$$P(k_i) = \text{UniformDiscrete}(1, K). \quad (6)$$

These pick out one copy for each latent variable, and the tuple containing all random indices is denoted  $k$ ,

$$k = (k_1, \dots, k_n) \in \mathcal{K}^n. \quad (7)$$

Now, we can write the probability density,

$$P(x, z|k) = \pi(x, z^k) \prod_i q(z_i'^{k_i}; x, z^k, z_{\text{qa}(i)}^k). \quad (8)$$

There's a lot of notation here, so let's unpack it. First,  $z^k$  is a tuple containing a single copy of each latent variable, indexed by  $k$ ,

$$z^k = (z_1^{k_1}, \dots, z_n^{k_n}) \in \mathcal{Z}. \quad (9)$$

Next,  $z_i'^{k_i}$  is the other, non-indexed copies,

$$z_i'^{k_i} = (z_i^1, \dots, z_i^{k_i-1}, z_i^{k_i+1}, \dots, z_i^K) \in \mathcal{Z}_i^{K-1}. \quad (10)$$

and a tuple containing all non-indexed copies of all latents is,

$$z'^k = (z_1'^{k_1}, \dots, z_n'^{k_n}) \in \mathcal{Z}^{K-1}, \quad (11)$$

i.e. all copies of the  $i$ th latent variable except  $z_i^{k_i}$ . Now we can understand the implied generative process:

- First, generate the indices,  $k_i$  from uniform discrete distributions.
- Next, generate the indexed latent variables,  $z^k$ , and data,  $x$ , from the underlying model with density  $\pi$ .
- Next, generate the other non-indexed latent variables,  $z'^k$ .

Now comes the question of what is the most general distribution over the non-indexed latent variables,  $q(z_i^{/k_i}; x, z^k, z_{\text{qa}(i)}^{/k})$ . We generate  $z^k$  and  $x$  first, so clearly  $z_i^{/k_i}$  can depend on them. But we can also introduce conditional dependencies between the non-indexed latent variables,  $z_i^{/k_i}$ . The structure of these conditional dependencies is expressed by  $z_{\text{qa}(i)}^{/k}$ , which selects out those latent variables that are direct parents of  $z_i^{k_i}$  under  $q$ . Importantly, this conditional dependency structure can be different from that in the underlying model, which is why the parents for the distribution over the non-indexed latent variables is denoted  $\text{qa}(i)$ , rather than  $\text{pa}(i)$  for the underlying model.

## 1.1 MCMC

In MCMC, we typically make only a single proposal, then compare the probability of the original sample and proposal under the posterior. This is very problematic because in larger state-spaces, it becomes increasingly difficult to make a proposal for all latent variables jointly that has high probability under the posterior. While we can make multiple proposals (e.g. 20 proposals) [1], this doesn't make a material difference as we move to high-dimensions. Here, we introduce a tensorised MCMC method that makes exponentially proposals,  $n^K$ , and is therefore able to mix quickly and efficiently even in high-dimensional state-spaces.

In tensorised MCMC, the idea is to use the extended state-space construction described in the first Section. One iteration then involves:

- Start with random variables,  $k$ ,  $z^k$  and  $z^{/k}$  from the previous iteration.
- Gibbs sample for  $z^{/k}$  by sampling from the density,  $\prod_i q(z_i^{/k_i}; x, z^k, z_{\text{qa}(i)}^{/k})$ , which acts as a proposal.
- Gibbs sample  $k$ , with fixed  $z$ .

This is an extension of the parallel Metropolis Hastings [1], where we introduce structured conditional independencies in the latent variables, and consider one index for each latent variable.

Resampling  $z^{/k}$  is straightforward. The only difficult part is in sampling  $k$ .

$$P(k|z, x) \propto \pi(x, z^k) \prod_i q(z_i^{/k_i}; x, z^k, z_{\text{qa}(i)}^{/k}). \quad (12)$$

To obtain an efficient MCMC algorithm, we need a proposal that simplifies this conditional. One approach would be to generalise the notion of a symmetric proposal to multiple variables, following [1]. In particular, we consider a proposal which is independent of other latent variables,

$$q(z_i^{/k_i}; x, z^k, z_{\text{qa}(i)}^{/k}) = q(z_i^{/k_i}; z_i^{k_i}) \quad (13)$$

We restrict this proposal to be symmetric in the sense that its probability density should be the same for any choice of  $k_i$  and  $k'_i \neq k_i$ ,

$$q(z_i^{/k_i}; z_i^{k_i}) = q(z_i^{/k'_i}; z_i^{k'_i}). \quad (14)$$

Thus, the conditional distribution for Gibbs sampling becomes proportional to the probability density for the underlying model,

$$P(k|z, x) \propto \pi(x, z^k). \quad (15)$$

What does this distribution look like? Well, remember that the underlying model factorises,

$$\pi(x, z^k) = \pi_x(x; z^k) \prod_{i=1}^n \pi_i(z_i^k; z_{\text{pa}(i)}^k). \quad (16)$$

and that  $z$  and  $x$  are fixed, and we are interested only in the  $k$  dependence. This is a product of factors,  $\pi_x(x; z^k)$  and  $\pi_i(z_i^k; z_{\text{pa}(i)}^k)$ . Each factor depends only on a subset of the  $k_i$ 's. We could thus write each factor as a tensor,

$$f_i^{k_i, k_{\text{pa}(i)}} = \pi_i(z_i^k; z_{\text{pa}(i)}^k) \quad f_i \in \mathbb{R}^{K^{1+|\text{pa}(i)|}} \quad (17)$$

$$f_x^{k_{\text{pa}(x)}} = \pi_x(z_{\text{pa}(x)}^k) \quad f_x \in \mathbb{R}^{K^{|\text{pa}(x)|}} \quad (18)$$

where  $i$  and  $x$  is a label indicating which tensor we are referring to and  $|\text{pa}(i)|$  is the number of parents of the  $i$ th latent variable. Thus,  $k_i, k_{\text{pa}(i)}$  and  $k_{\text{pa}(x)}$  are a list of elements of  $k$ , which act as indices into  $f_i$  and  $f_x$  respectively. The conditional distribution for Gibbs sampling can thus be written as a product of  $n + 1$  factors,

$$P(k|z, x) \propto \pi(x, z^k) = f_x^{k_{\text{pa}(x)}} \prod_{i=1}^n f_i^{k_i, k_{\text{pa}(i)}}. \quad (19)$$

And this is a discrete factor-graph, which can be treated using standard message-passing techniques.

### 1.1.1 Symmetric MCMC proposals with $K > 2$

**Draft: not ready for consumption!**

Unsurprisingly, it is non-trivial to build symmetric proposals for  $K > 2$ . Calderhead [1] had complex proposals involving auxiliary variables that involved are difficult to build into a general probabilistic programming framework. Instead, we consider other approaches to building symmetric proposals. In particular, we would like  $P(Z^{-K}|Z^K)$  to be independent of  $K$ . This is straightforward if we set  $\kappa = 2$ , and use a factorised, symmetric proposal,

$$P(Z^{-K}|Z^K) = \prod_i P(Z_i^{-K_i}|Z_i^{K_i}) \quad (20)$$

$$P(Z_i^{-K_i} = z'_i | Z_i^{K_i} = z_i, X) = P(Z_i^{-K_i} = z_i | Z_i^{K_i} = z'_i, X) \quad (21)$$

It should be possible to obtain analogues of symmetric proposals for higher values of  $K$  by considering spatial structure. For instance, if the individual proposals here are symmetric, and we generate the proposals as a chain,

$$Z_i^1 \leftarrow \dots \leftarrow Z_i^{K_i-1} \leftarrow Z_i^{K_i} \rightarrow Z_i^{K_i+1} \rightarrow \dots \rightarrow Z_i^n. \quad (22)$$

As this requires many tensor operations, it may be possible to put samples on a hyper-cube.

With a symmetric proposal, we have,

$$P(K|Z^K, Z^{-K}) \propto P(K) P(Z^K) P\left(X|Z_{\text{pa}(X)}^{K_{\text{pa}(X)}}\right) \quad (23)$$

$$= P(K) P\left(X|Z_{\text{pa}(X)}^{K_{\text{pa}(X)}}\right) \prod_i P\left(Z_i|Z_{\text{pa}(i)}^{K_{\text{pa}(i)}}\right) \quad (24)$$

Importantly, the right-hand-side can be understood as a factor graph for  $K$ .

## 1.2 VI

Just as the effectiveness of MCMC is determined by the quality of the proposal, the effectiveness of VI is determined by the quality of the approximate posterior. Training a very accurate approximate posterior is difficult, and not always possible (e.g. due to the presence of complex conditional independencies). An alternative approach to improving the variational bound is to draw multiple samples from the approximate posterior and do importance sampling. However, past work on this topic (IWAE) drew a small number of samples (e.g. 20) from the full joint state-space. Here, we seek an approach that is equivalent to drawing exponentially many samples,  $K^n$ , which should dramatically improve the quality of the variational bound.

We construct a variational approximate posterior just over  $z$  from the extended state-space in the first section. To achieve this, we marginalise  $k$  out of the generative model, then compute the usual ELBO (albeit in an extended state-space),

$$\mathcal{L} = E_{Q(z|x)} \left[ \log \frac{\sum_k P(z, x, k)}{Q(z|x)} \right] = E_{Q(z|x)} \left[ \log \frac{\sum_{k_1, \dots, k_n} P(z, x, k_1, \dots, k_n)}{Q(z|x)} \right]. \quad (25)$$

[WHY DOES THIS LOOK LIKE IWAE] where the sums over the  $k_i$ 's run from 1 to  $K$ . The original TMC paper showed that such a marginalisation is closely related to using an importance weighted bound (e.g. as in IWAE), but with an exponential number of importance samples. This exponential number of importance samples enables a much tighter variational bound, and much more accurate posteriors.

We consider an approximate posterior which independently samples each “copy” of the full latent space, indexed  $\kappa$ ,

$$Q(z^\kappa|x) = \prod_{i=1}^n Q(z_i^\kappa|z_{\text{qa}(i)}^\kappa, x) \quad (26)$$

where  $z^\kappa$  is the  $\kappa$ th copy of each latent variable,

$$z^\kappa = (z_1^\kappa, \dots, z_n^\kappa) \in \mathcal{Z}, \quad (27)$$

and where  $\text{qa}(i)$  denotes the parents of variable  $i$  under the approximate posterior. Overall, we have,

$$Q(z|x) = \prod_{\kappa \in \mathcal{K}} Q(z^\kappa|x) = \prod_{\kappa \in \mathcal{K}} \prod_{i=1}^n Q(z_i^\kappa|z_{\text{qa}(i)}^\kappa, x) \quad (28)$$

where, remember,  $\mathcal{K} = \{1, \dots, K\}$ . Note: it is this choice of dependency structure in the approximate posterior that distinguishes our approach from past work in tensor Monte-Carlo and particle filtering. In particular, in all of these approaches, the  $z_i^\kappa$  depends on all copies of the previous latent variables,  $z_{\text{qa}(i)}$ . For instance, the proposal might be a mixture, with one component corresponding to each copy of the previous latent variables.

Now, we need to pick  $q(z_i^{/k_i}; x, z^k, z_{\text{qa}(i)}^{/k})$  in the generative model (Eq. 8) to ensure as much cancellation as possible in the ELBO (Eq. 25). We choose,

$$q(z_i^{/k_i}; x, z^k, z_{\text{qa}(i)}^{/k}) = \prod_{\kappa \in \mathcal{K}/k_i} Q(z_i^\kappa|z_{\text{qa}(i)}^\kappa, x). \quad (29)$$

and where  $\mathcal{K}/k_i$  is all indicies except  $k_i$ ,

$$\mathcal{K}/k_i = \{1, \dots, k_i - 1, k_i + 1, \dots, K\}. \quad (30)$$

Importantly,  $z_{\text{qa}(i)}^\kappa$  is always available in  $q(z_i^{/k_i}; x, z^k, z_{\text{qa}(i)}^{/k})$ . Consider one parent,  $j \in \text{qa}(i)$ . If  $\kappa = k_j$ , then  $z_j^\kappa$  is included in  $z^k$ . Alternatively, if  $\kappa \neq k_j$ , then  $z_j^\kappa$  is included in  $z_{\text{qa}(i)}^{/k} = z_j^{/k}$ .

Substituting this choice of generative model for the extended state space, and approximate posterior into the ELBO, we get,

$$\sum_k \frac{P(z, k, x)}{Q(z|x)} = \sum_k P(k) \frac{\pi(x, z^k) \prod_{i=1}^n \prod_{\kappa \in \mathcal{K}/k_i} Q(z_i^\kappa|z_{\text{qa}(i)}^\kappa, x)}{\prod_{i=1}^n \prod_{\kappa \in \mathcal{K}} Q(z_i^\kappa|z_{\text{qa}(i)}^\kappa, x)} \quad (31)$$

Now, the  $Q(z_i^\kappa|z_{\text{qa}(i)}^\kappa, x)$  terms for all  $\kappa \in \mathcal{K}/k_i$  cancel,

$$\sum_k \frac{P(z, k, x)}{Q(z|x)} = \sum_k P(k) \frac{\pi(x, z^k)}{\prod_{i=1}^n Q(z_i^{k_i}|z_{\text{qa}(i)}^{k_i}, x)}. \quad (32)$$

Remembering that  $\pi(x, z^k)$  could be written as a factor graph, we can write our estimator as,

$$\sum_{k \in \mathcal{K}^n} \frac{P(z, k, x)}{Q(z|x)} = \sum_{k_1, \dots, k_n} P(k) f_x^{k_{\text{pa}(x)}} \prod_{i=1}^n \frac{f_i^{k_i, k_{\text{pa}(i)}}}{Q(z_i^{k_i} | z_{\text{qa}(i)}^{k_i}, x)}. \quad (33)$$

this factor graph has the same structure as the underlying graphical model.

### 1.3 Advantages of this construction over TMC/particle filtering

This approach has several advantages over TMC/particle filtering.

First, this approach has minimal asymptotic complexity. In particular, the factors, which dominate asymptotic compute and memory costs have size  $K^{1+|\text{pa}(i)|}$ . Our particular choice of approximate posterior ensures that the asymptotic costs depend only on the structure of the prior, through  $\text{pa}(i)$ ; they do not depend on the structure of the approximate posterior, expressed in  $\text{qa}(i)$ . We can therefore have arbitrarily complex approximate posteriors, without increasing the computational complexity. This only arises under our particular choice of approximate posterior. Other methods such as particle filtering and the original TMC paper have an approximate posterior that depends on all particles of all previous latent variables, and the usual choice is to use a mixture over all parent particles. This is fine if there is only one parent variable, in which case this mixture has  $K$  components. However, the number of components grows exponentially with the number of parents,  $K^{|\text{qa}(i)|}$ . This exponential growth can be prohibitive in many settings. For instance, consider doing inference over latent variables,  $z = (\theta, y_1, \dots, y_N)$  (there is one  $\theta$  for all datapoints, and one  $y_i$  associated with each datapoint), with the following generative model and approximate posterior.

$$P(x_1, \dots, x_N, y_1, \dots, y_N, \theta) = P(\theta) \prod_{i=1}^{n-1} P(y_i | \theta) P(x_i | y_i) \quad (34)$$

$$Q(y_1, \dots, y_N, \theta | x_1, \dots, x_N) = Q(\theta | y_1, \dots, y_{n-1}) \prod_{i=1}^{n-1} Q(y_i | x_i). \quad (35)$$

Note that each term in the generative model has only one or two latent variables ( $P(\theta)$  and  $P(x_i | y_i)$  have one latent variable, and  $P(y_i | \theta)$  has two). As such, if we implemented this model in our framework, the largest factor would be  $K^2$ , which is eminently tractable. But problems arise if we use a mixture proposal for  $Q(\theta | y_1, \dots, y_{n-1})$ , as in past work from TMC and particle filtering. The issue is that  $\theta$  depends on  $N$  latent variables, so there are  $K^N$  mixture components, which is intractable in all but the smallest models.

Second, our approach offers increased parallelisation as, when sampling from the approximate posterior we simply sample  $K$  independent copies from the full

joint latent space. In contrast, sampling is much more complicated under the mixture proposals in TMC/particle filtering, as  $z_i^k$  might depend on any copy of the previous latent variables. These complex dependencies impede effective GPU parallelisation.

Third, our method is straightforwardly differentiable using autodiff build in to e.g. PyTorch. In contrast, high dimensional mixture proposals are difficult to differentiate (as they cannot easily be reparameterised; see Justin’s paper).

## 1.4 Marginal inference

TPP not only provides a tighter lower bound for estimating the likelihood, it can also benefit the inference over the marginal distribution of latent variables, especially those with multiple modes.

# 2 Implementation notes: don’t look at this! Its not relevant to the underlying algorithm

## 2.1 Summing and sampling in factor graphs

Take the case of MCMC. Here, we start with the full joint, written as a product of factors. We progressively sum out each variable, (omitting conditioning on  $Z$ ),

$$P(K_1, K_2, K_3, K_4) = \sum_{K_5} P(K_1, K_2, K_3, K_4, K_5) \quad (36)$$

$$P(K_1, K_2, K_3) = \sum_{K_4} P(K_1, K_2, K_3, K_4) \quad (37)$$

$$P(K_1, K_2) = \sum_{K_3} P(K_1, K_2, K_3) \quad (38)$$

$$P(K_1) = \sum_{K_2} P(K_1, K_2) \quad (39)$$

Each joint distribution in this list is represented as a product of factors. When we sum out a factor, we:

- Remove all tensors/factors from the list that depend on that variable.
- Take the product of all factors, to make a single big factor.
- Sum out  $k_i$  for that variable.
- Put the single resulting tensor back in the list.

When we do a sum, we multiply all factors in the list that depend on that variable, sum it out, and put it back in the list.



Sampling is now straightforward. We begin by,

$$K_1 \sim P(K_1) \tag{40}$$

$$K_2 \sim P(K_2|K_1) \propto P(K_1, K_2) \tag{41}$$

$$K_3 \sim P(K_3|K_1, K_2) \propto P(K_1, K_2, K_3) \tag{42}$$

$$K_4 \sim P(K_4|K_1, K_2, K_3) \propto P(K_1, K_2, K_3, K_4) \tag{43}$$

$$K_5 \sim P(K_5|K_1, K_2, K_3, K_4) \propto P(K_1, K_2, K_3, K_4, K_5) \tag{44}$$

Importantly, sampling these distributions is easy, because we have all the joints on the right-hand-side, represented as lists of factors. Even if that list is long, we can still sample easily because we have a fixed sample for all the other  $K_i$ 's in that list.

## 2.2 IID sampling in plates

We keep the list of factors as a flat list. When we combine factors, we can only do so if:

- They have the same plates
- There are no lower-layer plates

When there is only one factor left under a plate, we can sum across those dimensions.

This can be achieved, and should be efficient if we just sum out plates in reverse order.

Algorithm:

- Extract plate order, by looking for the first tensor with that plate.
- Sum out plates in reverse order of definition.
  - Find all tensors in plate, remove them from underlying list of tensors, and add them to a new list.
  - Sum out all sample indexes within the plate.
  - Sum out the plate
  - Put the resulting tensor back in the original list.

## References

- [1] Ben Calderhead. A general construction for parallelizing metropolis-hastings algorithms. *Proceedings of the National Academy of Sciences*, 111(49):17408–17413, 2014.