
A Universal Marginalizer for Amortized Inference in Generative Models

Laura Douglas^{*†}, Iliyan Zarov^{*†}, Konstantinos Gourgoulis[†], Chris Lucas[†],
Chris Hart[†], Adam Baker[†], Maneesh Sahani[‡], Yura Perov[†], Saurabh Johri[†]

Abstract

We consider the problem of inference in a causal generative model where the set of available observations differs between data instances. We show how combining samples drawn from the graphical model with an appropriate masking function makes it possible to train a single neural network to approximate all the corresponding conditional marginal distributions and thus amortize the cost of inference. We further demonstrate that the efficiency of importance sampling may be improved by basing proposals on the output of the neural network. We also outline how the same network can be used to generate samples from an approximate joint posterior via a chain decomposition of the graph.

1. Introduction

Graphical models provide a natural framework for expressing probabilistic relationships between random variables, and are widely used to facilitate reasoning and decision-making. Bayesian networks (BN), a directed form of probabilistic graphical model (PGM), have been used extensively in medicine to capture the causal relationships between medical entities such as diseases and symptoms, and through inference, enable diagnosis of unobserved disease states. In sensitive domains such as health care, the penalty for errors in inference is potentially severe. This risk can be mitigated by increasing the complexity of the model of the underlying process; however, this will also increase the cost of inference, limiting the feasibility of most algorithms.

In complex models, exact inference is often computationally intractable. We therefore must resort to approximate methods such as variational inference [Wainwright et al., 2008] and Monte Carlo methods, e.g., importance sampling [Cheng and Druzdzel, 2000, Neal, 2001]. Variational inference methods are fast but inexact; Monte Carlo inference is asymptotically exact, but can be slow. For this reason, we focus on Amortized Inference, techniques which speed up sampling by allowing us to “flexibly reuse inferences so as to answer a variety of related queries” [Gershman and Goodman, 2014]. Amortized inference has been popular for Sequential Monte Carlo and has been used to learn in advance either parameters [Gu et al., 2015] or a discriminative model which provides conditional density estimates [Morris, 2001, Paige and Wood, 2016]. These conditional density estimates can be used as proposals for importance sampling (see Appendix A), an approach also explored in [Le et al., 2017], using a fixed sequential density estimator MADE [Germain et al., 2015]. We propose a related technique with a more general density estimator to allow arbitrary evidence.

Notation: Consider the set of random variables, $\mathbf{X} = \{X_1, \dots, X_N\}$. A BN is a combination of a Directed Acyclic Graph (DAG), with X_i as nodes, and a joint distribution P of the X_i . The distribution P factorizes according to the structure of the DAG with $P(X_i | \text{Pa}(X_i))$ being the conditional distribution of X_i given its parents, $\text{Pa}(X_i)$. We denote by $\mathbf{X}_\mathcal{O}$ a set of instantiated nodes, $\mathbf{X}_\mathcal{O} \subset \mathbf{X}$,

^{*}The authors equally contributed to the manuscript.

[†]babylon, London, UK. Email for correspondence: saurabh.johri@babylonhealth.com

[‡]Gatsby Computational Neuroscience Unit, University College London.

and their instantiations $\mathbf{x}_\mathcal{O}$. To conduct Bayesian inference when provided with a set of unobserved nodes, $\mathbf{X}_\mathcal{U} \subset \mathbf{X} \setminus \mathbf{X}_\mathcal{O}$, we need to compute the posterior marginal, $P(\mathbf{X}_\mathcal{U}|\mathbf{X}_\mathcal{O})$.

2. Universal Marginalizer

We consider a function approximator, such as a neural network (NN), trained to return the marginal posterior distributions $P(X_i|\mathbf{X}_\mathcal{O} = \mathbf{x}_\mathcal{O})$ for each node $X_i \in \mathbf{X}$ given an instantiation $\mathbf{x}_\mathcal{O}$ of *any* set of observations $\mathbf{X}_\mathcal{O} \subset \mathbf{X}$. We write $\tilde{\mathbf{x}}_\mathcal{O}$ for an encoding of the instantiation that specifies both which nodes are observed, and what their values are – some suitable encodings are discussed in Section 5 below. Then, for binary X_i , the desired network maps $\tilde{\mathbf{x}}_\mathcal{O}$ to a vector in $\{0, 1\}^N$ representing the probabilities $p_i := P(X_i = 1|\mathbf{X}_\mathcal{O} = \mathbf{x}_\mathcal{O})$:

$$\mathbf{Y} = \text{UM}(\tilde{\mathbf{x}}_\mathcal{O}) \approx (p_1, \dots, p_N). \quad (1)$$

To approximate any possible posterior marginal distribution (i.e., given any possible set of evidence $\mathbf{X}_\mathcal{O}$), only one model is needed. Due to this we describe this discriminative model as a *Universal Marginalizer* (UM). The existence of such a network is a direct consequence of the universal function approximation theorem (UFAT) [Hornik et al., 1989]. This is simply illustrated by considering marginalization in a BN as a function and that, by UFAT, any continuous function can be arbitrarily approximated by a NN.

2.1. Training a UM

1. Generating Data. The UM can be trained off-line by generating unbiased samples from the BN using ancestral sampling [Koller and Friedman, 2009, Algorithm 12.2]. Each sample is a binary vector which are the values the classifier will learn to predict.

2. Masking. For the purpose of prediction, a subset of the nodes in the sample must be hidden, or *masked*. This masking can be deterministic, i.e., always masking specific nodes, or probabilistic over nodes. We choose to probabilistically mask a sample in an unbiased way by defining a masking probability, $p \sim U[0, 1]$, which is applied to each node. This will create a dataset whose number of observed nodes is uniformly distributed. There is some analogy here to dropout in the input layer and so could work well as a regularizer, independently of this problem [Srivastava et al., 2014].

3. Representation of the Unobserved/Masked Nodes. Masked nodes are created for the purpose of mimicking unobserved nodes and so the representation of masked nodes must be consistent with the unobserved nodes at the time of inference. Different representations will be investigated and tested in Section 5.

4. Training a neural net with Cross Entropy Loss. By [Saerens et al., 2002], the output of a NN with any reasonable loss function can be mapped to a probability estimate, however the cross entropy loss is the most obvious choice as the output is exactly the predicted probability distribution. We train the network using a binary cross entropy loss function in a multi-label classification setting to predict the state of all observed and unobserved nodes.

5. Outputs: Posterior Marginals. The desired posterior marginals are the output of the sigmoid layer. This result can already be used as a rough posterior estimate, however results can be further improved by combining with Importance Sampling. This is discussed in Sections 3 and 4 and is empirically verified in Section 6.

3. Sequential UM for Importance Sampling

We now have a discriminative model which, given a set of observations $\mathbf{X}_\mathcal{O}$, will approximate all the posterior marginals. While useful on its own, the estimated marginals are not guaranteed to be unbiased. To obtain a guarantee of asymptotic unbiasedness, while making use of the speed of the approximate solution, we use the UM for proposals in importance sampling. A naive approach might be to sample each $X_i \in \mathbf{X}_\mathcal{U}$ independently from $\text{UM}(\tilde{\mathbf{x}}_\mathcal{O})_i$, where $\text{UM}(\tilde{\mathbf{x}}_\mathcal{O})_i$ is i -th element of the vector $\text{UM}(\tilde{\mathbf{x}}_\mathcal{O})$. However, the product of the (approximate) posterior marginals may be very different to the true posterior joint, even if the marginal approximations are good. A problematic example of this, where the variance of weights becomes very large, is highlighted in Appendix B.

The universality of the UM makes possible a scheme we call *Sequential Universal Marginalizer Importance Sampling* (SUM-IS). A single proposal sample \mathbf{x}_S is generated sequentially as follows. First, introduce a new partially observed state $\tilde{\mathbf{x}}_{S \cup \mathcal{O}}$ initialized to $\tilde{\mathbf{x}}_{\mathcal{O}}$. Sample $[\mathbf{x}_S]_1 \sim \text{UM}(\tilde{\mathbf{x}}_{\mathcal{O}})_1$, and update $\tilde{\mathbf{x}}_{S \cup \mathcal{O}}$ so that X_1 is now observed with this value. Now we repeat the process, at each step sampling $[\mathbf{x}_S]_i \sim \text{UM}(\tilde{\mathbf{x}}_{S \cup \mathcal{O}})_i$, and updating $\tilde{\mathbf{x}}_{S \cup \mathcal{O}}$ to include the new sampled value. Thus, we can approximate the conditional marginal for a node i given the current sampled state \mathbf{X}_S and evidence $\mathbf{X}_{\mathcal{O}}$, to get the optimal proposal Q_i^* as:

$$Q_i^* = P(X_i | \{X_1, \dots, X_{i-1}\} \cup \mathbf{X}_{\mathcal{O}}) \approx \text{UM}(\tilde{\mathbf{x}}_{S \cup \mathcal{O}})_i = Q_i. \quad (2)$$

The full sample \mathbf{x}_S is thus drawn from an implicit encoding by the UM of the (approximate) posterior *joint* distribution, as can be seen by observing the product of sample probabilities (Equation 3), so may be expected to yield low variance importance weights when used as the proposal distribution.

$$Q = \text{UM}(\tilde{\mathbf{x}}_{\mathcal{O}})_1 \prod_{i=2}^N \text{UM}(\tilde{\mathbf{x}}_{S \cup \mathcal{O}})_i \approx P(X_1 | \mathbf{X}_{\mathcal{O}}) \prod_{i=2}^N P(X_i | X_1, \dots, X_{i-1}, \mathbf{X}_{\mathcal{O}}). \quad (3)$$

The process by which we sample from these approximately optimal proposals is illustrated in Algorithm 1 and in Figure 3 in Appendix C. This procedure requires that nodes are sampled sequentially, using the UM to provide a conditional probability estimate at each step. This can affect computation time, depending on the parallelization scheme used for sampling. However, some parallelization efficiency can be recovered by increasing the number of samples, or batch size, for all steps. Alternatively a hybrid method which approximates the joint but requires only one call of the UM is proposed in Section 4.

4. Hybrid Proposals

The full SUM-IS process requires sequential sampling and many evaluations of the UM, which may be costly. We also explored a heuristic scheme by which a single UM output of all marginals may be combined with ancestral sampling, when nodes are sampled in topological order.

The proposal distribution for each node X_i is a mixture of the UM marginal $\text{UM}(\tilde{\mathbf{x}}_{\mathcal{O}})_i$, and the conditional $P(X_i | \mathbf{x}_{S \cap \text{Pa}(X_i)})$, where $\mathbf{x}_{S \cap \text{Pa}(X_i)}$ encodes the (sampled or evidential) observations of all ancestors of X_i . Note that this conditional can be calculated directly from the graph. The scheme uses a mixture model, with

$$Q(X_i) = \beta \cdot \text{UM}(\tilde{\mathbf{x}}_{\mathcal{O}})_i + (1 - \beta) \cdot P(X_i | \mathbf{x}_{S \cap \text{Pa}(X_i)}).$$

Here, each node in the proposal is drawn either from the UM approximate marginal given the observed evidence, independently of previously sampled nodes, or according to its prior dependence on previously sampled nodes (and any ancestral evidence), independently of evidence nodes that fall later in the topological sequence. This approach expects to blend these two forms of dependence, generating a reasonable IS proposal - described in Algorithm 2 in the Appendix D. We compare different fixed β values in Section 6. However, β can also be a function of the currently sampled state and the observations. This is left for future work.

5. Methods

We trial feed-forward NN architectures with a hyperparameter search on the number and size of hidden layers. The quality of conditional marginals is measured using a test set of posterior marginals computed for multiple sets of evidence via ancestral sampling with 300 million samples. Two main metrics are considered - overall mean absolute error (absolute difference between the true and predicted node posterior) and mean maximum absolute error (maximum absolute difference averaged across the evidence sets). For importance sampling results we also examine the Pearson correlation of the true and predicted marginal vectors, as well as Effective Sample Size (ESS). Kish's ESS is defined as $(\sum_{j=1}^M w_j)^2 / \sum_{j=1}^M w_j^2$. The best performing UM is used for subsequent experiments using the hybrid proposals scheme proposed in Section 4.

We use ReLU non-linearities, apply a dropout with a probability of 0.5 after each hidden layer and use the Adam optimization method [Kingma and Ba, 2014]. We consider two encoding schemes

for unobserved and observed nodes: **2-bit Representation:** Two binary values representation. One binary value represents whether the node is observed, the other represents (if observed) whether it is True or False. **33-bit Representation (1-bit + 32-bit Continuous):** One binary node represents whether the node is observed and another continuous node is in $\{0, 1\}$ if observed and the prior probability if not observed.

6. Results

UM Architecture and Representation Search. We run a hyperparameter search on network size and unobserved representation, reporting the results in Table 1. The largest one layer network performs the best. The difference between the representations is not large, but the results suggest that providing the priors may help improve performance.

Units per hidden layer	2-bit		33-bit (priors)	
	$ e $	$\max(e)$	$ e $	$\max(e)$
(2048)	0.0063	0.3425	0.0060	0.3223
(4096)	0.0053	0.2982	0.0052	0.2951
(1024, 1024)	0.0081	0.4492	0.0083	0.4726
(2048, 2048)	0.0071	0.4076	0.0071	0.4264

Table 1: Average mean absolute error ($|e|$) and average maximum absolute error ($\max(|e|)$) of the UM evaluated on the test set after training on different sized one- and two-layer NN architectures for 20,000 iterations. Best values are highlighted in bold.

Hybrid Importance Sampling using the UM. We assess the change in performance on the evidence sets with increasing number of samples. An increase in the maximum achieved correlation is observed, as well as higher ESS (Table 2 in Appendix E). Figure 1 indicates standard IS ($\beta = 0$) reaches 92% correlation after 2 million samples, whereas hybrid proposals with $\beta = 0.25$ exceed 95% after only 250,000 samples, ultimately achieving 96% correlation in 2 million samples. We achieve both a higher accuracy and a significant reduction in computational cost per inference.

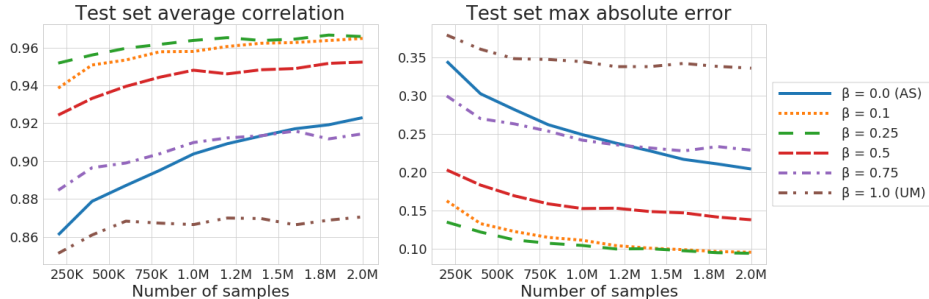


Figure 1: Hybrid importance sampling performance for various values of the mixing parameter β between pure ancestral sampling proposals ($\beta = 0$) and UM marginals independent of the sampled state of unobserved nodes ($\beta = 1$). When $\beta \in [0.1, 0.5]$ we see better marginal estimates in 250k samples than obtained in all 2 million samples when not mixing in the predictions from the UM.

7. Conclusion

This paper introduces a Universal Marginalizer, a neural network which can approximate all conditional marginal distributions of a BN. We have shown that a UM can be used via a chain decomposition of the BN to estimate the joint posterior and thus the optimal proposal distribution for importance sampling. While this process is more computationally intensive, a first-order approximation can be used requiring only a single evaluation of a UM per set of evidence. Our experiments show that the hybrid IS procedure delivers significant improvements in sampling efficiency.

References

- Jian Cheng and Marek J. Druzdzel. AIS-BN: An adaptive importance sampling algorithm for evidential reasoning in large Bayesian networks. *Journal of Artificial Intelligence Research*, 2000.
- Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. MADE: masked autoencoder for distribution estimation. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 881–889, 2015.
- Samuel Gershman and Noah Goodman. Amortized inference in probabilistic reasoning. In *Proceedings of the Cognitive Science Society*, volume 36, 2014.
- Shixiang Gu, Zoubin Ghahramani, and Richard E Turner. Neural adaptive sequential monte carlo. In *Advances in Neural Information Processing Systems*, pages 2629–2637, 2015.
- Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- Tuan Anh Le, Atilim Gunes Baydin, Robert Zinkov, and Frank Wood. Using synthetic data to train neural networks is model-based reasoning. *arXiv preprint arXiv:1703.00868*, 2017.
- Quaid Morris. Recognition networks for approximate inference in bn20 networks. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence, UAI’01*, pages 370–377, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc. ISBN 1-55860-800-1. URL <http://dl.acm.org/citation.cfm?id=2074022.2074068>.
- Radford M Neal. Annealed importance sampling. *Statistics and computing*, 11(2):125–139, 2001.
- Brooks Paige and Frank Wood. Inference networks for sequential Monte Carlo in graphical models. In *International Conference on Machine Learning*, pages 3040–3049, 2016.
- Marco Saerens, Patrice Latinne, and Christine Decaestecker. Any reasonable cost function can be used for a posteriori probability approximation. *IEEE transactions on neural networks*, 13(5):1204–1210, 2002.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014. URL <http://jmlr.org/papers/v15/srivastava14a.html>.
- Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.

A. Importance Sampling

We may use Importance Sampling (IS) to provide posterior marginal estimates, $P(\mathbf{X}_U|\mathbf{X}_O)$ in BN inference. To do so, we draw samples \mathbf{x}_U from a distribution $Q(\mathbf{X}_U|\mathbf{X}_O)$, known as the *proposal* distribution, which we can both sample and evaluate efficiently. Then, assuming that the numerator can be evaluated (which requires that \mathbf{X}_U contain the Markov boundary of \mathbf{X}_O along with all its ancestors), we have:

$$\begin{aligned} P(\mathbf{X}_U = \mathbf{x}_U | \mathbf{X}_O = \mathbf{x}_O) &= \frac{Q(\mathbf{x}_O)}{P(\mathbf{x}_O)} \int 1_{\mathbf{x}_U}(\mathbf{x}) \frac{P(\mathbf{x}, \mathbf{x}_O)}{Q(\mathbf{x}, \mathbf{x}_O)} Q(\mathbf{x}|\mathbf{x}_O) d\mathbf{x} \\ &= \lim_{n \rightarrow \infty} \sum_{i=1}^n 1_{\mathbf{x}_U}(\mathbf{x}_i) \frac{w_i}{\sum_{j=1}^n w_j}, \end{aligned} \tag{4}$$

where $\mathbf{x}_i \sim Q$ and $w_i = P(\mathbf{x}_i, \mathbf{x}_O)/Q(\mathbf{x}_i, \mathbf{x}_O)$ are the *importance sampling weights* and $1_{\mathbf{x}_U}(\mathbf{x})$ is an indicator function for \mathbf{x}_U .

The simplest proposal distribution is the prior, $P(\mathbf{X}_U)$. However, as the prior and the posterior may be very different this is often an inefficient approach. An alternative is to use an estimate of the posterior distribution as a proposal. This is the approach we develop.

B. Sampling from the Posterior Marginals: A Problematic Example

Take a BN with Bernoulli nodes and of arbitrary size and shape. Consider 2 specific nodes, X_i and X_j , such that X_j is caused only and always by X_i :

$$\begin{aligned} P(X_j = 1|X_i = 1) &= 1, \\ P(X_j = 1|X_i = 0) &= 0. \end{aligned}$$

Given evidence E , we assume that $P(X_i|E) = 0.001 = P(X_j|E)$. We will now illustrate that using the posterior distribution $P(X|E)$ as a proposal will not necessarily yield the best result.

Say we have been given evidence, E , and the true conditional probability of $P(X_i|E) = 0.001$, therefore also $P(X_j|E) = 0.001$. We naively would expect $P(X|E)$ to be the optimal proposal distribution. However we can illustrate the problems here by sampling with $Q = P(X|E)$ as the proposal.

Each node $k \in N$ will have a weight $w_k = P(X_k)/Q(X_k)$ and the total weight of the sample will be

$$w = \prod_{k=0}^N w_k.$$

The weights should be approximately 1 if Q is close to P . However, consider the w_j . There are four combinations of X_i and X_j . We will sample $X_i=1, X_j=1$ only, in expectation, one every million samples, however when we do the weight w_j will be $w_j = P(X_j = 1)/Q(X_j = 1) = 1/0.001 = 1000$. This is not a problem in the limit, however if it happens for example in the first 1000 samples then it will outweigh all other samples so far. As soon as we have a network with many nodes whose conditional probabilities are much greater than their marginal proposals this becomes almost inevitable. A further consequence of these high weights is that, since the entire sample is weighted by the same weight, every node probability will be effected by this high variance.

C. Process Diagrams

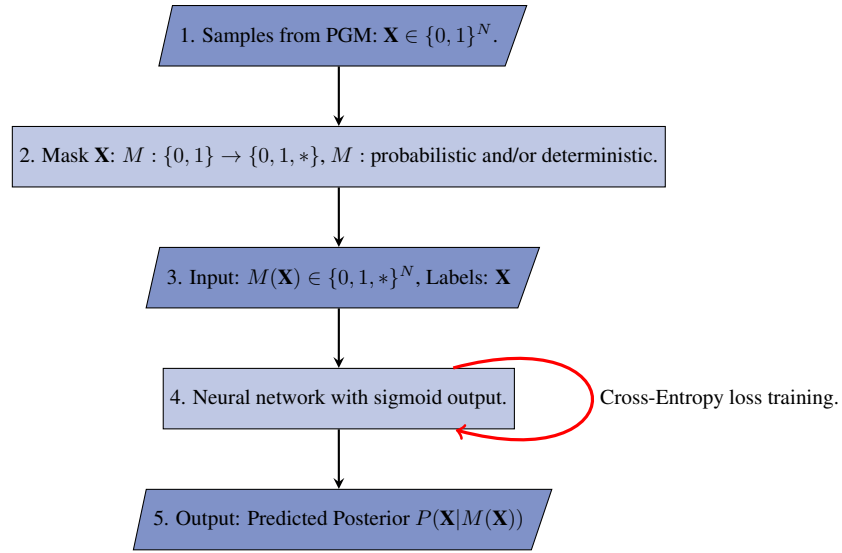


Figure 2: The process to train a Universal Marginalizer using binary data generated from a Bayesian Network

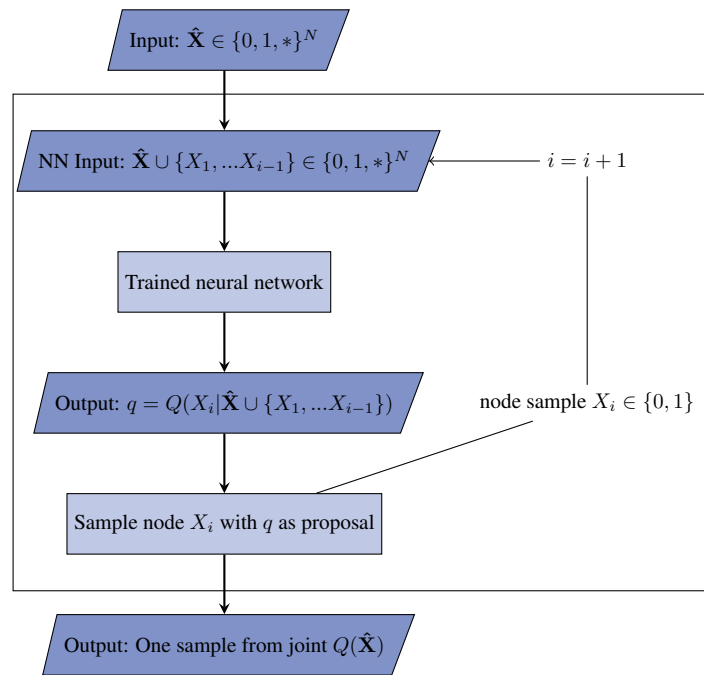


Figure 3: Importance Sampling + UM: The part in the box is repeated N times, for each node i in topological order

D. Algorithms

Algorithm 1 Sequential Universal Marginalizer importance sampling

- 1: Order the nodes topologically X_1, \dots, X_N , where N is the total number of nodes.
 - 2: **for** j in $[1, \dots, M]$ (where M is the total number of samples): **do**
 - 3: $\tilde{\mathbf{x}}_{\mathcal{S}} = \emptyset$
 - 4: **for** i in $[1, \dots, N]$: **do**
 - 5: sample node x_i from $Q(X_i) = \text{UM}(\tilde{\mathbf{x}}_{\mathcal{S} \cup \mathcal{O}})_i \approx P(X_i | \mathbf{X}_{\mathcal{S}}, \mathbf{X}_{\mathcal{O}})$
 - 6: add x_i to $\tilde{\mathbf{x}}_{\mathcal{S}}$
 - 7: $[\mathbf{x}_{\mathcal{S}}]_j = \tilde{\mathbf{x}}_{\mathcal{S}}$
 - 8: $w_j = \prod_{i=1}^N \frac{P_i}{Q_i}$ (where P_i is the likelihood, $P_i = P(X_i = x_i | \mathbf{x}_{\mathcal{S} \cap \text{Pa}(X_i)})$ and $Q_i = Q(X_i = x_i)$)
 - 9: $E_p[X] = \frac{\sum_{j=1}^M X_j w_j}{\sum_{j=1}^M w_j}$ (as in standard IS)
-

Algorithm 2 Hybrid UM-IS

- 1: Order the nodes topologically X_1, \dots, X_N , where N is the total number of nodes.
 - 2: **for** j in $[1, \dots, M]$ (where M is the total number of samples): **do**
 - 3: $\tilde{\mathbf{x}}_{\mathcal{S}} = \emptyset$
 - 4: **for** i in $[1, \dots, N]$: **do**
 - 5: sample node x_i from $Q(X_i) = \beta \text{UM}(\tilde{\mathbf{x}}_{\mathcal{O}})_i + (1 - \beta)P(X_i = x_i | \mathbf{x}_{\mathcal{S} \cap \text{Pa}(X_i)})$
 - 6: add x_i to $\tilde{\mathbf{x}}_{\mathcal{S}}$
 - 7: $[\mathbf{x}_{\mathcal{S}}]_j = \tilde{\mathbf{x}}_{\mathcal{S}}$
 - 8: $w_j = \prod_{i=1}^N \frac{P_i}{Q_i}$ (where P_i is the likelihood, $P_i = P(X_i = x_i | \mathbf{x}_{\mathcal{S} \cap \text{Pa}(X_i)})$ and $Q_i = Q(X_i = x_i)$)
 - 9: $E_p[X] = \frac{\sum_{j=1}^M X_j w_j}{\sum_{j=1}^M w_j}$ (as in standard IS)
-

E. Additional Results

Table 2: Effective sample size for hybrid UM-IS scheme with 2 million samples

β	0.0	0.1	0.25	0.5	0.75	1.0
ESS	7678	15458	11779	1218	171	92