
Conditionally Tractable Density Estimation using Neural Networks

Abstract

Tractable models such as cutset networks and sum-product networks (SPNs) have become increasingly popular as they admit polynomial time inference in some cases. Among them, cutset networks, which model the mechanics of Pearl’s cutset conditioning algorithm, demonstrate great scalability and prediction accuracy. Existing research on cutset networks has mainly focused on discrete domains, and the best mechanism to extend cutset networks to continuous domains is unclear. We propose one possible alternative to cutset networks that models the full joint distribution as the product of a general, complex distribution over a small subset of variables and a fully tractable conditional distribution whose parameters are controlled by a neural network. This model admits exact inference when all variables in the general distribution are observed, and although the model is not fully tractable in general, we show that “cutset” sampling can be employed to efficiently generate accurate predictions in practice. We show that our model performs comparably or better than existing competitors on a variety of real datasets.

1 INTRODUCTION

Tractable probabilistic models such as cutset networks Rahman et al. [2014], sum-product networks Poon and Domingos [2011], arithmetic circuits Darwiche [2003] and probabilistic sentential decision diagrams Kisa et al. [2014] have gained popularity in recent years, primarily because of their ability to accurately answer various reasoning queries in linear time in the size of the model while maintaining other desirable properties such as high expressivity and flexibility. However, to date, with a few exceptions Poon and Domingos [2011], Molina et al. [2018], Uria et al. [2013], a majority of

work on learning tractable models has focused on discrete random variables and not on continuous variables. This is primarily due to the fact that, unlike discrete domains, guaranteeing tractability of models in continuous domains while also maintaining their high expressivity and flexibility is very challenging.

Existing tractable models for continuous variables operate under the assumption that the underlying data can be modeled using a Gaussian mixture. For instance, sum-product networks Poon and Domingos [2011], Molina et al. [2018], Uria et al. [2013] use latent discrete (variable) architectures distributed over sum and product nodes via the *distributive law* to efficiently represent a Gaussian mixture. Unfortunately, these architectures often exhibit poor predictive performance despite the fact that their average test set log-likelihood scores can be quite high. This is a well-known phenomena in continuous models where test set likelihood scores are unable to discriminate low-performing models from high-performing ones. This is because unlike discrete likelihood which always lies between 0 and 1, the continuous likelihood can be potentially infinite or close to it, and thus even one test point that has high likelihood may skew the numbers.

Our aim then, in this work, is to design a tractable model for continuous domains that is flexible enough to be broadly applicable while also yielding good predictive performance in practice. To this end, we propose a novel tractable model for continuous domains motivated by cutset networks (CNs) Rahman et al. [2014]. CNs and their subsequent extensions Roy et al. [2021], Rahman and Gogate [2016] make use of cutset conditioning Pearl [2014]: select a set of variables to condition on (namely cutset variables), such that the distribution over the rest of the variables (leaf variables) can be well-represented by tractable tree structured Bayesian Networks (BNs) Chow and Liu [1968]. Cutset networks are expressive and have shown great scalability and predictive accuracy in high-dimensional discrete domains. A natural extension of cutset networks to continuous domains is to condition each cutset variable over a finite

number of ranges, similar to decision tree regressors, and then model the leaf distributions using conditional linear Gaussian (CLGs) Bayesian networks Grzegorzczak [2010]. Although simple and appealing, in this representation it is hard to ensure that the number of cutset variables is bounded by a small number in order to guarantee tractability without incurring a significant loss in expressive power, and partitioning the space into rectangular blocks is not natural in many practical applications.

To overcome these limitations of cutset networks in continuous domains, and motivated by the work of Rahman et al. [2019] on cutset Bayesian networks, we propose to model the joint distribution as a product of two distributions: (1) a general, complex unconditional distribution defined over a small subset \mathbf{X} of variables and (2) a fully tractable conditional distribution defined over the remaining variables \mathbf{Y} whose parameters are controlled by a neural network. Like cutset Bayesian networks, this model admits exact inference when all variables in the set \mathbf{X} are observed, and although the model is not fully tractable, we show that “cutset” sampling can be employed to efficiently generate accurate predictions in practice. We demonstrate the competitive performance of our approach against two well-known models in the literature: SPNs Lorenzo [2020], Molina et al. [2019], which use latent discrete sum-product architectures and RNADE Uria et al. [2013], which uses neural density estimators, on a prediction task generated from real datasets.

To summarize our contributions, we present (1) a novel tractable architecture, inspired by cutset networks, (2) a novel neural network architecture for choosing the parameters of a conditional linear Gaussian, (3) and an extensive study (90 different experiments) to evaluate the (marginal MAP) prediction accuracy of our model, SPNs, and RNADE on a variety of real datasets.

2 CONDITIONALLY TRACTABLE DENSITY ESTIMATION

In this section, we describe our approach to building tractable models for continuous domains. Note that we use bold uppercase letters for a set of random variables, e.g., \mathbf{X} , while a single random variable is denoted using uppercase letters, e.g., A . The instantiation (configuration) of random variables are denoted as lowercase letters. For example, \mathbf{x} is one possible configuration for all variables in \mathbf{X} and a is a possible value that the random variable A can take. All random variables considered in this paper are assumed to be defined over the real domain \mathbb{R} unless otherwise noted.

Given a set of random variables \mathbf{Z} , we model the full joint distribution over \mathbf{Z} as the product of a general, complex distribution $p(\mathbf{X})$ over a small subset of variables $\mathbf{X} \subset \mathbf{Z}$ and a fully tractable conditional distribution $p(\mathbf{Y}|\mathbf{X})$ over variables $\mathbf{Y} = \mathbf{Z} \setminus \mathbf{X}$ whose parameters are controlled by

a neural network. Specifically, we use a mixture of multivariate Gaussian (MixMG) distribution to model $p(\mathbf{X})$, and a fully connected Gaussian Bayesian network (fGBN) to model the distribution $p(\mathbf{Y}|\mathbf{X})$. In addition, we assume that the parameters of the fGBN are given by a neural network (NN) that takes the values of \mathbf{X} as input and outputs all the parameters inside the fGBN model, i.e., $p(\mathbf{Y}|\mathbf{X} = \mathbf{x}) \sim fGBN(\mathbf{Y}; \theta)$ where $\theta = NN(\mathbf{x})$. Note that any GBN can be converted into an equivalent Multivariate Gaussian (MG) distribution with full covariance matrix Koller and Friedman [2009]. Therefore, the conditional distribution $p(\mathbf{Y}|\mathbf{X})$ is fully tractable, which means our model admits tractable inference when all variables in the general distribution $p(\mathbf{X})$ are observed. In the cases where \mathbf{X} is not fully observed, cutset sampling can be employed to efficiently generate accurate predictions in practice.

In the following sections, we will first introduce the detailed architecture of our model, this includes how we select the conditional variables \mathbf{X} from \mathbf{Z} , as well as the architecture of the NN, and then we will discuss learning and inference algorithms.

2.1 DESIGN OF PARAMETER GENERATION NEURAL NETWORK

In this work, we propose a neural network architecture that we dub parameter generation neural networks (PGNNs) that can be used to generate the parameters for other models (specifically GBNs), see Figure 1. PGNNs are composed of a series of building blocks that are used to extract a shared feature vector stage by stage and a collection of headers that are used to compute the parameters based on the feature vector. The architecture consists of a basic building block (BBlock) that is used to extract non-linear features. It consists of three layers: (1) a (fully connected) linear layer that produces linear features, (2) a one-dimension batch normal layer, and (3) finally a ReLU layer. The Header blocks are used to refine the shared features and include a linear layer that is used to compute the parameters. We stack $s + 1$ building blocks to extract a feature vector of size k that is shared by all Headers. Given n raw input features, the first building block creates a large and expressive feature vector of size $k \cdot r^s$, and the following s building blocks compress the large vector stage by stage until we reach a feature vector with size k . The size of the feature vector is reduced by a factor of r each time it passes through a building block.

Once the final feature vector is built, the individual headers (output layers) generate the parameters for the underlying model. Note that we do not use a single giant header for all parameters to encourage diversity among the different types of generated parameters. Similarly, the model does not include a header for each parameter as, due to the large number of parameters, the computation complexity can be pro-

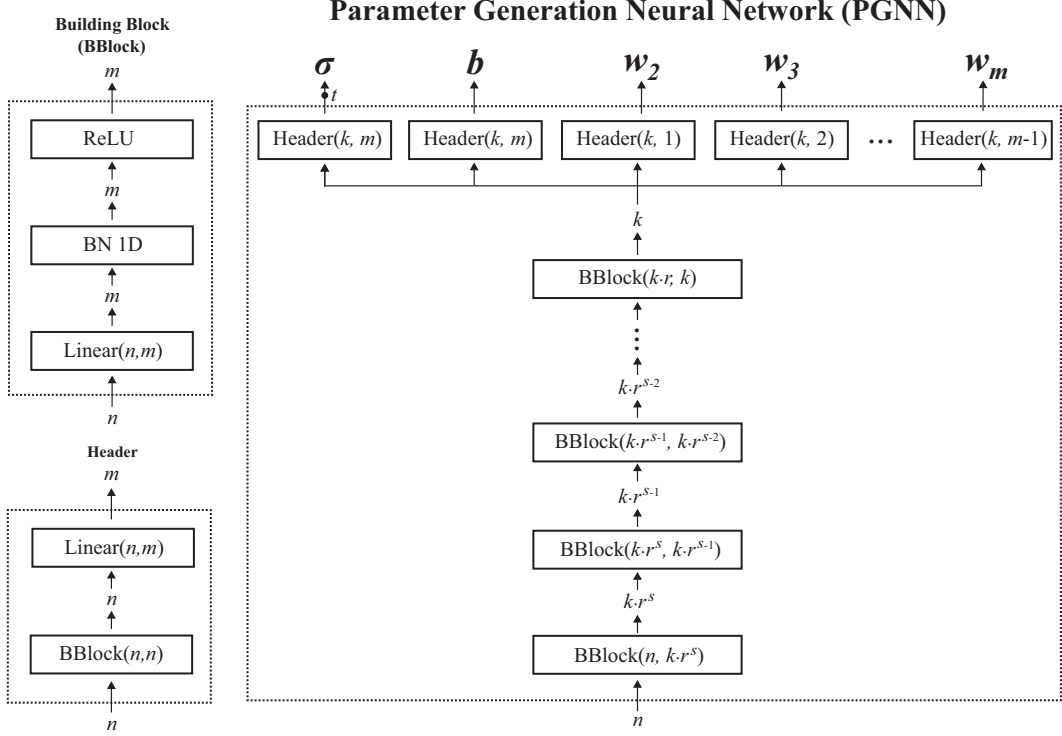


Figure 1: Structure of Parameter Generation Neural Network (PGNN), n is the size of input feature, k is the size of extracted feature, r is the compress ratio and s controls number of stages.

hibitive and this could lead to overfitting. The compromise, which we adopt here, is to use multiple smaller headers, one for each group of parameters. Note that the grouping of parameters is highly dependent on the underlying model. Here, we are generating the parameters of a fully connected GBN over m variables and the conditional distribution of X_j given its parents $\mathbf{X}_{1..j-1}$ is characterized by a conditional linear Gaussian (CLG), i.e.,

$$p(X_j | \mathbf{X}_{1..j-1} = \mathbf{x}_{1..j-1}) \sim N(X_j; \mu = \mathbf{w}_j^T \cdot \mathbf{x}_{1..j-1} + b_j, \sigma_j^2).$$

The mean of the normal distribution over X_j is a linear function of its parents controlled by the weight vector \mathbf{w}_j^T and a bias term b_j and σ_j to denotes the conditional standard deviation, which must be greater than zero. We first create a group $\{\sigma_i | i = 1..m\}$ and use a single header to predict all standard deviations. Then, we do the same for all bias terms $\{b_i | i = 1..m\}$. And finally, we construct one header for each of the weight vectors \mathbf{w}_i . Note that there is no \mathbf{w}_1 since X_1 has no parents and the μ is determined solely by the bias term b_1 . In addition, the output of the σ header clipped to a predefined threshold $t > 0$ in order to guarantee validity of σ as well as to ensure numerical stability.

2.2 VARIABLE SELECTION METHODS

As was the case for cutset networks, the selection of the conditional variables, X , can have a significant impact on

the quality of the final model. We need to determine both the number of conditional variables to select as well as which specific variables. In general, the expressiveness of the model increases as the number of conditional variables and the number of mixture components in $p(X)$ increases. As picking too many or too few random variables to condition on can lead to overfitting or underfitting, we use a validation set to automatically determine the best number of variables to condition on.

Once the size of X is fixed, we also need to determine which variables to pick in order to achieve the the best practical performance. Exhaustive search is prohibitive as there are $\binom{n}{k}$ possible choices if we want to choose $k < n$ variables out of a total of n variables. Conversely, randomly selecting k variables to be the conditional variables is efficient, but some random variables may not be informative for the prediction task. In some settings domain expertise can be used to select the conditional variables, but this expertise is not always available in practice. Consequently, we propose two sets of heuristic methods for selecting the conditional variables X .

The first set of heuristic methods select the variables based on the amount of variance explained by each features in the dataset. Specifically, we use Principal Components Analysis (PCA) as a feature selection method to determine the amount of variance explained by each feature (variable).

After that, we rank those features based on how much variance they explain from highest to lowest. Intuitively, the more variance a feature explains, the more informative it is. One reasonable heuristic would be select variables that have the largest amount of information, which gives us the first heuristic method `MostVar` that selects the first k variables based on the rank. Note that different variables might share some similar information, especially among variables that explain the most variance. Hence, we also consider an alternative heuristic method called `MixVar` that selects k evenly spaced variables based on the rank order. Compared to `MostVar`, `MixVar` produces a more diverse set of variables that have both high and low variance, but the total useful information might be higher than that of `MostVar`. Another alternative to `MixVar` is to randomly select k variables proportional to their contribution to the total variance.

The second set of heuristic methods selects conditional variables based on the correlations between each pair of variables. Recall that, the parameters of the distribution $p(\mathbf{Y}|\mathbf{X})$ are controlled by a neural network that takes \mathbf{x} as input. If the \mathbf{X} random variables have no correlation with the \mathbf{Y} variables, the neural network will be unable to produce an appropriate guess for the model parameters. In other words, we want all $Y_i \in \mathbf{Y}$ to be highly correlated to some or all variables in the set \mathbf{X} . With this idea in mind, we designed our third variable selection heuristic `MaxMincorrHard`. This heuristic initializes \mathbf{X} as an empty set \emptyset and initializes $\mathbf{Y} = \mathbf{Z}$. Then it iteratively selects a variable $A \in \mathbf{Y}$ and adds it into \mathbf{X} such that the minimum correlation between those two sets is maximized. Specifically, the algorithm first calculates the Pearson correlation coefficient $R_{A,B}$ to measure the degree of correlation between each pair of random variables A and B . Then, for each variable $V \in \mathbf{Y}$, it calculates the minimal correlation of all variable $Y_i \in \{\mathbf{Y} \setminus V\}$ to the set $\mathbf{S} = \mathbf{X} \cup V$. Note that the correlation of a single variable Y_i to a set of variables \mathbf{S} is defined as $R_{Y_i, \mathbf{S}} = \sum_{S_i \in \mathbf{S}} R_{Y_i, S_i}$. After that, the variable A that maximizes the minimal correlation will be picked as conditional the variable and added into the set \mathbf{X} before the next iteration starts, i.e.,

$$A \in \operatorname{argmax}_V \left(\min_{Y_i \in \{\mathbf{Y} \setminus V\}} R_{Y_i, \mathbf{X} \cup V} \right).$$

Finally, we also consider another variable picking strategy named `MaxMincorrSoft`. It follows the same iterative strategy as `MaxMincorrHard` but uses a slightly different criteria to select conditional variables. Specifically, in addition to maximizing the minimum correlation, `MaxMincorrSoft` also takes the correlation of V to all variables $Y_i \in \{\mathbf{Y} \setminus V\}$ into consideration. Obviously, if $R_{V, \mathbf{Y} \setminus V}$ is high, it means V contains a lot of useful information about other variables in \mathbf{Y} , and this should encourage the algorithm to select V as the next conditional variable. In short, `MaxMincorrSoft` selects variable B to be the

conditional variable where

$$B \in \operatorname{argmax}_V \left\{ R_{V, \mathbf{Y} \setminus V} * \left(\min_{Y_i \in \{\mathbf{Y} \setminus V\}} R_{Y_i, \mathbf{X} \cup V} \right) \right\}.$$

We have introduced four heuristic methods for selecting the conditional variables \mathbf{X} . Again, as the best strategy may be dataset dependent, we make it a tunable hyperparameter and use a validation set to automatically determine the variable picking method.

2.3 INFERENCE

In this paper, we focus on marginal Max-A-Posterior (MMAP) inference, which reflects the model’s predictive performance in the case of missing values. Given evidence $\mathbf{Z}_e = \mathbf{z}_e$, variables of interest (query variables) \mathbf{Z}_q and missing variables \mathbf{Z}_m , MMAP inference aims to find the best assignment \mathbf{z}_q^* such that $p(\mathbf{Z}_q|\mathbf{z}_e) = \int \mathbf{Z}_m p(\mathbf{Z}_q, \mathbf{Z}_m|\mathbf{z}_e) d\mathbf{Z}_m$ is maximized. As our model has partitioned the variables into two sets, we denote the query, missing, and evidence variables from the general distribution as $\mathbf{X}_q, \mathbf{X}_m, \mathbf{X}_e$ and those in the conditional part as $\mathbf{Y}_q, \mathbf{Y}_m, \mathbf{Y}_e$. The inference task can be formulated as finding \mathbf{x}_q^* and \mathbf{y}_q^* such that

$$\begin{aligned} \mathbf{x}_q^*, \mathbf{y}_q^* &\in \operatorname{argmax}_{\mathbf{x}_q, \mathbf{y}_q} \int_{\mathbf{x}_m} \int_{\mathbf{y}_m} p(\mathbf{x}_{q,m}, \mathbf{y}_{q,m}|\mathbf{x}_e, \mathbf{y}_e) d\mathbf{y}_m d\mathbf{x}_m \\ &\in \operatorname{argmax}_{\mathbf{x}_q, \mathbf{y}_q} \int_{\mathbf{x}_m} \int_{\mathbf{y}_m} p(\mathbf{x}_{q,m,e}, \mathbf{y}_{q,m,e}) d\mathbf{y}_m d\mathbf{x}_m \\ &\in \operatorname{argmax}_{\mathbf{x}_q, \mathbf{y}_q} \int_{\mathbf{x}_m} \int_{\mathbf{y}_m} p_{\mathbf{X}}(\mathbf{x}_{q,m,e}) p_{\mathbf{Y}}(\mathbf{y}_{q,m,e}|\mathbf{x}_{q,m,e}) d\mathbf{y}_m d\mathbf{x}_m \\ &\in \operatorname{argmax}_{\mathbf{x}_q, \mathbf{y}_q} \int_{\mathbf{x}_m} p_{\mathbf{X}}(\mathbf{x}_{q,m,e}) p'_{\mathbf{Y}}(\mathbf{y}_{q,e}|\mathbf{x}_{q,m,e}) d\mathbf{x}_m \end{aligned}$$

where $P'_{\mathbf{Y}}$ is the marginal distribution obtained by integrating out \mathbf{Y}_m from distribution $P_{\mathbf{Y}}$ and $\mathbf{x}_{q,m}$ denotes the union of \mathbf{x}_q and \mathbf{x}_m .

The difficulty of the above MMAP problem varies depends on whether \mathbf{X}_q and \mathbf{X}_m are empty sets.

Case 1: All \mathbf{X} variables are observed. This means $\mathbf{X}_q, \mathbf{X}_m = \emptyset$, and the distribution over the remaining variables \mathbf{Y} is an *fGBN* with parameters $\theta = NN(\mathbf{x})$. In this case, we can do *exact* MMAP inference by converting the *fGBN* into a multivariate Gaussian distribution.

Case 2: $\mathbf{X}_q \neq \emptyset$ while $\mathbf{X}_m = \emptyset$. In this case, since the joint distribution is not tractable, strategy similar to cutset sampling can be employed to efficiently answer the above query. Compared to standard sampling techniques that sample all the variables, cutset sampling only samples a subset of variables and does exact inference over the remaining

variables, which makes the prediction more accurate Rahman et al. [2019]. In our case, we can first generate n samples $\mathbf{x}_q^{(1)}, \dots, \mathbf{x}_q^{(n)}$ for \mathbf{X}_q from the distribution $p_{\mathbf{X}}$ given evidence \mathbf{x}_e . After that, for each $\mathbf{x}_q^{(i)}$, we do exact MMAP inference on $p'_{\mathbf{Y}}$ to find the best assignment $\mathbf{y}_q^{(i)}$ for \mathbf{Y}_q such that $P'_{\mathbf{Y}}(\mathbf{y}_q^{(i)}, \mathbf{y}_e | \mathbf{x}_q^{(i)}, \mathbf{x}_e)$ is maximized. Finally, we evaluate the joint density for each pair $(\mathbf{x}_q^{(i)}, \mathbf{y}_q^{(i)})$ along with the evidence and treat the one with maximum density as the MMAP result. Note that the size of \mathbf{X} is usually small in our case. Therefore, the number of variables being sampled is small, which enables the cutset sampling inference to generate high quality approximations in practice.

Case 3: There are missing variables in \mathbf{X} , which means $\mathbf{X}_m \neq \emptyset$. In this case, we treat those missing variables as part of the query variables and use the cutset sampling method described in Case 2 to find best assignment $\mathbf{x}_{q,m}^*$ for both \mathbf{X}_q and \mathbf{X}_m . Then we simply discard \mathbf{x}_m^* and return \mathbf{x}_q^* as the MMAP result for \mathbf{X}_q . This is a popular method in practice used for approximating MMAP by using the partial assignment from the MAP results Liu and Ihler [2013], Poon and Domingos [2011].

2.4 LEARNING

Our model can be learned with two steps given a dataset D over variables \mathbf{Z} and conditional variables $\mathbf{X} \subset \mathbf{Z}$. First, we learn a mixture of Gaussians over just the \mathbf{X} variables (using the expectation maximization algorithm), where the number of components can be a fixed number or automatically determined through simple tuning. Then, we train a PGNN for generating the parameters of the underlying $fGBN$ model defined over random variables $\mathbf{Y} = \mathbf{Z} \setminus \mathbf{X}$, using the negative log-likelihood of the dataset D as the loss function.

3 EXPERIMENTS

We evaluated our model’s MMAP inference prediction performance, as measured with respect to root mean square error (RMSE), on ten publicly available datasets from the UCI Machine Learning Repository Dua and Graff [2017]. Following Uria et al. [2013], we pre-process the datasets by eliminating discrete valued features and one of the attributes from every pair of attributes whose Pearson correlation coefficient is greater than 0.98. The number of instances and features for each dataset after pre-processing are shown in Table 1. The number of instances range from 5,875 to 150,000 while the number of features range from 12 to 117. All datasets were normalized by subtracting the mean and dividing by the standard deviation.

We consider four competitors in our experiments.

Chow-Liu GBNs (CLGBNs) are GBNs where the underlying structure is a tree learned using the Chow-Liu algorithm. This is considered as a baseline model.

Multivariate Gaussians (MGs) with full covariance matrix. Another baseline that is equivalent to a fully connected GBN.

Real-valued Neural Autoregressive Density Estimators (RNADEs) Uria et al. [2013, 2016] are equivalent to a fully connected Bayesian Network where each of the conditional distributions that define the model are given by a mixture of multivariate Gaussians whose parameters are controlled by neural networks. For MMAP inference, we used two approximate inference algorithms for RNADE. The first one is based on sampling, where we employ likelihood weighted sampling to generate N samples for both the query and the missing variables given the evidence, and then return the sample that achieves the highest likelihood as the prediction result. We denote this algorithm as RNADE-LW. The second algorithm is based on gradient ascent, which is denoted as RNADE-GA. It initializes all query and missing variables to a fixed value, iteratively optimizes the assignment through gradient ascent with the log-likelihood as the objective function, and returns the iteration with the highest log-likelihood. Note that the gradient is approximated using numerical methods since a closed form gradient is not available for RNADE. In addition, backtracking line search is used to determine an appropriate step size in each iteration. For the tuning of hyperparameters, we follow the same strategy as Uria et al. [2013]: we tune the number of components, $\{2, 5, 10, 20\}$; the weight-decay, $\{2.0, 1.0, 0.1, 0.01, 0.001, 0\}$; and the learning rate $\{0.1, 0.05, 0.025, 0.0125\}$. The remaining hyperparameters, where possible, are set using the discussion in Uria et al. [2013], left at the default value determined by existing code, or set to a value that is equivalent to our method if we have the same hyperparameter. For example, both our model and RNADE use the same number of pre-trainings.

Sum-Product Networks (SPNs) Poon and Domingos [2011], Molina et al. [2018] admit tractable MAP inference (under certain constraints). However for MMAP inference, SPNs are not tractable, but MMAP can be approximated by first approximately marginalizing out the missing variables and then applying MAP inference over the marginalized SPN. We also consider another inference method called SPN-S to estimate the MMAP results using the same sampling technique as described above for RNADE. We implement SPNs using existing libraries Molina [2019], Molina et al. [2019], Lorenzo [2020], and tune three different hyperparameters: the instance threshold for creating leaf nodes, $\{50, 100, 200, 300, 500, 1000\}$; the row splitting method, $\{kmeans, gmm, rdc\}$; and the column splitting method,

Table 1: Datasets information

| Name | #instance | #feature |
|--------------|-----------|----------|
| airquality | 9357 | 12 |
| cropmapping | 50000 | 117 |
| energy | 19735 | 24 |
| hepmass | 150000 | 21 |
| miniboone | 36488 | 43 |
| onlinenews | 39644 | 32 |
| parkinson | 5875 | 15 |
| sdd | 58509 | 29 |
| superconduct | 21263 | 68 |
| wec-sydney | 72000 | 49 |

$\{rdc, gvs\}$. Note that we restrict the leaf nodes in SPNs to univariate Gaussian distributions and thus we do not need to tune the feature threshold for creating product nodes.

PGNNs (our model) is trained using the Adam optimizer from pyTorch Paszke et al. [2019] with fixed batch size 100, and weight decay $1e^{-6}$ for 150 epochs. For each of the experiments, the number of stages $s = 3$, the reduction factor $r = 2.5$, and the standard deviation threshold $t = \sqrt{5}e^{-3}$ are fixed. As for the parameter k , which is the size of the feature vector generated by the building blocks, we make it 2 times the size of the input, i.e., $k = 2n$. The learning rate update is scheduled using a linear warm up of 5 epochs starting from 1% of the maximum learning rate, followed by cosine decay. Note that the complexity of the PGNN varies depending on the parameter k , and we decided to scale the max learning rate to encourage the convergence of training when k is relatively large¹. Given a base max learning rate α , we scale it to $\beta = \frac{2\alpha}{1+\sqrt{k}}$, and use β as the maximum learning rate in the training of the PGNN. At last, following Uria et al. [2013, 2016], we pre-train 4 times for 15 epochs each and select the model that achieves best loss for subsequent training. For hyperparameter tuning, we considered base maximum learning rate, $\{0.3, 0.16, 0.08, 0.04\}$; maximum number of conditional variables, $\{10\%, 20\%, 30\%, 40\%, 50\%\}$; and the variable selection method from the four heuristics we discussed in Section 2.2.

For each dataset, we randomly selected 200 instances for both validation and test sets and the remaining instances are used as the training split. All models (SPN, RNADE, OURS) are tuned based on the best RMSE achieved on the validation set. In addition, for sample based approximate inference schemes, the sample size is based on the number of query and missing variables during inference, with a minimum number of 200 samples.

¹Learning rate tuning can also be used to guarantee the convergence of training. However, the range of LR being tuned is usually roughly divided and is not appropriate for every case.

Table 2: One-to-One comparison of win/tie/lose achieved by one model (row) over the other (column).

| | CLGBN | MG | RNADE | SPN | OURS |
|-------|----------|----------|----------|----------|----------|
| CLGBN | - | 00/13/77 | 65/08/17 | 10/13/67 | 00/04/86 |
| MG | 77/13/00 | - | 90/00/00 | 45/07/38 | 00/11/79 |
| RNADE | 17/08/65 | 00/00/90 | - | 01/10/79 | 00/00/90 |
| SPN | 67/13/10 | 38/07/45 | 79/10/01 | - | 15/09/66 |
| OURS | 86/04/00 | 79/11/00 | 90/00/00 | 66/09/15 | - |

All experiments were conducted on a workstation with a 16-core Intel Xeon Gold 6130 CPU and two Quadro P5000 GPUs. Note that running on a GPU is not required for our model as the PGNNs are typically small compared to modern deep neural networks. For a small dataset like parkinson, CPU only performance is 2x faster than using GPUs. For larger datasets like cropmapping, running our model on GPUs can achieve a 3x to 5x speedup over CPUs. All code will be made publicly available on Github post acceptance.

3.1 PREDICTIVE PERFORMANCE

We conducted a comprehensive experiment using nine different settings of the query and missing variables where the percentage of query and missing variables is chosen from the set $\{10\%, 20\%, 30\%\}$. For each setting, we randomly assign query and missing variables for *each* of the instances in test set and we report the mean RMSE over all query variables averaged across all instances in test set. Table 3 shows the average RMSE of our model along with the other four competitors on the ten different datasets under the above nine settings. Numbers in bold denote the lowest RMSE, we consider two RMSEs within 5% difference as a tie. In addition, Table 2 presents a one-to-one comparison of win/tie/lose for one model (row) over another (column) based on the results in Table 3.

From the above results, we have the following observations. (1) Our model achieves the lowest average RMSE over all datasets, regardless of the setting of the query and missing variables. Further, it achieves the best or a comparable RMSE on seven of the datasets for all nine settings. For the other three datasets (airquality, energy and sdd), we lose 2, 5, and 8 out of 9 settings respectively, to only the SPN model. The poorer performance on the sdd dataset is likely the result of overfitting. (2) SPNs are the second best model in terms of the average RMSE over all datasets. However, their predictive performance is not stable across datasets: they produce very accurate predictions for some datasets (sdd) while the predictions on other datasets can be worse

than that of the multivariate Gaussian baseline (cropmapping), which it fails to outperform in roughly half of the cases. Our model, on the other hand, is only outperformed by SPNs and only in roughly 17% of the cases. (3) RNADE, somewhat surprisingly, seems to be the worst model among all models tested. For example, it fails to beat the baseline CLGBN model in over 70% of the cases. We believe there are two likely reasons for this. First, RNADE is a complicated model and is intractable everywhere (even the gradient is numerically approximated). This might result in worse performance of the approximate MMAP procedures used in this case. Further, these MMAP techniques are quite general, i.e., they are not specifically hand-tailored to the RNADE model. Second, the number of hyperparameters that need to be tuned in RNADE is extensive, Uria et al. [2013] manually tuned RNADE on a relatively large dataset. In our case, we only tuned a fraction of the possible hyperparameters as otherwise the search would have been computationally expensive. As such, with additional tuning, the performance of RNADE would likely be improved (though the same could be said of more extensive hyperparameter tuning for SPNs and our model as well). (4) MG produces acceptable or even very good results in many of the cases. This is not surprising as MGs admit exact MMAP inference, so no approximations are necessary.

We also evaluated the log-likelihood of each model on the respective test sets (see Appendix A of the supplementary material). In summary, SPNs achieve the highest average log-likelihood on six out of the ten datasets while our model achieves the highest log-likelihood on the remaining four. However, given the observations in Table 3, it does not appear that the highest test set log-likelihood translates into the best predictive performance (at least not on the MMAP task). In addition, if the hyperparameter tuning was done in order to maximize the log-likelihood score on the validation data instead of minimizing RMSE, then the predictive performance of all of the non-baseline models would likely decrease significantly.

4 CONCLUSION AND FUTURE WORK

We described a flexible family of tractable models for marginal MAP prediction tasks in continuous domains. Experimentally, we verified that our approach outperforms existing approaches such as SPNs and RNADE on a variety of real-world prediction tasks. In particular, we found that (1) although SPNs typically produced models with higher test set log-likelihoods, our model typically resulted in the best performance on the prediction task and (2) the number of hyperparameters required to fit RNADE in practice was computationally prohibitive and/or much more careful hand-tuning may be required, which likely limits its practical utility.

We focused primarily on relatively low dimensional datasets

in this work, and more detailed studies need to be conducted to assess performance differences between these models on high dimensional datasets. We are also interested in exploring applications of this approach to temporal data or more generally in situations where there may be additional graphical structure that can be exploited. In addition, many more neural network structures are possible for parameter selection, and in different applications, it may be of interest to explore different function characterizations of both general distribution $p(\mathbf{X})$ and conditional distribution $p(\mathbf{Y}|\mathbf{X})$. We leave these extensions for future work.

References

- CKCN Chow and Cong Liu. Approximating discrete probability distributions with dependence trees. *IEEE transactions on Information Theory*, 14(3):462–467, 1968.
- Adnan Darwiche. A differential approach to inference in bayesian networks. *Journal of the ACM (JACM)*, 50(3): 280–305, 2003.
- Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- Marco Grzegorzczuk. An introduction to gaussian bayesian networks. In *Systems Biology in Drug Discovery and Development*, pages 121–147. Springer, 2010.
- Doga Kisa, Guy Van den Broeck, Arthur Choi, and Adnan Darwiche. Probabilistic sentential decision diagrams. In *Proceedings of the 14th international conference on principles of knowledge representation and reasoning (KR)*, pages 1–10, 2014.
- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- Qiang Liu and Alexander Ihler. Variational algorithms for marginal map. *The Journal of Machine Learning Research*, 14(1):3165–3200, 2013.
- Loconte Lorenzo. Sum-product networks and normalizing flows for tractable density estimation [mit license]. <https://github.com/loreloc/spnflow>, 2020.
- Alejandro Molina. Sum product flow: An easy and extensible library for sum-product networks [apache license v2.0]. <https://github.com/SPFlow/SPFlow>, 2019.
- Alejandro Molina, Antonio Vergari, Nicola Di Mauro, Sri-raam Natarajan, Floriana Esposito, and Kristian Kersting. Mixed sum-product networks: A deep architecture for hybrid domains. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

Alejandro Molina, Antonio Vergari, Karl Stelzner, Robert Peharz, Pranav Subramani, Nicola Di Mauro, Pascal Poupart, and Kristian Kersting. Spflow: An easy and extensible library for deep probabilistic learning using sum-product networks, 2019.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.

Judea Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier, 2014.

Hoifung Poon and Pedro Domingos. Sum-product networks: A new deep architecture. In *2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*, pages 689–690. IEEE, 2011.

Tahrira Rahman and Vibhav Gogate. Learning ensembles of cutset networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016.

Tahrira Rahman, Prasanna Kothalkar, and Vibhav Gogate. Cutset networks: A simple, tractable, and scalable approach for improving the accuracy of chow-liu trees. In *Joint European conference on machine learning and knowledge discovery in databases*, pages 630–645. Springer, 2014.

Tahrira Rahman, Shasha Jin, and Vibhav Gogate. Cutset bayesian networks: A new representation for learning rao-blackwellised graphical models. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*, 2019.

Chiradeep Roy, Tahrira Rahman, Hailiang Dong, Nicholas Ruozi, and Vibhav Gogate. Dynamic cutset networks. In *International Conference on Artificial Intelligence and Statistics*, pages 3106–3114. PMLR, 2021.

Benigno Uria, Iain Murray, and Hugo Larochelle. Rnade: the real-valued neural autoregressive density-estimator. In *Proceedings of the 26th International Conference on Neural Information Processing Systems (NeurIPS)*, pages 2175–2183, 2013.

Benigno Uria, Marc-Alexandre Côté, Karol Gregor, Iain Murray, and Hugo Larochelle. Neural autoregressive distribution estimation. *The Journal of Machine Learning Research*, 17(1):7184–7220, 2016.

A TEST SET LOG-LIKELIHOOD

Table 4 shows the train and test loglikelihood for all of the models compared in the paper. In addition to the average loglikelihood, we also show the 25%, 50% (median) and 75% percentile of train/test loglikelihood over all train/test instances.

Table 4: Train (marked in grey)/test Loglikelihood, and best test LL is marked in bold.

| Models | Datasets | | | | | | | | | |
|----------------------|---------------|----------------|---------------|----------------|----------------|---------------|---------------|---------------|----------------|----------------|
| | airquality | cropmapping | energy | hepmass | miniboone | onlinenews | parkinson | sdd | superconduct | wec sydney |
| 25 percentile | | | | | | | | | | |
| CLGBN | -11.383 | -72.587 | -22.489 | -29.888 | -50.113 | -37.358 | -10.935 | -31.542 | -55.350 | -65.904 |
| | -12.460 | -70.940 | -21.627 | -30.428 | -51.549 | -37.241 | -11.026 | -31.109 | -60.863 | -64.181 |
| MG | -8.613 | -21.946 | -16.294 | -29.642 | -38.912 | -26.496 | -9.778 | -28.944 | -11.033 | -58.507 |
| | -9.012 | -20.699 | -15.314 | -30.254 | -40.705 | -26.453 | -9.727 | -28.511 | -15.001 | -56.181 |
| RNADE | 8.110 | -151.363 | -12.605 | -27.548 | -32.791 | -24.116 | -8.871 | -17.871 | -46.538 | -65.660 |
| | 7.738 | -148.388 | -11.969 | -28.100 | -32.791 | -25.218 | -8.913 | -16.645 | -43.908 | -64.999 |
| SPN | 17.239 | 15.530 | -0.160 | -28.887 | -35.789 | 37.651 | -7.683 | 12.295 | 1.841 | -54.269 |
| | 17.314 | 10.884 | -2.745 | -29.640 | -35.580 | 29.474 | -8.694 | 12.100 | -3.045 | -54.230 |
| OURS | -1.331 | 63.827 | 0.342 | -26.729 | -24.271 | -9.646 | -4.999 | -12.213 | 26.403 | -58.135 |
| | -1.524 | 64.601 | -0.218 | -27.119 | -28.240 | -10.994 | -6.011 | -11.579 | 25.645 | -57.579 |
| Median | | | | | | | | | | |
| CLGBN | -8.741 | -55.268 | -18.906 | -27.457 | -42.106 | -32.656 | -8.800 | -27.351 | -38.429 | -60.996 |
| | -8.875 | -54.845 | -18.505 | -27.861 | -42.877 | -32.500 | -8.791 | -26.955 | -40.621 | -61.379 |
| MG | -5.871 | -13.910 | -12.642 | -27.195 | -30.665 | -21.348 | -7.634 | -24.708 | 7.662 | -47.397 |
| | -5.949 | -14.178 | -12.554 | -27.726 | -31.127 | -21.057 | -7.657 | -24.075 | 5.400 | -48.223 |
| RNADE | 17.277 | 269.283 | 5.637 | -23.267 | -13.650 | 5.278 | -0.151 | -19.875 | 125.006 | -30.385 |
| | 17.416 | 266.802 | 4.548 | -23.774 | -16.033 | 4.874 | -0.269 | -20.079 | 113.721 | -33.907 |
| SPN | 20.433 | 53.421 | 9.641 | -26.342 | -27.436 | 63.851 | -3.355 | 17.940 | 86.134 | -37.666 |
| | 20.403 | 50.126 | 8.028 | -27.249 | -28.500 | 61.093 | -4.259 | 17.043 | 84.018 | -38.811 |
| OURS | 0.876 | 82.298 | 5.487 | -23.972 | -15.288 | -0.939 | -1.546 | -7.190 | 62.223 | -39.593 |
| | 1.057 | 82.919 | 4.791 | -24.701 | -17.369 | -2.318 | -2.041 | -6.892 | 50.711 | -41.495 |
| 75 percentile | | | | | | | | | | |
| CLGBN | -7.093 | -43.213 | -16.626 | -25.571 | -37.537 | -29.948 | -7.390 | -25.052 | -26.318 | -57.738 |
| | -7.359 | -43.092 | -16.141 | -25.999 | -38.651 | -29.691 | -7.449 | -25.260 | -27.085 | -58.090 |
| MG | -4.444 | -8.425 | -10.331 | -25.312 | -26.020 | -18.561 | -6.269 | -22.397 | 19.288 | -38.185 |
| | -4.583 | -8.217 | -10.285 | -25.731 | -27.158 | -18.082 | -6.280 | -22.439 | 19.122 | -38.550 |
| RNADE | 13.070 | -105.694 | -4.115 | -22.049 | -16.695 | -9.410 | -1.611 | -7.520 | 48.924 | -15.120 |
| | 13.127 | -104.361 | -4.336 | -22.232 | -17.689 | -9.010 | -1.869 | -7.167 | 47.275 | -15.219 |
| SPN | 25.283 | 92.417 | 20.264 | -24.278 | -19.996 | 90.897 | 0.045 | 24.068 | 255.095 | 1.171 |
| | 27.054 | 92.352 | 18.648 | -24.539 | -21.470 | 87.679 | -0.722 | 23.833 | 290.550 | 6.218 |
| OURS | 2.812 | 94.529 | 10.769 | -21.590 | -7.345 | 5.259 | 1.366 | -3.372 | 91.351 | 2.496 |
| | 3.021 | 95.279 | 9.513 | -22.133 | -8.480 | 5.284 | 1.196 | -2.729 | 88.547 | 5.954 |
| Average | | | | | | | | | | |
| CLGBN | -10.156 | -63.346 | -20.252 | -28.178 | -47.836 | -36.688 | -11.614 | -30.517 | -46.466 | -62.348 |
| | -10.596 | -62.214 | -19.564 | -28.584 | -48.657 | -37.501 | -11.713 | -29.021 | -49.418 | -62.054 |
| MG | -7.454 | -20.564 | -14.210 | -27.925 | -37.167 | -27.551 | -10.650 | -31.893 | -2.042 | -49.815 |
| | -7.639 | -17.089 | -13.627 | -28.348 | -37.182 | -27.195 | -10.566 | -26.362 | -4.764 | -48.627 |
| RNADE | 10.484 | -131.016 | -8.591 | -25.089 | -25.942 | -17.803 | -5.938 | -13.537 | 2.619 | -39.457 |
| | 10.242 | -128.895 | -8.370 | -25.583 | -27.512 | -18.257 | -6.371 | -12.555 | 2.560 | -38.361 |
| SPN | 27.953 | 51.723 | 11.080 | -27.003 | -28.953 | 65.680 | -4.634 | 18.856 | 180.669 | -24.174 |
| | 29.579 | 46.459 | 8.657 | -27.562 | -30.970 | 59.444 | -5.543 | 18.526 | 184.080 | -23.590 |
| OURS | 0.746 | 73.029 | 5.312 | -24.498 | -17.023 | -3.078 | -2.634 | -10.124 | 55.991 | -26.813 |
| | 0.763 | 74.106 | 4.455 | -25.466 | -19.347 | -4.174 | -4.169 | -7.427 | 52.627 | -26.559 |