
Tractable Regularization of Probabilistic Circuits

Anji Liu*¹

Guy Van den Broeck^{†1}

¹Computer Science Department, University of California, Los Angeles, USA

Abstract

Probabilistic Circuits (PCs) are a promising avenue for probabilistic modeling. They combine advantages of probabilistic graphical models (PGMs) with those of neural networks (NNs). Crucially, however, they are tractable probabilistic models, supporting efficient and exact computation of many probabilistic inference queries, such as marginals and MAP. Further, since PCs are structured computation graphs, they can take advantage of deep-learning-style parameter updates, which greatly improves their scalability. However, this innovation also makes PCs prone to overfitting, which has been observed in many standard benchmarks. Despite the existence of abundant regularization techniques for both PGMs and NNs, they are not effective enough when applied to PCs. Instead, we re-think regularization for PCs and propose two intuitive techniques, *data softening* and *entropy regularization*, that both take advantage of PCs' tractability and still have an efficient implementation as a computation graph. Specifically, data softening provides a principled way to add uncertainty in datasets in closed form, which implicitly regularizes PC parameters. To learn parameters from a softened dataset, PCs only need linear time by virtue of their tractability. In entropy regularization, the exact entropy of the distribution encoded by a PC can be regularized directly, which is again infeasible for most other density estimation models. We show that both methods consistently improve the generalization performance of a wide variety of PCs. Moreover, when paired with a simple PC structure, we achieved state-of-the-art results on 10 out of 20 standard discrete density estimation benchmarks.

*liuanji@cs.ucla.edu

[†]guyvdb@cs.ucla.edu

1 INTRODUCTION

Probabilistic Circuits (PCs) [Choi et al., 2020c, Dang et al., 2021] are considered to be the lingua franca for Tractable Probabilistic Models (TPMs) as they offer a unified framework to abstract from a wide variety of TPM circuit representations, such as arithmetic circuits (ACs) [Darwiche, 2003], sum-product networks (SPNs) [Poon and Domingos, 2011], and probabilistic sentential decision diagrams (PS-DDs) [Kisa et al., 2014]. PCs are a successful combination of classic probabilistic graphical models (PGMs) and neural networks (NNs). Moreover, by enforcing various structural properties, PCs permit efficient and exact computation of a large family of probabilistic inference queries [Vergari et al., 2021, Khosravi et al., 2019, Shen et al., 2016]. The ability to answer these queries leads to successful applications in areas such as model compression [Liang and Van den Broeck, 2017] and model bias detection [Choi et al., 2020a,b]. At the same time, PCs are analogous to NNs since their evaluation is also carried out using computation graphs. By exploiting the parallel computation power of GPUs, dedicated implementations [Dang et al., 2021, Molina et al., 2019] can train a complex PC with millions of parameters in minutes. These innovations have made PCs much more expressive and scalable to richer datasets that are beyond the reach of “older” TPMs [Peharz et al., 2020a].

However, such advances make PCs more prone to overfitting. Although parameter regularization has been extensively studied in both the PGM and NN communities [Srivastava et al., 2014, Ioffe and Szegedy, 2015], we find that existing regularization techniques for PGMs and NNs are either not suitable or not effective enough when applied to PCs. For example, parameter priors or Laplace smoothing typically used in PGMs, and often used in PC learning as well [Liang et al., 2017, Dang et al., 2020, Gens and Pedro, 2013], incur unwanted bias when learning PC parameters – we will illustrate this point in Sec. 3. Classic NN methods

such as L1 and L2 regularization are not always suitable since PCs often use either closed-form or EM-based parameter updates.

This paper designs parameter regularization methods that are directly tailored for PCs. We propose two regularization techniques, *data softening* and *entropy regularization*. Both formulate the regularization objective in terms of distributions, regardless of their representation and parameterization. Yet, both leverage the tractability and structural properties of PCs. Specifically, data softening injects noise into the dataset by turning hard evidence in the samples into soft evidence [Chan and Darwiche, 2005, Pan et al., 2006]. While learning with such softened datasets is infeasible even for simple machine learning models, with their tractability, PCs can learn the maximum-likelihood estimation (MLE) parameters given a softened dataset in $\mathcal{O}(|p| \cdot |\mathcal{D}|)$ time, where $|p|$ is the size of the PC and $|\mathcal{D}|$ is the size of the (original) dataset. Additionally, the entropy of the distribution encoded by a PC can be tractably regularized. Although the entropy regularization objective for PC is multi-modal and a global optimum cannot be found in general, we propose an algorithm that is guaranteed to converge monotonically towards a stationary point.

We show that both proposed approaches consistently improve the test set performance over standard density estimation benchmarks. Furthermore, we observe that when data softening and entropy regularization are properly combined, even better generalization performance can be achieved. Specifically, when paired with a simple PC structure, this combined regularization method achieves state-of-the-art results on 10 out of 20 standard discrete density estimation benchmarks.

Notation We denote random variables by uppercase letters (e.g., X) and their assignments by lowercase letters (e.g., x). Analogously, we use bold uppercase letters (e.g., \mathbf{X}) and bold lowercase letters (e.g., \mathbf{x}) for sets of variables and their joint assignments, respectively.

2 TWO INTUITIVE IDEAS FOR REGULARIZING DISTRIBUTIONS

A common way to prevent overfitting in machine learning models is to regularize the syntactic representation of the distribution. For example, L1 and L2 losses add mutually independent priors to all parameters of a model; other approaches such as Dropout [Srivastava et al., 2014] and Bayesian Neural Networks (BNNs) [Goan and Fookes, 2020] incorporate more complex and structured priors into the model [Gal and Ghahramani, 2016]. In this section, we ask the question: how would we regularize an arbitrary distribution, regardless of the model at hand, and the way it is parameterized? Such global, model-agnostic regularizers appear to be under-explored. Next, we introduce two intuitive

ideas for regularizing distributions, and study how they can be practically realized in the context of probabilistic circuits in the remainder of this paper.

Data softening Data augmentation is a common technique to improve the generalization performance of machine learning models [Perez and Wang, 2017, Szegedy et al., 2016]. A simple yet effective type of data augmentation is to inject noise into the samples, for example by randomly corrupting bits or pixels [Vincent et al., 2008]. This can greatly improve generalization as it renders the model more robust to such noise. While current noise injection methods are implemented as a sequence of sampled transformations, we stress that some noise injection can be done in closed form: we will be considering all possible corruptions, each with their own probability, as a function of how similar they are to a training data point.

Consider boolean variables¹ as an example: after noise injection, a sample $X = 1$ is represented as a distribution over all possible assignments (i.e., $X = 1$ and $X = 0$), where the instance $X = 1$, which is “similar” to the original sample, gets a higher probability: $P(X = 1) = \beta$. Here $\beta \in (0.5, 1]$ is a hyperparameter that specifies the regularization strength — if $\beta = 1$, no regularization is added; if β approaches 0.5, the regularized sample represents an (almost) uniform distribution. For a sample \mathbf{x} with K variables $\mathbf{X} := \{X_i\}_{i=1}^K$, where the k th variable takes value x_k , we can similarly ‘soften’ \mathbf{x} by independently injecting noise into each variable, resulting in a *softened distribution* $P_{\mathbf{x},\beta} (\forall \mathbf{x}' \in \text{val}(\mathbf{X}))$:

$$\begin{aligned} P_{\mathbf{x},\beta}(\mathbf{X} = \mathbf{x}') &:= \prod_{i=1}^K P_{x_i,\beta}(X_i = x'_i) \\ &= \prod_{i=1}^K \left(\beta \cdot \mathbb{1}[x'_i = x_i] + (1-\beta) \cdot \mathbb{1}[x'_i \neq x_i] \right). \end{aligned}$$

For a full dataset $\mathcal{D} := \{\mathbf{x}^{(i)}\}_{i=1}^N$, this softening of the data can also be represented through a new, *softened dataset* \mathcal{D}_β . Its empirical distribution is the average softened distribution of its data. It is a weighted dataset $\mathcal{D}_\beta := \{\mathbf{x} | \mathbf{x} \in \text{val}(\mathbf{X})\}$, where the weight $\text{weight}(\mathcal{D}_\beta, \mathbf{x})$ of sample \mathbf{x} is:

$$\text{weight}(\mathcal{D}_\beta, \mathbf{x}) = \frac{1}{N} \sum_{i=1}^N P_{\mathbf{x}^{(i)},\beta}(\mathbf{X} = \mathbf{x}). \quad (1)$$

This softened dataset ensures that each possible assignment has a small but non-zero weight in the training data. Consequently, any distribution learned on the softened data must assign a small probability everywhere as well. Of course, materializing this dataset, which contains all possible training examples, is not practical. Regardless, we will think of data softening as implicitly operating on this softened dataset. We remark that data softening is related to soft evidence [Jeffrey, 1990] and virtual evidence [Pearl, 2014],

¹We postpone the discussion on regularizing samples with non-boolean variables in Appendix A.1.

which both define a framework to incorporate uncertain evidence into a distribution.

Entropy regularization Shannon entropy is an effective indicator for overfitting. For a dataset \mathcal{D} with N distinct samples, a perfectly overfitting model that learns the exact empirical distribution has entropy $\log(N)$. A distribution that generalizes well should have a much larger entropy, since it assigns positive probability to exponentially more assignments near the training samples. Concretely, for the protein sequence density estimation task [Russ et al., 2020] that we will experiment with in Sec. 4.3, the perfectly overfitting empirical distribution has entropy 3, a severely overfitting learned model has entropy 92, yet a model that generalizes well has entropy 177. Therefore, directly controlling the entropy of the learned distribution will help mitigate overfitting. Given a model P_θ parametrized by θ and a dataset $\mathcal{D} := \{\mathbf{x}^{(i)}\}_{i=1}^N$, we define the following entropy regularization objective:

$$\text{LL}_{\text{ent}}(\theta; \mathcal{D}, \tau) := \frac{1}{N} \sum_{i=1}^N \log P_\theta(\mathbf{x}^{(i)}) + \tau \cdot \text{ENT}(P_\theta), \quad (2)$$

where $\text{ENT}(P_\theta) := -\sum_{\mathbf{x} \in \text{val}(\mathbf{X})} P_\theta(\mathbf{x}) \log P_\theta(\mathbf{x})$ denotes the entropy of distribution P_θ , and τ is a hyperparameter that controls the regularization strength. Various forms of entropy regularization have been used in the training process of deep learning models. Different from Eq. (2), these methods regularize the entropy of a parametric [Grandvalet and Bengio, 2006, Zhu et al., 2017] or non-parametric [Feng et al., 2017] output space of the model.

Although both ideas for regularizing distributions are rather intuitive, it is surprisingly hard to implement them in practice since they are intractable even for the simplest machine learning models.

Theorem 1. *Computing the likelihood of a distribution represented as a exponentiated logistic regression (or equivalently, a single neuron) given softened data is #P-hard.*

Theorem 2. *Computing the Shannon entropy of a normalized logistic regression model is #P-hard.*

Proofs of all theoretical results are provided in an extended version of the paper due to space limitation.² Although data softening and entropy regularization are infeasible for many models, we will show in the following sections that they are tractable to use when applied to Probabilistic Circuits (PCs) [Choi et al., 2020c], a class of expressive TPMs.

3 BACKGROUND AND MOTIVATION

Probabilistic Circuits (PCs) are a collective term for a wide variety of TPMs. They present a unified set of notations

²<https://tinyurl.com/sfbb67vc>

that provides succinct representations for TPMs such as Probabilistic Sentential Decision Diagrams (PSDDs) [Kisa et al., 2014], Sum-Product Networks (SPNs) [Poon and Domingos, 2011], and Arithmetic Circuits (ACs) [Darwiche, 2003]. We proceed by introducing the syntax and semantics of a PC.

Definition 1 (Probabilistic Circuits). A PC p that represents a probability distribution over variables \mathbf{X} is defined by a parametrized directed acyclic graph (DAG) with a single root node, denoted n_r . The DAG comprises three kinds of units: *input*, *sum*, and *product*. Each leaf node n in the DAG corresponds to an input unit; each inner node n (i.e., sum and product units) receives inputs from its children, denoted $\text{in}(n)$. Each unit n encodes a probability distribution p_n , defined as follows:

$$p_n(\mathbf{x}) := \begin{cases} f_n(\mathbf{x}) & \text{if } n \text{ is an input unit,} \\ \sum_{c \in \text{in}(n)} \theta_{n,c} \cdot p_c(\mathbf{x}) & \text{if } n \text{ is a sum unit,} \\ \prod_{c \in \text{in}(n)} p_c(\mathbf{x}) & \text{if } n \text{ is a product unit,} \end{cases}$$

where f_n is a univariate input distribution (e.g., boolean, categorical or Gaussian), and $\theta_{n,c}$ represents the parameter corresponds to edge (n, c) . Intuitively, a sum unit models a weighted mixture distribution over its children, and a product unit encodes a factored distribution over its children. We assume w.l.o.g. that all parameters are positive and the parameters associated with any sum unit n sum up to 1 (i.e., $\sum_{c \in \text{in}(n)} \theta_{n,c} = 1$). We further assume w.l.o.g. that a PC alternates between sum and product layers [Vergari et al., 2015]. The size of a PC p , denoted $|p|$, is the number of edges in its DAG.

This paper focuses on two classes of PCs that support different types of queries: (i) PCs that allow linear-time computation of marginal (MAR) and maximum-a-posterior (MAP) inferences (e.g., PSDDs [Kisa et al., 2014], selective SPNs [Peharz et al., 2014]); (ii) PCs that only permit linear-time computation of MAR queries (e.g., SPNs [Poon and Domingos, 2011]). The borders between these two types of PCs are defined by their *structural properties*, i.e., constraints imposed on a PC. First, in order to compute MAR queries in linear time, both classes of PCs should be decomposable (Def. 2) and smooth (Def. 3) [Choi et al., 2020c]. These are properties of the (variable) scope $\phi(n)$ of PC units n , that is, the collection of variables defined by all its descendent input nodes.

Definition 2 (Decomposability). A PC is decomposable if for every product unit n , its children have disjoint scopes: $\forall c_1, c_2 \in \text{in}(n) (c_1 \neq c_2), \phi(c_1) \cap \phi(c_2) = \emptyset$.

Definition 3 (Smoothness). A PC is smooth if for every sum unit n , its children have the same scope: $\forall c_1, c_2 \in \text{in}(n), \phi(c_1) = \phi(c_2)$.

Next, *determinism* is required to guarantee efficient computation of MAP inference [Mei et al., 2018].

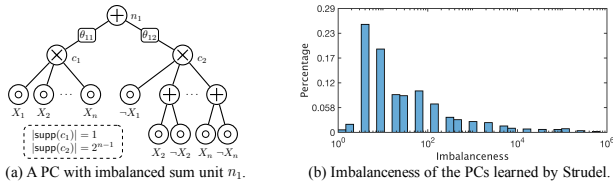


Figure 1: A Problem of Laplace smoothing. (a) Laplace smoothing cannot properly regularize this PC as the sum unit n_1 is imbalanced, i.e., its two children have drastically different support sizes. (b) A large fraction of sum units learned by a PC structure learning algorithm [Dang et al., 2020] are imbalanced.

Definition 4 (Determinism). Define the support $\text{supp}(n)$ of a PC unit n as the set of complete variable assignments $\mathbf{x} \in \text{val}(\mathbf{X})$ for which $p_n(\mathbf{x})$ has non-zero probability (density): $\text{supp}(n) = \{\mathbf{x} \mid \mathbf{x} \in \text{val}(\mathbf{X}), p_n(\mathbf{x}) > 0\}$. A PC is deterministic if for every sum unit n , its children have disjoint support: $\forall c_1, c_2 \in \text{in}(n) (c_1 \neq c_2), \text{supp}(c_1) \cap \text{supp}(c_2) = \emptyset$.

Since the only difference in the structural properties of both PCs classes is determinism, we denote members in the first PC class as deterministic PCs, and members in the second PC class as non-deterministic PCs. Interestingly, both PC classes not only differ in their tractability, which is characterized by the set of queries that can be computed within $\text{poly}(|p|)$ time [Vergari et al., 2021], they also exhibit drastically different expressive efficiency. Specifically, abundant empirical [Dang et al., 2020, Peharz et al., 2020a] and theoretical [Choi and Darwiche, 2017] evidences suggest that non-deterministic PCs are more expressive than their deterministic counterparts. Due to their differences in terms of tractability and expressive efficiency, this paper studies parameter regularization on deterministic and non-deterministic PCs separately.

Motivation Laplace smoothing is widely adopted as a PC regularizer [Liang et al., 2017, Dang et al., 2020]. Since it is also the default regularizer for classical probabilistic models such as Bayesian Networks (BNs) [Heckerman, 2008] and Hierarchical Bayesian Models (HBMs) [Allenby and Rossi, 2006], this naturally raises the following question: *are there differences between a good regularizer for classical probabilistic models such as BNs and HBMs and effective regularizers for PCs?* The question can be answered affirmatively — while Laplace smoothing provides good priors to BNs and HBMs, its uniform prior could add unwanted bias to PCs. Specifically, for every sum unit n , Laplace smoothing assigns the same prior to all its child parameters (i.e., $\{\theta_{n,c} \mid c \in \text{in}(n)\}$), while in many practical PCs, these parameters should be given drastically different priors. For example, consider the PC shown in Fig. 1(a). Since c_2 has an exponentially larger support than c_1 , it should be assumed as prior that θ_{12} will be much larger than θ_{11} .

Algorithm 1 Forward pass

```

1: Input: A deterministic PC  $p$ ; sample  $\mathbf{x}$ 
2: Output:  $\text{value}[n] := (\mathbf{x} \in \text{supp}(n))$  for each unit  $n$ 
3: foreach  $n$  traversed in postorder do
4:   if  $n$  isa input unit then  $\text{value}[n] \leftarrow f_n(\mathbf{x})$ 
5:   elif  $n$  isa product unit then
6:      $\text{value}[n] \leftarrow \prod_{c \in \text{in}(n)} \text{value}[c]$ 
7:   else  $n$  isa sum unit
8:      $\text{value}[n] \leftarrow \sum_{c \in \text{in}(n)} \text{value}[c]$ 

```

Algorithm 2 Backward pass

```

1: Input: A deterministic PC  $p$ ;  $\forall n, \text{value}[n]$ 
2: Output:  $\text{flow}[n, c] := (\mathbf{x} \in (\gamma_n \cap \gamma_c))$  for each pair  $(n, c)$ ,
   where  $n$  is a sum unit and  $c \in \text{in}(n)$ 
3:  $\forall n, \text{context}[n] \leftarrow 0$ ;  $\text{context}[n_r] \leftarrow \text{value}[n_r]$ 
4: foreach sum unit  $n$  traversed in preorder do
5:   foreach  $m \in \text{pa}(n)$  do (denote  $g \leftarrow \text{pa}(m)$ )
6:      $\mathbf{f} \leftarrow \frac{\text{value}[m]}{\text{value}[g]} \cdot \text{context}[g]$ 
7:      $\text{context}[n] += \mathbf{f}$ ;  $\text{flow}[g, m] = \mathbf{f}$ 

```

We highlight the significance of the above issue by examining the fraction of sum units with imbalanced child support sizes in PCs learned by Strudel, a state-of-the-art structure learning algorithm for deterministic PCs [Kisa et al., 2014]. We examine 20 PCs learned from the 20 density estimation benchmarks [Van Haaren and Davis, 2012], respectively. All sum units with ≥ 3 children and with a support size ≥ 128 are recorded. We measure “imbalanceness” of a sum unit n by the fraction of the maximum and minimum support size of its children (i.e., $\frac{\max_{c_1 \in \text{in}(n)} |\text{supp}(c_1)|}{\min_{c_2 \in \text{in}(n)} |\text{supp}(c_2)|}$). As demonstrated in Fig. 1(b), more than 20% of the sum units have imbalanceness $\geq 10^2$, which suggests that the inability of Laplace smoothing to properly regularize PCs with imbalanced sum units could lead to severe performance degradation in practice.

4 HOW IS THIS TRACTABLE AND PRACTICAL?

In this section, we first provide additional background about the parameter learning algorithms for deterministic and non-deterministic PCs (Sec. 4.1). We then demonstrate how the two intuitive ideas for regularizing distributions (Sec. 2), i.e., data softening and entropy regularization, can be efficiently implemented for deterministic (Sec. 4.2) and non-deterministic (Sec. 4.3) PCs.

4.1 LEARNING THE PARAMETERS OF PCS

Deterministic PCs Given a deterministic PC p defined on variables \mathbf{X} and a dataset $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$, the maximum likelihood estimation (MLE) parameters $\theta_{\mathcal{D}}^* := \text{argmax}_{\theta} \sum_{i=1}^N \log p(\mathbf{x}^{(i)}; \theta)$ can be learned in closed-

form. To formalize the MLE solution, we need a few extra definitions.

Definition 5 (Context). The context γ_n of every unit n in a PC p is defined in a top-down manner: for the base case, context of the root node n_r is defined as its support: $\gamma_{n_r} := \text{supp}(n_r)$. For every other node n , its context is the intersection of its support and the union of its parents’ ($\text{pa}(n)$) contexts:

$$\gamma_n := \bigcup_{m \in \text{pa}(n)} \gamma_m \cap \text{supp}(n).$$

Intuitively, if an assignment \mathbf{x} is in the context of unit n , then there exists a path on the PC’s DAG from n to the root unit n_r such that for any unit m in the path, we have $\mathbf{x} \in \text{supp}(m)$. Circuit flow extends the notation of context to indicate whether a sample \mathbf{x} is in the context of an edge.

Definition 6 (Flows). The flow $F_{n,c}(\mathbf{x})$ of any edge (n, c) in a PC given variable assignments $\mathbf{x} \in \text{val}(\mathbf{X})$ is defined as $\mathbb{1}[\mathbf{x} \in \gamma_n \cap \gamma_c]$, where $\mathbb{1}[\cdot]$ is the indicator function. The flow $F_{n,c}(\mathcal{D})$ w.r.t. dataset $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ is the sum of the flows of all its samples: $F_{n,c}(\mathcal{D}) := \sum_{i=1}^N F_{n,c}(\mathbf{x}^{(i)})$.

The flow $F_{n,c}(\mathbf{x})$ for *all* edges (n, c) in a PC p w.r.t. sample \mathbf{x} can be computed through a forward and backward path that both take $\mathcal{O}(|p|)$ time. The forward path, as shown in Alg. 1, starts from the leaf units and traverses the PC in postorder to compute $\forall n, \text{value}[n] := \mathbb{1}[\mathbf{x} \in \text{supp}(n)]$; afterwards, the backward path illustrated in Alg. 2 begins at the root unit n_r and traverses the PC in preorder to compute $\forall n, \text{context}[n] := \mathbb{1}[\mathbf{x} \in \gamma_n]$ as well as $\forall (n, c), \text{flow}[n, c] := F_{n,c}(\mathbf{x})$. By Def. 6, the time complexity for computing $F_{n,c}(\mathcal{D})$ with respect to all edges (n, c) in p is $\mathcal{O}(|p| \cdot |\mathcal{D}|)$, where $|\mathcal{D}|$ is the size of dataset \mathcal{D} .

The MLE parameters $\theta_{\mathcal{D}}^*$ given dataset \mathcal{D} can be computed using the flows [Kisa et al., 2014]:

$$\forall (n, c), \quad \theta_{n,c}^* = F_{n,c}(\mathcal{D}) / \sum_{c \in \text{in}(n)} F_{n,c}(\mathcal{D}). \quad (3)$$

Define hyperparameter $\alpha \geq 0$, for every sum unit n , Laplace smoothing regularizes its child parameters (i.e., $\{\theta_{n,c} \mid c \in \text{in}(n)\}$) by adding a *pseudocount* $\alpha/|\text{in}(n)|$ to every child branch of n , which is equivalent to adding $\alpha/|\text{in}(n)|$ to the numerator of Eq. (3) and α to its denominator.

Non-deterministic PCs As justified by Peharz et al. [Peharz et al., 2016], every non-deterministic PC can be augmented as a deterministic PC with additional hidden variables. For example, in Fig. 2, the left PC is not deterministic since the support of both children of n_1 (i.e., n_2 and n_3) contains $x_1 \bar{x}_2$. The right PC augments the left one by adding input units correspond to hidden variable Z_1 , which retains determinism by “dividing” the overlapping support $x_1 \bar{x}_2$ into $x_1 \bar{x}_2 z_1 \in \text{supp}(n_2)$ and $x_1 \bar{x}_2 \bar{z}_1 \in \text{supp}(n_3)$. Under this interpretation, parameter learning of non-deterministic PCs

is equivalent to learning the parameters of deterministic PCs given incomplete data (we never observe the hidden variables), which can be solved by Expectation-Maximization (EM) [Darwiche, 2009, Dempster et al., 1977]. In fact, EM is the default parameter learning algorithm for non-deterministic PCs [Peharz et al., 2020a, Choi et al., 2020a].

Under the latent variable model view of a non-deterministic PC, its EM updates can be computed using *expected flows* [Choi et al., 2020a]. Specifically, given observed variables \mathbf{X} and (implicit) hidden variables \mathbf{Z} , the expected flow of edge (n, c) given dataset \mathcal{D} is defined as

$$\text{EF}_{n,c}(\mathcal{D}; \theta) := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{z} \sim p_c(\cdot | \mathbf{x}; \theta)} [F_{n,c}(\mathbf{x}, \mathbf{z})],$$

where θ is the set of parameters, and $p_c(\cdot | \mathbf{x}; \theta)$ is the conditional probability over hidden variables \mathbf{Z} given \mathbf{x} specified by the PC rooted at unit c . Similar to flows, the expected flows can be computed via a forward and backward pass of the PC (Alg. 5 and 6 in the Appendix). As shown by Choi et al. [Choi et al., 2020a], for a non-deterministic PC, its parameters for the next EM iteration are given by

$$\theta_{n,c}^{(new)} = \text{EF}_{n,c}(\mathcal{D}; \theta) / \sum_{c \in \text{in}(n)} \text{EF}_{n,c}(\mathcal{D}; \theta). \quad (4)$$

This paper uses a hybrid EM algorithm, which uses mini-batch EM updates to initiate the training process, and switch to full-batch EM updates afterwards. Specifically, in mini-batch EM, $\theta^{(new)}$ are computed using mini-batches of samples, and the parameters are updated towards the target with a step size η : $\theta^{(k+1)} \leftarrow (1 - \eta)\theta^{(k)} + \eta\theta^{(new)}$; when using full-batch EM, we iteratively compute the updated parameters $\theta^{(new)}$ using the whole dataset. Fig. 3 demonstrates that this hybrid approach offers faster convergence speed compared to using full-batch or mini-batch EM only.

4.2 REGULARIZING DETERMINISTIC PCS

We demonstrate how the intuitive ideas for regularizing distributions presented in Sec. 2 (i.e., data softening and entropy regularization) can be efficiently applied to deterministic PCs.

Data softening As hinted by Eq. (1), we need exponentially many samples to represent a softened dataset, which makes parameter learning intractable even for the simple logistic regression model (Thm. 1), let alone more complex probabilistic models such as VAEs [Kingma and Welling, 2013] and GANs [Goodfellow et al., 2014]. Despite this negative result, the MLE parameters of a PC p w.r.t. \mathcal{D}_β can be computed in time $\mathcal{O}(|p| \cdot |\mathcal{D}|)$, which is linear w.r.t. the model size as well as the size of the *original* dataset.

Theorem 3. *Let $f_n(\mathbf{x}) = \beta \cdot \mathbb{1}[\mathbf{x} \in \text{supp}(n)] + (1 - \beta) \cdot \mathbb{1}[\mathbf{x} \notin \text{supp}(n)]$ in Alg. 1. Given a deterministic PC p , a boolean dataset \mathcal{D} , and hyperparameter $\beta \in (0.5, 1]$, the set of all flows $\{F_{n,c}(\mathcal{D}_\beta) \mid \forall \text{ edge } (n, c)\}$ w.r.t. the softened dataset \mathcal{D}_β can be computed by Alg. 1 and 2 within $\mathcal{O}(|p| \cdot |\mathcal{D}|)$ time.*

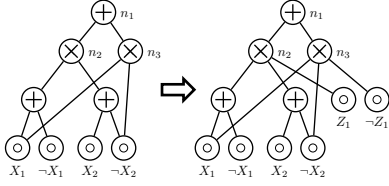


Figure 2: A non-deterministic PC can be modified as an equivalent deterministic PC with hidden variables.

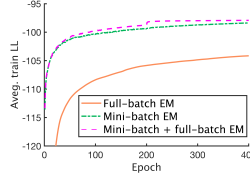


Figure 3: Average train LL on MNIST using different EM updates.

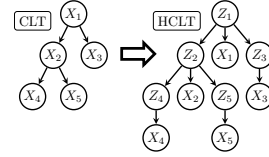


Figure 4: HCLT is constructed by adding hidden variables in a CLT.

Since the MLE parameters (Eq. (3)) w.r.t. \mathcal{D}_β can be computed in $\mathcal{O}(|p|)$ time using the flows, the overall time complexity to compute the MLE parameters is again $\mathcal{O}(|p| \cdot |\mathcal{D}|)$.

Entropy regularization The hope for tractable PC entropy regularization comes from the fact that the entropy of a deterministic PC p can be exactly computed in $\mathcal{O}(|p|)$ time [Vergari et al., 2021, Shih and Ermon, 2020]. However, it is still unclear whether the entropy regularization objective $\text{LL}_{\text{ent}}(\theta; \mathcal{D}, \tau)$ (Eq. (2)) can be tractably maximized. We answer this question with a mixture of positive and negative results: while the objective is multi-modal and the global optimal is hard to find, we propose an efficient algorithm that (i) guarantees convergence to a stationary point, and (ii) achieves high convergence rate in practice. We start with the negative result.

Proposition 1. *There exists a deterministic PC p , a hyperparameter τ , and a dataset \mathcal{D} such that $\text{LL}_{\text{ent}}(\theta; \mathcal{D}, \tau)$ (Eq. (2)) is non-concave and has multiple local maximas.*

Although global optimal solutions are generally infeasible, we propose an efficient algorithm that guarantees to find a stationary point of $\text{LL}_{\text{ent}}(\theta; \mathcal{D}, \tau)$. Specifically, Alg. 3 takes as input a deterministic PC p and all its edge flows w.r.t. \mathcal{D} , and returns a set of learned log-parameters that correspond to a stationary point of the objective.³ In its main loop (lines 4-10), the algorithm alternates between two procedures: (i) compute the entropy of the distribution encoded by every node w.r.t. the current parameters (line 5), and (ii) update PC parameters with regard to the computed entropies (lines 6-10). Specifically, in the parameter update phase (i.e., the second phase), the algorithm traverses every sum unit n in preorder and updates its child parameters by maximizing the entropy regularization objective ($\text{LL}_{\text{ent}}(\theta; \mathcal{D}, \tau)$) with all other parameters fixed. This is done by solving the set of equations in Eq. (5) using Newton’s method (lines 7-8).⁴ In addition to the child nodes’ entropy computed in the first phase, Eq. (5) uses the top-down probability of every unit n (i.e., $\text{node_prob}[n]$), which is progressively updated in lines 9-10.

Theorem 4. *Alg. 3 converges monotonically to a stationary point of $\text{LL}_{\text{ent}}(\theta; \mathcal{D}, \tau)$ (Eq. (2)).*

³Log-domain parameters are used for numerical stability.

⁴Details for solving Eq. (5) is given in Appendix A.2.

Proof. The high-level idea of the proof is to show that the parameter update phase (lines 6-10) optimizes a concave surrogate objective of $\text{LL}_{\text{ent}}(\theta; \mathcal{D}, \tau)$, which is determined by the entropies computed in line 5. Specifically, we show that whenever the surrogate objective is improved, $\text{LL}_{\text{ent}}(\theta; \mathcal{D}, \tau)$ is also improved. Since the surrogate objective is concave, it can be easily optimized. Therefore, Alg. 3 converges to a stationary point of $\text{LL}_{\text{ent}}(\theta; \mathcal{D}, \tau)$. Please refer to the extended version of this paper for a detailed proof. \square

Alg. 3 can be regarded as a EM-like algorithm, where the E-step is the entropy computation phase (line 5) and the M-step is the parameter update phase (lines 6-10). Specifically, the E-step constructs a concave surrogate of the true objective ($\text{LL}_{\text{ent}}(\theta; \mathcal{D}, \tau)$), and the M-step updates all parameters by maximizing the concave surrogate function. Although Thm. 4 provides no convergence rate analysis, the outer loop typically takes 3-5 iterations to converge in practice. Furthermore, Eq. (5) can be solved with high precision in a few (< 10) iterations. Therefore, compared to the computation of all flows w.r.t. \mathcal{D} , which takes $\mathcal{O}(|p| \cdot |\mathcal{D}|)$ time, Alg. 3 takes a negligible $\mathcal{O}(|p|)$ time.

In response to the motivation in Sec. 3, we show that both proposed methods can overcome the imbalanced regularization problem of Laplace smoothing. Again consider the example PC in Fig. 1(a), we conceptually demonstrate that both data softening and entropy regularization will not over-regularize θ_{11} compared to θ_{12} . First, data softening essentially add no prior to the parameters, and only soften the evidences in the dataset. Therefore, it will not over-regularize children with small support sizes. Second, entropy regularization will add a much higher prior to θ_{12} . Suppose $n = 10$, consider maximizing Eq. (2) with an empty dataset (i.e., we maximize $\text{ENT}(p_{n_1})$ directly), the optimal parameters would be $\theta_{11} \approx 0.002$ and $\theta_{12} \approx 0.998$. Therefore, entropy regularization will tend to add a higher prior to children with large support sizes. More fundamentally, the reason why both proposed approaches do not add biased priors to PCs is that they are designed to be model-agnostic, i.e., their definitions as shown in Sec. 2 are independent with the model they apply to.

Empirical evaluation We empirically evaluate both proposed regularization methods on the twenty density estima-

Algorithm 3 PC Entropy regularization

```

1: Input: A deterministic PC  $p$ ; flow  $F_{n,c}(\mathcal{D})$  for every edge  $(n, c)$  in  $p$ ; hyperparameter  $\tau$ .
2: Output: A set of log-parameters,  $\{\varphi_{n,c} : (n, c) \in p\}$ , which are the solution of Eq. (2).
3:  $\forall n$ ,  $\text{node\_prob}[n] \leftarrow 0$ ;  $\text{node\_prob}[n_r] \leftarrow 1$  //  $n_r$  is the root node of  $p$ 
4: while not converge do
5:    $\forall n$ ,  $\text{entropy}[n] \leftarrow$  The entropy of the sub-PC rooted at  $n$  (see Alg. 4)
6:   foreach sum unit  $n$  traversed in preorder (parent before children) do
7:      $d_i \leftarrow F_{n,c_i}(\mathcal{D})/|\mathcal{D}|$ ;  $b = \tau \cdot \text{node\_prob}[n]$  //  $\{c_i\}_{i=1}^{\text{in}(n)}$  is the set of children of  $n$ 
8:     Solve for  $\{\varphi_{n,c_i}\}_{i=1}^{|\text{in}(n)|}$  in the following set of equations ( $y$  is a variable):
           
$$\begin{cases} d_i e^{-\varphi_{n,c_i}} - b \cdot \varphi_{n,c_i} + b \cdot \text{entropy}[c_i] = y & (\forall i \in \{1, \dots, |\text{in}(n)|\}) \\ \sum_{i=1}^{|\text{in}(n)|} e^{\varphi_{n,c_i}} = 1 \end{cases} \quad (5)$$

9:     for each  $c \in \text{in}(n)$  and each  $m \in \text{in}(c)$  do // Update  $\text{node\_prob}$  of grandchildren
10:     $\text{node\_prob}[m] \leftarrow \text{node\_prob}[m] + e^{\varphi_{n,c}} \cdot \text{node\_prob}[n]$ 

```

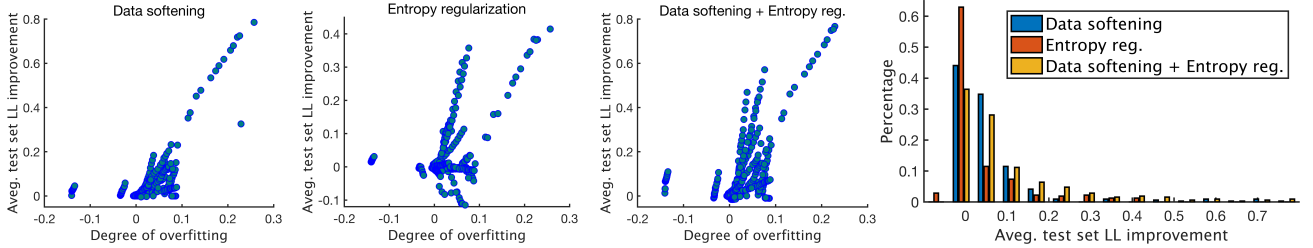


Figure 5: Both data softening and entropy regularization effectively improve the test set log-likelihood (LL) across various datasets [Van Haaren and Davis, 2012] and PC structures [Dang et al., 2020]. LL improvement (higher is better) represents the gain of test set LL compared to Laplace smoothing.

tion datasets [Van Haaren and Davis, 2012]. Since we are only concerned with parameter learning, we adopt PC structures (defined by its DAG) learned by Strudel [Dang et al., 2020]. 16 PCs with different sizes were selected for each of the 20 datasets. For all experiments, we performed a hyperparameter search for all three regularization approaches (Laplace smoothing, data softening, and entropy regularization)⁵ using the validation set and report results on the test set. Please refer to Appendix A.3 for more details.

Results are summarized in Fig. 5. First look at the scatter plots on the left. The x -axis represents the degree of overfitting, which is computed as follows: denote LL_{train} and LL_{val} as the average train and validation log-likelihood under the MLE estimation with Laplace smoothing ($\alpha = 1.0$), the degree of overfitting is defined as $(\text{LL}_{\text{val}} - \text{LL}_{\text{train}})/\text{LL}_{\text{val}}$, which roughly captures how much the dataset/model pair suffers from overfitting. The y -axis represents the improvement on the average test set log-likelihood compared to Laplace smoothing. As demonstrated by the scatter plots, despite a few outliers, both proposed regularization methods steadily improve the test set LL over various datasets and PC structures, and the LL improvements are positively correlated with the degree of overfitting. Further-

more, as shown by the last scatter plot and the histogram plot, when combining data softening and entropy regularization, the LL improvement becomes much higher compared to using the two regularizers individually.

4.3 REGULARIZING NON-DETERMINISTIC PCS

By viewing every non-deterministic PC as a deterministic PC with additional hidden variables (Sec. 4.1), the regularization techniques developed in Sec. 4.2 can be directly adapted. Specifically, data softening can be regarded as injecting noise in both observed and hidden variables. Since the dataset provides no information about the hidden variables anyway, data softening essentially still “perturbs” the observed variables only. On the other hand, entropy regularization will have different behaviors when applied to non-deterministic PCs. Specifically, since it is coNP-hard to compute the entropy of a non-deterministic PC [Vergari et al., 2021], it is infeasible to optimize the entropy regularization objective $\text{LL}_{\text{ent}}(\theta; \mathcal{D}, \tau)$ (Eq. (2)). However, we can still regularize the entropy of the distribution encoded by a non-deterministic PC over both of its observed and hidden variables, since explicitly representing the hidden variables renders the PC deterministic (Sec. 4.1).

On the implementation side, data softening is performed

⁵Specifically, $\alpha \in \{0.1, 0.4, 1.0, 2.0, 4.0, 10.0\}$, $\beta \in \{0.9996, 0.999, 0.996\}$, $\tau \in \{0.001, 0.01, 0.1\}$.

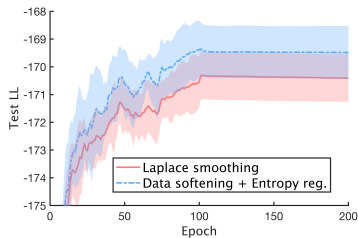


Figure 6: Test set LL over 5 trials on the protein dataset.

by modifying the forward pass of the algorithm used to compute expected flows (i.e., Alg. 5 and 6 in the Appendix). Entropy regularization is again performed by Alg. 3 at the M-step of each min-batch/full-batch EM update, except that the input flows (i.e., F) are replaced by the corresponding expected flows (i.e., EF).

Empirical evaluation We use a simple yet effective PC structure, hidden Chow-Liu Tree (HCLT), as demonstrated in Fig. 4. Specifically, on the left is a Bayesian network representation of a Chow-Liu Tree (CLT) [Chow and Liu, 1968] over 5 variables. For any CLT over variables $\{X_i\}_{i=1}^k$, we can modify it as a HCLT through the following steps. First, we introduce a set of k latent variables $\{Z_i\}_{i=1}^k$. Next, we replace all observed variables in the CLT with its corresponding latent variable (i.e., $\forall i, X_i$ is replaced by Z_i). Finally, we add an edge from every latent variable to its corresponding observed variable (i.e., $\forall i$, add an edge $Z_i \rightarrow X_i$). The HCLT structure is then compiled into a PC that encodes the same probability distribution. We used the hybrid mini-batch + full-batch EM as described in Sec. 4.1. For all experiments, we trained the PCs with 100 mini-batch EM epochs and 100 full-batch EM epochs. Please refer to Appendix A.4 for hyperparameters related to the HCLT structure. Similar to Sec. 4.2, we perform hyperparameter search for all methods using the validation set, and report results on the test set.

We first examine the performance on a protein sequence dataset [Russ et al., 2020] that suffers from severe overfitting. Specifically, the training LL is typically above -100 while the validation and test set LL are around -170 . Fig. 6 shows the test LL for Laplace smoothing and the hybrid regularization approach as training progresses. With the help of data softening and entropy regularization, we were able to obtain consistently higher test set LL. Next, we compare our HCLT model (with regularization) with the state-of-the-art PSDD (Strudel [Dang et al., 2020] and LearnPSDD [Liang et al., 2017]) and SPN (EinSumNet [Peharz et al., 2020a], LearnSPN [Gens and Pedro, 2013], ID-SPN [Rooshenas and Lowd, 2014], and RAT-SPN [Peharz et al., 2020b]) learning algorithms. With proper regularization, HCLT outperformed all baselines in 10 out of 20 datasets. Comparing with individual baselines, HCLT outperforms both PSDD learners on all datasets; HCLT achieved higher log-likelihood on 18, 19, 10, and 17 datasets compared to EinSumNet, LearnSPN, ID-SPN, and RAT-SPN, respectively.

Table 1: Test set log-likelihood in 20 density estimation benchmarks. We compare our method (HCLT) with the best performance (Best PSDD) over 2 deterministic PC learner: Strudel [Dang et al., 2020] and LearnPSDD [Liang et al., 2017] as well as the best performance (Best SPN) over 4 SPN learning algorithms: EinSumNet [Peharz et al., 2020a], LearnSPN [Gens and Pedro, 2013], ID-SPN [Rooshenas and Lowd, 2014], and RAT-SPN [Peharz et al., 2020b]. With the help of data softening and entropy regularization ($\alpha = 0.1$, $\beta = 0.002$, and $\tau = 0.001$), HCLT achieved the best performance over 10 out of 20 datasets. All experiments for HCLT were repeated 5 times, and the average and standard deviation are reported.

Dataset	HCLT	Best PSDD	Best SPN
accidents	-26.74 \pm 0.03	-28.29	-26.98
jester	-52.46 \pm 0.01	-54.63	-52.56
ad	-16.07 \pm 0.06	-16.52	-19.00
kdd	-2.18 \pm 0.00	-2.17	-2.12
baudio	-39.77 \pm 0.01	-41.51	-39.79
kosarek	-10.66 \pm 0.01	-10.98	-10.60
bbc	-251.04 \pm 1.19	-258.96	-248.33
msnbc	-6.05 \pm 0.01	-6.04	-6.03
bnetflix	-56.27 \pm 0.01	-58.53	-56.36
msweb	-9.98 \pm 0.05	-9.93	-9.73
book	-33.83 \pm 0.01	-35.77	-34.14
nltes	-5.99 \pm 0.01	-6.03	-6.01
c20ng	-153.40 \pm 3.83	-160.43	-151.47
plants	-14.26 \pm 0.16	-13.49	-12.54
cr52	-86.26 \pm 3.67	-92.38	-83.35
pumbs*	-23.64 \pm 0.25	-25.28	-22.40
cwebkb	-152.77 \pm 1.07	-160.5	-151.84
tmovie	-50.81 \pm 0.12	-55.41	-51.51
dna	-79.05 \pm 0.17	-82.03	-81.21
retail	-10.84 \pm 0.01	-10.90	-10.85

5 CONCLUSIONS

This paper proposes two model-agnostic distribution regularization techniques: data softening and entropy regularization. While both methods are infeasible for many machine learning models, we theoretically show that they can be efficiently implemented when applied to probabilistic circuits. On the empirical side, we show that both proposed regularizers consistently improve the generalization performance over a wide variety of PC structures and datasets.

ACKNOWLEDGEMENTS

This work is partially supported by NSF grants #IIS-1943641, #IIS-1633857, #CCF-1837129, DARPA grant #N66001-17-2-4032, a Sloan Fellowship, Intel, and Facebook.

References

Greg M Allenby and Peter E Rossi. Hierarchical bayes models. *The handbook of marketing research: Uses, misuses, and future advances*, pages 418–440, 2006.

- Hei Chan and Adnan Darwiche. On the revision of probabilistic beliefs using uncertain evidence. *Artificial Intelligence*, 163(1):67–90, 2005.
- Arthur Choi and Adnan Darwiche. On relaxing determinism in arithmetic circuits. In *International Conference on Machine Learning*, pages 825–833. PMLR, 2017.
- YooJung Choi, Meihua Dang, and Guy Van den Broeck. Group fairness by probabilistic modeling with latent fair decisions. *arXiv preprint arXiv:2009.09031*, 2020a.
- YooJung Choi, Golnoosh Farnadi, Behrouz Babaki, and Guy Van den Broeck. Learning fair naive bayes classifiers by discovering and eliminating discrimination patterns. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10077–10084, 2020b.
- YooJung Choi, Antonio Vergari, and Guy Van den Broeck. Probabilistic circuits: A unifying framework for tractable probabilistic models. 2020c.
- CKCN Chow and Cong Liu. Approximating discrete probability distributions with dependence trees. *IEEE transactions on Information Theory*, 14(3):462–467, 1968.
- Meihua Dang, Antonio Vergari, and Guy Van den Broeck. Strudel: Learning structured-decomposable probabilistic circuits. *arXiv preprint arXiv:2007.09331*, 2020.
- Meihua Dang, Pasha Khosravi, Yitao Liang, Antonio Vergari, and Guy Van den Broeck. Juice: A julia package for logic and probabilistic circuits. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence (Demo Track)*, 2021.
- Adnan Darwiche. A differential approach to inference in bayesian networks. *Journal of the ACM (JACM)*, 50(3):280–305, 2003.
- Adnan Darwiche. *Modeling and reasoning with Bayesian networks*. Cambridge university press, 2009.
- Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
- Yihao Feng, Dilin Wang, and Qiang Liu. Learning to draw samples with amortized stein variational gradient descent. *arXiv preprint arXiv:1707.06626*, 2017.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- Robert Gens and Domingos Pedro. Learning the structure of sum-product networks. In *International conference on machine learning*, pages 873–880. PMLR, 2013.
- Ethan Goan and Clinton Fookes. Bayesian neural networks: An introduction and survey. In *Case Studies in Applied Bayesian Data Science*, pages 45–87. Springer, 2020.
- Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014.
- Yves Grandvalet and Yoshua Bengio. Entropy regularization., 2006.
- David Heckerman. A tutorial on learning with bayesian networks. *Innovations in Bayesian networks*, pages 33–82, 2008.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- Richard C Jeffrey. *The logic of decision*. University of Chicago press, 1990.
- Pasha Khosravi, YooJung Choi, Yitao Liang, Antonio Vergari, and Guy Van den Broeck. On tractable computation of expected predictions. In *Advances in Neural Information Processing Systems 32 (NeurIPS)*, dec 2019.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- 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.
- Yitao Liang and Guy Van den Broeck. Towards compact interpretable models: Shrinking of learned probabilistic sentential decision diagrams. In *IJCAI 2017 Workshop on Explainable Artificial Intelligence (XAI)*, August 2017. URL <http://starai.cs.ucla.edu/papers/LiangXAI17.pdf>.
- Yitao Liang, Jessa Bekker, and Guy Van den Broeck. Learning the structure of probabilistic sentential decision diagrams. In *Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence (UAI)*, 2017.
- Jun Mei, Yong Jiang, and Kewei Tu. Maximum a posteriori inference in sum-product networks. 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. *arXiv preprint arXiv:1901.03704*, 2019.

- Rong Pan, Yun Peng, and Zhongli Ding. Belief update in bayesian networks using uncertain evidence. In *2006 18th IEEE International Conference on Tools with Artificial Intelligence (ICTAI'06)*, pages 441–444. IEEE, 2006.
- Judea Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier, 2014.
- Robert Peharz, Robert Gens, and Pedro Domingos. Learning selective sum-product networks. In *LTPM workshop*, volume 32, 2014.
- Robert Peharz, Robert Gens, Franz Pernkopf, and Pedro Domingos. On the latent variable interpretation in sum-product networks. *IEEE transactions on pattern analysis and machine intelligence*, 39(10):2030–2044, 2016.
- Robert Peharz, Steven Lang, Antonio Vergari, Karl Stelzner, Alejandro Molina, Martin Trapp, Guy Van den Broeck, Kristian Kersting, and Zoubin Ghahramani. Einsum networks: Fast and scalable learning of tractable probabilistic circuits. In *International Conference on Machine Learning*, pages 7563–7574. PMLR, 2020a.
- Robert Peharz, Antonio Vergari, Karl Stelzner, Alejandro Molina, Xiaoting Shao, Martin Trapp, Kristian Kersting, and Zoubin Ghahramani. Random sum-product networks: A simple and effective approach to probabilistic deep learning. In *Uncertainty in Artificial Intelligence*, pages 334–344. PMLR, 2020b.
- Luis Perez and Jason Wang. The effectiveness of data augmentation in image classification using deep learning. *arXiv preprint arXiv:1712.04621*, 2017.
- 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.
- Amirmohammad Rooshenas and Daniel Lowd. Learning sum-product networks with direct and indirect variable interactions. In *International Conference on Machine Learning*, pages 710–718. PMLR, 2014.
- William P Russ, Matteo Figliuzzi, Christian Stocker, Pierre Barrat-Charlaix, Michael Socolich, Peter Kast, Donald Hilvert, Remi Monasson, Simona Cocco, Martin Weigt, et al. An evolution-based model for designing chorismate mutase enzymes. *Science*, 369(6502):440–445, 2020.
- Yujia Shen, Arthur Choi, and Adnan Darwiche. Tractable operations for arithmetic circuits of probabilistic models. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 3943–3951. Citeseer, 2016.
- Andy Shih and Stefano Ermon. Probabilistic circuits for variational inference in discrete graphical models. In *Advances in Neural Information Processing Systems 33 (NeurIPS)*, december 2020. URL <https://cs.stanford.edu/~andyshih/assets/pdf/SEneurips20.pdf>.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016.
- Jan Van Haaren and Jesse Davis. Markov network structure learning: A randomized feature generation approach. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 26, 2012.
- Antonio Vergari, Nicola Di Mauro, and Floriana Esposito. Simplifying, regularizing and strengthening sum-product network structure learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 343–358. Springer, 2015.
- Antonio Vergari, YooJung Choi, Anji Liu, Stefano Teso, and Guy Van den Broeck. A compositional atlas of tractable circuit operations: From simple transformations to complex information-theoretic queries. *arXiv preprint arXiv:2102.06137*, 2021.
- Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103, 2008.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.

Algorithm 5 Forward pass (expected flows)

```
1: Input: A non-deterministic PC  $p$ ; sample  $\mathbf{x}$ 
2: Output:  $\text{value}[n] := (\mathbf{x} \in \text{supp}(n))$  for each unit  $n$ 
3: foreach  $n$  traversed in postorder do
4:   if  $n$  isa input unit then  $\text{value}[n] \leftarrow f_n(\mathbf{x})$ 
5:   elif  $n$  isa product unit then
6:      $\text{value}[n] \leftarrow \prod_{c \in \text{in}(n)} \text{value}[c]$ 
7:   else  $n$  is a sum unit
8:      $\text{value}[n] \leftarrow \sum_{c \in \text{in}(n)} \theta_{n,c} \cdot \text{value}[c]$ 
```

Algorithm 6 Backward pass (expected flows)

```
1: Input: A non-deterministic PC  $p$ ;  $\forall n, \text{value}[n]$ 
2: Output:  $\text{eflow}[n, c] := \mathbb{E}_{\mathbf{z} \in p_c(\cdot | \mathbf{x}; \theta)}((\mathbf{x}, \mathbf{z}) \in (\gamma_n \cap \gamma_c))$  for
   each pair  $(n, c)$ , where  $n$  is a sum unit and  $c \in \text{in}(n)$ 
3:  $\forall n, \text{context}[n] \leftarrow 0$ ;  $\text{context}[n_r] \leftarrow \text{value}[n_r]$ 
4: foreach sum unit  $n$  traversed in preorder do
5:   foreach  $m \in \text{pa}(n)$  do (denote  $g \leftarrow \text{pa}(m)$ )
6:      $\mathbf{f} \leftarrow \frac{\text{value}[m]}{\text{value}[g]} \cdot \text{context}[g] \cdot \theta_{g,m}$ 
7:      $\text{context}[n] += \mathbf{f}$ ;  $\text{flow}[g, m] = \mathbf{f}$ 
```

A METHOD AND EXPERIMENT DETAILS

A.1 SOFTEN NON-BOOLEAN DATASETS

As a direct extension of softening boolean datasets, datasets with categorical variables can be similarly softened. Suppose X is a categorical variable with k categories. For an assignment $x = j$, we can soften it as follows

$$\begin{cases} P(x = i) = \frac{1-\beta}{k} & (i \neq j), \\ P(x = j) = \beta. \end{cases}$$

To compute the flow $F_{n,c}(\mathcal{D}_\beta)$ w.r.t. a softened categorical dataset, we can again adopt Alg. 1 and 2 by choosing

$$f_n(\mathbf{x}) = \beta \cdot \mathbb{1}[\mathbf{x} \in \text{supp}(n)] + \frac{1-\beta}{k} \cdot \mathbb{1}[\mathbf{x} \notin \text{supp}(n)].$$

A.2 SOLVING EQUATION 5

Denote $\gamma_{c_i} := \text{entropy}[c_i]$, our goal is to solve the following set of equations:

$$\begin{cases} d_i e^{-\varphi_{n,c_i}} - b \cdot \varphi_{n,c_i} + b \cdot \gamma_{c_i} = y & (\forall i \in \{1, \dots, |\text{in}(n)|\}), \\ \sum_{i=1}^{|\text{in}(n)|} e^{\varphi_{n,c_i}} = 1. \end{cases}$$

We break down the problem by iteratively solve for $\{\varphi_{n,c_i}\}_{i=1}^{|\text{in}(n)|}$ and y , respectively.

- Solve for y . Given variables $\{\varphi_{n,c_i}\}_{i=1}^{|\text{in}(n)|}$, we update y as

$$y = \frac{1}{|\text{in}(n)|} \sum_{i=1}^{|\text{in}(n)|} d_i e^{-\varphi_{n,c_i}} - b \cdot \varphi_{n,c_i} + b \cdot \gamma_{c_i}.$$

- Solve for $\{\varphi_{n,c_i}\}_{i=1}^{|\text{in}(n)|}$. Given y , we first update each φ_{n,c_i} individually by solving the equation

$$d_i e^{-\varphi_{n,c_i}} - b \cdot \varphi_{n,c_i} + b \cdot \gamma_{c_i} = y.$$

Specifically, this is done by iterative Newton method update:

$$\varphi_{n,c_i} += \frac{\frac{d_i}{\varphi_{n,c_i}} + b \cdot (\gamma_{c_i} - \varphi_{n,c_i}) + y}{\frac{d_i}{\varphi_{n,c_i}} + b}$$

After one Newton method update step for every parameter in $\{\varphi_{n,c_i}\}_{i=1}^{|\text{in}(n)|}$, we enforce the constraint $\sum_{i=1}^{|\text{in}(n)|} e^{\varphi_{n,c_i}} = 1$ by

$$\varphi_{n,c_i} -= \log \left(\sum_{i=1}^{|\text{in}(n)|} e^{\varphi_{n,c_i}} \right).$$

A.3 DETAILS OF THE EXPERIMENTS ON DETERMINISTIC PCS

PC structures For each dataset, we adopt 16 PCs by running Strudel [Dang et al., 2020] for $\{1000, 1200, 1400, \dots, 4000\}$ iterations except for the dataset “dna”, which we ran Strudel for $\{50, 100, 150, \dots, 800\}$ iterations since the learning algorithm takes significantly longer for this dataset.

Hyperparameters We always perform hyperparameter search using the validation set, and report the final performance on the test set. Whenever we use data softening or entropy regularization, we also add pseudocount $\alpha = 1$ since it yields better performance.

Server specifications All our experiments were run on a server with 72 CPUs, 512G Memory, and 2 TITAN RTX GPUs.

A.4 DETAILS OF THE EXPERIMENTS ON NON-DETERMINISTIC PCS

The HCLT structure For the experiments on the twenty datasets, we set the hidden size of the HCLT structure as 12, i.e., every latent variable Z is a categorical variable with 12 categories. Additionally, following Dang et al. [2020], Liang et al. [2017], we learn a mixture of 4 HCLTs to achieve better performance. For the protein sequence dataset, we adopted a mixture of 2 HCLTs with hidden size 32.

Detailed results As an extension of Table 1, Table 2 provides the average test set log-likelihood for all adopted baselines.

Algorithm 4 PC entropy

1: **Input:** A deterministic PC p
2: **Output:** $\text{entropy}[n] := \text{ENT}(p_n)$ for every unit n
3: **foreach** n traversed in postorder **do**
4: **if** n **isa** input unit **then** $\text{entropy}[n] = \text{ENT}(p_n)$ //entropy of the input distribution
5: **elif** n **isa** product unit **then** $\text{entropy}[n] = \sum_{c \in \text{in}(n)} \text{entropy}[c]$
6: **else** // n is a sum unit **then**

$$\text{entropy}[n] = - \sum_{c \in \text{in}(n)} \theta_{n,c} \log \theta_{n,c} + \sum_{c \in \text{in}(n)} \theta_{n,c} \cdot \text{entropy}[c]$$

Table 2: Full results on the 20 density estimation benchmarks. As an extension of Table 1, we report the average test-set log-likelihood of all baselines: Strudel [Dang et al., 2020], LearnPSDD [Liang et al., 2017], EinSumNet [Peharz et al., 2020a], LearnSPN [Gens and Pedro, 2013], ID-SPN [Rooshenas and Lowd, 2014], and RAT-SPN [Peharz et al., 2020b].

Dataset	HCLT	EiNet	LearnSPN	ID-SPN	RAT-SPN	Strudel	LearnPSDD
accidents	-26.78	-35.59	-40.50	-26.98	-35.48	-29.46	-28.29
ad	-16.04	-26.27	-19.73	-19.00	-48.47	-16.52	-20.13
baudio	-39.77	-39.87	-40.53	-39.79	-39.95	-42.26	-41.51
bbc	-250.07	-248.33	-250.68	-248.93	-252.13	-258.96	-260.24
bnetflix	-56.28	-56.54	-57.32	-56.36	-56.85	-58.68	-58.53
book	-33.84	-34.73	-35.88	-34.14	-34.68	-35.77	-36.06
c20ng	-151.92	-153.93	-155.92	-151.47	-152.06	-160.77	-160.43
cr52	-84.67	-87.36	-85.06	-83.35	-87.36	-92.38	-93.30
cwebkb	-153.18	-157.28	-158.20	-151.84	-157.53	-160.50	-161.42
dna	-79.33	-96.08	-82.52	-81.21	-97.23	-87.10	-83.02
jester	-52.45	-52.56	-75.98	-52.86	-52.97	-55.30	-54.63
kdd	-2.18	-2.18	-2.18	-2.13	-2.12	-2.17	-2.17
kosarek	-10.66	-11.02	-10.98	-10.60	-10.88	-10.98	-10.99
msnbc	-6.05	-6.11	-6.11	-6.04	-6.03	-6.05	-6.04
msweb	-9.90	-10.02	-10.25	-9.73	-10.11	-10.19	-9.93
nltes	-6.00	-6.01	-6.11	-6.02	-6.01	-6.06	-6.03
plants	-14.31	-13.67	-12.97	-12.54	-13.43	-13.72	-13.49
pumbs*	-23.32	-31.95	-24.78	-22.40	-32.53	-25.28	-25.40
tmovie	-50.69	-51.70	-52.48	-51.51	-53.63	-59.47	-55.41
treail	-10.84	-10.91	-11.04	-10.85	-10.91	-10.90	-10.92