

Scalable Learning of Probabilistic Circuits

Renato Lui Geh

THESIS PRESENTED TO THE
INSTITUTE OF MATHEMATICS AND STATISTICS
OF THE UNIVERSITY OF SÃO PAULO
IN PARTIAL FULFILLMENT
OF THE REQUIREMENTS
FOR THE DEGREE OF
MASTER OF SCIENCE

Program: Computer Science

Advisor: Professor Denis Deratani Mauá

This work was supported by CNPq grant #133787/2019-2,
CAPES grant #88887.339583/2019-00 and EPECLIN FM-USP.

São Paulo
November 1, 2021

Scalable Learning of Probabilistic Circuits

Renato Lui Geh

This is the original version of the thesis
prepared by candidate Renato Lui Geh, as
submitted to the Examining Committee.

I hereby authorize the total or partial reproduction and publishing of this work for educational ou research purposes, as long as properly cited.

Acknowledgements

Abstract

Renato Lui Geh. **Scalable Learning of Probabilistic Circuits**. Thesis (Master's). Institute of Mathematics and Statistics, University of São Paulo, São Paulo, 2021.

The rising popularity of generative models together with the growing need for flexible and exact inferences has motivated the machine learning community to look for expressive yet tractable probabilistic models. Probabilistic circuits (PCs) are a family of tractable probabilistic models capable of answering a wide range of queries exactly and in polynomial time. Their operational syntax in the form of a computational graph and their principled probabilistic semantics allow their parameters to be estimated by the highly scalable and efficient optimization techniques used in deep learning. Importantly, tractability is tightly linked to constraints on their underlying graph: by enforcing certain structural assumptions, queries like marginals, *maximum a posteriori* or entropy become linear time computable while still retaining great expressivity. While inference is usually straightforward, learning PCs that both obey the needed structural restrictions and exploit their expressive power has proven a challenge. Current state-of-the-art structure learning algorithms for PCs can be roughly divided into three main categories. Most learning algorithms seek to generate a usually tree-shaped circuit from recursive decompositions on data, often through clustering and costly statistical (in)dependence tests, which can become prohibitive in higher dimensional data. Alternatively, other approaches involve constructing an intricate network by growing an initial circuit through structural preserving iterative methods. Besides depending on a sufficiently expressive initial structure, these can possibly take several minutes per iteration and many iterations until visible improvement. Lastly, other approaches involve randomly generating a probabilistic circuit by some criterion. Although usually less performant compared to other methods, random PCs are orders of magnitude more time efficient. With this in mind, this dissertation aims to propose fast and scalable random structure learning algorithms for PCs from two different standpoints: from a logical point of view, we efficiently construct a highly structured binary PC that takes certain knowledge in the form of logical constraints and scalably translate them into a probabilistic circuit; from the viewpoint of data guided structure search, we propose hierarchically building PCs from random hyperplanes. We empirically show that either approach is competitive against state-of-the-art methods of the same class, and that their performance can be further boosted by simple ensemble strategies.

Keywords: Probabilistic circuits. Machine learning. Probabilistic models.

Resumo

Renato Lui Geh. **Aprendizado Escalável de Circuitos Probabilísticos**. Dissertação (Mestrado). Instituto de Matemática e Estatística, Universidade de São Paulo, São Paulo, 2021.

Palavras-chave: Circuitos probabilísticos. Aprendizado de máquina. Modelos probabilísticos.




Nomenclature

List of Symbols

σ Sigmoid function

List of Figures

- | | | |
|-----|--|----|
| 2.1 | A probabilistic circuit (a) and 3 of the 12 possible induced subcircuits (b). | 9 |
| 2.2 | Decomposable but non-smooth (a), smooth but non-decomposable (b), and smooth and decomposable (c) circuits. | 10 |
| 2.3 | A vtree (a) defining an order (A, B, C, D) , a 2-standard structure decomposable probabilistic circuit that respects the vtree (b), and a 2-standard decomposable probabilistic circuit that does not (c). | 14 |
| 2.4 | Two smooth, structure decomposable and deterministic logic circuits encoding the same logic constraint $\phi \equiv (A \wedge B) \vee (\neg C \wedge D)$ for a balanced (a) and a right-linear (b) vtree. In (a), a circuit evaluation for an assignment, with each node value in the bottom-up evaluation pass shown inside nodes. | 17 |
| 3.1 | LEARNSPN assigns either rows (a) or columns (b) for sum and product nodes respectively. For sums, their edge weights are set proportionally to the assignments. For product children, scopes are defined by which columns are assigned to them. | 24 |
| 3.2 | The pairwise (in)dependence graph where each node is a variable. In (a) we show the full graph, computing independence tests for each pair of variables in $\mathcal{O}(m^2)$. However, it suffices to compute for only the connected components (b), saving up pairwise computation time for reachable nodes. The resulting product node and scope partitioning is shown in (c). | 25 |

3.3	Two iterations of ID-SPN, where the contents inside the dashed line are Markov networks. The red color indicates that a node has been chosen as the best candidate for an extension with EXTENDID. Although here we only extend input nodes, inner nodes can in fact be extended as well. . . .	27
3.4	The fully connected correlation graph (a) with weights as the pairwise correlation measurements for each pair of variables; the maximum spanning tree for determining decompositions (b); and the mixture of decompositions (c). Colors in (b) match their partitionings in (c).	28
3.5	Snapshots of four iterations from running the vtree top-down learning strategy with pairwise mutual information. Each iteration shows a variable partitioning, the cut-set that minimizes the average pairwise mutual information as black edges, and the subsequent (partial) vtree. The algorithm finishes when all partitions are singletons.	32
3.6	Snapshots from running the vtree bottom-up learning strategy with pairwise mutual information. Snapshots show pairings of two vtrees, with edges between partitions joined into a single edge whose weight is the average pairwise mutual information of all collapsed edges. In black are edges that correspond to the matchings that maximize the average pairwise mutual information. The algorithm finishes when all vtrees have been joined together into a single tree.	33
3.7	SPLIT (left) and CLONE (right) operations for growing a circuit when $m = 1$. Nodes and edges highlighted in red show the modified structure. In both cases smoothness, (structure) decomposability and determinism are evidently preserved.	34
3.8	A vtree (middle) and probabilistic circuit (right) compiled from a Chow-Liu Tree (left). Each conditional probability $p(Y X)$ is encoded as a (deterministic) sum node where each of the two children sets Y to 0 or 1. Colors in the CLT indicate the variables in the PC, while vtree inner node colors match with product nodes that respect them. Edges in red indicate the induced subcircuit activated on assignment $\{A = 1, B = 0, C = 1, D = 0\}$	37
3.9	A RAT-SPN generated from parameters $d = 3$, $r = 2$, $s = 2$ and $l = 2$. Nodes within a  belong to inner region nodes,  to partitions and  to leaf regions; dashed lines (and node colors) indicate different PC layers. Scope of region nodes are shown in curly braces.	45

List of Algorithms

1	EVI	11
2	MAR	11
3	MAP	12
4	ARGMAP	12
5	EXP	20
6	LEARNSPN	24
7	EXTENDID	26
8	ID-SPN	27
9	PROMETHEUS	29
10	LEARNSDD	35
11	INITIALSTRUDEL	39
12	STRUDEL	40
13	COMPILEREGIONGRAPH	43
14	RAT-SPN	44

List of Examples

2.1.1	Gaussians as probabilistic circuits	5
2.1.2	Gaussian mixture models as probabilistic circuits	6
2.1.3	Factors as probabilistic circuits	7
2.2.1	Naïve Bayes as probabilistic circuits	13
2.2.2	Hidden Markov models as probabilistic circuits	15
2.3.1	BDDs as logic circuits	18
2.3.2	Embedding certain knowledge in probabilistic circuits	19
2.3.3	Computing the probability of logical events	20

List of Remarks

2.1.1	On operators and tractability	8
2.3.1	On applications of probabilistic circuits	21
3.1.1	On variations of divide-and-conquer learning	30
3.2.1	On the choice of initial circuits	41
3.3.1	On parameter learning of probabilistic circuits	45

Notation

We use the following notation throughout the work. Random variables are written in upper case (e.g. X , Y) and their values in lower case (e.g. x , y). We write \mathcal{X} as the sample space and for independence we use \perp to indicate that two variables X and Y are statistically independent, i.e. $X \perp Y$. We identify propositional variables with 0/1-valued random variables, and use them interchangeably. Sets of variables and their joint values are written in boldface (e.g. \mathbf{X} , \mathbf{x}). Given a Boolean formula f , we write $\langle f \rangle$ to denote its semantics, i.e. the Boolean function represented by f . For Boolean formulas f and g , we write $f \equiv g$ if they are logically equivalent, that is, if $\langle f \rangle = \langle g \rangle$; we abuse notation and write $\phi \equiv f$ to indicate that $\phi = \langle f \rangle$ for a Boolean function ϕ . We use the notation $[a..b]$, with $b \geq a$ to denote the integer set $\{a, a+1, \dots, b\} \subset \mathbb{Z}$. Similarly, we use $[b]$ as an equivalent for $[1..b]$. We denote the Iverson bracket as $\llbracket \phi \rrbracket$, i.e. a function that returns 1 if ϕ is true and 0 otherwise. In the context of graph theory, we use sans serif letters for graph nodes (e.g. N , S , P , L) and bold variants for sets of nodes (e.g. \mathbf{N} , \mathbf{S} , \mathbf{P} , \mathbf{L}). We call $\text{Ch}(\mathbf{N})$ the set of all children of a node N , $\text{Pa}(\mathbf{N})$ as the set of all parents, and $\text{Desc}(\mathbf{N})$ the set of all descendants.

Contents

1	Introduction	1
1.1	Contributions and Dissertation Outline	3
2	Probabilistic Circuits	5
2.1	Distributions as Computational Graphs	5
2.2	Deciding What to Constraint	10
2.2.1	Basic Queries	10
2.2.2	Complex Queries	13
2.3	Probabilistic Circuits as Knowledge Bases	16
2.3.1	From Certainty...	16
2.3.2	...to Uncertainty	18
3	Learning Probabilistic Circuits	23
3.1	Divide-and-Conquer Learning	23
3.1.1	LEARNSPN	24
3.1.2	ID-SPN	26
3.1.3	PROMETHEUS	28
3.2	Incremental Learning	31
3.2.1	LEARNPSDD	31
3.2.2	STRUDEL	37
3.3	Random Learning	42
3.3.1	RAT-SPN	42
3.3.2	XPCs	45
4	Scalable Learning of Probabilistic Circuits	47
4.1	A Logical Perspective	47
4.1.1	SAMPLEPSDD	47
4.2	A Data Perspective	47
4.2.1	LEARNRP	47

Appendices

A Appendix	49
A.1 Proofs	49

Annexes

References	55
-------------------	-----------

1

Introduction

When reasoning about the world, rarely can we find a realistic model that perfectly subsumes all of the needed relationships for flawless prediction. As such, the presence of a reliable uncertainty quantifier in intelligent systems is essential in developing performant yet diagnostible agents. This is made explicitly clear in the case of high-risk settings, such as autonomous vehicles or automated power plant systems, where a wrong prediction could cause disastrous consequences. A well known example in the case of the former is obstacle avoidance: while the agent should be capable of accurately identifying obstructions in its way during normal conditions, so should it be able to identify its own lack of confidence in high uncertainty situations like ones brought by environmental factors, such as severe blizzard or heavy rain. In these situations where predictions are highly unreliable, the safest option might be for the agent to first identify its uncertainty, and second to reach out for human help. Other interesting applications of uncertainty quantification include out-of-distribution detection, which in our previous example could be visualized as the agent identifying a human's driving as abnormally irregular (due to inebriation, infirmity, etc.) and appropriately taking control of the vehicle before a potential accident takes place.

One popular approach to quantifying uncertainty is through probability theory. By abstracting the world as a probability distribution with a finite number of observable random variables that encode a possibly incomplete knowledge of the environment, we are, in theory, able to answer a diverse set of complex queries as long as we have access to the (approximate) true data distribution. Practice is far different from theory, however, as most machine learning models either lie within a very limited range in the tractability spectrum in terms of inference (Y. CHOI, VERGARI, and VAN DEN BROECK, 2020) or are too simplistic for complex real-world problems. Besides, although the majority of recent advances in deep learning claim some probabilistic meaning from the model's output, they are often uncalibrated distributions, a result of focusing on maximizing predictive accuracy at the expense of predictive uncertainty (GUO *et al.*, 2017; OVADIA *et al.*, 2019; CHERNIKOVA *et al.*, 2019), ultimately producing overconfident and peculiar results (SZEGEDY *et al.*, 2013; WEI *et al.*, 2018; SU *et al.*, 2019; CHERNIKOVA *et al.*, 2019). Crucially, mainstream deep models (i.e. standard neural networks) usually optimize a conditional distribution over the to-be-predicted random variables – and are thus often called *discriminative* models – and do not model the actual joint distribution of the data, limiting inference capabilities and

uncertainty estimation.

In contrast, *generative* models seek to extract information from the joint (in varying capacities), and have lately seen a sharp rise in interest within deep learning. Despite this, most popular models do not admit either exact or tractable querying of key inference scenarios. For instance, although Generative Adversarial Networks (GANs) allow for efficient sampling (GOODFELLOW *et al.*, 2014), basic queries such as likelihood or marginals are outside of their capabilities. Similarly, Normalizing Flows (NF) also permit access to efficient sampling, with the added feature of computing likelihoods, but are severely limited by the base distribution when it comes to discrete data (REZENDE and MOHAMED, 2015; PAPAMAKARIOS *et al.*, 2021) albeit recent works on discretizing NFs have shown empirically good results (LIPPE and GAVVES, 2021; ZIEGLER and RUSH, 2019). Variational Auto-Encoders (VAEs) are (under certain conditions) a generalization of NFs (YU, 2020; GRITSENKO *et al.*, 2019) with known extensions for categorical data (ROLFE, 2017; VAHDAT, MACREADY, *et al.*, 2018; VAHDAT, ANDRIYASH, *et al.*, 2018), but only permit access to sampling and an upper-bound on the likelihood, with the latter available only after solving a complex optimization task (KINGMA and WELLING, 2014).

Despite the impressive achievements of the aforementioned generative models on realistically producing samples consistent with evidence, in none of the previous models are complex queries like structured prediction under partial observations, *maximum a posteriori* (MAP), conditional or marginal probabilities tractable. An obvious alternative would be Probabilistic Graphical Models (PGMs), although they too suffer from intractability when dealing with high treewidth networks (R. DECHTER, 1998; KOLLER and FRIEDMAN, 2009), severely limiting expressivity. Instead, we draw our attention to an expressive class of models that subsumes several families of probabilistic models with tractable inference capabilities.

Probabilistic Circuits (PCs) define a superclass of probabilistic models distinctly specified by recursive compositions of distributions through graphical formalisms. Vaguely speaking, PCs are computational graphs akin to neural networks, but whose network structure and computational units abide by special constraints. Within these specific conditions span a wide range of subclasses, each establishing a distinct set of restrictions on their structure in order to enable different segments within the tractability spectrum. As an example, Sum-Product Networks (SPNs, POON and P. DOMINGOS, 2011) are usually loosely defined over a couple of constraints: namely *smoothness* and *decomposability*, which in turn enables likelihood, marginal and conditional computations. Arithmetic Circuits (ACs, DARWICHE, 2003) add *determinism* to the mix, allowing for tractable computation of MAP probabilities. Similarly, Cutset Networks (C Nets, RAHMAN, KOTHALKAR, *et al.*, 2014) employ the same constraints as ACs, but accept more expressive distributions as part of their computational units. Probabilistic Sentential Decision Diagrams (PSDDs, KISA *et al.*, 2014), Probabilistic Decision Graphs (PDGs, JAEGER, 2004) and And/Or-Graphs (AOGs, Rina DECHTER and MATEESCU, 2007) all require *smoothness* and *determinism*, but also call for a stronger version of *decomposability*, permitting all queries previously mentioned as well as computation of the Kullback-Leibler divergence and expectation between two circuits (Y. CHOI, VERGARI, and BROECK, 2020). Usually, PCs represent the joint distribution of the data, although they are sufficiently expressive for generative *and* discriminative modeling (KHOSRAVI *et al.*, 2019; RASHWAN *et al.*, 2018; ROOSHENAS and LOWD, 2016; GENS and P. DOMINGOS, 2012;

SHAO *et al.*, 2020). In this dissertation though, we shall focus on the generative side of PCs.

While inference is usually straightforward, as we shall see in [Section 2.2](#), learning the structure of PCs so that they obey the needed structural restrictions requires either careful handcrafted architectures (POON and P. DOMINGOS, 2011; CHENG *et al.*, 2014; NATH and P. M. DOMINGOS, 2016) or usually involves running costly (in)dependence tests over most (if not all) variables (GENS and P. DOMINGOS, 2013; JAINI, GHOSE, *et al.*, 2018; VERGARI, MAURO, *et al.*, 2015; DI MAURO *et al.*, 2017), which can become prohibitive in higher dimension data. Alternatively, some learning algorithms resort to structure preserving iterative methods to grow a PC that already initially satisfies desired constraints, adding complexity to the underlying distribution at each iteration (LIANG, BEKKER, *et al.*, 2017; DANG *et al.*, 2020). However, these can take several iterations until visible improvement and often take several minutes for each iteration when the circuit is big. Common techniques used in deep learning for generating scalable architectures for neural network also pose a problem, as the nature of the needed structural constraints make for sparse computational graphs. To circumvent these issues, work on scaling PCs to higher dimensions has focused mainly on random architectures, with competitive results (PEHARZ, VERGARI, *et al.*, 2020; MAURO *et al.*, 2021; R. L. GEH and Denis Deratani MAUÁ, 2021; PEHARZ, LANG, *et al.*, 2020). Apart from the scalability side of random structure generation, usual structure learning algorithms often require grid-search for hyperparameter tuning to achieve top quality performance, which is usually not the case for random algorithms. For the usual data scientist or machine learning practitioner, hyperparameter tuning can become exhaustive, especially if the goal is to analyze and infer from large data, and not to achieve top tier performance on benchmark datasets.

In this dissertation, we propose two scalable structural learning algorithms for probabilistic circuits that are especially suited for large data and fast deployment. They both take advantage of random network generation to quickly construct PCs with little to no need for hyperparameters. The first is effective for constructing PCs from binary data with a highly constrained structure, and thus appropriate when complex querying is needed. The second builds less constrained random PCs, but supports both discrete and continuous data.

1.1 Contributions and Dissertation Outline

We organize this dissertation as follows. We begin [Chapter 2](#) by formally defining probabilistic circuits, conducting a review of some of the structural constraints that we might impose on PCs, as well as what we may gain from them in terms of tractability. We then list existing formalisms that may be viewed as instances of PCs, and what their structure entail in terms of inference power. In [Chapter 3](#), we address existing PC structure learning algorithms, and which guarantees in terms of tractability each give. We cover the two new structure learners in [Chapter 4](#), providing empirical results on their performance. The final chapter is dedicated to summarizing our research contributions and pointing to potential future work in learning PCs.

Our contributions in this dissertation address the following research topics.

Scalably learning PCs directly from background knowledge

In [R. L. GEH and Denis Deratani MAUÁ \(2021\)](#), we provide a learning algorithm for PSDDs that learns a PC directly from background knowledge in the form of logical constraints. The algorithm samples a structure from a distribution of possible PSDDs that are weakly consistent with the logical formula. How weak consistency is depends on a parameter that trades permission of false statements as non zero probability events with circuit complexity. We provide the algorithm and empirical results in Section 4.??.

Using ensembles to strengthen consistency

The PC sampler given by [R. L. GEH and Denis Deratani MAUÁ \(2021\)](#) produces competitive probabilistic models (in terms of likelihood), albeit weak logical models in the sense that it possibly assigns non-zero probability to false variable assignments – as we discuss in Section 4.??, it never assigns zero probability to true statements. By producing many weak models, we not only gain in terms of data fitness, but also consistency: if any one component in the ensemble returns an assignment to be impossible, the whole model should return false.

Random projections to efficiently learn PCs

Usual methods often employ clustering algorithms for constructing convex combinations of computational units. These can take many iterations to converge or require space quadratic in the number of data points. Instead, in Section 4.?? we present linear alternatives based on random projections ([FREUND *et al.*, 2008](#); [DASGUPTA and FREUND, 2008](#)).

2

Probabilistic Circuits

As we briefly mentioned in the last chapter, Probabilistic Circuits (PCs) are conceptualized as computational graphs under special conditions. In this chapter, [Section 2.1](#) to be more precise, we formally define PCs and give an intuition on their syntax, viewing other probabilistic models through the lenses of the PC framework. In [Section 2.2](#), we describe the special structural constraints that give PCs their inference power over other generative models and state which queries (as far as we know) are enabled from each constraint.

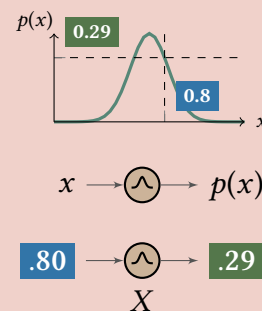
2.1 Distributions as Computational Graphs

Probabilistic circuits are directed acyclic graphs usually recursively defined in terms of their computational units. In its simplest form a PC is a single unit with no outgoing edges whose value corresponds to the result of a function. These are often called *input* nodes, and can take any form as long as its value is tractably computable. More concretely, input nodes typically represent probability density (or mass) functions, although they also support inputs as joint probability density functions of complex non-parametric models as well. To simplify notation, from here on out we shall use the term distribution and probability density (resp. mass) function interchangeably and argue that the input node represents a probability distribution.

Example 2.1.1: Gaussians as probabilistic circuits

Let L_p an input node and $p(X) = \mathcal{N}(X; \mu, \sigma^2)$ a univariate Gaussian distribution. Computing any query on L_p is straightforward: any query on L_p directly translates to p . As an example, suppose $\mu = 0$ and $\sigma^2 = 1$ and we wish to compute $L_p(x = 0.8)$. The probability of this input shall then be

$$L_p(x = 0.8) = \mathcal{N}(x = 0.8; \mu = 0, \sigma^2 = 1) = 0.29.$$



Let L_p a PC input node and denote p as its inherent probability distribution. By definition, any query $f : \mathcal{X} \rightarrow \mathcal{Y}$ which is tractable on p is tractable on L_p . We shall denote $L_p(\mathbf{X}) = p(\mathbf{X})$, and often omit p when its explicit form is not needed. Evidently, a single input node lacks the expressivity for modeling complex models, otherwise we would have just used the input distribution as a standalone model. The expressiveness of PCs comes from recursively combining distributions into complex functions. This can be done through computational units that either compute convex combinations or products of their children. Let us first look at convex combinations, known in the literature as *sum* nodes.

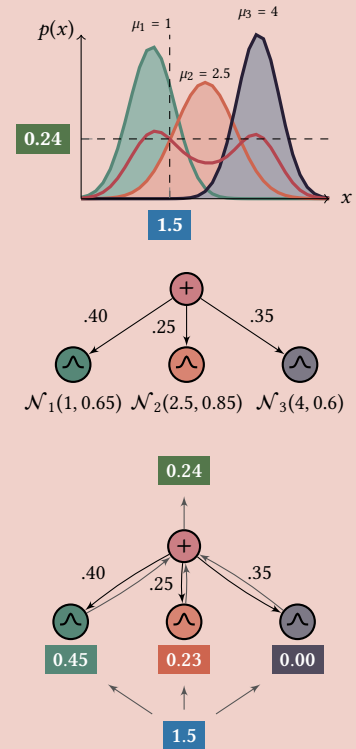
Let S be a PC sum node, and denote by $\text{Ch}(S)$ the children nodes of S . For every edge \overrightarrow{SC} coming out of S and going to C , we attribute a weight $w_{S,C} > 0$, such that $\sum_{C \in \text{Ch}(S)} w_{S,C} = 1$. A sum node semantically defines a mixture model over its children, essentially acting as a latent variable over the component distributions (POON and P. DOMINGOS, 2011; PEHARZ, GENS, PERNKOPF, *et al.*, 2016). Its value is the weighted sum of its children $p_S(\mathbf{X} = \mathbf{x}) = \sum_{C \in \text{Ch}(S)} w_{S,C} \cdot p_C(\mathbf{X} = \mathbf{x})$. To simplify notation, we shall use $N(\mathbf{X} = \mathbf{x})$ as an alias for $p_N(\mathbf{X} = \mathbf{x})$, that is, the probability function given by N 's induced distribution. We often omit the variable assignment $\mathbf{X} = \mathbf{x}$ to \mathbf{x} if the notation is unambiguous.

Example 2.1.2: Gaussian mixture models as probabilistic circuits

A Gaussian Mixture Model (GMM) defines a mixture over Gaussian components. Say we wish to compute the probability of $X = x$ for a GMM \mathcal{G} with three components $\mathcal{N}_1(\mu_1 = 1, \sigma_1^2 = 0.65)$, $\mathcal{N}_2(\mu_2 = 2.5, \sigma_2^2 = 0.85)$ and $\mathcal{N}_3(\mu_3 = 4, \sigma_3^2 = 0.6)$, and suppose we have weights set to $\phi = (0.4, 0.25, 0.35)$. Computing the probability of G amounts to the weighted summation

$$\mathcal{G}(X = x) = 0.4 \cdot \mathcal{N}_1(x; \mu_1, \sigma_1^2) + 0.25 \cdot \mathcal{N}_2(x; \mu_2, \sigma_2^2) + 0.35 \cdot \mathcal{N}_3(x; \mu_3, \sigma_3^2),$$

which is equivalent to a computational graph (i.e. a PC) with a sum node whose weights are set to ϕ and children are the components of the mixture. The figure on the right shows \mathcal{G} (top) and its corresponding PC (middle). Given $x = 1.5$ (in blue), input nodes are computed following the inference flow (bottom, gray edges) up to the root sum node (in red), where a weighted summation is computed to output the probability (in green).



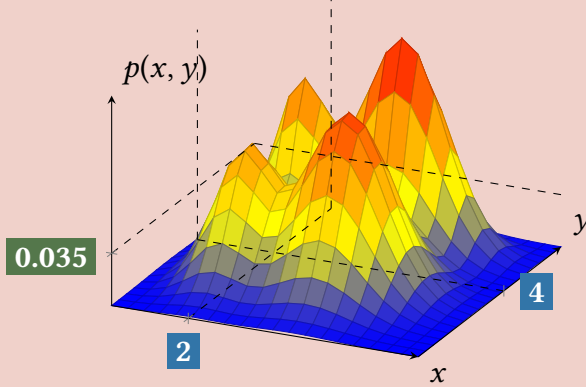
So far, the only nonlinearities present in PCs come from the internal computations of input nodes. In fact, a PC that only contains sums inputs can always be reduced to a sum node rooted PC with a single layer, i.e. a mixture model (see Theorem A.1.1). Adding *product nodes* as another form of nonlinearity increases expressivity sufficiently for PCs

to be capable of representing any (discrete) probability distribution (DARWICHE, 2003; MARTENS and MEDABALIMI, 2014; PEHARZ, TSCHIATSCHKEK, *et al.*, 2015). More importantly, products semantically act as factorizations of their children, indicating an independence relationship between variables from different children. In practice, product nodes are simply products of their childrens' distribution: if P is a PC product node, then its value is given by $P(X = \mathbf{x}) = \prod_{C \in \text{Ch}(P)} C(X = \mathbf{x})$.

Example 2.1.3: Factors as probabilistic circuits

Say we have two GMMs \mathcal{G}_1 and \mathcal{G}_2 . The first is a mixture model over variable X , with component weights $\phi_1 = (0.3, 0.7)$ and gaussians $\mathcal{N}_1(\mu_1 = 2, \sigma_1 = 0.5)$ and $\mathcal{N}_2(\mu_2 = 4, \sigma_2 = 0.8)$. The second is composed of $\mathcal{N}_3(\mu_3 = 3, \sigma_2 = 0.7)$ and $\mathcal{N}_4(\mu_4 = 5, \sigma_2 = 0.4)$, both distributions over variable Y and with weights $\phi_2 = (0.6, 0.4)$.

Suppose $X \perp\!\!\!\perp Y$, yet we wish to compute the joint probability of both x and y . If $X \perp\!\!\!\perp Y$, then $p(x, y) = p(x)p(y) = \mathcal{G}_1(x)\mathcal{G}_2(y)$, which corresponds to a factoring of mixtures. This is represented as a product node (in green) over the two mixture models (in red and purple). The resulting joint of this circuit is shown below.



Now that we have introduced the three most important computational units in PCs, we are finally ready to formally define probabilistic circuits.

Definition 2.1.1 (Probabilistic circuit). A probabilistic circuit C is a rooted connected DAG whose nodes compute any tractable operation of their children, usually either convex combinations, known as sum nodes, or products. Nodes with no outgoing edges, i.e. input nodes, are tractable nonnegative functions whose integrals exist and equal to one. Computing a value from C amounts to a bottom-up feedforward pass from input nodes to root.

While we assume that *tractable* operation or function is acceptable, we are usually interested in $\mathcal{O}(1)$ time computable operations, and often assume the same of input functions to simplify analysis. Further, in this dissertation we are only interested in convex combinations and products, and as such only these operations are considered. When a

probabilistic circuit C contains no consecutive sums or products (i.e. for every sum all of its children are either inputs or products and respectively for products) then it is said to be a *standard* form circuit. Any PC can be transformed into a *standardized* circuit, a process we call *standardization* (see [Theorem A.1.2](#)).

Remark 2.1.1: On operators and tractability

Throughout this work we consider only products and convex combinations (apart from the implicit operations contained within input nodes) as potential computational units. The question of whether any other operator could be used to gain expressivity without loss of tractability is without a doubt an interesting research question, and one that is actively being pursued. However, this is certainly out of the scope of this dissertation, and so we restrict discussion on this topic and only give a brief comment on operator tractability here, pointing to existing literature in this area of research.

[A. FRIESEN and P. DOMINGOS \(2016\)](#) formalize the notion of replacing sums and products in PCs with any pair of operators in a commutative semiring, giving results on the conditions for marginalization to be tractable. They provide examples of common semirings and to which known formalisms they correspond to. One such example are PCs under the Boolean semiring $(\{0, 1\}, \vee, \wedge, 0, 1)$ for logical inference, which are equivalent to Negation Normal Form (NNF, [BARWISE, 1982](#)) and constitute an instance of Logic Circuits (LCs), of which Sentential Decision Diagrams (SDDs, [DARWICHE, 2011](#)) and Binary Decision Diagrams (BDDs, [AKERS, 1978](#)) are a part of. Another less common semiring in PCs is the real min-sum semiring $(\mathbb{R}_\infty, \min, +, \infty, 0)$ for nonconvex optimization ([Abram L. FRIESEN and P. DOMINGOS, 2015](#)).

Recently, [VERGARI, Y. CHOI, et al. \(2021\)](#) extensively covered tractability conditions and complexity bounds for convex combinations, products, exp (and more generally powers in both naturals and reals), quotients and logarithms, even giving results for complex information-theoretic queries, such as entropies and divergences. Notably, they analyze whether structural constraints (and thus, in a sense, tractability) under these conditions are preserved.

Up to now, we have only considered summations as nonnegative weighted sums. Indeed, in most literature the sum node is defined as a convex combination. However, negative weights have appeared in Logistic Circuits ([LIANG and VAN DEN BROECK, 2019](#)) for discriminative modeling; and in Probabilistic Generating Circuits ([ZHANG et al., 2021](#)), a class of tractable probabilistic models that subsume PCs. [Denis D. MAUÁ et al. \(2017\)](#) extend (nonnegative) weights in sum nodes with probability intervals, effectively inducing a credal set ([Fabio G. COZMAN, 2000](#)) for measuring imprecision.

Other works include PCs with quotients ([SHARIR and SHASHUA, 2018](#)), transformations ([PEVNÝ et al., 2020](#)), max ([MELIBARI, POUPART, and DOSHI, 2016](#)), and einsum ([PEHARZ, LANG, et al., 2020](#)) operations.

Before we address the key components that make PCs interesting tractable probabilistic models, we must first discuss some important concepts that often come up in PC literature.

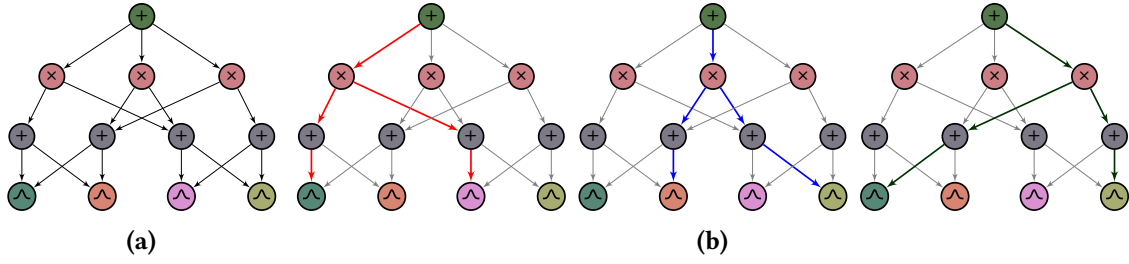


Figure 2.1: A probabilistic circuit (a) and 3 of the 12 possible induced subcircuits (b).

Mainly, we are interested in defining two notions here: the scope of a unit and induced subcircuits.

In simple terms, the scope of a computational unit N of a PC is merely the set of all variables that appear in the descendants of N . More formally, denote by $\text{Sc}(N)$ the set of all variables that appear in N . We inductively compute the scope of circuit by a bottom-up approach: the scope of an input node L_p is the set of variables that appear in p 's distribution¹, and the scope of any other node is the union of all of its children's scopes. The notion of scope is essential to the structural constraints seen in [Section 2.2](#).

As an example, take the circuit from [Example 2.1.3](#). The scope of input nodes $\text{Sc}(\text{green})$ and $\text{Sc}(\text{red})$ are $\text{Sc}(\text{green}) = \text{Sc}(\text{red}) = \{X\}$, while $\text{Sc}(\text{blue}) = \text{Sc}(\text{green}) = \{Y\}$. Consequentially, their parent sum nodes will have the same scope as their children $\text{Sc}(\text{blue}) = \{X\}$ and $\text{Sc}(\text{blue}) = \{Y\}$, yet the root node's scope is $\text{Sc}(\text{green}) = \{X, Y\}$, since its children's scopes are distinct. The size of a probabilistic circuit is the number of nodes and edges of its computational graph. We use $|C|$ to denote the size of a probabilistic circuit C .

Let C a probabilistic circuit and node $N \in C$. We say that S_N is a subcircuit of C rooted at N if S_N 's root is N , all nodes and edges in S_N are also in C and S is also a probabilistic circuit. We now introduce the concept of induced subcircuits ([CHAN and DARWICHE, 2006](#); [DENNIS and VENTURA, 2015](#); [PEHARZ, GENS, and P. DOMINGOS, 2014](#)).

Definition 2.1.2 (Induced subcircuit). *Let C a probabilistic circuit. An induced subcircuit S of C is a subcircuit of C rooted at C 's root such that all edges coming out of product nodes in C are also in S , and of all edges coming out of sum nodes in C , only one is in S .*

Examples of induced subcircuits are visualized in [Figure 2.1](#). When the induced subcircuit is a tree, as is the case in [Figure 2.1](#), they are referred to as induced tree ([H. ZHAO, MELIBARI, et al., 2015](#); [H. ZHAO, ADEL, et al., 2016](#)).

So far, by [Definition 2.1.1](#) a PC does not yet necessarily represent a probability distribution with tractable (marginal) inference. In the next section, we formally define sufficient conditions for tractability in the form of structural constraints that we have only previously mentioned in passing.

¹ Although we previously defined input nodes as mere functions, here we are explicitly associating a random variable to a *probability* function. Indeed, if we are being rigorous, we should define input nodes as a pair of random variable and function. To save space we instead assume, as previously stated, that the function is seen as both the probability (density) function as well as the distribution itself, and thus its scope is the scope of its distribution, i.e. the random variables that come into play in a probability distribution.

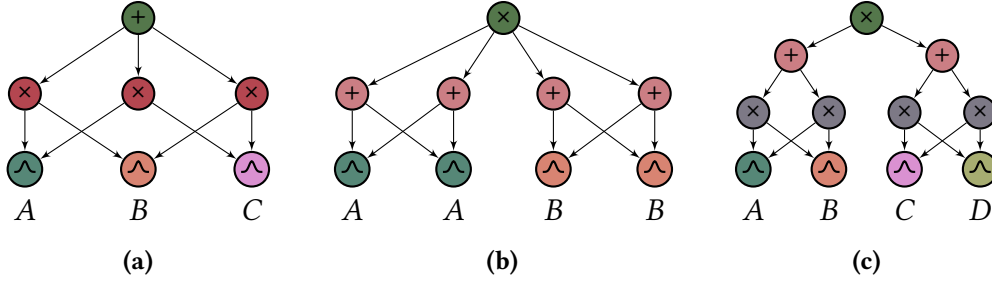


Figure 2.2: *Decomposable but non-smooth (a), smooth but non-decomposable (b), and smooth and decomposable (c) circuits.*

2.2 Deciding What to Constraint

In this section we are interested in studying how the structural constraints in PCs enable different inference tasks. We shall first cover the more basic queries, namely *probability of evidence* (EVI), *marginal probability* (MAR), *conditional probability* (CON) and *maximum a posteriori probability* (MAP). After that we briefly address more complex queries such as mutual information, entropies and *expectation* (EXP).

2.2.1 Basic Queries

The most basic inference task we are interested in computing is the probability of evidence. To unlock this, we must introduce two structural constraints known as *smoothness* and *decomposability*.

Definition 2.2.1 (Smoothness). *A probabilistic circuit C is said to be smooth if for every sum node S in C , $\text{Sc}(C_1) = \text{Sc}(C_2)$ for $C_1, C_2 \in \text{Ch}(S)$.*

Definition 2.2.2 (Decomposability). *A probabilistic circuit C is said to be decomposable if for every product node P in C , $\text{Sc}(C_1) \cap \text{Sc}(C_2) = \emptyset$ for $C_1, C_2 \in \text{Ch}(P)$.*

When a PC is *decomposable*, all of its induced subcircuits are induced trees. For any *smooth* and *decomposable* PC, computing EVI is done in linear time in the number of edges. In fact this is true for MAR and CON as well. To compute marginals, it is sufficient to compute the corresponding marginals with respect to each input node and proceed to propagate values bottom-up. For conditionals, we simply compute two passes: one where we marginalize the conditional variables and the other any other variables that are not present in our query. These procedures are formalized in the theorem below.

Theorem 2.2.1 (POON and P. DOMINGOS, 2011; Y. CHOI, VERGARI, and VAN DEN BROECK, 2020; VERGARI, Y. CHOI, *et al.*, 2021). *Let C a smooth and decomposable PC. Any one of EVI, MAR or CON can be computed in linear time (in the number of edges of C).*

Importantly, EVI and CON are both special cases of MAR in the sense that they are marginalizations over different intervals (see page 51). More specifically, EVI is a marginalization over an empty set and CON is simply the quotient between two marginalization passes. Algorithmically, this means that the only distinction between these three queries is on what to do on the input nodes. We say that a variable assignment \mathbf{x} for circuit C is *complete* if for every variable $X \in \text{Sc}(C)$, \mathbf{x} assigns a value to X ; otherwise it is said to be a

Algorithm 1 EVI**Input** A probabilistic circuit C and complete assignment \mathbf{x} **Output** Probability $C(\mathbf{x})$

- 1: Let v a hash function mapping a node to its probability
- 2: **for** each N in reverse topological order **do**
- 3: **if** N is an input **then** $v_N \leftarrow N(\mathbf{x})$
- 4: **else if** N is a sum **then** $v_N \leftarrow \sum_{C \in \text{Ch}(N)} w_{N,C} v_C$
- 5: **else if** N is a product **then** $v_N \leftarrow \prod_{C \in \text{Ch}(N)} v_C$
- 6: **return** v_R , where R is C 's root

Algorithm 2 MAR**Input** A probabilistic circuit C and partial assignment \mathbf{x} **Output** Probability $\int C(\mathbf{x}, \mathbf{y}) d\mathbf{y}$

- 1: Let v a hash function mapping a node to its probability
- 2: **for** each N in reverse topological order **do**
- 3: **if** N is an input **then** $v_N \leftarrow \int N(\mathbf{x}, \mathbf{y}) d\mathbf{y}$
- 4: **else if** N is a sum **then** $v_N \leftarrow \sum_{C \in \text{Ch}(N)} w_{N,C} v_C$
- 5: **else if** N is a product **then** $v_N \leftarrow \prod_{C \in \text{Ch}(N)} v_C$
- 6: **return** v_R , where R is C 's root

partial assignment. [Algorithm 1](#) and [Algorithm 2](#) show the exact algorithmic procedures to extract EVI and MAR queries from C . For CON, it suffices to run two MARs.

Because EVI, MAR and CON are the most basic forms of querying in PCs, we shall refer to these as *base queries*. Although smoothness and decomposability are sufficient conditions for tractable base queries, they are not necessary conditions. As a matter of fact, decomposability can be replaced with a third weaker constraint known as *consistency*. Denote by $\text{Desc}(N)$ the set of all descendants of N .

Definition 2.2.3 (Consistency). *A probabilistic circuit C is said to be consistent if for any product node P in C , it holds that for every two children $C_1, C_2 \in \text{Ch}(P)$ there does not exist two input nodes $L_p^1 \in \text{Desc}(C_1)$ and $L_q^2 \in \text{Desc}(C_2)$ whose scope is the same and $p(\mathbf{x}) \neq q(\mathbf{x})$ for any $\mathbf{x} \in \mathcal{X}$.*

In [PEHARZ, TSCHIATSCHEK, et al. \(2015\)](#), the authors show that although smooth and consistent PCs are more succinct ([DARWICHE and MARQUIS, 2002](#)) compared to smooth and decomposable circuits, the gain is only mild, further proving that any smooth and consistent PC can be polynomially translated to a decomposable equivalent. In practice, because constructing decomposable circuits (and verifying decomposability) is much easier compared to doing the same with consistency, we instead focus on smooth and decomposable PCs.

Suppose we wish to compute the most probable assignment of a variable, say for classification or image reconstruction. To do so, we must compute the conditional query

$$\max_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}) = \max_{\mathbf{y}} \frac{p(\mathbf{y}, \mathbf{x})}{p(\mathbf{x})} = \frac{\max_{\mathbf{y}} p(\mathbf{y}, \mathbf{x})}{p(\mathbf{x})}, \quad (2.1)$$

Algorithm 3 MAP**Input** A probabilistic circuit C and evidence assignment \mathbf{x} **Output** Probability $\max_{\mathbf{y}} C(\mathbf{y}|\mathbf{x})$

- 1: Let v a hash function mapping a node to its probability
- 2: **for** each N in reverse topological order **do**
- 3: **if** N is an input **then** $v_N \leftarrow \max_{\mathbf{y}} N(\mathbf{y}, \mathbf{x})$
- 4: **else if** N is a sum **then** $v_N \leftarrow \max_{C \in \text{Ch}(N)} w_{N,C} v_C$
- 5: **else if** N is a product **then** $v_N \leftarrow \prod_{C \in \text{Ch}(N)} v_C$
- 6: **return** $v_R/C(\mathbf{x})$, where R is C 's root

Algorithm 4 ARGMAP**Input** A probabilistic circuit C and evidence assignment \mathbf{x} **Output** The most probable state $\arg \max_{\mathbf{y}} C(\mathbf{y}|\mathbf{x})$

- 1: Compute $\max_{\mathbf{y}} C(\mathbf{y}|\mathbf{x})$ and store values in v
- 2: $\mathbf{N} \leftarrow \text{MAXINDUCEDTREE}(C, v)$
- 3: Call \mathbf{y} the set of initially empty arg max states
- 4: **for** each input node $L \in \mathbf{N}$ **do**
- 5: $\mathbf{y} \leftarrow \mathbf{y} \cup \arg \max_{\mathbf{z}} N(\mathbf{z}, \mathbf{x})$
- 6: **return** \mathbf{y}

often called the *maximum a posteriori* (MAP) probability. For this dissertation, we shall only consider the case of full MAP, as computing partial MAP is provably hard in most PCs (PEHARZ, GENS, PERNKOPF, *et al.*, 2016; DE CAMPOS, 2011). Unless explicitly stated, MAP shall mean full MAP, i.e. when $\mathbf{x} \cup \mathbf{y}$ forms a complete assignment of the scope. Full MAP also goes by the name of *most probable explanation* (MPE) in literature. Although at first it seems like MAP is no harder than computing a CON, it turns out that for smooth and decomposable PCs, MAP is unfortunately intractable (CONATY *et al.*, 2017; MEI *et al.*, 2018). To unlock access to MAP we must make the circuit *deterministic*.

Definition 2.2.4 (Determinism). *A probabilistic circuit C is said to be deterministic if for every sum node $S \in C$ only one child of S has nonnegative value at a time for any complete assignment.*

At this point, we must introduce a graphical notation for *indicator* nodes. An indicator node is an input node whose function is the characteristic function $f(x) = \llbracket x = k \rrbracket$, i.e. a degenerate function with all of its mass on k and zero anywhere else. A special case is when X is binary and $k = 1$, in which case we say the input node is a literal node, denoting by the usual propositional notation X for when $k = 1$ and $\neg X$ for $k = 0$. Graphically, we shall use \odot for indicators and the textual propositional notation for literals.

With determinism, we now have access to MAP.

Theorem 2.2.2 (PEHARZ, GENS, PERNKOPF, *et al.*, 2016). *Let C a smooth, decomposable and deterministic PC. MAP is computable in linear time (in the number of edges of C).*

When a PC is smooth, decomposable and deterministic, the MAP is easily computable by simply replacing sum nodes with a max operation and performing a feedforward

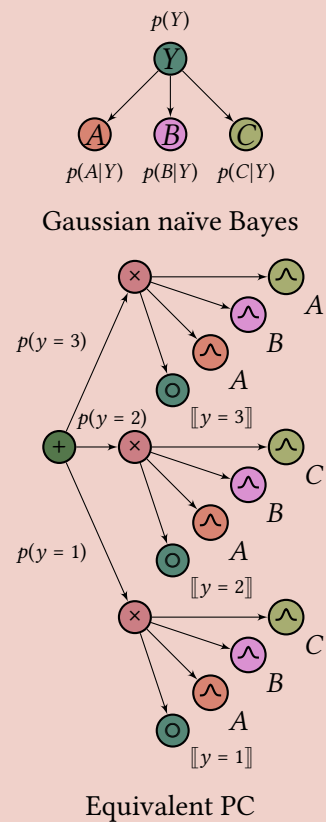
EVI pass. This is commonly called the Max-Product algorithm, shown more formally in [Algorithm 3](#). To find the states that maximize [Equation \(2.1\)](#) in a given circuit C , we first compute the MAP probabilities through the usual bottom-up pass, and then find the maximum (in terms of probability) induced tree \mathcal{M} rooted at C . This maximum induced tree can be retrieved by a top-down pass selecting the most probable sum child nodes according to the probabilities set by MAP. Since C is decomposable, there cannot exist a node in \mathcal{M} with more than one parent, meaning it is by construction a tree whose leaves are input nodes with scopes whose union is the scope of C . This reduces the problem to a divide-and-conquer approach where each input node is individually maximized (see [Algorithm 4](#) and [page 52](#) for a formal proof on its validity).

Example 2.2.1: Naïve Bayes as probabilistic circuits

Suppose we have samples of per capita census measurements on three different features, say age A , body mass index B and average amount of cheese consumed daily C from three different cities Y . Assuming A , B and C are independent, given a sample $x = (a, b, c)$ we can use Gaussian naïve Bayes to predict x 's class

$$p(y|a, b, c) = p(y)p(a|y)p(b|y)p(c|y).$$

In PC terms, $p(y)$ are the prior probabilities, i.e. sum weights, for each class and $p(z|y)$ are Gaussian input nodes corresponding to the distributions of each feature in each city. To make sure that these are in fact conditional distributions, we introduce indicator variables “selecting” Y 's state. Since the PC resulting is deterministic, we can compute the MAP for classification in linear time by simply replacing the root node with a max, which is exactly equivalent to finding the highest likelihood of x for each city y .



2.2.2 Complex Queries

Although base queries cover most of the needs of the usual data scientist, more complex tasks involving information-theoretic measures, logical queries or distributional divergences are also (tractably) computable under the right conditions. Particularly, we are interested in a key component for tractability in all these tasks: the notion of *vtrees* and *structure decomposability*, a stronger variant of decomposability where variable partitionings on product nodes follow a hierarchy. This hierarchy is easily visualized through a *vtree* (variable tree), a data structure that defines a (partial) ordering of variables.

Definition 2.2.5 (Vtree). A variable tree \mathcal{V} , or vtree, over a set of variables X is a binary

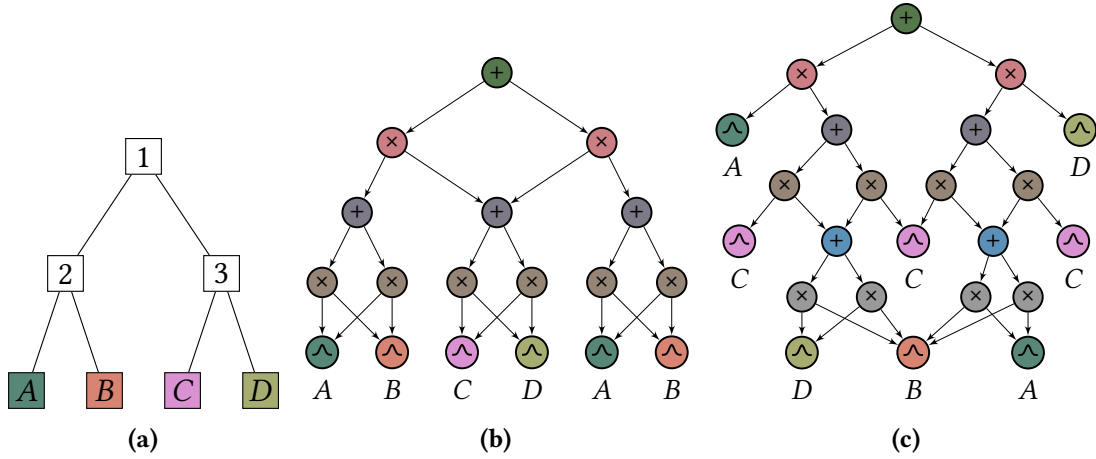


Figure 2.3: A vtree (a) defining an order (A, B, C, D) , a 2-standard structure decomposable probabilistic circuit that respects the vtree (b), and a 2-standard decomposable probabilistic circuit that does not (c).

tree whose leaf nodes have a one-to-one and onto mapping $\phi_{\mathcal{V}}$ with X .

We shall adopt the same scope definition and notation $\text{Sc}(\cdot)$ for vtrees as in PCs. Let v a vtree node from a vtree \mathcal{V} . If v is a leaf node, its scope is $\phi_{\mathcal{V}}(v)$, i.e. the leaf's assigned variable; otherwise its scope is the union of the scope of its children. For an inner node v , we shall call its left child v^{\leftarrow} and right child v^{\rightarrow} . Every inner node v of a vtree \mathcal{V} defines a variable *partitioning* of the scope $(\text{Sc}(v^{\leftarrow}), \text{Sc}(v^{\rightarrow}))$, while the leaves of \mathcal{V} define a partial ordering of $\text{Sc}(\mathcal{V})$. We are especially interested in the scope partitioning aspect of vtrees.

Definition 2.2.6. A product node P respects a vtree node v if P contains only two children $\text{Ch}(P) = \{C_1, C_2\}$, and $\text{Sc}(C_1) = \text{Sc}(v^{\leftarrow})$ and $\text{Sc}(C_2) = \text{Sc}(v^{\rightarrow})$.

Obviously, the above definition is vague with regards to which child (i.e. graphically, left or right) of P should respect the scope of which v child. We therefore assume a fixed order for P 's children and say that the (graphically) left child is called the *prime* and (graphically) right child the *sub*, and refer to P as an *element*. This ultimately means that the scope of the prime (resp. sub) of P must be the same as the scope of the left (resp. right) child of v . Although the graphical concept of left and right is needed for easily visualizing the scope partitioning of a product with respect to a vtree node, we do not use it strictly. In fact, when the situation is unambiguous, we compactly represent the computational graph without adhering to the left-right convention in favor of readability.

We say that a vtree is linear, if either it is left-linear or right-linear. A left- (resp. right) linear vtree is a vtree whose inner nodes all have leaf nodes on their right (resp. left) child. Similarly, a vtree is said to be left- (resp. right) leaning if the number of leaf nodes as right (resp. left) children is much higher than left (resp. right) children. Otherwise, it is a balanced vtree. The variable order of a vtree is the sequence of leaf nodes (i.e. variables) read from left to right. Figure 2.3a shows a balanced vtree with order (A, B, C, D) .

Now that we understand what a vtree is, we can properly introduce *structure decomposability*, a stronger variant of decomposability. We say that a PC is 2-standard if it is standard and all of its product nodes have exactly two children.

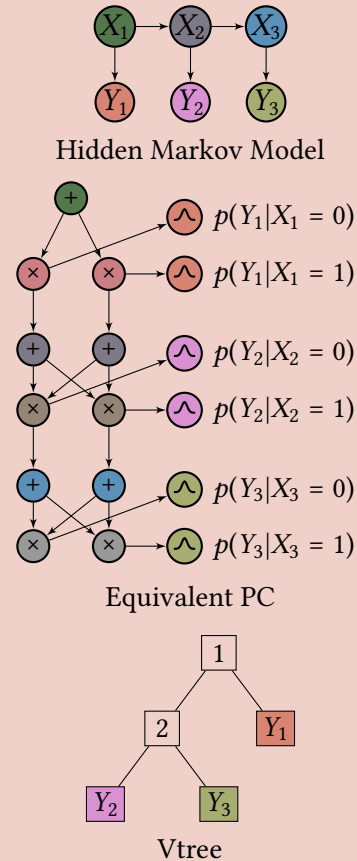
Definition 2.2.7 (Structure decomposability). Let C a 2-standard probabilistic circuit and \mathcal{V} a vtree with same scope as C . C is said to be structure decomposable if every i -th product layer of C respects every i -th inner node layer of \mathcal{V} .

Example 2.2.2: Hidden Markov models as probabilistic circuits

Say we wish to model a temporal structured dependence between three latent binary variables, for example the presence of a subject X_1 , verb X_2 and object X_3 in a natural language phrase. Each observation Y_i is a fragment (of X_i) taken from a complete sentence $y = (y_1, y_2, y_3)$. The first-order Hidden Markov Model (HMM) (on the right) models the joint probability of sentences

$$p(X_{1..3}, Y_{1..3}) = p(X_1) \prod_{i=2}^3 p(X_i | X_{i-1}) \prod_{i=1}^3 p(Y_i | X_i).$$

This is computationally equivalent to the PC on the right. Each input node $p(Y_i | X_i)$ is a conditional distribution model (possibly another PC) for each assignment (here two) of X_i , meaning that if $p(Y_i | X_i = 0) > 0$, then $p(Y_i | X_i = 1) = 0$ and vice-versa. Root weights are exactly $p(X_1)$, and each $p(X_i | X_{i-1})$ translates into the other matched color sum weights. Further, every product (except for the redundant \otimes nodes) follows the partitionings imposed by the vtree. This means that this PC is not only smooth, but structure decomposable and deterministic.



Although we assume a 2-standard PC in Definition 2.2.7, its assumption was only for convenience, and does not imply in a loss of expressivity. As a matter of fact, any PC can be 2-standardized, i.e. to standardize the circuit such that every product node only has two children (see Theorem A.1.3). Intuitively, structure decomposability merely states that every two product nodes whose scopes are the same must partition their scopes (between their two children) exactly the same (and according to their corresponding vtree node). Semantically speaking, a vtree's inner node v defines an independence relationship between $\text{Sc}(v^{\leftarrow})$ and $\text{Sc}(v^{\rightarrow})$ under the distribution encoded by its PC.

Figure 2.3 shows a vtree \mathcal{V} and two probabilistic circuits, say C_1 for the one in the middle and C_2 for the one on the right. Notice how C_1 respects \mathcal{V} , as each \otimes respects the split at vtree node 1 (namely $\{A, B\}$). The primes are then \oplus whose scopes are $\{A, B\}$, while the sub is the one with two parents and scope $\{C, D\}$. For each of these, their children \otimes also respect \mathcal{V} : they either encode the same split as 2 or as 3, depending on whether they are descendants from the sub or prime of 1. Although C_2 is decomposable, it does *not* respect \mathcal{V} , as \otimes encode different variable partitionings ($(\{A\}, \{B, C, D\})$ and $(\{A, B, C\}, \{D\})$). In fact, it is not structure decomposable, as it does not respect any vtree.

Despite our structure decomposability definition relying on a vtree, there is at least one alternative definition that defines it in terms of *circuit compatibility*. Essentially, a circuit C_1 is *compatible* with C_2 if they can be 2-standardized and rearranged such that any two products with same scope, one from C_1 and the other C_2 , partition the scope the same way (VERGARI, Y. CHOI, *et al.*, 2021). A structure decomposable PC is then defined as a PC that is compatible with a copy of itself. In summary, the two definitions of structure decomposability are equivalent, except compatibility implicitly assumes an arrangement of product scopes that is analogous to a vtree.

The notion of a vtree (or compatibility for that matter) is key to more complex queries. For instance, given two probabilistic circuits C_1 and C_2 , computing cross entropy between the two is $\mathcal{O}(|C_1||C_2|)$ as long as both have the same vtree and at least the circuit that needs to be log-computable is also deterministic. Likewise, computing the Kullback-Leibler (KL) divergence between C_1 and C_2 requires that the two share the same vtree and both be deterministic. Mutual Information (MI), in turn, calls for the circuit to be smooth, structure decomposable and an even stronger version of determinism known as *marginal determinism* where sums can only have one nonnegative valued child for any *partial* assignment at a time. In fact, when a PC is smooth, decomposable and marginal deterministic, marginal MAP, i.e. MAP over partial assignments becomes linear time computable. For a more detailed insight on the tractability of these (and other) queries, as well as proofs on these results, we point to the comprehensive study of VERGARI, Y. CHOI, *et al.* (2021).

A particularly interesting class of queries that becomes tractable when circuits are structure decomposable is expectation (EXP). One notable example from this class is computing the probability of logical events. This leads us to logic circuits, a parallel version of probabilistic circuits for logical reasoning.

2.3 Probabilistic Circuits as Knowledge Bases

We superficially mentioned in Remark 2.1.1 that PCs under a Boolean semiring with conjunctions and disjunctions as operators are known as Logic Circuits (LCs). In this section, we formally yet briefly define LCs and more precisely show the connection between PCs and LCs.

2.3.1 From Certainty...

Logic circuits are computational graphs just like PCs, but whose input are always Booleans (and as such the scope is over propositional variables) and computational units define either a conjunction, disjunction or literal of their inputs. While the computational graph in PCs encodes uncertainty as a probability distribution, in LCs their computational graph encodes certain knowledge as a propositional language. Similar to PCs, computing the satisfiability of an assignment is done by a bottom-up feedforward evaluation of the circuit. In terms of notation, we shall use \bigvee for disjunction nodes, \bigwedge for conjunction nodes, and X and $\neg X$ for literal nodes.

Definition 2.3.1 (Logic circuit). *A logic circuit \mathcal{L} is a rooted connected DAG whose nodes compute either a conjunction or a disjunction of their children. Nodes with no outgoing*

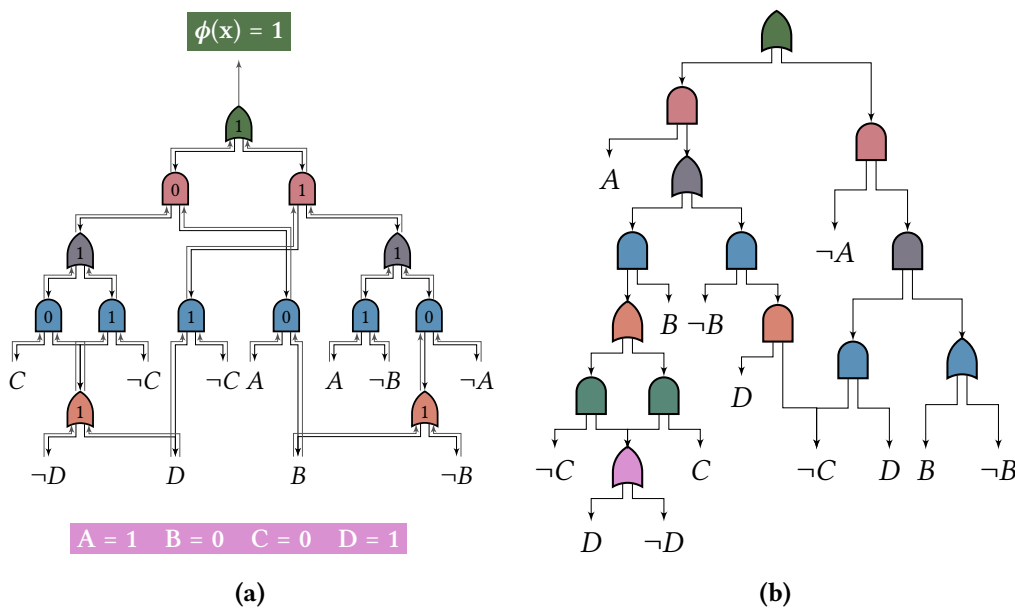


Figure 2.4: Two smooth, structure decomposable and deterministic logic circuits encoding the same logic constraint $\phi \equiv (A \wedge B) \vee (\neg C \wedge D)$ for a balanced (a) and a right-linear (b) vtree. In (a), a circuit evaluation for an assignment, with each node value in the bottom-up evaluation pass shown inside nodes.

edges, i.e. literal nodes, are indicator functions of either a positive (true) or negative (false) assignment. Computing the satisfiability of \mathcal{L} amounts to a bottom-up feedforward pass from literal nodes to root.

Evidently, logic circuits are closely related to probabilistic circuits. The strikingly similar definitions are not coincidence: much of the literature on PCs have their origins on LCs. In fact, most structural constraints in PCs are the exact same (up to even their names²) as their LC analogues. In this dissertation, we are particularly interested in the specific subset of smooth, structure decomposable and deterministic LCs, known as Sentential Decision Diagrams (SDDs, DARWICHE, 2011). Figure 2.4 shows two SDDs encoding the same knowledge base but under different vtrees. The one on the left respects a balanced vtree and the one on the right a right-linear vtree.

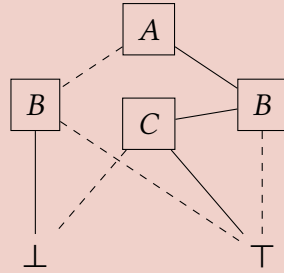
Logic circuits have appeared in literature under the name of Negated Normal Form (NNF) (DARWICHE, 2001; DARWICHE, 1999), a superset of other propositional compilation languages such as Binary Decision Diagrams (BDDs, BRYANT, 1986), Propositional DAGs (PDAGs, WACHTER and HAENNI, 2006), Sentential Decision Diagrams (SDDs, DARWICHE, 2011), DNFs and CNFs. Although no structural constraint is required for correctly evaluating satisfiability in logic circuits, the same properties defined in past sections have come up in LCs to enable other more complex queries like equivalence, implicates, sentential entailment and model counting, as well as transformations such as closed conditionings, forgetting, conjunctions and disjunctions (DARWICHE and MARQUIS, 2002). The succinctness (i.e. expressive efficiency) of LCs are also impacted by these structural restrictions

² There are some exceptions. Smoothness is sometimes referred to as *completeness* in PCs, while determinism has the alternative name of *selectivity*.

(GOGIC *et al.*, 1995; PAPADIMITRIOU, 1994; DARWICHE and MARQUIS, 2002). See DARWICHE and MARQUIS (2002) and DARWICHE (2020) for more on logic circuits.

Example 2.3.1: BDDs as logic circuits

A Binary Decision Diagram (BDD, BRYANT, 1986) defines a Boolean function over binary variables as a rooted DAG. BDDs are a subset of smooth, structure decomposable and deterministic LCs (i.e. SDDs) whose vtrees are always right-linear.



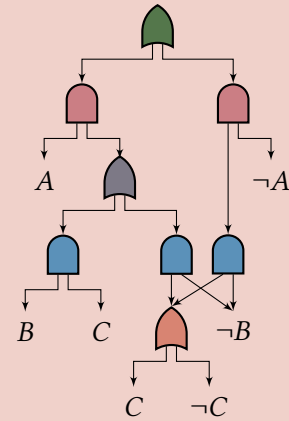
In their usual notation, inner nodes are variables and leaves are constants \perp and \top . Evaluating an assignment $\mathbf{x} \subseteq \{0, 1\}^n$ is equivalent to a path from the root to a leaf where each variable X determines a decision to go through the dashed line when $x = 0$ or solid line when $x = 1$. If the path ends at a \top , the function returns 1, otherwise it must end at a \perp and therefore returns 0. The BDD above encodes the following logic formula

$$\phi(A, B, C) = (A \vee \neg B) \wedge (\neg B \vee C),$$

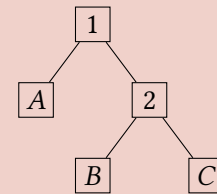
also shown as a truth table on the right, together with a logic circuit that encodes the same truth table and its vtree.

A	B	C	$\phi(\mathbf{x})$
0	0	0	1
1	0	0	1
0	1	0	0
1	1	0	0
0	0	1	1
1	0	1	1
0	1	1	0
1	1	1	1

Truth Table



Equivalent LC



Vtree

2.3.2 ...to Uncertainty

Logic circuits are easily extensible to probabilistic circuits. In fact, if we think of an LC as the support of a PC the connections between the two come naturally. Suppose a 2-standard smooth, structure decomposable and deterministic probabilistic circuit C over binary variables. We can construct an identically structured logic circuit (up to input nodes) \mathcal{L} with same vtree as C whose underlying Boolean function encodes $\phi(\mathbf{x}) = \llbracket C(\mathbf{x}) > 0 \rrbracket$. Since sums act exactly like disjunctions and products like conjunctions under the Boolean semiring, products in \mathcal{L} are replaced with conjunctions in \mathcal{L} , and sums with disjunctions. Input nodes from C are replaced with a literal node if the function is degenerate, or with a disjunction over positive and negative literals otherwise. This makes sure \mathcal{L} acts as the

support of \mathcal{C} , as each disjunction node $S_{\mathcal{L}}$ of \mathcal{L} defines

$$S_{\mathcal{L}}(\mathbf{x}) = \bigvee_{C \in \text{Ch}_C(S_{\mathcal{L}})} \llbracket C_p(\mathbf{x}) > 0 \rrbracket \wedge \llbracket C_s(\mathbf{x}) > 0 \rrbracket,$$

where $\text{Ch}_C(S)$ retrieves the children of S 's corresponding sum node in \mathcal{C} , with $C_p(\mathbf{x})$ and $C_s(\mathbf{x})$ the probabilities of C 's prime and sub respectively. The corresponding sum node S_C in \mathcal{C} then only attributes a weight (i.e. probability) to each positive element as usual

$$S_C(\mathbf{x}) = \sum_{C \in \text{Ch}(S_C)} w_{S_C, C} \cdot C_p(\mathbf{x}) \cdot C_s(\mathbf{x}).$$

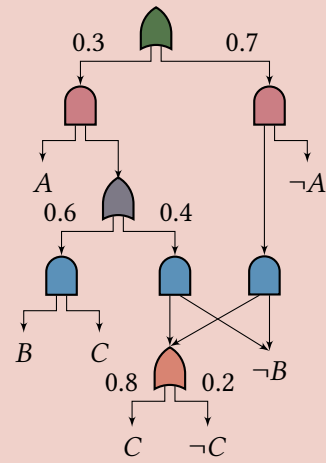
Example 2.3.2: Embedding certain knowledge in probabilistic circuits

Recall the logic circuit \mathcal{L} from [Example 2.3.1](#). Imagine we wish to model the uncertainty coming from all assignments where $\mathcal{L}(\mathbf{x}) = 1$. In other words, we want to assign a positive probability to all true entries in the previous example's truth table, turning it into a probability table. The table on the right shows the chosen probabilities for each instance. Naturally, they all sum to one, with logically impossible assignments set to zero.

Compiling an LC into a PC is straightforward: replace conjunctions with product nodes and disjunctions with sum nodes. Input nodes are left untouched, as literal nodes are just degenerate probability distributions. Sum weights are what ultimately define the probabilities in the probability table. The PC on the right is the result of the compilation of \mathcal{L} into a probabilistic circuit whose distribution is defined by the probability table above it. When we mean to say that a PC has its support defined by its underlying LC, then we use the logic gate notation with the added weights on edges coming out from sum nodes.

A	B	C	$\phi(\mathbf{x})$	$p(\mathbf{x})$
0	0	0	1	0.140
1	0	0	1	0.024
0	1	0	0	0.000
1	1	0	0	0.000
0	0	1	1	0.560
1	0	1	1	0.096
0	1	1	0	0.000
1	1	1	1	0.180

Probability Table



Probabilistic Circuit

This compatibility between logic and probabilistic circuits allows certain knowledge to be embedded into an uncertain model by constructing a computational graph whose underlying logic circuit correctly attributes positive values only to the support of the distribution. When such a computational graph is also smooth, structure decomposable and deterministic, then it belongs to a special subclass of PCs called Probabilistic Sentential Decision Diagrams (PSDDs, [Kisa et al., 2014](#)). An alternative use case for logic circuits within the context of probabilistic reasoning is querying for the probability of logical events, i.e. the expectation of a logic query with respect to a distribution. This kind of

Algorithm 5 EXP

Input A smooth, structure decomposable PC \mathcal{C} and LC \mathcal{L} both with vtree \mathcal{V}

Output The expectation $\mathbb{E}_{\mathcal{C}}[\mathcal{L}] = \int \mathcal{C}(\mathbf{x})\mathcal{L}(\mathbf{x}) d\mathbf{x}$

- 1: Let v a hash function mapping a pair of nodes to an expectation
- 2: **for** each pair of nodes (N_C, N_L) in reverse topological order **do**
- 3: **if** N_C is an input **then** $v_{N_C, N_L} \leftarrow \mathbb{E}_{N_C}[N_L]$
- 4: **else if** N_C is a product **then** $v_{N_C, N_L} \leftarrow \prod_{C \in \text{Ch}(N_C)} v_{C, C_L}$
- 5: **else if** N_C is a sum **then** $v_{N_C, N_L} \leftarrow \sum_{C' \in \text{Ch}(N_C)} \sum_{C'' \in \text{Ch}(N_L)} w_{S, C'} \cdot v_{C', C''}$
- 6: **return** v_R , where R is \mathcal{C} 's root

query is enabled by the following result.

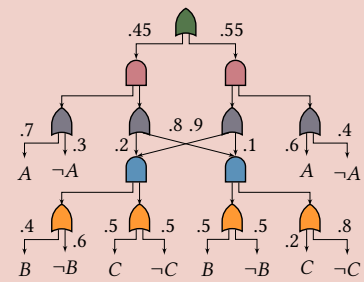
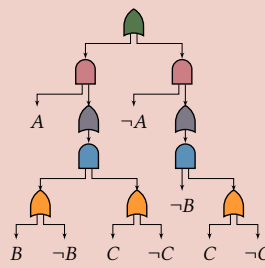
Theorem 2.3.1 (Y. CHOI, VERGARI, and BROECK, 2020). *If \mathcal{C} is a smooth and structure decomposable probabilistic circuit with vtree \mathcal{V} , and \mathcal{L} a structure decomposable logic circuit also respecting \mathcal{V} , then $\mathbb{E}_{\mathcal{C}}[\mathcal{L}]$ is polynomial time computable (in the number of edges).*

Let \mathcal{C} a smooth and structure decomposable circuit with vtree \mathcal{V} , and \mathcal{L} a logic circuit representing a logical query whose vtree is also \mathcal{V} . Computing the probability of \mathcal{L} with respect to the distribution encoded by \mathcal{C} is done by a bottom-up evaluation over both circuits at the same time. Algorithm 5 shows the procedure algorithmically. Importantly, the procedure relies on evaluating the expectation for each node in a reverse topological fashion according to line 2. This paired reverse topological sorting is done by a recursive assignment. For product nodes, primes are paired with primes, and subs with subs; for sums, each combination of PC child with LC child is paired up.

Example 2.3.3: Computing the probability of logical events

Say we have a PC encoding the distribution shown in the table on the right. Suppose we wish to compute the probability of a logical event, say $\phi \equiv A \vee \neg B$. A naïve approach would be to go over each assignment \mathbf{x}

where $\phi(\mathbf{x}) = 1$, and compute their sum. This is obviously exponential on the number of variables. Instead, we may compile ϕ into the logic circuit above and run the EXP algorithm to (in polynomial time) compute this otherwise intractable marginalization. Running the EXP gives us a probability of $0.6415 = 1.0 - (0.1365 + 0.2220)$, which is exactly the desired probability.



A	B	C	$\phi(\mathbf{x})$	$p(\mathbf{x})$
0	0	0	1	.1000
1	0	0	1	.0580
0	1	0	0	.1365
1	1	0	1	.0805
0	0	1	1	.1860
1	0	1	1	.0970
0	1	1	0	.2220
1	1	1	1	.1195

Remark 2.3.1: On applications of probabilistic circuits

So far, we have not yet addressed real-world applications of probabilistic circuits. Although their usage has not yet gained popularity among the data science crowd, they have been successfully employed in a multitude of interdisciplinary tasks. Here, we give a brief survey on the different use cases of probabilistic circuits present in literature.

Computer vision is perhaps the most popular application for deep learning, and this could not be different for probabilistic circuits. PCs have been used for image classification (GENS and P. DOMINGOS, 2012; SGUERRA and F. G. COZMAN, 2016; LLERENA and DERATANI MAUÁ, 2017; R. GEH and D. MAUÁ, 2019; PEHARZ, VERGARI, *et al.*, 2020), image reconstruction and sampling (POON and P. DOMINGOS, 2011; DENNIS and VENTURA, 2017; PEHARZ, LANG, *et al.*, 2020; Cory J BUTZ *et al.*, 2019), image segmentation and scene understanding (Abram L FRIESEN and P. DOMINGOS, 2017; YUAN *et al.*, 2016; Abram L FRIESEN and P. M. DOMINGOS, 2018; RATHKE *et al.*, 2017), and activity recognition (J. WANG and G. WANG, 2018; AMER and TODOROVIC, 2012; NOURANI *et al.*, 2020; AMER and TODOROVIC, 2016).

Probabilistic circuits have also been used for sequential data (MELIBARI, POUPART, DOSHI, and TRIMPONIAS, 2016), such as speech recognition and reconstruction (PEHARZ, KAPELLER, *et al.*, 2014; RATAJCZAK *et al.*, 2014; RATAJCZAK *et al.*, 2018) and natural language processing (CHENG *et al.*, 2014).

Remarkably, probabilistic circuits have seen a recent boom in hardware-aware research (SHAH, ISABEL GALINDEZ OLASCOAGA, *et al.*, 2019; OLASCOAGA *et al.*, 2019) and dedicated hardware for PCs in embedded systems (SOMMER *et al.*, 2018; SHAH, OLASCOAGA, MEERT, *et al.*, 2020; SHAH, OLASCOAGA, S. ZHAO, *et al.*, 2021).

Some other applications include robotics (SGUERRA and F. G. COZMAN, 2016; R. GEH and D. MAUÁ, 2019; ZHENG *et al.*, 2018; PRONOBIS *et al.*, 2017), biology (Cory J. BUTZ, SANTOS, *et al.*, 2018; Abram L. FRIESEN and P. DOMINGOS, 2015), probabilistic programming (STUHL"ULLER and GOODMAN, 2012; SAAD *et al.*, 2021), and fault localization (NATH and P. M. DOMINGOS, 2016).

Now that we have addressed the (arguably) most important aspects of probabilistic circuits, we are now ready to understand how to build them in a principled way.

3

Learning Probabilistic Circuits

As we have seen in [Chapter 2](#), inference in probabilistic circuits is, for the most part, straightforward. This is not so much the case when *learning* PCs. Despite the uncomplicated syntax, learning sufficiently expressive PCs in a principled way is comparatively harder to, say the usual neural network. For a start, we are usually required to comply with smoothness and decomposability to ensure marginalization at the least. This restriction excludes the possibility of adopting any of the most popular neural network patterns or architectures used in deep learning today. To make matters worse, constructing a PC graph more often than not involves costly statistical tests that make learning their structure a challenge for high dimensional data.

In this chapter, we review the most popular PC structure learning algorithms, their pros and cons, and more importantly, what can we learn from them to efficiently build scalable probabilistic circuits. We broadly divide existing structure learners into three main categories: divide-and-conquer (DIV, [Section 3.1](#)), incremental methods (INCR, [Section 3.2](#)) and random approaches (RAND, [Section 3.3](#)).

3.1 Divide-and-Conquer Learning

Arguably the most popular approach to learning the structure of probabilistic circuits are algorithms that follow a *divide-and-conquer* scheme. This class of PC learning algorithms, which here we denote by DIV, are characterized by recursive calls over (usually mutually exclusive) subsets of data in true divide-and-conquer fashion. This kind of procedure is more clearly visualized by LEARNSPN, the first, most well-known, and perhaps most archetypal of its class.

Before we start however, we must first address how we denote data. Data is commonly represented as a matrix where rows are assignments (of all variables), and columns are the values that each variable takes at each assignment. Let $\mathbf{D} \in \mathbb{R}^{m \times n}$ a matrix with m rows and n columns. We use $\mathbf{D}_{i,j}$ to access an element of \mathbf{D} at the i -th row, j -th column of matrix \mathbf{D} . We denote by $\mathbf{D}_{i,j}$, where $\mathbf{i} \subseteq [1..n]$ and $\mathbf{j} \subseteq [1..m]$ are sets of indices, a submatrix from the extraction of the \mathbf{i} rows and \mathbf{j} columns of \mathbf{D} . We use a colon as a shorthand for selecting all rows or columns, e.g. $\mathbf{D}_{:,j} = \mathbf{D}$, $\mathbf{D}_{:,j}$ is the j -th column and $\mathbf{D}_{i,:}$ is the i -th row.

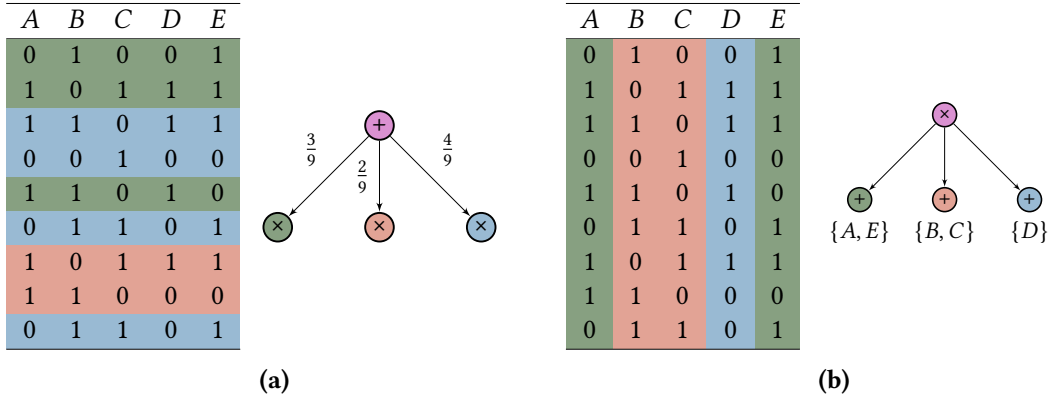


Figure 3.1: *LEARNSPN* assigns either rows (a) or columns (b) for sum and product nodes respectively. For sums, their edge weights are set proportionally to the assignments. For product children, scopes are defined by which columns are assigned to them.

Algorithm 6 LEARNSPN

Input Data D , whose columns are indexed by variables X

Output A smooth and decomposable probabilistic circuit learned from D

- 1: **if** $|X| = 1$ **then return** an input node learned from D
 - 2: **else**
 - 3: Find scope partitions $X_1, \dots, X_t \subseteq X$ st $X_i \perp\!\!\!\perp X_j$ for $i \neq j$
 - 4: **if** $k > 1$ **then return** $\prod_{j=1}^t \text{LEARNSPN}(D_{:X_j}, X_j)$
 - 5: **else**
 - 6: Find subsets of data $x_1, \dots, x_k \subseteq D$ st all assignments within x_i are all similar
 - 7: **return** $\sum_{i=1}^k \frac{|x_i|}{|D|} \cdot \text{LEARNSPN}(x_i, X)$
-

3.1.1 LEARNSPN

Recall the semantics of sum and product nodes in a smooth and decomposable probabilistic circuit. A sum is a mixture of distributions $p(X) = \sum_{i=1}^m w_i \cdot p_i(X)$ whose children scopes are all the same. A product is a factorization $p(X_1, \dots, X_n) = \prod_{i=1}^n p(X_i)$, implying that $X_i \perp\!\!\!\perp X_j$ for $i, j \in [n]$ and $i \neq j$. LEARNSPN (GENS and P. DOMINGOS, 2013) exploits these semantics in an intuitive and uncomplicated manner: sum children are defined by sub-PCs learned from similar (by some arbitrary metric) assignments, and product children are sub-PCs learned from data conditioned on the variables defined by their scope. In practice, this means that, for a dataset $D \in \mathbb{R}^{m \times n}$, sums assign rows to their children, while product children are assigned columns. This procedure continues recursively until data are reduced to a $k \times 1$ matrix, in which case a univariate distribution acting as input node is learned from it. This recursive procedure is shown more formally in Algorithm 6.

Notably, (GENS and P. DOMINGOS, 2013) purposely does not strictly specify which techniques should be used for assigning rows and columns, although they do provide empirical results on a particular form of LEARNSPN where row assignments are computed through EM clustering and products by pairwise G-testing. Instead, they call the algorithm a *schema* that incorporates several actual learning algorithms whose concrete form depends on the choice of how to split data.

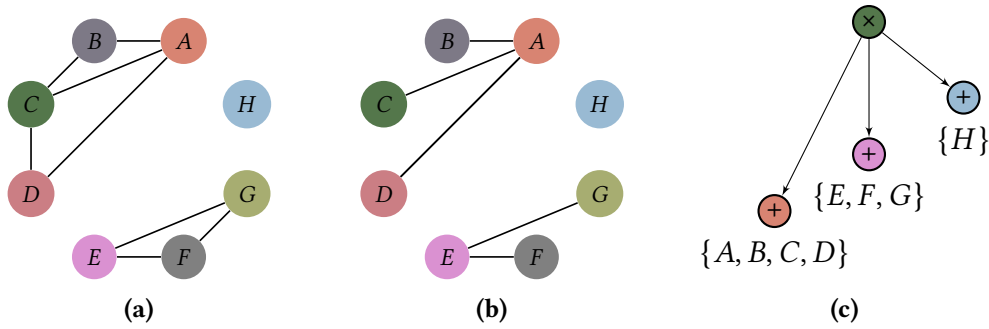


Figure 3.2: The pairwise (in)dependence graph where each node is a variable. In (a) we show the full graph, computing independence tests for each pair of variables in $\mathcal{O}(m^2)$. However, it suffices to compute for only the connected components (b), saving up pairwise computation time for reachable nodes. The resulting product node and scope partitioning is shown in (c).

Complexity

To be able to analyze the complexity of LEARNSPN, we assume a common implementation where sums are learned from k -means clustering, and products through pairwise G-testing. We know learning sums is efficient: k -means takes $\mathcal{O}(n \cdot k \cdot m \cdot c)$ time, where k is the number of clusters and c the number of iterations to be run. Products, on the other hand, are much more costly. The naïve approach would be to compute whether $X_i \perp\!\!\!\perp X_j$ for every possible combination. This is clearly quadratic on the number of variables $\mathcal{O}\left(\binom{m}{2} = \frac{m!}{2(m-2)!}\right)$ assuming an $\mathcal{O}(1)$ oracle for independence testing. In reality, G-test takes $\mathcal{O}(n \cdot m)$ time, as we must compute a ratio of observed versus expected values for each cell in the contingency table. This brings the total runtime for products to a whopping $\mathcal{O}(n \cdot m^3)$, prohibitive to any reasonably large dataset. In terms of space, independence tests most commonly used require either a correlation (for continuous data) or contingency (for discrete data) matrix that takes up $\mathcal{O}(m^2)$ space, another barrier for scaling up to high dimensional data.

Alternatively, instead of computing the G-test for every possible combination of variables, (GENS and P. DOMINGOS, 2013) constructs an independence graph \mathcal{G} whose nodes are variables and edges indicate whether two variables are statistically dependent. Within this context, the variable partitions we attribute to product children are exactly the connected components of \mathcal{G} , meaning it suffices testing only some combinations. Even so, this heuristic is still quadratic worst case. Figure 3.2 shows \mathcal{G} , the spanning forest resulted from the connected component heuristic, and the equivalent product node from this decomposition.

Pros and cons

Pros. Perhaps the main factor for LEARNSPN's popularity is how easily implementable, intuitive and modular it is. Even more remarkably, it is an empirically competitive PC learning algorithm despite its age, serving as a baseline for most subsequent works in PC literature. Lastly, the fact that each recursive call from LEARNSPN is completely independent from each the other makes it an attractive candidate for CPU parallelization.

Algorithm 7 EXTENDID**Input** Data \mathbf{D} , whose columns are indexed by variables \mathbf{X} , and memoization function \mathcal{M} **Output** A smooth and decomposable probabilistic circuit learned from \mathbf{D}

```

1: Find scope partitions  $\mathbf{X}_1, \dots, \mathbf{X}_t \subseteq \mathbf{X}$  st
2: if  $k > 1$  then
3:   for each  $j \in [t]$  do
4:      $N_j \leftarrow \text{LEARNMARKOV}(\mathbf{D}_{:, \mathbf{X}_j}, \mathbf{X}_j)$ 
5:     Associate  $\mathcal{M}(N_j)$  with  $\mathbf{D}_{:, \mathbf{X}_j}$  and  $\mathbf{X}_j$ 
6:   return  $\prod_{j=1}^t N_j$ 
7: else
8:   Find subsets of data  $\mathbf{x}_1, \dots, \mathbf{x}_k \subseteq \mathbf{D}$  st all assignments within  $\mathbf{x}_i$  are all similar
9:   for each  $i \in [k]$  do
10:     $N_i \leftarrow \text{LEARNMARKOV}(\mathbf{x}_i, \mathbf{X})$ 
11:    Associate  $\mathcal{M}(N_i)$  with  $\mathbf{x}_i$  and  $\mathbf{X}$ 
12:   return  $\sum_{i=1}^k \frac{|\mathbf{x}_i|}{|\mathbf{D}|} \cdot N_i$ 

```

Cons. Debatably, one of the key weakness of LEARNSPN is its tree-shaped computational graph, meaning that they are strictly less succinct compared to non-tree DAG PCs (MARTENS and MEDABALIMI, 2014). In terms of runtime efficiency, the algorithm struggles on high dimensional data due to the complexity involved in computing costly statistical tests. Despite Algorithm 6 giving the impression that no hyperparameter tuning is needed for LEARNSPN, in practice the modules for learning sums and products often take many parameters, most of which (if not all) are exactly the same for every recursive call. This can have a negative impact on the algorithm’s performance, since the same parameters are repeatedly used even under completely different data.

3.1.2 ID-SPN

A subtle yet effective way of improving the performance of LEARNSPN is to consider tractable probabilistic models over many variables as input nodes instead of univariate distributions. ID-SPN (ROOSHENAS and LOWD, 2014) does so by assuming that input nodes are Markov networks. Further, instead of blindly applying the recursion over subsequent sub-data, it attempts to compute some metric of quality from each node. The worst scored node is then replaced with a LEARNSPN-like tree. This is repeated until no significant increase in likelihood is observed. Algorithm 8 shows the ID-SPN pipeline, where EXTENDID is used in line 7 to grow the circuit in a divide-and-conquer fashion. The name ID-SPN comes from *direct* variable interactions, meaning the relationships modeled through the Markov networks as input nodes; and *indirect* interactions brought from the latent variable interpretation of sum nodes.

With respect to its implementation, ID-SPN is as modular as LEARNSPN in the sense that the data partitioning is left as a subroutine. Indeed, even the choice of input distributions is customizable: although ROOSHENAS and LOWD recommend Markov networks, any tractable distribution will do. Despite this seemingly small change compared to the original LEARNSPN algorithm, ID-SPN seems to perform better compared to its counterpart most

Algorithm 8 ID-SPN**Input** Data \mathbf{D} , whose columns are indexed by variables \mathbf{X} **Output** A smooth and decomposable probabilistic circuit learned from \mathbf{D}

- 1: Create a single-node PC: $C \leftarrow \text{LEARNMARKOV}(\mathbf{D}, \mathbf{X})$
- 2: Let \mathcal{M} a memoization function associating a node with a dataset and scope
- 3: Call C' a copy of C
- 4: **while** improving C yields better likelihood **do**
- 5: Pick worse node N from C'
- 6: Extract sub-data \mathbf{D}' and sub-scope \mathbf{X}' from $\mathcal{M}(N)$
- 7: Replace N with $\text{EXTENDID}(\mathbf{D}', \mathbf{X}', \mathcal{M})$
- 8: **if** C' has better likelihood than C **then** $C \leftarrow C'$
- 9: **return** C

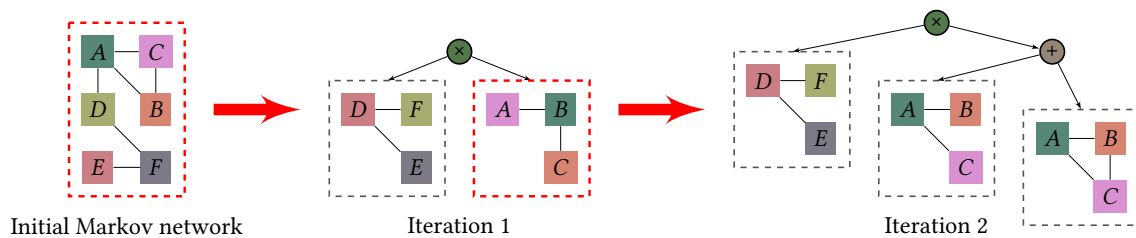


Figure 3.3: Two iterations of ID-SPN, where the contents inside the dashed line are Markov networks. The red color indicates that a node has been chosen as the best candidate for an extension with *EXTENDID*. Although here we only extend input nodes, inner nodes can in fact be extended as well.

of the time (ROOSHENAS and LOWD, 2014; JAINI, GHOSE, *et al.*, 2018), although at a cost to learning speed. Further, because of the enormous parameter space brought by having to learn Markov networks as inputs *and* perform the optimizations from sums and products, grid search hyperparameter tuning is infeasible. (ROOSHENAS and LOWD, 2014) recommend random search (BERGSTRA and BENGIO, 2012) as an alternative.

Complexity

As ID-SPN is a special case of LEARNSPN, the analysis for the sums and products subroutines holds. The only difference is on the runtime complexity for learning input nodes and the convergence rate for ID-SPN. Assuming input nodes are learned from the method suggested by ROOSHENAS and LOWD (2014), which involves learning a probabilistic circuit from a Markov network (LOWD and ROOSHENAS, 2013), then each “input” node takes time $\mathcal{O}(i \cdot c(k \cdot n + m))$, where i is the number of iterations to run, c is the size of the generated PC, and constant k is a bound on the number of candidate improvements to the circuit, which can grow exponentially for multi-valued variables. Importantly, opposite from LEARNSPN where we only learn input nodes once per call *if* data is univariate, ID-SPN requires learning multiple multivariate inputs for *every* *EXTENDID* call.

Pros and Cons

Pros. If we assume any multivariate distribution in place of Markov networks, PCs learned from ID-SPN are strictly more expressive than ones learned from LEARNSPN, as

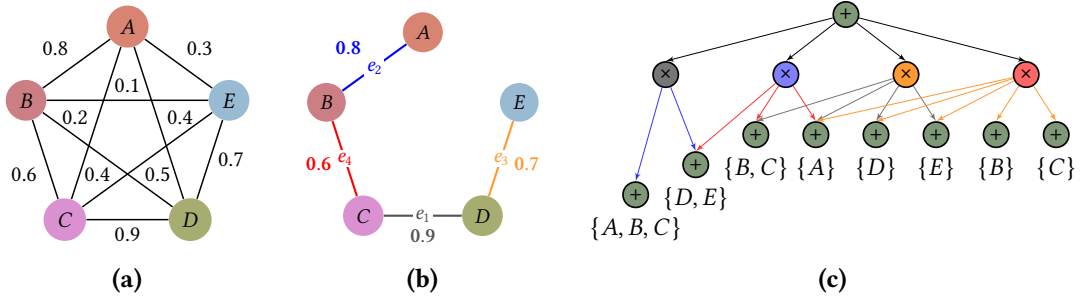


Figure 3.4: The fully connected correlation graph (a) with weights as the pairwise correlation measurements for each pair of variables; the maximum spanning tree for determining decompositions (b); and the mixture of decompositions (c). Colors in (b) match their partitionings in (c).

input nodes could potentially be replaced with LEARNSPN distributions. Additionally, the modularity inherited from LEARNSPN allows ID-SPN to adapt to data according to expert knowledge, bringing some flexibility to the algorithm.

Cons. Unfortunately, most of the disadvantages from LEARNSPN also apply to ID-SPN. Just like LEARNSPN, independence tests are more often than not a bottleneck for most executions with reasonably large number of variables. However, ID-SPN relies on a likelihood improvement for the computational graph to be extended, which ends up curbing the easy parallelization aspect of LEARNSPN. Besides, the complexity involved in learning Markov networks (or any other complex multivariate distribution as input node) carries a heavy weight during learning. This, coupled with the fact that hyperparameter tuning in the huge parameter space of ID-SPN must be done by a random search method, can take a heavy price in terms of learning time.

3.1.3 PROMETHEUS

So far, we have only considered structure learning algorithms that produce tree-shaped circuits. Even though ID-SPN *might* produce non-tree graphs at the input nodes depending on the choice of families of multivariate distributions, it does not do so as a rule. We now turn our attention to a PC learner that *does* generate non-tree computational graphs in a divide-and-conquer manner.

Recall that in both LEARNSPN and ID-SPN the scope partitioning is done greedily; we define a graph encoding the pairwise (in)dependencies of variables and greedily search for connected components by comparing independence test results with some correlation threshold, adding an edge if the correlation is sufficiently high. The choice of this threshold is often arbitrary and subject to hyperparameter tuning during learning, which is especially worrying when dealing with high dimensional data. In this section we review PROMETHEUS (JAINI, GHOSE, *et al.*, 2018), a divide-and-conquer LEARNSPN-like PC learning algorithm with two main features that stand out compared to the last two methods we have seen so far: (1) it requires no hyperparameter tuning for variable partitionings, and (2) accepts a more scalable alternative to computing all pairwise correlations.

Let \mathcal{G} the independence graph for scope $\mathbf{X} = \{X_1, X_2, \dots, X_m\}$. Remember that \mathcal{G} 's vertices are \mathbf{X} and each (undirected) edge $\overline{X_i X_j}$ coming from X_i to X_j means that $X_i \perp\!\!\!\perp X_j$.

Algorithm 9 PROMETHEUS**Input** Data D , whose columns are indexed by variables X **Output** A smooth and decomposable probabilistic circuit learned from D

```

1: if  $|X|$  is sufficiently small then return an input node learned from  $D$ 
2: else
3:   Find subsets of data  $x_1, \dots, x_k \subseteq D$  st all assignments within  $x_i$  are all similar
4:   Create a sum node  $S$  with initially no children and uniform weights
5:   for each  $x_i$  do
6:      $\mathcal{T} \leftarrow \text{CORRELATIONMST}(x_i, X)$ 
7:     for each weighted edge  $e_j$  in  $\mathcal{T}$  in decreasing order do
8:       Remove edge  $e_i$  from  $\mathcal{T}$ 
9:       Call  $S_1, \dots, S_t$  the scopes of each component in  $\mathcal{T}$ 
10:      Create product node  $P_j$  and associate it with  $S_1, \dots, S_t$ 
11:      Associate  $P_j$  with dataset  $x_i$ 
12:      Add  $P_j$  as a child of  $S$ 
13:   Let  $\mathcal{H}$  a hash table (initially empty) associating scopes to sum nodes
14:   for each  $P \in \text{Ch}(S)$  do
15:     for each scope  $S$  associated with  $P$  do
16:       if  $S \notin \mathcal{H}$  then
17:         Let  $x$  the dataset associated with  $P$ 
18:          $N \leftarrow \text{PROMETHEUS}(x_{\setminus S}, S)$ 
19:         Add  $N$  as a child of  $P$ 
20:          $\mathcal{H}_S \leftarrow N$ 
21:       else
22:         Add  $\mathcal{H}_S$  as a child of  $P$ 
23:   return  $S$ 

```

Previously, we constructed \mathcal{G} by comparing the output of an independence test (such as the G-test) against a threshold (e.g. a sufficiently low p -value). Instead, suppose \mathcal{G} is fully connected and that we attribute weights corresponding to a correlation metric of X_i against X_j for each edge (e.g. Pearson's correlation coefficient). The *maximum spanning tree* (MST) of \mathcal{G} , here denoted by \mathcal{T} , defines a graph where the removal of any edge in \mathcal{T} partitions the component into two subcomponents. Let e_i the i -th lowest (weight) valued edge; PROMETHEUS obtains a set of decompositions by iteratively removing edges from e_1 to $e_{|X|-1}$. In other words, the algorithm constructs a product node for each decomposition, assigning the scope of each child as the scope of each component at each edge removal. These products are then joined together by a parent sum node that acts as a mixture of decompositions. Figure 3.4 shows an example of \mathcal{T} , the subsequent decompositions, and the resulting mixture of decompositions.

Sum nodes are learned by clustering data into similar instances, just like in previous cases. Since the previously mentioned procedure involving products creates a mixture of decompositions (and thus a sum node), we can simply collapse the consecutive sum layers into a single sum node. Algorithm 9 shows the algorithm in its entirety. CORRELATIONMST computes the (fully connected) correlation graph, returning its MST. It is worth mentioning

that PROMETHEUS makes sure each recursive call shares subcircuits whenever scopes are the same (this is when the hash table \mathcal{H} in [Algorithm 9](#) comes into play). This avoids an exponential growth from the $k \cdot (|X| - 1)$ potential recursive calls.

Complexity

Up to now, the computation of decompositions is done by a $\mathcal{O}(m^2)$ construction of a fully connected correlation graph. This gives PROMETHEUS no asymptotic advantage over neither LEARNSPN nor ID-SPN. To change this, JAINI, GHOSE, *et al.* propose a more scalable alternative: in place of constructing the entire correlation graph, sample $m \log m$ variables and construct a correlation graph where only $\log m$ edges are added for each of these sampled variables instead, bringing down complexity to $\mathcal{O}(m(\log m)^2)$.

The analysis of sum nodes is exactly the same as LEARNSPN if we assume the same clustering method. If PROMETHEUS is implemented with the same multivariate distributions as ID-SPN at the input nodes, the analysis for those also holds.

Pros and Cons

Pros. The notable achievements of PROMETHEUS are evidently the absence of parameters for computing scope partitionings, reducing the dimension of hyperparameters to tune; a scalable alternative to partitionings that runs in sub-quadratic time; and (more debatably) the fact that the algorithm produces non-tree shaped computational graphs. Further, since product nodes are learned through correlation metrics, PROMETHEUS is easily adaptable to continuous data. To some extent, PROMETHEUS also inherits the modularity of LEARNSPN, as the choice of how to cluster and what input nodes to use is open to the the user.

Cons. Although the construction of the correlation graph in PROMETHEUS is not done greedily (at least in the quadratic version), selecting the decompositions (i.e. partitioning the graph into maximal components) is; of course, this is not exactly a drawback but a compromise, as graph partitioning is a known NP-hard problem ([FELDMANN and FOSCHINI, 2015](#)). Because PROMETHEUS accounts for all decompositions yielded from components after the removal of each edge from the MST, the circuit can grow considerably, even if we reuse subcircuits at each recursive call. An alternative would be to globally reuse subcircuits (i.e. share \mathcal{H} among different recursive calls) throughout learning, although this curbs expressivity somewhat, as these subcircuits are learned from possibly (completely) different data. Another option would be to bound the number of decompositions, or in other words remove only a bounded number of edges from the MST.

Remark 3.1.1: On variations of divide-and-conquer learning

Because of LEARNSPN's simplicity and modularity, there is a lot of room for improvement. This is reflected in the many works in literature on refining LEARNSPN to specific data, choosing the right parameters, producing non-tree shaped circuits, and choice of input nodes. In this remark segment, we briefly discuss other advances in divide-and-conquer PC learning.

As we have previously mentioned, one of the drawbacks of LEARNSPN is the

sheer volume of hyperparameters involved. VERGARI, MAURO, *et al.* (2015) suggests simplifying clustering to only binary row splits, while Y. LIU and LUO (2019) proposes clustering methods that automatically decide the number of clusters from data. Together with PROMETHEUS, the space of hyperparameters to tune is greatly reduced.

We again go back to the issue of reducing the cost of learning variable partitions. Apart from PROMETHEUS, DI MAURO *et al.* (2017) also investigate more efficient decompositions, proposing two approximate sub-quadratic methods to producing variable splits: one by randomly sampling pairs of variables and running G-test, and the other by a linear time entropy criterion.

VERGARI, MAURO, *et al.* (2015) proposes the use of Chow-Liu Trees as input nodes instead of univariate distributions, while MOLINA, NATARAJAN, *et al.* (2017) recommend Poisson distributions for modeling negative dependence. BUEFF *et al.* (2018) combines LEARNSPN with weighted model integration by learning polynomials as input nodes for continuous variables and counts for discrete data. MOLINA, VERGARI, *et al.* (2018) adapts LEARNSPN to hybrid domains by employing the randomized dependence coefficient for both clustering and variable partitioning, with pairwise polynomial approximations for input nodes.

Other contributions include adapting LEARNSPN to relational data (NATH and P. DOMINGOS, 2015), an empirical study comparing different techniques for clustering and partitioning in LEARNSPN (Cory J. BUTZ, OLIVEIRA, *et al.*, 2018), and LEARNSPN post-processing strategies for deriving non-tree graphs (RAHMAN and GOGATE, 2016).

3.2 Incremental Learning

Learning algorithms from the DIV class heavily rely on recursively constructing a probabilistic circuit in a top-down fashion. This facilitates learning, as we need only to greedily optimize at a local level. We now draw our attention to incremental algorithms that iteratively grow an initial circuit. These usually require a search over possible candidate nodes to be extended, and as such involve evaluating the entire circuit to determine best scores. In this section, we look at two examples of INCR class learning algorithms: LEARNPSDD and STRUDEL.

3.2.1 LEARNPSDD

As the name suggests, LEARNPSDD (LIANG, BEKKER, *et al.*, 2017) learns a smooth, structure decomposable and deterministic probabilistic circuit (see Section 2.3.2), meaning its computational graph must respect a vtree. We therefore must address the issue of learning the vtree before we turn to the PC learning algorithm *per se*.

Recall that for a vtree \mathcal{V} , every inner node $v \in \mathcal{V}$ with $\mathbf{X} = \text{Sc}(v^{\leftarrow})$ and $\mathbf{Y} = \text{Sc}(v^{\rightarrow})$ determines that \mathbf{X} and \mathbf{Y} are independent, i.e. $p_C(\mathbf{X}, \mathbf{Y}) = p_C(\mathbf{X})p_C(\mathbf{Y})$ for a PC C . This means that a PC's vtree is pivotal in embedding the independencies of the circuit's distribution.

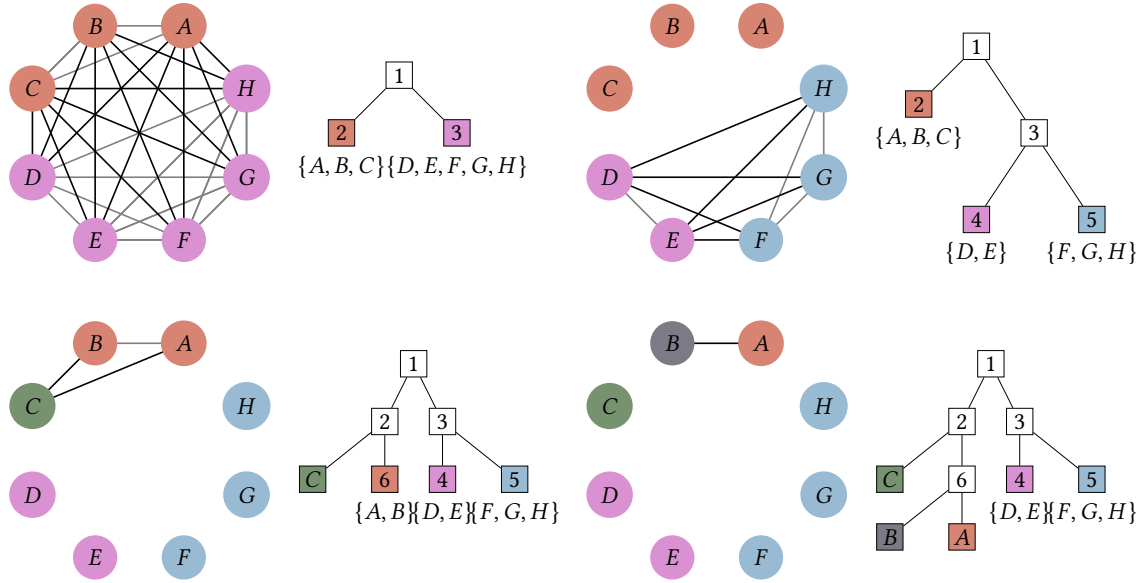


Figure 3.5: Snapshots of four iterations from running the vtree top-down learning strategy with pairwise mutual information. Each iteration shows a variable partitioning, the cut-set that minimizes the average pairwise mutual information as black edges, and the subsequent (partial) vtree. The algorithm finishes when all partitions are singletons.

With this in mind, [LIANG, BEKKER, et al. \(2017\)](#) propose two approaches to inducing vtrees from data, both of which use mutual information

$$\text{MI}(\mathbf{X}, \mathbf{Y}) = \sum_{\mathbf{x} \in \mathbf{X}} \sum_{\mathbf{y} \in \mathbf{Y}} p(\mathbf{x}, \mathbf{y}) \log \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})}$$

for deciding independence. To avoid computing an exponential number of MI terms, an approximation based on the average pairwise MI is computed instead

$$\text{pMI}(\mathbf{X}, \mathbf{Y}) = \frac{1}{|\mathbf{X}| |\mathbf{Y}|} \cdot \sum_{\mathbf{x} \in \mathbf{X}} \sum_{\mathbf{y} \in \mathbf{Y}} \text{MI}(\mathbf{x}, \mathbf{y}).$$

The first approach learns vtrees in a top-down fashion, starting with a full scope and recursively partitioning down to the unit set. The second learns bottom-up, starting with singletons and joining sets of variables up to full scope.

Top-down vtree learning. Let \mathcal{G} a fully connected weighted graph where variables are nodes. For each edge \overrightarrow{XY} , attribute its weight as $\text{MI}(\mathbf{x}, \mathbf{y})$. Learning the vtree top-down amounts to partitioning \mathcal{G} such that the cut-set that divides the two partitions \mathbf{X} and \mathbf{Y} is minimal with respect to pMI. [LIANG, BEKKER, et al. \(2017\)](#) further argue that balanced vtrees produce smaller PCs, and so they reduce learning to a balanced min-cut bipartition problem. Although this is known to be NP-complete ([GAREY and JOHNSON, 1990](#)), optimized solvers are able to produce high quality bipartitions efficiently ([KARYPIS and KUMAR, 1998](#)). In a nutshell, the vtree construction goes as follows: find a balanced min-cut bipartition (\mathbf{X}, \mathbf{Y}) in \mathcal{G} minimizing the pMI of the edges; add a vtree inner node representing this bipartition and connect it to the two vtrees produced by the recursive calls over \mathbf{X} and

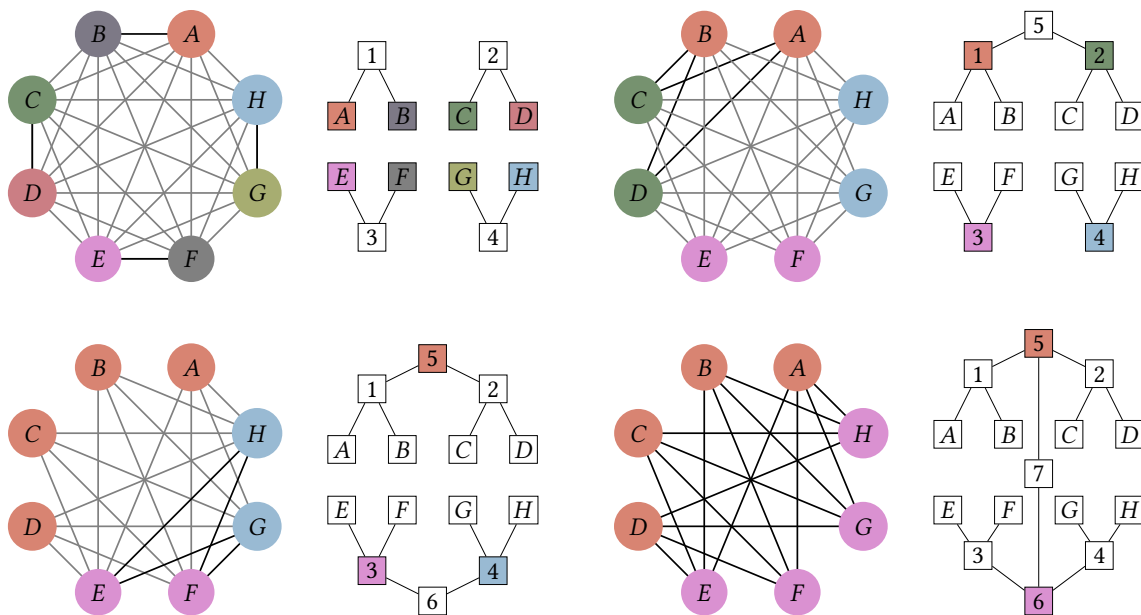


Figure 3.6: Snapshots from running the vtree bottom-up learning strategy with pairwise mutual information. Snapshots show pairings of two vtrees, with edges between partitions joined into a single edge whose weight is the average pairwise mutual information of all collapsed edges. In black are edges that correspond to the matchings that maximize the average pairwise mutual information. The algorithm finishes when all vtrees have been joined together into a single tree.

Y ; if $X = \{X\}$ (resp. $|Y| = \{Y\}$), produce a leaf node X (resp. Y). Figure 3.5 shows four iterations of this procedure.

Bottom-up vtree learning. Again, take \mathcal{G} as the fully connected weighted graph from computing the pairwise mutual information of variables. Now consider that every node of \mathcal{G} is a vtree whose only node is the variable itself. To learn a vtree bottom-up is to find pairings of vtrees such that the mutual information between them is high, meaning that the partitionings at higher levels are minimized (and so determine the “true” independence relationships between subsets of variables). To produce balanced vtrees, the algorithm attempts to join vtrees of same height whose pMI is maximal; this is equivalent to min-cost perfect matching, which can be solved, in our case, in $\mathcal{O}(m^4)$, where m is the number of variables (EDMONDS, 1965; KOLMOGOROV, 2009). Figure 3.6 exemplifies the algorithm.

LEARNPSDD is an incremental learning algorithm. This means that it takes an existing PC and incrementally grows the circuit by some criterion, preserving the structural constraints from the PC in the process. Once a vtree \mathcal{V} has been learned from data, we use it to construct an initial circuit that respects \mathcal{V} . The choice of circuit initialization is dependent on our task. For example, within the context of PSDDs, we are mostly interested in starting out with a PC induced from an LC encoding a certain knowledge base (see Section 2.3); this is usually done in a case-by-case basis, where LCs are compiled for a particular task and then promoted to PCs (see Remark 3.2.1). However, if one does not require specifying the distribution’s support, any PC will do.

How and where the circuit is grown – once we have acquired a vtree and an initial circuit – are the main topics of interest now. We first address the matter of *how*, i.e. how

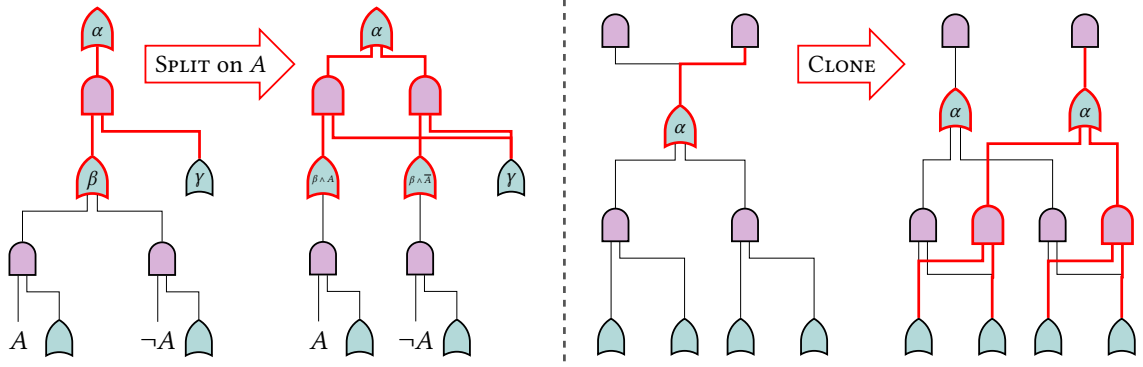


Figure 3.7: *SPLIT* (left) and *CLONE* (right) operations for growing a circuit when $m = 1$. Nodes and edges highlighted in red show the modified structure. In both cases smoothness, (structure) decomposability and determinism are evidently preserved.

can we increase a PC's expressivity such that we preserve a desired set of structural constraints; and later of *where*, i.e. which portions of the circuit are eligible for growth and how do we know they are good candidates.

LIANG, BEKKER, *et al.* propose two local transformations for growing a circuit C : *SPLIT* and *CLONE*. The first acts by multiplying a sum node's product child P into P_1, \dots, P_k products such that π_1, \dots, π_k (primes of P_1, \dots, P_k respectively) are mutually exclusive. This is done by attributing all possible values of a variable in $\text{Sc}(P)$, say A , to each prime, meaning that π_i will contain the assignment $A = i$ for every $i \in [k]$. This attribution is done by partially copying C_P into k circuits $C_P^{(1)}, \dots, C_P^{(k)}$ up to some depth m and then conditioning $C_P^{(i)}$ on $\llbracket A = i \rrbracket$. This is straightforward for the discrete case: at the appropriate vtree node (i.e. one that contains A as a leaf), replace the input node whose scope is A into an indicator node, setting it to the appropriate assignment of A . Although LIANG, BEKKER, *et al.* (2017) only considers the binary case, the transformation can be extended to the continuous if we consider k piecewise distributions whose support is over only a set interval. Naturally, input nodes must then have their support truncated to the appropriate i -th interval, which is no easy feat in the general case. The left side of Figure 3.7 shows *SPLIT* for the binary case.

The other proposed transformation, *CLONE*, does something similar for sum nodes. Pick a sum node S whose children are C_1, \dots, C_k and parents P_1 and P_2 ; double S and C_1, \dots, C_k , producing clones S' and C'_1, \dots, C'_k . Disconnect the edge coming from P_2 to S and instead connect it to S' . Connect all C'_1, \dots, C'_k to the same children as their original counterparts. This operation is visualized on the right side of Figure 3.7. One can further extend *CLONE* to apply this operation cloning nodes up to some depth m and then joining the last remaining deepest nodes similar to what was described for C'_1, \dots, C'_k .

It is easy to see that, in both cases, smoothness, structure decomposability and determinism are preserved. In fact, if the original circuit encodes a particular support (i.e. a knowledge base), the PC resulting from applying any of the two transformations must also encode the same support, since we have only made the underlying logic circuit more redundant. Probabilistically though, this “redundancy” only increases the parameterization space and as such increases the expressiveness of the PC. However, not all applications of

Algorithm 10 LEARNPSDD**Input** Data D , vtree \mathcal{V} , initial PC C , max depth m , scope X **Output** A smooth, structure decomposable and deterministic PC learned from D

```

1: while there is score improvement or has not reached the iteration/time limit do
2:    $s_S \leftarrow -\infty$ 
3:   Let  $(S^*, P^*)$  the best SPLIT candidate seen so far, initially empty
4:   for each candidate  $(S, P)$  of all possible SPLIT candidates do
5:      $C' \leftarrow \text{SPLIT}(C, S, P, \mathcal{V}, m)$ 
6:      $s' \leftarrow \text{Score}(D, C, C')$ 
7:     if  $s' > s_S$  then  $s_S \leftarrow s'$  and  $S^*, P^* \leftarrow S, P$ 
8:    $s_C \leftarrow -\infty$ 
9:   Let  $C^*$  the best CLONE candidate seen so far, initially empty
10:  for each candidate  $C$  of all possible CLONE candidates do
11:     $C' \leftarrow \text{CLONE}(C, C, \mathcal{V}, m)$ 
12:     $s' \leftarrow \text{Score}(D, C, C')$ 
13:    if  $s' > s_C$  then  $s_C \leftarrow s'$  and  $C^* \leftarrow C$ 
14:  if  $s_S > s_C$  then  $C \leftarrow \text{SPLIT}(C, S^*, P^*, \mathcal{V}, m)$ 
15:  else  $C \leftarrow \text{CLONE}(C, C^*, \mathcal{V}, m)$ 
16: return  $C$ 

```

SPLIT or CLONE are equal in terms of performance. While it is true that the application of SPLIT to any product node or CLONE to any sum node strictly increases expressivity, it is more meaningful to choose candidates whose growth carries a bigger impact on the overall fit relative to the training data. LEARNPSDD searches for reasonable candidates by computing

$$\text{Score}(D, C, C') = \frac{\log C'(D) - \log C(D)}{|C'| - |C|},$$

where C and C' are, respectively, the PCs before and after the application of any of the two operations. In other words, the algorithm randomly evaluates applying SPLIT and/or CLONE and ultimately chooses the one candidate that maximizes the log-likelihood of training data penalized by the size of the resulting PC, iteratively growing the circuit until there is no more improvement or reaches an iteration step or time limit, as Algorithm 10 shows.

Complexity

Although learning the vtree top-down reduces to an NP-complete min-cut graph partitioning problem, there are approximate algorithms that provide high quality partitionings in $\mathcal{O}(|X|^2)$ (KARYPIS and KUMAR, 1998). Learning bottom-up is reduced to min-cost perfect matching, which can be done in $\mathcal{O}(|X|^4)$ via the Edmonds Blossom algorithm (EDMONDS, 1965; KOLMOGOROV, 2009).

SPLIT runs, for a given variable X , in $\mathcal{O}(v \cdot |C|)$ if unbounded by m , where v is $|\text{Val}(X)|$, the number of possible assignments to X if X is discrete; or the number of intervals to fragment $\text{Val}(X)$ if X is continuous. CLONE's runtime is $\mathcal{O}(|C|)$ when m is unbounded, as it

needs to produce an almost exact copy of the circuit. We say that a local transformation, such as SPLIT or CLONE, is *minimal* when the copy depth is $m = 0$. When SPLIT and CLONE are minimal and X is binary, then the transformation is done in constant time. In fact, any non-minimal transformation can be composed out of minimal transformations (LIANG, BEKKER, *et al.*, 2017).

Perhaps the most costly routine of LEARNPSDD is its score function. Although log-likelihood is linear time computable on the number of edges of the circuit, C can grow substantially as transformations pile up. Each score evaluation requires four passes on the circuit: log-likelihoods and circuit sizes for both C and its updated circuit C' . However, since transformations are local, log-likelihood and circuit sizes only change for the nodes affected in the transformation and their ancestors, allowing LEARNPSDD to cache values. The overall complexity of LEARNPSDD at each iteration is therefore $\mathcal{O}(|C|^2)$ if we assume $m = 0$, with the first $|C|$ coming from the search of all candidates in C , and the second from the computation of Score. Each iteration further increases $|C|$, slowing down the algorithm's runtime.

Pros and Cons

Pros. The fact that LEARNPSDD preserves smoothness, structural decomposability, determinism *and* any logical semantic coming from its underlying LC is remarkable. On top of that, in theory and under minor modifications to SPLIT and CLONE, any PC is eligible as an initial circuit, even ones which do not respect any vtree. Besides, computing variable splits beforehand through a separate process of learning the vtree relieves the learning algorithm from having to compute costly statistical tests at each product node. Where LEARNPSDD really shines (and perhaps more fittingly PSDDs in general) is when the support is explicitly defined through the initial circuit's LC; because the PC attributes non-zero probability only to events where the LC does not return false, the circuit wastes no mass on impossible events.

Cons. In practice, LEARNPSDD is very slow even with caching; even worse, it may take several hours for only a minor (if any) improvement. (LIANG, BEKKER, *et al.*, 2017) suggests improving performance by producing ensembles of LEARNPSDDs, although this negates determinism in the final model (as well as structure decomposability if different vtrees are used at each component), denying the access to tractably computing queries like divergences, MI and entropies, not to mention the time cost to learn all components. Another issue is with the choice of the initial circuit. As previously mentioned, any circuit will do, however the performance (and efficiency) of LEARNPSDD is highly dependent on it. Within the context of PSDDs and encoding their support, LEARNPSDD requires that a separate algorithm compiles an LC for a specific task without looking at data. Although there are many ways of doing so, they are often not task agnostic (see Remark 3.2.1). More importantly, because the process of learning the circuit (from data) is decoupled from the task of encoding logical constraints imposed by a knowledge base, all variables that do not appear in the logic formula are compiled into a trivial form (e.g. fully factorized circuit). Lastly, although decoupling the process of learning the vtree from learning the PC helps with scalability, the ability of identifying the proper vtree for the most expressive PC given data is certainly desirable, and one which might be hindered by this separated

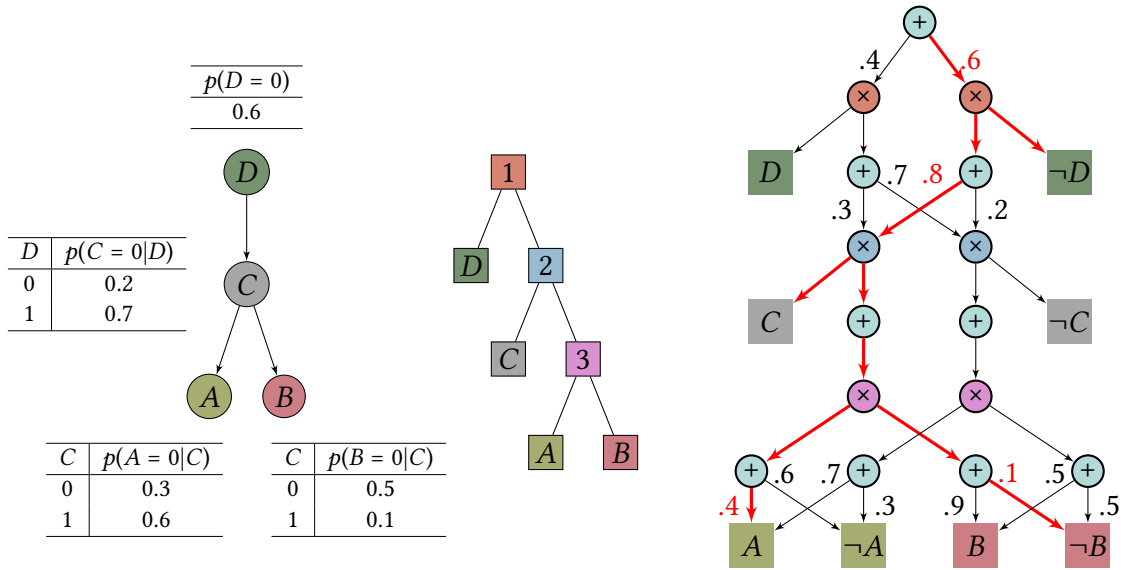


Figure 3.8: A vtree (middle) and probabilistic circuit (right) compiled from a Chow-Liu Tree (left). Each conditional probability $p(Y|X)$ is encoded as a (deterministic) sum node where each of the two children sets Y to 0 or 1. Colors in the CLT indicate the variables in the PC, while vtree inner node colors match with product nodes that respect them. Edges in red indicate the induced subcircuit activated on assignment $\{A = 1, B = 0, C = 1, D = 0\}$.

process.

3.2.2 STRUDEL

DANG *et al.* (2020) build upon the work of LEARNPSDD and propose STRUDEL, which mainly improves LEARNPSDD on two fronts: (1) by providing a simple algorithm for generating an initial circuit and vtree from data, and (2) proposing a heuristic for efficiently searching for good transformation candidates.

We first address how to construct the initial circuit from data. DANG *et al.* suggests doing so by compiling both a vtree and linear sized PC (in the number of variables) from a Chow-Liu Tree (CLT, CHOW and C. LIU, 1968). Let \mathcal{T} a CLT over variables $\mathbf{X} = \{X_1, \dots, X_m\}$. A vtree \mathcal{V} is extracted from \mathcal{T} by traversing \mathcal{T} top-down. For each node $X_i \in \mathcal{T}$, if X_i is a leaf node in \mathcal{T} , then create a vtree leaf node of X_i ; otherwise create an inner vtree node v , attach a vtree leaf node of X_i as v^{\leftarrow} and assign v^{\rightarrow} as a vtree built over all the vtrees coming from the children of X_i . The construction of v^{\leftarrow} depends on how balanced one wishes the vtree to be: if we want a more right-leaning vtree, it suffices to construct a right-linear vtree connecting all vtrees from each child $X_j \in \text{Ch}(X_i)$. Likewise, a balanced vtree is built by balancing the vtree connecting the recursive vtree calls from each X_j . Note that this does not necessarily mean that v^{\rightarrow} is completely right-linear or balanced, only that it is somewhat close to it, as the rest of the structure depends on the recursive calls of each CLT node.

STRUDEL compiles an initial circuit by looking at the vtree bottom-up and caching subcircuits. Let v a vtree node and $Y \in \mathcal{T}$ a CLT node with conditional probability $p(Y|X)$, where X is the parent of Y . If v is a leaf node in \mathcal{V} and v 's variable is also a leaf

node in \mathcal{T} , two sum nodes S_0 and S_1 over literal nodes $\neg Y$ and Y are created, each with weights $w_{S_0, \neg Y} = p(Y = 0|X = 0)$, $w_{S_0, Y} = p(Y = 1|X = 0)$ and $w_{S_1, \neg Y} = p(Y = 0|X = 1)$, $w_{S_1, Y} = p(Y = 1|X = 1)$. The two sum nodes connecting B and $\neg B$ in the PC shown on the right of Figure 3.8 show this exact case. The left sum node encodes $p(B|C = 1)$ and the right one $p(B|C = 0)$. These circuits are then cached by associating them with v . When Y is not a leaf node in \mathcal{T} but v is, we simply return literal nodes. If v is an inner node, we must define a scope partition, splitting $\mathbf{X} = \text{Sc}(v^{\leftarrow})$ and $\mathbf{Y} = \text{Sc}(v^{\rightarrow})$ into product nodes P_1, \dots, P_k , one for each value cached value in v . Each prime is set to the cached circuits from v^{\leftarrow} and each sub the cached circuits from v^{\rightarrow} . Finally, if two variables $X \in \mathbf{X}$ and $Y \in \mathbf{Y}$ are such that their parents are the same variable, say Z , then X and Y are independent when Z is given (because of a divergent connection in \mathcal{T}) and thus cannot be merged together into a single sum because of the context-specific independence set by Z (BOUTILIER *et al.*, 1996). This is visualized in the \otimes nodes; in this situation, A and B are siblings coming from C , and so $A \perp\!\!\!\perp B|C$ (redundant sum nodes are added for standardization). When the prior situation is not true, then not only X is the only variable in \mathbf{X} , but X must also be the parent of Y and so we must model $p(Y|X)$. This is the case for \otimes , where C is the parent of B and so we have to be join the two by sum nodes attributing the conditional probabilities $p(A, B|C = 0)$ for the right-most \otimes and $p(A, B|C = 1)$ for the left-most sibling. This procedure is shown more formally in Algorithm 11.

Now that we have an initial PC constructed from INITIALSTRUDEL, we are ready to discuss STRUDEL's second contribution. To do so, we must first understand the notion of *circuit flows* introduced in DANG *et al.* (2020). In short, the circuit flow of a deterministic probabilistic circuit C with respect to a variable assignment \mathbf{x} is the induced tree (see Definition 2.1.2) whose edges are all non-zero when C is evaluated under \mathbf{x} . Such an induced tree is unique in deterministic PCs because every sum node admits only one non-zero valued child for \mathbf{x} (or any assignment for that matter). Note how circuit flows are more specific in the sense they are intrinsically linked to an assignment, while induced subcircuits specify a deterministic subcircuit within its supercircuit.

The circuit flow of deterministic PCs helps us understand how to efficiently compute inference in circuits of that nature. As we briefly mentioned before, for any assignment \mathbf{x} in a smooth, decomposable and deterministic PC C , there exists a unique circuit flow \mathcal{F} that encodes the log-likelihood computation

$$C(\mathbf{x}) = \mathcal{F}_C(\mathbf{x}) = \prod_{(S,C) \in \text{Edges}(\mathcal{F}_C)} w_{S,C} \prod_{L \in \text{Inputs}(\mathcal{F}_C)} p_L(\mathbf{x}),$$

where $\text{Inputs}(\cdot)$ returns the set of input nodes of a circuit. When inputs are all binary, then one might encode \mathcal{F}_C as a mapping $f_C : \mathcal{X} \rightarrow \{0, 1\}^{|\mathcal{W}_C|}$, here \mathcal{W}_C denoting the set of all parameters (i.e. sum node weights) of C , which “activates” edge $w \in \mathcal{W}_C$ under assignment \mathbf{x} . With this, the above operation under log-space is reduced to a vector multiplication

$$\log C(\mathbf{x}) = f_C(\mathbf{x})^\top \cdot \log(\mathcal{W}_C).$$

Importantly, by aggregating circuit flows through counting the number of activations of each parameter $w_{S,C}$ in the entire training dataset \mathbf{D} , we get a sense of the number of samples $w_{S,C}$ impacts over \mathbf{D} , and thus a sense of how meaningful is that edge on the

Algorithm 11 INITIALSTRUDEL**Input** Data \mathbf{D} , whose columns are indexed by variables \mathbf{X} **Output** A smooth, structure decomposable and deterministic initial PC and vtree

```

1:  $\mathcal{T} \leftarrow \text{LEARNCLT}(\mathbf{D}, \mathbf{X})$ 
2:  $\mathcal{V} \leftarrow \text{COMPILEVTREE}(\mathcal{T})$ 
3: Let  $\mathcal{M}$  a hash table for caching circuits, initially empty
4: for each vtree node  $v \in \mathcal{V}$  in reverse topological order do
5:   if  $v$  is a leaf node then
6:     Let  $X \in \mathcal{T}$  the variable represented by  $v$ , and  $Y$  its parent
7:     if  $X$  is a leaf node in  $\mathcal{T}$  then
8:        $S_j \leftarrow \sum_{i \in \text{Val}(X)} p(X = i | Y = j) \cdot \llbracket X = i \rrbracket$  for each  $j \in \text{Val}(Y)$ 
9:        $\mathcal{M}(v) \leftarrow \mathcal{M}(v) \cup \{S_j | \forall j \in \text{Val}(Y)\}$ 
10:    else
11:       $\mathcal{M}(v) \leftarrow \mathcal{M}(v) \cup \{\llbracket X = i \rrbracket | i \in \text{Val}(X)\}$ 
12:    else
13:      Attribute  $\mathbf{X} \leftarrow \text{Sc}(v^{\leftarrow})$  and  $\mathbf{Y} \leftarrow \text{Sc}(v^{\rightarrow})$ 
14:      Let  $X \in \mathbf{X}$  and  $Y \in \mathbf{Y}$ 
15:      Attribute  $\mathbf{N}^{\leftarrow} \leftarrow \mathcal{M}(v^{\leftarrow})$  and  $\mathbf{N}^{\rightarrow} \leftarrow \mathcal{M}(v^{\rightarrow})$ 
16:       $k \leftarrow |\text{Val}(X)|$ 
17:      Construct product nodes  $\mathbf{P} = \{\mathbf{N}_i^{\leftarrow} \times \mathbf{N}_i^{\rightarrow} | \forall i \in [k]\}$ 
18:      if  $\text{Pa}(X) = \text{Pa}(Y)$  then
19:        Create sum nodes  $S_i$  each with only a single child  $P_i \in \mathbf{P}$ , for each  $i \in [k]$ 
20:         $\mathcal{M}(v) \leftarrow \mathcal{M}(v) \cup \{S_1, \dots, S_k\}$ 
21:      else
22:         $S_j \leftarrow \sum_{i \in \text{Val}(X)} p(\mathbf{Y} | X = i) \cdot \llbracket X = i \rrbracket$ , for each  $j \in \text{Val}(Y)$ 
23:         $\mathcal{M}(v) \leftarrow \mathcal{M}(v) \cup \{S_j | \forall j \in \text{Val}(Y)\}$ 
24: return  $\mathcal{M}(v_r)$ , where  $v_r$  is  $\mathcal{V}$ 's root node

```

fitness of data. As we shall see briefly, this aggregated circuit flow shall then be used as a score for a greedy search over the space of candidates for local transformations.

To overcome the scalability limitations of LEARNPSDD, STRUDEL proposes using only SPLIT to reduce the search space, looking at performing the search greedily instead of exhaustively and exploiting the efficiency of aggregate circuit flows as a fast heuristic in place of computing the whole likelihood. Searching is done by finding the edge to SPLIT whose aggregate circuit flow is maximal

$$\text{Score}_{\text{eFlow}}(w_{S,C} | \mathbf{C}, \mathbf{D}) = \sum_{\mathbf{x} \in \mathbf{D}} f_{\mathbf{C}}(\mathbf{x}) [w_{S,C}],$$

while the choice of which variable to condition SPLIT on is done by selecting the variable X that shares the most dependencies (and thus the higher pairwise mutual information) with other variables within the scope of that edge, estimated from the aggregate flows

$$\text{Score}_{\text{vMI}}(X, w_{S,C} | \mathbf{C}, \mathbf{D}) = \sum_{\substack{Y \in \text{Sc}(S) \\ Y \neq X}} \text{MI}(X, Y).$$

Algorithm 12 STRUDEL**Input** Data \mathbf{D} , max depth m , scope \mathbf{X} **Output** A smooth, structure decomposable and deterministic PC learned from \mathbf{D}

- 1: $\mathcal{C}, \mathcal{V} \leftarrow \text{INITIALSTRUDEL}(\mathbf{D}, \mathbf{X})$
- 2: **while** there is score improvement of has not reached the iteration/time limit **do**
- 3: Compute the aggregate flow over all edges
- 4: $w_{S,C}^* \leftarrow \arg \max_{w \in \mathcal{W}_C} \text{Score}_{\text{eFlow}}(w|\mathcal{C}, \mathbf{D})$
- 5: $X^* \leftarrow \arg \max_{X \in \text{Sc}(\mathcal{S})} \text{Score}_{\text{vMI}}(X, w_{S,C}^*|\mathcal{C}, \mathbf{D})$
- 6: $\mathcal{C} \leftarrow \text{SPLIT}(\mathcal{C}, \mathcal{S}, \mathcal{C}, \mathcal{V}, m)$
- 7: **return** \mathcal{C}

The entire algorithm for STRUDEL is showcased in [Algorithm 12](#).

Complexity

Learning the Chow-Liu Tree is done in $\mathcal{O}(|\mathbf{X}|^2 \cdot |\mathbf{D}|)$ through Chow-Liu’s algorithm ([CHOW and C. LIU, 1968](#)), while the vtree is compiled in time linear to the size of the CLT, i.e. $\mathcal{O}(|\mathbf{X}|)$ since the Bayesian network is a tree. Consequentially, INITIALSTRUDEL runs in $\mathcal{O}(|\mathbf{X}| \cdot |\text{Val}(\mathbf{X})|)$, or linear on $|\mathbf{X}|$ if we assume binary variables as originally intended. The bulk of the computation falls under STRUDEL, which runs in $\mathcal{O}(|\mathbf{X}|^2 \cdot |\mathbf{D}| + i(|\mathcal{C}| \cdot |\mathbf{D}| + |\mathbf{X}|^2))$ assuming a bounded max depth m and binary variables. Term $|\mathbf{X}|^2 \cdot |\mathbf{D}|$ corresponds to learning the CLT, $|\mathcal{C}| \cdot |\mathbf{D}|$ to the computation of the aggregate circuit flows, $|\mathbf{X}|^2$ to the computation of $\text{Score}_{\text{vMI}}$ which involves the pairwise mutual informations of \mathbf{X} , and i the number of iterations of STRUDEL.

Pros and Cons

Pros. Arguably, the most valuable contribution of STRUDEL is its improvement on LEARNPSDD’s scalability. Compared to LEARNPSDD, STRUDEL can take orders of magnitude less time per iteration, which in practice means a higher number of transformations accomplished in the same range of time. In addition, the nature of circuit flows allows for easy vectorization and thus CPU or GPU parallelization. In terms of data fitness, [DANG et al. \(2020\)](#) empirically shows that initial circuits constructed from STRUDELINITIAL greatly improve performance compared to fully factorized initial PCs from LEARNPSDD. Similar to LEARNPSDD, one can learn an ensemble of STRUDELS to further boost performance at the cost of losing determinism. Opposite to LEARNPSDD however, [DANG et al.](#) employ structure-sharing components so that the act of learning the circuit’s structure is done once, greatly reducing learning time. Parameters are then learned through closed form EM (see [Remark 3.3.1](#)) and bagging.

Cons. Although STRUDEL’s greedy heuristic search strategy translates into possibly more accurate PCs, it also produces more sizable circuits when compared to the exhaustive search of LEARNPSDD. Indeed, [DANG et al. \(2020\)](#)’s empirical evaluation shows STRUDEL PCs up to 12 times bigger than LEARNPSDD’s with the two somewhat tied in terms of fitness. This is especially worrying given that STRUDEL’s complexity grows with its circuit size. In fact, experiments show a sharp increase in seconds per iterations for the two

INCR algorithms, with both reaching multiple digits for each iteration even in smaller sized datasets; though LEARNPSDD much sharper and sooner comparatively (DANG *et al.*, 2020).

Remark 3.2.1: On the choice of initial circuits

We only briefly mentioned in Section 3.2.1 how we might want to start out with an initial PC conveying a specific support and then run an INCR class algorithm to further boost its probabilistic expressiveness without changing the underlying knowledge base. We devote this remark segment to discussing several works in literature that construct a so-called *canonical* (i.e. minimal with respect to their size without sacrificing its logical semantics) logic circuit, becoming perfect candidates to be used as an initial circuit in INCR learners.

Just like in probabilistic reasoning, the field of knowledge compilation and symbolic reasoning is often interested in finding succinct representations capable of tractably computing queries, a subject which we briefly touched in Section 2.3.1. For this reason, smooth, structure decomposable and deterministic logic circuits, who usually go by the name of Sentential Decision Diagrams (SDDs, DARWICHE, 2011) have proven to be a useful tool in several applications (VLASSELAER, RENKENS, *et al.*, 2014; VLASSELAER, BROECK, *et al.*, 2015; LOMUSCIO and PAQUET, 2015; HERRMANN and BARROS, 2013). Fortunately, both LEARNPSDD and STRUDEL preserve all the necessary structural constraints for both logical *and* probabilistic queries in (P)SDDs. With this in mind, we highlight compilation of SDDs in this short remark.

For most cases, SDDs can be compiled directly from CNFs and DNFs. (A. CHOI and DARWICHE, 2013) constructs SDDs bottom-up by first compiling C/DNF clauses and then combining smaller SDDs by either conjoining or disjoining them. In contrast, (OZTOK and DARWICHE, 2015) presents a faster compilation process which recursively breaks down C/DNFs by decomposing the formula into components according to a vtree, and then combines them into an SDD.

Although CNFs and DNFs are the most common form of encoding propositional knowledge bases, they struggle under specific logical constraints such as cardinality constraints (NISHINO *et al.*, 2016; SINZ, 2005). Interesting alternatives include BDDs (see Example 2.3.1) which are also widely used in formal methods and program verification, and for which efficient compilation from cardinality constraints are available (EÉN and SÖRENSON, 2006). Because BDDs are special case SDDs whose vtrees are always right-linear (DARWICHE, 2011; BOVA, 2016), their reduced representations (BRYANT, 1986) are natural initial circuit candidates (see Section 4.1.1).

We now cover some of the existing literature on producing task specific (P)SDDs. A. CHOI, TAVABI, *et al.* (2016) analyzes the feasibility of compiling LCs (and subsequently producing a PC by parameterizing disjunction edges) from tic-tac-toe game traces, and route planning within a city. Both involve exhaustively disjoining all permutations of valid conjoined configurations and compiling the resulting DNF through previously cited SDD compilers. A. CHOI, SHEN, *et al.* (2017) further studies route planning by compiling them into SDDs, but analyze the feasibility of

(P)SDDs in route planning in larger scale maps. A. CHOI, BROECK, *et al.* (2015) explore (P)SDDs in preference learning and rankings, providing an algorithm for compiling an SDD from total or partial rankings. Similarly, SHEN *et al.* (2017) investigates (P)SDDs in probabilistic cardinality constraint tasks, also known as subset selection or n -choose- k models.

3.3 Random Learning

We now explore random approaches to constructing probabilistic circuits, which we classify as RAND class learning algorithms. Essentially, RAND class circuits are constructed either by a completely random procedure (Section 3.3.1), or guided by data (Section 3.3.2). We first look at RAT-SPN (PEHARZ, VERGARI, *et al.*, 2020), a connectionist PC structure learning algorithm for randomly generating tensorized smooth and decomposable probabilistic circuits. We then address XPCs (MAURO *et al.*, 2021), a flexible algorithm capable of learning smooth and decomposable PCs as well as structure decomposable and/or deterministic circuits through simple modifications to their method.

3.3.1 RAT-SPN

A key ingredient to RAT-SPN (PEHARZ, VERGARI, *et al.*, 2020) is the concept of *region graphs*. First introduced in PC literature in DENNIS and VENTURA (2012), region graphs are tensorized templates for PC construction. Informally, a region graph is composed out of *region nodes* and *partition nodes*; the former a set of sum or input nodes, and the latter of products. Regions can be thought of sets of computational units explaining the same interactions among variables (DENNIS and VENTURA, 2012), for instance semantically similar pixel regions in an image; while partitions define independencies between these regions. Edges coming out of region (resp. partition) nodes must necessarily connect to a partition (resp. region) node.

Definition 3.3.1 (Region graph). *A region graph is a rooted connected DAG whose nodes are either regions or partitions. Children of regions are partitions, and children of partitions are regions. The root is always a region node.*

Region graphs simplify the process of constructing PCs by ensuring that they are *at the least* smooth and decomposable. Call \mathcal{G} a region graph; \mathcal{G} is easily translated to a PC by applying the procedure described in Algorithm 13. Every region is compiled into a set of sums or inputs, fully connecting children; every partition into a set of products, producing a distinct permutation of children. Evidently, the resulting PC is exponential on s and l , as products must ensure that they encode different permutations. To deal with this blow-up, this number is often restricted to only two children per partition.

RAT-SPN works by randomly generating a region graph in a top-down divide-and-conquer approach similar to LEARN-SPN, except that the learned structure eventually produces non-tree shaped circuits and the procedure is done *completely* random (see Algorithm 14). In fact, PEHARZ, VERGARI, *et al.* (2020) argue that the parameterization of the circuit by means of sum weights is as important as its structure, looking at probabilistic

Algorithm 13 COMPILEREGIONGRAPH**Input** A region graph \mathcal{G} , parameters s for sums and l for inputs**Output** A smooth and decomposable probabilistic circuit

```

1: Let  $\mathcal{M}$  a mapping of region nodes to PC nodes
2: for each node  $N$  in  $\mathcal{G}$  except the root in reverse topological order do
3:   if  $N$  is a region then
4:     if  $N$  is a leaf node in  $\mathcal{G}$  then
5:       Construct  $\mathbf{L} = \{L_1, \dots, L_l\}$  input nodes over variables  $\text{Sc}(N)$ 
6:       Associate  $N$  with  $\mathbf{L}$ 
7:     else
8:       Construct  $\mathbf{S} = \{S_1, \dots, S_s\}$  sum nodes
9:       for each partition node  $P \in \text{Ch}(N)$  do
10:        Every sum in  $\mathbf{S}$  connects with every product in  $\mathcal{M}(P)$ 
11:     else if  $N$  is a partition then
12:       Let  $\text{Ch}(N) = \{R_1, \dots, R_k\}$  regions and  $q = \prod_{i=1}^k |\mathcal{M}(R_i)|$ 
13:       Construct  $\mathbf{P} = \{P_1, \dots, P_q\}$  product nodes
14:       for every product  $P \in \mathbf{P}$  do
15:        Connect  $P$  with a distinct combination of sums in  $\mathcal{M}(R_1), \dots, \mathcal{M}(R_k)$ 
16: Construct a root node  $R$ 
17: Connect  $R$  to all products in every child of  $\mathcal{G}$ 's root
18: return  $R$ 

```

circuits as a specific subclass of neural networks. Indeed, they show that this connectionist approach heavily inspired by traditional deep learning produces very competitive PCs. However, to do so requires extensive optimization of the circuit's weights, which unsurprisingly is where RAT-SPN shines: because of the tensorized nature of region graphs, the resulting PC is able to exploit the advantages of known deep learning frameworks, making the most of efficient stochastic gradient descent optimizers and GPU parallelization.

To ensure that the compiled PC is smooth and decomposable, the region graph must also be so. We extend the definition of scope function to region graphs. As long as leaf region nodes are assigned the correct scope, the PC is by construction smooth (every region is fully connected to their children) and decomposable (every partition splits variables into two regions). Function `CREATELAYER` in [Algorithm 14](#) does exactly that, making sure each partition decomposes into two distinct variable splits down to leaf region nodes. A PC is then constructed by following `COMPILEREGIONGRAPH` and passing the random region graph and number of nodes per region and partition nodes as parameters. The function ultimately produces a dense probabilistic circuit from the region graph blueprint, as [Figure 3.9](#) exemplifies.

Once the PC structure is successfully generated, [PEHARZ, VERGARI, et al. \(2020\)](#) suggest Expectation-Maximization (EM, [DEMPSTER et al., 1977](#); [PEHARZ, GENS, PERNKOPF, et al., 2016](#); [H. ZHAO, POUPART, et al., 2016](#)) for optimizing the circuit parameters (i.e. sum weights and input node distributions). Although parameter learning of PCs is not the focus of this dissertation, we briefly touch the subject in [Remark 3.3.1](#).

Algorithm 14 RAT-SPN**Input** Data \mathbf{D} , variables \mathbf{X} , max depth d , r # subcircuits, s # sums, and l # inputs**Output** A smooth and decomposable probabilistic circuit

```

1: function CREATELAYER( $R, d, \mathbf{X}$ )
2:   Assign  $\mathbf{X}$  as  $\text{Sc}(R)$ 
3:   Sample a variable split  $(Y, Z)$  from  $\mathbf{X}$ 
4:   Create a partition  $P$  and add it as a child of  $R$ 
5:   if  $d > 1$  then
6:     if  $|Y| > 1$  then
7:       Create a region  $R_1$ 
8:       CREATELAYER( $R_1, d - 1, Y$ )
9:     if  $|Z| > 1$  then
10:      Create a region  $R_2$ 
11:      CREATELAYER( $R_2, d - 1, Z$ )
12: Start with a root region node  $R$ 
13: for each  $i \in [r]$  do
14:   CREATELAYER( $R, d, \mathbf{X}$ )
15:  $C \leftarrow \text{COMPILEREGIONGRAPH}(R, s, l)$ 
16: return  $C$ 

```

Worthy of note is a discriminative version of RAT-SPN where instead of a single root sum node, k roots are learned, each connecting to every CREATELAYER subcircuit. Each i -th root describes the conditional probability $p(\mathbf{X}|Y = i)$, where Y is the query variable and \mathbf{X} evidence. Classification follows directly from Bayes Rule $p(Y|\mathbf{X}) = \frac{p(\mathbf{X}|Y)p(Y)}{\sum_{i=1}^k p(\mathbf{X}|Y=i)p(Y=i)}$, where $p(Y)$ is either estimated from the training data or assumed to be fixed. Accordingly, a discriminative objective function is proposed involving cross-entropy and log-likelihood for hybrid generative-discriminative optimization (BOUCHARD and TRIGGS, 2004) instead of running EM.

Complexity

Although the procedure described in Algorithm 14 is $\mathcal{O}(r \cdot d(s + l))$ and therefore extremely fast, it does not paint the whole picture. The main bulk of the complexity when learning RAT-SPN comes from parameter learning. PEHARZ, VERGARI, *et al.* (2020) calculate the number of sum weights to be

$$|\mathcal{W}_C| = \begin{cases} r \cdot k \cdot l^2 & \text{if } d = 1, \\ r \cdot (k \cdot s^2 + (2^{d-1} - 2)s^3 + 2^{d-1} \cdot s \cdot l^2) & \text{if } d > 1; \end{cases}$$

meaning that learning only the non-input parameters of RAT-SPN takes time $\mathcal{O}((r \cdot 2^d \cdot (s^3 + s \cdot l^2) + r \cdot k \cdot s^2) \cdot |\mathbf{D}|)$. However, given that all other structure learning algorithms covered in this dissertation also require learning circuit parameters, one might argue that the true cost of structure learning in RAT-SPN is indeed subquadratic on the parameters.

4

Scalable Learning of Probabilistic Circuits

4.1 A Logical Perspective

4.1.1 SAMPLEPSDD

4.2 A Data Perspective

4.2.1 LEARNRP



Appendix

A.1 Proofs

Theorem A.1.1. *Let C a probabilistic circuit whose first l layers are composed solely of sum nodes. Call \mathbf{N} the set of all nodes in layer $l + 1$. C is equivalent to a PC C' whose root is a sum node with \mathbf{N} as children.*

Proof. We adapt a similar proof due to [JAINI, POUPART, et al. \(2018\)](#). Every sum node is of the form

$$S(\mathbf{x}) = \sum_{C \in \text{Ch}(S)} w_{S,C} \cdot C(\mathbf{x}).$$

Particularly, every child C in a sum node in layer $1 \leq i \leq l - 1$, is a sum node, and so for the first layer we have that

$$\begin{aligned} S(\mathbf{x}) &= \sum_{C_1 \in \text{Ch}(S)} w_{S,C_1} \sum_{C_2 \in \text{Ch}(C_1)} w_{C_1,C_2} C_2(\mathbf{x}) \\ &= \sum_{C_1 \in \text{Ch}(S)} \sum_{C_2 \in \text{Ch}(C_1)} w_{S,C_1} w_{C_1,C_2} C_2(\mathbf{x}). \end{aligned}$$

Define a one-to-one mapping that takes a tuple (C_1, C_2) where $C_1 \in \text{Ch}(S)$ and $C_2 \in \text{Ch}(C_1)$ and returns a (unique) path from S to every grandchild C_2 of S . Call \mathbf{K} the set of all paths, and w_{S,C_1} and w_{C_1,C_2} the weights for one such path. We can merge these two weights into a single weight $w'_{S,C_2} = w_{S,C_1} \cdot w_{C_1,C_2}$, yielding

$$S(\mathbf{x}) = \sum_{(w_{S,C_1}, w_{C_1,C_2}) \in \mathbf{K}} w'_{S,C_2} C_2(\mathbf{x}).$$

This ensures that two consecutive sum layers can be collapsed into a single layer. Particularly, for the first (root) and second layers, the above transformation generates a circuit with one fewer layer and whose root has $\mathcal{O}(nm)$ edges, where n and m are the number of edges coming from the original root and its children respectively. We can apply this

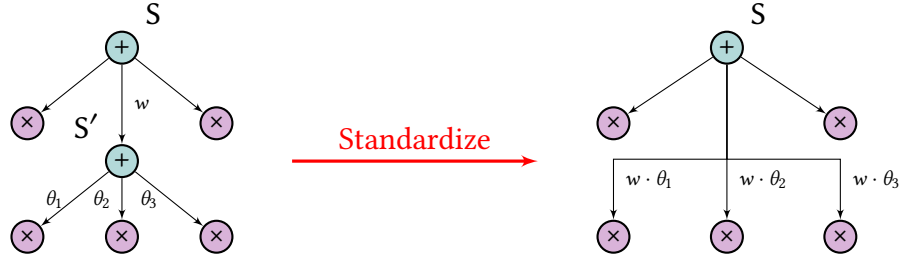
procedure until there are no more consecutive sum nodes. This results in a PC of the form

$$S(\mathbf{x}) = \sum_{C \in \text{Ch}(S)} w_{S,C} N(\mathbf{x}),$$

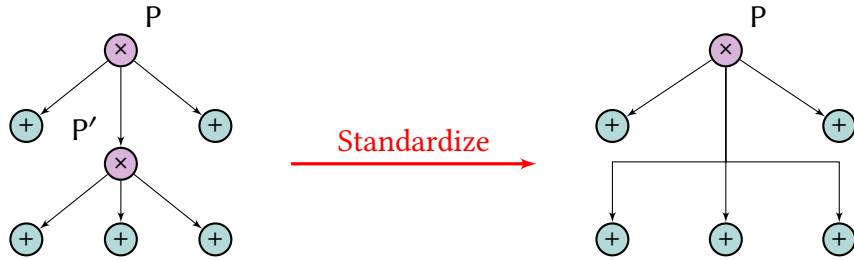
where $N \in \mathbf{N}$. The number of children of the resulting root sum node will be exponential on the number of edges of its children. \square

Theorem A.1.2 (Standardization). *Any probabilistic circuit C can be reduced to a circuit where every sum node contains only products or inputs and every product node contains only sums or inputs.*

Proof. If C is already standard we are done. Otherwise, there exists either (i) a sum node S with a sum S' as child; or (ii) a product node P with a product P' as child. We first address (i): let w be the weight of edge $\overrightarrow{S S'}$ and θ_i the weights from all edges coming out from S' .



Connect S with every child of S' , assigning as weight $w \cdot \theta_i$ for each child i . Delete S' and all edges coming out from it. The resulting circuit is computationally equivalent but now without a consecutive pair of sums. This transformation is visualized by the figure above. We do a similar procedure in (ii), but now instead remove P' and connect all children of P' to P , as we show below.



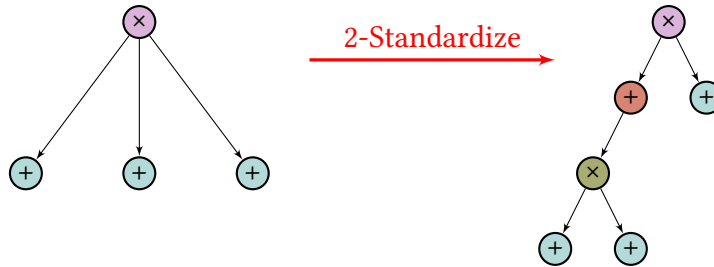
\square

Theorem A.1.3 (2-Standardization). *Any probabilistic circuit C can be transformed into a circuit where every sum node contains only products or inputs and every product node contains only two sums or inputs.*

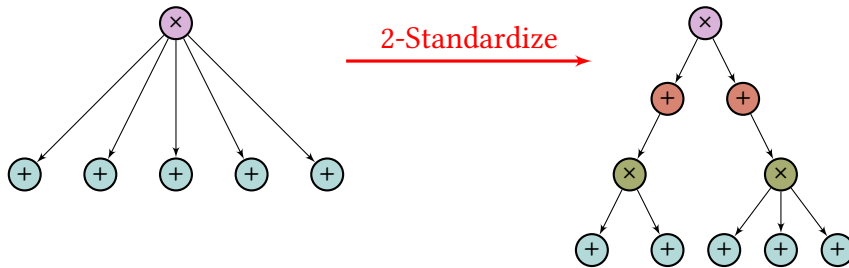
Proof. For sums, apply the same standardization procedure as [theorem A.1.2](#). Let P a product and call $n = |\text{Ch}(P)|$. If $n = 1$ and $\text{Ch}(P)$ is a product, then remove P and connect

all previous parents of P with its child. If $n = 1$ and $\text{Ch}(P)$ is not a product, remove P and apply the standardization procedure for sums on all of $\text{Pa}(P)$.

For $n > 2$, we simply need to split into 2-products recursively. We prove this by induction. The base case is when $n = 2$, which is already done, or $n = 3$, in which case we need apply the transformation below.



Where \oplus and \otimes are newly introduced nodes. When $n > 3$, we create two products P_1 and P_2 , each connected with a sum and product, and with $\lfloor \frac{n}{2} \rfloor$ and $\lceil \frac{n}{2} \rceil$ potential children. By the induction hypothesis, we can recursively binarize the subsequent grandchildren products.



As an example, we have $n = 5$ in the figure above. We introduce the sums \oplus and products \otimes and then recursively apply the transformation again on the \otimes s.

When $\text{Ch}(P)$ are product nodes we do the same procedure as before, but with the added post-process addition of a sum node connecting \otimes to every $\text{Ch}(P)$. \square

Theorem 2.2.1 (POON and P. DOMINGOS, 2011; Y. CHOI, VERGARI, and VAN DEN BROECK, 2020; VERGARI, Y. CHOI, *et al.*, 2021). *Let C a smooth and decomposable PC. Any one of EVI, MAR or CON can be computed in linear time (in the number of edges of C).*

Proof. For a sum node S , we have the following marginalization query

$$\begin{aligned} \int S(\mathbf{x}, \mathbf{y}) \, d\mathbf{y} &= \int \sum_{C \in \text{Ch}(S)} w_{S,C} C(\mathbf{x}, \mathbf{y}) \, d\mathbf{y} \\ &= \sum_{C \in \text{Ch}(S)} w_{S,C} \int C(\mathbf{x}, \mathbf{y}) \, d\mathbf{y}. \end{aligned}$$

Analogously, for a product node

$$\begin{aligned} \int P(\mathbf{x}, \mathbf{y}) d\mathbf{y} &= \int \prod_{C \in \text{Ch}(P)} C(\mathbf{x}, \mathbf{y}) d\mathbf{y} \\ &= \prod_{C \in \text{Ch}(P)} \int C(\mathbf{x}, \mathbf{y}) d\mathbf{y}. \end{aligned}$$

This ensures that marginals are pushed down to children. This can be done recursively until C is an input node L_p , in which case we marginalize \mathbf{y} according to p , which by definition should be tractable and here we assume can be done in $\mathcal{O}(1)$. We have proved the case for MAR. For EVI, we simply assign $\mathbf{y} = \emptyset$ with input nodes acting as probability density functions. Conditionals can easily be computed by an EVI or MAR followed by a second pass marginalizing the conditional variables $p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{y})}$ which are both done in linear time as we have seen here. \square

Theorem 2.2.2 (PEHARZ, GENS, PERNKOPF, *et al.*, 2016). *Let C a smooth, decomposable and deterministic PC. MAP is computable in linear time (in the number of edges of C).*

Proof. For a sum node S , we want to compute the following query

$$\max_{\mathbf{y}} S(\mathbf{y}|\mathbf{x}) = \frac{1}{S(\mathbf{x})} \max_{\mathbf{y}} S(\mathbf{y}, \mathbf{x}) = \frac{1}{S(\mathbf{x})} \max_{\mathbf{y}} \sum_{C \in \text{Ch}(S)} w_{S,C} C(\mathbf{y}, \mathbf{x}),$$

yet notice that for any assignment of \mathbf{x} and \mathbf{y} only one $C \in \text{Ch}(S)$ must have a nonnegative value by the definition of determinism, so we may replace the summation with a maximization over the children, giving

$$\max_{\mathbf{y}} S(\mathbf{y}|\mathbf{x}) = \frac{1}{S(\mathbf{x})} \max_{\mathbf{y}} \max_{C \in \text{Ch}(S)} w_{S,C} C(\mathbf{y}, \mathbf{x}) = \frac{1}{S(\mathbf{x})} \max_{C \in \text{Ch}(S)} \max_{\mathbf{y}} w_{S,C} C(\mathbf{y}, \mathbf{x}).$$

For a product node P , we compute

$$\max_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) = \frac{1}{P(\mathbf{x})} \max_{\mathbf{y}} P(\mathbf{y}, \mathbf{x}) = \frac{1}{P(\mathbf{x})} \max_{\mathbf{y}} \prod_{C \in \text{Ch}(P)} C(\mathbf{y}, \mathbf{x}) = \frac{1}{P(\mathbf{x})} \prod_{C \in \text{Ch}(P)} \max_{\mathbf{y}} C(\mathbf{y}, \mathbf{x}).$$

This is equivalent to an inductive top-down pass where we maximize instead of sum until we reach all input nodes, in which case we simply maximize the supposedly tractable functions. Once these are computed, we unroll the induction, maximizing over all values. \square

Theorem 2.3.1 (Y. CHOI, VERGARI, and BROECK, 2020). *If C is a smooth and structure decomposable probabilistic circuit with vtree \mathcal{V} , and \mathcal{L} a structure decomposable logic circuit also respecting \mathcal{V} , then $\mathbb{E}_C[\mathcal{L}]$ is polynomial time computable (in the number of edges).*

Proof. For completeness, we show the proof of this claim as stated in Y. CHOI, VERGARI, and BROECK, 2020. We assume, without loss of generality, that the layers of both C and \mathcal{L} are compatible, i.e. they both have the same number of layers and if the i -th layer of C is made out of sums (resp. products), then the i -th layer of \mathcal{L} is made out of disjunctions

(resp. conjunctions). The expectation $\mathbb{E}_C [\mathcal{L}]$ has the following form when the root of C is a product

$$\begin{aligned} \mathbb{E}_C [\mathcal{L}] &= \int C(\mathbf{x}) \mathcal{L}(\mathbf{x}) d\mathbf{x} = \int (C_p(\mathbf{x}) C_s(\mathbf{x})) (\mathcal{L}_p(\mathbf{x}) \mathcal{L}_s(\mathbf{x})) d\mathbf{x} \\ &= \int (C_p(\mathbf{x}) \mathcal{L}_p(\mathbf{x})) (C_s(\mathbf{x}) \mathcal{L}_s(\mathbf{x})) d\mathbf{x} = \int (C_p(\mathbf{x}) \mathcal{L}_p(\mathbf{x})) d\mathbf{x} \int (C_s(\mathbf{x}) \mathcal{L}_s(\mathbf{x})) d\mathbf{x} \\ &= \mathbb{E}_{C_p} [\mathcal{L}_p] \cdot \mathbb{E}_{C_s} [\mathcal{L}_s], \end{aligned}$$

where the subscript p and s indicate the prime and sub of a node. When the root is a sum

$$\begin{aligned} \mathbb{E}_C [\mathcal{L}] &= \int C(\mathbf{x}) \mathcal{L}(\mathbf{x}) d\mathbf{x} = \int \left(\sum_{C' \in \text{Ch}(C)} w_{C,C'} C'(\mathbf{x}) \right) \left(\sum_{C'' \in \text{Ch}(\mathcal{L})} C''(\mathbf{x}) \right) d\mathbf{x} \\ &= \int \sum_{C' \in \text{Ch}(C)} \sum_{C'' \in \text{Ch}(\mathcal{L})} w_{C,C'} \cdot C'(\mathbf{x}) \cdot C''(\mathbf{x}) d\mathbf{x} = \sum_{C' \in \text{Ch}(C)} \sum_{C'' \in \text{Ch}(\mathcal{L})} w_{C,C'} \int C'(\mathbf{x}) C''(\mathbf{x}) d\mathbf{x} \\ &= \sum_{C' \in \text{Ch}(C)} \sum_{C'' \in \text{Ch}(\mathcal{L})} w_{C,C'} \cdot \mathbb{E}_{C'} [C'']. \end{aligned}$$

Therefore, if expectation is tractable for input nodes, then expectation is tractable for the whole circuit. \square

References

- [AKERS 1978] S. B. AKERS. “Binary decision diagrams”. In: *IEEE Transactions on Computers* 27.6 (1978), pp. 509–516 (cit. on p. 8).
- [AMER and TODOROVIC 2012] Mohamed R. AMER and Sinisa TODOROVIC. “Sum-product networks for modeling activities with stochastic structure”. In: *2012 IEEE Conference on Computer Vision and Pattern Recognition*. 2012, pp. 1314–1321. DOI: [10.1109/CVPR.2012.6247816](https://doi.org/10.1109/CVPR.2012.6247816) (cit. on p. 21).
- [AMER and TODOROVIC 2016] Mohamed R. AMER and Sinisa TODOROVIC. “Sum product networks for activity recognition”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 38.4 (2016), pp. 800–813. DOI: [10.1109/TPAMI.2015.2465955](https://doi.org/10.1109/TPAMI.2015.2465955) (cit. on p. 21).
- [BARWISE 1982] Jon BARWISE. *Handbook of Mathematical Logic*. 1982 (cit. on p. 8).
- [BERGSTRA and BENGIO 2012] James BERGSTRA and Yoshua BENGIO. “Random search for hyper-parameter optimization”. In: *Journal of Machine Learning Research* 13.10 (2012), pp. 281–305. URL: <http://jmlr.org/papers/v13/bergstra12a.html> (cit. on p. 27).
- [BOUCHARD and TRIGGS 2004] Guillaume BOUCHARD and Bill TRIGGS. “The Tradeoff Between Generative and Discriminative Classifiers”. In: *16th IASC International Symposium on Computational Statistics (COMPSTAT '04)*. Prague, Czech Republic, Aug. 2004, pp. 721–728. URL: <https://hal.inria.fr/inria-00548546> (cit. on p. 44).
- [BOUTILIER *et al.* 1996] Craig BOUTILIER, Nir FRIEDMAN, Moises GOLDSZMIDT, and Daphne KOLLER. “Context-specific independence in bayesian networks”. In: *Proceedings of the Twelfth International Conference on Uncertainty in Artificial Intelligence*. UAI'96. Portland, OR: Morgan Kaufmann Publishers Inc., 1996, pp. 115–123. ISBN: 155860412X (cit. on p. 38).
- [BOVA 2016] Simone BOVA. “Sdds are exponentially more succinct than obdds”. In: *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*. AAAI'16. Phoenix, Arizona: AAAI Press, 2016, pp. 929–935 (cit. on p. 41).
- [BRYANT 1986] Randal E. BRYANT. “Graph-based algorithms for boolean function manipulation”. In: *IEEE Transactions on Computers* 35.8 (1986), pp. 677–691 (cit. on pp. 17, 18, 41).
- [BUEFF *et al.* 2018] Andreas BUEFF, Stefanie SPEICHERT, and Vaishak BELLE. “Tractable querying and learning in hybrid domains via sum-product networks”. In: *CoRR* abs/1807.05464 (2018). arXiv: [1807.05464](https://arxiv.org/abs/1807.05464). URL: <http://arxiv.org/abs/1807.05464> (cit. on p. 31).

- [Cory J BUTZ *et al.* 2019] Cory J BUTZ, Jhonatan S OLIVEIRA, André E dos SANTOS, and André L TEIXEIRA. “Deep convolutional sum-product networks”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 33. 2019, pp. 3248–3255 (cit. on p. 21).
- [Cory J. BUTZ, OLIVEIRA, *et al.* 2018] Cory J. BUTZ, Jhonatan S. OLIVEIRA, *et al.* “An empirical study of methods for spn learning and inference”. In: *Proceedings of the Ninth International Conference on Probabilistic Graphical Models*. Ed. by Václav KRATOCHVÍL and Milan STUDENÝ. Vol. 72. Proceedings of Machine Learning Research. PMLR, Nov. 2018, pp. 49–60. URL: <https://proceedings.mlr.press/v72/butz18a.html> (cit. on p. 31).
- [Cory J. BUTZ, SANTOS, *et al.* 2018] Cory J. BUTZ, André E. dos SANTOS, Jhonatan S. OLIVEIRA, and John STAVRINIDES. “Efficient examination of soil bacteria using probabilistic graphical models”. In: *Recent Trends and Future Technology in Applied Intelligence*. Ed. by Malek MOUHOUB, Samira SADAoui, Otmane AIT MOHAMED, and Moonis ALI. Cham: Springer International Publishing, 2018, pp. 315–326. ISBN: 978-3-319-92058-0 (cit. on p. 21).
- [CHAN and DARWICHE 2006] Hei CHAN and Adnan DARWICHE. “On the robustness of most probable explanations”. In: *Proceedings of the 22nd Conference in Uncertainty in Artificial Intelligence*. 2006 (cit. on p. 9).
- [CHENG *et al.* 2014] Wei-Chen CHENG, Stanley KOK, Hoai Vu PHAM, Hai Leong CHIEU, and Kian Ming A. CHAI. “Language modeling with sum-product networks”. In: *Fifteenth Annual Conference of the International Speech Communication Association*. 2014 (cit. on pp. 3, 21).
- [CHERNIKOVA *et al.* 2019] A. CHERNIKOVA, A. OPREA, C. NITA-ROTARU, and B. KIM. “Are self-driving cars secure? evasion attacks against deep neural networks for steering angle prediction”. In: *2019 IEEE Security and Privacy Workshops*. 2019, pp. 132–137 (cit. on p. 1).
- [A. CHOI, BROECK, *et al.* 2015] Arthur CHOI, Guy Van den BROECK, and Adnan DARWICHE. “Tractable learning for structured probability spaces: A case study in learning preference distributions”. In: *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence*. 2015, pp. 2861–2868 (cit. on p. 42).
- [A. CHOI and DARWICHE 2013] Arthur CHOI and Adnan DARWICHE. “Dynamic minimization of sentential decision diagrams”. In: *Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence*. 2013, pp. 187–194 (cit. on p. 41).
- [A. CHOI, SHEN, *et al.* 2017] Arthur CHOI, Yujia SHEN, and Adnan DARWICHE. “Tractability in structured probability spaces”. In: *Advances in Neural Information Processing Systems*. Vol. 30. 2017, pp. 3477–3485 (cit. on p. 41).

REFERENCES

- [A. CHOI, TAVABI, *et al.* 2016] Arthur CHOI, Nazgol TAVABI, and Adnan DARWICHE. “Structured features in naive Bayes classification”. In: *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*. 2016, pp. 3233–3240 (cit. on p. 41).
- [Y. CHOI, VERGARI, and BROECK 2020] YooJung CHOI, Antonio VERGARI, and Guy Van den BROECK. “Probabilistic circuits: a unifying framework for tractable probabilistic models”. In: (2020). In preparation (cit. on pp. 2, 20, 52).
- [Y. CHOI, VERGARI, and VAN DEN BROECK 2020] YooJung CHOI, Antonio VERGARI, and Guy VAN DEN BROECK. “Lecture notes: probabilistic circuits: representation and inference”. In: (Feb. 2020). URL: <http://starai.cs.ucla.edu/papers/LecNoAAAI20.pdf> (cit. on pp. 1, 10, 51).
- [CHOW and C. LIU 1968] C. CHOW and C. LIU. “Approximating discrete probability distributions with dependence trees”. In: *IEEE Transactions on Information Theory* 14.3 (1968), pp. 462–467. DOI: [10.1109/TIT.1968.1054142](https://doi.org/10.1109/TIT.1968.1054142) (cit. on pp. 37, 40).
- [CONATY *et al.* 2017] Diarmaid CONATY, Cassio Polpo de CAMPOS, and Denis Deratani MAUÁ. “Approximation complexity of maximum A posteriori inference in sum-product networks”. In: *Proceedings of the Thirty-Third Conference on Uncertainty in Artificial Intelligence*. 2017 (cit. on p. 12).
- [Fabio G. COZMAN 2000] Fabio G. COZMAN. “Credal networks”. In: *Artificial Intelligence* 120.2 (2000), pp. 199–233. ISSN: 0004-3702. DOI: [https://doi.org/10.1016/S0004-3702\(00\)00029-1](https://doi.org/10.1016/S0004-3702(00)00029-1). URL: <https://www.sciencedirect.com/science/article/pii/S0004370200000291> (cit. on p. 8).
- [DANG *et al.* 2020] Meihua DANG, Antonio VERGARI, and Guy Van den BROECK. “Strudel: learning structured-decomposable probabilistic circuits”. In: *Proceedings of the 10th International Conference on Probabilistic Graphical Models*. PGM. 2020 (cit. on pp. 3, 37, 38, 40, 41).
- [DARWICHE 1999] Adnan DARWICHE. “Compiling knowledge into decomposable negation normal form”. In: *IJCAI*. Vol. 99. 1999, pp. 284–289 (cit. on p. 17).
- [DARWICHE 2001] Adnan DARWICHE. “Decomposable negation normal form”. In: *J. ACM* 48.4 (July 2001), pp. 608–647. ISSN: 0004-5411. DOI: [10.1145/502090.502091](https://doi.org/10.1145/502090.502091). URL: <https://doi.org/10.1145/502090.502091> (cit. on p. 17).
- [DARWICHE 2003] Adnan DARWICHE. “A differential approach to inference in bayesian networks”. In: *Journal of the ACM* 50.3 (2003), pp. 280–305 (cit. on pp. 2, 7).
- [DARWICHE 2011] Adnan DARWICHE. “SDD: a new canonical representation of propositional knowledge bases”. In: *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence*. 2011, pp. 819–826 (cit. on pp. 8, 17, 41).

- [DARWICHE 2020] Adnan DARWICHE. “Three modern roles for logic in ai”. In: *Proceedings of the 39th Symposium on Principles of Database Systems (PODS)*. 2020 (cit. on p. 18).
- [DARWICHE and MARQUIS 2002] Adnan DARWICHE and Pierre MARQUIS. “A knowledge compilation map”. In: *J. Artif. Int. Res.* 17.1 (Sept. 2002), pp. 229–264. ISSN: 1076-9757 (cit. on pp. 11, 17, 18).
- [DASGUPTA and FREUND 2008] Sanjoy DASGUPTA and Yoav FREUND. “Random projection trees and low dimensional manifolds”. In: *Proceedings of the fortieth annual ACM symposium on Theory of computing*. STOC. 2008, pp. 537–546 (cit. on p. 4).
- [DE CAMPOS 2011] Cassio P. DE CAMPOS. “New complexity results for map in bayesian networks”. In: *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence - Volume Volume Three*. IJCAI’11. Barcelona, Catalonia, Spain, 2011, pp. 2100–2106. ISBN: 9781577355151 (cit. on p. 12).
- [R. DECHTER 1998] R. DECHTER. “Bucket elimination: a unifying framework for probabilistic inference”. In: *Learning in Graphical Models*. Ed. by Michael I. JORDAN. Dordrecht: Springer Netherlands, 1998, pp. 75–104. ISBN: 978-94-011-5014-9. DOI: [10.1007/978-94-011-5014-9_4](https://doi.org/10.1007/978-94-011-5014-9_4). URL: https://doi.org/10.1007/978-94-011-5014-9_4 (cit. on p. 2).
- [Rina DECHTER and MATEESCU 2007] Rina DECHTER and Robert MATEESCU. “And/or search spaces for graphical models”. In: *Artificial Intelligence* 171.2 (2007), pp. 73–106. ISSN: 0004-3702. DOI: <https://doi.org/10.1016/j.artint.2006.11.003>. URL: <https://www.sciencedirect.com/science/article/pii/S000437020600138X> (cit. on p. 2).
- [DEMPSTER *et al.* 1977] A. P. DEMPSTER, N. M. LAIRD, and D. B. RUBIN. “Maximum likelihood from incomplete data via the em algorithm”. In: *Journal of the Royal Statistical Society. Series B (Methodological)* 39.1 (1977), pp. 1–38. ISSN: 00359246. URL: <http://www.jstor.org/stable/2984875> (cit. on p. 43).
- [DENNIS and VENTURA 2012] Aaron DENNIS and Dan VENTURA. “Learning the architecture of sum-product networks using clustering on variables”. In: *Advances in Neural Information Processing Systems* 25. NIPS, 2012, pp. 2033–2041 (cit. on p. 42).
- [DENNIS and VENTURA 2015] Aaron DENNIS and Dan VENTURA. “Greedy structure search for sum-product networks”. In: *Proceedings of the 24th International Conference on Artificial Intelligence*. 2015, pp. 932–938 (cit. on p. 9).
- [DENNIS and VENTURA 2017] Aaron DENNIS and Dan VENTURA. “Autoencoder-enhanced sum-product networks”. In: *2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA)*. 2017, pp. 1041–1044. DOI: [10.1109/ICMLA.2017.00-13](https://doi.org/10.1109/ICMLA.2017.00-13) (cit. on p. 21).

REFERENCES

- [DI MAURO *et al.* 2017] Nicola DI MAURO, Floriana ESPOSITO, Fabrizio G. VENTOLA, and Antonio VERGARI. “Alternative variable splitting methods to learn sum-product networks”. In: *AI*IA 2017 Advances in Artificial Intelligence*. Ed. by Floriana ESPOSITO, Roberto BASILI, Stefano FERILLI, and Francesca A. LISI. Cham: Springer International Publishing, 2017, pp. 334–346. ISBN: 978-3-319-70169-1 (cit. on pp. 3, 31).
- [EDMONDS 1965] Jack EDMONDS. “Paths, trees, and flowers”. In: *Canadian Journal of Mathematics* 17 (1965), pp. 449–467. DOI: [10.4153/CJM-1965-045-4](https://doi.org/10.4153/CJM-1965-045-4) (cit. on pp. 33, 35).
- [EÉN and SÖRENSON 2006] Niklas EÉN and Niklas SÖRENSON. “Translating pseudo-boolean constraints into SAT”. In: *Journal on Satisfiability, Boolean Modeling and Computation* (2006) (cit. on p. 41).
- [FELDMANN and FOSCHINI 2015] Andreas Emil FELDMANN and Luca FOSCHINI. “Balanced partitions of trees and applications”. In: *Algorithmica* 71.2 (Feb. 2015), pp. 354–376. ISSN: 0178-4617. DOI: [10.1007/s00453-013-9802-3](https://doi.org/10.1007/s00453-013-9802-3). URL: <https://doi.org/10.1007/s00453-013-9802-3> (cit. on p. 30).
- [FREUND *et al.* 2008] Yoav FREUND, Sanjoy DASGUPTA, Mayank KABRA, and Nakul VERMA. “Learning the structure of manifolds using random projections”. In: *Advances in Neural Information Processing Systems*. Vol. 20. NeurIPS. 2008 (cit. on p. 4).
- [A. FRIESEN and P. DOMINGOS 2016] Abram FRIESEN and Pedro DOMINGOS. “The sum-product theorem: a foundation for learning tractable models”. In: *Proceedings of The 33rd International Conference on Machine Learning*. Ed. by Maria Florina BALCAN and Kilian Q. WEINBERGER. Vol. 48. Proceedings of Machine Learning Research. New York, New York, USA: PMLR, June 2016, pp. 1909–1918. URL: <https://proceedings.mlr.press/v48/friesen16.html> (cit. on p. 8).
- [Abram L FRIESEN and P. DOMINGOS 2017] Abram L FRIESEN and Pedro DOMINGOS. “Unifying sum-product networks and submodular fields”. In: *Proceedings of the Workshop on Principled Approaches to Deep Learning at ICML*. 2017 (cit. on p. 21).
- [Abram L FRIESEN and P. M. DOMINGOS 2018] Abram L FRIESEN and Pedro M DOMINGOS. “Submodular field grammars: representation, inference, and application to image parsing”. In: *Advances in Neural Information Processing Systems*. Ed. by S. BENGIO *et al.* Vol. 31. Curran Associates, Inc., 2018. URL: <https://proceedings.neurips.cc/paper/2018/file/c5866e93cab1776890fe343c9e7063fb-Paper.pdf> (cit. on p. 21).
- [Abram L. FRIESEN and P. DOMINGOS 2015] Abram L. FRIESEN and Pedro DOMINGOS. “Recursive decomposition for nonconvex optimization”. In: *Proceedings of the 24th International Conference on Artificial Intelligence*. IJCAI’15. Buenos Aires, Argentina: AAAI Press, 2015, pp. 253–259. ISBN: 9781577357384 (cit. on pp. 8, 21).

- [GAREY and JOHNSON 1990] Michael R. GAREY and David S. JOHNSON. *Computers and Intractability; A Guide to the Theory of NP-Completeness*. USA: W. H. Freeman Co., 1990. ISBN: 0716710455 (cit. on p. 32).
- [R. GEH and D. MAUÁ 2019] Renato GEH and Denis MAUÁ. “End-to-end imitation learning of lane following policies using sum-product networks”. In: *Anais do XVI Encontro Nacional de Inteligência Artificial e Computacional*. Salvador: SBC, 2019, pp. 297–308. DOI: [10.5753/eniac.2019.9292](https://doi.org/10.5753/eniac.2019.9292). URL: <https://sol.sbc.org.br/index.php/eniac/article/view/9292> (cit. on p. 21).
- [R. L. GEH and Denis Deratani MAUÁ 2021] Renato Lui GEH and Denis Deratani MAUÁ. “Learning probabilistic sentential decision diagrams under logic constraints by sampling and averaging”. In: *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*. Ed. by Cassio de CAMPOS and Marloes H. MAATHUIS. Vol. 161. Proceedings of Machine Learning Research. PMLR, July 2021, pp. 2039–2049. URL: <https://proceedings.mlr.press/v161/geh21a.html> (cit. on pp. 3, 4).
- [GENS and P. DOMINGOS 2012] Robert GENS and Pedro DOMINGOS. “Discriminative learning of sum-product networks”. In: *Advances in Neural Information Processing Systems 25*. NIPS, 2012, pp. 3239–3247 (cit. on pp. 2, 21).
- [GENS and P. DOMINGOS 2013] Robert GENS and Pedro DOMINGOS. “Learning the structure of sum-product networks”. In: *Proceedings of the 30th International Conference on Machine Learning*. ICML. 2013, pp. 873–880 (cit. on pp. 3, 24, 25).
- [GOGIC *et al.* 1995] Goran GOGIC, Henry KAUTZ, Christos PAPADIMITRIOU, and Bart SELMAN. “The comparative linguistics of knowledge representation”. In: *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 1*. IJCAI’95. Montreal, Quebec, Canada, 1995, pp. 862–869. ISBN: 1558603638 (cit. on p. 18).
- [GOODFELLOW *et al.* 2014] Ian GOODFELLOW *et al.* “Generative adversarial nets”. In: *Advances in Neural Information Processing Systems 27*. Ed. by Z. GHAHRAMANI, M. WELING, C. CORTES, N. D. LAWRENCE, and K. Q. WEINBERGER. Curran Associates, Inc., 2014, pp. 2672–2680. URL: <http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf> (cit. on p. 2).
- [GRITSENKO *et al.* 2019] Alexey A. GRITSENKO, Jasper SNOEK, and Tim SALIMANS. “On the relationship between normalising flows and variational- and denoising autoencoders”. In: *Deep Generative Models for Highly Structured Data, ICLR 2019 Workshop, New Orleans, Louisiana, United States, May 6, 2019*. OpenReview.net, 2019. URL: https://openreview.net/forum?id=HklKEUY%5C_E (cit. on p. 2).
- [GUO *et al.* 2017] Chuan GUO, Geoff PLEISS, Yu SUN, and Kilian Q WEINBERGER. “On calibration of modern neural networks”. In: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org. 2017, pp. 1321–1330 (cit. on p. 1).

REFERENCES

- [HERRMANN and BARROS 2013] Ricardo G. HERRMANN and Leliane N. de BARROS. “Algebraic sentential decision diagrams in symbolic probabilistic planning”. In: *2013 Brazilian Conference on Intelligent Systems*. 2013, pp. 175–181. DOI: [10.1109/BRACIS.2013.37](https://doi.org/10.1109/BRACIS.2013.37) (cit. on p. 41).
- [JAEGER 2004] Manfred JAEGER. “Probabilistic decision graphs-combining verification and ai techniques for probabilistic inference”. In: *Int. J. Uncertain. Fuzziness Knowl.-Based Syst.* 12.1 supp (Jan. 2004), pp. 19–42. ISSN: 0218-4885. DOI: [10.1142/S0218488504002564](https://doi.org/10.1142/S0218488504002564). URL: <https://doi.org/10.1142/S0218488504002564> (cit. on p. 2).
- [JAINI, GHOSE, *et al.* 2018] Priyank JAINI, Amur GHOSE, and Pascal POUPART. “Prometheus: directly learning acyclic directed graph structures for sum-product networks”. In: *International Conference on Probabilistic Graphical Models*. PGM. 2018, pp. 181–192 (cit. on pp. 3, 27, 28, 30).
- [JAINI, POUPART, *et al.* 2018] Priyank JAINI, Pascal POUPART, and Yaoliang YU. “Deep homogeneous mixture models: representation, separation, and approximation”. In: *Advances in Neural Information Processing Systems*. Ed. by S. BENGIO *et al.* Vol. 31. Curran Associates, Inc., 2018. URL: <https://proceedings.neurips.cc/paper/2018/file/c5f5c23be1b71adb51ea9dc8e9d444a8-Paper.pdf> (cit. on p. 49).
- [KARYPIS and KUMAR 1998] George KARYPIS and Vipin KUMAR. “A fast and high quality multilevel scheme for partitioning irregular graphs”. In: *SIAM Journal on Scientific Computing* 20.1 (1998), pp. 359–392. DOI: [10.1137/S1064827595287997](https://doi.org/10.1137/S1064827595287997) (cit. on pp. 32, 35).
- [KHOSRAVI *et al.* 2019] Pasha KHOSRAVI, Yitao LIANG, YooJung CHOI, and Guy VAN DEN BROECK. “What to expect of classifiers? reasoning about logistic regression with missing features”. In: *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*. International Joint Conferences on Artificial Intelligence Organization, July 2019, pp. 2716–2724. DOI: [10.24963/ijcai.2019/377](https://doi.org/10.24963/ijcai.2019/377). URL: <https://doi.org/10.24963/ijcai.2019/377> (cit. on p. 2).
- [KINGMA and WELLING 2014] Diederik P. KINGMA and Max WELLING. “Auto-encoding variational bayes”. In: *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*. Ed. by Yoshua BENGIO and Yann LECUN. 2014. URL: <http://arxiv.org/abs/1312.6114> (cit. on p. 2).
- [KISA *et al.* 2014] Doga KISA, Guy Van den BROECK, Arthur CHOI, and Adnan DARWICHE. “Probabilistic sentential decision diagrams”. In: *Knowledge Representation and Reasoning Conference* (2014) (cit. on pp. 2, 19).
- [KOLLER and FRIEDMAN 2009] Daphne KOLLER and Nir FRIEDMAN. *Probabilistic Graphical Models: Principles and Techniques - Adaptive Computation and Machine Learning*. The MIT Press, 2009. ISBN: 0262013193 (cit. on p. 2).

- [KOLMOGOROV 2009] Vladimir KOLMOGOROV. “Blossom v: a new implementation of a minimum cost perfect matching algorithm”. In: *Mathematical Programming Computation* 1.1 (July 2009), pp. 43–67. ISSN: 1867-2957. DOI: [10.1007/s12532-009-0002-8](https://doi.org/10.1007/s12532-009-0002-8). URL: <https://doi.org/10.1007/s12532-009-0002-8> (cit. on pp. 33, 35).
- [LIANG, BEKKER, *et al.* 2017] Yitao LIANG, Jessa BEKKER, and Guy Van den BROECK. “Learning the structure of probabilistic sentential decision diagrams”. In: *Proceedings of the Thirty-Third Conference on Uncertainty in Artificial Intelligence*. 2017 (cit. on pp. 3, 31, 32, 34, 36).
- [LIANG and VAN DEN BROECK 2019] Yitao LIANG and Guy VAN DEN BROECK. “Learning logistic circuits”. In: *Proceedings of the AAAI Conference on Artificial Intelligence* 33.01 (July 2019), pp. 4277–4286. DOI: [10.1609/aaai.v33i01.33014277](https://doi.org/10.1609/aaai.v33i01.33014277). URL: <https://ojs.aaai.org/index.php/AAAI/article/view/4336> (cit. on p. 8).
- [LIPPE and GAVVES 2021] Phillip LIPPE and Efstratios GAVVES. “Categorical normalizing flows via continuous transformations”. In: *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL: <https://openreview.net/forum?id=-GLNZeVDuik> (cit. on p. 2).
- [Y. LIU and LUO 2019] Yang LIU and Tiejian LUO. “The optimization of sum-product network structure learning”. In: *Journal of Visual Communication and Image Representation* 60 (2019), pp. 391–397. ISSN: 1047-3203. DOI: <https://doi.org/10.1016/j.jvcir.2019.02.012>. URL: <https://www.sciencedirect.com/science/article/pii/S1047320319300653> (cit. on p. 31).
- [LLERENA and DERATANI MAUÁ 2017] Julissa Villanueva LLERENA and Denis DERATANI MAUÁ. “On using sum-product networks for multi-label classification”. In: *2017 Brazilian Conference on Intelligent Systems (BRACIS)*. 2017, pp. 25–30. DOI: [10.1109/BRACIS.2017.34](https://doi.org/10.1109/BRACIS.2017.34) (cit. on p. 21).
- [LOMUSCIO and PAQUET 2015] Alessio LOMUSCIO and Hugo PAQUET. “Verification of multi-agent systems via sdd-based model checking”. In: *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*. AAMAS ’15. Istanbul, Turkey: International Foundation for Autonomous Agents and Multiagent Systems, 2015, pp. 1713–1714. ISBN: 9781450334136 (cit. on p. 41).
- [LOWD and ROOSHENAS 2013] Daniel LOWD and Amirmohammad ROOSHENAS. “Learning markov networks with arithmetic circuits”. In: *Proceedings of the Sixteenth International Conference on Artificial Intelligence and Statistics*. Ed. by Carlos M. CARVALHO and Pradeep RAVIKUMAR. Vol. 31. Proceedings of Machine Learning Research. Scottsdale, Arizona, USA: PMLR, Apr. 2013, pp. 406–414. URL: <https://proceedings.mlr.press/v31/lowd13a.html> (cit. on p. 27).
- [MARTENS and MEDABALIMI 2014] James MARTENS and Venkatesh MEDABALIMI. “On the expressive efficiency of sum product networks”. In: *CoRR* abs/1411.7717 (2014). arXiv: [1411.7717](https://arxiv.org/abs/1411.7717). URL: <http://arxiv.org/abs/1411.7717> (cit. on pp. 7, 26).

- [Denis D. MAUÁ *et al.* 2017] Denis D. MAUÁ, Fabio G. COZMAN, Diarmaid CONATY, and Cassio P. CAMPOS. “Credal sum-product networks”. In: *Proceedings of the Tenth International Symposium on Imprecise Probability: Theories and Applications*. Ed. by Alessandro ANTONUCCI, Giorgio CORANI, Inés COUSO, and Sébastien DESTERCKE. Vol. 62. Proceedings of Machine Learning Research. PMLR, Oct. 2017, pp. 205–216. URL: <https://proceedings.mlr.press/v62/mau%C3%A117a.html> (cit. on p. 8).
- [MAURO *et al.* 2021] Nicola Di MAURO, Gennaro GALA, Marco IANNOTTA, and Teresa M. A. BASILE. “Random probabilistic circuits”. In: *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*. 2021 (cit. on pp. 3, 42).
- [MEI *et al.* 2018] Jun MEI, Yong JIANG, and Kewei TU. “Maximum a posteriori inference in sum-product networks”. In: *AAAI Conference on Artificial Intelligence*. 2018 (cit. on p. 12).
- [MELIBARI, POUPART, and DOSHI 2016] Mazen MELIBARI, Pascal POUPART, and Prashant DOSHI. “Sum-product-max networks for tractable decision making”. In: *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*. IJCAI’16. New York, New York, USA: AAAI Press, 2016, pp. 1846–1852. ISBN: 9781577357704 (cit. on p. 8).
- [MELIBARI, POUPART, DOSHI, and TRIMPONIAS 2016] Mazen MELIBARI, Pascal POUPART, Prashant DOSHI, and George TRIMPONIAS. “Dynamic sum product networks for tractable inference on sequence data”. In: *Probabilistic Graphical Models*. Vol. 52. JMLR Workshop and Conference Proceedings. JMLR.org, 2016, pp. 345–355 (cit. on p. 21).
- [MOLINA, NATARAJAN, *et al.* 2017] Alejandro MOLINA, Sriraam NATARAJAN, and Kristian KERSTING. “Poisson sum-product networks: a deep architecture for tractable multivariate poisson distributions”. In: *Proceedings of the AAAI Conference on Artificial Intelligence* 31.1 (Feb. 2017). URL: <https://ojs.aaai.org/index.php/AAAI/article/view/10844> (cit. on p. 31).
- [MOLINA, VERGARI, *et al.* 2018] Alejandro MOLINA, Antonio VERGARI, *et al.* “Mixed sum-product networks: a deep architecture for hybrid domains”. In: (2018). URL: <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16865> (cit. on p. 31).
- [NATH and P. DOMINGOS 2015] Aniruddh NATH and Pedro DOMINGOS. “Learning relational sum-product networks”. In: *Proceedings of the AAAI Conference on Artificial Intelligence* 29.1 (Feb. 2015). URL: <https://ojs.aaai.org/index.php/AAAI/article/view/9538> (cit. on p. 31).
- [NATH and P. M. DOMINGOS 2016] Aniruddh NATH and Pedro M DOMINGOS. “Learning tractable probabilistic models for fault localization”. In: *Thirtieth AAAI Conference on Artificial Intelligence*. 2016 (cit. on pp. 3, 21).

- [NISHINO *et al.* 2016] Masaaki NISHINO, Norihito YASUDA, Shin-ichi MINATO, and Masaaki NAGATA. “Zero-suppressed sentential decision diagrams”. In: *AAAI Conference on Artificial Intelligence*. 2016 (cit. on p. 41).
- [NOURANI *et al.* 2020] Mahsan NOURANI *et al.* *Don’t Explain without Verifying Veracity: An Evaluation of Explainable AI with Video Activity Recognition*. 2020. arXiv: 2005.02335 [cs.HC] (cit. on p. 21).
- [OLASCOAGA *et al.* 2019] Laura Isabel Galindez OLASCOAGA, Wannes MEERT, Nimish SHAH, Marian VERHELST, and Guy VAN DEN BROECK. “Towards hardware-aware tractable learning of probabilistic models”. In: *NeurIPS*. 2019, pp. 13726–13736 (cit. on p. 21).
- [OVADIA *et al.* 2019] Yaniv OVADIA *et al.* “Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift”. In: *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. 2019 (cit. on p. 1).
- [OZTOK and DARWICHE 2015] Umut OZTOK and Adnan DARWICHE. “A top-down compiler for sentential decision diagrams”. In: *Proceedings of the 24th International Conference on Artificial Intelligence*. 2015, pp. 3141–3148 (cit. on p. 41).
- [PAPADIMITRIOU 1994] C.H. PAPADIMITRIOU. *Computational Complexity*. Theoretical computer science. Addison-Wesley, 1994. ISBN: 9780201530827 (cit. on p. 18).
- [PAPAMAKARIOS *et al.* 2021] George PAPAMAKARIOS, Eric NALISNICK, Danilo Jimenez REZENDE, Shakir MOHAMED, and Balaji LAKSHMINARAYANAN. “Normalizing flows for probabilistic modeling and inference”. In: *Journal of Machine Learning Research* 22.57 (2021), pp. 1–64. URL: <http://jmlr.org/papers/v22/19-1028.html> (cit. on p. 2).
- [PEHARZ, GENS, and P. DOMINGOS 2014] Robert PEHARZ, Robert GENS, and Pedro DOMINGOS. “Learning selective sum-product networks”. In: *Workshop on Learning Tractable Probabilistic Models*. 2014 (cit. on p. 9).
- [PEHARZ, GENS, PERNKOPF, *et al.* 2016] Robert PEHARZ, Robert GENS, Franz PERNKOPF, and Pedro DOMINGOS. “On the latent variable interpretation in sum-product networks”. In: *IEEE transactions on pattern analysis and machine intelligence* 39.10 (2016), pp. 2030–2044 (cit. on pp. 6, 12, 43, 52).
- [PEHARZ, KAPPELLER, *et al.* 2014] Robert PEHARZ, Georg KAPPELLER, Pejman MOWLAEE, and Franz PERNKOPF. “Modeling speech with sum-product networks: application to bandwidth extension”. In: *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. 2014, pp. 3699–3703 (cit. on p. 21).

REFERENCES

- [PEHARZ, LANG, *et al.* 2020] Robert PEHARZ, Steven LANG, *et al.* “Einsum networks: fast and scalable learning of tractable probabilistic circuits”. In: *Proceedings of the 37th International Conference on Machine Learning*. Ed. by Hal Daumé III and Aarti SINGH. Vol. 119. Proceedings of Machine Learning Research. PMLR, July 2020, pp. 7563–7574. URL: <https://proceedings.mlr.press/v119/peharz20a.html> (cit. on pp. 3, 8, 21).
- [PEHARZ, TSCHIATSCHKE, *et al.* 2015] Robert PEHARZ, Sebastian TSCHIATSCHKE, Franz PERKOPF, and Pedro DOMINGOS. “On theoretical properties of sum-product networks”. In: *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*. 2015, pp. 744–752 (cit. on pp. 7, 11).
- [PEHARZ, VERGARI, *et al.* 2020] Robert PEHARZ, Antonio VERGARI, *et al.* “Random sum-product networks: a simple and effective approach to probabilistic deep learning”. In: *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*. Ed. by Ryan P. ADAMS and Vibhav GOGATE. Vol. 115. Proceedings of Machine Learning Research. PMLR, July 2020, pp. 334–344. URL: <https://proceedings.mlr.press/v115/peharz20a.html> (cit. on pp. 3, 21, 42–44).
- [PEVNÝ *et al.* 2020] Tomáš PEVNÝ, Václav SMÍDL, Martin TRAPP, Ondřej POLÁČEK, and Tomáš OBERHUBER. “Sum-product-transform networks: exploiting symmetries using invertible transformations”. In: *Proceedings of the 10th International Conference on Probabilistic Graphical Models*. Ed. by Manfred JAEGER and Thomas Dyhre NIELSEN. Vol. 138. Proceedings of Machine Learning Research. PMLR, Sept. 2020, pp. 341–352. URL: <https://proceedings.mlr.press/v138/pevny20a.html> (cit. on p. 8).
- [POON and P. DOMINGOS 2011] Hoifung POON and Pedro DOMINGOS. “Sum-product networks: a new deep architecture”. In: *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*. 2011, pp. 337–346 (cit. on pp. 2, 3, 6, 10, 21, 51).
- [PRONOBIS *et al.* 2017] Andrzej PRONOBIS, Francesco RICCIO, and Rajesh PN RAO. “Deep spatial affordance hierarchy: spatial knowledge representation for planning in large-scale environments”. In: *ICAPS 2017 Workshop on Planning and Robotics*. 2017 (cit. on p. 21).
- [RAHMAN and GOGATE 2016] Tahrima RAHMAN and Vibhav GOGATE. “Merging strategies for sum-product networks: from trees to graphs”. In: *Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence*. UAI’16. Jersey City, New Jersey, USA: AUAI Press, 2016, pp. 617–626. ISBN: 9780996643115 (cit. on p. 31).
- [RAHMAN, KOTHALKAR, *et al.* 2014] Tahrima RAHMAN, Prasanna KOTHALKAR, and Vibhav GOGATE. “Cutset networks: a simple, tractable, and scalable approach for improving the accuracy of chow-liu trees”. In: *Proceedings of the 2014th European Conference on Machine Learning and Knowledge Discovery in Databases*. 2014, pp. 630–645 (cit. on p. 2).

- [RASHWAN *et al.* 2018] Abdullah RASHWAN, Pascal POUPART, and Chen ZHITANG. “Discriminative training of sum-product networks by extended baum-welch”. In: *Proceedings of the Ninth International Conference on Probabilistic Graphical Models*. Vol. 72. Proceedings of Machine Learning Research. 2018, pp. 356–367 (cit. on p. 2).
- [RATAJCZAK *et al.* 2014] Martin RATAJCZAK, Sebastian TSCHIATSCHEK, and Franz PERNKOPF. “Sum-product networks for structured prediction: context-specific deep conditional random fields”. In: *International Conference on Machine Learning (ICML) Workshop on Learning Tractable Probabilistic Models Workshop*. 2014 (cit. on p. 21).
- [RATAJCZAK *et al.* 2018] Martin RATAJCZAK, Sebastian TSCHIATSCHEK, and Franz PERNKOPF. “Sum-product networks for sequence labeling”. In: *arXiv preprint arXiv:1807.02324* (2018) (cit. on p. 21).
- [RATHKE *et al.* 2017] Fabian RATHKE, Mattia DESANA, and Christoph SCHNÖRR. “Locally adaptive probabilistic models for global segmentation of pathological oct scans”. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer. 2017, pp. 177–184 (cit. on p. 21).
- [REZENDE and MOHAMED 2015] Danilo REZENDE and Shakir MOHAMED. “Variational inference with normalizing flows”. In: *Proceedings of the 32nd International Conference on Machine Learning*. Ed. by Francis BACH and David BLEI. Vol. 37. Proceedings of Machine Learning Research. Lille, France: PMLR, July 2015, pp. 1530–1538. URL: <https://proceedings.mlr.press/v37/rezende15.html> (cit. on p. 2).
- [ROLFE 2017] Jason Tyler ROLFE. “Discrete variational autoencoders”. In: *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL: <https://openreview.net/forum?id=ryMxXPfex> (cit. on p. 2).
- [ROOSHENAS and LOWD 2014] Amirmohammad ROOSHENAS and Daniel LOWD. “Learning sum-product networks with direct and indirect variable interactions”. In: *Proceedings of the 31st International Conference on Machine Learning*. Ed. by Eric P. XING and Tony JEBA. Vol. 32. Proceedings of Machine Learning Research 1. Beijing, China: PMLR, June 2014, pp. 710–718. URL: <https://proceedings.mlr.press/v32/rooshenas14.html> (cit. on pp. 26, 27).
- [ROOSHENAS and LOWD 2016] Amirmohammad ROOSHENAS and Daniel LOWD. “Discriminative structure learning of arithmetic circuits”. In: *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*. Ed. by Arthur GRETTON and Christian C. ROBERT. Vol. 51. Proceedings of Machine Learning Research. Cadiz, Spain: PMLR, Sept. 2016, pp. 1506–1514. URL: <https://proceedings.mlr.press/v51/rooshenas16.html> (cit. on p. 2).

REFERENCES

- [SAAD *et al.* 2021] Feras A. SAAD, Martin C. RINARD, and Vikash K. MANSINGHKA. “SPPL: probabilistic programming with fast exact symbolic inference”. In: *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*. Virtual, Canada: Association for Computing Machinery, 2021, pp. 804–819. ISBN: 9781450383912. DOI: [10.1145/3453483.3454078](https://doi.org/10.1145/3453483.3454078) (cit. on p. 21).
- [SGUERRA and F. G. COZMAN 2016] B. M. SGUERRA and F. G. COZMAN. “Image classification using sum-product networks for autonomous flight of micro aerial vehicles”. In: *2016 5th Brazilian Conference on Intelligent Systems (BRACIS)*. 2016, pp. 139–144 (cit. on p. 21).
- [SHAH, ISABEL GALINDEZ OLASCOAGA, *et al.* 2019] Nimish SHAH, Laura Isabel GALINDEZ OLASCOAGA, Wannes MEERT, and Marian VERHELST. “Problp: a framework for low-precision probabilistic inference”. In: *DAC 2019*. ACM, 2019, p. 190. DOI: [10.1145/3316781.3317885](https://doi.org/10.1145/3316781.3317885). URL: <https://doi.org/10.1145/3316781.3317885> (cit. on p. 21).
- [SHAH, OLASCOAGA, MEERT, *et al.* 2020] Nimish SHAH, Laura Isabel Galindez OLASCOAGA, Wannes MEERT, and Marian VERHELST. “Acceleration of probabilistic reasoning through custom processor architecture”. In: *DATE*. IEEE, 2020, pp. 322–325 (cit. on p. 21).
- [SHAH, OLASCOAGA, S. ZHAO, *et al.* 2021] Nimish SHAH, Laura Isabel Galindez OLASCOAGA, Shirui ZHAO, Wannes MEERT, and Marian VERHELST. “Piu: a 248gop-s/w stream-based processor for irregular probabilistic inference networks using precision-scalable posit arithmetic in 28nm”. In: *2021 IEEE International Solid-State Circuits Conference (ISSCC)*. Vol. 64. 2021, pp. 150–152. DOI: [10.1109/ISSCC42613.2021.9366061](https://doi.org/10.1109/ISSCC42613.2021.9366061) (cit. on p. 21).
- [SHAO *et al.* 2020] Xiaoting SHAO *et al.* “Conditional sum-product networks: imposing structure on deep probabilistic architectures”. In: *Proceedings of the 10th International Conference on Probabilistic Graphical Models*. Ed. by Manfred JAEGER and Thomas Dyhre NIELSEN. Vol. 138. Proceedings of Machine Learning Research. PMLR, Sept. 2020, pp. 401–412. URL: <https://proceedings.mlr.press/v138/shao20a.html> (cit. on p. 3).
- [SHARIR and SHASHUA 2018] Or SHARIR and Amnon SHASHUA. “Sum-product-quotient networks”. In: *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*. Ed. by Amos STORKEY and Fernando PEREZ-CRUZ. Vol. 84. Proceedings of Machine Learning Research. PMLR, Sept. 2018, pp. 529–537. URL: <https://proceedings.mlr.press/v84/sharir18a.html> (cit. on p. 8).
- [SHEN *et al.* 2017] Yujia SHEN, Arthur CHOI, and Adnan DARWICHE. “A tractable probabilistic model for subset selection”. In: *Proceedings of the Thirty-Third Conference on Uncertainty in Artificial Intelligence*. 2017 (cit. on p. 42).

- [SINZ 2005] Carsten SINZ. “Towards an optimal CNF encoding of boolean cardinality constraints”. In: *Proceedings of the 11th International Conference on Principles and Practice of Constraint Programming*. 2005 (cit. on p. 41).
- [SOMMER *et al.* 2018] Lukas SOMMER *et al.* “Automatic mapping of the sum-product network inference problem to fpga-based accelerators”. In: *ICCD*. IEEE Computer Society, 2018, pp. 350–357 (cit. on p. 21).
- [STUHLMULLER and GOODMAN 2012] Andreas STUHLMULLER and Noah D GOODMAN. “A dynamic programming algorithm for inference in recursive probabilistic programs”. In: *Workshop of Statistical and Relational AI (StarAI)* (2012) (cit. on p. 21).
- [SU *et al.* 2019] Jiawei SU, Danilo Vasconcellos VARGAS, and Kouichi SAKURAI. “One pixel attack for fooling deep neural networks”. In: *IEEE Transactions on Evolutionary Computation* 23.5 (2019), pp. 828–841 (cit. on p. 1).
- [SZEGEDY *et al.* 2013] Christian SZEGEDY *et al.* “Intriguing properties of neural networks”. In: *arXiv preprint arXiv:1312.6199* (2013) (cit. on p. 1).
- [VAHDAT, ANDRIYASH, *et al.* 2018] Arash VAHDAT, Evgeny ANDRIYASH, and William MACREADY. “Dvae#: discrete variational autoencoders with relaxed boltzmann priors”. In: *Advances in Neural Information Processing Systems*. Ed. by S. BENGIO *et al.* Vol. 31. Curran Associates, Inc., 2018. URL: <https://proceedings.neurips.cc/paper/2018/file/9f53d83ec0691550f7d2507d57f4f5a2-Paper.pdf> (cit. on p. 2).
- [VAHDAT, MACREADY, *et al.* 2018] Arash VAHDAT, William G. MACREADY, Zhengbing BIAN, Amir KHOSHAMAN, and Evgeny ANDRIYASH. “DVAE++: discrete variational autoencoders with overlapping transformations”. In: *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*. Ed. by Jennifer G. DY and Andreas KRAUSE. Vol. 80. Proceedings of Machine Learning Research. PMLR, 2018, pp. 5042–5051. URL: <http://proceedings.mlr.press/v80/vahdat18a.html> (cit. on p. 2).
- [VERGARI, Y. CHOI, *et al.* 2021] 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”. In: *CoRR* abs/2102.06137 (2021). arXiv: [2102.06137](https://arxiv.org/abs/2102.06137) (cit. on pp. 8, 10, 16, 51).
- [VERGARI, MAURO, *et al.* 2015] Antonio VERGARI, Nicola Di MAURO, and Floriana ESPOSITO. “Simplifying, regularizing and strengthening sum-product network structure learning”. In: *ECML/PKDD*. 2015 (cit. on pp. 3, 31).
- [VLASSELAER, BROECK, *et al.* 2015] Jonas VLASSELAER, Guy Van den BROECK, Angelika KIMMIG, Wannes MEERT, and Luc De RAEDT. “Anytime inference in probabilistic logic programs with tp-compilation”. In: (2015). URL: <https://www.aaai.org/ocs/index.php/IJCAI/IJCAI15/paper/view/11005> (cit. on p. 41).

REFERENCES

- [VLASSELAER, RENKENS, *et al.* 2014] Jonas VLASSELAER, Joris RENKENS, Guy VAN DEN BROECK, and Luc DE RAEDT. “Compiling probabilistic logic programs into sentential decision diagrams”. eng. In: 2014, pp. 1–10. URL: [https://lirias.kuleuven.be/retrieve/275254\\$Dplp_sdd_paper.pdf%20\[freely%20available\]](https://lirias.kuleuven.be/retrieve/275254$Dplp_sdd_paper.pdf%20[freely%20available]) (cit. on p. 41).
- [WACHTER and HAENNI 2006] Michael WACHTER and Rolf HAENNI. “Propositional dags: a new graph-based language for representing boolean functions”. In: *Proceedings of the Tenth International Conference on Principles of Knowledge Representation and Reasoning*. KR’06. Lake District, UK, 2006, pp. 277–285. ISBN: 9781577352716 (cit. on p. 17).
- [J. WANG and G. WANG 2018] Jinghua WANG and Gang WANG. “Hierarchical spatial sum-product networks for action recognition in still images”. In: *IEEE Transactions on Circuits and Systems for Video Technology* 28.1 (2018), pp. 90–100. DOI: [10.1109/TCSVT.2016.2586853](https://doi.org/10.1109/TCSVT.2016.2586853) (cit. on p. 21).
- [WEI *et al.* 2018] Wenqi WEI *et al.* “Adversarial examples in deep learning: characterization and divergence”. In: *arXiv preprint arXiv:1807.00051* (2018) (cit. on p. 1).
- [YU 2020] Ronald YU. *A Tutorial on VAEs: From Bayes’ Rule to Lossless Compression*. 2020. arXiv: [2006.10273](https://arxiv.org/abs/2006.10273) [cs.LG] (cit. on p. 2).
- [YUAN *et al.* 2016] Zehuan YUAN *et al.* “Modeling spatial layout for scene image understanding via a novel multiscale sum-product network”. In: *Expert Syst. Appl.* 63.C (Nov. 2016), pp. 231–240. ISSN: 0957-4174. DOI: [10.1016/j.eswa.2016.07.015](https://doi.org/10.1016/j.eswa.2016.07.015). URL: <https://doi.org/10.1016/j.eswa.2016.07.015> (cit. on p. 21).
- [ZHANG *et al.* 2021] Honghua ZHANG, Brendan JUBA, and Guy VAN DEN BROECK. “Probabilistic generating circuits”. In: *Proceedings of the 38th International Conference on Machine Learning*. Ed. by Marina MEILA and Tong ZHANG. Vol. 139. Proceedings of Machine Learning Research. PMLR, July 2021, pp. 12447–12457. URL: <https://proceedings.mlr.press/v139/zhang21i.html> (cit. on p. 8).
- [H. ZHAO, ADEL, *et al.* 2016] Han ZHAO, Tameem ADEL, Geoff GORDON, and Brandon AMOS. “Collapsed variational inference for sum-product networks”. In: *Proceedings of The 33rd International Conference on Machine Learning*. Ed. by Maria Florina BALCAN and Kilian Q. WEINBERGER. Vol. 48. Proceedings of Machine Learning Research. New York, New York, USA: PMLR, June 2016, pp. 1310–1318. URL: [http://proceedings.mlr.press/v48/zhaoa16.html](https://proceedings.mlr.press/v48/zhaoa16.html) (cit. on p. 9).
- [H. ZHAO, MELIBARI, *et al.* 2015] Han ZHAO, Mazen MELIBARI, and Pascal POUPART. “On the relationship between sum-product networks and Bayesian networks”. In: *Proceedings of the 32nd International Conference on Machine Learning*. Vol. 37. Proceedings of Machine Learning Research. 2015, pp. 116–124 (cit. on p. 9).
- [H. ZHAO, POUPART, *et al.* 2016] Han ZHAO, Pascal POUPART, and Geoffrey J GORDON. “A unified approach for learning the parameters of sum-product networks”. In: *Advances in Neural Information Processing Systems*. NeurIPS. 2016 (cit. on p. 43).

- [ZHENG *et al.* 2018] Kaiyu ZHENG, Andrzej PRONOBIS, and Rajesh PN RAO. “Learning graph-structured sum-product networks for probabilistic semantic maps”. In: *Thirty-Second AAAI Conference on Artificial Intelligence*. 2018 (cit. on p. 21).
- [ZIEGLER and RUSH 2019] Zachary M. ZIEGLER and Alexander M. RUSH. “Latent normalizing flows for discrete sequences”. In: *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*. Ed. by Kamalika CHAUDHURI and Ruslan SALAKHUTDINOV. Vol. 97. Proceedings of Machine Learning Research. PMLR, 2019, pp. 7673–7682. URL: <http://proceedings.mlr.press/v97/ziegler19a.html> (cit. on p. 2).