On the Latent Variable Interpretation in Sum-Product Networks

Robert Peharz, Robert Gens, Franz Pernkopf, Senior Member, IEEE, and Pedro Domingos

Abstract—One of the central themes in Sum-Product networks (SPNs) is the interpretation of sum nodes as marginalized latent variables (LVs). This interpretation yields an increased syntactic or semantic structure, allows the application of the EM algorithm and to efficiently perform MPE inference. In literature, the LV interpretation was justified by explicitly introducing the indicator variables corresponding to the LVs' states. However, as pointed out in this paper, this approach is in conflict with the completeness condition in SPNs and does not fully specify the probabilistic model. We propose a remedy for this problem by modifying the original approach for introducing the LVs, which we call SPN augmentation. We discuss conditional independencies in augmented SPNs, formally establish the probabilistic interpretation of the sum-weights and give an interpretation of augmented SPNs as Bayesian networks. Based on these results, we find a sound derivation of the EM algorithm for SPNs. Furthermore, the Viterbi-style algorithm for MPE proposed in literature was never proven to be correct. We show that this is indeed a correct algorithm, when applied to selective SPNs, and in particular when applied to augmented SPNs. Our theoretical results are confirmed in experiments on synthetic data and 103 real-world datasets.

Index Terms—Sum-product networks, latent variables, mixture models, expectation-maximization, MPE inference

1 INTRODUCTION

S UM-PRODUCT Networks (SPNs) are a promising type of probabilistic model, combining the domains of deep learning and graphical models [1], [2]. One of their main advantages is that many interesting inference scenarios are expressed as single forward and/or backward passes, i.e., these inference scenarios have a computational cost linear in the SPN's representation size. SPNs have shown convincing performance in applications such as image completion [1], [3], [4], computer vision [5], classification [6] and speech and language modeling [7], [8], [9]. Since their proposition [1], one of the central themes in SPNs has been their interpretation as hierarchically structured latent variable (LV) models. This is essentially the same approach as the LV interpretation in mixture model. Consider for example a Gaussian mixture model (GMM) with K components over a set of random variables (RVs) **X**

$$p(\mathbf{X}) = \sum_{k=1}^{K} w_k \mathcal{N}(\mathbf{X} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \qquad (1)$$

where $\mathcal{N}(\cdot | \cdot)$ is the Gaussian PDF, $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$ are the means and covariances of the *k*th component, and w_k are the mixture

 F. Pernkopf is with the Signal Processing and Speech Communication Lab, Graz University of Technology, Graz 8010, Austria. E-mail: pernkopf@tugraz.at.

Manuscript received 12 Nov. 2015; revised 20 June 2016; accepted 30 Sept. 2016. Date of publication 17 Nov. 2016; date of current version 12 Sept. 2017. Recommended for acceptance by G. Elidan.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TPAMI.2016.2618381 weights with $w_k \ge 0$, $\sum w_k = 1$. The GMM can be interpreted in two ways: i) It is a convex combination of PDFs and thus itself a PDF, or ii) it is a marginal distribution of a distribution $p(\mathbf{X}, Z)$ over \mathbf{X} and a latent, marginalized variable Z, where $p(\mathbf{X} | Z = k) = \mathcal{N}(\mathbf{X} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ and $p(Z = k) = w_k$. The second interpretation, the LV interpretation, yields a syntactically wellstructured model. For example, following the LV interpretation, it is clear how to draw samples from $p(\mathbf{X})$ by using ancestral sampling. This structure can also be of semantic nature, for instance when Z represents a clustering of \mathbf{X} or when Z is a class variable. Furthermore, the LV interpretation allows the application of the EM algorithm—which is essentially maximum-likelihood learning under missing data [10], [11]—and enables advanced Bayesian techniques [12], [13].

Mixture models can be seen as a special case of SPNs with a single sum node, which corresponds to a single LV. More generally, SPNs can have arbitrarily many sum nodes, each corresponding to its own LV, leading to a hierarchically structured model. In [1], the LV interpretation in SPNs was justified by explicitly introducing the LVs in the SPN model, using the so-called *indicator variables* (IVs) corresponding to the LVs' states. However, as shown in this paper, this justification is actually too simplistic, since it is potentially in conflict with the completeness condition [1], leading to an incompletely specified model. As a remedy we propose the *augmentation* of an SPN, which additionally to the IVs also introduces the so-called twin sum nodes, in order to completely specify the LV model. We further investigate the independency structure of the LV model resulting from augmentation and find a parallel to the local independence assertions in Bayesian networks (BNs) [14], [15]. This allows us to define a BN representation of the augmented SPN. Using our BN interpretation and the differential approach to inference [16], [17] in augmented SPNs, we give a sound derivation of the (soft) EM algorithm for SPNs.

This work is licensed under a Creative Commons Attribution 3.0 License. For more information, see http://creativecommons.org/licenses/by/3.0/

R. Peharz is with the Institute of Physiology (iDN), Medical University of Graz, and BioTechMed—Graz, Graz 8036, Austria. E-mail: robert.peharz@gmail.com.

R. Gens and P. Domingos are with the Department of Computer Science and Engineering, University of Washington, Seattle, WA 98105. E-mail: {rcg, pedrod}@cs.washington.edu.

Closely related to the LV interpretation is the inference scenario of finding the most-probable-explanation (MPE), i.e., finding a probability maximizing assignment for all RVs. Using results from [18], [19], we first point out that this problem is generally NP-hard for SPNs. In [1] it was proposed that an MPE solution can be found efficiently when maximizing over both model RVs (i.e., non-latent RVs) and LVs. The proposed algorithm replaces sum nodes by max nodes and recovers the solution by using Viterbi-style backtracking. However, it was not shown that this algorithm delivers a correct MPE solution. In this paper, we show that this algorithm is indeed correct, when applied to selective SPNs [20]. In particular, since augmented SPNs are selective, this algorithm obtains an MPE solution in augmented SPNs. However, when applied to non-augmented SPNs, the algorithm still returns an MPE solution of the augmented SPN, but implicitly assumes that the weights for all twin sums are deterministic, i.e., they are all 0 except a single 1. This leads to a phenomenon in MPE inference which we call low-depth bias, i.e., more shallow parts of the SPN are preferred during backtracking.

The main contribution in this paper is to provide a sound theoretical foundation for the LV interpretation in SPNs and related concepts, i.e., the EM algorithm and MPE inference. Our theoretical findings are confirmed in experiments on synthetic data and 103 real-world datasets.

The paper is organized as follows: In the remainder of this section we introduce notation, review SPNs and discuss related work. In Section 2 we propose the augmentation of SPNs, show its soundness as hierarchical LV model and give an interpretation as BN. Furthermore, we discuss independency properties in augmented SPNs and the interpretation of sum-weights as conditional probabilities. The EM algorithm for SPNs is derived in Section 3. In Section 4 we discuss MPE inference for SPNs. Experiments are presented in Sections 5 and 6 concludes the paper. Proofs for our theoretical findings are deferred to the Appendix.

1.1 Background and Notation

RVs are denoted by upper-case letters W, X, Y and Z. The set of values of an RV X is denoted by val(X), where corresponding lower-case letters denote elements of val(X), e.g., x is an element of val(X). Sets of RVs are denoted by boldface letters W, X, Y and Z. For RV set $X = \{X_1, \ldots, X_N\}$, we define $val(X) = \times_{n=1}^N val(X_n)$ and use corresponding lower-case boldface letters for elements of val(X), e.g., x is an element of val(X). For a sub-set $Y \subseteq X$, x[Y] denotes the projection of x onto Y.

The elements of $\mathbf{val}(X)$ can be interpreted as *complete* evidence, assigning each RV in X a fixed value. *Partial* evidence about X is represented as a subset $\mathcal{X} \subseteq \mathbf{val}(X)$, which is an element of the sigma-algebra \mathcal{A}_X induced by RV X. For all RVs we use $\mathcal{A}_X = \{\mathcal{X} \in \mathcal{B} \mid \mathcal{X} \subseteq \mathbf{val}(X)\}$, \mathcal{B} being the Borelsets over \mathbb{R} . For discrete RVs, this choice yields the powerset $\mathcal{A}_X = 2^{\mathbf{val}(X)}$. For example, partial evidence $\mathcal{X} = \{1, 3, 5\}$ for a discrete RV X with $\mathbf{val}(X) = \{1, \ldots, 6\}$ represents evidence that X takes one of the states 1, 3 or 5, and $\mathcal{Y} = [-\infty, \pi]$ for a real-valued RV Y represents evidence that Y takes a value smaller than π . Formally speaking, partial evidence is used to express the domain of *marginalization* for a particular RV.

For sets of RVs $\mathbf{X} = \{X_1, \ldots, X_N\}$, we use the product sets $\mathcal{H}_{\mathbf{X}} := \{\times_{n=1}^N \mathcal{X}_n \mid \mathcal{X}_n \in \mathcal{A}_{X_n}\}$ to represent partial evidence about \mathbf{X} . Elements of $\mathcal{H}_{\mathbf{X}}$ are denoted using boldface notation, e.g., \mathcal{X} . When $\mathbf{Y} \subseteq \mathbf{X}$ and $\mathcal{X} \in \mathcal{H}_{\mathbf{X}}$, we define $\mathcal{X}[\mathbf{Y}] := \{\mathbf{x}[\mathbf{Y}] \mid \mathbf{x} \in \mathcal{X}\}$. Furthermore, we use \mathbf{e} to symbolize any combination of complete and partial evidence, i.e., for RVs \mathbf{X} we have some complete evidence \mathbf{x}' for $\mathbf{X}' \subseteq \mathbf{X}$ and some partial evidence $\mathcal{X}'' \in \mathcal{H}_{\mathbf{X}''}$ for $\mathbf{X}'' = \mathbf{X} \setminus \mathbf{X}'$.

Given a node N in some directed graph \mathcal{G} , let ch(N) and pa(N) be the set of children and parents of N, respectively. Furthermore, let desc(N) be the set of descendants of N, recursively defined as the set containing N itself and any child of a descendant. Similarly, we define anc(N) as the ancestors of N, recursively defined as the set containing N itself and any parent of an ancestor. SPNs are defined as follows.

Definition 1 (Sum-Product Network). A Sum-Product network S over a set of RVs **X** is a tuple (\mathcal{G}, w) where \mathcal{G} is a connected, rooted and acyclic directed graph, and w is a set of nonnegative parameters. The graph *G* contains three types of nodes: distributions, sums and products. All leaves of *G* are distributions and all internal nodes are either sums or products. A distribution node (also called input distribution or simply distribution) $D_{\mathbf{Y}} : \mathbf{val}(\mathbf{Y}) \mapsto [0, \infty]$ is a distribution function over a subset of RVs $\mathbf{Y} \subseteq \mathbf{X}$, i.e., either a PMF (discrete RVs), a PDF (continuous RVs), or a mixed distribution function (discrete and continuous RVs mixed). A sum node S computes a weighted sum of its children, i.e., $S = \sum_{C \in ch(S)} \mathit{w}_{S,C} \, C$, where $w_{S,C}$ is a non-negative weight associated with edge $S \rightarrow C$, and w contains the weights for all outgoing sum-edges. A product node P computes the product over its children, i.e., $\mathsf{P} = \prod_{\mathsf{C} \in \mathsf{ch}(\mathsf{P})} \mathsf{C}$. The sets $\mathsf{S}(\mathcal{S})$ and $\mathsf{P}(\mathcal{S})$ contain all sum nodes and all product nodes in *S*, respectively.

The size |S| of the SPN is defined as the number of nodes and edges in G. For any node N in G, the scope of N is defined as

$$\mathbf{sc}(\mathsf{N}) = \begin{cases} \mathsf{Y} & \text{if } \mathsf{N} \text{ is a distribution } \mathsf{D}_{\mathsf{Y}} \\ \bigcup_{\mathsf{C} \in \mathsf{ch}(\mathsf{N})} \mathbf{sc}(\mathsf{C}) & \text{otherwise.} \end{cases}$$
(2)

The function computed by S is the function computed by its root and denoted as $S(\mathbf{x})$, where without loss of generality we assume that the scope of the root is \mathbf{X} .

We use symbols D, S, P, N, C and F for nodes in SPNs, where D denotes a distribution, S denotes a sum, and P denotes a product. Symbols N, C and F denote generic nodes, where C and F indicate a child or parent relationship to another node, respectively. The *distribution* p_S of an SPN S is defined as the normalized output of S, i.e., $p_S(\mathbf{x}) \propto S(\mathbf{x})$. For each node N, we define the *sub-SPN* S_N rooted at N as the SPN defined by the graph induced by the descendants of N and the corresponding parameters.

Inference in unconstrained SPNs is generally intractable. However, efficient inference in SPNs is enabled by two structural constraints, *completeness* and *decomposability* [1]. An SPN is *complete* if for all sums **S** it holds that

$$\forall \mathbf{C}', \mathbf{C}'' \in \mathbf{ch}(\mathbf{S}) : \mathbf{sc}(\mathbf{C}') = \mathbf{sc}(\mathbf{C}''). \tag{3}$$

An SPN is *decomposable* if for all products P it holds that

$$\forall \mathbf{C}', \mathbf{C}'' \in \mathbf{ch}(\mathsf{P}), \mathbf{C}' \neq \mathbf{C}'' : \mathbf{sc}(\mathbf{C}') \cap \mathbf{sc}(\mathbf{C}'') = \emptyset.$$
(4)

Furthermore, a sum node S is called *selective* [20] if for all choices of sum-weights w and all possible inputs x it holds that at most one child of S is non-zero. An SPN S is called selective if all its sum nodes are selective.

As shown in [17], [19], integrating $S(\mathbf{x})$ over arbitrary sets $\mathcal{X} \in \mathcal{H}_{\mathbf{X}}$, i.e., *marginalization* over \mathcal{X} , reduces to the corresponding integrals at the input distributions and evaluating sums and products in the usual way. This property is known as *validity* of the SPNs [1], and key for efficient inference. In this paper we only consider complete and decomposable SPNs. Without loss of generality [17], [21], we assume *locally normalized* sum-weights, i.e., for each sum node S we have $\sum_{C \in ch(S)} w_{S,C} = 1$, and thus $p_S \equiv S$, i.e., the SPN's normalization constant is 1.

For RVs with finitely many states, we will use so-called *indicator variables* as input distributions [1]. For a finitestate RV X and state $x \in val(X)$, we introduce the IV $\lambda_{X=x}(x') := \mathbb{1}(x = x')$, assigning all probability mass to x. A complete and decomposable SPN represents the *(extended) network polynomial* of p_S , which can be used in the *differential approach to inference* [1], [16], [17]. Assume any evidence **e** which is evaluated in the SPN. The derivatives of the SPN function with respect to the IVs (by interpreting the IVs as real-valued variables, see [16], [17] for details) yield

$$\frac{\partial \mathcal{S}(\mathbf{e})}{\partial \lambda_{X=x}} = \mathcal{S}(X = x, \mathbf{e} \setminus X), \tag{5}$$

representing the inference scenario of *modified evidence*, i.e., evidence **e** is modified such that *X* is set to *x*. The computationally attractive feature of the differential approach is that (5) can be evaluated for *all* $X \in \mathbf{X}$ and *all* $x \in \mathbf{val}(X)$ simultaneously using a *single* back-propagation pass in the SPN, after evidence has been evaluated. Similarly, for the second (and higher) derivatives, we get

$$\frac{\partial^2 \mathcal{S}(\mathbf{e})}{\partial \lambda_{X=x} \lambda_{Y=y}} = \begin{cases} \mathcal{S}(X=x, Y=y, \mathbf{e} \setminus \{X, Y\}) & \text{if } X \neq Y \\ 0 & \text{otherwise.} \end{cases}$$
(6)

Furthermore, the differential approach can be generalized to SPNs with arbitrary input distributions, i.e., SPNs over RVs with countably infinite or uncountably many states (cf. [17] for details).

1.2 Related Work

SPNs are related to negation normal forms (NNFs), a potentially deep network representation of propositional theories [22], [23], [24]. Like in SPNs, structural constraints in NNFs enable certain polynomial-time queries in the represented theory. In particular, the notions of smoothness, decomposability and determinism in NNFs translate to the notions of completeness, decomposability and selectivity in SPNs, respectively. The work on NNFs led to the concept of network polynomials as a multilinear representation of BNs over finitely many states [16], [25]. BNs were cast into an intermediate deterministic decomposable NNF (d-DNNF) representation in order to generate an arithmetic circuit (AC), representing the BN's network polynomial. ACs, when restricted to sums and products, are equivalent to SPNs but have a slightly different syntax. In [26], ACs were learned by optimizing an objective trading off the loglikelihood on the training set and the inference cost of the AC, measured as the worst-case number of arithmetic operations required for inference (i.e., the number of edges in the AC). The learned models still represent BNs with context-specific independencies [27]. A similar approach learning Markov networks represented by ACs is followed in [28]. SPNs were the first time proposed in [1], where the represented distribution was not defined via a background graphical model any more, but directly as the normalized output of the network. In this work, SPNs were applied to image data, where a generic architecture reminiscent to convolutional neural networks was proposed. Structure learning algorithms not restricted to the image domain were proposed in [2], [3], [4], [29], [30], [31]. Discriminative learning of SPNs, optimizing conditional likelihood, was proposed in [6]. Furthermore, there is a growing body of literature on theoretical aspects of SPNs and their relationship to other types of probabilistic models. In [32] two families of functions were identified which are efficiently representable by deep, but not by shallow SPNs, where an SPN is considered as shallow if it has no more than three layers. In [17] it was shown that SPNs can w.l.o.g. be assumed to be locally normalized and that the notion of consistency does not allow exponentially more compact models than decomposability. These results were independently found in [21]. Furthermore, in [17], a sound derivation of inference mechanisms for generalized SPNs was given, i.e., SPNs over RVs with (uncountably) infinitely many states. In [21], a BN representation of SPNs was found, where LVs associated with sum nodes and the model RVs are organized in a two layer bipartite structure. The actual SPN structure is captured in structured conditional probability tables (CPTs) using algebraic decision diagrams. Recently, the notion of SPNs was generalized to sumproduct functions over arbitrary semirings [33]. This yields a general unifying framework for learning and inference, subsuming, among others, SPNs for probabilistic modeling, NNFs for logical propositions and function representations for integration and optimization.

2 LATENT VARIABLE INTERPRETATION

As pointed out in [1], each sum node in an SPN can be interpreted as a marginalized LV, similar as in the GMM example in Section 1. For each sum node S, one postulates a discrete LV *Z* whose states correspond to the children of S. For each state, an IV and a product is introduced, such that the children are switched on/off by the corresponding IVs, as illustrated in Fig. 1.¹ When all IVs in Fig. 1b are set to 1, S still computes the same value as in Fig. 1a. Since setting all IVs of *Z* to 1 corresponds to marginalizing *Z*, the sum S should be interpreted as a latent, marginalized RV.

1. In graphical representations of SPNs, IVs are depicted as nodes containing a small circle, general distributions as nodes containing a Gaussian-like PDF, and sum and products as nodes with + and \times symbols. Empty nodes are of arbitrary type.



Fig. 1. Problems occurring when IVs of LVs are introduced. (a): Excerpt of SPN containing a sum S, corresponding to LV Z. (b): Introducing IVs for Z renders S' incomplete, assuming that $S \notin desc(N)$. (c): Remedy by extending SPN further, introducing twin sum node \overline{S} .

However, when we regard a larger structural context in Fig. 1b, we recognize that this justification is actually too simplistic. Explicitly introducing the IVs renders the ancestor S' incomplete, when S is no descendant of N, and Z is thus not in the scope of N. Note that setting all IVs to 1 in an *incomplete* SPN generally does *not* correspond to marginalization. Furthermore, note that also S' corresponds to an LV, say Z'. While we know the probability distribution of Z if Z' is in the state corresponding to P, namely the weights of S, we do not know this distribution when Z' is in the state corresponding to N. Intuitively, we recognize that the state of Z is irrelevant in this case, since it does not influence the resulting distribution over the model RVs **X**. Nevertheless, the probabilistic model is not completely specified, which is unsatisfying.

A remedy for these problems is shown in Fig. 1c. We introduce the twin sum node S whose children are the IVs corresponding to Z. The twin \bar{S} is connected as child of an additional product node, which is interconnected between S' and N. Since this new product node has scope $sc(N) \cup \{Z\}, S'$ is rendered complete now. Furthermore, if Z' takes the state corresponding to N (or actually the state corresponding to the new product node), we now have a specified conditional distribution for Z, namely the weights of the twin sum node. Clearly, given that all IVs of Z are set to 1, the network depicted in Fig. 1c still computes the same function as the network in Fig. 1a (or Fig. 1b), since S constantly outputs 1, as long as we use normalized weights for it. Which weights should be used for the twin sum node S? Basically, we can assume arbitrary normalized weights, which will cause S to constantly output 1, where, however, a natural choice would be to use uniform weights for \bar{S} (maximizing the entropy of the resulting LV model). Although the choice of weights is not crucial for *evaluating* evidence in the SPN, it plays a role in MPE inference, see Section 4. For now, let us formalize the explicit introduction of LVs, denoted as augmentation.

2.1 Augmentation of SPNs

Let *S* be an SPN over **X**. For each $S \in S(S)$ we assume an arbitrary but fixed ordering of its children $ch(S) = \{C_S^1, \ldots, C_S^{K_S}\}$, where $K_S = |ch(S)|$. Let Z_S be an RV on the same probability space as **X**, with $val(Z_S) = \{1, \ldots, K_S\}$, where state *k* corresponds to child C_S^k . We call Z_S the *LV* associated with **S**. For sets of sum nodes **S** we define $Z_S = \{Z_S | S \in S\}$. To distinguish **X** from the LVs, we will

1: **procedure** AUGMENTSPN(S)

 $S' \leftarrow S$ 2: $\forall \mathsf{S} \in \mathsf{S}(\mathcal{S}'), \ \forall k \in \{1, \ldots, K_{\mathsf{S}}\}:$ 3. let $w_{\mathsf{S},k} = w_{\mathsf{S},\mathsf{C}^k_\mathsf{c}}$, $\overline{w}_{\mathsf{S},k} = \overline{w}_{\mathsf{S},\mathsf{C}^k_\mathsf{c}}$ for $S \in \mathbf{S}(\mathcal{S}')$ do 4: 5: for $k = 1 \dots K_S$ do Introduce a new product node $\mathsf{P}^k_{\mathsf{S}}$ in $\mathsf{S}(\mathcal{S}')$ 6: Disconnect C_{S}^{k} from S 7: Connect C_{S}^{k} as child of P_{S}^{k} 8: Connect $\mathsf{P}^{\overline{k}}_{\mathsf{S}}$ as child of S with weight $w_{\mathsf{S},k}$ 9: 10:end for 11: end for for $S \in \mathbf{S}(S')$ do $12 \cdot$ 13: for $k \in \{1, ..., K_{S}\}$ do Connect new IV $\lambda_{Z_S=k}$ as child of P^k_S 14: 15: end for 16: if $\mathbf{S}^{c}(S) \neq \emptyset$ then $17 \cdot$ Introduce a twin sum node S in S'18: $\forall k \in \{1, \ldots, K_{\mathsf{S}}\}$: connect $\lambda_{Z_{\mathsf{S}}=k}$ as child of S , and let $w_{\bar{S},\lambda_{Z_{S}=k}} = \bar{w}_{S,k}$ for $S^c \in S^c(S)$ do 19. for $k \in \{k \mid \mathsf{S} \notin \operatorname{desc}(\mathsf{P}^k_{\mathsf{S}^c})\}$ do $20 \cdot$ Connect \overline{S} as child of P_{Sc}^k 21: 22. end for end for 23. 24: end if 25: end for return S'26: 27: end procedure

Fig. 2. Pseudo-code for augmentation of an SPN.

refer to the former as *model RVs*. For node N, we define the *sum ancestors/descendants* as

$$\operatorname{anc}_{\mathbf{S}}(\mathsf{N}) := \operatorname{anc}(\mathsf{N}) \cap \mathbf{S}(\mathcal{S}),$$
 (7)

$$\operatorname{desc}_{\mathsf{S}}(\mathsf{N}) := \operatorname{desc}(\mathsf{N}) \cap \mathsf{S}(\mathcal{S}). \tag{8}$$

For each sum node S we define the *conditioning sums* as

$$\mathbf{S}^{c}(\mathbf{S}) := \{\mathbf{S}^{c} \in \mathbf{anc}_{\mathbf{S}}(\mathbf{S}) \setminus \{\mathbf{S}\} \mid \exists \mathbf{C} \in \mathbf{ch}(\mathbf{S}^{c}) : \mathbf{S} \notin \mathbf{desc}(\mathbf{C})\}.$$
(9)

Furthermore, we assume a set of locally normalized *twinweights* \bar{w} , containing a twin-weight $\bar{w}_{S,C}$ for each weight $w_{S,C}$ in the SPN. We are now ready to define the *augmentation* of an SPN.

Definition 2 (Augmentation of SPN). Let S be an SPN over X, \bar{w} be a set of twin-weights and S' be the result of algorithm AUGMENTSPN, shown in Fig. 2. S' is called the augmented SPN of S, denoted as S' =: aug(S). Within the context of S', C_S^k is called the kth former child of S. The introduced product node P_S^k is called link of S, C_S^k and $\lambda_{Z_S=k}$, respectively. The sum node \bar{S} , if introduced, is called the twin sum node of S. With respect to S', we denote S as the original SPN.

In steps 4–11 of AUGMENTSPN we introduce the links P_{S}^{k} which are interconnected between sum node S and its *k*th child. Each link P_{S}^{k} has a single parent, namely S, and simply copies the former child C_{S}^{k} . In steps 13–15, we introduce IVs corresponding to the associated LV Z_{S} , as



Fig. 3. Augmentation of an SPN. (a): Example SPN over $X = \{X_1, X_2, X_3\}$, containing sum nodes S^1 , S^2 , S^3 and S^4 . (b): Augmented SPN, containing IVs corresponding to Z_{S^1} , Z_{S^2} , Z_{S^3} , Z_{S^4} , links and twin sum nodes \bar{S}^2 , \bar{S}^3 , \bar{S}^4 . For nodes introduced by augmentation, smaller circles are used.

proposed in [1]. As we saw in Fig. 1 and the discussion above, this can render other sum nodes incomplete. These sums are clearly the conditioning sums $S^{c}(S)$. Thus, when necessary, we introduce a twin sum node in steps 17–23, to treat this problem. The following proposition states the soundness of augmentation.

Proposition 1. Let S be an SPN over X, $S' = \operatorname{aug}(S)$ and $Z := Z_{S(S)}$. Then S' is a complete and decomposable SPN over $X \cup Z$ with $S'(X) \equiv S(X)$.

Proposition 1 states that the marginal distribution over X in the augmented SPN is the same distribution as represented by the original SPN, while being a *completely specified probabilistic model* over X and Z. Thus, augmentation provides a sound way to generalize the LV interpretation from mixture models to more general SPNs. An example of augmentation is shown in Fig. 3.

Note that we understand the augmentation mainly as a *theoretical* tool to establish and work with the LV interpretation in SPNs. In most cases, it will be neither necessary nor advisable to *explicitly* construct the augmented SPN.

An interesting question is how the sizes of the original SPN and the augmented SPN relate to each other. A lower bound is $|S'| \in \Omega(|S|)$, holding, e.g., for SPNs with a single sum node. An asymptotic upper bound is $|S'| \in \mathcal{O}(|S|^2)$. To see this, note that the introduction of links, IVs and twin sums cause at most a linear increase of the SPN's size. The number of edges introduced when

connecting twins to the links of conditioning sums is bounded by $|\mathcal{S}|^2$, since the number of twins and links are both bounded by $|\mathcal{S}|$. Therefore, we have $|\mathcal{S}'| \in \mathcal{O}(|\mathcal{S}|^2)$. This asymptotic upper bound is indeed achieved by certain types of SPNs: Consider, e.g., a chain consisting of K sum nodes and K + 1 distribution nodes. For k < Kthe *k*th sum is the parent of the (k+1)th sum and the *k*th distribution, and the *K*th sum is the parent of the last two distributions. For the *k*th sum, all preceding sums are conditioning sums, yielding k - 1 introduced edges. In total this gives $\sum_{k=2}^{K} (k-1) = \frac{K(K-1)}{2} = \frac{K^2 - K}{2}$ edges, i.e., in this case $|\mathcal{S}'|$ indeed grows quadratically in $|\mathcal{S}|$.

2.2 Conditional Independencies in Augmented SPNs and Probabilistic Interpretation of Sum-Weights

It is helpful to introduce the notion of configured SPNs, which takes a similar role as conditioning in the literature on DNNFs [22], [23], [24].

Definition 3 (Configured SPN). Let S be an SPN over X, $\mathbf{Y} \subseteq \mathbf{Z}_{\mathbf{S}(S)}$ and $\mathbf{y} \in \mathbf{val}(\mathbf{Y})$. The configured SPN S^y is obtained by deleting the IVs $\lambda_{Y=y}$ and their corresponding link for each $Y \in \mathbf{Y}, y \neq \mathbf{y}[Y]$ from $\mathbf{aug}(S)$, and further deleting all nodes which are rendered unreachable from the root.

Intuitively, the configured SPN isolates the computational structure selected by **y**. All sum edges which "survive" in the configured SPN are equipped with the same weights as in the augmented SPN. Therefore, a configured SPN is in general not locally normalized. We note the following properties of configured SPNs.

Proposition 2. Let S be an SPN over X, $Y \subseteq Z_{S(S)}$ and $Z = Z_{S(S)} \setminus Y$. Let $y \in val(Y)$ and let S' = aug(S). It holds that

- 1) Each node in S^y has the same scope as its corresponding node in S'.
- 2) $S^{\mathbf{y}}$ is a complete and decomposable SPN over $\mathbf{X} \cup \mathbf{Y} \cup \mathbf{Z}$.
- 3) For any node N in $S^{\mathbf{y}}$ with $\mathbf{sc}(N) \cap \mathbf{Y} = \emptyset$, we have that $S_{N}^{\mathbf{y}} = S_{N}'$.
- 4) For $\mathbf{y}' \in \mathbf{val}(\mathbf{Y})$ it holds that

$$S^{\mathbf{y}}(\mathbf{X}, \mathbf{Z}, \mathbf{y}') = \begin{cases} S'(\mathbf{X}, \mathbf{Z}, \mathbf{y}') & \text{if } \mathbf{y}' = \mathbf{y} \\ 0 & \text{otherwise.} \end{cases}$$
(10)

The next theorem shows certain conditional independencies in the augmented SPN. For ease of discussion, we make the following definitions.

Definition 4. Let **S** be a sum node in an SPN and Z_S its associated LV. All other RVs (model RVs and LVs) are divided into three sets:

- Parents \mathbf{Z}_p , which are all LVs "above" S, i.e., $\mathbf{Z}_p = \mathbf{Z}_{\operatorname{anc}_{S(S)} \setminus Z_S}$.
- Children Y_c, which are all model RVs and LVs "below"
 S, i.e., Y_c = sc(S) ∪ Z_{desc_S(S)}\Z_S.
- Non-descendants Y_n, which are the remaining RVs,
 i.e., Y_n = (X ∪ Z_S(S)) \ (Z_p ∪ Y_c ∪ Z_S).



Fig. 4. Dependency structure of augmented SPN from Fig. 3, represented as BN.

We will show that the *parents*, *children* and *nondescendants* play the likewise role as for independencies in BNs [14], [15], i.e., Z_S is independent of \mathbf{Y}_n given \mathbf{Z}_p . We will further show that the sum-weights of \mathbf{S} are the conditional distribution of Z_S , conditioned on the event that " \mathbf{Z}_p select a path to \mathbf{S} ". One problem in the original LV interpretation [1] was, that no conditional distribution of Z_S was specified for the complementary event. Here, we will show that the twin-weights are precisely this conditional distribution. This requires that the event " \mathbf{Z}_p select a path to the twin $\overline{\mathbf{S}}$ " is indeed the complementary event to " \mathbf{Z}_p select a path to \mathbf{S} ". This is shown in following lemma.

Lemma 1. Let S be an SPN over X, let S be a sum node in S and Z_p be the parents of Z_S . For any $z \in val(Z_p)$, the configured SPN S^z contains either S or its twin \overline{S} , but not both.

We are now ready to state our theorem concerning conditional independencies in augmented SPNs.

Theorem 1. Let S be an SPN over X and S' = $\operatorname{aug}(S)$. Let S be an arbitrary sum in S and $w_k = w_{S,C_S^k}$, $\overline{w}_k = \overline{w}_{S,C_S^k}$, $k = 1, \ldots, K_S$. With respect to S, let \mathbb{Z}_p be the parents, Y_c be the children and Y_n be the non-descendants, respectively. Then there exists a two-partition of $\operatorname{val}(\mathbb{Z}_p)$, i.e., $Z, \overline{Z} : Z \cup \overline{Z}$ = $\operatorname{val}(\mathbb{Z}_p), Z \cap \overline{Z} = \emptyset$, such that

$$\forall \mathbf{z} \in \boldsymbol{\mathcal{Z}} : \boldsymbol{\mathcal{S}}'(Z_{\mathsf{S}} = k, \mathbf{Y}_n, \mathbf{z}) = w_k \boldsymbol{\mathcal{S}}'(\mathbf{Y}_n, \mathbf{z}), and \quad (11)$$

$$\forall \mathbf{z} \in \bar{\mathbf{Z}} : \mathcal{S}'(Z_{\mathsf{S}} = k, \mathbf{Y}_n, \mathbf{z}) = \bar{w}_k \mathcal{S}'(\mathbf{Y}_n, \mathbf{z}).$$
(12)

From Theorem 1 it follows that the weights and twinweights of a sum node S can be interpreted as *conditional probability tables* of Z_S , conditioned on Z_p and that Z_S is conditionally independent of Y_n given Z_p , i.e.,

$$\mathcal{S}'(Z_{\mathsf{S}} = k \,|\, \mathbf{Y}_n, \mathbf{z}) = \mathcal{S}'(Z_{\mathsf{S}} = k \,|\, \mathbf{z}) = \begin{cases} w_k & \text{if } \mathbf{z} \in \mathbf{Z} \\ \bar{w}_k & \text{if } \mathbf{z} \in \mathbf{\bar{Z}}. \end{cases}$$
(13)

Using this result, we can define a BN representing the augmented SPN as follows: For each sum node S, connect Z_p as parents of Z_S , and all RVs sc(S) as children of Z_S . By doing this for each LV, we obtain our BN representation of the augmented SPN, serving as a useful tool to understand SPNs in the context of probabilistic graphical models. An example of the BN interpretation is shown in Fig. 4.

Note that the BN representation by Zhao et al. [21] can be recovered from the BN representation of augmented SPNs. They proposed a BN representation of SPNs using a bipartite structure, where an LV is a parent of a model RV if it is contained in the scope of the corresponding sum node. The model RVs and LVs are unconnected among each other, respectively. When we constrain the twin-weights to be equal to the sum-weights, we can see in (13) that Z_S becomes independent of Z_p . This special choice of twin weights effectively removes all edges between LVs, recovering the BN structure in [21]. In the next section, we use the augmented SPN and the BN interpretation to derive the EM algorithm for SPNs.

3 EM ALGORITHM

The EM algorithm is a general scheme for maximum likelihood learning, when for some RVs complete evidence is missing [10], [11]. Thus, augmented SPNs are amenable for EM due to the LVs associated with sum nodes. Moreover, the twin-weights can be kept fixed, so that EM applied to augmented SPNs actually optimizes the weights of the original SPN. This approach was already pointed out in [1], where it was suggested that for evidence e and for any LV Z_S , the marginal posteriors should be given as $p(Z_{S} = k | \mathbf{e}) \propto w_{S,C_{e}^{k}} \frac{\partial S(\mathbf{e})}{\partial S(\mathbf{e})}$, which should be used for EM updates. These updates, however, cannot be the correct ones, as they actually leave the weights unchanged. Here, using augmented SPNs, we formally derive the standard EM updates for sum-weights and the input distributions, when they are chosen from an exponential family.

3.1 Updates for Weights

Assume a dataset $\mathcal{D} = \{\mathbf{e}^{(1)}, \dots, \mathbf{e}^{(L)}\}\$ of *L* i.i.d. samples, where each $\mathbf{e}^{(l)}$ is any combination of complete and partial evidence for the model RVs **X**, cf. Section 1.1. Let $\mathbf{Z} = \mathbf{Z}_{\mathbf{S}(S)}$ be the set of all LVs and consider an arbitrary sum node **S**. Eq. (13) shows that the weights can be interpreted as conditional probabilities in our BN interpretation, where

$$\mathcal{S}'(Z_{\mathsf{S}} = k \,|\, \mathbf{Z}_p = \mathbf{z}) = \begin{cases} w_k & \text{if } \mathbf{z} \in \mathcal{Z} \\ \bar{w}_k & \text{if } \mathbf{z} \in \bar{\mathcal{Z}}. \end{cases}$$
(14)

As mentioned above, the twin-weights \bar{w}_k are kept fixed. Using the well-known EM-updates in BNs over discrete RVs [10], [15], the updates for sum-weight w_k are given by summing over the expected statistics

$$\mathcal{S}'(Z_{\mathsf{S}} = k, \mathbf{Z}_p \in \mathcal{Z} \,|\, \mathbf{e}^{(l)}),\tag{15}$$

followed by renormalization. We make the event $Z_p \in Z$ explicit, by introducing a *switching parent* Y_S of Z_S : When the twin sum of S exists, Y_S assumes the two states $val(Y_S) = \{y_S, y_S\}$, where $Y_S = y_S \Leftrightarrow Z_p \in Z$ and $Y_S = y_S \Leftrightarrow Z_p \in \overline{Z}$. When the twin sum does not exist, Y_S just takes the single value $val(Y_S) = \{y_S\}$. Clearly, when observed, Y_S renders Z_S independent from Z_p . The switching parent can be explicitly introduced in the augmented SPN, as depicted in Fig. 5.

Here we simply introduce two new IVs $\lambda_{Y_S=y_S}$ and $\lambda_{Y_S=y_{\bar{S}}}$, which switch on/off the output of **S** and \bar{S} , respectively. It is easy to see that when these IV are



Fig. 5. Explicitly introducing a switching parent $Y_{\rm S}$ in an augmented SPN. (a): Part of an augmented SPN containing a sum node with three children and its twin. (b): Explicitly introduced switching parent $Y_{\rm S}$ using IVs $\lambda_{Y_{\rm S}=y_{\rm S}}$ and $\lambda_{Y_{\rm S}=y_{\rm S}}$.

constantly set to 1, i.e., when Y_S is marginalized, the augmented SPN performs exactly the same computations as before. It is furthermore easy to see that completeness and decomposability of the augmented SPN are maintained when the switching parent is introduced. Using the switching parent, the required expected statistics (15) translate to

$$\mathcal{S}'(Z_{\mathsf{S}} = k, Y_{\mathsf{S}} = y_{\mathsf{S}} \mid \mathbf{e}^{(l)}). \tag{16}$$

To compute (16), we use the differential approach, [16], [17], [19], cf. also Section 1.1. First note that

$$\mathcal{S}'(Z_{\mathsf{S}} = k, Y_{\mathsf{S}} = y_{\mathsf{S}}, \mathbf{e}^{(l)}) = \frac{\partial^2 \mathcal{S}'(\mathbf{e}^{(l)})}{\partial \lambda_{Y_{\mathsf{S}} = y_{\mathsf{S}}} \partial \lambda_{Z_{\mathsf{S}} = k}}.$$
 (17)

The first derivative is given as

$$\frac{\partial \mathcal{S}'(\mathbf{e}^{(l)})}{\partial \lambda_{Y_{\mathbf{S}}=y_{\mathbf{S}}}} = \frac{\partial \mathcal{S}'(\mathbf{e}^{(l)})}{\partial \mathbf{P}} \, \mathbf{S}(\mathbf{e}^{(l)}) \tag{18}$$

$$= \frac{\partial \mathcal{S}'(\mathbf{e}^{(l)})}{\partial \mathbf{P}} \sum_{k=1}^{K_{\mathbf{S}}} \lambda_{Z_{\mathbf{S}}=k} w_k \mathbf{C}_{\mathbf{S}}^k(\mathbf{e}^{(l)}), \qquad (19)$$

where P is the common product parent of S and $\lambda_{Y_S=y_S}$ in the augmented SPN (see Fig. 5b). Differentiating (19) after $\lambda_{Z_S=k}$ yields the second derivative

$$\frac{\partial^2 \mathcal{S}'(\mathbf{e}^{(l)})}{\partial \lambda_{Y_{\mathsf{S}}=y_{\mathsf{S}}} \partial \lambda_{Z_{\mathsf{S}}=k}} = \frac{\partial \mathcal{S}'(\mathbf{e}^{(l)})}{\partial \mathsf{P}} w_k \, \mathsf{C}^k_{\mathsf{S}}(\mathbf{e}^{(l)}), \tag{20}$$

delivering the required posteriors

$$\mathcal{S}'(Z_{\mathsf{S}} = k, Y_{\mathsf{S}} = y_{\mathsf{S}} | \mathbf{e}^{(l)}) = \frac{1}{\mathcal{S}'(\mathbf{e}^{(l)})} \frac{\partial \mathcal{S}'(\mathbf{e}^{(l)})}{\partial \mathsf{P}} w_k \mathsf{C}^k_{\mathsf{S}}(\mathbf{e}^{(l)}).$$
(21)

We do not want to construct the augmented SPN explicitly, so we express (21) in terms of the original SPN. Since all LVs are marginalized, it holds that $S'(\mathbf{e}^{(l)}) = S(\mathbf{e}^{(l)})$ and $\frac{\partial S'(\mathbf{e}^{(l)})}{\partial \mathbf{P}} = \frac{\partial S(\mathbf{e}^{(l)})}{\partial \mathbf{S}}$, yielding

$$\mathcal{S}'(Z_{\mathsf{S}} = k, Y_{\mathsf{S}} = y_{\mathsf{S}} \,|\, \mathbf{e}^{(l)}) = \frac{1}{\mathcal{S}(\mathbf{e}^{(l)})} \frac{\partial \mathcal{S}(\mathbf{e}^{(l)})}{\partial \mathsf{S}} \,w_k \, \mathbf{C}^k_{\mathsf{S}}(\mathbf{e}^{(l)}), \quad (22)$$

delivering the required statistics for updating the sumweights. We now turn to the updates of the input distributions.

3.2 Updates for Input Distributions

For simplicity, we derive updates for univariate input distributions, i.e., for all distributions D_Y we have $|sc(D_Y)| = 1$. Similar updates can rather easily be derived also for multivariate input distributions. In [17], the socalled distribution selectors (DSs) were introduced to derive the differential approach for generalized SPNs. Similar as the switching parents for (twin) sum nodes, the DSs are RVs which render the respective model RVs independent from the remaining RVs. More formally, for each $X \in \mathbf{X}$, let \mathbf{D}_X be the set of all input distributions which have scope $\{X\}$. Assume an arbitrary but fixed ordering of \mathbf{D}_X and let $[\mathbf{D}_X]$ be the index of \mathbf{D}_X in this ordering. Let the DS W_X be a discrete RV with $|\mathbf{D}_X|$ states. The so-called gated SPN S^g is obtained by replacing each distribution by the product node

$$\mathsf{D}_X \to \mathsf{D}_X \times \lambda_{W_X = [\mathsf{D}_X]}.$$
 (23)

The introduced product is denoted as gate. As shown in [17], *X* is rendered independent from all other RVs in the SPN when conditioned on W_X . Moreover, D_X is the conditional distribution of *X* given $W_X = [D_X]$. Therefore, each *X* and its DS W_X can be incorporated as a two RV family in our BN interpretation. When each input distribution D_X is chosen from an exponential family with natural parameters θ_{D_X} , the M-step is given by the expected sufficient statistics

$$\theta_{\mathsf{D}_X} \leftarrow \frac{\sum_l \mathcal{S}^g(W_X = k \,|\, \mathbf{e}^{(l)}) \int \mathsf{D}_X(x \,|\, \mathbf{e}^{(l)}) \theta_{\mathsf{D}_X}(x) \mathrm{d}x}{\sum_l \mathcal{S}^g(W_X = k \,|\, \mathbf{e}^{(l)})}, \quad (24)$$

where $k = [D_X]$. When $\mathbf{e}^{(l)}$ contains complete evidence x' for X, then the integral $\int D_X(x | \mathbf{e}^{(l)}) \theta_{D_X}(x) dx$ reduces to $\theta_{D_X}(x')$. When $\mathbf{e}^{(l)}$ contains partial evidence \mathcal{X} , then

$$\int \mathsf{D}_X(x \,|\, \mathbf{e}^{(l)}) \theta_{\mathsf{D}_X}(x) \mathrm{d}x = \frac{\int_{\mathcal{X}} \mathsf{D}_X(x) \theta_{\mathsf{D}_X}(x) \mathrm{d}x}{\int_{\mathcal{X}} \mathsf{D}_X(x) \mathrm{d}x}.$$
 (25)

Depending on *X* and the the type of D_X , evaluating (25) can be more or less demanding. A simple but practical case is when D_X is Gaussian and \mathcal{X} is some interval, permitting a closed form solution for integrating the Gaussian's statistics $\theta(x) = (x, x^2)$, using truncated Gaussians [34].

To obtain the posteriors $S^{g}(W_{X} = k | \mathbf{e}^{(l)})$ required in (24), we again use the differential approach. Note that

$$\mathcal{S}^{g}(W_{X}=k,\mathbf{e}^{(l)}) = \frac{\partial \mathcal{S}^{g}(\mathbf{e}^{(l)})}{\partial \lambda_{W_{X}=k}} = \frac{\partial \mathcal{S}^{g}(\mathbf{e}^{(l)})}{\partial \mathbf{P}} \mathbf{D}_{X}(\mathbf{e}^{(l)}), \qquad (26)$$

1: **procedure** EXPECTATION-MAXIMIZATION(S)

2:	Initialize w and input distributions
3:	while not converged do
4:	$\forall S \in S(\mathcal{S}), \forall C \in \mathbf{ch}(S) \colon n_{S,C} \leftarrow 0$
5:	$\forall X \in \mathbf{X}, \forall D_X \in \mathbf{D}_X : \theta_{D_X} \leftarrow 0, n_{D_X} \leftarrow 0$
6:	for $l = 1 \dots L$ do
7:	Input $\mathbf{e}^{(l)}$ to $\mathcal S$
8:	Evaluate $\mathcal S$ (upward-pass)
9:	Backprop ${\cal S}$ (backward-pass)
10:	for $S \in \mathbf{S}(\mathcal{S}), C \in \mathbf{ch}(S)$ do
11:	$n_{S,C} \leftarrow n_{S,C} + \frac{1}{\mathcal{S}} \frac{\partial \mathcal{S}}{\partial S} C w_{S,C}$
12:	end for
13:	for $X \in \mathbf{X}$, $D_X \in \mathbf{D}_X$ do
14:	if $e^{(l)}$ is complete w.r.t. X then
15:	$x \leftarrow \text{complete evidence for } X$
16:	$\theta \leftarrow \theta(x)$
17:	else
18:	$\mathcal{X} \leftarrow \text{partial evidence for } X$
19:	$\theta \leftarrow \frac{\int_{\mathcal{X}} D_X(x)\theta(x) dx}{\int D_X(x) dx}$
20:	end if
21:	$p \leftarrow \frac{1}{S} \frac{\partial S}{\partial D} D_X$
22:	$\theta_{D_{\mathbf{X}}} \leftarrow \theta_{D_{\mathbf{X}}} + p \theta$
23:	$n_{D_X} \leftarrow n_{D_X} + p$
24:	end for
25:	end for
26:	$\forall S \in S(\mathcal{S}), \forall C \in \mathbf{ch}(S): \ w_{S,C} \leftarrow \frac{n_{S,C}}{\sum_{C' \in \mathbf{ch}(S)} n_{S,C'}}$
27:	$\forall X \in \mathbf{X}, \forall D_X \in D_X$: set parameters to $\frac{\theta_{D_X}}{\theta_{D_X}}$
28:	end while n_{D_X}
29:	return S
30:	end procedure
	1

Fig. 6. Pseudo-code for EM algorithm in SPNs.

where $k = [D_X]$ and P is the gate of D_X , cf. (23). If we do not want to construct the gated SPN explicitly, we can use the identity $\frac{\partial S^g(\mathbf{e}^{(l)})}{\partial D_X} = \frac{\partial S(\mathbf{e}^{(l)})}{\partial D_X}$. Thus the required posteriors are given as

$$\mathcal{S}^{g}(W_{X} = k \,|\, \mathbf{e}^{(l)}) = \frac{1}{\mathcal{S}(\mathbf{e}^{(l)})} \frac{\partial \mathcal{S}(\mathbf{e}^{(l)})}{\partial \mathbf{D}_{X}} \mathbf{D}_{X}(\mathbf{e}^{(l)}).$$
(27)

The EM algorithm for SPNs, both for sum-weights and input distributions, is summarized in Fig. 6. In Section 5.1 we empirically verify our derivation of EM and show that standard EM successfully trains SPNs when a suitable structure is at hand.

Note that recently Zhao and Poupart [35] derived a concave-convex procedure (CCCP) which yield the same sumweight updates as the EM algorithm presented here and in [19]. This result is surprising, as EM and CCCP are rather different approaches in general.

4 MOST PROBABLE EXPLANATION

In [1], [4], [7], SPNs were applied for reconstructing data using MPE inference. Given some distribution p over **X** and evidence **e**, MPE can be formalized as finding $\arg \max_{x \in \mathcal{A}} p(\mathbf{x})$,

where we assume that p actually has a maximum in e.

MPE is a special case of MAP, defined as finding $\underset{\mathbf{y} \in \mathbf{e}[\mathbf{Y}]}{\operatorname{arg\,max}} \int_{\mathbf{e}[\mathbf{Z}]} p(\mathbf{y}, \mathbf{z}) \, \mathrm{d}\mathbf{z}, \text{ for some two-partition of } \mathbf{X}, \text{ i.e.,}$ $X = Y \cup Z, Y \cap Z = \emptyset$. Both MPE and MAP are generally NP-hard in BNs [36], [37], [38], and MAP is inherently harder than MPE [37], [38]. Using the result in [18], it follows that MAP inference is NP-hard also in SPNs. In particular, Theorem 5 in [18] shows that the decision version of MAP is NP-complete for a Naive Bayes model, when the class variable is marginalized. Naive Bayes is represented by the augmentation of an SPN with a single sum node, the LV representing the class variable. Therefore, MAP in SPNs is generally NP-hard. Since MAP in the augmented SPN representing the Naive Bayes model corresponds to MPE inference in the original SPN, i.e., a mixture model, it follows that also MPE inference is generally NP-hard in SPNs. A proof tailored to SPNs can be found in [19].

However, when considering the the sub-class of *selective* SPNs (cf. Section 1.1 and [20]), an MPE solution can be obtained using a Viterbi-style backtracking algorithm in *max-product networks* (MPN).

Definition 5 (Max-Product Network). Let S be an SPN over

X. We define the max-product network \hat{S} , by replacing each distribution node D by a maximizing distribution node

$$\hat{\mathsf{D}}: \mathcal{H}_{\mathsf{sc}(\mathsf{D})} \mapsto [0, \infty], \hat{\mathsf{D}}(\mathcal{Y}) := \max_{\mathbf{y} \in \mathcal{Y}} \mathsf{D}(\mathbf{y}), \tag{28}$$

and each sum node S by a max node

$$\hat{\mathbf{S}} := \max_{\hat{\mathbf{C}} \in \mathbf{ch}(\hat{\mathbf{S}})} w_{\hat{\mathbf{S}},\hat{\mathbf{C}}} \hat{\mathbf{C}}.$$
(29)

A product node P in S corresponds to a product node P in \hat{S} .

Theorem 2. Let *S* be a selective SPN over **X** and let \hat{S} the corresponding MPN. Let N be some node in *S* and \hat{N} its corresponding node in \hat{S} . Then, for every $\mathcal{X} \in \mathcal{H}_{sc(N)}$ we have $\hat{N}(\mathcal{X}) = \max_{\mathbf{x} \in \mathcal{X}} N(\mathbf{x})$.

Theorem 2 shows that the MPN maximizes the probability in its corresponding selective SPN. The proof (see appendix) also shows how to actually find a maximizing assignment. For a product, a maximizing assignment is given by combining the maximizing assignments of its children. For a sum, a maximizing assignment is given by the maximizing assignment of a single child, whose weighted maximum is maximal among all children. Here the children's maxima are readily given by the upwards pass in the MPN. Thus, finding a maximizing assignment of any node in an selective SPN recursively reduces to finding maximizing assignments for the children of this node; this can be accomplished by a Viterbi-like backtracking procedure. This algorithm, denoted as MPESELECTIVE, is shown in Fig. 7. Here Q denotes a queue of nodes, where $Q \cap N$ and $N \cap Q$ denote the en-queue and de-queue operations, respectively. Note that Theorem 2 has already been derived for a special case, namely for arithmetic circuits representing network polynomials of BNs over discrete RVs [39].

A direct corollary of Theorem 2 is that MPE inference is tractable in augmented SPNs, since augmented SPNs are

1: **procedure** MPESELECTIVE(S, **e**) Initialize zero-vector \mathbf{x}^* of length $|\mathbf{X}|$ 2: Evaluate e in corresponding MPN \hat{S} (upwards pass) 3: $Q \succ$ root node of MPN 4: 5: while Q not empty do $\hat{\mathsf{N}} \bowtie Q$ 6: if \hat{N} is a max node then 7: $Q \sim \arg \max \left\{ w_{\hat{\mathsf{N}},\hat{\mathsf{C}}} \mathsf{C} \right\}$ 8: $\hat{C} \in \mathbf{ch}(\hat{N})$ 9: else if N is a product node then $\forall \hat{\mathsf{C}} \in \mathbf{ch}(\hat{\mathsf{N}}) : Q \backsim \hat{\mathsf{C}}$ $10 \cdot$ else if \hat{N} is a maximizing distribution node then 11: $N \leftarrow$ corresponding distribution node 12: $\mathbf{x}^*[\mathbf{sc}(\mathsf{N})] = \arg \max \mathsf{N}(\mathbf{x})$ 13. $\mathbf{x} \in \mathbf{e}[\mathbf{sc}(\mathsf{N})]$ end if $14 \cdot$ end while $15 \cdot$ return x* 16: 17: end procedure

Fig. 7. Pseudo-code for MPE inference in selective SPNs.

selective SPNs over X and Z. This can easily be seen in AUG-MENTSPN, as for any z and any sum S, exactly one IV of Z_S is set to 1, causing that at most one child of S (or \overline{S}) can be non-zero. Therefore, we can use MPESELECTIVE in augmented SPNs, in order to find an MPE solution over *both* model RVs and LVs. Note that an MPE solution for the augmented SPN does in general *not* correspond to an MPE solution for the original SPN, when discarding the states of the LVs. However, this procedure is a frequently used approximation for models where MPE is tractable for both model RVs and LVs, but not for model RVs alone.

In [1], MPESELECTIVE was applied to *original* SPNs, not to *augmented* SPNs, but also with the goal to recover an MPE solution over both model RVs and LVs. The states of the LVs were assigned during max-backtracking, as sum-children and LV states are in one-to-one correspondence. The states of the LVs whose sums are *not visited* during backtracking, are not assigned—again, this causes some confusion, since some LVs appear to be undefined in some contexts, cf. the illustrations in Section 2. However, since this algorithm was used as approximation for MPE over model RVs by discarding the states of the LVs, this situation was not paid any further attention.

Nevertheless, as we show here, applying MPESELEC-TIVE to original (non-selective) SPNs effectively "simulates" MPESELECTIVE in the corresponding augmented SPN. Thereby, however, deterministic twin-weights are implicitly assumed, i.e., twin-weights which are 0, except a single 1. To see this, let us modify MPESELECTIVE, such that it can be applied to an original SPN, but returning an MPE solution for the corresponding augmented SPN. First note that in the augmented MPN, every twin node simply outputs the maximal twin-weight among all children whose states are contained in evidence e. For twin node S, let this maximal weight be denoted by $\hat{w}_{\bar{S}}$. The effect of the twin nodes can now be simulated in the original SPN by replacing each weight $w_{S,C}$ in the original SPN by $w_{S,C} \times \tilde{w}_{S,C}$. Here $\tilde{w}_{S,C}$ is a correction factor and given as $\tilde{w}_{S,C} = \prod_{\bar{S}} \hat{w}_{\bar{S}}$, where the product runs over all twins of those sums for which S is a conditioning sum. By using



Fig. 8. Illustration of the low-depth bias using an SPN over RVs $\{X_1, X_2, X_3\}$. The structure introduced by augmentation is depicted by small nodes and edges. When deterministic twin-weights are used, the state of Z_{S^1} corresponding to P¹ is preferred over P² and P³, since their probabilities are "dampened" by the weights of S² and S³, respectively.

these corrected weights, each max node in the corresponding MPN gets the same input as in the MPN of the augmented SPN, i.e., the twin nodes are simulated. We can identify the maximizing states of those LVs whose sums are visited during backtracking, as in [1]. The states of the sums which are not visited are given by the child which correspond to the maximal twin-weight $\hat{w}_{\tilde{S}}$. Pseudo-code for this somewhat technical modification of MPESELECTIVE can be found in [19].

We see that the algorithm used in [1] is essentially equivalent to MPESELECTIVE in augmented SPNs when $\tilde{w}_{S,C} = 1$ for all sum nodes, which implies that the twin-weights are deterministic. Therefore, although the LV model in [1] is not completely specified and it was not shown that the Viterbilike algorithm recovers an MPE solution, it nevertheless corresponds to MPE inference in the augmented SPN for special twin-weights, i.e., deterministic weights.

However, using deterministic twin-weights is a rather unnatural choice, since this prefers one arbitrary state over the others in cases where this LV is actually "rendered irrelevant". In this case, MPE inference also has a bias towards less structured sub-models, which we call low-depth bias. This is illustrated in Fig. 8, which shows an SPN over three RVs X_1, X_2, X_3 . The augmented SPN has two twin sum nodes \bar{S}^2 and \bar{S}^3 , corresponding to S^2 and S^3 , respectively. When their twin-weights are deterministic, the selection of the state of Z_{S^1} is *biased* towards the state corresponding to P^1 , which is a distribution assuming independence among X_1 , X_2 and X_3 . This comes from the fact, that the values of P^2 and P^3 are dampened by the weights of S^2 and S^3 , respectively, which are generally smaller than 1. Therefore, when using deterministic weights for twin sum nodes, we introduce a bias towards the selection of sub-SPNs that are less deep and less structured. Using uniform weights for twin sum nodes is somewhat "fairer", since in this case P^1 gets dampened by \overline{S}^2 and \overline{S}^3 , P^2 by S^2 and \overline{S}^3 , and P^3 by \overline{S}^2 and S^3 . Uniform weights are to some extend the opposite choice to deterministic twin-weights: the former represent

the strongest possible dampening via twin-weights and therefore actually *penalize* less structured distributions. Investigating these effects further is subject to future work.

5 EXPERIMENTS

5.1 Experiments with EM Algorithm

In [1], [40] SPNs were applied to image data, where a generic architecture reminiscent to convolutional neural networks was proposed. We refer to this architecture as PD architecture. Standard EM was not used in experiments for two reasons: First, explicitly constructing the proposed structure and to train it with standard EM is hardly possible with current hardware, since the number of nodes grows $O(l^3)$, where *l* is the square-length of the modeled image domain in pixels [19]. Instead, a sparse hard EM algorithm was used, which virtualizes the PD structure, i.e., sum and products are generated on the fly (see [40] for details). Second, using standard EM seemed unsuited to train large and dense SPNs, either because it is trapped in local optima or due to the gradient vanishing phenomenon.

In our experiments,² we investigated three questions:

- 1) Is our derivation of EM correct, both for complete and missing data?
- 2) Can the result of hard EM [1] be improved by standard EM?
- 3) Given a suited sparse structure, does EM yield a good solution for parameters?

Question 1) is important since the original derivation contained an error. Questions 2) and 3) are concerned with the general applicability of EM for training SPN.

We used the same datasets and SPN structures as in [1], obtainable from [40]. The datasets comprise Caltech-101 (inclusive background class) [43] and the ORL face images [44], i.e., in total 103 datasets. The input distributions in these SPNs are single-dimensional Gaussians (four for each pixel), where means were set to the averages of the fourquantiles and variances were constantly 1. We ran EM (Fig. 6) for 30 iterations, with various settings:

- Update any combination of the three different types of parameters, i.e., sum-weights, Gaussian means and Gaussian variances. Each set of parameters types is encoded by a string of letters W (weights), M (means) and V (variances). (seven combinations)
- Use original parameters for initialization, obtained from [40], or use three random initialization, where sum-weights are drawn from a Dirichlet distribution with uniform $\alpha = 1$ hyper-parameter (i.e., uniform distribution on the standard simplex), Gaussian means are uniformly drawn from [-1, 1] and Gaussian variances from [0.01, 1]. Only parameters which are actually updated are initialized randomly; otherwise the original parameters [1] are used and kept fixed. (four combinations)
- Use complete data or missing training data, randomly discarding 33 or 66 percent of the observations, independently for each sample. (three combinations)



Fig. 9. Normalized log-likelihood over EM-iterations, averaged over all 103 datasets and three random initializations. (a): Training set. (b): Test set; Curves for V and WV are outside the displayed region, for better readability of the other curves. They start at approximately -8,000 nats and decreased to approximately -11,000 nats.

Thus, in total we ran EM $7 \times 4 \times 3 \times 103 = 8,652$ times, yielding 259,560 EM-iterations. To avoid pathological solutions we used a lower bound of 0.01 for the Gaussian variances. In *no iteration* we observed a decreasing likelihood on the training set,³ i.e., our derived EM algorithm showed monotonicity in our experiments. Moreover, as can be seen in Fig. 9a, the training log-likelihood actually increased over iterations. The curves for the missing data scenarios are similar. This gives affirmative evidence for question 1).

Fig. 9b shows the log-likelihood on the test set. Note that optimizing the parameter sets V and WV led to severe overfitting: while achieving extremely high likelihoods on the training set, they achieved extremely poor likelihoods on the test set. Also the parameter sets MV and WMV tend to overfit, although not as strong as V and WV.

Regarding question 2), we closer inspected the test loglikelihood when the original parameters are used for initialization, i.e., when the parameters obtained by [40] are post-trained using EM. Table 1 summarizes the results. When parameter sets not including Gaussian variances are optimized (i.e., W, M, and WM), the test log-likelihood increased most of the time, i.e., for 83.5 percent (M) to up to 92.23 percent (WM) of the datasets. Furthermore, having oracle knowledge about the ideal number of iterations (i.e., column best), the average log-likelihood increased by 0.58 percent (M) to up to 1.39 percent (WM) relative to the original parameters. Most of this improvement happens in the first iteration, yielding 0.52 percent (M) up to 1.05 percent

2. Code available under [41], and [42].

3. Except for tiny occasional decreases (always $<10^{-8}$) after EM had converged, which can be attributed to numerical artifacts.

TABLE 1 Changes in Test Log-Likelihoods When Original Parameters Are Post-Trained Using EM

		after	after 1st iteration			best		
	% inc.	% all.	% pos.	% neg.	% all	% pos.	% neg.	
W	91.26	0.55	0.61	-0.03	0.87	0.96	-0.03	
Μ	83.50	0.52	0.67	-0.21	0.58	0.73	-0.21	
WM	92.23	1.06	1.18	-0.30	1.39	1.53	-0.30	
V	39.81	-13.47	14.44	-31.93	-13.45	14.51	-31.93	
WV	39.81	-13.41	14.79	-32.06	-13.33	14.98	-32.06	
MV	38.83	-17.24	14.27	-37.25	-17.21	14.35	-37.25	
WMV	38.83	-17.18	14.63	-37.37	-17.12	14.78	-37.37	

% inc.: percentage of datasets where log-likelihood increased in the first iteration. % all, % pos., % neg.: relative change of log-likelihood, averaged over all datasets, datasets with positive change, datasets with negative change, respectively.

(WM) improvement. These results indicate that the parameters obtained by [40] slightly underfit the given datasets. Similar as in Fig. 9, we see that parameter sets including the Gaussian variances (V, WV, MV, WMV) are prone to overfitting: more than 60 percent of the datasets decreased their test log-likelihood during EM. However, in the remaining 40 percent of the datasets, the test log-likelihood could be improved *substantially* by at least 14 percent on average.

We now turn to question 3). As pointed out above, a hard EM variant was used in [1], [40] which at the same time finds the effective SPN structure. Optimizing W using the three random initialization amounts to using the oracle structure obtained by [1], [40], discarding the learned parameters. For each dataset we selected the random initialization which yielded the highest likelihood on the training set in iteration 30. For this run, we compared the log-likelihoods with the log-likelihoods obtained by the original parameters. The results are summarized in Table 2.

We see that on all data sets the log-likelihood on the training set is larger than for the original parameters. This is also the case for each individual random start (not just best one)—every random restart always yielded a higher training log-likelihood than the original parameters. Thus, by considering the actual optimization objective—the likelihood on the training set—EM successfully trains SPNs, given a suited oracle structure. Furthermore, as can be seen in Table 2, EM is also not more prone to overfitting than the algorithm in [1]: on 67.96 percent of the datasets, EM delivered a higher test log-likelihood than the original parameters, when using oracle knowledge about the ideal number of iterations (column best).

TABLE 2 Log-Likelihoods When Sum-Weights (W) Are Trained, Using Random Initialization

	after 1st iteration				best			
	%>	% all.	% pos.	% neg.	% >	% all	% pos.	% neg.
train test	70.87 41.75	$0.68 \\ -0.11$	1.38 0.40	$-1.00 \\ -0.48$	100.00 67.96	3.97 0.46	3.97 0.76	-0.18

% >: percentage of data sets, where log-likelihood is larger than for original parameters. % all, % pos., % neg.: relative log-likelihood w.r.t. original parameters, for all data sets, data sets where relative log-likelihood is positive/negative, respectively.

TABLE 3 Differences of Log-Likelihood to the Ground-Truth MPE Solution Found by Exhaustive Enumeration, Averaged over 100 Independent Draws of Sum-Weights

		MPEDET	MPEUNI
4 RVs	$\begin{aligned} \alpha &= 0.5\\ \alpha &= 1.0\\ \alpha &= 2.0 \end{aligned}$	0.00 (100) 0.00 (100) 0.00 (100)	0.00 (100) 0.00 (100) 0.00 (100)
9 RVs	$\begin{aligned} \alpha &= 0.5\\ \alpha &= 1.0\\ \alpha &= 2.0 \end{aligned}$	-0.10(70) -0.10(68) -0.11(62)	0.00 (100) 0.00 (100) 0.00 (100)
16 RVs	$\begin{aligned} \alpha &= 0.5\\ \alpha &= 1.0\\ \alpha &= 2.0 \end{aligned}$	-0.63 (19) -0.85 (12) -0.82 (12)	0.00 (100) 0.00 (100) 0.00 (100)

Numbers in parentheses are the number of times where an MPE solution was found. Results for augmented SPNs using uniform twin-weights.

5.2 Experiments with MPE Inference

To illustrate correctness of MPESELECTIVE (Fig. 7) when applied to augmented SPNs, we generated SPNs using the PD architecture [1], arranging 4, 9 and 16 binary RVs in a 2×2 , 3×3 and 4×4 grid, respectively. As inputs we used two indicator variables for each RV representing their two states. The sum-weights were drawn from a Dirichlet distribution with uniform α -parameters, where $\alpha \in \{0.5, 1, 2\}$. For all networks we drew 100 independent parameters sets. We ran MPESELECTIVE on the augmented SPN, once equipped with uniform twin-weights and once with deterministic twin-weights. For uniform twin-weights, we denote the result obtained by MPESELECTIVE as MPEUNI. For deterministic twin-weights, we denote the result as MPEDET. As described in Section 4, MPEDET corresponds essentially to the result when MPESELECTIVE is applied to the original SPN [1]. For each assignment, the loglikelihoods were evaluated in the augmented SPN with deterministic weights, the augmented SPN with uniform weights and in the original SPN (discarding the states of the LVs). Additionally, we found ground truth MPE assignments in the two augmented SPNs and the original SPN using exhaustive enumeration. The results relative to the ground truth MPE solutions are shown in Tables 3, 4, and 5. As can be seen, MPEUNI always finds an MPE solution in the augmented SPN with uniform twin-weights and

TABLE 4 Similar as in Table 3

		MPEDET	MPEUNI
4 RVs	$\begin{array}{l} \alpha = 0.5 \\ \alpha = 1.0 \\ \alpha = 2.0 \end{array}$	0.00 (100) 0.00 (100) 0.00 (100)	0.00 (100) 0.00 (100) 0.00 (100)
9 RVs	$\begin{array}{l} \alpha = 0.5 \\ \alpha = 1.0 \\ \alpha = 2.0 \end{array}$	0.00 (100) 0.00 (100) 0.00 (100)	-0.10 (70) -0.12 (68) -0.15 (62)
16 RVs	$\begin{aligned} \alpha &= 0.5\\ \alpha &= 1.0\\ \alpha &= 2.0 \end{aligned}$	0.00 (100) 0.00 (100) 0.00 (100)	$\begin{array}{c} -0.89(19)\\ -1.11(12)\\ -1.01(12)\end{array}$

Results for augmented SPNs using deterministic twin-weights.

Similar as in Table 3				
		MPEDET	MPEUNI	
4 RVs	$\begin{aligned} \alpha &= 0.5\\ \alpha &= 1.0\\ \alpha &= 2.0 \end{aligned}$	-0.06 (72) -0.09 (59) -0.10 (52)	-0.06 (72) -0.09 (59) -0.10 (52)	
9 RVs	$\begin{aligned} \alpha &= 0.5\\ \alpha &= 1.0\\ \alpha &= 2.0 \end{aligned}$	-0.31 (32) -0.47 (12) -0.40 (6)	-0.38 (27) -0.48 (12) -0.37 (7)	
16 RVs	$\begin{aligned} \alpha &= 0.5\\ \alpha &= 1.0\\ \alpha &= 2.0 \end{aligned}$	-0.76 (5) -0.76 (3) -0.67 (1)	-1.04 (4) -1.18 (2) -0.92 (0)	

TABLE 5

Results for original SPNs.

MPEDET always finds an MPE solution in augmented SPNs with deterministic twin-weights. This gives empirical evidence for the correctness of MPESELECTIVE for MPE inference in augmented SPNs.

Furthermore, we wanted to investigate the quality of both algorithms when serving as approximation for MPE inference in the original SPNs. For the SPNs considered here, MPEDET delivered on average slightly better approximations than MPEUNI. However, these results should be interpreted with caution, due to the rather similar nature of the distributions considered here. Closer investigating approximate MPE for (original) SPNs is an interesting direction and will be subject to future research.

6 CONCLUSION

In this paper we revisited the interpretation of SPNs as hierarchically structured LV models. We pointed out that the original approach to explicitly incorporate LVs does not produce a sound probabilistic model. As a remedy we proposed the augmentation of SPNs and proved its soundness as LV model. Within augmented SPNs, we investigated the independency structure represented as BN, and showed that the sum-weights can be interpreted as structured CPTs within this BN. Using augmented SPNs, we derived the EM algorithm for sum-weights and single-dimensional input distributions from exponential families. While MPE-inference is generally NP-hard in SPNs, we showed that a Viterbi-style backtracking algorithm recovers an MPE solution in selective SPNs, and in particular in augmented SPNs. In experiments we give empirical evidence supporting our theoretical results. We furthermore showed that standard EM can successfully train generative SPNs, given a suitable network structure at hand.

APPENDIX A PROOFS

A.1 Proof of Proposition 1

If S' is a complete and decomposable SPN over $\mathbf{X} \cup \mathbf{Z}$, then $S'(\mathbf{X}) \equiv S(\mathbf{X})$ is immediate: Computing $S'(\mathbf{x})$ for any $\mathbf{x} \in \mathbf{val}(\mathbf{X})$ is done by marginalizing \mathbf{Z} , i.e., setting all $\lambda_{Z_{\mathbf{S}}=k} = 1$. In this case, it is easy to see that none of the structural changes modifies the output of the SPN, i.e., the outputs of S and S' agree for each \mathbf{x} , i.e., $S'(\mathbf{X}) \equiv S(\mathbf{X})$.

It remains to show that S' is complete and decomposable, and that the root's scope is $X \cup Z$. Steps 4–11 in AUG-MENTSPN introduce the links, representing "private copies" of the sum's children, and clearly leave the SPN complete and decomposable. In steps 13–15 the LV Z_S is introduced in the scope of S and thus in the scope of the root. Since this is done for all sum nodes, all Z are introduced in the root's scope. Steps 13–15 cannot render products non-decomposable, since this would imply that S is reachable by two distinct children of this product—a contradiction to the fact that the SPN was decomposable before. However, as shown in Fig. 1, steps 13-15 can render ancestor sums incomplete. These are treated in steps 17–23. The twin sum S, if introduced, is clearly complete and has scope $\{Z\}$. Furthermore, incompleteness of any conditioning sum S^c can only be caused by links not having $Z_{\rm S}$ in their scope. The scope of these links is augmented by $Z_{\rm S}$ in step 21. These links clearly remain decomposable and moreover, S^c is rendered complete now.

A.2 Proof of Proposition 2

ad 1.) When deleting the IVs and their links, the scopes of any (twin) sum remains the same, since it is complete and is left with one child. Thus also the scope of any ancestor remains the same.

ad 2.) The graph of S^{y} is rooted and acyclic, since the root cannot be a link and deleting nodes and edges cannot introduce cycles. When an IV $\lambda_{Y=y}$ is deleted, also the link $\mathsf{P}_{\mathsf{S}_{Y}}^{y}$ is deleted, so no internal nodes are left as leaves. The roots in S^{y} and S' are the same, and by point $1., \mathsf{X} \cup \mathsf{Y} \cup \mathsf{Z}$ is the scope of the root. S^{y} is also complete and decomposable: Whenever an IV and its link are deleted, the corresponding sum node and twin sum node remain trivially complete, since they are left with a single child. Furthermore, completeness and decomposability of any ancestor of S_{Y} or $\bar{\mathsf{S}}_{Y}$ is left intact, since neither S_{Y} nor $\bar{\mathsf{S}}_{Y}$ changes its scope.

ad 3.) According to point 1., the scope of N is the same in S' and S^{y} . Since $sc(N) \cap Y = \emptyset$, the disconnected IVs and deleted links are no descendants of N, i.e., no descendants of N are disconnected during configuration. Since N is present in S^{y} , it must still be reachable from the root. Therefore also all descendants of N are reachable, i.e., $S_{N}^{y} = S'_{N}$.

ad 4.) When the input is fixed to $\mathbf{x}, \mathbf{z}, \mathbf{y}$, all IVs and links which are deleted from the configured SPN $S^{\mathbf{y}}$ evaluate to zero in the augmented SPN S'. The outputs of all sums and twin sums are therefore the same in S' and $S^{\mathbf{y}}$. Therefore, also the output of all other nodes remains the same. This includes the root and therefore $S^{\mathbf{y}}(\mathbf{x}, \mathbf{z}, \mathbf{y}) = S'(\mathbf{x}, \mathbf{z}, \mathbf{y})$, for any \mathbf{x}, \mathbf{z} .

When $\mathbf{y}' \neq \mathbf{y}$, then there must be a $Y \in \mathbf{Y}$ such that the IV $\lambda_{Y=\mathbf{y}'[Y]}$ has been deleted, i.e., $\lambda_{Y=\mathbf{y}'[Y]} \notin \mathbf{desc}(\mathbf{N})$, where N is the root of $S^{\mathbf{y}}$. Using Lemma 1 in [17], it follows that $S^{\mathbf{y}}(\mathbf{x}, \mathbf{z}, \mathbf{y}') = 0$.

A.3 Proof of Lemma 1

 S^{z} must contain either **S** or \overline{S} , since Z_{S} is in the scope of the root by Proposition 2. To show that *not both* are in S^{z} , let Π_{k} denote the set of paths of length k from the root to any node N with $Z_{S} \in sc(N)$. For k > 1, all paths in Π_{k} can be constructed by extending each path in Π_{k-1} with each child of this path's

last node, if it has Z_S in its scope. Let K be the smallest number such that there is a path in Π_k containing S or \overline{S} .

We show by induction, that $|\Pi_k| = 1$, k = 1, ..., K. Note that Π_1 contains a single path (N), where N is the root, therefore the induction basis holds.

For the induction step, we show that given $|\Pi_{k-1}| = 1$, then also $|\Pi_k| = 1$. Let $(N_1, ..., N_{k-1})$ be the single path in Π_{k-1} . If N_{k-1} is a product node, then it has a single child C with $Z_{\mathbf{S}} \in \mathbf{sc}(\mathbf{C})$, due to decomposability. If N_{k-1} is a sum node, then it must be in $\mathbf{anc}_{\mathbf{S}}(\mathbf{S}) \setminus \{\mathbf{S}\}$, and therefore has a single child in the configured SPN. Therefore, there is a single way to extend the path and therefore $|\Pi_k| = 1$, k = 1, ..., K. This single path does either lead to S or \overline{S} . Since $\mathbf{S} \notin \mathbf{desc}(\overline{S})$ and $\overline{\mathbf{S}} \notin \mathbf{desc}(\mathbf{S})$, $S^{\mathbf{z}}$ contains a single path to one of them, but not to both.

A.4 Proof of Theorem 1

By Lemma 1, for each $\mathbf{z} \in \mathbf{val}(\mathbf{Z}_p)$ the configured SPN $S^{\mathbf{z}}$ contains either S or \overline{S} , but not both. Let \mathcal{Z} be the subset of $\mathbf{val}(\mathbf{Z}_p)$ such that S is in $S^{\mathbf{z}}$ and $\overline{\mathcal{Z}}$ be the subset of $\mathbf{val}(\mathbf{Z}_p)$ such that \overline{S} is in $S^{\mathbf{z}}$.

Fix $Z_{S} = k$ and $z \in \mathbb{Z}$. We want to compute $S'(Z_{S} = k, \mathbf{Y}_{n}, \mathbf{z})$, i.e., we marginalize \mathbf{Y}_{c} . According to Proposition 2 (4.), this equals $S^{\mathbf{z}}(Z_{S} = k, \mathbf{Y}_{n}, \mathbf{z})$. According to Proposition 2 (3.), the sub-SPN rooted at former child \mathbf{C}_{S}^{k} is the same in S' and $S^{\mathbf{z}}$. Since S' is locally normalized, this sub-SPN is also locally normalized in $S^{\mathbf{z}}$. Since the scope of the former child \mathbf{C}_{S}^{k} is a sub-set of \mathbf{Y}_{c} , which is marginalized, and $\lambda_{Z_{S}=k} = 1$, the link P_{S}^{k} outputs 1. Since $\lambda_{Z_{S}=k'} = 0$ for $k' \neq k$, the sum S outputs w_{k} .

Now consider the set of nodes in S^z which have Z_S in their scope, not including $\lambda_{Z_S=k}$ and P^k_S . Clearly, since \bar{S} is not in S^z , this set must be $\operatorname{anc}(S)$. Let $\mathsf{N}_1, \ldots, \mathsf{N}_L$ be a topologically ordered list of $\operatorname{anc}(S)$, where S is N_1 and N_L is the root. Let $\mathbf{Y}_{n,l} := \operatorname{sc}(\mathsf{N}_l) \cap \mathbf{Y}_n$ and $\mathbf{Z}_l := \operatorname{sc}(\mathsf{N}_l) \cap \mathbf{Z}_p$. We show by induction that for $l = 1, \ldots, L$, we have

$$\mathsf{N}_l(Z_\mathsf{S} = k, \mathbf{Y}_{n,l}, \mathbf{z}[\mathbf{Z}_l]) = w_k \,\mathsf{N}_l(\mathbf{Y}_{n,l}, \mathbf{z}[\mathbf{Z}_l]). \tag{30}$$

Since $\mathbf{Y}_{n,1} = \emptyset$ and $\mathbf{Z}_1 = \emptyset$, and $\mathsf{N}_1(Z_\mathsf{S} = k) = w_k$, the induction basis holds. Assume that (30) holds for all $\mathsf{N}_1, \ldots, \mathsf{N}_{l-1}$. If N_l is a sum, we have due to completeness

$$\mathsf{N}_{l}(Z_{\mathsf{S}} = k, \mathbf{Y}_{n,l}, \mathbf{z}[\mathbf{Z}_{l}]) = \sum_{\mathsf{C} \in \mathsf{ch}(\mathsf{N}_{l})} w_{\mathsf{N}_{l},\mathsf{C}} w_{k} \,\mathsf{C}(\mathbf{Y}_{n,l}, \mathbf{z}[\mathbf{Z}_{l}]) \quad (31)$$

$$= w_k \,\mathsf{N}_l(\mathbf{Y}_{n,l}, \mathbf{z}[\mathbf{Z}_l]), \tag{32}$$

i.e., the induction step holds for sums. When N_l is a product, due to decomposability, it must have a single child with Z_S in its scope. Hence, this child must be a node $N_m \in anc(S)$ We have

$$\mathsf{N}_l(Z_\mathsf{S} = k, \mathbf{Y}_{n,l}, \mathbf{z}[\mathbf{Z}_l]) \tag{33}$$

$$= w_k \mathsf{N}_m(\mathbf{Y}_{n,m}, \mathbf{z}[\mathbf{Z}_m]) \prod_{\mathsf{C} \in \mathbf{ch}(\mathsf{N}_l) \setminus \mathsf{N}_m} \mathsf{C}(\mathbf{Y}_{n,l} \cap \mathbf{sc}(\mathsf{C}))$$
(34)

$$= w_k \mathsf{N}_l(\mathbf{Y}_{n,l}, \mathbf{z}[\mathbf{Z}_l]), \tag{35}$$

i.e., the induction step holds for products. Therefore, by induction, (30) also holds for the root, and (11) follows.

Now we show (12). If the twin sum \overline{S} does not exist, \overline{Z} is empty and (12) holds trivially. Otherwise, fix the input to $Z_{S} = k$ and $z \in \overline{Z}$. Clearly, \overline{S} outputs \overline{w}_{k} and (12) can be shown in similar way as (11).

A.5 Proof of Theorem 2

We prove the theorem using an inductive argument. The theorem clearly holds for any \hat{D} by definition. Consider a product \hat{P} and assume the theorem holds for all $ch(\hat{P})$. Then the theorem also holds for \hat{P} , since

$$\hat{\mathsf{P}}(\boldsymbol{\mathcal{X}}) = \prod_{\mathsf{C}\in ch(\hat{\mathsf{P}})} \max_{\mathbf{x}\in\boldsymbol{\mathcal{X}}} \mathsf{C}(\mathbf{x}) = \max_{\mathbf{x}\in\boldsymbol{\mathcal{X}}} \prod_{\mathsf{C}\in ch(\hat{\mathsf{P}})} \mathsf{C}(\mathbf{x}) = \max_{\mathbf{x}\in\boldsymbol{\mathcal{X}}} \mathsf{P}(\mathbf{x}),$$
(36)

where the max and the product can be switched due to decomposability.

Now consider a max node \hat{S} and its corresponding sum node S. Let the *support* of an SPN-node N be the set $\sup_N := \{x \mid N(x) > 0\}$. Since S is selective, its support is partitioned by the supports of its children, i.e., $\sup_S = \bigcup_{C \in ch(S)} \sup_{C'}$, $\sup_{C'} \bigcap \sup_{C''} = \emptyset$, for $C' \neq C''$. Assuming that the theorem holds for all $ch(\hat{S})$, we have

$$\hat{\mathsf{S}}(\mathcal{X}) = \max_{\mathsf{C} \in \mathsf{ch}(\mathsf{S})} w_{\mathsf{S},\mathsf{C}} \max_{\mathsf{x} \in \mathcal{X}} \mathsf{C}(\mathsf{x})$$
(37)

$$= \max_{\mathbf{C} \in \mathbf{ch}(\mathbf{S})} w_{\mathbf{S},\mathbf{C}} \max_{\mathbf{x} \in \sup_{\mathbf{C}} \cap \mathcal{X}} \mathbf{C}(\mathbf{x})$$
(38)

$$= \max_{\mathbf{C} \in \mathbf{ch}(\mathbf{S})} \max_{\mathbf{x} \in \sup_{\mathbf{C}} \cap \mathcal{X}} w_{\mathbf{S},\mathbf{C}} \mathbf{C}(\mathbf{x})$$
(39)

$$= \max_{\mathbf{x} \in \sup_{\mathbf{S}} \cap \mathcal{X}} \mathbf{S}(\mathbf{x}) = \max_{\mathbf{x} \in \mathcal{X}} \mathbf{S}(\mathbf{x}).$$
(40)

In (38) we have a slight abuse of notation, as we actually should use suprema over the sets $\sup_{C} \cap \mathcal{X}$ and define the supremum over the empty set as 0. In (39) we used the fact that the support of the sum node is partitioned by the supports of its children and that for selective sums we have $S = w_{S,C} C$ whenever we have single child with C > 0.

We see that the induction step also holds for S. Therefore, the theorem holds for all nodes.

ACKNOWLEDGMENTS

We would like to thank the anonymous reviewers for their constructive comments. This work was supported by the Austrian Science Fund (FWF): P25244-N15 and Austrian Science Fund (FWF): P27803-N15. This research was partly funded by ONR grant N00014-16-1-2697 and AFRL contract FA8750-13-2-0019.

REFERENCES

- H. Poon and P. Domingos, "Sum-product networks: A new deep architecture," in *Proc. Uncertainty Artif. Intell.*, 2011, pp. 337–346.
- [2] R. Gens and R. Domingos, "Learning the structure of sum-product networks," in Proc. 30th Int. Conf. Mach. Learn., 2013, pp. 873–880.

- [3] A. Dennis and D. Ventura, "Learning the architecture of sumproduct networks using clustering on variables," in *Proc. 25th Int. Conf. Neural Inf. Process. Syst.*, 2012, pp. 2042–2050.
 [4] R. Peharz, B. Geiger, and F. Pernkopf, "Greedy part-wise learning
- [4] R. Peharz, B. Geiger, and F. Pernkopf, "Greedy part-wise learning of sum-product networks," in *Proc. Eur. Conf. Mach. Learn. Knowl. Discovery Databases*, 2013, pp. 612–627.
 [5] M. Amer and S. Todorovic, "Sum-product networks for modeling
- [5] M. Amer and S. Todorovic, "Sum-product networks for modeling activities with stochastic structure," in *Proc. Comput. Vis. Pattern Recognit.*, 2012, pp. 1314–1321.
- [6] R. Gens and P. Domingos, "Discriminative learning of sumproduct networks," in Proc. Advances Neural Inf. Process. Syst., 2012, pp. 3248–3256.
- [7] R. Peharz, G. Kapeller, P. Mowlaee, and F. Pernkopf, "Modeling speech with sum-product networks: Application to bandwidth extension," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, 2014, pp. 3699–3703.
- [8] W. C. Cheng, S. Kok, H. V. Pham, H. L. Chieu, and K. M. A. Chai, "Language modeling with sum-product networks," in *Proc. INTERSPEECH*, 2014, pp. 2098–2102.
- [9] M. Zöhrer, R. Peharz, and F. Pernkopf, "Representation learning for single-channel source separation and bandwidth extension," *IEEE/ACM Trans. Audio Speech Language Process.*, vol. 23, no. 12, pp. 2398–2409, Dec. 2015.
- [10] A. Dempster, N. Laird, and D. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," J. Royal Statistical Soc. Series B, vol. 39, no. 1, pp. 1–38, 1977.
- [11] Z. Ghahramani and M. Jordan, "Supervised learning from incomplete data via an EM approach," in *Proc. Advances Neural Inf. Process. Syst.*, 1994, pp. 120–127.
- [12] S.-W. Lee, C. Watkins, and B. Zhang, "Non-parametric Bayesian sum-product networks," in *Proc. Workshop Learn. Tractable Probabilistic Models*, (2014). [Online]. Available: https://sites.google. com/site/ltpm2014/accepted-papers
- [13] M. Trapp, R. Peharz, M. Skowron, T. Madl, F. Pernkopf, and R. Trappl, "Structure inference in sum-product networks using infinite sum-product trees," in *Proc. NIPS Workshop Practical Bayesian Nonparametrics*, 2016.
- [14] J. Pearl, Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference. San Francisco, CA, USA: Morgan Kaufmann, 1988.
- [15] D. Koller and N. Friedman, Probabilistic Graphical Models: Principles and Techniques. Cambridge, MA, USA: MIT Press, 2009.
- [16] A. Darwiche, "A differential approach to inference in Bayesian networks," J. ACM, vol. 50, no. 3, pp. 280–305, 2003.
- [17] R. Peharz, S. Tschiatschek, F. Pernkopf, and P. Domingos, "On theoretical properties of sum-product networks," in *Proc. 18th Int. Conf. Artif. Intell. Statist.*, 2015, pp. 744–752.
- [18] C. de Campos, "New complexity results for MAP in Bayesian networks," in *Proc. 22nd Int. Joint Conf. Artif. Intell.*, 2011, pp. 2100–2106.
 [19] R. Peharz, "Foundations of sum-product networks for probabilis-
- [19] R. Peharz, "Foundations of sum-product networks for probabilistic modeling," Ph.D. dissertation, Signal Processing and Speech Communication Lab, Graz Univ. Technol., Graz, Austria, 2015.
- [20] R. Peharz, R. Gens, and P. Domingos, "Learning selective sumproduct networks," in *Proc. ICML Workshop Learn. Tractable Probabilistic Models*, (2014). [Online]. Available: https://sites.google. com/site/ltpm2014/accepted-papers
- [21] H. Zhao, M. Melibari, and P. Poupart, "On the relationship between sum-product networks and Bayesian networks," in *Proc.* 32nd Int. Conf. Mach. Learn., 2015, pp. 116–124.
- [22] A. Darwiche, "Compiling knowledge into decomposable negation normal form," in *Proc. 16th Int. Joint Conf. Artif. Intell.*, 1999, pp. 284–289.
- [23] A. Darwiche, "Decomposable negation normal form," J. ACM, vol. 48, no. 4, pp. 608–647, 2001.
- [24] A. Darwiche and P. Marquis, "A knowledge compilation map," J. Artif. Intell. Res., vol. 17, no. 1, pp. 229–264, 2002.
- [25] A. Darwiche, "A logical approach to factoring belief networks," in Proc. 8th Int. Conf. Principles Knowl. Representation Reasoning, 2002, pp. 409–420.
- [26] D. Lowd and P. Domingos, "Learning arithmetic circuits," in Proc. 24th Conf. Uncertainty Artif. Intell., 2008, pp. 383–392.
- [27] C. Boutilier, N. Friedman, M. Goldszmidt, and D. Koller, "Context-specific independence in Bayesian networks," in *Proc.* 12th Int. Conf. Uncertainty Artif. Intell., 1996, pp. 115–123.
 [28] D. Lowd and A. Rooshenas, "Learning Markov networks with
- [28] D. Lowd and A. Rooshenas, "Learning Markov networks with arithmetic circuits," in Proc. 16th Int. Conf. Artif. Intell. Statist., 2013, pp. 406–414.

- [29] A. Rooshenas and D. Lowd, "Learning sum-product networks with direct and indirect variable interactions," in *Proc. 31st Int. Conf. Mach. Learn.*, 2014, pp. 710–718.
- [30] T. Ádel, D. Balduzzi, and A. Ghodsi, "Learning the structure of sum-product networks via an SVD-based algorithm," in Proc. Uncertainty Artif. Intell., 2015, pp. 32–41.
- [31] A. Vergari, N. Di Mauro, and F. Esposito, "Simplifying, regularizing and strengthening sum-product network structure learning," in *Proc. Eur. Conf. Mach. Learn. Knowl. Discovery Databases*, 2015, pp. 343–358.
- [32] O. Delalleau and Y. Bengio, "Shallow vs. deep sum-product networks," in Proc. Advances Neural Inf. Process. Syst. 24, 2011, pp. 666–674.
- [33] A. Friesen and P. Domingos, "The sum-product theorem: A foundation for learning tractable models," in *Proc. 33rd Int. Conf. Mach. Learn.*, 2016, pp. 1909–1918.
- [34] N. Johnson, S. Kotz, and N. Balakrishnan, Continuous Univariate Distributions. Hoboken, NJ, USA: Wiley, 1994.
- [35] H. Zhao and P. Poupart, "A unified approach for learning the parameters of sum-product networks," (2016). [Online]. Available: http://arxiv.org/abs/1601.00318
- [36] H. Bodlaender, F. van den Eijkhof, and L. van der Gaag, "On the complexity of the MPA problem in probabilistic networks," in *Proc. 15th Eur. Conf. Artif. Intell.*, 2002, pp. 675–679.
- [37] J. D. Park and A. Darwiche, "Complexity results and approximation strategies for MAP explanations," J. Artif. Intell. Res., vol. 21, no. 1, pp. 101–133, 2004.
- [38] J. Kwisthout, "Most probable explanations in Bayesian networks: Complexity and tractability," *Int. J. Approximate Reasoning*, vol. 52, no. 9, pp. 1452–1469, 2011.
- [39] A. Darwiche, Modeling and Reasoning with Bayesian Networks. Cambridge, U.K.: Cambridge Univ. Press, 2014.
- [40] (2011). [Online]. Available: http://alchemy.cs.washington.edu/ spn
- [41] (2016). [Online]. Available: http://spn.cs.washington.edu/pubs. shtml
- [42] (2016). [Online]. Available: https://www.spsc.tugraz.at/tools
- [43] L. Fei-Fei, R. Fergus, and R. Perona, "Learning generative visual models from few training examples: An incremental Bayesian approach tested on 101 object categories," *Comput. Vis. Image Understanding*, vol. 106, no. 1, pp. 59–70, 2007.
- [44] F. Samaria and A. Harter, "Parameterisation of a stochastic model for human face identification," in *Proc. 2nd IEEE Workshop Appl. Comput. Vis.*, 1994, pp. 138–142.



Robert Peharz received the MSc degree in computer engineering and the PhD degree in electrical engineering from Graz University of Technology. His main research interest lies in machine learning, in particular probabilistic modeling, with applications to signal processing, speech and audio processing, and computer vision. Currently, he is with the research unit interdisciplinary developmental neuroscience (iDN), Medical University of Graz, applying machine learning techniques to detect early markers of neurological conditions in

infants. He is funded by the BioTechMed-Graz cooperation, an interdisciplinary network of the three major universities in Graz with a focus on basic bio-medical research, technological development, and medical applications.



Robert Gens received the SB degree in electrical engineering and computer science from the Massachusetts Institute of Technology, Cambridge, Massachusetts, in 2009, and the MSc degree in computer science and engineering from the University of Washington, Seattle, Washington, in 2012. He is currently working toward the PhD degree in computer science and engineering at the University of Washington, Seattle, Washington. During the Summer of 2014, he was a research intern with Microsoft Research, Red-

mond, Washington. He is supported by the 2014 Google Ph.D. Fellowship in Deep Learning. He received an Outstanding Student Paper Award at the Neural Information Processing Systems conference in 2012.



Franz Pernkopf received the MSc (Dipl. Ing.) degree in electrical engineering from Graz University of Technology, Austria, in summer 1999. He received the PhD degree from the University of Leoben, Austria, in 2002. In 2002, he was awarded the Erwin Schrödinger Fellowship. He was a research associate in the Department of Electrical Engineering, University of Washington, Seattle, from 2004 to 2006. Currently, he is an associate professor in the Laboratory of Signal Processing and Speech Communication, Graz

University of Technology, Austria. His research interests include machine learning, discriminative learning, graphical models, feature selection, finite mixture models, and image- and speech processing applications. He is a senior member of the IEEE.



Pedro Domingos received the PhD degree from the University of California at Irvine. He is a professor of computer science with the University of Washington and the author of *The Master Algorithm*. He is a winner of the SIGKDD Innovation Award, the highest honor in data science. He has received a Fulbright Scholarship, a Sloan Fellowship, the National Science Foundations CAREER Award, and numerous best paper awards. He is the author or co-author of more than 200 technical publications. He has held visiting positions at

Stanford, Carnegie Mellon, and MIT. He co-founded the International Machine Learning Society in 2001. His research spans a wide variety of topics in machine learning, artificial intelligence, and data science, including scaling learning algorithms to big data, maximizing word of mouth in social networks, unifying logic and probability, and deep learning. He is a fellow of the Association for the Advancement of Artificial Intelligence.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.