

Random Probabilistic Circuits

Nicola Di Mauro¹

Gennaro Gala¹

Marco Iannotta¹

Teresa M.A. Basile^{2,3}

¹Department of Computer Science, University of Bari, Bari, Italy

²Department of Physics, University of Bari, Bari, Italy

³National Institute for Nuclear Physics (INFN), Bari Division, Bari, Italy

Abstract

Density estimation could be viewed as a core component in machine learning, since a good estimator could be used to solve many tasks such as classification, regression, and imputing missing values. The main challenge of density estimation is balancing the model expressiveness and its learning and inference complexity. Probabilistic circuits (PCs) model a probability distribution as a computational graph. By imposing specific structural properties on such models many inference tasks become tractable. However, learning PCs usually relies on greedy and time consuming procedures. In this paper we propose a new unified approach to efficiently learn PCs having several structural properties. We introduce extremely randomized PCs (XPCs), PCs with a random structure. We show their advantage on standard density estimation benchmarks when compared to other density estimators.

1 INTRODUCTION

Density estimation is the unsupervised task of reconstructing the joint probability density function (pdf)—learning an estimator—underlying a set of observed samples over specific random variables (RVs). Once such an estimator is learned, it can be used to make *inference*—computing the probability of *queries* about certain RVs states, such as complete evidence queries, marginal queries or conditional queries. A perfect estimator would allow to solve many classical machine learning tasks (e.g., regression, classification, clustering and unsupervised prediction) by casting them in specific types of inference. Indeed, density estimation is recognized as one of the most general and powerful task in Machine Learning (ML).

The main challenge of density estimation is balancing the

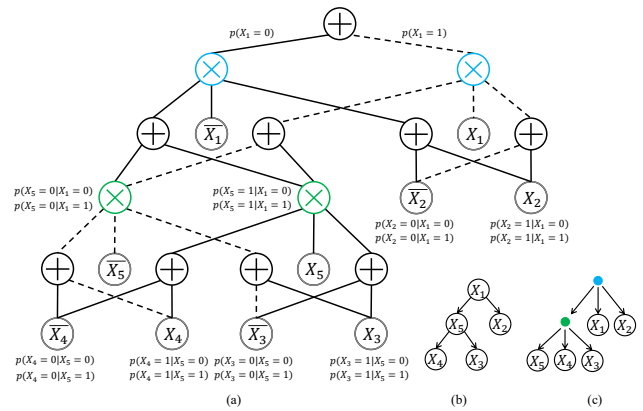


Figure 1: A PC, its vtree and corresponding CLT. A smooth, structured-decomposable, and deterministic PC over five RVs $\mathbf{X} = \{X_1, X_2, X_3, X_4, X_5\}$ in (a), where dotted edges denote the *flowing* for the input sample $\{X_1 = 1, X_2 = 0, X_3 = 0, X_4 = 1, X_5 = 0\}$. Its equivalent CLT in (b), and its vtree in (c). Conditional probabilities in (a) represent the weights of the left and right incoming edges.

model representation expressiveness (i.e., the ability to model complex distributions) against its learning and inference complexity. For instance, *Probabilistic Graphical Models* (PGMs), like Bayesian Networks (BNs), can model highly complex probability distributions but exact inference with them is generally intractable—non-polynomial in the size of the network. This is the reason why, a relatively recent probabilistic ML research area focuses on designing and exploiting models that can theoretically guarantee reliable and efficient probabilistic inference. These models belong to the family of *Tractable Probabilistic Models* (TPMs), compact representations for rich probability distributions which allow complex inference routines to be computed exactly and in polynomial time, i.e., *tractably*.

However, it should be noticed that tractability is associated to classes of queries: computing exact marginals on a TPM may be tractable, while MAP (Maximum A Posteriori)

may not be. TPMs like bounded tree-width graphical models [Meila and Jordan, 2000], Arithmetic Circuits [Darwiche, 2009], Sum-Product Networks (SPNs) [Poon and Domingos, 2011], Cut-Set Networks (CNETs) [Rahman et al., 2014] and Probabilistic Sentential Decision Diagrams (PSDDs) [Kisa et al., 2014], promise a good compromise between expressive power and tractable inference. The recently proposed *Probabilistic Circuits* (PCs) framework [Vergari et al., 2020, Choi et al., 2020] allow us to describe, learn, and reason about these TPMs and to categorize them under this single and unified framework.

PCs (see Figure 1) are computational graphs that define a joint probability distribution as recursive mixtures (sum units) and factorizations (product units) of simpler distributions (e.g., parametric distributions such as Gaussians or Bernoullis). Learning PCs can be naturally organized into *structure learning* and *parameter learning*, following the same dichotomy as in PGMs. Differently from other neural density estimators such as NADE [Uribe et al., 2016] and autoregressive flows [Papamakarios et al., 2017], PCs enable large classes of tractable inference with little or no compromise in terms of model expressiveness.

However, even if learning TPMs may be done in polynomial time, thanks to several recent algorithmic schemes, making these algorithms scale to high dimensional data is still an issue. In particular, various maximum likelihood based approaches have been proposed using either gradient-based optimisation or expectation-maximization (and related schemes), but unfortunately, most structure learners proposed so far are tedious to tune. For instance, tractable SPN learners as LEARNSPN [Gens and Domingos, 2013] and ID-SPN [Rooshenas and Lowd, 2014] spend a substantial amount of their execution time on partitioning RVs into conditionally independent components. On the contrary, Peharz et al. [2019] proposed RAT-SPN, a simple random approach to employ SPNs for deep learning, and demonstrated that tractable models like SPNs can get surprisingly far even without sophisticated structure learning. Also, Di Mauro et al. [2017] showed how mixtures of CNETs whose OR tree is learned by performing random conditioning outperform state-of-the-art density estimators on a series of standard benchmark datasets.

Following these trends, we introduce *eXtremely randomized Probabilistic Circuits* (XPCs), smooth and decomposable PCs which can be easily forced to be deterministic and/or structured decomposable without increasing learning complexity. We are proposing a new learning scheme for PCs that can be forced to build networks with specific structural constraints. XPCs leverage on random conditionings and employ naïve factorizations (*aka* PoB, Product of Bernoullis) and Chow-Liu Trees (CLTs) [Chow and Liu, 1968] as multivariate leaf distributions. To satisfy structured decomposability, we propose a simple yet effective approach to normalize several CLTs for the same vtree. While the

log-likelihood of an XPC is slightly worse than other state-of-the-art (SOTA) competitors, ensembles of XPC (EXPCs) perform as SOTA density estimators on a series of standard benchmark datasets, yet employing only a fraction of the time needed to learn the competitors.

2 RELATED WORKS

Randomly learning the structure of PCs is a problem already tackled in some recent works. The first kind of random PCs have been proposed in [Di Mauro et al., 2017] where the authors introduced random CNETs, a deterministic PC, whose OR tree is learned by performing random conditioning. Here we do not limit to single variable conditioning, and furthermore our proposed learning scheme is not limited to deterministic circuits as CNETs.

In [Peharz et al., 2019] the authors proposed a simple and scalable method to build random and tensorized SPNs (RAT-SPNs), avoiding the necessity of a structure learner. In particular, they first construct a random region graph [Dennis and Ventura, 2012] subsequently populated with arrays of SPN nodes, hence leading to a random hierarchical tensorial decomposition. For the generative case, the parameters of the network are learned using the classical expectation-maximization (EM) algorithm [Dempster et al., 1977] derived for SPNs as in [Peharz et al., 2017]. In this paper the construction of random region graph is data-driven and we avoid the time consuming EM algorithm for parameter optimization. Furthermore, our learning scheme can impose other structural constraints to the obtained PCs than that imposed on SPNs.

Finally, in [Ventola et al., 2019] the authors proposed random sum-product forests (RSPFs), an ensemble approach for mixing multiple randomly generated SPNs, with an approach similar to that proposed in [Di Mauro et al., 2017]. Even in this case the parameters of the networks are learned using the time consuming EM algorithm and the approach is restricted to SPNs.

3 PROBABILISTIC CIRCUITS

Notation. Upper-case letters are used to denote random variables (RVs), while lower-case ones for their assignments. Upper-case and lower-case bold letters are used, resp., to denote set of RVs and their joint values. In this paper we consider Boolean RVs, i.e., having values in $\{0, 1\}$.

3.1 STRUCTURAL PROPERTIES

A probabilistic circuit \mathcal{C} , defined over a set of RVs \mathbf{X} , is a computational graph, represented with a directed acyclic graph (DAG), encoding a probability distribution $P_{\mathcal{C}}(\mathbf{X})$.

PCs have three kinds of nodes only: input distributions (leaves), product nodes (\otimes -node) and sum nodes (\oplus -node). Given a node n of the DAG, let C_n denote the sub-circuit rooted at n , $ch(n)$ its child nodes, and P_{C_n} its encoding distribution. An input distribution n encodes a tractable probability distribution P_n over some RVs $\phi(n) \subseteq \mathbf{X}$. A product node defines a factorized distribution $P_n(\mathbf{X}) = \prod_{c \in ch(n)} P_c(\mathbf{X})$. A sum node n represents a mixture distribution $P_n(\mathbf{X}) = \sum_{c \in ch(n)} \theta_{nc} P_c(\mathbf{X})$. The RVs over which an input distribution n is defined is named its *scope* $\phi(n)$. The scope of a product or a sum node n is defined as $\phi(n) = \cup_{c \in ch(n)} \phi(c)$. In this paper we consider univariate input distributions.

A PC C over RVs \mathbf{X} supports computing the likelihood $p_C(\mathbf{x})$ given a complete configuration \mathbf{x} (a complete evidence query) by evaluating the circuit bottom up: starting from the input distributions and computing the output of children before parents. Additional properties of the PC extend the set of probabilistic queries that are guaranteed to be answered exactly and in time linear in the size of the PC—its number of edges.

Definition 3.1. A probabilistic circuit is *smooth* if for every sum node, the children have the same scope.

Definition 3.2. A probabilistic circuit is *decomposable* if for each \otimes -node its children have disjoint scope.

Smooth and decomposable PCs enable the tractable computation of any marginal query [Darwiche, 2000]. In this way, each node in a PC recursively defines a distribution over its scope: a) leaves are distribution by definition, b) \oplus -nodes are mixtures of their child distributions, and c) \otimes -nodes are factorized distributions.

Definition 3.3. A circuit is *deterministic* if for every sum node n and assignment \mathbf{x} , at most one of the children of n have a non-zero output.

In a deterministic circuit \oplus -nodes define a mixture model whose components have disjoint support, thus enabling tractable MAP inference [Chan and Darwiche, 2006] and closed-form parameter learning [Kisa et al., 2014].

Even if a vtree has been defined to be a binary tree [Kisa et al., 2014], its definition could be extended to general n-ary trees.

Definition 3.4. A *vtree* over a set of \mathbf{X} RVs is a n-ary tree encoding a hierarchical decomposition of the RVs. Each leaf in a vtree denotes a RV, while an internal node indicates how to decompose a set of RVs into many subsets.

A PC is said to be *normalized* for a vtree if the scope of every product node decomposes over its children as its corresponding node in the vtree (see Figure 1 for an example of a vtree and a corresponding normalized structured decomposable PC. Figure 3 in Appendix A is another example of structured decomposable PC).

A PC is *structured decomposable* (SD) if its product nodes cannot decompose in arbitrary ways, but must agree on a *vtree*. In particular, a structured-decomposable PC encodes a probability distribution in a computational graph by recursively decomposing it into smaller distributions according to a hierarchical partitioning of the random variables, also called vtree. A structured decomposable PC provides several classes of advanced probabilistic queries computable exactly and efficiently [Kisa et al., 2014].

3.2 TRACTABLE INFERENCE

A probabilistic model is tractable when it provides exact inference—answers to queries are not approximated—and the query computation can be obtained efficiently—in time polynomial in the size of the probabilistic model [Choi et al., 2020].

Definition 3.5. A class of queries \mathbf{Q} is *tractable* on a family of probabilistic models \mathcal{M} iff any query $q \in \mathbf{Q}$ on a model $m \in \mathcal{M}$ can be computed in time $\mathcal{O}(\text{poly}(|m|))$.

The concept of efficiency translates to polytime complexity w.r.t. the size of models in a class, $|m|$. For models represented as computational graphs, such as our probabilistic circuits, model size will directly translate to the number of edges in the graph. The complexity of answering queries in the above definition depends only on one model size.

3.3 STRUCTURE LEARNING

Structure learning for PCs corresponds to learning both its structure and parameters approximating the data distribution. For advanced inference task is necessary to require structured decomposability and determinism. Two structure learner for those PCs have been recently proposed: LEARNPSDD [Liang et al., 2017] and STRUDEL [Dang et al., 2020].

LEARNPSDD performs a local search over the space of possible structured-decomposable PCs, given a vtree as input. A first hierarchical clustering is performed over the RVs in order to learn a vtree. Next the local search start from a fully factored PC conforming to a vtree—all RVs are considered to be independent. Each local step changes the circuit while preserving its semantics and structural properties of smoothness, determinism and structured decomposability. Candidates are proposed using two structural transformations—split and clone—to all nodes in the circuits. LEARNPSDD does not use the dependencies discovered during the learning of the vtree, and is very slow making difficult the learning of a mixture model.

In order to overcome these shortcomings, in STRUDEL the authors proposed to extract a vtree from the best graphical model that can be learned in tractable time, and then

compile it into a structured-decomposable PC as a more informative starting point. The complexity of the learning algorithm has been reduced adopting a greedier local search employing a single transformation, and using circuit *flows* to speed up parameter learning and likelihood computation. The resulting algorithm is a simpler, faster structure learning scheme, enabling learning of large mixtures.

4 EXTREMELY RANDOMIZED CIRCUITS

Here we present our proposed PC structure learning method exploiting a random conditioning method without optimizing the internal parameters of the circuit, thus leading to a fast learning algorithm.

4.1 CHOW-LIU TREES

A Chow-Liu Tree (see Figure 1) is a direct tree-structured Bayesian network minimizing the Kullback-Leibler divergence with the data distribution and supporting linear time marginals and MAP inference. More formally, a CLT is a first-order dependency tree \mathcal{T} over $\mathbf{X} = \{X_i\}_{i=1}^d$ equipped with parameters $\theta_{X_i|X_{\tau(i)}}$ which encodes a factorized probability distribution of the form:

$$q(\mathbf{X}) = \prod_{i=1}^d p(X_i|X_{\tau(i)}), \quad (1)$$

where $X_{\tau(i)}$ is the parent of X_i in \mathcal{T} and p is the data probability distribution¹. A CLT is learned according to the classic Chow-Liu algorithm [Chow and Liu, 1968]: the CLT structure is based on the maximum spanning tree derived by the empirical pairwise mutual information (MI) matrix of \mathbf{X} , then, the related parameters θ are estimated from the data. CLTs can be quickly compiled into smooth, deterministic and structured-decomposable PCs [Dang et al., 2020].

4.2 REGION GRAPH

In this paper we restrict to the class of PCs following a tree-shaped *region graph*. Region graphs can be viewed as a vectorised representation of PCs, and have been already used in many SPN learners [Dennis and Ventura, 2012, Peharz et al., 2013, 2019, Trapp et al., 2019].

Definition 4.1. Given a set of RVs \mathbf{X} , a *region graph* is a couple (\mathcal{R}, ϕ) where \mathcal{R} is a connected DAG containing both *regions* nodes (\ominus) and *partitions* nodes (\boxplus). Let \mathbf{R} be the set of all \ominus -nodes and \mathbf{P} be the set of all \boxplus -nodes. The scope

¹Note that, if X_i is the root of \mathcal{T} then $\tau(i) = 0$ and $p(X_i|X_0) = p(X_i)$.

function is defined as a function $\phi : \mathbf{R} \cup \mathbf{P} \rightarrow 2^{\mathbf{X}}$, assigning each node in \mathcal{R} a subset of \mathbf{X} ($2^{\mathbf{X}}$ denotes the power set of \mathbf{X}). A region graph has the following properties: i) the root R is a \ominus -node and $\phi(R) = \mathbf{X}$; ii) all leaves are \ominus -nodes; iii) \mathcal{R} is bipartite, i.e., all children of \ominus -node are \boxplus -nodes and vice versa; iv) if Q is either a \ominus -node with children or a \boxplus -node, then $\phi(Q) = \bigcup_{Q' \in \text{ch}(Q)} \phi(Q')$; v) for a \boxplus -node P we have $\forall R, R' \in \text{ch}(P) : \phi(R) \cap \phi(R') = \emptyset$; vi) for a \ominus -node R we have $\forall P \in \text{ch}(R) : \phi(R) = \phi(P)$.

Given a region graph (\mathcal{R}, ϕ) , we can obtain its corresponding PC structure \mathcal{C} as follows. The root \ominus -node in \mathcal{R} is replaced with a \oplus -node as a root in \mathcal{C} . Each leaf \ominus -node in \mathcal{R} introduces a leaf node in \mathcal{C} , while each other \ominus -node corresponds to a \oplus -node in \mathcal{C} . For each \boxplus -node in \mathcal{R} , we introduce a \otimes -node in \mathcal{C} . All the nodes introduced in \mathcal{C} have the same scope as the nodes in \mathcal{R} . It is immediate to observe that such a PC satisfies the smoothness (Definition 3.1) and decomposability (Definition 3.2) structural properties.

4.3 RANDOM REGION GRAPH

Here, we propose a random conditioning approach to build a region graph based on the satisfaction of logical constraints. The construction of a random region graph has been already proposed in [Peharz et al., 2019]. Our approach differs from that since the construction is data driven, thus leading to a fast method for learning both the structure and the parameters of a random PC without relying on parameter optimization. Indeed, while in [Peharz et al., 2019] the graph is randomly defined, here its structure is guided exploiting a random partitioning of the data.

The notation $\mathbf{x} \models \Gamma$ denotes that the assignment \mathbf{x} satisfies the logical constraint Γ , and $|\Gamma|$ the length of the constraint. For instance, given the assignment $\mathbf{x} = \{X_1 = 0, X_2 = 1, X_3 = 1, X_4 = 0\}$ and two constraints $\Gamma_1 = \{X_1 = 0 \wedge X_3 = 1\}$ and $\Gamma_2 = \{X_2 = 0 \wedge X_3 = 1\}$, then it follows that $\mathbf{x} \models \Gamma_1$ and $\mathbf{x} \not\models \Gamma_2$.

Smooth and decomposable PC. Each region or partition node Q in \mathcal{R} refers to a data slice \mathcal{S} of the entire dataset \mathcal{D} used to learn \mathcal{R} , in the following denoted as $Q_{\mathcal{S}}$. The process to build the random region graph is reported in Algorithm 1. The root region R in \mathcal{R} should be a \ominus -node, and it refers to the entire dataset \mathcal{D} (lines 1-2). The algorithm iteratively extends the graph replacing a leaf region with a sub-graph in the following manner. In order to extend a randomly selected leaf region $R_{\mathcal{S}}$ (line 4), we horizontally split its corresponding data slice \mathcal{S} into non overlapping sub-slices as described in Algorithm 2. In particular, given the $R_{\mathcal{S}}$ region, l RVs \mathbf{C} from its scope $\phi(R_{\mathcal{S}})$ are randomly chosen. Then, k logical constraints Γ_i , among the 2^l possible ones, over the selected \mathbf{C} RVs are chosen at random.

The k chosen logical constraints are used to split the \ominus -node

Algorithm 1 RandomRegionGraph($\mathcal{D}, \delta, l, k, \sigma$)

Input: a set of samples \mathcal{D} over a set of RVs \mathbf{X} , a minimum number of examples per partition δ , a split arity k , a conjunction length l and an ordering σ of \mathbf{X}

Output: a random region graph \mathcal{G}

```
1:  $\mathcal{G} \leftarrow \text{INSROOT}(\Xi_{\mathcal{D}})$ 
2:  $\mathcal{P} \leftarrow \{\text{ROOT}(\mathcal{G})\}$ 
3: while  $\mathcal{P} \neq \emptyset$  do
4:    $P \leftarrow$  a random region  $\Xi_{\mathcal{R}}$  from  $\mathcal{P}$ 
5:    $\mathbf{C} \leftarrow$  the first RVs from  $\sigma(\mathbf{X}_P)$ 
6:    $\mathcal{S} \leftarrow \text{RandRegions}(P, \delta, \mathbf{C}, k)$ 
7:   if  $\mathcal{S} \neq \emptyset$  then
8:     for each  $Q \in \mathcal{S}$  do
9:       AddSub( $P, \Xi_Q$ )
10:    for each  $Q \in \text{sub}(P)$  do
11:      AddSub( $Q, \Xi_{Q[\mathbf{C}]}$ )
12:      AddSub( $Q, \Xi_{Q[\mathbf{X}_P - \mathbf{C}]}$ )
13:     $\mathcal{P} \leftarrow \mathcal{P} \cup \{\Xi_{Q[\mathbf{X}_P - \mathbf{C}]}\}$ 
14: return  $\mathcal{G}$ 
```

Algorithm 2 RandRegions($\Xi_{\mathcal{S}}, \delta, \mathbf{C}, k$)

Input: a region $\Xi_{\mathcal{S}}$ over a data slice \mathcal{S} , a minimum number of instances per partition δ , a set of RVs \mathbf{C} , a split arity k

Output: a set of sub-slices $\mathcal{H} \cup \mathcal{R}$ of \mathcal{S}

```
1:  $\mathcal{R} \leftarrow \mathcal{S}$ 
2:  $\mathcal{H} \leftarrow \emptyset$ 
3: while  $\perp \neq \Delta \leftarrow \text{RFORM}(\mathbf{C})$  and  $|\mathcal{H}| \neq k - 1$  do
4:    $\mathcal{Q} \leftarrow \{\mathbf{x} \in \mathcal{R} : \mathbf{x} \neq \Delta\}$ 
5:   if  $|\mathcal{Q}| \geq \delta$  and  $|\mathcal{R} - \mathcal{Q}| \geq \delta$  then
6:      $\mathcal{H} \leftarrow \mathcal{H} \cup \{\mathcal{Q}\}$ 
7:      $\mathcal{R} \leftarrow \mathcal{R} - \mathcal{Q}$ 
8:   if  $\mathcal{H} \neq \emptyset$  then
9:     return  $\mathcal{H} \cup \{\mathcal{R}\}$ 
10: else
11:   return  $\emptyset$ 
```

$R_{\mathcal{S}}$ into $k + 1$ \boxplus -nodes. Each constraint $\Gamma_i, i = 1, \dots, k$, leads to a \boxplus -node $R'_{\mathcal{S}_i}$ if at least δ samples \mathcal{S}'_i of the data slice \mathcal{S} satisfy Γ_i . A remaining \boxplus -node contains all the samples in \mathcal{S} not satisfying any constraint. The corresponding \boxplus -node in the circuit to the \boxplus -node $R_{\mathcal{S}}$ in \mathcal{R} has $k + 1$ branches whose weights are estimated as $|\mathcal{S}'_i|/|\mathcal{S}|$.

Now, each obtained \boxplus -node $R'_{\mathcal{S}_i}$ becomes a child of the \boxplus -node $R_{\mathcal{S}}$, and a parent of two new \boxplus -nodes (a vertical split of the sub-slice): the first region $R''_{\mathcal{S}'_i[\mathbf{C}]}$ on the sub-slice $\mathcal{S}'_i[\mathbf{C}]$ and the second one $R''_{\mathcal{S}'_i[\phi(R_{\mathcal{S}}) - \mathbf{C}]}$ on the remaining variables, where with $\mathcal{S}[\mathbf{C}]$ we denote the selection of the columns \mathbf{C} from the slice \mathcal{S} . The first k obtained $R''_{\mathcal{S}'_i[\mathbf{C}]}$ regions are called \mathcal{Q} -regions, while the last one is called \mathcal{R} -region. All the $R''_{\mathcal{S}'_i[\phi(R_{\mathcal{S}}) - \mathbf{C}]}$ regions are called \mathcal{S} -regions. During the iterative process, only \mathcal{S} -regions are considered for the graph extension.

At the end of the iterative process all the leaf \boxplus -nodes corresponding to \mathcal{Q} -regions or \mathcal{R} -regions are represented in the corresponding circuit as a PoB with Laplace correction, while those corresponding to \mathcal{S} -regions are modeled as CLTs with Laplace correction. This process provides a random region graph that leads to a smooth and decomposable PC (see Figure 2 in Appendix A).

The parameters of every sub-circuit are estimated from the data slice corresponding to the related region. Differently from [Peharz et al., 2019], where the graph is randomly constructed, here, however, its construction is guided by the chosen variable ordering and by the corresponding data partitioning, conditioned on the variable assignments.

Deterministic PC. In order to impose determinism, given an assignment, the *flow* for each sum node in the circuit must be unique, where with flows we indicate the 'activated' edges in the circuit for an input assignment. Hence, first of all, we remove the Laplace corrections from the sub-circuits corresponding to \mathcal{Q} -regions and \mathcal{S} -regions. Furthermore, \mathcal{R} -regions are modeled as a deterministic PC which evaluates to 0 only for samples having an assignment satisfying the chosen constraints for its sibling regions, thus ensuring determinism for the whole circuit.

Structured Decomposable PC. In order to impose structural decomposability we cannot chose l RVs at random for each region to split since the circuit must be normalized for a vtree (see Section 3.1). Hence, a fixed random ordering $\sigma(\mathbf{X})$ of the RVs is chosen at the beginning of the iterative process and the variables are always pushed from that ordering. However, providing an ordering $\sigma(\mathbf{X})$ of the RVs is necessary but not sufficient to learn structured decomposable PCs. Indeed, we cannot independently learn the CLTs in each \mathcal{S} -region, because there is not guarantee they share the same vtree when learned on their corresponding data slice. Therefore, it is necessary to learn for each

\mathcal{S} -region a CLT sharing the same structure. The following theorem proves what is the best first-order dependency tree approximating many distributions.

Theorem 1. *Let $\{p_i\}_{k=1}^n$ be probability distributions over $\mathbf{X} = \{X_i\}_{i=1}^d$ and $\{q_i\}_{k=1}^n$ their corresponding approximations based on a same first-order dependency tree \mathcal{T} . It is possible to prove that the first-order dependency tree $\widehat{\mathcal{T}}$ minimizing $\sum_{k=1}^n \mathbb{KL}(p_k || q_k)$ is:*

$$\widehat{\mathcal{T}} = \arg \max_{\mathcal{T}} \sum_{k=1}^n \sum_{i=1}^d \text{MI}_k(X_i; X_{\tau(i)}), \quad (2)$$

where MI_k is the mutual information on p_k and $X_{\tau(i)}$ is the parent of X_i in \mathcal{T} .

Theorem 1 proves that the best first-order dependency tree for many \mathcal{S} -regions—over the same scope—is based on the maximum spanning tree derived by the matrix obtained by adding the empirical MI computed on each data slice (proof in Appendix B).

However, it may happen to have \mathcal{S} -regions over different scopes (see Figure 3 in Appendix A) thus avoiding to straightforwardly apply the Theorem 1. Furthermore, we can note that the conditioning process implicitly imposes a partial structure to the vtree. In particular, let $\mathbb{C} = \{\mathbf{C}_i\}_{i=1}^s$ be the ordered set in which \mathbf{C}_i contains the RVs of the i -th conditioning traversing top-down the region graph. Let $\mathbf{F} = \mathbf{X} - \bigcup_{i=1}^s \mathbf{C}_i$ be the set of variables never involved in a conditioning. Hence, only on those variables \mathbf{F} we can learn the best dependence tree using the result in Theorem 1.

To complete the vtree we can proceed as follows. For each \mathcal{S} -partition, we accumulate the estimated MI just for those variables in \mathbf{F} . Let $\mathcal{T}_{\mathbf{C}_1}, \mathcal{T}_{\mathbf{C}_2}, \dots, \mathcal{T}_{\mathbf{C}_s}$ be the dependence trees resp. over the variables in $\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_s$, and $\mathcal{T}_{\mathbf{F}}$ that over \mathbf{F} . We build the complete dependence tree over \mathbf{X} joining the dependence trees such that, the root of $\mathcal{T}_{\mathbf{C}_{i+1}}$ becomes a child of $\mathcal{T}_{\mathbf{C}_i}$ for $i = 1, \dots, s-1$, and the root of $\mathcal{T}_{\mathbf{F}}$ becomes a child of $\mathcal{T}_{\mathbf{C}_s}$. Now, each CLT for a \mathcal{S} -region is compiled using the sub-tree in the vtree associated to the variables in its scope.

4.4 MIXTURES OF XPCS

Recently, a lot of attention has been devoted to learn mixtures of PCs to greatly improve their performance as density estimators, encoding a distribution as the following:

$$m(\mathbf{X}) = \sum_{i=1}^k w_i q_i(\mathbf{X}),$$

where k is the number of components $q_i(\mathbf{X})$ and w_i are the weights such that $\sum_{i=1}^k w_i = 1$.

Putting many XPC in a mixture results in a non-deterministic circuit we call EXPC, since it introduces a sum node over many XPC marginalizing a latent variable [Peharz et al., 2017]—all the branches of this sum node are ‘activated’ for an input assignment.

An ensemble of structured decomposable XPCs remains not structured decomposable unless all of them are normalized for the same vtree [Dang et al., 2020], a constraint very difficult to fulfill in general. However, in our learning scheme we propose to just learn a structured decomposable mixture using the same ordering σ for each involved XPC. Therefore, while non-structured decomposable EXPCs can rely on higher randomness, structured decomposable ones are highly dependent on such a fixed ordering, and their randomness only lies in choosing the logical constraints for horizontal splits. Moreover, since this ordering is shared among all XPCs, it directly affects the ensemble performance.

We noticed that choosing the ordering at random provides poor results. Instead, we propose the following greedy procedure. The algorithm iteratively selects blocks of l variables from \mathbf{X} . A first variable is selected having the highest MI w.r.t. the other ones; the remaining $l-1$ ones are those having the highest MI w.r.t. the first one. After the selection block of selected variables, the process continues selecting the next block on the remaining variables. In this way the conditioning is done on strong conditionally dependent variables thus providing better results. All the experiments for structured decomposable PCs use this procedure for selecting the ordering.

It is possible to note that vertical splits are automatically introduced due to the random conditioning proposed approach. Furthermore, the proposed heuristic for variable ordering (and consequently for a vtree structure) imposes a correlation among the variables in the conjunction, thus providing effective product nodes.

The weights of each component are simply set to be uniform among the ensemble, i.e., $w_i = 1/k$ for each $i = 1, \dots, k$. Even if this is not the optimal choice, however in this way we avoid to optimize the using time consuming procedures like EM.

5 RESULTS

In this section we empirically evaluate our proposed approach² [Gala, 2021]. We evaluate our learner on 20 real-world benchmark datasets [Haaren and Davis, 2012] (reported in Appendix C, Table 3), already used to evaluate different tractable density estimators. In particular, we aim to answer the following research questions: **Q1**) how much

²<https://github.com/gengala/Random-Probabilistic-Circuits>.

accurate are single XPCs when compared to other density estimators? And what is the effect of increasing the length of the conditioning on XPCs. **Q2)** What is the learning time and circuit size obtained with the proposed approach. **Q3)** Are ensemble of XPCs competitive to those learned using other approaches?

In our grid searches some hyper-parameters have been tuned, while keeping fixed to 0.01, for computational reasons, the smoothing parameter α used in the CLTs construction for the Laplace correction. We stopped to grow the region graph as soon as its leaves are more than 200. Experiments have been executed on a 8-core Intel Xeon CPU E5-1620 v3 @ 3.50GHz with 16 GB RAM.

Let us denote a deterministic XPC as XPC_{Det} and a structured decomposable XPC as XPC^{SD} .

5.1 Q1) SINGLE MODELS

We evaluate different XPCs and we investigate how they perform by varying the values for l (the length of the logical constraint) and k (the number of constraints for the horizontal splits) in order to asses whether the proposed random region graph method works. In particular, we are interested in XPCs with $l > 1$, since XPCs with $l = 1$ are basically CNets, extensively studied in literature.

Due to the randomness of XPCs, we report the results averaged on 40 different runs for each kind of model in the following grid search space: $l = \{1, 2, 3\}$, $k = \{2, 3, 4, 8\}$ and $\delta = \{16, 32, 64, 128, 256, 512\}$ ³ (the minimum number of samples per slice for splitting).

The results of such a grid search (reported in Appendix D, Table 4), shows that XPCs with $l > 1$ obtain better results, i.e., just 33 XPCs with $l = 1$ out of 120 are marked as best ones, proving that conditioning on more than one variable at the same time is more effective. The results of XPC wins over XPC_{Det} 15 out of 20 times, and XPC^{SD} wins 13 out of 20 over XPC_{Det} . This obviously reflects the higher expressiveness of non-deterministic PCs over deterministic ones. Furthermore, XPC^{SD} wins 15 out of 20 over XPC and XPC_{Det} wins 15 out of 20 over XPC_{Det} , even though structured decomposable PCs are more restrictive than non-structured decomposable ones, since their structure have to fulfill a vtree. However, as reported in Section 4.4, this could be explained by the adopted greedy heuristic used to select the variables on which is chosen the conditioning, and thus providing a correct conditional probability estimation. A random vtree does not represent the conditional dependencies among the variables the greedy heuristic instead provides.

Table 1 compares, in the first two columns, the results of a

³Configurations with $k \notin [2, 2^l]$ are discarded.

single XPC_{Det} , averaged by 40 runs, against those of another deterministic circuit such a XCNet, as reported in [Di Mauro et al., 2017]. The results show how the two models rank the same. In the remaining columns we evaluate $XPC_{\text{Det}}^{\text{SD}}$ sharing the same structural constraints of the two competitors LEARNPSDD [Liang et al., 2017] and STRUDEL. $XPC_{\text{Det}}^{\text{SD}}$ log-likelihoods are in line with those of its competitors despite its random structure and the extremely fast learning time.

5.2 Q2) LEARNING TIME AND CIRCUIT SIZE

XPCs learning time depends mainly on learning CLTs and regardless of whether or not we force SD and/or determinism. More specifically, under the same region graph, every XPC type requires the same learning time. In terms of circuit sizes, deterministic XPCs need a greater number of edges in order to model \mathcal{R} -partitions, but nevertheless the size increase is negligible for small values of l . In general, learning times and circuit sizes are obviously inversely proportional to δ given that smaller δ values generate deeper region graphs and then bigger circuits.

In Appendix D, Table 5 reports the average circuit size, the average training time (in seconds) and the δ value of the 40 best models (bold ones) of Table 4. The statistics are quite stable unless few outliers such as the XPC^{SD} over the AD dataset.

Table 6, in Appendix D, reports the learning times (in seconds) and circuit sizes for $XPC_{\text{Det}}^{\text{SD}}$ (averaged over 40 runs) and STRUDEL best models. STRUDEL models are larger PCs than $XPC_{\text{Det}}^{\text{SD}}$ thus leading to efficient inference. Furthermore, the learning time required to learn $XPC_{\text{Det}}^{\text{SD}}$ models is only a tiny fraction of the time required by the competitor so as to extremely speed up learning. To the best of our knowledge, XPCs are the fastest TPMs to learn among those available in literature.

5.3 Q3) MIXTURE MODELS

To investigate the performance of ensembles of XPCs (EXPCs) we run a grid search in the space formed by $l = \{1, 2, 3\}$, $k = \{2, 3, 4, 8\}$, $\delta = \{16, 32, 64, 128, 256, 512\}$ and $M = \{2, 5, 10, 15, 20, 25, 30, 40\}$ (the number of mixtures) for every EXPC type.

Learning times and circuit sizes of EXPCs scale linearly in M . Next, we compare EXPCs to other state-of-the-art ensembling techniques and much more sophisticated single models such as ID-SPN.

Table 7, in Appendix E, shows that most of best test-set log-likelihoods is associated to EXPCs having $l > 1$ and this, once again, justifies the conjunctions-based approach. In particular, just 13 EXPCs having $l = 1$ out of 80 (i.e. 16.25%)

Table 1: Average test-set log-likelihoods for XCNet, XPC_{Det}, XPC_{Det}^{SD}, LEARNPSDD and STRUDEL.

DATASET	XPC _{Det}	XCNet	XPC _{Det} ^{SD}	LEARNPSDD	STRUDEL
NLTCS	-6.10	-6.06	-6.09	-6.03	-6.06
MSNBC	-6.18	-6.09	-6.21	-6.04	-6.05
KDD	-2.22	-2.19	-2.20	-2.17	-2.17
PLANTS	-13.96	-13.43	-14.59	-13.49	-13.72
AUDIO	-42.65	-42.66	-41.97	-41.51	-42.26
JESTER	-56.00	-56.10	-54.94	-54.63	-55.30
NETFLIX	-59.28	-59.21	-58.73	-58.53	-58.68
ACCIDENTS	-31.88	-31.58	-31.03	-28.29	-29.46
RETAIL	-10.95	-11.44	-10.98	-10.92	-10.90
PUMSB-STAR	-25.90	-25.55	-26.56	-25.40	-25.28
DNA	-87.75	-87.67	-87.46	-83.02	-87.10
KOSAREK	-11.29	-11.70	-10.99	-10.99	-10.98
MSWEB	-10.19	-10.47	-10.12	-9.93	-10.19
BOOK	-37.51	-42.36	-36.83	-36.06	-35.77
EACHMOVIE	-62.62	-60.71	-59.99	-55.41	-59.47
WEBKB	-163.33	-167.45	-161.26	-161.42	-160.50
ROUTERS-52	-94.29	-99.52	-88.92	-93.30	-92.38
20NEWS-GRP	-163.95	-172.60	-159.37	-160.43	-160.77
BBC	-261.96	-261.79	-260.62	-260.24	-258.96
AD	-16.40	-18.70	-16.39	-20.13	-16.52
AVG. RANK	1.5	1.5	2.45	1.55	2

Table 2: Average test-set log-likelihoods for structured decomposable mixture models.

DATASET	EXPC ^{SD}	EXPC _{Det} ^{SD}	STRUDEL
NLTCS	-6.05	-6.05	-6.08
MSNBC	-6.18	-6.17	-6.04
KDD	-2.15	-2.16	-2.16
PLANTS	-14.19	-14.21	-13.73
AUDIO	-40.91	-40.97	-41.48
JESTER	-53.43	-53.51	-55.04
NETFLIX	-57.58	-57.69	-58.25
ACCIDENTS	-31.02	-30.99	-29.07
RETAIL	-10.94	-10.90	-10.90
PUMSB-STAR	-26.06	-26.05	-24.17
DNA	-86.61	-85.09	-87.21
KOSAREK	-10.77	-10.81	-10.89
MSWEB	-9.93	-9.94	-9.76
BOOK	-34.75	-35.14	-35.89
EACHMOVIE	-54.82	-55.26	-55.76
WEBKB	-153.67	-154.23	-159.85
ROUTERS-52	-84.70	-85.09	-90.12
20NEWS-GRP	-153.75	-154.21	-158.79
BBC	-248.34	-248.79	-257.40
AD	-15.50	-15.59	-15.39
AVG. RANK	1.6	2	2.25

are marked as best ones. By comparing non-rounded test-set log-likelihoods, EXPC wins over EXPC_{Det} 12 out of 20 times, and EXPC^{SD} wins 14/20 over EXPC_{Det}^{SD}. Therefore, regardless of whether we use deterministic or non-deterministic components, EXPCs do not show a significant difference in log-likelihood. In fact, an ensemble can be expressive enough even employing deterministic components. Furthermore, we can therefore reach competitive log-

likelihoods by employing non-structured decomposable ensemble with structured decomposable components. Finally, non-structured decomposable EXPC wins over structured decomposable ones 19 out of 20 times, the only exception occurs for ROUTERS-52. This proves the limited expressiveness of structured decomposable circuits compared to non-structured decomposable ones.

Table 8, in Appendix E, reports the log-likelihoods for best EXPCs, SOTA ensembling techniques and single model competitors. Among mixture models, EXPCs provide the best average rank; while, when we include single models competitors, they provide the second best average rank. However, it should be notice that EXPCs are obviously bigger PCs than single model competitors but nevertheless can be learned employing only a fraction of the time.

Table 2 reports the log-likelihoods of structured decomposable EXPCs and their natural competitor STRUDEL. As we can see our proposed approach ranks better than STRUDEL, thus providing its validity.

6 CONCLUSIONS

We introduced XPCs, simple and customizable probabilistic circuits based on random logical constraints-based conditioning. XPCs are smooth and decomposable by default but can be easily forced to be deterministic and/or structured decomposable without increasing the learning complexity. When learned in ensembles, XPCs achieve the state-of-

the-art results for density estimation on several benchmark datasets. Due to their simplicity to implement, fast learning times and accurate inference performances XPCs are promising tractable density estimators.

Many lines of work are possible for XPCs. Firstly, we plan to employ greedy orderings also for non-structured decomposable circuits just to improve their log-likelihoods as single models. Secondly, we plan to extend XPCs for continuous RVs and to employ a DAG rather than a tree structure to partition the data. Thirdly, we plan to investigate how EXPCs perform using bagged datasets, as done by mixtures of LEARNPSDD like in [Dang et al., 2020]. Finally, due to their fast learning, we plan to investigate how sophisticated learners as LEARNPSDD and STRUDEL perform if we provide XPCs as initial PCs.

Acknowledgements

The authors would like to thank Antonio Vergari for inspiring discussions and valuable feedback.

References

Hei Chan and Adnan Darwiche. On the robustness of most probable explanations. In *UAI*, 2006.

YooJung Choi, Antonio Vergari, and Guy Van den Broeck. Probabilistic circuits: A unifying framework for tractable probabilistic models. Technical report, 2020.

C. K. Chow and C. N. Liu. Approximating discrete probability distributions with dependence trees. *IEEE Trans. Inf. Theory*, 14(3):462–467, 1968.

Meihua Dang, Antonio Vergari, and Guy Van den Broeck. Strudel: Learning structured-decomposable probabilistic circuits. In *PGM*, 2020.

Adnan Darwiche. A differential approach to inference in bayesian networks. In Craig Boutilier and Moisés Goldszmidt, editors, *UAI*, pages 123–132, 2000.

Adnan Darwiche. *Modeling and Reasoning with Bayesian Networks*. Cambridge University Press, 2009.

A. Dempster, N. Laird, and D. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B*, 1977.

Aaron W. Dennis and Dan Ventura. Learning the architecture of sum-product networks using clustering on variables. In *NIPS*, pages 2042–2050, 2012.

Nicola Di Mauro, Antonio Vergari, Teresa Maria Altomare Basile, and Floriana Esposito. Fast and accurate density estimation with extremely randomized cutset networks. In *ECML*, pages 203–219, 2017.

Gennaro Gala. gengala/Random-Probabilistic-Circuits: First release, May 2021. URL <https://doi.org/10.5281/zenodo.4775258>.

Robert Gens and Pedro Domingos. Learning the structure of sum-product networks. In *ICML*, pages 873–880, 2013.

Jan Van Haaren and Jesse Davis. Markov network structure learning: A randomized feature generation approach. In *AAAI*, 2012.

Doga Kisa, Guy Van den Broeck, Arthur Choi, and Adnan Darwiche. Probabilistic sentential decision diagrams. In *ICPKRR*, 2014.

Yitao Liang, Jessa Bekker, and Guy Van den Broeck. Learning the structure of probabilistic sentential decision diagrams. In *UAI*, 2017.

Marina Meila and Michael I. Jordan. Learning with mixtures of trees. *Journal of Machine Learning Research*, 1:1–48, 2000.

George Papamakarios, Iain Murray, and Theo Pavlakou. Masked autoregressive flow for density estimation. In *NIPS*, pages 2338–2347, 2017.

Robert Peharz, Bernhard C. Geiger, and Franz Pernkopf. Greedy part-wise learning of sum-product networks. In *ECML/PKDD*, pages 612–627, 2013.

Robert Peharz, Robert Gens, Franz Pernkopf, and Pedro M. Domingos. On the latent variable interpretation in sum-product networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(10):2030–2044, 2017.

Robert Peharz, Antonio Vergari, Karl Stelzner, Alejandro Molina, Martin Trapp, Xiaoting Shao, Kristian Kersting, and Zoubin Ghahramani. Random sum-product networks: A simple and effective approach to probabilistic deep learning. In *UAI*, volume 115, pages 334–344, 2019.

Hoifung Poon and Pedro M. Domingos. Sum-product networks: A new deep architecture. In *UAI*, 2011.

Tahrima Rahman and Vibhav Gogate. Learning ensembles of cutset networks. In *AAAI*, pages 3301–3307, 2016.

Tahrima Rahman, Prasanna Kothalkar, and Vibhav Gogate. Cutset networks: A simple, tractable, and scalable approach for improving the accuracy of chow-liu trees. In *ECML/PKDD*, pages 630–645, 2014.

Amirmohammad Rooshenas and Daniel Lowd. Learning sum-product networks with direct and indirect variable interactions. In *ICML*, pages 710–718, 2014.

Martin Trapp, Robert Peharz, Hong Ge, Franz Pernkopf, and Zoubin Ghahramani. Bayesian learning of sum-product networks. In *NIPS*, pages 6344–6355, 2019.

Benigno Uribe, Marc-Alexandre Côté, Karol Gregor, Iain Murray, and Hugo Larochelle. Neural autoregressive distribution estimation. *Journal of Machine Learning Research*, 17:205:1–205:37, 2016.

Fabrizio Ventola, Karl Stelzner, Alejandro Molina, and Kristian Kersting. Random sum-product forests with residual links. *CoRR*, abs/1908.03250, 2019.

Antonio Vergari, YooJung Choi, Robert Peharz, and Guy Van den Broeck. Probabilistic circuits: Representations, inference, learning and applications. In *Tutorial at the The 34th AAAI Conference on Artificial Intelligence*, 2020.

Supplementary Materials

A AN EXAMPLE OF DETERMINISTIC AND STRUCTURED DECOMPOSABLE XPC

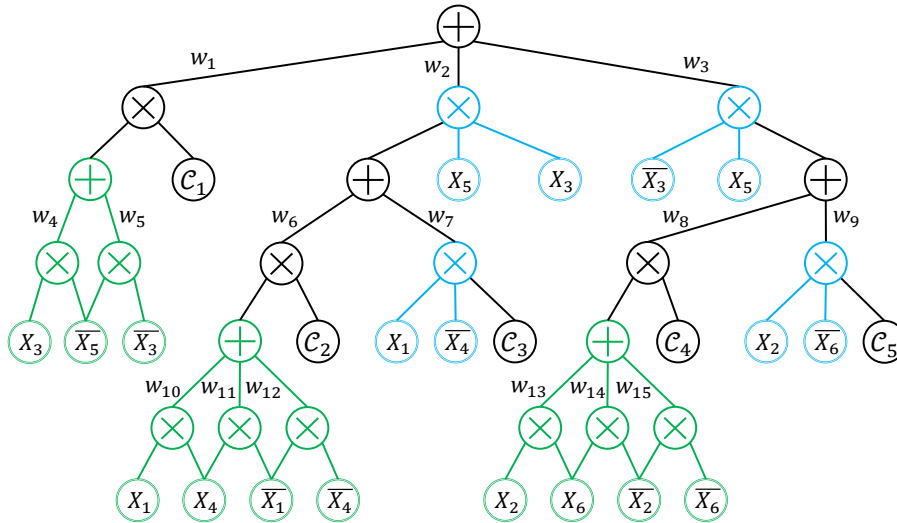


Figure 2: A deterministic and non-SD XPC learned setting $l = 2$ and $k = 3$ over RVs $\mathbf{X} = \{X_i\}_{i=1}^7$. The sub-circuits in blue (resp., green) encode \mathcal{Q} -regions (resp., \mathcal{R} -regions). Every CLT (represented as a multivariate node \mathcal{C}_i) is learned by following the traditional Chow-Liu tree algorithm [Chow and Liu, 1968].

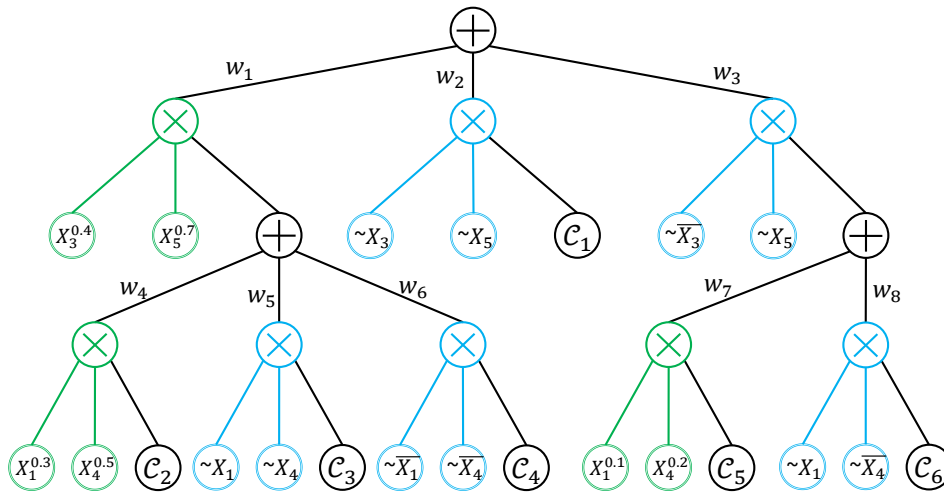


Figure 3: A non-deterministic and SD XPC learned setting $l = 2$ and $k = 3$ over RVs $\mathbf{X} = \{X_i\}_{i=1}^7$. The sub-circuits in blue (resp., green) encode \mathcal{Q} -regions (resp., \mathcal{R} -regions). $\sim X_i$ denotes the smoothed indicator of X_i and X_i^p denotes a Bernoulli node of X_i with parameter p . The CLTs (represented as multivariate nodes \mathcal{C}_i) have the same dependency tree—are normalized for the same vtree.

B PROOF OF THEOREM 1.

We recall from [Chow and Liu, 1968] that the best first-order dependency tree $\widehat{\mathcal{T}}$ over a set of RVs $\mathbf{X} = \{X_i\}_{i=1}^d$ is:

$$\widehat{\mathcal{T}} = \arg \max_{\mathcal{T}} \sum_{i=1}^d \text{MI}(X_i, X_{\tau(i)}), \quad (3)$$

where $X_{\tau(i)}$ is the parent of X_i in \mathcal{T} . Let $\{p_i\}_{i=1}^n$ be probability distributions on \mathbf{X} and $\{q_i\}_{i=1}^n$ their corresponding approximations based on a same first-order dependency tree \mathcal{T} . We want to find $\widehat{\mathcal{T}}$ such that:

$$\widehat{\mathcal{T}} = \arg \min_{\mathcal{T}} \sum_{k=1}^n \mathbb{KL}(p_k || q_k) = \arg \max_{\mathcal{T}} \sum_{k=1}^n \mathbb{E}_{\mathbf{x} \sim p_k} [\log q_k(\mathbf{x})], \quad (4)$$

where \mathbb{KL} denotes the Kullback-Leibler divergence. Thus considering a dependency tree \mathcal{T} :

$$\sum_{k=1}^n \mathbb{E}_{\mathbf{x} \sim p_k} [\log q_k(\mathbf{x})] = \sum_{k=1}^n \mathbb{E}_{\mathbf{x} \sim p_k} \left[\sum_{i=1}^d \log p_k(X_i | X_{\tau(i)}) \right] \quad (5)$$

$$= \sum_{k=1}^n \left(\sum_{i=1}^d \text{MI}_k(X_i; X_{\tau(i)}) + \mathbb{E}_{\mathbf{x} \sim p_k} \left[\sum_{i=1}^d \log p_k(X_i) \right] \right), \quad (6)$$

where $\text{MI}_k(X, Y)$ is the mutual information of X and Y on p_k . Therefore:

$$\widehat{\mathcal{T}} = \arg \min_{\mathcal{T}} \sum_{k=1}^n \mathbb{KL}(p_k || q_k) = \arg \max_{\mathcal{T}} \sum_{k=1}^n \sum_{i=1}^d \text{MI}_k(X_i; X_{\tau(i)}). \quad (7)$$

C TWENTY DENSITY ESTIMATION DATASETS

We report in Table 3 the adopted 20 real-world density estimation datasets [Haaren and Davis, 2012] and their statistics, already used to evaluate different tractable density estimators.

Table 3: Standard benchmark datasets and their statistics.

DATASET	RVS	TRAIN	VALID	TEST
NLTCS	16	16181	2157	3236
MSNBC	17	291326	38843	58265
KDD	64	180092	19907	34955
PLANTS	69	17412	2321	3482
AUDIO	100	15000	2000	3000
JESTER	100	9000	1000	4116
NETFLIX	100	15000	2000	3000
ACCIDENTS	111	12758	1700	2551
RETAIL	135	22041	2938	4408
PUMSB-STAR	163	12262	1635	2452
DNA	180	1600	400	1186
KOSAREK	190	33375	4450	6675
MSWEB	294	29441	3270	5000
BOOK	500	8700	1159	1739
EACHMOVIE	500	4525	1002	591
WEBKB	839	2803	558	838
ROUTERS-52	889	6532	1028	1540
20NEWS-GRP	910	11293	3764	3764
BBC	1058	1670	225	330
AD	1556	2461	327	491

D SINGLE MODELS

D.1 DENSITY ESTIMATION PERFORMANCES

We compare all the variants of XPCs. The average test-set log-likelihood is reported in Table 4

Table 4: Means and standard deviations of the best average test-set log-likelihoods over 40 runs for each variant of XPCs type in the space formed by $l = \{1, 2, 3\}$, $k = \{2, 3, 4, 8\}$ and $\delta = \{16, 32, 64, 128, 256, 512\}$. For every dataset, each row corresponds (from top to down) to XPC, XPC_{Det}, XPC^{SD} and XPC_{Det}^{SD}. For each model (row) we report the best value in bold.

	$l = 1$ $k = 2$	$l = 2$ $k = 2$	$l = 2$ $k = 3$	$l = 2$ $k = 4$	$l = 3$ $k = 2$	$l = 3$ $k = 3$	$l = 3$ $k = 4$	$l = 3$ $k = 8$
NLTCS	-6.10±0.03	-7.03±0.36	-6.35±0.23	-6.14±0.02	-7.75±0.57	-6.92±0.41	-6.55±0.33	-6.15±0.02
	-6.10±0.03	-6.76±0.23	-6.26±0.13	-6.11±0.02	-7.16±0.31	-6.59±0.21	-6.37±0.17	-6.12±0.02
	-6.09±0.00	-7.03±0.58	-6.42±0.26	-6.12±0.01	-7.57±0.84	-6.69±0.35	-6.36±0.06	-6.14±0.02
	-6.09±0.00	-6.70±0.31	-6.27±0.13	-6.10±0.00	-6.77±0.31	-6.40±0.13	-6.27±0.03	-6.10±0.01
MSNBC	-6.18±0.04	-6.68±0.06	-6.38±0.07	-6.21±0.03	-6.78±0.06	-6.70±0.07	-6.54±0.10	-6.22±0.03
	-6.18±0.04	-6.51±0.06	-6.29±0.06	-6.19±0.03	-6.62±0.05	-6.49±0.08	-6.36±0.08	-6.19±0.03
	-6.21±0.02	-6.72±0.08	-6.45±0.10	-6.25±0.02	-6.75±0.07	-6.73±0.09	-6.62±0.08	-6.26±0.04
	-6.21±0.02	-6.50±0.06	-6.33±0.06	-6.23±0.02	-6.57±0.04	-6.49±0.06	-6.40±0.06	-6.24±0.04

KDD	-2.21±0.01	-2.23±0.02	-2.21±0.01	-2.21±0.00	-2.24±0.02	-2.22±0.01	-2.21±0.01	-2.21±0.01
	-2.22±0.01	-2.23±0.02	-2.22±0.01	-2.22±0.01	-2.23±0.02	-2.21±0.01	-2.21±0.01	-2.22±0.01
	-2.20±0.00	-2.24±0.02	-2.21±0.01	-2.19±0.00	-2.28±0.04	-2.24±0.03	-2.21±0.01	-2.20±0.00
	-2.20±0.00	-2.24±0.02	-2.21±0.01	-2.20±0.00	-2.26±0.03	-2.22±0.02	-2.20±0.01	-2.20±0.00
PLANTS	-13.83±0.15	-15.15±0.63	-14.02±0.26	-13.67±0.17	-16.50±0.90	-14.98±0.60	-14.33±0.49	-13.56±0.09
	-14.11±0.14	-15.26±0.34	-14.26±0.21	-13.97±0.09	-16.22±0.71	-14.94±0.45	-14.41±0.26	-13.96±0.13
	-14.93±0.00	-15.55±0.59	-14.75±0.25	-14.55±0.04	-15.55±0.81	-14.91±0.31	-14.67±0.07	-14.52±0.05
	-14.93±0.01	-15.37±0.47	-14.82±0.16	-14.67±0.03	-15.25±0.34	-14.86±0.15	-14.73±0.05	-14.59±0.05
AUDIO	-42.54±0.10	-42.36±0.24	-42.30±0.21	-42.38±0.12	-42.98±0.44	-42.33±0.29	-42.20±0.20	-42.32±0.13
	-42.63±0.11	-43.05±0.24	-42.78±0.13	-42.67±0.15	-43.44±0.35	-42.97±0.22	-42.85±0.25	-42.65±0.14
	-41.96±0.00	-42.31±0.47	-42.13±0.16	-42.08±0.03	-43.07±0.71	-42.58±0.38	-42.22±0.26	-42.11±0.03
	-41.97±0.00	-42.67±0.44	-42.29±0.16	-42.16±0.03	-43.00±0.51	-42.73±0.31	-42.49±0.24	-42.25±0.03
JESTER	-55.94±0.21	-55.52±0.34	-55.62±0.23	-55.73±0.14	-56.26±0.49	-55.65±0.34	-55.47±0.29	-55.61±0.21
	-56.09±0.22	-56.59±0.38	-56.25±0.25	-56.06±0.15	-56.85±0.45	-56.37±0.37	-56.19±0.33	-56.00±0.23
	-54.89±0.00	-55.15±0.52	-54.86±0.22	-54.84±0.02	-56.11±0.69	-55.20±0.32	-54.95±0.33	-54.67±0.00
	-54.94±0.00	-55.72±0.45	-55.16±0.23	-55.01±0.02	-56.12±0.52	-55.59±0.29	-55.42±0.31	-54.93±0.01
NETFLIX	-59.28±0.10	-59.23±0.20	-59.22±0.13	-59.21±0.13	-60.08±0.31	-59.44±0.19	-59.22±0.17	-59.22±0.13
	-59.29±0.10	-59.79±0.17	-59.43±0.15	-59.28±0.13	-60.26±0.24	-59.82±0.14	-59.57±0.16	-59.33±0.12
	-58.72±0.00	-59.17±0.51	-58.88±0.15	-58.79±0.02	-60.29±0.72	-59.27±0.27	-58.99±0.31	-58.87±0.01
	-58.73±0.00	-59.37±0.36	-58.99±0.17	-58.84±0.01	-59.85±0.41	-59.34±0.20	-59.17±0.21	-58.97±0.01
ACCIDENTS	-31.90±0.27	-32.76±0.38	-32.03±0.34	-31.84±0.24	-33.83±0.63	-32.77±0.40	-32.31±0.34	-31.92±0.27
	-31.88±0.37	-32.76±0.36	-32.09±0.34	-31.91±0.25	-33.77±0.56	-32.72±0.35	-32.32±0.33	-32.03±0.28
	-31.01±0.00	-32.02±0.37	-31.35±0.01	-31.35±0.01	-34.40±1.02	-32.06±0.28	-31.38±0.22	-31.17±0.21
	-31.03±0.00	-31.88±0.24	-31.35±0.01	-31.35±0.01	-32.78±0.44	-31.73±0.13	-31.30±0.15	-31.15±0.08
RETAIL	-10.95±0.00	-10.96±0.02	-10.96±0.02	-10.96±0.02	-11.01±0.04	-11.02±0.04	-11.02±0.03	-11.02±0.03
	-10.95±0.00	-10.96±0.01	-10.96±0.01	-10.96±0.01	-10.99±0.04	-11.00±0.03	-11.00±0.03	-11.00±0.03
	-10.97±0.00	-11.17±0.09	-11.04±0.07	-10.98±0.00	-11.24±0.09	-11.11±0.07	-11.06±0.03	-11.01±0.01
	-10.98±0.00	-11.06±0.09	-10.99±0.01	-10.99±0.00	-11.07±0.07	-11.02±0.05	-11.00±0.01	-11.00±0.01
PUMSB	-25.95±0.48	-27.32±0.76	-26.00±0.47	-25.69±0.37	-28.37±0.74	-27.03±0.65	-26.18±0.57	-25.62±0.27
	-26.14±0.52	-27.42±0.88	-26.18±0.46	-25.90±0.41	-28.29±0.68	-27.04±0.62	-26.27±0.58	-25.92±0.37
	-27.06±0.00	-27.19±0.06	-27.19±0.06	-27.19±0.06	-26.65±0.47	-26.55±0.50	-26.55±0.50	-26.55±0.50
	-27.06±0.00	-27.19±0.06	-27.19±0.06	-27.19±0.06	-26.66±0.48	-26.56±0.51	-26.56±0.51	-26.56±0.51
DNA	-87.74±0.12	-87.79±0.11	-87.79±0.11	-87.79±0.11	-87.92±0.11	-87.92±0.11	-87.92±0.11	-87.92±0.11
	-87.75±0.12	-87.79±0.08	-87.79±0.08	-87.79±0.08	-87.95±0.10	-87.95±0.10	-87.95±0.10	-87.95±0.10
	-87.46±0.00	-87.53±0.00	-87.53±0.00	-87.53±0.00	-87.54±0.00	-87.54±0.00	-87.54±0.00	-87.54±0.00
	-87.46±0.00	-87.53±0.00	-87.53±0.00	-87.53±0.00	-87.47±0.00	-87.47±0.00	-87.47±0.00	-87.47±0.00
KOSAREK	-11.25±0.07	-11.11±0.06	-11.09±0.09	-11.12±0.08	-11.08±0.08	-11.03±0.05	-11.04±0.06	-11.04±0.06
	-11.38±0.08	-11.35±0.10	-11.31±0.08	-11.34±0.08	-11.34±0.10	-11.29±0.08	-11.29±0.09	-11.29±0.09
	-10.97±0.00	-11.23±0.22	-11.05±0.16	-10.92±0.02	-11.40±0.15	-11.16±0.14	-11.04±0.06	-10.96±0.01
	-11.01±0.00	-11.25±0.14	-11.06±0.09	-10.99±0.01	-11.42±0.14	-11.20±0.14	-11.08±0.05	-11.03±0.00
MSWEB	-10.19±0.02	-10.19±0.02	-10.19±0.02	-10.19±0.02	-10.19±0.02	-10.19±0.02	-10.19±0.02	-10.19±0.02
	-10.19±0.01	-10.19±0.02	-10.19±0.02	-10.19±0.02	-10.20±0.02	-10.20±0.02	-10.20±0.03	-10.20±0.03
	-10.12±0.00	-10.43±0.13	-10.13±0.03	-10.12±0.01	-10.65±0.14	-10.44±0.13	-10.27±0.11	-10.13±0.02
	-10.12±0.00	-10.31±0.07	-10.15±0.02	-10.15±0.00	-10.38±0.06	-10.29±0.06	-10.20±0.05	-10.13±0.01
BOOK	-37.66±0.32	-37.14±0.32	-37.34±0.40	-37.34±0.40	-36.73±0.37	-36.68±0.34	-36.68±0.34	-36.68±0.34
	-37.79±0.19	-37.69±0.21	-37.69±0.21	-37.69±0.21	-37.51±0.29	-37.50±0.21	-37.50±0.21	-37.50±0.21
	-36.35±0.00	-36.00±0.17	-36.00±0.17	-36.00±0.17	-35.90±0.00	-35.90±0.00	-35.90±0.00	-35.90±0.00
	-36.97±0.00	-36.83±0.10	-36.83±0.10	-36.83±0.10	-36.60±0.00	-36.60±0.00	-36.60±0.00	-36.60±0.00

EMOVIE	-62.19±2.13	-60.61±1.69	-60.98±1.87	-61.38±1.99	-59.78±1.95	-59.64±1.92	-59.41±2.01	-59.20±1.21
	-63.91±1.24	-63.37±1.41	-63.46±1.31	-63.46±1.31	-62.86±1.54	-62.62±1.58	-62.72±1.30	-62.74±1.31
	-59.07±0.01	-58.44±0.54	-57.94±0.39	-57.80±0.15	-57.90±0.89	-58.06±0.20	-57.71±0.57	-57.98±0.00
	-60.56±0.00	-60.86±0.24	-60.49±0.06	-60.49±0.06	-59.99±0.65	-60.47±0.22	-60.35±0.00	-60.35±0.00
WEBKB	-163.31±0.35	-162.82±0.88	-162.88±0.88	-162.88±0.88	-162.42±1.10	-162.53±0.92	-162.44±1.00	-162.44±1.00
	-163.43±0.23	-163.33±0.29	-163.33±0.29	-163.33±0.29	-163.33±0.34	-163.33±0.34	-163.33±0.34	-163.33±0.34
	-159.56±0.00	-158.44±0.35	-159.11±0.41	-159.11±0.41	-158.30±0.09	-158.52±0.63	-158.48±0.00	-158.48±0.00
	-161.30±0.00	-161.19±0.22	-161.91±0.09	-161.91±0.09	-160.93±0.29	-161.26±0.00	-161.26±0.00	-161.26±0.00
ROUTERS	-94.05±0.69	-93.45±1.15	-93.69±1.13	-93.71±1.11	-93.00±1.34	-93.13±0.80	-93.12±0.84	-93.12±0.84
	-94.32±0.33	-94.32±0.51	-94.29±0.37	-94.29±0.37	-94.29±0.48	-94.29±0.48	-94.29±0.48	-94.29±0.48
	-88.13±0.00	-87.94±0.52	-87.95±0.45	-87.95±0.45	-88.89±0.80	-88.18±0.61	-87.91±0.52	-87.91±0.29
	-88.92±0.00	-89.71±0.58	-89.33±0.42	-89.33±0.42	-90.66±0.55	-90.58±0.65	-89.47±0.35	-89.43±0.18
20NEWS	-162.96±1.12	-161.58±0.77	-162.04±0.96	-162.04±0.96	-160.73±0.56	-161.06±0.67	-161.06±0.63	-161.06±0.63
	-164.26±0.61	-163.81±0.53	-163.95±0.44	-163.95±0.44	-163.86±0.84	-163.95±0.85	-164.03±0.87	-164.03±0.87
	-158.38±0.00	-157.31±0.58	-157.55±0.43	-157.44±0.16	-157.79±0.80	-157.86±0.41	-157.71±0.52	-158.16±0.01
	-159.58±0.00	-160.61±1.07	-159.73±0.46	-159.37±0.17	-161.12±0.65	-160.68±0.44	-160.41±0.42	-160.20±0.01
BBC	-261.88±0.34	-261.81±0.28	-261.81±0.28	-261.81±0.28	-261.96±0.48	-261.96±0.48	-261.96±0.48	-261.96±0.48
	-262.04±0.77	-261.96±0.58	-261.96±0.58	-261.96±0.58	-262.24±0.98	-262.24±0.98	-262.24±0.98	-262.24±0.98
	-262.33±0.00	-260.52±0.00	-260.52±0.00	-260.52±0.00	-259.45±1.26	-258.79±0.00	-258.79±0.00	-258.79±0.00
	-262.91±0.00	-261.41±0.00	-261.41±0.00	-261.41±0.00	-260.62±0.00	-260.62±0.00	-260.62±0.00	-260.62±0.00
AD	-16.40±0.04	-16.38±0.11	-16.38±0.11	-16.38±0.11	-16.40±0.12	-16.40±0.12	-16.40±0.12	-16.40±0.12
	-16.62±0.13	-16.61±0.13	-16.61±0.13	-16.61±0.13	-16.66±0.26	-16.66±0.26	-16.66±0.26	-16.66±0.26
	-16.02±0.00	-16.02±0.00	-16.02±0.00	-16.02±0.00	-16.26±0.10	-16.26±0.10	-16.05±0.20	-15.97±0.08
	-16.28±0.00	-16.28±0.00	-16.28±0.00	-16.28±0.00	-16.39±0.07	-16.39±0.07	-16.39±0.07	-16.39±0.07

D.2 TRAINING TIME AND CIRCUIT SIZES

The training time and the circuit size for our models are reported in Table 5.

Table 5: Each row reports δ -values/training times in seconds (top) and circuit sizes (bottom) obtained for learning for the best XPC models reported in Table 4.

DATASET	XPC	XPC _{Det}	XPC ^{SD}	XPC ^{SD} _{Det}
NLTCS	64 / 0.06 2765±458	64 / 0.08 2765±458	64 / 0.05 4335±67	64 / 0.05 4401±67
MSNBC	512 / 0.36 3575±546	512 / 0.71 3575±546	512 / 0.31 4887±125	512 / 0.31 4887±125
KDD	64 / 1.29 21331±1591	128 / 2.09 16089±1312	128 / 0.64 12951±109	128 / 1.02 13040±109
PLANTS	64 / 0.25 31403±1301	128 / 0.18 19138±2135	64 / 0.12 22242±532	128 / 0.1 13960±644
AUDIO	256 / 0.18 20897±1232	512 / 0.12 11396±740	128 / 0.17 29275±42	128 / 0.17 29317±42
JESTER	512 / 0.06 6530±722	256 / 0.11 12584±870	128 / 0.13 25351±69	128 / 0.11 20273±29
NETFLIX	256 / 0.17 20239±1655	512 / 0.12 12510±976	128 / 0.21 39811±57	128 / 0.21 39868±57
ACCIDENTS	256 / 0.15 16252±2742	256 / 0.09 8994±2417	128 / 0.09 11906±15	128 / 0.09 11921±15
RETAIL	512 / 0.04 1098±416	512 / 0.04 1098±416	512 / 0.09 6645±7	512 / 0.09 6651±7
PUMSB-STAR	128 / 0.45 43737±6494	128 / 0.35 32356±7266	128 / 0.1 8703±1442	128 / 0.12 8866±1480
DNA	256 / 0.04 3927±888	256 / 0.04 3927±888	512 / 0.02 2615±1	512 / 0.02 2616±1
KOSAREK	256 / 0.28 13730±3381	256 / 0.29 13617±2872	256 / 0.27 20898±34	256 / 0.3 20938±34
MSWEB	512 / 0.13 2298±714	512 / 0.13 2298±714	512 / 0.27 12130±5	512 / 0.27 12135±5
BOOK	256 / 0.37 11610±4205	256 / 0.35 11227±4306	256 / 0.39 13505±4	256 / 0.38 13678±28
EACHMOVIE	128 / 0.53 19913±8846	256 / 0.31 11375±5785	64 / 1.58 80294±9009	256 / 0.5 21369±3124
WEBKB	256 / 0.89 17026±6503	512 / 0.42 7719±3334	128 / 2.59 61915±13804	512 / 0.87 17122±2
ROUTERS-52	128 / 2.88 55112±18036	512 / 0.65 10479±5049	256 / 2.91 67594±2067	512 / 1.52 36440±5
20NEWS-GRP	512 / 2.11 37138±10983	512 / 2.2 38557±10433	512 / 3.28 65838±2880	512 / 3.34 65881±2884
BBC	512 / 0.49 7755±1930	512 / 0.51 7755±1932	128 / 1.86 36578±6	256 / 0.77 14578±1
AD	64 / 0.94 11661±2170	64 / 1.01 12368±3157	16 / 21.59 298553±10013	128 / 1.69 22093±20

D.3 LEARNING TIME AND CIRCUIT SIZE COMPETITORS

The learning time and circuit size for $\text{XPC}_{\text{Det}}^{\text{SD}}$ compared to those for STRUDEL are reported in Table 6. As reported in [Dang et al., 2020] the learning time and circuit size have been already compared to those obtained with Learn-PSDD showing the superiority of STRUDEL.

Table 6: Learning times (in seconds) and circuit sizes for $\text{XPC}_{\text{Det}}^{\text{SD}}$ (averaged over 40 runs) and STRUDEL best models.

DATASET	$\text{XPC}_{\text{Det}}^{\text{SD}}$	STRUDEL	$\text{XPC}_{\text{Det}}^{\text{SD}}$	STRUDEL
NLTCS	0.05	172.5	4401	13827
MSNBC	0.31	3182.7	4887	35629
KDD	1.02	532.7	13040	16984
PLANTS	0.1	10805.6	13960	454141
AUDIO	0.17	930.9	29317	87090
JESTER	0.11	748.5	20273	78342
NETFLIX	0.21	530.1	39868	45328
ACCIDENTS	0.09	2213.3	11921	99277
RETAIL	0.09	50.3	6651	6609
PUMSB-STAR	0.12	216.8	8866	17395
DNA	0.02	211.2	2616	24106
KOSAREK	0.3	657.4	20938	73086
MSWEB	0.27	224.5	12135	3177
BOOK	0.38	1589.9	13678	140351
EACHMOVIE	0.5	10686.1	21369	988953
WEBKB	0.87	1631.0	17122	84299
ROUTERS-52	1.52	3875.6	36440	198905
20NEWS-GRP	3.34	4425.4	65881	204983
BBC	0.77	317.3	14578	22962
AD	1.69	113.2	22093	18151

E MIXTURE MODELS

E.1 EXPCS

The best average test-set log-likelihoods for every EXPC type are reported in Table 7.

Table 7: Best average test-set log-likelihoods for every EXPC type in the space formed by $l = \{1, 2, 3\}$, $k = \{2, 3, 4, 8\}$, $\delta = \{16, 32, 64, 128, 256, 512\}$ and $M = \{2, 5, 10, 15, 20, 25, 30, 40\}$. For every dataset, each row corresponds (from top to down) to: EXPC, EXPC_{Det} , $\text{EXPC}_{\text{Det}}^{\text{SD}}$ and $\text{EXPC}_{\text{Det}}^{\text{SD}}$.

	$l = 1$		$l = 2$		$l = 3$			
	$k = 2$	$k = 2$	$k = 3$	$k = 4$	$k = 2$	$k = 3$	$k = 4$	$k = 8$
NLTCS	-6.00	-6.30	-6.04	-6.00	-6.54	-6.20	-6.10	-6.01
	-6.00	-6.26	-6.04	-6.00	-6.48	-6.18	-6.09	-6.01
	-	-6.43	-6.09	-6.05	-6.50	-6.29	-6.15	-6.10
	-	-6.38	-6.08	-6.05	-6.42	-6.23	-6.11	-6.08
MSNBC	-6.12	-6.40	-6.19	-6.12	-6.56	-6.38	-6.25	-6.12
	-6.12	-6.35	-6.17	-6.12	-6.46	-6.33	-6.21	-6.12
	-	-6.43	-6.22	-6.18	-6.54	-6.44	-6.33	-6.19
	-	-6.38	-6.20	-6.17	-6.48	-6.38	-6.29	-6.19
KDD	-2.13	-2.14	-2.13	-2.13	-2.13	-2.13	-2.13	-2.13
	-2.13	-2.13	-2.13	-2.14	-2.13	-2.13	-2.13	-2.13
	-	-2.15	-2.17	-2.19	-2.17	-2.15	-2.16	-2.18
	-	-2.15	-2.17	-2.20	-2.17	-2.16	-2.16	-2.19
PLANTS	-12.41	-13.14	-12.53	-12.35	-13.91	-13.03	-12.59	-12.24
	-12.41	-13.11	-12.51	-12.35	-13.78	-12.96	-12.57	-12.24
	-	-14.73	-14.19	-14.55	-14.82	-14.45	-14.39	-14.46
	-	-14.61	-14.21	-14.62	-14.76	-14.46	-14.43	-14.50

		-39.84	-40.03	-39.80	-39.78	-40.46	-40.06	-39.89	-39.75
	AUDIO	-39.84	-40.12	-39.86	-39.80	-40.53	-40.14	-39.97	-39.77
		-	-40.96	-40.96	-41.74	-41.66	-40.91	-40.96	-41.91
		-	-41.00	-41.05	-41.76	-41.65	-40.97	-41.07	-42.01
		-52.72	-52.89	-52.65	-52.67	-53.27	-52.89	-52.73	-52.60
	JESTER	-52.72	-52.96	-52.75	-52.68	-53.35	-52.97	-52.84	-52.63
		-	-53.43	-53.52	-54.59	-54.01	-53.49	-53.41	-54.68
		-	-53.51	-53.68	-54.63	-54.03	-53.61	-53.59	-54.76
		-56.43	-56.79	-56.50	-56.38	-57.39	-56.77	-56.64	-56.41
	NETFLIX	-56.43	-56.93	-56.58	-56.40	-57.41	-56.88	-56.75	-56.43
		-	-57.66	-57.68	-58.46	-58.18	-57.67	-57.58	-58.60
		-	-57.68	-57.80	-58.46	-58.14	-57.73	-57.69	-58.62
		-29.63	-30.63	-29.76	-29.54	-31.33	-30.35	-29.97	-29.55
	RETAILACCIDENTS	-29.63	-30.58	-29.74	-29.53	-31.28	-30.31	-29.93	-29.54
		-	-31.71	-31.36	-31.36	-33.12	-31.45	-31.02	-31.08
		-	-31.60	-31.36	-31.36	-31.99	-31.32	-30.99	-31.08
		-10.83	-10.82	-10.82	-10.82	-10.81	-10.81	-10.81	-10.81
	RETAILACCIDENTS	-10.83	-10.81	-10.81	-10.81	-10.81	-10.80	-10.79	-10.79
		-	-10.94	-10.96	-10.97	-11.00	-10.94	-10.94	-10.96
		-	-10.91	-10.95	-10.98	-10.90	-10.91	-10.95	-10.97
		-23.78	-24.87	-23.86	-23.58	-25.61	-24.56	-23.90	-23.55
	PUMSB	-23.78	-24.78	-23.83	-23.57	-25.49	-24.44	-23.85	-23.53
		-	-27.21	-27.21	-27.21	-26.16	-26.06	-26.06	-26.06
		-	-27.21	-27.21	-27.21	-26.16	-26.05	-26.05	-26.05
		-84.83	-85.75	-85.47	-85.47	-86.25	-85.69	-85.62	-85.53
	DNA	-84.83	-85.78	-85.48	-85.48	-86.24	-85.73	-85.62	-85.54
		-	-87.25	-87.53	-87.53	-86.72	-86.61	-87.54	-87.54
		-	-86.85	-87.25	-87.25	-85.09	-86.30	-87.47	-87.47
		-10.68	-10.64	-10.66	-10.66	-10.65	-10.63	-10.61	-10.62
	KOSAREK	-10.68	-10.66	-10.66	-10.66	-10.67	-10.64	-10.63	-10.63
		-	-10.87	-10.77	-10.89	-11.04	-10.85	-10.81	-10.96
		-	-10.86	-10.81	-10.95	-11.01	-10.87	-10.85	-11.03
		-10.09	-9.94	-9.90	-9.90	-9.84	-9.78	-9.78	-9.77
	MSWEB	-10.09	-9.94	-9.90	-9.90	-9.84	-9.78	-9.78	-9.78
		-	-10.09	-9.93	-10.07	-10.24	-10.05	-9.97	-10.05
		-	-10.02	-9.94	-10.12	-10.06	-9.97	-9.93	-10.09
		-35.72	-35.16	-35.23	-35.23	-34.85	-35.02	-35.02	-35.02
	BOOK	-35.63	-35.40	-35.40	-35.40	-34.98	-35.16	-35.16	-35.16
		-	-34.86	-35.37	-35.68	-34.75	-35.08	-35.41	-35.90
		-	-35.25	-36.28	-36.28	-35.14	-35.71	-36.60	-36.60
		-54.50	-53.78	-53.78	-53.87	-53.08	-53.23	-53.15	-52.93
	EMOIE	-54.65	-54.06	-54.13	-54.30	-53.35	-53.65	-53.44	-53.45
		-	-55.35	-56.05	-57.27	-55.23	-54.82	-55.42	-57.98
		-	-55.88	-57.39	-60.15	-55.26	-55.70	-56.47	-60.35
		-153.97	-153.92	-155.16	-155.23	-154.03	-153.93	-154.17	-154.17
	WEBKB	-154.29	-154.36	-155.33	-155.49	-154.42	-154.49	-154.66	-154.66
		-	-153.62	-157.22	-157.22	-153.67	-154.17	-156.28	-157.89
		-	-154.45	-159.88	-159.88	-154.23	-155.58	-158.22	-161.26
		-85.61	-85.32	-85.71	-85.28	-85.07	-84.82	-85.31	-85.04
	ROUTERS	-85.89	-85.96	-86.12	-85.72	-85.53	-85.27	-85.74	-85.65
		-	-85.19	-85.17	-86.48	-85.38	-84.70	-85.02	-86.94
		-	-85.51	-85.93	-87.25	-85.49	-85.09	-85.96	-88.61
		-154.48	-153.98	-154.47	-154.35	-153.85	-154.13	-154.45	-154.48
	2ONEWS	-154.60	-154.54	-154.80	-154.71	-154.64	-154.91	-155.09	-155.07
		-	-153.78	-155.10	-156.86	-153.75	-154.02	-154.52	-157.26
		-	-154.35	-156.26	-158.63	-154.21	-154.56	-155.49	-160.10
		-237.96	-242.58	-243.68	-243.72	-250.91	-251.07	-251.33	-251.33
	BBC	-239.01	-247.14	-246.98	-247.43	-251.35	-251.64	-251.85	-251.85
		-	-248.43	-253.19	-260.52	-248.34	-250.16	-251.50	-258.79
		-	-249.46	-255.95	-261.41	-248.79	-252.21	-255.59	-260.62
		-15.90	-14.39	-14.23	-14.23	-13.95	-13.61	-13.60	-13.60
	AD	-15.94	-14.28	-14.14	-14.14	-13.78	-13.45	-13.45	-13.45
		-	-16.03	-16.03	-16.03	-15.66	-15.59	-15.50	-15.67
		-	-16.31	-16.31	-16.31	-15.59	-15.73	-15.68	-15.87

E.2 COMPETITORS

The average test-set log-likelihoods for best EXPCs are reported in Table 8.

Table 8: Average test-set log-likelihoods for best EXPCs, SOTA ensembling techniques and single model competitors. In the first half of the table the best results for ensembles of bagged (CNet⁴⁰) and boosted (CNet⁴⁰_{boost}) entropy based CNets taken from [Rahman and Gogate, 2016], ensembles of bagged CNets learned with dCSN and ensemble of extremely randomized CNets as in [Di Mauro et al., 2017]. The remaining columns report the comparison to other models employing more sophisticated models such as ID-SPN [Rooshenas and Lowd, 2014], RAT-SPN [Peharz et al., 2019], and Learn-SPN [Gens and Domingos, 2013].

DATASET	ensembles					competitors		
	EXPC	XCNet ⁴⁰	CNet ⁴⁰	CNet ⁴⁰ _{boost}	dCSN ⁴⁰	RAT	ID-SPN	Learn-SPN
NLTCS	-6.00	-6.00	-6.00	-6.01	-6.00	-6.01	-6.02	-6.11
MSNBC	-6.12	-6.06	-6.08	-6.15	-6.05	-6.04	-6.04	-6.11
KDD	-2.12	-2.13	-2.14	-2.15	-2.15	-2.13	-2.13	-2.18
PLANTS	-12.24	-11.99	-12.32	-12.67	-12.59	-13.44	-12.54	-12.99
AUDIO	-39.75	-39.77	-40.09	-39.84	-40.19	-39.96	-39.79	-40.50
JESTER	-52.51	-52.65	-52.88	-52.82	-52.99	-52.97	-52.86	-75.98
NETFLIX	-56.38	-56.38	-56.55	-56.44	-56.69	-56.85	-56.36	-57.33
ACCIDENTS	-29.54	-29.31	-29.88	-29.45	-29.27	-35.49	-26.98	-30.04
RETAIL	-10.79	-10.93	-10.84	-10.81	-11.17	-10.91	-10.85	-11.04
PUMSB-STAR	-23.53	-23.44	-23.98	-23.46	-23.78	-32.53	-22.41	-24.78
DNA	-84.83	-84.96	-81.07	-85.67	-85.95	-97.23	-81.21	-82.52
KOSAREK	-10.61	-10.72	-10.74	-10.60	-10.97	-10.89	-10.60	-10.99
MSWEB	-9.78	-9.66	-9.77	-9.74	-9.93	-10.12	-9.73	-10.25
BOOK	-34.14	-36.35	-35.55	-34.46	-37.38	-34.68	-34.14	-35.89
EACHMOVIE	-52.38	-51.72	-53.00	-51.53	-54.14	-53.63	-51.51	-52.49
WEBKB	-152.44	-153.01	-153.12	-152.53	-155.47	-157.53	-151.84	-158.21
ROUTERS-52	-84.70	-84.05	-83.71	-83.69	-86.19	-87.37	-83.35	-85.07
20NEWS-GRP	-152.26	-153.89	-156.09	-153.12	-156.46	-152.06	-151.47	-155.93
BBC	-237.96	-238.47	-237.42	-247.01	-248.84	-252.14	-248.93	-250.69
AD	-13.45	-14.20	-15.28	-14.36	-15.55	-48.47	-19.05	-19.73
AVG. RANK	2.00	2.25	3.25	2.85	4.25			
	2.75	3.15	4.3	3.8	5.85	6	2.5	6.9