

SUBTRACTIVE MIXTURE MODELS VIA SQUARING: REPRESENTATION AND LEARNING

Lorenzo Loconte^{1*} Aleksanteri M. Sladek² Stefan Mengel³
 Martin Trapp² Arno Solin² Nicolas Gillis⁴ Antonio Vergari¹

¹ School of Informatics, University of Edinburgh, UK

² Department of Computer Science, Aalto University, Finland

³ University of Artois, CNRS, Centre de Recherche en Informatique de Lens (CRIL), France

⁴ Department of Mathematics and Operational Research, Université de Mons, Belgium

ABSTRACT

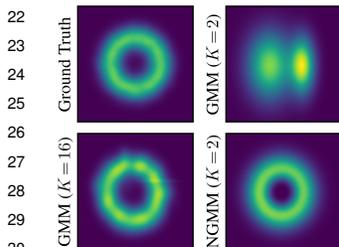
Mixture models are traditionally represented and learned by *adding* several distributions as components. Allowing mixtures to *subtract* probability mass or density can drastically reduce the number of components needed to model complex distributions. However, learning such subtractive mixtures while ensuring they still encode a non-negative function is challenging. We investigate how to learn and perform inference on deep subtractive mixtures by *squaring* them. We do this in the framework of probabilistic circuits, which enable us to represent tensorized mixtures and generalize several other subtractive models. We theoretically prove that the class of squared circuits allowing subtractions can be exponentially more expressive than traditional additive mixtures; and, we empirically show this increased expressiveness on a series of real-world distribution estimation tasks.

1 INTRODUCTION

Finite mixture models (MMs) are a staple in probabilistic machine learning, as they offer a simple and elegant solution to model complex distributions by blending simpler ones in a linear combination (McLachlan et al., 2019). The classical recipe to design MMs is to compute a *convex combination* over input components. That is, a MM representing a probability distribution p over a set of random variables $\mathbf{X} = \{X_1, X_2, \dots, X_D\}$ is usually defined as

$$p(\mathbf{X}) = \sum_{i=1}^K w_i p_i(\mathbf{X}), \quad \text{with } w_i \geq 0, \quad \sum_{i=1}^K w_i = 1, \quad (1)$$

where w_i are the mixture parameters and each component p_i is a mass or density function. This is the case for widely-used MMs such as Gaussian mixture models (GMMs) and hidden Markov models (HMMs) but also mixtures of generative models such as normalizing flows (Papamakarios et al., 2021) and deep mixture models such as probabilistic circuits (PCs, Vergari et al., 2019b).



The convexity constraint in Eq. (1) is the simplest *sufficient* condition to ensure that p is a non-negative function and integrates to 1,¹ i.e., is a valid probability distribution, and is often assumed in practice. However, this implies that the components p_i can only be combined in an additive manner and as such it can greatly impact their ability to estimate a distribution efficiently. For instance, consider approximating distributions having “holes” in their domain, such as the simple 2-dimensional ring distribution on the left (ground truth). A classical additive MM such a GMM would ultimately recover it, as it is a universal approximator of density functions (Nguyen et al., 2019), but only by employing an unnecessarily high number of components. A MM allowing negative mixture weights, i.e., $w_i < 0$, would instead require only two components, as it can *subtract* one outer Gaussian density from an inner

*Corresponding author, l.loconte@sms.ed.ac.uk

¹Across the paper we will abuse the term integration to also refer to summation in case of discrete variables.

34 one (NGMM). We call these MMs subtractive or *non-monotonic* MMs (NMMs), as opposed to their
35 classical additive counterpart, called *monotonic* MMs (Shpilka & Yehudayoff, 2010).

36 The challenge with NMMs is ensuring that the modeled $p(\mathbf{X})$ is a valid distribution, as the convex-
37 ity constraint does not hold anymore. This problem has been investigated in the past in a number
38 of ways, in its simplest form by imposing ad-hoc constraints over the mixture parameters w_i , de-
39 rived for simple components such as Gaussian and Weibull distributions (Zhang & Zhang, 2005;
40 Rabusseau & Denis, 2014; Jiang et al., 1999). However, different families of components would
41 require formulating different constraints, whose closed-form existence is not guaranteed.

42 In this paper, we study a more general principle to design NMMs that circumvents the aforemen-
43 tioned limitation while ensuring non-negativity of the modeled function: *squaring the encoded lin-*
44 *ear combination*. For example, the NGMM above is a squared combination of Gaussian densities
45 with negative mixture parameters. We theoretically investigate the expressive efficiency of squared
46 NMMs, i.e., their expressiveness w.r.t. their model size, and show how to effectively represent and
47 learn them in practice. Specifically, we do so in the framework of PCs, tractable models general-
48 izing classical shallow MMs into deep MMs represented as structured neural networks. Deep PCs
49 are already more expressive efficient than shallow MMs as they compactly encode a mixture with
50 an exponential number of components (Vergari et al., 2019b; Choi et al., 2020). However, they are
51 classically represented with non-negative parameters, hence being restricted to encode deep but ad-
52 ditive MMs. Instead, as a main theoretical contribution we prove that *our squared non-monotonic*
53 *PCs (NPC²s) can be exponentially more parameter-efficient than their monotonic counterparts.*

54 **Contributions.** **i)** We introduce a general framework to represent NMMs via squaring (Sec. 2),
55 within the language of tensorized PCs (Mari et al., 2023), and show how NPC²s can be effectively
56 learned and used for tractable inference (Sec. 3). **ii)** We show how NPC²s generalize not only mono-
57 tonic PCs but other apparently different models allowing negative parameters that have emerged in
58 different literatures, such as square root of density models in signal processing (Pinheiro & Vi-
59 dakovic, 1997), positive semi-definite (PSD) models in kernel methods (Rudi & Ciliberto, 2021),
60 and Born machines from quantum mechanics (Orús, 2013) (Sec. 4). This allows us to understand
61 why they lead to tractable inference via the property-oriented framework of PCs. **iii)** We derive an
62 exponential lower bound over the size of monotonic PCs to represent functions that can be com-
63 pactly encoded by one NPC² (Sec. 4.1), hence showing that NPC²s (and thus the aforementioned
64 models) can be more expressive for a given size. Finally, **iv)** we provide empirical evidence (Sec. 5)
65 that NPC²s can approximate distributions better than monotonic PCs for a variety of experimen-
66 tal settings involving learning from real-world data and distilling intractable models such as large
67 language models to unlock tractable inference (Zhang et al., 2023).

68 2 SUBTRACTIVE MIXTURES VIA SQUARING

69 We start by formalizing how to represent *shallow* NMMs by *squaring* non-convex combinations of
70 K simple functions. Like exponentiation in energy-based models (LeCun et al., 2006), squaring
71 ensures the non-negativity of our models, but differently from it, allows to tractably renormalize
72 them. A squared NMM encodes a (possibly unnormalized) distribution $c^2(\mathbf{X})$ over variables \mathbf{X} as

$$c^2(\mathbf{X}) = \left(\sum_{i=1}^K w_i c_i(\mathbf{X}) \right)^2 = \sum_{i=1}^K \sum_{j=1}^K w_i w_j c_i(\mathbf{X}) c_j(\mathbf{X}), \quad (2)$$

73 where c_i are the learnable components and the mixture parameters $w_i \in \mathbb{R}$ are unconstrained, as
74 opposed to Eq. (1). Squared NMMs can therefore represent $\binom{K+1}{2}$ components within the same pa-
75 rameter budget of K components of an additive MM. Each component of a squared NMM computes
76 a *product of experts* $c_i(\mathbf{X}) c_j(\mathbf{X})$ (Hinton, 2002) allowing negative parameters $2w_i w_j$ if $i \neq j$, and
77 $c_i^2(\mathbf{X})$ with w_i^2 otherwise. Fig. 1 shows a concrete example of this construction, which constitutes
78 the simplest NPC² we can build (see Sec. 3), i.e., comprising a single layer and having depth one.

79 **Tractable marginalization.** Analogously to traditional MMs, squared NMMs support tractable
80 marginalization and conditioning, if their component distributions do as well. The distribution en-
81 coded by $c^2(\mathbf{X})$ can be normalized to compute a valid probability distribution $p(\mathbf{X}) = c^2(\mathbf{X})/Z$,
82 by computing its partition function Z as

$$Z = \int c^2(\mathbf{x}) d\mathbf{x} = \sum_{i=1}^K \sum_{j=1}^K w_i w_j \int c_i(\mathbf{x}) c_j(\mathbf{x}) d\mathbf{x}. \quad (3)$$

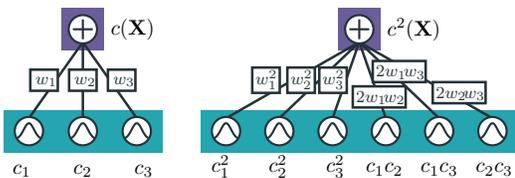


Figure 1: **Shallow MMs and squared NMMs represented as PCs**, mapped to a computational graph having input components and a weighted sum unit as output. Squaring a mixture with $K = 3$ components (left) can yield more components that share parameters (right).

83 Computing Z translates to evaluating $\binom{K+1}{2}$ integrals over products of components $c_i(\mathbf{X})c_j(\mathbf{X})$.
 84 More generally, marginalizing any subset of variables in \mathbf{X} can be done in $\mathcal{O}(K^2)$. This how-
 85 ever implies that the components c_i are chosen from a family of functions such that their product
 86 $c_i(\mathbf{X})c_j(\mathbf{X})$ can be tractably integrated, and Z is non-zero and finite. This is true for many para-
 87 metric families, including exponential families (Seeger, 2005). For instance, the product of two
 88 Gaussian or two categorical distributions is another Gaussian (Rasmussen & Williams, 2005) or
 89 categorical up to a multiplicative factor, which can be computed in polynomial time.

90 **A wider choice of components.** Note that we do not require each c_i to model a probability distri-
 91 bution, e.g., we might have $c_i(\mathbf{x}) < 0$. This allows us to employ more expressive tractable functions
 92 as base components in squared NMMs such as splines (see App. E for details) or potentially small
 93 neural networks (see discussion in App. G). However, if the components are already flexible enough
 94 there might not be an increase in expressiveness when mixing them in a linear combination or squar-
 95 ing them. E.g., a simple categorical distribution can already capture any discrete distribution with
 96 finite support and a (subtractive) mixture thereof might not yield additional benefits besides being
 97 easier to learn. An additive mixture of Binomials is instead more expressive than a single Binomial,
 98 but expected to be less expressive than its subtractive version (as illustrated in Sec. 5).

99 **Learning squared NMMs.** The canonical way to learn traditional MMs (Eq. (1)) is by maximum-
 100 likelihood estimation (MLE), i.e., by maximizing $\sum_{\mathbf{x} \in \mathcal{D}} \log p(\mathbf{x})$ where \mathcal{D} is a set of independent
 101 and identically distributed (i.i.d.) samples. For squared NMMs, the MLE objective is

$$\sum_{\mathbf{x} \in \mathcal{D}} \log (c^2(\mathbf{x})/Z) = -|\mathcal{D}| \log Z + 2 \sum_{\mathbf{x} \in \mathcal{D}} \log |c(\mathbf{x})|, \tag{4}$$

102 where $c(\mathbf{x}) = \sum_{i=1}^K w_i c_i(\mathbf{x})$. Unlike other NMMs mentioned in Sec. 1, we do not need to derive
 103 additional closed-form constraints for the parameters to preserve non-negativity. Although mate-
 104 rializing the squared mixture having $\binom{K+1}{2}$ components is required to compute Z as in Eq. (3),
 105 evaluating $\log |c(\mathbf{x})|$ is linear in K . Hence, we can efficiently perform batched stochastic gradient-
 106 based optimization and compute Z just once per batch.

107 3 SQUARING DEEP MIXTURE MODELS

108 So far, we dealt with mixtures that are shallow, i.e., that can be represented as simple computational
 109 graphs with a single weighted sum unit (e.g., Fig. 1). We now generalize them in the framework
 110 of PCs (Vergari et al., 2019b; Choi et al., 2020; Darwiche, 2001) as they offer a property-driven
 111 language to model structured neural networks which allow tractable inference. PCs enable us to
 112 encode an exponential number of mixture components in a compact but deep computational graph.

113 PCs are usually defined in terms of scalar computational units: sum, product and input (see App. A).
 114 Following Vergari et al. (2019a); Mari et al. (2023), we instead formalize them as tensorized compu-
 115 tational graphs. That is, we group several computational units together in layers, whose advantage is
 116 twofold. First, we are able to derive a simplified tractable algorithm for squaring that requires only
 117 linear algebra operations and benefits from GPU acceleration (Alg. 1). Second, we can more easily
 118 generalize many recent PC architectures (Peharz et al., 2020b;a; Liu & Van den Broeck, 2021), as
 119 well as other tractable tensor representations (Sec. 4). Fig. A.1 illustrates how scalar computational
 120 units are mapped to tensorized layers. We start by defining deep computational graphs that can
 121 model possibly negative functions, simply named *circuits* (Vergari et al., 2021).

122 **Definition 1** (Tensorized circuit). A *tensorized circuit* c is a parameterized computational graph
 123 encoding a function $c(\mathbf{X})$ and comprising of three kinds of layers: *input*, *product* and *sum*. Each
 124 layer comprises computational units defined over the same set of variables, also called its *scope*, and
 125 every non-input layer receives input from one or more layers. The scope of each non-input layer is

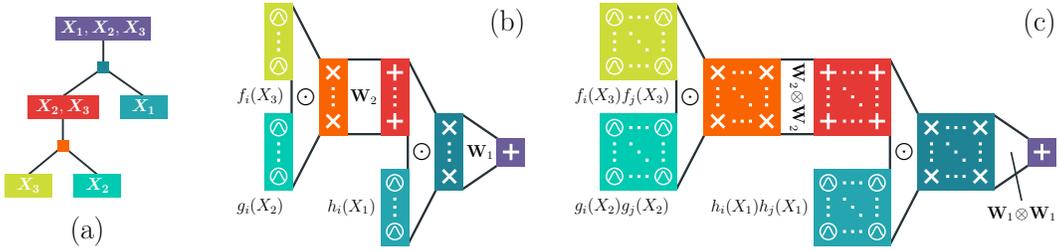


Figure 2: **Squaring tensorized structured-decomposable circuits reduces to squaring layers**, depicted as colored boxes of \odot (input), \times (product), and $+$ (sum). Connections to a sum layer are labeled by the matrix parameterizing the layer, while connections to product layers are labeled by the Hadamard product sign (see also Fig. A.1). A tensorized structured-decomposable circuit (b) over three variables defined from the RG in (a) is squared in (c) by recursively squaring each layer via Alg. 1. Squared layers contain a quadratic number of units, but still output vectors.

126 the union of the scope of its inputs, and the scope of the output layer computing $c(\mathbf{X})$ is \mathbf{X} . Each
 127 input layer ℓ has scope $\mathbf{Y} \subseteq \mathbf{X}$ and computes a collection of K functions $f_i(\mathbf{Y}) \in \mathbb{R}$, i.e., ℓ outputs
 128 a K -dimensional vector. Each product layer ℓ computes an Hadamard (or element-wise) product
 129 over the N layers it receives as input, i.e., $\ell = \odot_{i=1}^N \ell_i$. A sum layer with S sum units and receiving
 130 input from a previous layer $\ell \in \mathbb{R}^K$, is parameterized by $\mathbf{W} \in \mathbb{R}^{S \times K}$ and computes $\mathbf{W}\ell$.

131 Fig. 2b shows a deep circuit in tensorized form. To model a distribution via circuits we first require
 132 that the output of the computational graph is non-negative. We call such a circuit a PC. Similarly
 133 to shallow additive MM (Eq. (1)), a sufficient condition to ensure non-negativity of the output is
 134 make the PC monotonic, i.e., to parameterize all sum layers with non-negative matrices and to
 135 restrict input layers to encode non-negative functions (e.g., probability mass or density functions).
 136 So far, monotonic PCs have been the canonical way to represent and learn PCs (App. G). In Def. 1
 137 we presented product layers computing Hadamard products only, to simplify notation and as this
 138 implementation choice is commonly used in many existing PC architectures (Darwiche, 2009; Liu
 139 & Van den Broeck, 2021; Mari et al., 2023). We generalize our treatment of PCs in Def. A.6 to deal
 140 with another popular product layer implementation: Kronecker products (Peharz et al., 2020b;a;
 141 Mari et al., 2023). Our results still hold for both kinds of product layers, if not specified otherwise.

142 3.1 BUILDING TRACTABLE CIRCUITS FOR MARGINALIZATION

143 Deep PCs can be renormalized and marginalize out any subset of \mathbf{X} in a single feed-forward pass
 144 if they are *smooth* and *decomposable*, i.e., each sum layer receives inputs from layers whose units
 145 are defined over the same scopes, and each product layer receives inputs from layers whose scopes
 146 are pairwise disjoint, respectively. See Prop. A.1 for more background. Sum layers in our Def. 1
 147 guarantee smoothness by design as they have exactly one input. A simple way to ensure decompos-
 148 ability is to create a circuit that follows a *hierarchical scope partitioning of variables* \mathbf{X} , also called
 149 a *region graph*, which is formalized next.

150 **Definition 2** (Region graph (Dennis & Ventura, 2012)). Given a set of variables \mathbf{X} , a *region graph*
 151 (RG) is a bipartite and rooted graph whose nodes are either *regions*, denoting subsets \mathcal{R} of \mathbf{X} , or
 152 *partitions* specifying how a region is partitioned into other regions.

153 Fig. 2a shows an example of a RG. Given a RG, we can build a smooth and decomposable tensorized
 154 circuit as follows. First, we parameterize regions $\mathcal{R} \subseteq \mathbf{X}$ that are not further partitioned with an
 155 input layer encoding some functions over variables in \mathcal{R} . Then, we parameterize each partitioning
 156 $\{\mathcal{R}_i\}_{i=1}^N$ with a product layer having as inputs one layer for each \mathcal{R}_i . Each product layer is then
 157 followed by a sum layer. Figs. 2a and 2b illustrate such a construction by color-coding regions and
 158 corresponding layers. As we will show in Sec. 3.2, this also provides us a clean recipe to efficiently
 159 square a deep circuit. The literature on PCs provides several ways to build RGs (Peharz et al.,
 160 2020b;a; Mari et al., 2023). In our experiments (Sec. 5), we recursively partition sets of variables
 161 randomly until no further partitioning is possible (Peharz et al., 2020b). Moreover, we experiment
 162 with RGs that partitions variables one by one (e.g., the one in Fig. 2a), as they are related to other
 163 classes of models (see Sec. 4). App. F details how to construct RGs.

164 3.2 SQUARING DEEP TENSORIZED CIRCUITS

165 **(Squared negative) MMs as circuits.** It is easy to see that traditional shallow MMs (Eq. (1)) can
 166 be readily represented as tensorized smooth and decomposable PCs consisting of an input layer
 167 encoding K components p_i followed by a sum layer, which is parameterized by a non-negative row-
 168 vector $\mathbf{W} \in \mathbb{R}_+^{1 \times K}$ whose entries sum up to one. Squared NMMs (Eq. (2)) can be represented in a
 169 similar way, as they can be viewed as mixtures over an increased number of components (see Fig. 1
 170 and Fig. A.1), where the sum layer is parameterized by real entries, instead. Next, we discuss how
 171 to square *deep* tensorized circuits as to retrieve our NPC²s model class.

172 **Squaring (and renormalizing) tensorized circuits.** The challenge of squaring a tensorized non-
 173 monotonic circuit c (potentially encoding a negative function) is guaranteeing c^2 to be representable
 174 as a smooth and decomposable PC with polynomial size, as these two properties are necessary
 175 conditions to being able to renormalize c^2 efficiently and in a single feed-forward pass (Choi et al.,
 176 2020). In general, even squaring a decomposable circuit while preserving decomposability of the
 177 squared circuit is a #P-hard problem (Shen et al., 2016; Vergari et al., 2021). Fortunately, it is
 178 possible to obtain a decomposable representation of c^2 efficiently for circuits c that are *structured-*
 179 *decomposable* (Pipatsrisawat & Darwiche, 2008; Vergari et al., 2021). Intuitively, in a tensorized
 180 structured-decomposable circuit all product layers having the same scope $\mathbf{Y} \subseteq \mathbf{X}$ decompose \mathbf{Y}
 181 over their input layers in the exact same way. We formalize this property in the Appendix in Def. A.3.

182 Tensorized circuits satisfying this property by design can be easily constructed by stacking layers
 183 conforming to a RG, as discussed before, and requiring that such a RG is a *tree*, i.e., in which there is
 184 a single way to partition each region, and whose input regions do not have overlapping scopes. E.g.,
 185 the RG in Fig. 2a is a tree RG. From here on, w.l.o.g. we assume our tree RGs to be binary trees, i.e.,
 186 they partition each region into two other regions only. Given a tensorized structured-decomposable
 187 circuit c defined on such a tree RG, Alg. 1 efficiently constructs a smooth and decomposable ten-
 188 sorized circuit c^2 . Moreover, let L be the number of layers and M the maximum time required to
 189 evaluate one layer in c , then the following proposition holds.

190 **Proposition 1** (Tractable marginalization of squared circuits). *Let c be a tensorized structured-*
 191 *decomposable circuit where the products of functions computed by each input layer can be tractably*
 192 *integrated. Any marginalization of c^2 obtained via Alg. 1 requires time and space $\mathcal{O}(L \cdot M^2)$.*

193 See App. B.2 for a proof. In a nutshell, this is possible because Alg. 1 recursively squares each
 194 layer ℓ in c such as $\ell^2 = \ell \otimes \ell$ in c^2 , where \otimes denotes the Kronecker product of two vectors.² Our
 195 tensorized treatment of circuits allows for a much more compact version of the more general algo-
 196 rithm proposed in Vergari et al. (2021) which was defined in terms of squaring scalar computational
 197 units. At the same time, it lets us derive a tighter worst-case upper-bound than the one usually re-
 198 ported for squaring structured-decomposable circuits (Pipatsrisawat & Darwiche, 2008; Choi et al.,
 199 2015; Vergari et al., 2021), which is the squared number of computations in the whole computational
 200 graph, or $\mathcal{O}(L^2 \cdot M^2)$. Note that materializing c^2 is needed when we want to efficiently compute
 201 the normalization constant Z of c^2 or marginalizing any subset of variables. As such, when learning
 202 by MLE (Eq. (4)) and by batched gradient descent, we need to evaluate c^2 only once per batch, thus
 203 greatly amortizing its cost. In App. C, we investigate the time and memory costs of learning NPC²s
 204 having different size and on different data set dimensionalities. Finally, tractable marginalization
 205 enables tractable sampling from the distribution modeled by NPC²s, as we discuss in App. A.2.

206 3.3 NUMERICALLY STABLE INFERENCE AND LEARNING

207 Renormalizing deep PCs can easily lead to underflows and/or overflows. In monotonic PCs, this
 208 is usually addressed by performing computations in log-space and utilizing the log-sum-exp trick
 209 (Blanchard et al., 2021). However, this is not applicable to non-monotonic PCs as intermediate
 210 layers can compute negative values. Therefore, we instead evaluate circuits by propagating the log-
 211 arithm of absolute values and the sign values of the outputs of each layer. Then, sum layers are
 212 evaluated with a *sign-aware* version of the log-sum-exp trick. A similar idea has been already ap-
 213 plied to evaluate expectations of negative functions with monotonic PCs (Mauá et al., 2018; Correia
 214 & de Campos, 2019). App. D extends it to tensorized non-monotonic circuits.

²In Alg. B.2 we provide a generalization of Alg. 1 to square Kronecker product layers.

Algorithm 1 squareTensorizedCircuit(ℓ, \mathcal{R})**Input:** A tensorized circuit having output layer ℓ and defined on a tree RG rooted by \mathcal{R} .**Output:** The tensorized squared circuit defined on the same tree RG having ℓ^2 as output layer computing $\ell \otimes \ell$.

```

1: if  $\ell$  is an input layer then
2:    $\ell$  computes  $K$  functions  $f_i(\mathcal{R})$ 
3:   return An input layer  $\ell^2$  computing all  $K^2$ 
4:     product combinations  $f_i(\mathcal{R})f_j(\mathcal{R})$ 
5: else if  $\ell$  is a product layer then
6:    $\{(\ell_i, \mathcal{R}_i), (\ell_{ii}, \mathcal{R}_{ii})\} \leftarrow \text{getInputs}(\ell, \mathcal{R})$ 
7:    $\ell_i^2 \leftarrow \text{squareTensorizedCircuit}(\ell_i, \mathcal{R}_i)$ 
8:    $\ell_{ii}^2 \leftarrow \text{squareTensorizedCircuit}(\ell_{ii}, \mathcal{R}_{ii})$ 
9:   return  $\ell_i^2 \odot \ell_{ii}^2$ 
10: else  $\triangleright \ell$  is a sum layer
11:    $\{(\ell_i, \mathcal{R}_i)\} \leftarrow \text{getInputs}(\ell, \mathcal{R})$ 
12:    $\ell_i^2 \leftarrow \text{squareTensorizedCircuit}(\ell_i, \mathcal{R}_i)$ 
13:    $\mathbf{W} \in \mathbb{R}^{S \times K} \leftarrow \text{getParameters}(\ell)$ 
14:    $\mathbf{W}' \in \mathbb{R}^{S^2 \times K^2} \leftarrow \mathbf{W} \otimes \mathbf{W}$ 
15:   return  $\mathbf{W}' \ell_i^2$ 

```

215 4 EXPRESSIVENESS OF NPC²S AND RELATIONSHIP TO OTHER MODELS

216 Circuits have been used as the “lingua franca” to represent apparently different tractable model
217 representations (Darwiche & Marquis, 2002; Shpilka & Yehudayoff, 2010), and to investigate their
218 ability to exactly represent certain function families with only a polynomial increase in model size
219 – also called the *expressive efficiency* (Martens & Medabalimi, 2014), or *succinctness* (de Colnet
220 & Mengel, 2021) of a model class. This is because the size of circuits directly translates to the
221 computational complexity of performing inference. As we extend the language of monotonic PCs
222 to include negative parameters, here we provide polytime reductions from tractable probabilistic
223 model classes emerging from different application fields that can encode subtractions, to (deep)
224 non-monotonic PCs. By doing so, we not only shed light on why they are tractable, by explicitly
225 stating their structural properties as circuits, but also on why they can be more expressive than
226 classical additive MMs, as we prove that NPC²s can be exponentially more compact in Sec. 4.1.

227 **Simple shallow NMMs** have been investigated for a limited set of component families, as discussed
228 in Sec. 1. Notably, this can also be done by directly learning to approximate the square root of a
229 density function, as done in signal processing with wavelet functions as components (Daubechies,
230 1992; Pinheiro & Vidakovic, 1997) or RBF kernels, i.e., unnormalized Gaussians centered over data
231 points (Schölkopf & Smola, 2001), as in Hong & Gao (2021). As discussed in Sec. 3, we can readily
232 represent these NMMs as simple NPC²s where kernel functions are computed by input layers.

233 **Positive semi-definite (PSD) models** (Rudi & Ciliberto, 2021; Marteau-Ferey et al., 2020) are
234 a recent class of models from the kernel and optimization literature. Given a kernel function κ
235 (e.g., an RBF kernel as in Rudi & Ciliberto (2021)) and a set of d data points $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(d)}$ with
236 $\boldsymbol{\kappa}(\mathbf{x}) = [\kappa(\mathbf{x}, \mathbf{x}^{(1)}), \dots, \kappa(\mathbf{x}, \mathbf{x}^{(d)})]^\top \in \mathbb{R}^d$, and a real $d \times d$ PSD matrix \mathbf{A} , they define an unnor-
237 malized distribution as the non-negative function $f(\mathbf{x}; \mathbf{A}, \boldsymbol{\kappa}) = \boldsymbol{\kappa}(\mathbf{x})^\top \mathbf{A} \boldsymbol{\kappa}(\mathbf{x})$. Although apparently
238 different, they can be translated to NPC²s in polynomial time.

239 **Proposition 2** (Reduction from PSD models). *A PSD model with kernel function κ , defined over*
240 *d data points, and parameterized by a PSD matrix \mathbf{A} , can be represented as a mixture of squared*
241 *NMMs (hence NPC²s) in time $\mathcal{O}(d^3)$.*

242 We prove this in App. B.3. Note that while PSD models are *shallow* non-monotonic PCs, we can
243 stack them into deeper NPC²s that support tractable marginalization via structured-decomposability.

244 **Tensor networks and the Born rule.** Squaring a possibly negative function to retrieve an un-
245 normalized distribution is related to the Born rule in quantum mechanics (Dirac, 1930), used to
246 characterize the distribution of particles by squaring their wave function (Schollwoeck, 2010; Orús,
247 2013). These functions can be represented as a large D -dimensional tensor \mathcal{T} over discrete vari-
248 ables $\mathbf{X} = \{X_1, \dots, X_D\}$ taking value $\{1, \dots, m\}$, compactly factorized in a tensor network (TN)
249 such as a matrix-product state (MPS) (Pérez-García et al., 2007), also called tensor-train. Given an
250 assignment $\mathbf{x} = \langle x_1, \dots, x_D \rangle$ to \mathbf{X} , a rank r MPS compactly represents \mathcal{T} as

$$\mathcal{T}[x_1, \dots, x_D] = \sum_{i_1=1}^r \sum_{i_2=1}^r \cdots \sum_{i_{D-1}=1}^r \mathbf{A}_1[x_1, i_1] \mathbf{A}_2[x_2, i_1, i_2] \cdots \mathbf{A}_D[x_D, i_{D-1}], \quad (5)$$

251 where $\mathbf{A}_1, \mathbf{A}_D \in \mathbb{R}^{m \times r}$, $\mathbf{A}_j \in \mathbb{R}^{m \times r \times r}$ with $1 < j < D$, for indices $\{i_1, \dots, i_{D-1}\}$, and de-
252 noting indexing with square brackets. To encode a distribution $p(\mathbf{X})$, one can reparameterize ten-
253 sors \mathbf{A}_j to be non-negative (Novikov et al., 2021) or apply the Born rule and square \mathcal{T} to model

254 $p(\mathbf{x}) \propto (\mathcal{T}[x_1, \dots, x_D])^2$. Such a TN is called a Born machine (BM) (Glasser et al., 2019). Be-
 255 sides modeling complex quantum states, TNs such as BMs have also been explored as classical
 256 ML models to learn discrete distributions (Stoudenmire & Schwab, 2016; Han et al., 2018; Glasser
 257 et al., 2019; Cheng et al., 2019), in quantum ML (Liu & Wang, 2018; Huggins et al., 2018), and
 258 more recently extended to continuous domains by introducing sets of basis functions, called TTDE
 259 (Novikov et al., 2021). Next, we show they are a special case of NPC²s.

260 **Proposition 3** (Reduction from BMs). *A BM encoding D -dimensional tensor with m states by*
 261 *squaring a rank r MPS can be exactly represented as a structured-decomposable NPC² in $\mathcal{O}(D \cdot k^4)$*
 262 *time and space, with $k \leq \min\{r^2, mr\}$.*

263 We prove this in App. B.4 by showing an equivalent NPC² defined on linear tree RG (e.g., the one
 264 in Fig. 2a). This connection highlights how tractable marginalization in BMs is possible thanks to
 265 structured-decomposability (Proposition 1), a condition that to the best of our knowledge was not
 266 previously studied for TNs. Furthermore, as NPC²s we can now design more flexible tree RGs, e.g.,
 267 randomized tree structures (Peharz et al., 2020b; Di Mauro et al., 2017; Di Mauro et al., 2021),
 268 densely tensorized structures heavily exploiting GPU parallelization (Peharz et al., 2020a; Mari
 269 et al., 2023) or heuristically learn them from data (Liu & Van den Broeck, 2021).

270 4.1 EXPONENTIAL SEPARATION OF NPC²S AND STRUCTURED MONOTONIC PCs

271 Squaring via Alg. 1 can already make a tensorized (monotonic) PC more expressive, but only by a
 272 polynomial factor, as we quadratically increase the size of each layer, while keeping the same num-
 273 ber of learnable parameters (similarly to the increased number of components of squared NMMs
 274 (Sec. 2)). On the other hand, allowing negative parameters can provide an exponential advantage,
 275 as proven for certain circuits (Valiant, 1979), but understanding if this advantage carries over to
 276 our squared circuits is not immediate. In fact, we observe there cannot be any expressiveness ad-
 277 vantage in squaring certain classes of non-monotonic structured-decomposable circuits. These are
 278 the circuits that support tractable maximum-a-posteriori inference (Choi et al., 2020) and satisfy an
 279 additional property known as *determinism* (see Darwiche (2001), Def. A.5). Squaring these circuits
 280 outputs a PC of the same size and that is monotonic, as formalized next and proven in App. B.6.

281 **Proposition 4** (Squaring deterministic circuits). *Let c be a smooth, decomposable and deterministic*
 282 *circuit, possibly computing a negative function. Then, the squared circuit c^2 is monotonic and has*
 283 *the same structure (and hence size) of c .*

284 The NPC²s we considered so far, as constructed in Sec. 3, are *not* deterministic. Here we prove that
 285 some non-negative functions (hence probability distributions up to renormalization) can be com-
 286 puted by NPC²s that are exponentially smaller than any structured-decomposable monotonic PC.

287 **Theorem 1** (Expressive efficiency of NPC²s). *There is a class of non-negative functions \mathcal{F} over*
 288 *variables \mathbf{X} that can be compactly represented as a shallow squared NMM (hence NPC²s), but for*
 289 *which the smallest structured-decomposable monotonic PC computing any $F \in \mathcal{F}$ has size $2^{\Omega(|\mathbf{X}|)}$.*

290 We prove this in App. B.5 by showing a non-trivial lower bound on the size of structured-
 291 decomposable monotonic PCs for a variant of the unique disjointness problem (Fiorini et al., 2015).
 292 Intuitively, this tells us that, given a fixed number of parameters, NPC²s can potentially be much
 293 more expressive than structured-decomposable monotonic PCs (and hence shallow additive MMs).
 294 We conjecture that an analogous lower bound can be devised for decomposable monotonic PCs.
 295 Furthermore, as this result directly extends to PSD and BM models (Sec. 4), it opens up interesting
 296 theoretical connections in the research fields of kernel-based and TN models.

297 5 EXPERIMENTS

298 We aim to answer the following questions: (A) are NPC²s better distribution estimators than mono-
 299 tonic PCs? (B) how the increased model size given by squaring and the presence of negative pa-
 300 rameters independently influence the expressiveness of NPC²s? (C) how does the choice of input
 301 layers and the RG affect the performance of NPC²s? We perform several distribution estimation
 302 experiments on both synthetic and real-world data, and label the following paragraphs with letters
 303 denoting relevance to the above questions. Moreover, note that our comparisons between NPC²s
 304 and monotonic PCs are based on models having the same number of learnable parameters.

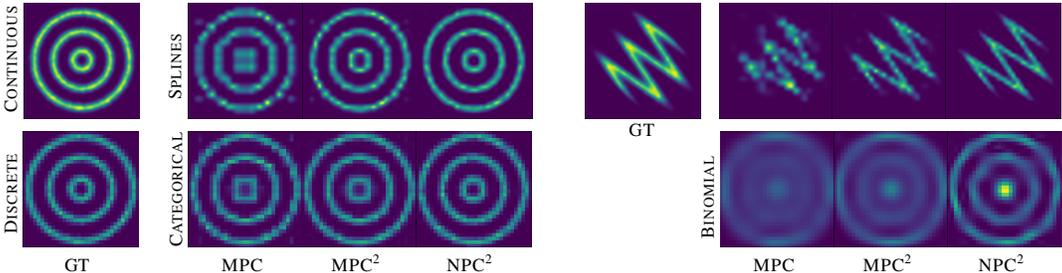


Figure 3: **NPC²s are better estimators, especially with parameter-efficient input layers.** Distribution estimated by monotonic PCs (MPC), squared monotonic PCs (MPC²) and NPC²s on 2D continuous (above) and discrete (below) data. On continuous data input layers compute splines (Eq. (11)), while on discrete data they compute either categoricals (for MPC and MPC²), embeddings (for NPC²s) or Binomials. Apps. H.1 and H.2 shows log-likelihoods on also additional data.

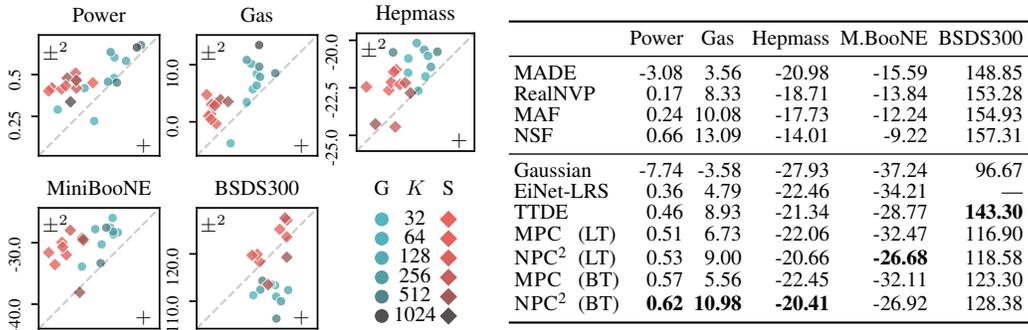


Figure 4: **NPC²s can be more expressive than monotonic PCs (MPCs).** Best average log-likelihoods achieved by monotonic PCs (+) and NPC²s (±²), built either from randomized linear tree (LT) or binary tree (BT) RGs (see App. H.3). The scatter plots (left) pairs log-likelihoods based on the number of units per layer K (the higher the darker), differentiating PCs with Gaussian (G:blue) and splines (S:red) input layers. Both axes of each scatter plot are on the same scale, thus the results above the diagonal are of NPC²s achieving higher log-likelihoods than MPCs at parity of model size. The table (right) shows our models’ best average test log-likelihoods and puts them in context with intractable (above) and tractable (below) models.

305 (A, B) **Synthetic continuous data.** Following Wenliang et al. (2019), we evaluate monotonic PCs
 306 and NPC²s on 2D density estimation tasks, as this allows us to gain an insight on the learned
 307 density functions. To disentangle the effect of squaring versus that of negative parameters, we also
 308 experiment with squared monotonic PCs. We build circuit structures from a trivial tree RG (see
 309 App. H.1 for details). We experiment with splines as input layers for NPC²s, and enforce their non-
 310 negativity for monotonic PCs (see App. E). Fig. 3 shows that, while squaring benefits monotonic
 311 PCs, negative parameters in NPC²s are needed to better capture complex target densities.

312 (C) **Synthetic discrete data.** We estimate the probability mass of the previous 2D data sets, now
 313 finitely-discretized (see App. H.2), to better understand when negative parameters might bring little
 314 to no advantage if input layers are already expressive enough. First, we experiment with (squared)
 315 monotonic PCs (resp. NPC²s) having input layers computing categoricals (resp. real-valued em-
 316 beddings). Second, we employ the less flexible but more parameter-efficient Binomials instead.
 317 App. H.2 reports the hyperparameters. Fig. 3 shows that, while there is no clear advantage for
 318 NPC²s equipped with embeddings instead of MPC² with categoricals, in case of Binomials they
 319 can better capture the target distribution. This is because categoricals (and embeddings) already
 320 have enough parameters to capture “holes” in the probability mass function. However, Binomials
 321 introduce a strong inductive bias that might hinder learning. We believe this is the reason why, ac-
 322 cording to some preliminary results, we did not observe an improvement of NPC²s with respect to
 323 monotonic PCs on estimating image distributions.

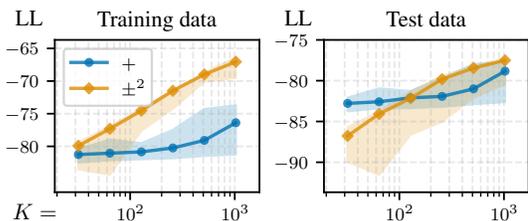


Figure 5: NPC²s (\pm^2) achieve higher log-likelihoods than monotonic PCs (+) on data sampled by GPT2. We report the median and the area including 80% of runs by varying the size of layers K and other hyperparameters (App. H.4). For comparison, the log-likelihood of GPT2 on the same training data is about -52 . Overfitting is observed for NPC²s, as it achieves much higher log-likelihood on the training data.

324 (A, B, C) **Multi-variate continuous data.** Following Papamakarios et al. (2017), we evaluate
 325 deeper PCs for density estimation on five multivariate data sets (statistics are reported in Table H.1).
 326 We evaluate monotonic PCs and NPC²s in tensorized form built out of randomized linear tree RGs.
 327 That is, for some variable permutation, we construct a tree RG where each partition splits a region
 328 into a set of only one variable and recursively factorizes the rest. By doing so, we recover architec-
 329 tures similar to a BMs or TTDEs (see Sec. 4). Following Peharz et al. (2020b), we also experiment
 330 with randomized binary tree RGs whose partitions randomly split regions in half. App. H.3 details
 331 these RGs, as well as the hyperparameters used. We compare against: a full covariance Gaussian,
 332 normalizing flows RealNVP (Dinh et al., 2017), MADE (Germain et al., 2015), MAF (Papamakari-
 333 os et al., 2017) and NSF (Durkan et al., 2019), a monotonic PC with input layers encoding flows
 334 (EiNet-LRS) (Sidheekh et al., 2023), and TTDE (Novikov et al., 2021). Fig. 4 shows that NPC²s
 335 with Gaussian input layers generally achieve higher log-likelihoods than monotonic PCs on four
 336 data sets. Fig. H.3 shows similar results when comparing to squared monotonic PCs, thus providing
 337 evidence that negative parameters other than squaring contribute to the expressiveness of NPC²s.
 338 Binary tree RGs generally deliver better likelihoods than linear tree ones, especially on Gas, where
 339 NPC²s using them outperform TTDE, which uses a sophisticated Riemanniann optimization scheme.

340 (A) **Distilling intractable models.** Monotonic PCs have been used to approximate intractable mod-
 341 els such as LLMs and perform exact inference in presence of logical constraints, such as for con-
 342 strained text generation (Zhang et al., 2023). As generation performance is correlated with how
 343 well the LLM is approximated by a tractable model, we are interested in how NPC²s can better be
 344 the distillation target of a LLM such as GPT2, rather than monotonic PCs. Following Zhang et al.
 345 (2023), we minimize the KL divergence between GPT2 and our PCs on a data set of sentences hav-
 346 ing bounded length (see App. H.4 for details). Since sentences are sequences of token variables, the
 347 architecture of tensorized circuits is built from a linear tree RG, thus corresponding to an inhomoge-
 348 neous HMM in case of monotonic PCs (see App. B.4.1) while resembling a BM for NPC²s. Fig. 5
 349 shows that NPC²s can scale and distill GPT2 more compactly than monotonic PCs, as they achieve
 350 log-likelihoods closer to the ones computed by GPT2. Moreover, we observe that NPC²s overfit the
 351 training data, which is however further evidence of the increased expressiveness of NPC²s. In the
 352 limit and by sampling enough sentences from GPT2, one can definitely reduce the risk of overfitting
 353 in the mentioned experimental setting. While regularization methods have been proposed for MPCs
 354 (Peharz et al., 2020b), regularizing NPC²s deserves future investigation.

355 6 DISCUSSION & CONCLUSION

356 With this work, we hope to popularize subtractive MMs via squaring as a simple and effective model
 357 class in the toolkit of tractable probabilistic modeling and reasoning that can rival traditional additive
 358 MMs. By casting them in the framework of circuits, we presented how to effectively represent
 359 and learn deep subtractive MMs such as NPC²s (Sec. 3) while showing how they can generalize
 360 other model classes such as PSD and tensor network models (Sec. 4). Our main theoretical result
 361 (Sec. 4.1) applies also to these models and justifies the increased performance we found in practice
 362 (Sec. 5). This work is the first to rigorously address representing and learning non-monotonic PCs
 363 in a general way, and opens up a number of future research directions. The first one is to retrieve a
 364 latent variable interpretation for NPC²s, as negative parameters in a non-monotonic PC invalidate
 365 the probabilistic interpretation of its sub-circuits (Peharz et al., 2017), making it not possible to learn
 366 its structure and parameters in classical ways (see App. G). Better ways to learn NPC²s, in turn, can
 367 benefit all applications in which PCs are widely used – from causal discovery (Wang et al., 2022)

368 to variational inference (Shih & Ermon, 2020) and neuro-symbolic AI (Ahmed et al., 2022) – by
 369 making more compact and expressive distributions accessible. Finally, by connecting circuits with
 370 tensor networks for the first time, we hope to inspire works that carry over the advancements of one
 371 community to the other, such as better learning schemes (Stoudenmire & Schwab, 2016; Novikov
 372 et al., 2021), and more flexible ways to factorize high-dimensional tensors (Mari et al., 2023).

373 REPRODUCIBILITY STATEMENT

374 In **App. H** we include all the details about the experiments we showed in **Sec. 5**. The source code,
 375 documentation, data sets and scripts needed to reproduce the results and figures, are available at
 376 <https://github.com/april-tools/squared-npcs>.

377 ACKNOWLEDGMENTS

378 AV was supported by the "UNREAL: Unified Reasoning Layer for Trustworthy ML" project
 379 (EP/Y023838/1) selected by the ERC and funded by UKRI EPSRC. NG acknowledges the support
 380 by the European Union (ERC consolidator, eLinoR, no 101085607). AMS acknowledges funding
 381 from the Helsinki Institute for Information Technology. MT acknowledges funding from the Re-
 382 search Council of Finland (grant number 347279). The authors acknowledge the computational
 383 resources provided by the CSC – IT Center for Science, Finland.

384 REFERENCES

- 385 Kareem Ahmed, Stefano Teso, Kai-Wei Chang, Guy Van den Broeck, and Antonio Vergari. Seman-
 386 tic probabilistic layers for neuro-symbolic learning. In *Advances in Neural Information Process-
 387 ing Systems 35 (NeurIPS)*, volume 35, pp. 29944–29959. Curran Associates, Inc., 2022.
- 388 Pierre Baldi, Kyle Cranmer, Taylor Faucett, Peter Sadowski, and Daniel Whiteson. Parameterized
 389 neural networks for high-energy physics. *European Physical Journal C*, 76(5):235, 2016. doi:
 390 10.1140/epjc/s10052-016-4099-4.
- 391 Pierre Blanchard, Desmond J. Higham, and Nicholas J. Higham. Accurately computing the log-sum-
 392 exp and softmax functions. *Institute of Mathematics and its Applications Journal of Numerical
 393 Analysis (IMAJNA)*, 41(4):2311–2330, 2021.
- 394 Song Cheng, Lei Wang, Tao Xiang, and Pan Zhang. Tree tensor networks for generative modeling.
 395 *Physical Review B*, 99(15):155131, 2019.
- 396 Arthur Choi, Guy Van den Broeck, Adnan Darwiche, Qiang Yang, and Michael Wooldridge.
 397 Tractable learning for structured probability spaces: A case study in learning preference distribu-
 398 tions. In *24th International Joint Conference on Artificial Intelligence (IJCAI)*, volume 2015, pp.
 399 2861–2868. IJCAI, 2015.
- 400 YooJung Choi, Antonio Vergari, and Guy Van den Broeck. Probabilistic circuits: A unifying frame-
 401 work for tractable probabilistic modeling. Technical report, University of California, Los Angeles
 402 (UCLA), 2020.
- 403 Alvaro H. C. Correia and Cassio P. de Campos. Towards scalable and robust sum-product networks.
 404 In *Scalable Uncertainty Management*, 2019.
- 405 Meihua Dang, Antonio Vergari, and Guy Van den Broeck. Strudel: A fast and accurate learner of
 406 structured-decomposable probabilistic circuits. *The International Journal of Approximate Rea-
 407 soning (IJAR)*, 140:92–115, 2021.
- 408 Adnan Darwiche. Decomposable negation normal form. *Journal of the ACM (JACM)*, 48:608–647,
 409 2001.
- 410 Adnan Darwiche. *Modeling and Reasoning with Bayesian Networks*. Cambridge University Press,
 411 2009.
- 412 Adnan Darwiche and Pierre Marquis. A knowledge compilation map. *Journal of Artificial Intelli-
 413 gence Research (JAIR)*, 17:229–264, 2002.

- 414 Ingrid Daubechies. Ten lectures on wavelets. *Computers in Physics*, 6:697–697, 1992.
- 415 Carl de Boor. Subroutine package for calculating with B-splines. Technical report, Los Alamos
416 National Lab. (LANL), 1971.
- 417 Alexis de Colnet and Stefan Mengel. A compilation of succinctness results for arithmetic circuits. In
418 *18th International Conference on Principles of Knowledge Representation and Reasoning (KR)*,
419 pp. 205–215, 2021.
- 420 Ronald De Wolf. Nondeterministic quantum query and communication complexities. *SIAM Journal*
421 *on Computing*, 32(3):681–699, 2003.
- 422 Aaron W. Dennis. *Algorithms for Learning the Structure of Monotone and Nonmonotone Sum-*
423 *Product Networks*. PhD thesis, Brigham Young University, 2016.
- 424 Aaron W. Dennis and Dan Ventura. Learning the architecture of sum-product networks using clus-
425 tering on variables. In *Advances in Neural Information Processing Systems 25 (NeurIPS)*, pp.
426 2033–2041. Curran Associates, Inc., 2012.
- 427 Nicola Di Mauro, Antonio Vergari, Teresa M. A. Basile, and Floriana Esposito. Fast and accu-
428 rate density estimation with extremely randomized cutset networks. In *Machine Learning and*
429 *Knowledge Discovery in Databases: ECML PKDD*, pp. 203–219. Springer, 2017.
- 430 Nicola Di Mauro, Gennaro Gala, Marco Iannotta, and Teresa Maria Altomare Basile. Random
431 probabilistic circuits. In *37th Conference on Uncertainty in Artificial Intelligence (UAI)*, volume
432 161, pp. 1682–1691. PMLR, 2021.
- 433 Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using Real NVP. In
434 *5th International Conference on Learning Representations (ICLR)*, 2017.
- 435 Paul Adrien Maurice Dirac. *The Principles of Quantum Mechanics*. Clarendon Press, Oxford,, 1930.
- 436 Dheeru Dua and Casey Graff. UCI Machine Learning Repository, 2017.
- 437 Conor Durkan, Artur Bekasov, Iain Murray, and George Papamakarios. Neural spline flows. In
438 *Advances in Neural Information Processing Systems 32 (NeurIPS)*, pp. 7511–7522. Curran Asso-
439 ciates, Inc., 2019.
- 440 Samuel Fiorini, Serge Massar, Sebastian Pokutta, Hans Raj Tiwary, and Ronald De Wolf. Exponen-
441 tial lower bounds for polytopes in combinatorial optimization. *Journal of the ACM (JACM)*, 62
442 (2):1–23, 2015.
- 443 Jordi Fonollosa, Sadique Sheik, Ramón Huerta, and Santiago Marco. Reservoir computing compen-
444 sates slow response of chemosensor arrays exposed to fast varying gas concentrations in contin-
445 uous monitoring. *Sensors and Actuators B: Chemical*, 215:618–629, 2015.
- 446 Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. MADE: Masked autoencoder
447 for distribution estimation. In *32nd International Conference on Machine Learning (ICML)*, pp.
448 881 – 889, 2015.
- 449 Nicolas Gillis. *Nonnegative Matrix Factorization*. Society for Industrial and Applied Mathematics
450 (SIAM), 2020.
- 451 Ivan Glasser, Ryan Sweke, Nicola Pancotti, Jens Eisert, and Ignacio Cirac. Expressive power of
452 tensor-network factorizations for probabilistic modeling. In *Advances in Neural Information Pro-*
453 *cessing Systems 32 (NeurIPS)*, pp. 1498–1510. Curran Associates, Inc., 2019.
- 454 Zhao-Yu Han, Jun Wang, Heng Fan, Lei Wang, and Pan Zhang. Unsupervised generative modeling
455 using matrix product states. *Physical Review X*, 8:031012, Jul 2018.
- 456 Georges Hebrail and Alice Berard. Individual household electric power consumption. UCI Machine
457 Learning Repository, 2012.
- 458 Geoffrey E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural*
459 *Computation*, 14:1771–1800, 2002.

- 460 Xia Hong and Junbin Gao. Estimating the square root of probability density function on Riemannian
461 manifold. *Expert Systems - The Journal of Knowledge Engineering*, 38(7), 2021.
- 462 Shlomo Hoory, Nathan Linial, and Avi Wigderson. Expander graphs and their applications. *Bulletin*
463 *of the American Mathematical Society*, 43(4):439–561, 2006.
- 464 William J. Huggins, Piyush S. Patil, Bradley K. Mitchell, K. Birgitta Whaley, and Edwin Miles
465 Stoudenmire. Towards quantum machine learning with tensor networks. *Quantum Science and*
466 *Technology*, 4, 2018.
- 467 Renyan Jiang, Ming J. Zuo, and Han-Xiong Li. Weibull and inverse weibull mixture models allowing
468 negative weights. *Reliability Engineering & System Safety*, 66(3):227–234, 1999.
- 469 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd Interna-*
470 *tional Conference on Learning Representations (ICLR)*, 2015.
- 471 Tamara G. Kolda and Brett W. Bader. Tensor decompositions and applications. *Society of Industrial*
472 *and Applied Mathematics (SIAM) Review*, 51(3):455–500, 2009.
- 473 Yann LeCun, Sumit Chopra, Raia Hadsell, Marc’Aurelio Ranzato, and Fugie Huang. A tutorial on
474 energy-based learning. *Predicting Structured Data*, 2006.
- 475 Anji Liu and Guy Van den Broeck. Tractable regularization of probabilistic circuits. In *Advances in*
476 *Neural Information Processing Systems 34 (NeurIPS)*, pp. 3558–3570. Curran Associates, Inc.,
477 2021.
- 478 Anji Liu, Honghua Zhang, and Guy Van den Broeck. Scaling up probabilistic circuits by latent
479 variable distillation. In *11th International Conference on Learning Representations (ICLR)*, 2023.
- 480 Jin-Guo Liu and Lei Wang. Differentiable learning of quantum circuit born machines. *Physical*
481 *Review A*, 98(6):062324, 2018.
- 482 Lorenzo Loconte, Nicola Di Mauro, Robert Peharz, and Antonio Vergari. How to turn your knowl-
483 edge graph embeddings into generative models via probabilistic circuits. In *Advances in Neural*
484 *Information Processing Systems 37 (NeurIPS)*. Curran Associates, Inc., 2023.
- 485 Antonio Mari, Gennaro Vessio, and Antonio Vergari. Unifying and understanding overparameter-
486 ized circuit representations via low-rank tensor decompositions. In *6th Workshop on Tractable*
487 *Probabilistic Modeling*, 2023.
- 488 Ulysse Marteau-Ferey, Francis Bach, and Alessandro Rudi. Non-parametric models for non-negative
489 functions. In *Advances in Neural Information Processing Systems 33 (NeurIPS)*, pp. 12816–
490 12826, 2020.
- 491 James Martens and Venkatesh Medabalimi. On the expressive efficiency of sum product networks.
492 *arXiv preprint arXiv:1411.7717*, 2014.
- 493 David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented
494 natural images and its application to evaluating segmentation algorithms and measuring ecological
495 statistics. In *8th International Conference on Computer Vision (ICCV)*, volume 2, pp. 416–423.
496 IEEE, 2001.
- 497 D. Mauá, Diarmaid Conaty, Fabio Gagliardi Cozman, Katja Poppenhaeger, and Cassio Polpo
498 de Campos. Robustifying sum-product networks. *International Journal of Approximate Rea-*
499 *soning*, 101:163–180, 2018.
- 500 Geoffrey J. McLachlan, Sharon X. Lee, and Suren I. Rathnayake. Finite mixture models. *Annual*
501 *Review of Statistics and its Application*, 6:355–378, 2019.
- 502 TrungTin Nguyen, Hien Duy Nguyen, Faicel Chamroukhi, and Geoffrey J. McLachlan. Approxima-
503 tion by finite mixtures of continuous density functions that vanish at infinity. *Cogent Mathematics*
504 *& Statistics*, 7, 2019.

- 505 Georgii S. Novikov, Maxim E. Panov, and Ivan V. Oseledets. Tensor-train density estimation. In
506 *37th Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 161 of *Proceedings of*
507 *Machine Learning Research*, pp. 1321–1331. PMLR, 2021.
- 508 Román Orús. A practical introduction to tensor networks: Matrix product states and projected
509 entangled pair states. *Annals of Physics*, 349:117–158, 2013.
- 510 George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density
511 estimation. In *Advances in Neural Information Processing Systems 30 (NeurIPS)*, pp. 2338–2347.
512 Curran Associates, Inc., 2017.
- 513 George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lak-
514 shminarayanan. Normalizing flows for probabilistic modeling and inference. *The Journal of*
515 *Machine Learning Research (JMLR)*, 22(1):2617–2680, 2021.
- 516 Robert Peharz, Robert Gens, Franz Pernkopf, and Pedro M. Domingos. On the latent variable
517 interpretation in sum-product networks. *IEEE Transactions on Pattern Analysis and Machine*
518 *Intelligence*, 39(10):2030–2044, 2017.
- 519 Robert Peharz, Steven Lang, Antonio Vergari, Karl Stelzner, Alejandro Molina, Martin Trapp, Guy
520 Van Den Broeck, Kristian Kersting, and Zoubin Ghahramani. Einsum networks: Fast and scalable
521 learning of tractable probabilistic circuits. In *37th International Conference on Machine Learning*
522 *(ICML)*, volume 119 of *Proceedings of Machine Learning Research*, pp. 7563–7574. PMLR,
523 2020a.
- 524 Robert Peharz, Antonio Vergari, Karl Stelzner, Alejandro Molina, Xiaoting Shao, Martin Trapp,
525 Kristian Kersting, and Zoubin Ghahramani. Random sum-product networks: A simple and ef-
526 fective approach to probabilistic deep learning. In *35th Conference on Uncertainty in Artificial*
527 *Intelligence (UAI)*, volume 115 of *Proceedings of Machine Learning Research*, pp. 334–344.
528 PMLR, 2020b.
- 529 David Pérez-García, F. Verstraete, Michael M. Wolf, and Juan Ignacio Cirac. Matrix product state
530 representations. *Quantum Information and Computing*, 7(5):401–430, 2007. ISSN 1533-7146.
- 531 Les A. Piegl and Wayne Tiller. The NURBS book. In *Monographs in Visual Communication*, 1995.
- 532 Aluisio Pinheiro and Brani Vidakovic. Estimating the square root of a density via compactly sup-
533 ported wavelets. *Computational Statistics and Data Analysis*, 25(4):399–415, 1997.
- 534 Knot Pipatsrisawat and Adnan Darwiche. New compilation languages based on structured decom-
535 posability. In *23rd Conference on Artificial Intelligence (AAAI)*, volume 8, pp. 517–522, 2008.
- 536 Hoifung Poon and Pedro Domingos. Sum-product networks: A new deep architecture. In *IEEE*
537 *International Conference on Computer Vision Workshops (ICCV Workshops)*, pp. 689–690. IEEE,
538 2011.
- 539 Guillaume Rabusseau and François Denis. Learning negative mixture models by tensor decomposi-
540 tions. *arXiv preprint arXiv:1403.4224*, 2014.
- 541 Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*.
542 Adaptive Computation and Machine Learning. MIT Press, 2005.
- 543 Byron P. Roe, Hai-Jun Yang, Ji Zhu, Yong Liu, Ion Stancu, and Gordon McGregor. Boosted decision
544 trees as an alternative to artificial neural networks for particle identification. *Nuclear Instruments*
545 *& Methods in Physics Research Section A-accelerators Spectrometers Detectors and Associated*
546 *Equipment*, 543:577–584, 2004.
- 547 Tim Roughgarden. Communication complexity (for algorithm designers). *Foundations and Trends®*
548 *in Theoretical Computer Science*, 11(3–4):217–404, 2016.
- 549 Alessandro Rudi and Carlo Ciliberto. PSD representations for effective probability models. In
550 *Advances in Neural Information Processing Systems 34 (NeurIPS)*, pp. 19411–19422. Curran
551 Associates, Inc., 2021.

- 552 Bernhard Schölkopf and Alex Smola. Learning with kernels: support vector machines, regulariza-
553 tion, optimization, and beyond. In *Adaptive Computation and Machine Learning Series*. MIT
554 Press, 2001.
- 555 Ulrich Schollwoeck. The density-matrix renormalization group in the age of matrix product states.
556 *Annals of Physics*, 326:96–192, 2010.
- 557 Matthias Seeger. Expectation propagation for exponential families. Technical report, Department of
558 EECS, University of California at Berkeley, 2005.
- 559 Yujia Shen, Arthur Choi, and Adnan Darwiche. Tractable operations for arithmetic circuits of prob-
560 abilistic models. In *Advances in Neural Information Processing Systems 29 (NeurIPS)*. Curran
561 Associates, Inc., 2016.
- 562 Andy Shih and Stefano Ermon. Probabilistic circuits for variational inference in discrete graphical
563 models. In *Advances in Neural Information Processing Systems 33 (NeurIPS)*. Curran Associates,
564 Inc., 2020.
- 565 Amir Shpilka and Amir Yehudayoff. Arithmetic circuits: A survey of recent results and open ques-
566 tions. *Foundations and Trends in Theoretical Computer Science*, 5:207–388, 2010.
- 567 Sahil Sidheekh, Kristian Kersting, and Sriraam Natarajan. Probabilistic flow circuits: Towards
568 unified deep models for tractable probabilistic inference. In *39th Conference on Uncertainty
569 in Artificial Intelligence (UAI)*, volume 216 of *Proceedings of Machine Learning Research*, pp.
570 1964–1973. PMLR, 2023.
- 571 Edwin Stoudenmire and David J Schwab. Supervised learning with tensor networks. In *Advances in
572 Neural Information Processing Systems 29 (NeurIPS)*, pp. 4799–4807. Curran Associates, Inc.,
573 2016.
- 574 Russell Tsuchida, Cheng Soon Ong, and Dino Sejdinovic. Squared neural families: A new class of
575 tractable density models. *arXiv preprint arXiv:2305.13552*, 2023.
- 576 Leslie G. Valiant. Negation can be exponentially powerful. In *11th Annual ACM Symposium on
577 Theory of Computing*, pp. 189–196, 1979.
- 578 Antonio Vergari, Nicola Di Mauro, and Floriana Esposito. Visualizing and understanding sum-
579 product networks. *Machine Learning*, 108(4):551–573, 2019a.
- 580 Antonio Vergari, Nicola Di Mauro, and Guy Van den Broeck. Tractable probabilistic models: Repre-
581 sentations, algorithms, learning, and applications. *Tutorial at the 35th Conference on Uncertainty
582 in Artificial Intelligence (UAI)*, 2019b.
- 583 Antonio Vergari, YooJung Choi, Anji Liu, Stefano Teso, and Guy Van den Broeck. A compositional
584 atlas of tractable circuit operations for probabilistic inference. In *Advances in Neural Information
585 Processing Systems 34 (NeurIPS)*, pp. 13189–13201. Curran Associates, Inc., 2021.
- 586 Allan H. Vermeulen, Richard H. Bartels, and Glenn R. Heppler. Integrating products of B-splines.
587 *SIAM Journal on Scientific and Statistical Computing*, 13:1025–1038, 1992.
- 588 Benjie Wang, Matthew R. Wicker, and Marta Kwiatkowska. Tractable uncertainty for structure
589 learning. In *39th International Conference on Machine Learning (ICML)*, pp. 23131–23150.
590 PMLR, 2022.
- 591 Li Wenliang, Danica J. Sutherland, Heiko Strathmann, and Arthur Gretton. Learning deep kernels
592 for exponential family densities. In *36th International Conference on Machine Learning (ICML)*,
593 volume 97 of *Proceedings of Machine Learning Research*, pp. 6737–6746. PMLR, 2019.
- 594 Baibo Zhang and Changshui Zhang. Finite mixture models with negative components. In *4th In-
595 ternational Conference on Machine Learning and Data Mining in Pattern Recognition (MLDM)*,
596 pp. 31–41. Springer, 2005.
- 597 Honghua Zhang, Brendan Juba, and Guy Van den Broeck. Probabilistic generating circuits. In
598 *International Conference on Machine Learning*, pp. 12447–12457. PMLR, 2021.

- 599 Honghua Zhang, Meihua Dang, Nanyun Peng, and Guy Van den Broeck. Tractable control for au-
600 toregressive language generation. In *40th International Conference on Machine Learning (ICML)*,
601 volume 202 of *Proceedings of Machine Learning Research*, pp. 40932–40945. PMLR, 2023.



Figure A.1: **Computational units can be grouped into layers as to build a tensorized circuit.** Sum units each parameterized by the rows of $\mathbf{W} \in \mathbb{R}^{3 \times 3}$ (left, in purple) form a sum layer parameterized by \mathbf{W} (right). Product units (left, in red) form an Hadamard product layer (right). Input units (left, in yellow) form an input layer computing the same functions (right)

602 A CIRCUITS

603 In [Sec. 3](#) we introduced circuits in a tensorized form. Here we instead present the definitions and
 604 properties of circuits as they are usually defined in the literature, which will be used in [App. B](#).

605 **Definition A.1** (Circuit ([Choi et al., 2020](#); [Vergari et al., 2021](#))). A *circuit* c is a parameterized
 606 computational graph over variables \mathbf{X} encoding a function $c(\mathbf{X})$ and comprising three kinds of
 607 computational units: *input*, *product*, and *sum*. Each product or sum unit n receives as inputs the
 608 outputs of other units, denoted with the set $\text{in}(n)$. Each unit n encodes a function c_n defined as: (i)
 609 $f_n(\text{sc}(n))$ if n is an input unit, where f_n is a function over variables $\text{sc}(n) \subseteq \mathbf{X}$, called its *scope*,
 610 (ii) $\prod_{i \in \text{in}(n)} c_i(\text{sc}(n_i))$ if n is a product unit, and (iii) $\sum_{i \in \text{in}(n)} w_i c_i(\text{sc}(n_i))$ if n is a sum unit, with
 611 $w_i \in \mathbb{R}$ denoting the weighted sum parameters. The scope of a product or sum unit n is the union
 612 of the scopes of its inputs, i.e., $\text{sc}(n) = \bigcup_{i \in \text{in}(n)} \text{sc}(i)$.

613 Note that tensorized circuits ([Def. 1](#)) are circuits where each input (resp. product and sum) layer
 614 consists of scalar input (resp. product and sum) units. For example, [Fig. A.1](#) shows how compu-
 615 tational units are grouped into layers. A *probabilistic circuit* (PC) is defined as a circuit encoding
 616 a non-negative function. PCs that are *smooth* and *decomposable* ([Def. A.2](#)) enable computing the
 617 partition function and, more in general, performing variable marginalization efficiently ([Prop. A.1](#)).

618 **Definition A.2** (Smoothness and decomposability ([Darwiche & Marquis, 2002](#))). A circuit is *smooth*
 619 if for every sum unit n , its input units depend all on the same variables, i.e., $\forall i, j \in \text{in}(n): \text{sc}(i) =$
 620 $\text{sc}(j)$. A circuit is *decomposable* if the inputs of every product unit n depend on disjoint sets of
 621 variables, i.e., $\forall i, j \in \text{in}(n) i \neq j: \text{sc}(i) \cap \text{sc}(j) = \emptyset$.

622 **Proposition A.1** (Tractability ([Choi et al., 2020](#))). Let c be a smooth and decomposable circuit over
 623 variables \mathbf{X} whose input units can be integrated efficiently. Then for any $\mathbf{Z} \subseteq \mathbf{X}$ and \mathbf{y} an assignment
 624 to variables in $\mathbf{X} \setminus \mathbf{Z}$, the quantity $\int c(\mathbf{y}, \mathbf{z}) d\mathbf{z}$ can be computed exactly in time and space $\Theta(|c|)$,
 625 where $|c|$ denotes the size of the circuit, i.e., the number of connections in the computational graph.

626 The size of circuits in tensorized form is obtained by counting the number of connections between
 627 the scalar computational units (as detailed in [App. A.1.1](#)). Squaring circuits or their tensorized rep-
 628 resentation efficiently such that the resulting PC is smooth and decomposable ([Def. A.2](#)) requires the
 629 satisfaction of *structured-decomposability*, as showed in ([Pipatsrisawat & Darwiche, 2008](#); [Vergari](#)
 630 [et al., 2021](#)).

631 **Definition A.3** (Structured-decomposability ([Pipatsrisawat & Darwiche, 2008](#); [Darwiche, 2009](#))). A
 632 circuit is *structured-decomposable* if (1) it is smooth and decomposable, and (2) any pair of product
 633 units n, m having the same scope decompose their scope at their input units in the same way.

634 Note that shallow MMs are both decomposable and structured-decomposable. As anticipated in
 635 [Sec. 3](#), the expressiveness of squared non-monotonic PCs that are also *deterministic* is the same as
 636 as monotonic deterministic PCs, which are used for tractable maximum-a-posteriori (MAP) infer-
 637 ence. We prove it formally in [App. B.6](#) by leveraging the definition of determinism that we show in
 638 [Def. A.5](#). Before that, we introduce the definition of *support* of a computational unit.

639 **Definition A.4** (Support ([Choi et al., 2020](#))). The *support* of a computational unit n over variables
 640 \mathbf{X} is defined as the set of value assignments to variables in \mathbf{X} such that the output of n is non-zero,
 641 i.e., $\text{supp}(n) = \{\mathbf{x} \in \text{val}(\mathbf{X}) \mid c_n(\mathbf{x}) \neq 0\}$, where $\text{val}(\mathbf{X})$ denotes the domain of variables \mathbf{X} .

642 **Definition A.5** (Determinism (Darwiche, 2001)). A circuit c is *deterministic* if for any sum unit
643 $n \in c$ its inputs have disjoint *support* (Def. A.4), i.e., $\forall i, j \in \text{in}(n), i \neq j: \text{supp}(i) \cap \text{supp}(j) = \emptyset$.

644 A.1 TENSORIZED CIRCUITS

645 **Def. 1** can be further generalized by introducing Kronecker product layers, which are the building
646 blocks of other tensorized circuit architectures, such as randomized and tensorized sum-product
647 networks (RAT-SPNs) (Peharz et al., 2020b), einsum networks (EiNets) (Peharz et al., 2020a).

648 **Definition A.6** (Tensorized circuit). A *tensorized circuit* c is a parameterized computational graph
649 encoding a function $c(\mathbf{X})$ and comprising of three kinds of layers: *input*, *product* and *sum*. Each
650 layer comprises computational units defined over the same set of variables, also called its *scope*, and
651 every non-input layer receives input from one or more layers. The scope of each non-input layer is
652 the union of the scope of its inputs, and the scope of the output layer computing $c(\mathbf{X})$ is \mathbf{X} . Each
653 input layer ℓ has scope $\mathbf{Y} \subseteq \mathbf{X}$ and computes a collection of K functions $f_i(\mathbf{Y}) \in \mathbb{R}$, i.e., ℓ outputs
654 a K -dimensional vector. Each product layer ℓ computes either an Hadamard (or element-wise) or
655 Kronecker product over the N layers it receives as input, i.e., $\ell = \odot_{i=1}^N \ell_i$ or $\otimes_{i=1}^N \ell_i$, respectively.
656 A sum layer with S sum units and receiving input from a previous layer $\ell \in \mathbb{R}^K$, is parameterized
657 by $\mathbf{W} \in \mathbb{R}^{S \times K}$ and computes $\mathbf{W}\ell$.

658 A.1.1 SIZE OF TENSORIZED CIRCUITS

659 The time and space complexity of evaluating a circuit is linear in its size. The size $|c|$ of a circuit c
660 (Def. A.1) is obtained by counting the number of input connections of each scalar product or sum
661 unit. In other words, it is the number of edges in its computational graph.

662 If c is a tensorized circuit, then its size is obtained by counting the number of connections in its non-
663 tensorized form. Fig. A.1 shows part of a tensorized circuit and its non-tensorized form. For sum
664 layers, the number of scalar input connections is the size of its parameterization matrix, i.e., $S \cdot K$
665 if it is parameterized by $\mathbf{W} \in \mathbb{R}^{S \times K}$. If ℓ is an Hadamard product layer computing $\ell = \odot_{i=1}^N \ell_i$,
666 where each ℓ_i outputs a K -dimensional vector, then the number of its scalar input connections is
667 $N \cdot K$. In case of Kronecker product layers as in the more general Def. A.6, i.e., $\ell = \otimes_{i=1}^N \ell_i$ where
668 each ℓ_i outputs a K -dimensional vector, then the number of its scalar input connections is K^{N+1} .

669 A.2 TRACTABLE EXACT SAMPLING

670 Each sum unit in a monotonic PC can be interpreted as a finitely discrete latent variable that can as-
671 sume as many values as the number of input connections (Peharz et al., 2017). As such, a monotonic
672 PC can be seen as a hierarchical MM. This allows us to sample exactly from the modeled distribution
673 by (1) recursively sampling latent variables until input units are reached, and (2) sampling observed
674 variables from the distributions modeled by input units (Vergari et al., 2019a).

675 Such probabilistic interpretation of *inner* sum units for NPC²s is not possible, as they can output
676 negative values. However, since NPC²s are smooth and decomposable (Def. A.2), they support
677 efficient marginalization and hence conditioning (Proposition 1). This allows us to still sample
678 exactly from the modeled distribution via *inverse transform sampling*. That is, we choose a variable
679 ordering X_1, X_2, \dots, X_D and sample them in an autoregressive fashion, i.e., $x_1 \sim p(X_1), x_2 \sim$
680 $p(X_2 | x_1), \dots, x_D \sim p(X_D | x_1, \dots, x_{D-1})$, which is still linear in the number of variables.

681 B PROOFS

682 B.1 SQUARING TENSORIZED CIRCUITS

683 **Proposition B.1** (Correctness of Alg. 1). Let c be a tensorized structured-decomposable circuit
684 (Def. 1 and Def. A.3), then Alg. 1 recursively constructs the layers of the squared tensorized PC c^2
685 such that c^2 is also structured-decomposable.

Proof. The proof is by induction on the structure of c . Let ℓ be a sum layer having as input ℓ_i and
computing $\mathbf{W}\ell_i$, with $\mathbf{W} \in \mathbb{R}^{S \times K}$ and ℓ_i computing an output in \mathbb{R}^S . If ℓ is the last layer of c (i.e.,

the output layer), then $S = 1$ since c outputs a scalar, and the squared layer ℓ^2 must compute

$$\ell^2 = (\mathbf{W}\ell_i) \cdot (\mathbf{W}\ell_i) = (\mathbf{W} \otimes \mathbf{W})(\ell_i \otimes \ell_i) = (\mathbf{W} \otimes \mathbf{W})\ell_i^2,$$

which requires squaring the input layer ℓ_i . By inductive hypothesis the squared circuit having ℓ_i^2 as output layer is structured-decomposable, hence also the squared circuit having ℓ^2 as output layer must be. If ℓ is a non-output sum layer, we still require computing the Kronecker product of its input layer. The squared layer ℓ^2 is again a sum layer that outputs a S^2 -dimensional vector, i.e.,

$$\ell^2 = \ell \otimes \ell = (\mathbf{W}\ell_i) \otimes (\mathbf{W}\ell_i) = (\mathbf{W} \otimes \mathbf{W})(\ell_i \otimes \ell_i) = (\mathbf{W} \otimes \mathbf{W})\ell_i^2$$

via mixed-product property (L11-15 in [Alg. 1](#)). Let ℓ be a binary³ Hadamard product layer computing $\ell_i \odot \ell_{ii}$ for input layers ℓ_i, ℓ_{ii} each computing a K -dimensional vector. Then, the squared layer ℓ^2 computes the Hadamard product between K^2 -dimensional vectors, i.e.,

$$\ell^2 = (\ell_i \odot \ell_{ii}) \otimes (\ell_i \odot \ell_{ii}) = (\ell_i \otimes \ell_i) \odot (\ell_{ii} \otimes \ell_{ii}) = \ell_i^2 \odot \ell_{ii}^2$$

686 via mixed-product property with respect to the Hadamard product. By inductive hypothesis ℓ_i^2 and
687 ℓ_{ii}^2 are the output layers of structured-decomposable circuits depending on a disjoint sets of variables.
688 As such, the circuit having ℓ^2 as output layer maintains structured-decomposability (L6-9 in [Alg. 1](#)).
689 For the base case we consider the squaring of an input layer ℓ that computes K functions f_i over
690 some variables $\mathbf{Y} \subseteq \mathbf{X}$. We replace ℓ with its squaring ℓ^2 which encodes the products $f_i(\mathbf{Y})f_j(\mathbf{Y})$,
691 $1 \leq i, j \leq K$, by introducing K^2 functions g_{ij} such that $g_{ij}(\mathbf{Y}) = f_i(\mathbf{Y})f_j(\mathbf{Y})$ (L2-4 in [Alg. 1](#)).

692 **Squaring Kronecker product layers.** In the case of ℓ being a binary Kronecker product layer
693 instead as in the more general [Def. A.6](#), then the squared layer ℓ^2 computes the Kronecker product
694 between K^2 -dimensional vectors up to a permutation of the entries, i.e.,

$$\ell^2 = (\ell_i \otimes \ell_{ii}) \otimes (\ell_i \otimes \ell_{ii}) = \mathbf{R}((\ell_i \otimes \ell_i) \otimes (\ell_{ii} \otimes \ell_{ii})) = \mathbf{R}(\ell_i^2 \otimes \ell_{ii}^2), \quad (6)$$

695 by introducing a $K^4 \times K^4$ permutation matrix \mathbf{R} whose rows are all zeros except for one entry
696 set to 1, which reorders the entries of $\ell_i^2 \otimes \ell_{ii}^2$ as to recover the equality in [Eq. \(6\)](#). Note that such
697 permutation maintains decomposability ([Def. A.2](#)), and its application can be computed by a sum
698 layer having \mathbf{R} as fixed parameters. Moreover, by inductive hypothesis, the squaring circuit having
699 ℓ^2 as output layer is still structured-decomposable. Finally, [Alg. B.2](#) generalizes [Alg. 1](#) as to support
700 the squaring of Kronecker product layers as showed above (L10-11 in [Alg. B.2](#)). \square

Algorithm B.2 squareTensorizedCircuit(ℓ, \mathcal{R})

Input: A tensorized circuit ([Def. A.6](#)) having output layer ℓ and defined on a tree RG rooted by \mathcal{R} .

Output: The tensorized squared circuit defined on the same tree RG having ℓ^2 as output layer computing $\ell \otimes \ell$.

1: if ℓ is an input layer then 2: ℓ computes K functions $f_i(\mathcal{R})$ 3: return An input layer ℓ^2 computing all K^2 4: product combinations $f_i(\mathcal{R})f_j(\mathcal{R})$ 5: else if ℓ is a product layer then 6: $\{(\ell_i, \mathcal{R}_i), (\ell_{ii}, \mathcal{R}_{ii})\} \leftarrow \text{getInputs}(\ell, \mathcal{R})$ 7: $\ell_i^2 \leftarrow \text{squareTensorizedCircuit}(\ell_i, \mathcal{R}_i)$ 8: $\ell_{ii}^2 \leftarrow \text{squareTensorizedCircuit}(\ell_{ii}, \mathcal{R}_{ii})$ 9: if $\ell = \ell_i \odot \ell_{ii}$ then return $\ell_i^2 \odot \ell_{ii}^2$	10: else return $\mathbf{R}(\ell_i^2 \otimes \ell_{ii}^2)$, where \mathbf{R} is 11: a permutation matrix (see proof of Prop. B.1) 12: else $\triangleright \ell$ is a sum layer 13: $\{(\ell, \mathcal{R})\} \leftarrow \text{getInputs}(\ell, \mathcal{R})$ 14: $\ell_i^2 \leftarrow \text{squareTensorizedCircuit}(\ell_i, \mathcal{R})$ 15: $\mathbf{W} \in \mathbb{R}^{S \times K} \leftarrow \text{getParameters}(\ell)$ 16: $\mathbf{W}' \in \mathbb{R}^{S^2 \times K^2} \leftarrow \mathbf{W} \otimes \mathbf{W}$ 17: return $\mathbf{W}'\ell_i^2$
--	---

701 B.2 TRACTABLE MARGINALIZATION OF NPC²S

702 **Proposition 1.** Let c be a tensorized structured-decomposable circuit where the products of func-
703 tions computed by each input layer can be tractably integrated. Any marginalization of c^2 obtained
704 via [Alg. 1](#) requires time and space $\mathcal{O}(L \cdot M^2)$.

705 *Proof.* Given c by hypothesis, [Prop. B.1](#) ensures that the PC built via [Alg. 1](#) computes c^2 and is
706 defined on the same tree RG ([Def. 2](#)) of c . As such, c^2 is structured-decomposable and hence also

³Without loss of generality, we assume product layers have exactly two layers as inputs.

smooth and decomposable (see [Def. A.3](#)). Now, we make an argument about c and c^2 in their non-tensorized form ([Def. A.1](#)) as to leverage [Prop. A.1](#) for tractable marginalization later. The size of c is $|c| \in \mathcal{O}(L \cdot M)$, where L is the number of layers and M the maximum number of scalar input connections of each layer in c (see [App. A.1.1](#) for details). The size of c^2 is therefore $|c^2| \in \mathcal{O}(L \cdot M^2)$, since [Alg. 1](#) squares the output dimension of each layer as well as the size of the parameterization matrix of each sum layer. Since c^2 is smooth and decomposable and the functions computed by its input layers can be tractably integrated, then [Prop. A.1](#) ensures we can marginalize any subset of variables in time and space $|c^2| \in \mathcal{O}(L \cdot M^2)$. \square

715 B.3 REPRESENTING PSD MODELS WITHIN THE LANGUAGE OF NPC²S

716 **Proposition 2.** A PSD model with kernel function κ , defined over d data points, and parameterized
717 by a PSD matrix \mathbf{A} , can be represented as a mixture of squared NMMs (hence NPC²s) in time
718 $\mathcal{O}(d^3)$.

Proof. The PSD model computes a non-negative function $f(\mathbf{x}; \mathbf{A}, \kappa) = \kappa(\mathbf{x})^\top \mathbf{A} \kappa(\mathbf{x})$, where $\kappa(\mathbf{x}) = [\kappa(\mathbf{x}, \mathbf{x}^{(1)}), \dots, \kappa(\mathbf{x}, \mathbf{x}^{(d)})] \in \mathbb{R}^d$, with data points $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(d)}$, and $\mathbf{A} \in \mathbb{R}^{d \times d}$ is PSD. Let $\mathbf{A} = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$ be the eigendecomposition of \mathbf{A} with rank r . Then we can rewrite $f(\mathbf{x}; \mathbf{A}, \kappa)$ as

$$f(\mathbf{x}; \mathbf{A}, \kappa) = \kappa(\mathbf{x})^\top \left(\sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^\top \right) \kappa(\mathbf{x}) = \sum_{i=1}^r \lambda_i (\mathbf{u}_i^\top \kappa(\mathbf{x}))^2,$$

719 where $\lambda_i > 0$ since \mathbf{A} is PSD. Therefore, such PSD model can be represented as a monotonic
720 mixture of $r \leq d$ squared NMMs ([Eq. \(2\)](#)), whose d components computing $\kappa(\mathbf{x})$ are shared. The
721 eigendecomposition of \mathbf{A} can be done in time $\mathcal{O}(d^3)$, and materializing each squared NMMs (e.g.,
722 as in [Fig. 1](#)) requires space $\mathcal{O}(d^2)$. Furthermore, note that if $\mathbf{A} = \mathbf{u} \mathbf{u}^\top$ is a rank-1 matrix, then
723 $f(\mathbf{x}; \mathbf{A}, \kappa) = (\mathbf{u}^\top \kappa(\mathbf{x}))^2$ is exactly a squared NMM whose d components compute $\kappa(\mathbf{x})$. \square

724 B.4 RELATIONSHIP WITH TENSOR NETWORKS

725 In this section, we detail the construction of a tensorized structured-decomposable circuit ([Def. 1](#))
726 that is equivalent to a matrix product state (MPS) tensor network ([Pérez-García et al., 2007](#)), as we
727 mention in [Sec. 4](#). As such, the application of the Born rule as to retrieve a probabilistic model called
728 Born machine (BM) ([Glasser et al., 2019](#)) is equivalent to squaring the equivalent circuit ([Sec. 3](#)).

729 **Proposition 3.** A BM encoding D -dimensional tensor with m states by squaring a rank r MPS
730 can be exactly represented as a structured-decomposable NPC² in $\mathcal{O}(D \cdot k^4)$ time and space, with
731 $k \leq \min\{r^2, mr\}$.

732 *Proof.* We prove it constructively, by using a similar transformation used by [Glasser et al. \(2019\)](#)
733 to represent a non-negative MPS factorization as an hidden Markov model (HMM). Let $\mathbf{X} =$
734 $\{X_1, \dots, X_D\}$ be a set of discrete variables each taking values in $\{1, \dots, m\}$. Let \mathcal{T} be a tensor
735 with D m -dimensional indices. Given an assignment $\mathbf{x} = \langle x_1, \dots, x_D \rangle$ to \mathbf{X} , we factorize \mathcal{T}
736 via a rank r MPS factorization, i.e.,

$$\mathcal{T}[x_1, \dots, x_D] = \sum_{i_1=1}^r \sum_{i_2=1}^r \cdots \sum_{i_{D-1}=1}^r \mathbf{A}_1[x_1, i_1] \mathbf{A}_2[x_2, i_1, i_2] \cdots \mathbf{A}_D[x_D, i_{D-1}] \quad (7)$$

where $\mathbf{A}_1, \mathbf{A}_D \in \mathbb{R}^{m \times r}$ and $\mathbf{A}_j \in \mathbb{R}^{m \times r \times r}$ with $1 < j < D$, for indices $\{i_1, \dots, i_{D-1}\}$ and denoting indexing with square brackets. To reduce \mathcal{T} to being computed by a tensorized structured-decomposable circuit c , i.e., such that $c(\mathbf{x}) = \mathcal{T}[x_1, \dots, x_D]$ for any \mathbf{x} , we perform the following construction. First, we perform a CANDECOMP/PARAFAC (CP) decomposition ([Kolda & Bader, 2009](#)) of each \mathbf{A}_j with $1 < j < D$, i.e.,

$$\mathbf{A}_j[x_j, i_{j-1}, i_j] = \sum_{s_j=1}^k \mathbf{B}_j[i_{j-1}, s_j] \mathbf{V}_j[x_j, s_j] \mathbf{C}_j[i_j, s_j]$$

where $k \leq \min\{r^2, mr\}$ is the maximum rank of the decomposition, and $\mathbf{V}_j \in \mathbb{R}^{m \times k}$, $\mathbf{B}_j \in \mathbb{R}^{r \times k}$, $\mathbf{C}_j \in \mathbb{R}^{r \times k}$. Then, we “contract” each \mathbf{C}_j with \mathbf{B}_{j+1} by computing

$$\mathbf{W}_j[s_j, s_{j+1}] = \sum_{i_j=1}^r \mathbf{C}_j[i_j, s_j] \mathbf{B}_{j+1}[i_j, s_{j+1}]$$

with $\mathbf{W}_j \in \mathbb{R}^{k \times k}$ for $1 < j < D - 1$. In addition, we “contract” \mathbf{C}_{D-1} with \mathbf{A}_D by computing

$$\mathbf{V}_D[x_D, s_{D-1}] = \sum_{i_{D-1}=1}^r \mathbf{C}_{D-1}[i_{D-1}, s_{D-1}] \mathbf{A}_D[x_D, i_{D-1}].$$

In addition, for notation clarity we rename \mathbf{B}_2 with \mathbf{W}_1 and \mathbf{A}_1 with \mathbf{V}_1 . By doing so, we can rewrite Eq. (7) as a sum with indices $\{s_2, \dots, s_{D-1}\}$ over products, i.e.,

$$\begin{aligned} \mathcal{T}[x_1, \dots, x_D] &= \sum_{i_1=1}^r \mathbf{V}_1[x_1, i_1] \sum_{s_2=1}^k \mathbf{W}_2[i_1, s_2] \mathbf{V}[x_2, s_2] \\ &\quad \cdots \sum_{s_{D-1}=1}^k \mathbf{W}_{D-2}[s_{D-2}, s_{D-1}] \mathbf{V}_{D-1}[x_{D-1}, s_{D-1}] \mathbf{V}_D[x_D, s_{D-1}] \end{aligned}$$

737 Fig. B.1 shows an example of such MPS factorization via CP decompositions. We see that we can
 738 encode the products over the same indices using a Hadamard product layers, and summations over
 739 indices $\{s_2, \dots, x_{D-1}\}$ with sum layers parameterized by the \mathbf{W}_j . More precisely, the sum layers
 740 that sum over s_2 and s_{D-1} are parameterized by matrices of ones. Each \mathbf{V}_j with $1 \leq j \leq D$
 741 is instead encoded by an input layer depending on the variable X_j and computing k functions $f_l(X_j)$
 742 such that $f_l(x_j) = \mathbf{V}_j[x_j, l]$ with $1 \leq j \leq k$. The tensorized circuit constructed in this way is
 743 structured-decomposable, as it is defined on a linear tree RG (e.g., Fig. 2a) induced by the variable
 744 ordering implicitly stated by the MPS factorization (Eq. (7), see App. B.4 for details). Fig. B.2
 745 shows the circuit representation corresponding to the MPS reported in Fig. B.1b.

746 Finally, note that the number of parameters of such tensorized circuit correspond to the size of
 747 all \mathbf{W}_j and \mathbf{V}_j introduced above, i.e., $\mathcal{O}(D \cdot k^2)$ where $k \leq \min\{r^2, mr\}$. Moreover, the CP
 748 decompositions at the beginning can be computed using iterative methods whose iterations require
 749 polynomial time (Kolda & Bader, 2009). To retrieve an equivalent BM, we can square the circuit
 750 constructed in this way using Alg. 1, which results in a circuit having size $\mathcal{O}(D \cdot k^4)$ (see Prop. B.1).

751

□

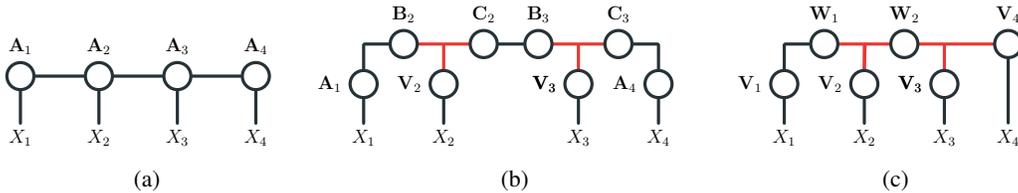


Figure B.1: **Further decomposing a matrix product state (MPS) via CP decompositions.** Tensor networks are represented using the Penrose graphical notation, where circles denote tensors and their connections denote summations over shared indices, and variables X_1, X_2, X_3, X_4 are input indices. Given a MPS (a), we perform a CP decomposition of \mathbf{A}_2 and \mathbf{A}_3 (b). Red edges denote additional indices given by the CP decompositions. Then, we rename \mathbf{A}_1 with \mathbf{V}_1 , \mathbf{B}_2 with \mathbf{W}_1 . Finally, we contract \mathbf{C}_2 with \mathbf{B}_3 , and \mathbf{C}_3 with \mathbf{A}_4 resulting in tensors \mathbf{W}_2 and \mathbf{V}_4 , respectively (c). Fig. B.2 shows the tensorized circuit corresponding to such tensor network, where $\mathbf{V}_1, \mathbf{V}_2, \mathbf{V}_3, \mathbf{V}_4$ and $\mathbf{W}_1, \mathbf{W}_2$ parameterize input layers and sum layers, respectively.

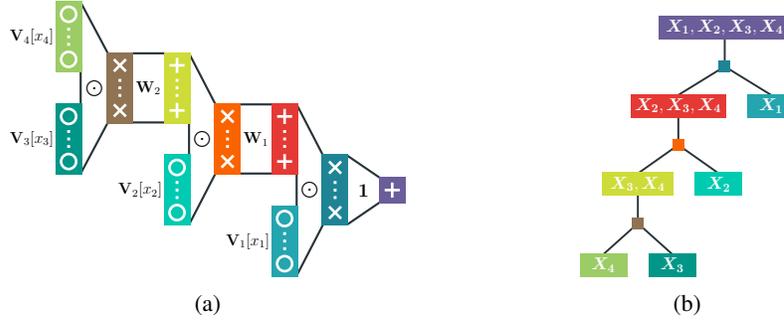


Figure B.2: **Matrix product states (MPS) as structured-decomposable circuits.** The decomposed MPS over three variables showed in Fig. B.1c can be immediately represented as a tensorized structured-decomposable circuit (a) defined on a linear tree RG (b, matching the colors of layers) having Hadamard product layers and sum layers parameterized by \mathbf{W}_1 , \mathbf{W}_2 and a row vector of ones $\mathbf{1}$. Each input layer maps x_1, x_2, x_3, x_4 to rows in $\mathbf{V}_1, \mathbf{V}_2, \mathbf{V}_3, \mathbf{V}_4$, respectively.

752 B.4.1 RELATIONSHIP WITH HIDDEN MARKOV MODELS

753 MPS tensor networks where each tensor \mathbf{A}_i is non-negative can be seen as inhomogeneous hidden
 754 Markov models (HMMs) as showed by Glasser et al. (2019), i.e., where latent state and emitting
 755 transitions do not necessarily share parameters. As such, the tensorized structured-decomposable
 756 circuit c that is equivalent to a MPS (see App. B.4) is also an inhomogenous HMM if c is monotonic.

757 In Sec. 5 we experiment with a tensorized monotonic PC that is an inhomogenous HMM to distill
 758 a large language model, as to leverage the sequential structure of the sentences. We compare it
 759 against a NPC² that is the squaring of a MPS (also called Born machine (Glasser et al., 2019)) or,
 760 equivalently, the squaring of an inhomogenous HMM-like whose parameters can be negative.

761 B.5 EXPONENTIAL SEPARATION

762 **Theorem 1.** There is a class of non-negative functions \mathcal{F} over variables \mathbf{X} that can be compactly
 763 represented as shallow squared NMMs (and hence squared non-monotonic PCs) but for which the
 764 smallest structured-decomposable monotonic PC computing any $F \in \mathcal{F}$ has size $2^{\Omega(|\mathbf{X}|)}$.

765 *Proof.* For the proof of Theorem 1, we start by constructing \mathcal{F} by introducing a variant of the *unique*
 766 *disjointness* (UDISJ) problem, which seems to have first been introduced by De Wolf (2003). The
 767 variant we consider here is defined over graphs, as detailed in the following definition.

768 **Definition B.1** (Unique disjointness function). Consider an undirected graph $G = (V, E)$, where V
 769 denotes its vertices and E its edges. To every vertex $v \in V$ we associate a Boolean variable X_v and
 770 let $\mathbf{X}_V = \{X_v \mid v \in V\}$ be the set of all these variables. The *unique disjointness function* of G is
 771 defined as

$$\text{UDISJ}_G(\mathbf{X}_V) := \left(1 - \sum_{uv \in E} X_u X_v\right)^2. \quad (8)$$

772 **The UDISJ function as a non-monotonic circuit.** We will construct \mathcal{F} as the class of functions
 773 UDISJ_G for graphs $G \in \mathcal{G}$, where \mathcal{G} is a family of graphs that we will choose later. Regardless
 774 of the way the class \mathcal{G} is picked, we can compactly represent UDISJ_G as a squared structured-
 775 decomposable (Def. A.3) and non-monotonic circuit as follows. First, we represent the function
 776 $c(\mathbf{X}_V) = 1 - \sum_{uv \in E} X_u X_v$ as sum unit computing $1 \cdot a(\mathbf{X}_V) + (-1) \cdot b(\mathbf{X}_V)$ where

- 777 • a is a circuit gadget that realizes an unnormalized uniform distribution over the domain
 778 of variables in \mathbf{X}_V , i.e., $a(\mathbf{X}_V) = \prod_{v \in V} (\mathbb{1}\{X_v = 0\} + \mathbb{1}\{X_v = 1\})$ where $\mathbb{1}\{X_v = 0\}$
 779 (resp. $\mathbb{1}\{X_v = 1\}$) is an indicator function that outputs 1 when X_v is set to 0 (resp. 1);

780 • b is another sum unit whose inputs are product units over the input units
 781 $\mathbb{1}\{X_u = 1\}, \mathbb{1}\{X_v = 1\}$ if there is an edge uv in G , i.e., $b(\mathbf{X}_V) = \sum_{uv \in E} \mathbb{1}\{X_u = 1\} \cdot$
 782 $\mathbb{1}\{X_v = 1\}$.

783 Note that b may not be smooth, but we can easily smooth it by adding to every product an additional
 784 input that is a circuit similar to a that outputs 1 for any input $\mathbf{X}_{\overline{uv}}$, where $\mathbf{X}_{\overline{uv}} = \mathbf{X}_V \setminus \{X_u, X_v\}$.
 785 Since c is structured-decomposable (Def. A.3), we can easily multiply it with itself to realize c^2
 786 that would be still a structured-decomposable circuit whose size is polynomially bounded as $|c^2| \in$
 787 $\mathcal{O}(|c|^2)$ (Vergari et al., 2021). In particular, in this case we have that $|c|$ is a polynomial in the
 788 number of variables (or vertices) $|\mathbf{X}_V|$ by the construction above. Furthermore, note that c^2 is non-
 789 monotonic as one of its sum unit has negative parameters (i.e., -1) to encode the subtraction in
 790 Eq. (8).

791 **The lower bound for monotonic circuits.** To prove the exponential lower bound for monotonic
 792 circuits in Theorem 1, we will use an approach that has been used in several other works (Martens
 793 & Medabalimi, 2014; de Colnet & Mengel, 2021). This approach is based on representing a decom-
 794posable circuit (and hence a structured-decomposable one) as a shallow mixture whose components
 795 are *balanced products*, as formalized next.

796 **Definition B.2.** Let \mathbf{X} be a set of variables. A *balanced decomposable product* over \mathbf{X} is a function
 797 from \mathbf{X} to \mathbb{R} that can be written as $f(\mathbf{Y}) \times h(\mathbf{Z})$ where (\mathbf{Y}, \mathbf{Z}) is a partitioning of \mathbf{X} , f and h are
 798 functions to \mathbb{R} and $|\mathbf{X}|/3 \leq |\mathbf{Y}| \leq 2|\mathbf{X}|/3$.

Theorem B.1 (Martens & Medabalimi (2014)). Let F be a non-negative function over Boolean
 variables \mathbf{X} computed by a smooth and decomposable circuit c . Then, F can be written as a sum of
 N balanced decomposable products (Def. B.2) over \mathbf{X} , with $N \leq |c|$ in the form⁴

$$F(\mathbf{X}) = \sum_{k=1}^N f_k(\mathbf{Y}_k) \times h_k(\mathbf{Z}_k),$$

799 where $(\mathbf{Y}_k, \mathbf{Z}_k)$ is partitioning of \mathbf{X} for $1 \leq k \leq N$. If c is structured-decomposable, the N
 800 partitions $\{(\mathbf{Y}_k, \mathbf{Z}_k)\}_{k=1}^N$ are all identical. Moreover, if c is monotonic, then all f_k, h_k only compute
 801 non-negative values.

802 Intuitively, Thm. B.1 tells us that to lower bound the size of c we can lower bound N . To this end, we
 803 first encode the UDISJ function (Eq. (8)) as a sum of N balanced products and show the exponential
 804 growth of N for a family of graphs. We start with a special case for a representation in the following
 805 proposition.

Proposition B.2. Let G_n be a matching of size n , i.e., a graph consisting of n edges none of which
 share any vertices. Assume that the UDISJ function (Eq. (8)) for G_n is written as a sum of products
 of balanced partitions

$$\text{UDISJ}_{G_n}(\mathbf{Y}, \mathbf{Z}) = \sum_{k=1}^N f_k(\mathbf{Y}) \times h_k(\mathbf{Z}),$$

806 where for every edge uv in G_n we have that $X_u \in \mathbf{Y}$ and $X_v \in \mathbf{Z}$. Then $N = 2^{\Omega(n)}$.

807 To prove the above results, we will make an argument on the rank of the so-called *communication*
 808 *matrix*, also known as the *value matrix*, for a function F and a fixed partition (\mathbf{Y}, \mathbf{Z}) .

809 **Definition B.3.** Let F be a function over (\mathbf{Y}, \mathbf{Z}) , its communication matrix M_F is a $2^{|\mathbf{Y}|} \times 2^{|\mathbf{Z}|}$
 810 matrix whose rows (resp. columns) are uniquely indexed by assignments to \mathbf{Y} (resp. \mathbf{Z}) such that
 811 for a pair of index⁵ $(i_{\mathbf{Y}}, j_{\mathbf{Z}})$, the entry at the row $i_{\mathbf{Y}}$ and column $j_{\mathbf{Z}}$ in M_F is $F(i_{\mathbf{Y}}, j_{\mathbf{Z}})$.

⁴In Martens & Medabalimi (2014), Theorem 38, this result is stated with $N \leq |c|^2$. The square materializes
 from the fact that they reduce their circuits to have all their inner units to have exactly two inputs, as we already
 assume, following de Colnet & Mengel (2021).

⁵An index $i_{\mathbf{Y}}$ (resp. $j_{\mathbf{Z}}$) is a complete assignment to Boolean variables in \mathbf{Y} (resp. \mathbf{Z}). See Example 1.

812 **Example 1.** Let us consider a simple matching on 6 vertices, where \mathbf{Y} correspond to the first 3
 813 vertices, and \mathbf{Z} to the last 3, and where there is an edge between the first, second and third vertices
 814 of \mathbf{Y} and \mathbf{Z} . The matrix M_F is an 8-by-8 matrix, a row and a column for each assignment of the 3
 815 binary variables associated to each vertex; it is given by

$\mathbf{Y} \setminus \mathbf{Z}$	000	100	010	001	110	101	011	111
000	1	1	1	1	1	1	1	1
100	1	0	1	1	0	0	1	0
010	1	1	0	1	0	1	0	0
001	1	1	1	0	1	0	0	0
110	1	0	0	1	1	0	0	1
101	1	0	1	0	0	1	0	1
011	1	1	0	0	0	0	1	1
111	1	0	0	0	1	1	1	4

817 Note that the name UDISJ comes from the fact that $M_F(i, j) = 0$ if and only if \mathbf{Y} and \mathbf{Z} share a
 818 single entry equal to 1.

819 In the following, we will rely on the following quantity.

820 **Definition B.4** (Non-negative rank). The non-negative rank of a non-negative matrix $A \in \mathbb{R}_+^{m \times n}$,
 821 denoted $\text{rank}_+(A)$, is the smallest k such that there exist k nonnegative rank-one matrices $\{A_i\}_{i=1}^k$
 822 such that $A = \sum_{i=1}^k A_i$. Equivalently, it is the smallest k such that there exists two non-negative
 823 matrices $B \in \mathbb{R}_+^{m \times k}$ and $C \in \mathbb{R}_+^{k \times n}$ such that $A = BC$.

824 Given a function F written as a sum over N decomposable products (see [Thm. B.1](#)) over a fixed par-
 825 tition (\mathbf{Y}, \mathbf{Z}) , we now show that the non-negative rank of its communication matrix M_F ([Def. B.3](#))
 826 is a lower bound of N .

Lemma B.1. Let $F(\mathbf{X}) = \sum_{k=1}^N f_k(\mathbf{Y}) \times h_k(\mathbf{Z})$ where f_k and h_k are non-negative functions and
 let M_F be the communication matrix ([Def. B.3](#)) of F for the partition (\mathbf{Y}, \mathbf{Z}) , then it holds that

$$\text{rank}_+(M_F) \leq N.$$

827 *Proof.* This proof is an easy extension of the proof of Lemma 13 from [de Colnet & Mengel \(2021\)](#).
 828 Assume w.l.o.g. that $f_k(\mathbf{Y}) \times h_k(\mathbf{Z}) \neq 0$ for any complete assignment to \mathbf{Y} and \mathbf{Z} .⁶ Let M_k
 829 denote the communication matrix of the function $f_k(\mathbf{Y}) \times h_k(\mathbf{Z})$. By construction, we have that
 830 $M_F = \sum_{k=1}^N M_k$. Furthermore, since all values in M_F are non-negative by definition, $\text{rank}_+(M_k)$
 831 is defined for all k and by sub-additivity of the non-negative rank we have that $\text{rank}_+(M_F) \leq$
 832 $\sum_{k=1}^N \text{rank}_+(M_k)$. To conclude the proof, it is sufficient to show that M_k are rank-1 matrices, i.e.,
 833 $\text{rank}_+(M_k) = 1$. To this end, consider an arbitrary k . Since $f_k(\mathbf{Y}) \times h_k(\mathbf{Z}) \neq 0$, there is a row in
 834 M_k that is not a row of zeros. Say it is indexed by $i_{\mathbf{Y}}$, then its entries are of the form $f_k(i_{\mathbf{Y}}) \times h_k(j_{\mathbf{Z}})$
 835 for varying $j_{\mathbf{Z}}$. In any other rows indexed by $i'_{\mathbf{Y}}$ we have $f_k(i'_{\mathbf{Y}}) \times h_k(j_{\mathbf{Z}}) = (f_k(i'_{\mathbf{Y}})/f_k(i_{\mathbf{Y}})) \times$
 836 $f_k(i_{\mathbf{Y}}) \times h_k(j_{\mathbf{Z}})$ for varying $j_{\mathbf{Z}}$. Consequently, all rows are non-negative multiples of the $i_{\mathbf{Y}}$ row,
 837 and therefore $\text{rank}_+(M_k) = 1$. \square

838 To complete the proof of [Prop. B.2](#), we leverage a known lower bound of the non-negative rank of
 839 the communication matrix of the UDISJ problem. The interested reader can find more information
 840 on this result in the books [Roughgarden \(2016\)](#), [Gillis \(2020\)](#) and the references therein.

Theorem B.2 ([Fiorini et al. \(2015\)](#)). Let a UDISJ function defined as in [Prop. B.2](#), and M_{UDISJ} be
 its communication matrix over a partition (\mathbf{Y}, \mathbf{Z}) , then it holds that

$$(3/2)^n \leq \text{rank}_+(M_{\text{UDISJ}}).$$

841 Using [Thm. B.2](#) and [Lem. B.1](#), we directly get [Prop. B.2](#). So we have shown that, for a fixed
 842 partition of variables (\mathbf{Y}, \mathbf{Z}) , every monotonic circuit c encoding the UDISJ function ([Eq. \(8\)](#)) of
 843 a matching of size n has size $|c| \geq 2^{\Omega(n)}$. However, the smallest non-monotonic circuit encoding

⁶If this were not the case we could simply drop the term from the summation, which would clearly reduce the number of summands.

844 the same function has polynomial size in n (see the construction of the UDISJ function as a circuit
845 above). Now, to complete the proof for the exponential lower bound in [Theorem 1](#), we need to
846 find a function class \mathcal{F} where this result holds for all possible partitions (\mathbf{Y}, \mathbf{Z}) . Such function
847 class consists of UDISJ functions over a particular family of graphs, as detailed in the following
848 proposition.

849 **Proposition B.3.** There is a family of graphs \mathcal{G} such that for every graph $G_n = (V_n, E_n) \in \mathcal{G}$ we
850 have $|V_n| = |E_n| = \mathcal{O}(n)$, and any monotonic structured-decomposable circuit representation of
851 UDISJ_{G_n} has size $2^{\Omega(n)}$.

852 *Proof.* We prove it by constructing a class of so-called *expander graphs*, which we introduce next.
853 We say that a graph $G = (V, E)$ has expansion ε if, for every subset V' of V of size at most $|V|/2$,
854 there are at least $\varepsilon|V'|$ edges from V' to $V \setminus V'$ in G . It is well-known, see e.g. [Hoory et al. \(2006\)](#),
855 that there are constants $\varepsilon > 0$ and $d \in \mathbb{N}$ and a family $(G_n)_{n \in \mathbb{N}}$ of graphs such that G_n has at least n
856 vertices, expansion ε and maximal degree d . We fix such a family of graphs in the remainder and
857 denote by V_n , resp. E_n , the vertex set, resp. the edge set, of G_n .

858 Let c be a monotonic structured-decomposable circuit of size N computing UDISJ_{G_n} . Then, by
859 using [Thm. B.1](#), we can write it as

$$\text{UDISJ}_{G_n}(\mathbf{Y}, \mathbf{Z}) = \sum_{k=1}^N f_k(\mathbf{Y}) \times h_k(\mathbf{Z}) \quad (9)$$

where (\mathbf{Y}, \mathbf{Z}) is a balanced partition of \mathbf{X}_V . Let $V_{\mathbf{Y}} = \{v \in V_n \mid X_v \in \mathbf{Y}\}$ and $V_{\mathbf{Z}} = \{v \in V_n \mid X_v \in \mathbf{Z}\}$. Then $(V_{\mathbf{Y}}, V_{\mathbf{Z}})$ form a balanced partition of V_n . By the expansion of G_n , it follows that there are $\Omega(n)$ edges from vertices in $V_{\mathbf{Y}}$ to vertices in $V_{\mathbf{Z}}$. By greedily choosing some of those edges and using the bounded degree of G_n , we can construct an edge set E'_n of size $\Omega(n)$ that is a matching between \mathbf{Y} and \mathbf{Z} , i.e., all edges in E'_n go from \mathbf{Y} to \mathbf{Z} and every vertex in V_n is incident to only one edge in E'_n . Let V'_n be the set of endpoints in E'_n and $\mathbf{X}_{V'_n} \subseteq \mathbf{X}_V$ be the variables associated to them. We construct a new circuit c' from c by substituting all input units for variables X_v that are not in $\mathbf{X}_{V'_n}$ by 0. Clearly, $|c'| \leq |c|$ and hence all the lower bounds for $|c'|$ are lower bounds for $|c|$. Let $\bar{\mathbf{Y}} = \mathbf{X}_{V'_n} \cap \mathbf{Y}$ and $\bar{\mathbf{Z}} = \mathbf{X}_{V'_n} \cap \mathbf{Z}$. By construction c' computes the function

$$\text{UDISJ}_{G'_n}(\bar{\mathbf{Y}}, \bar{\mathbf{Z}}) = \left(1 - \sum_{uv \in E'_n} X_u X_v \right)^2$$

which corresponds to solving the UDISJ problem over the graph $G'_n = (V'_n, E'_n)$. From [Eq. \(9\)](#) we get that

$$\text{UDISJ}_{G'_n}(\bar{\mathbf{Y}}, \bar{\mathbf{Z}}) = \sum_{k=1}^N f'_k(\bar{\mathbf{Y}}) \times h'_k(\bar{\mathbf{Z}}),$$

860 where f'_k (resp. h'_k) are obtained from f_k (resp. h_k) by setting all the variables not in $\mathbf{X}_{V'_n}$ to 0. Since
861 c' is monotonic by construction and $|E'_n| = \Omega(n)$, from [Prop. B.2](#) it follows that $N = 2^{\Omega(n)}$. \square

862 [Prop. B.3](#) concludes the proof of [Theorem 1](#), as we showed the existence of family of graphs for
863 which the smallest structured-decomposable monotonic circuit computing the UDISJ function over
864 n variables has size $2^{\Omega(n)}$. However, the smallest structured-decomposable non-monotonic circuit
865 has size polynomial in n , whose construction has been detailed at the beginning of our proof. \square

866 B.6 SQUARING DETERMINISTIC CIRCUITS

867 In [Sec. 4.1](#) we argued that squaring any non-monotonic, smooth, decomposable ([Def. A.2](#)), and
868 deterministic ([Def. A.5](#)) circuit yields a monotonic and deterministic PC. As a consequence, any
869 function computed by a NPC^2 that is deterministic can be computed by a monotonic and deter-
870 ministic PC. Therefore, we are interested in squaring structured-decomposable circuits that are *not*
871 deterministic. Below we formally prove [Proposition 4](#).

872 **Proposition 4.** Let c be a smooth, decomposable and deterministic circuit over variables \mathbf{X} possibly
 873 computing a negative function. Then, the squared circuit c^2 is monotonic and has the same structure
 874 (hence size) of c .

875 *Proof.* The proof is by induction. Let $n \in c$ be a product unit that computes $c_n(\mathbf{Z}) =$
 876 $\prod_{i \in \text{in}(n)} c_n(\mathbf{Z}_i)$, with $\mathbf{Z} \subseteq \mathbf{X}$ and $(\mathbf{Z}_1, \dots, \mathbf{Z}_{|\text{in}(n)|})$ forming a partitioning of \mathbf{Z} . Then its
 877 squaring computes $c_n^2(\mathbf{Z}) = \prod_{i \in \text{in}(n)} c_n^2(\mathbf{Z}_i)$. Now consider a sum unit $n \in c$ that computes
 878 $c_n(\mathbf{Z}) = \sum_{i \in \text{in}(n)} w_i c_i(\mathbf{Z})$ with $\mathbf{Z} \subseteq \mathbf{X}$ and $w_i \in \mathbb{R}$. Then its squaring computes $c_n^2(\mathbf{Z}) =$
 879 $\sum_{i \in \text{in}(n)} \sum_{j \in \text{in}(n)} w_i w_j c_i(\mathbf{Z}) c_j(\mathbf{Z})$. Since c is deterministic (Def. A.5), for any i, j with $i \neq j$
 880 either $c_i(\mathbf{Z})$ or $c_j(\mathbf{Z})$ is zero for any assignment to \mathbf{Z} . Therefore, we have that

$$c_n^2(\mathbf{Z}) = \sum_{i \in \text{in}(n)} w^2 c_i^2(\mathbf{Z}). \quad (10)$$

881 This implies that in deterministic circuits, squaring does not introduce additional components that
 882 encode (possibly negative) cross-products. The base case is defined on an input unit n that models
 883 a function f_n , and hence its squaring is an input unit that models f_n^2 . By induction c^2 is constructed
 884 from c by squaring the parameters of sum units w_i and squaring the functions f_n modeled by input
 885 units. Moreover, the number of inputs of each sum unit remains the same, as we observe in Eq. (10),
 886 and thus c^2 and c have the same size. \square

887 C EFFICIENT LEARNING OF NPC²S

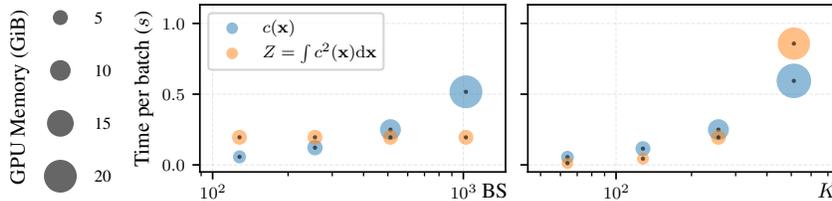


Figure C.1: **Evaluating the squared circuit representation adds little overhead during training.** By learning by MLE (Eq. (4)) and batched gradient descent, the time and space required to compute the partition function Z of c^2 is constant w.r.t. the batch size (BS) (left). By fixing the batch size to 512 and varying the output dimensionality (K) of each layer (right), the resources needed to compute Z are similar to the ones needed to evaluate c (i.e., $c(\mathbf{X})$). For the left figure, we fix $K = 256$ and vary the BS, while for the right figure we fix $\text{BS} = 512$ and vary K . The plots share the y-axis.

888 In this section, we investigate the computational cost of learning NPC²s with a series of benchmarks,
 889 showing that NPC²s add little computational overhead over traditional monotonic PCs (MPCs).

890 **Efficient renormalization in practice.** As suggested by the MLE objective (Eq. (4)), squaring the
 891 tensorized circuit c with Alg. 1 is only required to compute the partition function $Z = \int c^2(\mathbf{x}) dx$.
 892 In addition, we need to compute Z only once per parameter update via gradient ascent, as Z does
 893 not depend on the training data. For these reasons, the increased computational burden of evaluating
 894 a squared circuit (see Proposition 1) as to compute Z is negligible, and it is independent w.r.t. the
 895 batch size. Fig. C.1 illustrates this aspect by comparing the time needed to evaluate c on a batch of
 896 data and to compute the partition function Z . The results showed in Fig. C.1 are obtained by running
 897 benchmarks on NPC²s that are similar in size to the ones we experiment with in Sec. 5. That is,
 898 we benchmark a mixture of 32 NPC²s, each having an architecture built from a randomly-generated
 899 tree RG (see App. F for details) approximating the density function of BSDS300 (the data set with
 900 highest number of variables, see Table H.1). The input layers compute Gaussian distributions.

901 **Training efficiency on UCI data sets.** We benchmark the computational cost of learning NPC²s
 902 on UCI data sets (Table H.1). Fig. C.2 compares time and memory required to learn the best NPC²s
 903 and MPCs showed in Fig. 4, while Fig. C.3 compares time and memory required to learn NPC²s
 904 and MPCs in a worse scenario for NPC²s where the batch size is small and the layer dimension-
 905 ality is large, as NPC²s benefit from using large batch sizes as discussed above. NPC²s add very

906 little overhead during training in most configurations when compared to MPCs, as computing the
 907 partition function Z is comparable to evaluating MPCs on a batch of samples. In particular, on Gas
 908 ($|\mathbf{X}| = 8$), NPC² takes more time and memory to compute Z (times are 6ms and 121ms, while
 909 memory allocations are 0.6GiB and 5.8GiB), but it is only slightly more than the cost of computing
 910 c for MPCs (time 144ms and memory 4.4GiB). Moreover, note that NPC²s achieve about a $\times 2$
 911 improvement on the log-likelihood on Gas. On the much higher dimensional data set BSDS300
 912 ($|\mathbf{X}| = 63$) instead, we found that training NPC² is even cheaper as it requires fewer parameters
 913 while still achieving an higher log-likelihood (128.38 rather than 123.3).

914 **Hardware and significance of benchmarks.** The benchmarks mentioned above and illustrated in
 915 **Figs. C.1 to C.3** have been run on a single NVIDIA RTX A6000 with 48GiB of memory. The
 916 measured times are averaged over 50 independent circuit evaluations.

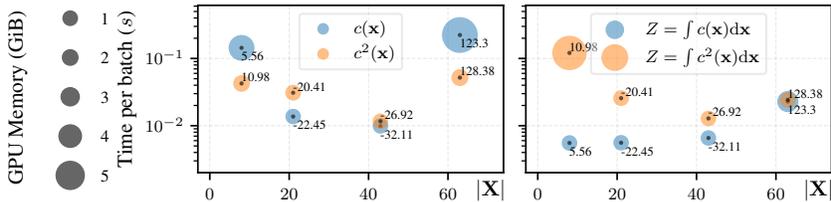


Figure C.2: NPC²s add little overhead during training on real-world data sets, while improving log-likelihoods. We evaluate time and memory required by monotonic PCs (MPCs) and NPC²s to perform one optimization step on UCI data sets (Gas, Hepmass, MiniBooNE, BSDS300) with number of variables $|\mathbf{X}|$ and using the best hyperparameters found (see App. H.3). We benchmark the computation of $c(\mathbf{x})$ by MPCs and $c^2(\mathbf{x})$ by NPC²s on a batch \mathbf{x} of data (left), as well as the partition functions Z for both models (right), and label the data points with the final log-likelihoods achieved by the corresponding models (as also reported in Fig. 4). The plots share the y-axis. For NPC²s, computing the partition function Z is more expensive both in time and memory (right), but it is still very similar to the cost of evaluating $c(\mathbf{x})$ or $c^2(\mathbf{x})$ (left).

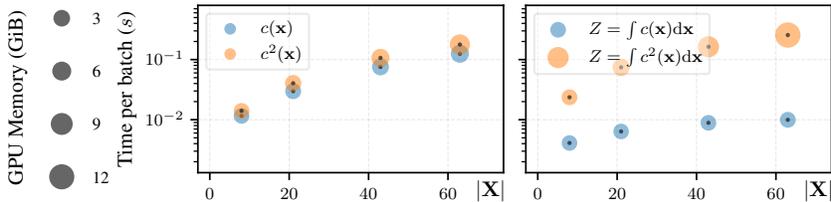


Figure C.3: NPC²s add little overhead during training even with relatively small batch sizes. We evaluate time and memory required by monotonic PCs (MPCs) and NPC²s to perform one optimization step on UCI data sets (Gas, Hepmass, MiniBooNE, BSDS300) with respect to the number of variables $|\mathbf{X}|$ and using the same hyperparameters (512 as batch size, 512 as layer dimensionality, and Gaussian input layers). The plots share the y-axis. The cost of computing $c^2(\mathbf{x})$ on a batch \mathbf{x} of data by NPC²s is only slightly higher than the cost of computing $c(\mathbf{x})$ by MPCs (left), while the cost of computing Z for NPC²s is comparable to evaluating $c^2(\mathbf{x})$ or $c(\mathbf{x})$ (right).

917 D THE SIGNED LOG-SUM-EXP TRICK

918 Scaling squared non-monotonic PCs to more than a few tens (resp. hundreds) of variables without
 919 performing computations in log-space is infeasible in 32-bit (resp. 64-bit) floating point arithmetic,
 920 as we illustrate in Fig. D.1. For this reason, we *must* perform computations in the log-space even
 921 in presence of negative values. The idea is to represent non-zero outputs $\mathbf{y} \in \mathbb{R}^S$ of each layer
 922 in terms of the element-wise logarithm of their absolute value $\log |\mathbf{y}|$ and their element-wise sign
 923 $\text{sign}(\mathbf{y}) \in \{-1, 1\}^S$, i.e., such that $\mathbf{y} = \text{sign}(\mathbf{y}) \odot \exp(\log |\mathbf{y}|)$.

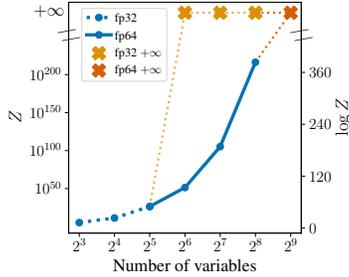


Figure D.1: **Squared non-monotonic PCs cannot scale without performing computations in log-space.** Partition functions (and their *natural* logarithm) of squared non-monotonic PCs having Gaussian input units, with increasing number of variables V and having depth $\lceil \log_2 V \rceil$ computed using 32-bit and 64-bit floating point arithmetic.

In practice, we evaluate product and sum layers according to the following evaluation rules. Given an Hadamard product layer ℓ , then it computes and propagates both $\log |\ell| = \sum_{i=1}^N \log |\ell_i|$ and $\text{sign}(\ell) = \odot_{i=1}^N \text{sign}(\ell_i)$ for some inputs $\{\ell_i\}_{i=1}^N$. Given a sum layer ℓ parameterized by $\mathbf{W} \in \mathbb{R}^{S \times K}$ and having ℓ' as input layer, then it computes and propagates both $\log |\ell| = \alpha + \log |s|$ and $\text{sign}(\ell) = \text{sign}(s)$ where α and s are defined as

$$\alpha = \mathbf{1} \cdot \max_{1 \leq j \leq S} \{\log |\ell'[j]|\} \quad \mathbf{s} = \mathbf{W} (\text{sign}(\ell') \odot \exp(\log |\ell'| - \alpha))$$

924 by assuming $s \neq \mathbf{0}, \mathbf{1}$ denoting a S -dimensional vector of ones, $\ell'[j]$ denoting the j -th entry of the
 925 output of ℓ' , and \exp being applied element-wise. We call *signed log-sum-exp trick* the evaluation
 926 rule above for sum layers, which generalizes the log-sum-exp trick (Blanchard et al., 2021) that is
 927 used to evaluate tensorized monotonic PC architectures (Peharz et al., 2020a).

928 For the more general definition of tensorized circuits instead (Def. A.6), given a Kronecker product
 929 layer ℓ , then it computes both $\log |\ell| = \bigoplus_{i=1}^N \log |\ell_i|$ and $\text{sign}(\ell) = \bigotimes_{i=1}^N \text{sign}(\ell_i)$, where \bigoplus
 930 denotes an operator similar to the Kronecker product but computing sums instead.

931 **E SPLINES AS EXPRESSIVE INPUT COMPONENTS**

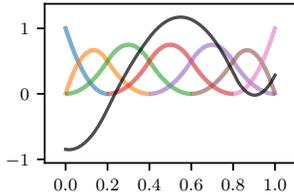


Figure E.1: **Splines represent a class of flexible non-linear functions.** A quadratic ($k = 2$) spline (in black) over $n = 4$ knots chosen uniformly in $(0, 1)$ (i.e. 0.2, 0.4, 0.6 and 0.8) is computed by a linear combination of $n + k + 1 = 7$ distinct basis functions (each colored differently).

932 Polynomials defined on fixed intervals are candidate functions to be modeled by components (resp.
 933 input layers) of squared NMMs (Sec. 2) (resp. NPC²s Sec. 3). This is because they can be negative
 934 function and their product can be tractably integrated. In particular, we experiment with piecewise
 935 polynomials, also called *splines*. An univariate spline function of order k is a piecewise polynomial
 936 defined on a variable X , and the n values of X where polynomials meet are called *knots*. *B-splines*
 937 of order k are basis functions for continuous spline functions of the same degree. In practice, we can
 938 represent any spline function f of order k defined over n knots inside an interval (a, b) as a linear
 939 combination of $n + k + 1$ basis functions, i.e.,

$$f(X) = \sum_{i=1}^{n+k+1} \alpha_i B_{i,k}(X) \tag{11}$$

940 where $\alpha_i \in \mathbb{R}$ are the parameters of the spline and $B_{i,k}(X)$ are polynomials of order k (i.e., the basis
 941 of f), which are unequivocally determined by the choice of the n knots. In particular, each $B_{i,k}(X)$
 942 is a non-negative polynomial that is recursively defined with the Cox-de-Boor formula (de Boor,
 943 1971; Piegel & Tiller, 1995). Given two splines f, g of order k defined over n knots and represented
 944 in terms of $n + k + 1$ basis functions as in Eq. (11), we can write their product integral as

$$\int_a^b f(X)g(X) dX = \sum_{i=1}^{n+k+1} \sum_{j=1}^{n+k+1} \alpha_i \beta_j \int_a^b B_{i,k}(X)B_{j,k}(X) dX \tag{12}$$

945 where $\alpha_i \in \mathbb{R}$ (resp. $\beta_j \in \mathbb{R}$) denote the parameters of f (resp. g). Therefore, integrating a
 946 product of splines requires integrating products of their basis functions. Among the various way of
 947 computing Eq. (12) exactly (Vermeulen et al., 1992), we can do it in time $\mathcal{O}(n^2 \cdot k^2)$ by representing
 948 the product $B_{i,k}(X)B_{j,k}(X)$ as the basis polynomial of another B-spline of order $2k+1$, and finally
 949 integrating it in the interval of definition. Fig. E.1 shows an example of a spline.

950 Since each $B_{i,k}$ is non-negative, we can use B-splines as components (resp. modeled by input layers)
 951 of traditional MMs (resp. monotonic PCs) by assuming each spline parameter α_i to be non-negative.
 952 This is the case of monotonic PCs we experimented with in Sec. 5.

953 F TREE REGION GRAPHS

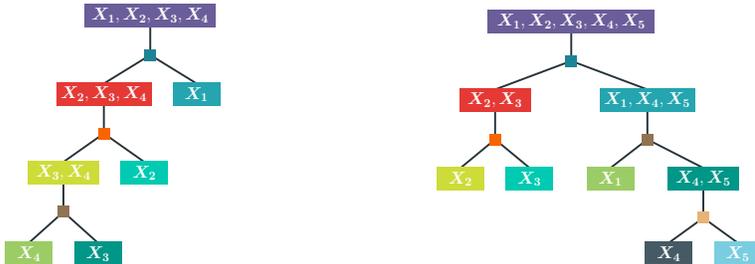


Figure F.1: **Different ways to construct region graphs.** The left figure illustrates a linear tree (LT) region graph (Def. 2) over four variables, which decomposes variables one by one. The right figure shows a possible binary tree (BT) region graph over five variables, which recursively splits them.

954 Since we require structured-decomposability to square circuits (see Sec. 3.2), we construct their
 955 architecture based on tree RGs (Def. 2). We choose to experiment with two kinds of tree RGs:
 956 *binary tree* (BT) and *linear tree* (LT). Following Peharz et al. (2020b), the BT is built by recursively
 957 partitioning variables evenly and randomly until regions with only one variable are obtained. The
 958 LT is built by (1) shuffling the variables randomly and then (2) recursively partitioning variables
 959 one by one, i.e., a set of variables $\{X_i, \dots, X_D\}$ is partitioned in $\{X_i\}$ and $\{X_{i+1}, \dots, X_D\}$ for
 960 $1 \leq i \leq D - 1$. Fig. F.1 shows examples of LT and BT RGs. Note that the LT is the same
 961 on which the circuit representation of matrix-product states (MPS) (Pérez-García et al., 2007) and
 962 TTDE (Novikov et al., 2021) depend on (see also Sec. 4 and App. B.4).

963 G ADDITIONAL RELATED WORKS

964 **Squared neural family (SNEFY)** (Tsuchida et al., 2023) have been concurrently proposed as a
 965 class of models squaring the 2-norm of the output of a single-hidden-layer neural network. Under
 966 certain parametric conditions, SNEFYs can be re-normalized as to model a density function, but
 967 they do not guarantee tractable marginalization of any subset of variables as our NPC²s do, unless
 968 they encode a fully-factorized distribution, which would limit their expressiveness. Hence, SNEFYs
 969 can be employed in our NPC²s to model multivariate units in input layers with bounded scopes.

970 The **rich literature of PCs** provides several algorithms to learn both the structure and the param-
 971 eters of circuits (Poon & Domingos, 2011; Peharz et al., 2017; Di Mauro et al., 2021; Dang et al.,
 972 2021; Liu & Van den Broeck, 2021; Liu et al., 2023). However, in these works circuits are always
 973 assumed to be monotonic. A first work considering subtractions is Dennis (2016) which generalizes
 974 the ad-hoc constraints over Gaussian NMMs (Zhang & Zhang, 2005) to deep PCs over Gaussian
 975 inputs by constraining their structure and reparameterizing their sum weights. Shallow NMM rep-
 976 resented as squared circuits have been investigated for low-dimensional categorical distributions
 977 in (Loconte et al., 2023). Circuit representations encoding probability generating functions allow
 978 negative coefficients, but in symbolic computational graphs (Zhang et al., 2021).

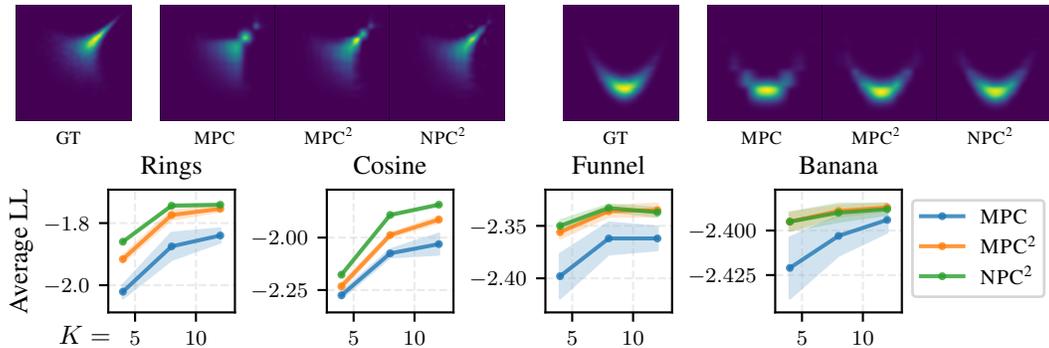


Figure H.1: **Negative parameters increases the expressiveness of NPC²s.** From left to right (above) and for each bivariate density, we show the ground truth (GT) and its estimation by a monotonic PC (MPC), a squared monotonic PC (MPC²), and a NPC² having input layers computing quadratic splines (App. E) and with the same number of parameters. Moreover, (below) we show the average log-likelihoods (and one standard deviation with 10 independent runs) on unseen data achieved by a monotonic MPC, a squared monotonic MPC², and a NPC² by increasing the dimensionality of input layers K .

979 H EXPERIMENTAL SETTINGS AND COMPLEMENTARY RESULTS

980 H.1 CONTINUOUS SYNTHETIC DATA

981 Following (Wenliang et al., 2019) we experiment with monotonic PCs, their squaring and NPC²s
 982 on synthetic continuous 2D data sets, named *rings*, *cosine*, *funnel* and *banana*. We generate each
 983 synthetic data set by sampling 10_000/1_000/2_000 training/validation/test samples. In these ex-
 984 periments, we are interested in studying whether NPC²s can be more expressive in practice, with-
 985 out making assumptions on the data distribution and therefore choosing parametric distributions as
 986 components. For this reason, we choose components computing the product of univariate spline
 987 functions (App. E) over 32 knots that are uniformly chosen in the data domain. In particular, for
 988 monotonic mixtures we restrict the spline coefficients to be non-negative.

989 **Learning and hyperparameters.** Since the data is bivariate, the tree on which PCs are defined
 990 on consists of just one region that is split in half. All models are learned by batched stochastic
 991 gradient descent using the Adam optimizer with default learning rate (Kingma & Ba, 2015) and a
 992 batch size of 256. The parameters of all mixtures are initialized by sampling uniformly between 0
 993 and 1. Furthermore, monotonicity in (squared) PCs is ensured by exponentiating the parameters.

994 Fig. 3 shows the density functions estimated from data sets *rings* and *cosine*, when using 8 and 12
 995 components, respectively. Moreover, Fig. H.1 report the log-likelihoods and other density functions
 996 learned from data sets *funnel* and *banana*, when using 4 components.

997 H.2 DISCRETE SYNTHETIC DATA

998 For our experiments investigating the flexibility of input layers of NPC²s (Sec. 2) in case of discrete
 999 data (Sec. 5), we quantize the bivariate continuous synthetic data sets reported in App. H.1. That is,
 1000 we discretize both continuous variables using 32 uniform bins each. The resulting target distribution
 1001 is therefore a probability mass function over two finitely discrete variables.

1002 We experiment with monotonic PCs, their squaring and NPC²s with two families of input layers.
 1003 First, we investigate very flexible input layers for finitely discrete data: categoricals for monotonic
 1004 PCs and embeddings for NPC²s. Second, we experiment with the less flexible but more parameter-
 1005 efficient Binomials. The learning and hyperparameters setting are the same used for the continuous
 1006 data (see App. H.1). Fig. H.2 shows that there is little advantage in subtracting probability mass
 1007 with respect to monotonic PCs having categorical components. However, in case of the less flex-
 1008 ible Binomial components, NPC²s capture the target distribution significantly better. This is also
 1009 confirmed by the log-likelihoods on unseen data, which we show in Fig. H.2.

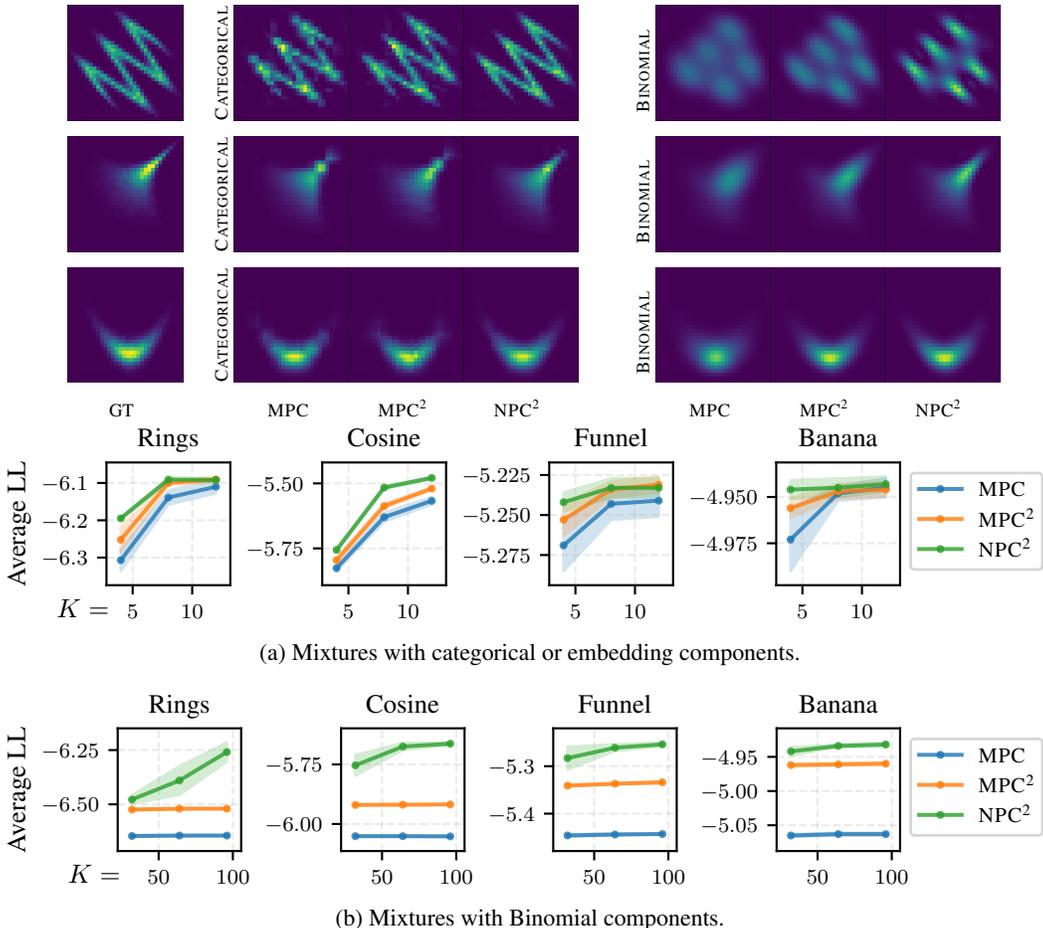


Figure H.2: **Negative parameters increases the expressiveness of NPC²s.** From left to right (above) and for each bivariate distribution, we show the ground truth (GT) and its estimation by a monotonic PC (MPC), a squared monotonic PC (MPC²), and a NPC² having input layers computing categoricals (embeddings for NPC²s) and with the same number of parameters. Moreover, we show the average log-likelihoods (and one standard deviation with 10 independent runs) on unseen data achieved by a monotonic MPC, a squared monotonic MPC², and a NPC² with either categorical (a) or Binomial (b) components and by increasing the dimensionality of input layers K .

1010 H.3 UCI CONTINUOUS DATA

1011 **Data sets.** In Sec. 5 we evaluate NPC²s for density estimation on five multivariate UCI data sets
 1012 (Dua & Graff, 2017): Power (Hebrail & Berard, 2012), Gas (Fonollosa et al., 2015), Hepmass (Baldi
 1013 et al., 2016), MiniBooNE (Roe et al., 2004) and BSDS300 patches (Martin et al., 2001) by following
 1014 the pre-processing by Papamakarios et al. (2017). Table H.1 reports their statistics.

	D	Number of samples		
		train	validation	test
Power	6	1,659,917	184,435	204,928
Gas	8	852,174	94,685	105,206
Hepmass	21	315,123	35,013	174,987
MiniBooNE	43	29,556	3,284	3,648
BSDS300	63	1,000,000	50,000	250,000

Table H.1: **UCI data set statistics.** Dimensionality D and number of samples of each data set split after the pre-processing by Papamakarios et al. (2017).

1015 **Models.** We compare monotonic PCs and NPC²s in tensorized form (Def. 1) for density estimation.
 1016 The tensorized architecture for both is constructed based on either the *binary tree* (BT) or *linear*

1017 *tree* (LT) RGs (see App. F). In addition, since both RGs are randomly-constructed, we instantiate
 1018 eight of them by changing the random seed. By doing so, our monotonic PCs consist of a mixture
 1019 of tensorized monotonic PCs each defined on a different RG. Conversely, our NPC²s consist of a
 1020 mixture (with non-negative parameters) of tensorized NPC²s, each constructed by squaring a circuit
 1021 defined on a different RG. To ensure a fair comparison, monotonic PCs and NPC²s have the exact
 1022 same structure, but NPC²s allow for negative parameters via the squaring mechanism (see Sec. 3).

1023 **Hyperparameters.** We search for hyperparameters by running a grid search with both monotonic
 1024 PCs and NPC²s. For each UCI data set, Tables H.2 and H.3 report the possible value of each
 1025 hyperparameter, depending on the chosen RG. In case of input layers modeling spline functions (see
 1026 App. E), we use quadratic splines and select 512 uniformly in the domain space.

1027 **Parameters initialization.** We found NPC²s to be more sensible to the choice of the initialization
 1028 method for parameters than monotonic PCs. The effect of initialization in monotonic PCs is not well
 1029 explored in the literature, and it is even more unclear for NPC²s as parameters are allowed to be
 1030 negative. In these experiments, we investigated initializing NPC²s by independently sampling the
 1031 parameters from a normal distribution. However, we found NPC²s to achieve higher log-likelihoods
 1032 if they are initialized with non-negative parameters only, i.e., by sampling uniformly between 0 and
 1033 1. Note that our work is a first attempt to learn non-monotonic PCs at scale, thus it opens interesting
 1034 future directions on how to initialize and learn NPC²s.

Table H.2: **Hyperparameter grid search space for each UCI data set (for BT experiments).** Each data set is associated to lists of hyperparameters: learning rate, the dimensionality of layers in tensorized PCs (K), batch size, and whether input layers compute Gaussian likelihoods or spline functions (see App. E).

Data set	Learning rate	K	Batch size	Input layer
Power	[0.01, 0.005]	[32, ..., 512]	[512, 1024, 2048]	[Gaussian, splines]
Gas		[32, ..., 1024]	[512, 1024, 2048, 4096]	
Hepmass		[32, ..., 512]	[512, 1024, 2048]	
MiniBooNE		[32, ..., 512]	[512, 1024, 2048]	
BSDS300		[32, ..., 256]	[512, 1024, 2048]	

Table H.3: **Hyperparameter grid search space for each UCI data set (for LT experiments).** Each data set is associated to lists of hyperparameters: learning rate, the dimensionality of layers in tensorized PCs (K), batch size, and whether input layers compute Gaussian likelihoods or spline functions (see App. E).

Data set	Learning rate	K	Batch size	Input layer
Power	[0.005, 0.001]	[32, ..., 512]	[512, 1024, 2048]	[Gaussian, splines]
Gas				
Hepmass				
MiniBooNE				
BSDS300				

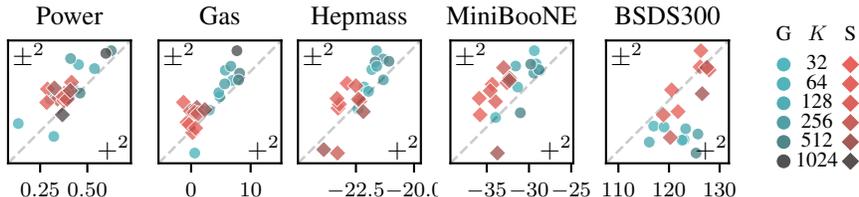


Figure H.3: **Negative parameters make squared non-monotonic PCs more expressive than squared monotonic PCs.** NPC²s (\pm^2 , vertical) generally achieve higher log-likelihoods than squared monotonic PCs ($+^2$, horizontal) when paired with the same number of units per layer K . as shown by the presence of more points in the upper triangle than in the lower triangle for most data sets. Blue circles \bullet and red diamonds \blacklozenge refer to runs with Gaussian (G) and spline (S) input layers respectively, and darker hues indicate larger K . The dashed grey line represents the points of equal log-likelihood for both the NPC² and the squared monotonic PC.

Table H.4: **Squared non-monotonic PCs can be more expressive than monotonic PCs.** Best average test log-likelihoods and two standard errors achieved by monotonic PCs (MPC) and NPC²s built either from randomized linear tree RGs (LT) or from randomized binary tree RGs (BT) (see App. H.3), when compared to baselines. MPC, MPC² and NPC² were experimented with both Gaussian (G) and spline (S) node input layers. † means no values were originally provided.

	Power		Gas		Hepmass		MiniBooNE		BSDS300	
MADE	-3.08 ±0.03		3.56 ±0.04		-20.98 ±0.02		-15.59 ±0.50		148.85 ±0.28	
RealNVP	0.17 ±0.01		8.33 ±0.14		-18.71 ±0.02		-13.84 ±0.52		153.28 ±1.78	
MAF	0.24 ±0.01		10.08 ±0.02		-17.73 ±0.02		-12.24 ±0.45		154.93 ±0.28	
NSF	0.66 ±0.01		13.09 ±0.02		-14.01 ±0.03		-9.22 ±0.48		157.31 ±0.28	
Gaussian	-7.74 ±0.02		-3.58 ±0.75		-27.93 ±0.02		-37.24 ±1.07		96.67 ±0.25	
EiNet-LRS	0.36 ±†		4.79 ±†		-22.46 ±†		-34.21 ±†		†	
TTDE	0.46 ±†		8.93 ±†		-21.34 ±†		-28.77 ±†		143.30 ±†	
	G	S	G	S	G	S	G	S	G	S
MPC (LT)	0.51 ±0.01	0.24 ±0.01	6.73 ±0.03	-2.05 ±0.02	-22.07 ±0.02	-23.09 ±0.02	-32.48 ±0.44	-37.53 ±0.46	123.15 ±0.28	116.90 ±0.28
MPC ² (LT)	0.49 ±0.01	0.39 ±0.01	7.06 ±0.03	0.95 ±0.01	-21.42 ±0.02	-22.24 ±0.02	-29.46 ±0.44	-32.81 ±0.47	—	—
NPC ² (LT)	0.53 ±0.01	0.43 ±0.01	9.00 ±0.02	3.03 ±0.02	-20.66 ±0.02	-21.53 ±0.02	-26.68 ±0.42	-29.36 ±0.42	112.99 ±0.29	120.11 ±0.29
MPC (BT)	0.57 ±0.01	0.32 ±0.01	5.56 ±0.03	-2.55 ±0.02	-22.45 ±0.02	-24.09 ±0.02	-32.11 ±0.43	-37.56 ±0.46	121.92 ±0.29	123.30 ±0.29
MPC ² (BT)	0.57 ±0.01	0.36 ±0.01	8.24 ±0.03	0.32 ±0.02	-21.47 ±0.02	-23.38 ±0.02	-29.46 ±0.43	-33.43 ±0.47	125.56 ±0.29	126.85 ±0.29
NPC ² (BT)	0.63 ±0.01	0.45 ±0.01	10.98 ±0.02	3.12 ±0.01	-20.41 ±0.02	-22.25 ±0.02	-26.92 ±0.44	-30.81 ±0.54	114.47 ±0.28	128.38 ±0.29

Table H.5: Table showing average test set log-likelihoods and one standard deviation achieved from running experiments 5 times with random parameters initialization, using the same hyperparameters that were used for achieving results showed in Table H.4.

	Power		Gas		Hepmass		MiniBooNE		BSDS300	
MPC (LT)	0.46 ±0.03	7.03 ±0.18	-22.07 ±0.02	-31.79 ±0.39	126.66 ±5.46					
MPC (BT)	0.53 ±0.03	6.16 ±0.56	-22.42 ±0.45	-33.30 ±0.98	122.77 ±0.71					
NPC ² (LT)	0.42 ±0.11	8.97 ±0.08	-20.67 ±0.05	-29.58 ±0.29	127.58 ±4.66					
NPC ² (BT)	0.62 ±0.01	10.55 ±0.39	-20.48 ±0.11	-27.64 ±0.44	128.45 ±0.52					

Table H.6: Table listing the hyperparameters combinations found via a grid search, which were used for achieving results showed in Table H.4. For input layers, G and S respectively denote Gaussian and spline.

Model	Data set	K	Batch size	Learning rate	Input layer
MPC (BT)	Power	512	512	0.01	G
	Gas	1024	4096	0.01	G
	Hepmass	128	512	0.01	G
	MiniBooNE	32	512	0.01	G
	BSDS300	512	512	0.01	S
MPC (LT)	Power	512	512	0.001	G
	Gas	512	1024	0.001	G
	Hepmass	512	512	0.005	G
	MiniBooNE	512	1024	0.005	G
	BSDS300	64	512	0.005	S
NPC ² (BT)	Power	512	512	0.01	G
	Gas	1024	512	0.01	G
	Hepmass	256	512	0.01	G
	MiniBooNE	32	512	0.01	G
	BSDS300	128	512	0.01	S
NPC ² (LT)	Power	512	512	0.001	G
	Gas	512	512	0.001	G
	Hepmass	256	512	0.001	G
	MiniBooNE	128	2048	0.005	G
	BSDS300	32	1024	0.001	S

1035 H.4 LARGE LANGUAGE MODEL DISTILLATION

1036 **Data set.** Given $p^*(\mathbf{x})$ the distribution modeled by GPT2 over sentences $\mathbf{x} = [x_1, \dots, x_D]$ having
1037 maximum length D , we aim to minimize the Kullback-Leibler divergence $\text{KL}[p^* \mid p]$, where p is
1038 modeled by a PC. Minimizing such divergence is equivalent to learn the PC by maximum-likelihood
1039 on data sampled by GPT2. Therefore, following the experimental setting by Zhang et al. (2023)
1040 we sample a data set of 4M sentences using GPT2 having bounded length $D = 32$, i.e., with a
1041 maximum of $D = 32$ tokens. However, differently from Zhang et al. (2023), we used half the
1042 number of sentences (i.e., 4M instead of 8M) due to training time limitations. A larger number of
1043 sampled sentences can be useful to reduce the overfitting we observed with NPC²s (see Fig. 5).

1044 **Models.** Then, we learn a monotonic PC and a NPC² as tensorized circuits whose architecture is
1045 determined by a linear tree RG (Def. 2), i.e., a region graph that recursively partitions each set of
1046 finitely-discrete variables $\{X_i, \dots, X_D\}$ into $\{X_i\}$ and $\{X_{i+1}, \dots, X_D\}$ for $1 \leq i \leq D - 1$ (e.g.,
1047 see Fig. 2a). This is because we are interested in exploiting the sequential dependencies between
1048 words in a sentence. By enforcing monotonicity, we recover that the monotonic PC is equivalent to
1049 an inhomogenous hidden Markov model (HMM), and that that NPC² corresponds to a Born machine
1050 (see App. B.4.1 for details).

1051 **Hyperparameters.** All PCs are learned by batched stochastic gradient descent using (Kingma &
1052 Ba, 2015) as optimizer using batch size 4096. We perform multiple runs by exploring combinations
1053 of learning rates and initialization. For monotonic PCs, we run experiments by choosing learning
1054 rates in $\{5 \cdot 10^{-3}, 10^{-2}, 5 \cdot 10^{-2}\}$ and initializing parameters by sampling uniformly in $(0, 1)$, by
1055 sampling from a standard log-normal distribution, and from a Dirichlet distribution with concentra-
1056 tion values set to 1. Similarly for NPC²s, we run experiments by choosing the same learning rates
1057 for monotonic PCs, but using different initialization. In addition to sampling uniformly in $(0, 1)$,
1058 we also initialize the parameters by sampling from a standard normal distribution. By doing so, we
1059 initialize an approximately even number of positive and negative parameters. Moreover, we also
1060 initialize parameters by sampling from a normal distribution with mean 1 and standard deviation 1,
1061 which initializes more parameters to be positive.