

Elevating Perceptual Sample Quality in Probabilistic Circuits through Differentiable Sampling

Steven Lang

TU Darmstadt

STEVEN.LANG@CS.TU-DARMSTADT.DE

Martin Mundt

TU Darmstadt, Hessian Center for AI (hessian.AI), Darmstadt

MARTIN.MUNDT@CS.TU-DARMSTADT.DE

Fabrizio Ventola

TU Darmstadt

VENTOLA@CS.TU-DARMSTADT.DE

Robert Peharz

TU Graz

R.PEHARZ@TUE.NL

Kristian Kersting

TU Darmstadt, Hessian Center for AI (hessian.AI), Darmstadt

KERSTING@CS.TU-DARMSTADT.DE

Abstract

Deep generative models have seen a dramatic improvement in recent years, due to the use of alternative losses based on perceptual assessment of generated samples. This improvement has not yet been applied to the model class of probabilistic circuits (PCs), presumably due to significant technical challenges concerning differentiable sampling, which is a key requirement for optimizing perceptual losses. This is unfortunate, since PCs allow a much wider range of probabilistic inference routines than main-stream generative models, such as exact and efficient marginalization and conditioning. Motivated by the success of loss reframing in deep generative models, we incorporate perceptual metrics into the PC learning objective. To this aim, we introduce a differentiable sampling procedure for PCs, where the central challenge is the non-differentiability of sampling from the categorical distribution over latent PC variables. We take advantage of the Gumbel-Softmax trick and develop a novel inference pass to smoothly interpolate child samples as a strategy to circumvent non-differentiability of sum node sampling. We initially hypothesized, that perceptual losses, unlocked by our novel differentiable sampling procedure, will elevate the generative power of PCs and improve their sample quality to be on par with neural counterparts like probabilistic auto-encoders and generative adversarial networks. Although our experimental findings empirically reject this hypothesis for now, the results demonstrate that samples drawn from PCs optimized with perceptual losses can have similar sample quality compared to likelihood-based optimized PCs and, at the same time, can express richer contrast, colors, and details. Whereas before, PCs were restricted to likelihood-based optimization, this work has paved the way to advance PCs with loss formulations that have been built around deep neural networks in recent years.

Keywords: Pre-registration, Machine Learning, Generative Models, Probabilistic Circuits

1. Introduction

The central task of approximating data-generating distributions by means of probabilistic models has experienced impressive improvements with the advent of deep learning. However, when considering the in-depth novelties that underlie the corresponding advancements,

it seems that several improvements are independent of the neural network and its parameters. The neural architecture of encoders and decoders involved in state-of-the-art deep generative models is predominantly shared among methods. Instead, it could be argued that major achievements were an immediate result of iterative stages of reframing the optimization and the particular employed objectives.

For instance, for natural images, the initial sampled generation on the basis of denoising auto-encoders (Bengio et al., 2013) (DAE), or variational inference in neural networks (Graves, 2011; Kingma and Welling, 2013) (VAE), was rapidly improved upon with the inclusion of a min-max objective that learns to separate real from fake data in a discriminator of a generative adversarial network (GAN) (Goodfellow et al., 2014). The latter approach, promoted to emphasize a “perceptual” measure over pixel values, has then been brought back into a hybrid VAEGAN model (Larsen et al., 2016). Subsequently, several concurrent realizations of adversarial training have argued that an encoder-decoder is sufficient to formulate a min-max adversarial objective (Ulyanov et al., 2018; Huang et al., 2018), without the presence of an extra discriminator. The loop has seemingly been closed with the proposal of the perceptual auto-encoder (Zhang et al., 2020), linking advancements back to the original neural starting point. Orthogonal developments have simultaneously improved generation quality through the addition of autoregressive sampling steps (Gulrajani et al., 2017; Chen et al., 2017) and losses motivated from a perspective of optimal transport with the incorporation of Wasserstein distances (Arjovsky et al., 2017; Tolstikhin et al., 2018).

Motivated by this evolution of deep generative models, we posit that the prevalence of (deep) neural networks is not primarily due to the nature of their neural architectural design and computational operations, but rather from objectives that directly relate samples from the approximated distribution to ground-truth data instances, in an attempt to introduce perceptual metrics into training. To date, the latter practice is absent from tractable probabilistic model alternatives, such as the family of probabilistic circuits (PCs) (Darwiche, 2002, 2003; Hoifung and Pedro, 2011; Rahman et al., 2014; Kisa et al., 2014; Vergari et al., 2020), which focus their training efforts on maximum likelihood estimation (MLE). Despite these models having a crucial advantage of providing general and computationally efficient inference, a large scale demonstration of their generative modeling capabilities comparable to their neural network counterparts still remains open. *Our paper’s central hypothesis is that the perceived gap to neural networks can be bridged and probabilistic circuits provide an equally adequate alternative, if their objective is posed from an analogous perspective. That is, we introduce to PCs a formulation of objectives over samples from the approximated distribution and ground-truth instances.*

To gather experimental evidence in support of our hypothesis, we propose to introduce differentiable sampling into PCs. This provides the currently missing component required for flexible loss formulation to include similarity measures between samples and data. Therefore, we present a novel differentiable sampling procedure for PCs, which provides continuous approximations to the non-differentiable operations in the standard sampling procedure. Finally, an evaluation of PCs on methods such as probabilistic auto-encoding (Kingma and Welling, 2013), adversarial training (Goodfellow et al., 2014), and maximum mean discrepancy optimization (Gretton et al., 2012) will be enabled.

2. Related Work

Probabilistic circuits (PCs) are a set of expressive probabilistic models that provide exact and tractable inference for a wide range of probabilistic queries. Well-known representatives of PCs are arithmetic circuits (Darwiche, 2002, 2003), cutset networks (Rahman et al., 2014), probabilistic sentential decision diagrams (Kisa et al., 2014), and sum-product networks (SPNs) (Hoifung and Pedro, 2011). PCs are graphical models which connect a set of random variables (RVs), modeled in leaf nodes, by repeatedly and hierarchically assuming independence between RVs (product node) and building convex mixtures over these local independencies (sum node). For our purpose, we concentrate on SPNs, although the findings of this work will be equally applicable to every member of PCs. SPNs constrain their graphs to be *smooth*, i.e. all children of a sum node must have equal scope, in other words, the distribution they encode should be defined over the same set of RVs. Moreover, SPNs are *decomposable*, i.e. all children of a product node must have pairwise disjoint scope. These constraints allow for efficient marginalization and conditioning routines, providing great representational power and efficiency for probabilistic queries. Recent work on the computational efficiency (Shah et al., 2020; Sommer et al., 2021), and specialized SPN structures such as EinsumNetworks (Peharz et al., 2020) and LibSPN (Pronobis et al., 2017), leverage modern GPU architectures to further scale SPN model complexity and reduce the gap to deep neural networks in terms of model capacity.

SPN parameters, i.e. sum node weights and leaf node distribution parameters, are optimized by means of maximum likelihood estimation. To date, SPNs cannot make use of flexible formulations of optimization objectives that include perceptual loss terms involving the generated samples. In contrast, due to the continuous nature of the way deep neural components are arranged, the generation of samples in deep probabilistic models admits direct formulation of objectives $\mathcal{L}(\mathbf{x}, \mathbf{x}^*)$ in the input domain, where \mathbf{x} is a data instance and \mathbf{x}^* is a sample generated by the model.

2.1. Flexible Loss Formulations in Neural Networks

As an early example, the generalization of denoising auto-encoders in Bengio et al. (2013) has formulated the denoising procedure from a probabilistic perspective. That is, given a data instance \mathbf{x} , a noisy sample $\hat{\mathbf{x}} \sim C(\hat{\mathbf{x}} | \mathbf{x})$ is generated from a corruption process C and the expected value of $-\log P_\theta(\mathbf{x} | \hat{\mathbf{x}})$, a reconstruction loss, is then minimized. A variational inference approach to auto-encoding has been suggested in Kingma and Welling (2013). Here, a variational approximation $q_\theta(\mathbf{z} | \mathbf{x})$, encoded through the parameters of a neural network, serves as the approximation to the true posterior $p(\mathbf{z} | \mathbf{x})$. A lower bound on the data likelihood $p(\mathbf{x})$, consisting of a reconstruction error between the data instance and a reconstruction $\mathbf{x}^* \sim p(\mathbf{x} | \mathbf{z})$ from the decoder, and a negative Kullback-Leibler divergence term that represents the distance between the variational posterior and a typically isotropic Gaussian prior $p(\mathbf{z}) = \mathcal{N}(0, 1)$, are optimized jointly with the help of a reparametrization trick.

A different strategy has been proposed by Goodfellow et al. (2014) who construct a min-max game objective. The introduced model consists of two components i.e. a generator G , generating samples, and a discriminator D , that discerns whether an input is generated by G (fake) or whether it is a true data instance (real) by incorporating a perceptual loss on the

quality of samples generated by G . The objectives are formulated such that D maximizes its accuracy, whereas G is pushed to generate samples that fool D , thus, reducing its accuracy. Since the sample generation by G and the fake-real discrimination by D are differentiable, the two components can be trained in an end-to-end fashion, allowing for direct optimization of perceived sample quality.

Therefore, a crucial component seems to be differentiability in sample generation, which is still missing from PC formulations. In this work, we address this challenge and thus bring a broader set of optimization objectives to PCs, similar to the one found in deep neural models. In Section 3, we first summarize the parts of the standard sampling in PCs that break differentiability and then continue to propose a novel procedure that permits gradient flow from model samples to parameters, therefore, making PCs amenable to the aforementioned flexible loss formulations.

2.2. Complementary Advances

Whereas the previous section has highlighted the role of differentiability in sampling and the consequent possibilities of flexible loss formulations as a key element of sample quality improvements, there naturally exist other avenues that have lead to further advances in recent years. For example, sophisticated data-augmentation techniques have been proposed and are commonly employed in practice. To name a few, geometric transformations, color space augmentations, image mixing, random erasing, feature space augmentation, all lead to respective model performance boosts, as summarized in several timely surveys (Shorten and Khoshgoftaar, 2019; Feng et al., 2021; Wen et al., 2021). Although interesting, these are fully complementary to perceptual loss formulations and the concept under investigation in this work. Therefore, the additional inclusion of these auxiliary efforts should be inspected in separate future work.

Another line of recent concurrent work are diffusion models (Kingma et al., 2021; Song et al., 2021; Ho et al., 2020). The latter sample by reversing a gradual noising process and, regarding sample quality, have been shown to be on par with GANs (Dhariwal and Nichol, 2021). Whereas diffusion models are conceptually different, we argue that they can also be perceived as a fundamental way of re-framing the optimization objective from a perceptual point of view, i.e. learning to produce a slightly less noisy \mathbf{x}_{t-1} from \mathbf{x}_t until reaching a final sample \mathbf{x}_0 .

3. Method

To gain comparable benefits from employing perceptual loss quantities as deep neural models do, we hypothesize that permitting gradient flow during sample generation w.r.t. the model parameters can elevate PCs to similar data generation quality as deep neural models.

3.1. Sampling in Sum-Product Networks

Sampling in SPNs is performed in a top-down fashion by starting at the root node, descending into either all of its children, in the case of a product node, or one of its children, in the case of a sum node. When a leaf node is reached, the distribution modeled by that

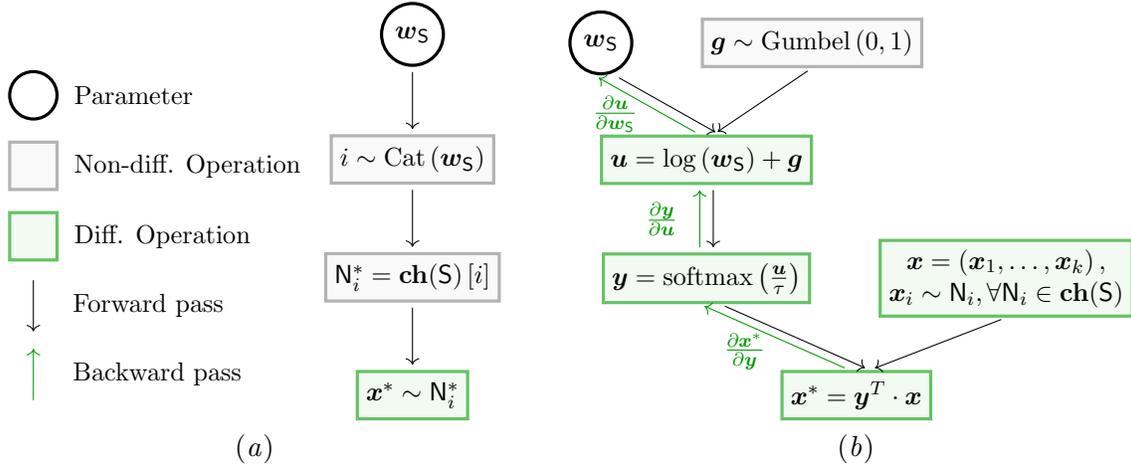


Figure 1: Computation graphs for the 1(a) standard and the 1(b) differentiable sampling procedure in an SPN sum node S . Gray nodes indicate non-differentiable operations which break the gradient flow on the path $w_S \rightarrow x^*$. Using the Gumbel-Softmax Trick combined with smooth child interpolation, we are able to construct a computation path $w_S \rightarrow x^*$, which consists of only smooth functions, allowing gradient computation $\partial x^* / \partial w_S$ between the generated sample x^* and the sum node parameters w_S using the chain-rule.

leaf node is sampled. The particular sampling procedures for sum and product nodes are defined in the following.

Product nodes in SPNs are decomposable and, thus, have children with pairwise non-overlapping scopes. Therefore, the sampling procedure for a product node P simply requires to sample from all of its children $\mathbf{ch}(P)$, i.e. $x^* = \{x \sim N \mid \forall N \in \mathbf{ch}(P)\}$. Figure 1(a) depicts the sampling procedure for a sum node S sampling procedure as a computation graph. Interpreting the sum node weights $w_S = \{w_{S,N} \mid N \in \mathbf{ch}(S)\}$ as categorical probabilities, we first sample from a categorical distribution $i \sim \text{Cat}(w_S)$, where i is used as an index to obtain the sampled sum node child $N_i^* = \mathbf{ch}(S)[i]$ and we finally obtain the actual sample $x^* \sim N_i^*$.

The sum node sampling steps break differentiability in two points as highlighted in Figure 1(a). The first operation is the drawing from a categorical distribution, which is a not differentiable due to its discrete nature. The second issue arises when obtaining the sampled child node N^* by indexing the set of all children $\mathbf{ch}(S)$ which is a discrete operation. Indexing prohibits gradients w.r.t. the index itself, thus, breaking differentiability. In the following, we explain how to tackle these challenges and formulate a differentiable sampling procedure for SPNs.

3.2. Differentiable Sampling in Sum-Product Networks

To leverage differentiability during sampling in SPNs, we propose a novel sampling procedure that differs from the standard sampling procedure in two ways. We first replace the discrete categorical sampling in sum nodes by a continuous approximation, making use of the well-known Gumbel-Softmax Trick (Jang et al., 2017; Maddison et al., 2017). Then, we perform smooth interpolation of sum node children to overcome the discretization introduced by indexing a single child.

Gumbel-Softmax Trick Let z be a categorical variable with probabilities π_1, \dots, π_k , the Gumbel-*Max* trick allows us to draw samples z from a categorical distribution $z = \arg \max_i \{\log(\pi_i) + g_i\}$, with gumbel distributed noise $g_i \sim \text{Gumbel}(0, 1)$. To relax the discretization introduced by the $\arg \max$ operation, Jang et al. (2017) have proposed to use the softmax function as a continuous, differentiable approximation to $\arg \max$, leading to a k -dimensional sample vector \mathbf{y} in the simplex Δ^{k-1} where $\mathbf{y} = \text{softmax}((\boldsymbol{\pi} + \mathbf{g})/\tau)$. The temperature τ controls the approximation precision. As τ goes towards zero, the Gumbel-Softmax Distribution (Jang et al., 2017) converges to the categorical distribution with probabilities π_i . For $\tau > 0$, the Gumbel-Softmax distribution is smooth and has a well-defined gradient $\partial y_i / \partial \pi_i$ with respect to its parameters π_i .

Smooth Child Interpolation For the sake of simplicity, let us first explore the issue of discretization induced by child node indexing for the case when $\tau \rightarrow 0$. After applying the Gumbel-Softmax Trick, sampling from a sum node S results in a one-hot encoded vector \mathbf{y} with $y_i = 1$ at exactly one position (the sampled child node index) and $y_j = 0$ for all $j \neq i$. Let’s assume a scenario in which we have already sampled from all children of S , i.e. we have a vector $\mathbf{x}^* = (\mathbf{x}_1^* \sim \mathbf{N}_1, \dots, \mathbf{x}_k^* \sim \mathbf{N}_k)^T$. In the conventional non-differentiable perspective, we are now required to index this vector with the sampled index i . Due to the one-hot encoding of the vector \mathbf{y} , we can perform the dot-product between \mathbf{x}^* and \mathbf{y} to obtain an interpolation between all samples $(\mathbf{x}_1^*, \dots, \mathbf{x}_k^*)^T \cdot \mathbf{y} = \mathbf{x}_i^*$. In the general case of $\tau > 0$, then, the dot product will result in a mixture over all possible child samples. Each child sample contributes to the final sample according to the sampled mixture weight y_i . The computation graph of the differentiable sampling procedure is shown in Figure 1(b). Alternatively, we can also use the Straight-Through Gumbel Estimator (Jang et al., 2017) by discretizing \mathbf{y} with $\arg \max$ in the forward pass, but use the continuous approximation $\nabla_{\boldsymbol{\pi}} \mathbf{z} \approx \nabla_{\boldsymbol{\pi}} \mathbf{y}$ in the backward pass.

Although this sampling procedure is differentiable, it is not identical to the sampling from the corresponding categorical distribution for $\tau > 0$. This leads to a trade-off during parameter learning, where small τ results in almost one-hot encoded samples and gradients with high variance, and large τ leads to smooth samples with low gradient variance. To allow for stable training, the temperature τ can either be annealed or learned as a parameter, which can be interpreted as entropy regularization (Szegedy et al., 2016; Pereyra et al., 2017). In this case, the Gumbel-Softmax distribution adaptively adjusts the “confidence” of the samples during the training process (Jang et al., 2017).

Moreover, by creating these smooth top-down evaluations through the network graph using the Gumbel-Softmax Trick, whereas out of scope for the investigation of this work, it is now possible to further consider variance reduction techniques to deal with discrete

distributions in differentiable programs, such as REBAR (Tucker et al., 2017), wake-sleep algorithms (Hinton et al., 1995), or REINFORCE (Williams, 1992).

4. Experimental Protocol

Our experimental protocol is devised to demonstrate the effectiveness of our contributions, namely, the introduced differentiable sampling for PCs to enable flexible objectives. Our main hypothesis is that PCs equipped with the latter can be on par with deep neural models like GANs and VAEs in sample generation. For this purpose, we leverage the implementation of PCs by Peharz et al. (2020), where we implement our proposed method, Section 3.1, to take advantage of the flexible loss formulations as seen in neural networks and investigate their impact in PCs. The source code is released at <https://github.com/ml-research/differentiable-sampling-pc>. Given that this is uncharted territory for PCs, as previous works did not prioritize perceived image sample quality, our focus is to investigate SPNs with differentiable sampling in the following three concrete settings.

4.1. Flexible Loss Experiments

1. Adversarial Training We adopt the min-max optimization objective of GANs (Goodfellow et al., 2014):

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log (1 - D(G(\mathbf{z})))] .$$

In our case, the generator G is an SPN \mathcal{S}_G and the second term of the objective becomes the expectation over samples drawn from \mathcal{S}_G : $\mathbb{E}_{\mathbf{x} \sim \mathcal{S}_G(\mathbf{x})} [\log (1 - D(\mathbf{x}))]$. The discriminator D can be either: 1) The generator SPN \mathcal{S}_G itself, similar in spirit to adversarial training with introspection (Huang et al., 2018), 2) a separate, discriminative, SPN \mathcal{S}_D (Gens and Domingos, 2012), 3) another arbitrary model, e.g. a neural network. We then compare the sample quality with the original GAN, where G and D are modeled by a neural network. For a fair comparison, we will also use the discriminative SPN with the neural generator.

2. Maximum Mean Discrepancy We train SPNs by using and minimizing the Maximum Mean Discrepancy (Gretton et al., 2012) distance:

$$\text{MMD}^2(\mathbf{x}, \mathbf{y}) = \frac{1}{\binom{n}{2}} \sum_{i \neq j} k(\mathbf{x}_i, \mathbf{x}_j) - \frac{2}{\binom{n}{2}} \sum_{i \neq j} k(\mathbf{x}_i, \mathbf{y}_j) + \frac{1}{\binom{n}{2}} \sum_{i \neq j} k(\mathbf{y}_i, \mathbf{y}_j) ,$$

as proposed in MMD-GANs (Li et al., 2017). The core idea is to employ the kernel two-sample test, denoted by k in above equation, instead of a discriminator with the goal of minimizing the distance between the target and the estimated distributions. In our evaluation, we contrast against MMD-GAN (Li et al., 2017) and employ adversarially learned kernels for both, the MMD-GAN and the SPN.

3. Probabilistic Auto-Encoding We extend SPNs by explicitly modeling the distribution of the latent space, similar to VAEs. That is, we learn an SPN $\mathcal{S}(\mathbf{x}, \mathbf{z})$ over the joint of the input \mathbf{x} and some latent space variables \mathbf{z} via auto-encoding. Given a data point $\mathbf{x}_i \sim p_{\text{data}}(\mathbf{x})$, we compute its latent representation with the most probable explanation

$\mathbf{z}_i = \arg \max_{\mathbf{z}} \mathcal{S}(\mathbf{z} | \mathbf{x}_i)$. Then, we generate its reconstruction by sampling $\hat{\mathbf{x}}_i \sim \mathcal{S}(\mathbf{x} | \mathbf{z}_i)$. We compare the results with a VAE (Kingma and Welling, 2013) that has a Gaussian prior. For a fair comparison, we also model the SPN latent variables \mathbf{z} as Gaussians.

4.2. Fair Experiment and Evaluation Protocol

We analyze the above three loss formulations and their respective comparison with neural networks on four commonly investigated image datasets: MNIST (LeCun et al., 1998), CIFAR-10 (Krizhevsky and Hinton, 2009), SVHN (Netzer et al., 2011), and CelebA (Liu et al., 2015). For each dataset we will use their default data splits for train, test, and validation sets and employ preprocessing, i.e. ensuring zero mean and unit variance for each feature across the dataset. Furthermore, for each proposed experiment we will also compare with a baseline SPN, optimized purely using maximum likelihood. As SPNs and deep neural networks deviate in their model structure, we practice particular caution to warrant a fair comparison that provides conclusive evidence with respect to our postulated hypothesis. For each experiment, we will perform five randomly seeded runs to measure statistical deviations. Specifically:

Optimization We use Adam (Kingma and Ba, 2015) as the optimization algorithm in all models. Since the aforementioned loss objectives are novel territory for SPNs, we conduct a grid-search over a discrete set of learning rates and the stochasticity induced by mini-batch size using a separate validation set. The same process is conducted for the neural counterparts. A set of hyper-parameters is then selected individually for each model to assure an adequate operating point.

Metrics For evaluation purposes we primarily employ the predominant Fréchet inception distance (FID). However, this metric is known to suffer from minor perturbation and resolution changes (Borji, 2019). Therefore, we also analyze the generative power of each model by training a separate classifier on the corresponding generations, and we assess its classification performance on the original test sets. For the auto-encoding setting, we will further report the reconstruction losses.

Model Capacity To keep the comparison between original models and the SPN variants fair, we employ the same amount of parameters for each model in direct comparisons.

Data Efficiency and Convergence Since deep neural networks and PCs have fundamentally different structures, they may have different data efficiency and convergence rate. For this purpose, we evaluate the aforementioned losses and metrics under the lens of different amounts of training instances available for optimization. In addition, to further underline the fairness w.r.t. the arguments in the previous points on optimization and metrics, we also test all models at different points over the course of training, in order to assess generation quality at potentially different convergence speeds.

Ablation Since the Gumbel-Softmax Trick exposes an additional hyper-parameter, i.e. the temperature τ , we will perform an ablation study between constant, annealed, and interactively learned temperature values. We perform this study on CIFAR-10 for each of the investigations in Section 4 with a fixed learning rate and mini-batch size, to keep the number of experiments practical.

5. Results

In the following, we report the results of SPNs trained and evaluated on MNIST (28×28), SVHN (32×32), CIFAR (32×32), and CelebA (64×64). We abbreviate the different optimization objectives as follows: log-likelihood maximization as LL, adversarial training as ADV, maximum mean discrepancy training as MMD, and probabilistic auto-encoding training as PAE. For each experiment, we use a random binary tree structure for the SPN with depth $D = 4$, number of repetitions $R = 10$, number of sum node and leaf node representations per random variable at each layer $K = 20$ and Binomial leaves with $N = 255$. To sample differentiably from Binomials, we approximate the Binomial distribution $B(N, p)$ with a Gaussian distribution $\mathcal{N}(\mu, \sigma)$ with $\mu = Np$ and $\sigma = Np(1 - p)$ and apply the well-known reparametrization trick.

To understand the generative capabilities of each model, we visualize samples of all datasets, drawn from the SPN optimized with the respective objective in Table 1. We further evaluate the FID and KID scores on the CelebA test split against 10,000 samples from each objective in Table 2. For the full grid-search results over learning-rate $\in \{0.1, 0.01, 0.001\}$ and batch-size $\in \{64, 128, 256\}$ between all datasets and all objectives, we refer the reader to Appendix A.

Adversarial Training The adversarial objective leads to samples of higher degrees of saturation and contrast, compared to LL-based samples. Even though on CelebA (Table 1, second row, fourth column) some faces are sampled twice, the general diversity between faces is much higher than those sampled with the LL optimized SPN. However, FID and KID scores are higher for the ADV objective (FID: 420.6 ± 24.3 , KID: 0.55 ± 0.06) compared to the LL objective (FID: 273.4 ± 13.9 , KID: 0.24 ± 0.02). On SVHN and CIFAR (Table 1, second row, second and third column), the LL-based SPN generates washed-out samples, whereas the ADV objective can generate samples that are richer in terms of features and colors. On MNIST (Table 1, second row, first column), the ADV setting is prone to mode collapse (see also different hyper-parameter combinations in Figure 1 of Appendix A).

Maximum Mean Discrepancy The SPN optimized with the MMD objective can generate samples with higher brightness on MNIST (Table 1, third row, first column), showing clearer boundaries between the white digit pixels and the black background, compared to the LL objective. Similarly, whereas facial details on CelebA (Table 1, third row, fourth column) are not as detailed as with other objectives, we can observe a richer variation in colors than the more average-looking ones from the LL or the PAE objective. On CelebA, the MMD objective achieves an FID score of 353.5 ± 5.7 and a KID score of 0.43 ± 0.02 . Unfortunately, we were not able to reproduce any stable results of MMD-GAN Li et al. (2017) with the official implementation.

Probabilistic Auto-Encoding The PAE objective can generate distinguishable digits on MNIST (Table 1, fourth row, first column) and approach the sample quality of CelebA faces compared to the LL objective with slightly fewer variations (Table 1, fourth row, fourth column). The FID and KID scores of the PAE objective come close to that of the vanilla VAE on CelebA: 354.2 ± 19.2 (PAE) compared to 316.3 ± 0.9 (VAE) for FID and 0.35 ± 0.02 (PAE) compared to 0.33 ± 0.1 (VAE). The probabilistic auto-encoding setting was further evaluated by measuring the reconstruction loss (as binary cross-entropy between

	MNIST	SVHN	CIFAR	CelebA
LL				
ADV				
MMD				
PAE				

Table 1: SPN samples on four different datasets after training with the maximum likelihood (LL), adversarial training (ADV), maximum mean discrepancy (MMD), and probabilistic auto-encoding (PAE) objectives.

FID	SPN	DNN	KID	SPN	DNN
LL	273.4 ± 13.9	–	LL	0.24 ± 0.02	–
ADV	420.6 ± 24.3	358.1 ± 1.3	ADV	0.55 ± 0.06	0.39 ± 0.01
MMD	353.5 ± 5.7	†	MMD	0.43 ± 0.02	†
PAE	354.2 ± 19.2	316.3 ± 0.9	PAE	0.35 ± 0.02	0.33 ± 0.01

Table 2: FID and KID scores on CelebA (lower is better). DNN indicates a deep neural model. The likelihood objective is only available to SPNs. We were not able to reproduce stable MMD-GAN results (†) with the official implementation (Li et al., 2017).

Model	MNIST	SVHN	CIFAR	CelebA
SPN	606.1 ± 44.2	2644.6 ± 423.7	4006.3 ± 1180.4	10945.4 ± 1493.0
VAE	75.3 ± 0.1	1817.9 ± 0.6	1807.7 ± 0.1	6352.1 ± 0.5

Table 3: Auto-encoding reconstruction error measured with the binary cross-entropy between pixels as sum over all pixels and mean over all images on the test splits of the respective datasets for SPNs trained in the PAE setting and vanilla VAEs.

pixel values in the range of $[0, 1]$) over all images in the test split of the respective datasets in Table 3 and compared to vanilla VAE models. Whereas on MNIST, the VAE reconstruction error (75.3 ± 0.1) is significantly lower than the SPN reconstruction error (606.1 ± 44.2), on SVHN (1817.9 ± 0.6 vs. 2644.6 ± 423.7), CIFAR (1807.7 ± 0.1 vs. 4006.3 ± 1180.4), and CelebA (6352.1 ± 0.5 vs. 10945.4 ± 1493.0) the SPN reconstruction is in the same order of magnitude, compared to the VAE. On all datasets, the SPN reconstruction error shows a larger standard deviation than the VAE across five differently seeded runs of the same setting.

Gumbel-Softmax Temperature Ablation Study The continuous approximation of the categorical sampling in sum nodes using the Gumbel-Softmax Trick introduces a temperature parameter τ that controls the softmax entropy. The choice of τ is a trade-off during parameter learning, where larger values of τ leads to smoother interpolations between child samples with low gradient variance and lower values of τ result in almost one-hot encoded child samples and high gradient variance. We evaluate the influence of the Gumbel-Softmax Temperature for all objectives with three schedules: A *constant* $\tau \in \{0.05, 0.5, 1.0\}$, an *annealed* $\tau = \max(0.5, \exp(-\frac{\text{current_epoch}}{\text{max_epochs}}))$ that starts at $\tau = 1.0$ and decays until $\tau = 0.5$, and a *learned* τ that is jointly optimized with the model parameters during training.

We train an SPN with each objective with each schedule for τ on CelebA and visualize their samples in Table 4. For a constant value of $\tau = 0.05$ that approximates a high samples uniqueness (i.e. less interpolation between child samples and closer to the argmax

τ	ADV	MMD	PAE
0.05			
0.5			
1.0			
annealed			
learned			

Table 4: Image reconstruction error of an SPN trained with the GAN, MMD, and PAE objective with different settings for τ . Values of 0.05, 0.5, and 1.0 provide a constant schedule, the “annealed” schedule anneals the value with $\tau = \max(0.5, \exp(-\frac{\text{current_epoch}}{\text{max_epochs}}))$, proposed in [Jang et al. \(2017\)](#), and “learned” jointly learns τ with the model parameters during the training.

operation), all objectives generate noisy samples. When balancing the interpolation and the uniqueness of child samples with $\tau = 0.5$, the different objectives produce colorful samples with a higher variance. For $\tau = 1.0$, the PAE objective visually improves over $\tau = 0.5$ slightly, which could be attributed to the fact, that the PAE objective leads to more average looking samples in most settings (see also hyper-parameter grid-search results in Appendix A). For the MMD objective, the constant $\tau = 1.0$ setting generates more washed-out samples with less saturation. Interestingly, the annealed and the learned schedule result in similar sample fidelity for all objectives. The ADV objective appears to generate stable results for larger τ , independent of constant, annealed or learned schedules.

5.1. Results on Synthetic Data

To further understand the previous results, we now analyze the model distributions learned with the ADV and MMD objectives on well known synthetic datasets and compare it to the standard SPN setting with the LL objective as an addition to the initial experimental protocol outlined in Section 4. Note, that we exclude the PAE setting from these experiments, as an encoding from the two to one-dimensional space will, in the optimal case, only be a reduction to the axis with the highest variance from which we will not be able to recover a meaningful two-dimensional representation anymore. The following experiments are performed with an SPN that consists of a single sum root node and K product child nodes that themselves have two Gaussian leaf nodes, one for each dimension. This is equivalent to a bi-variate Gaussian Mixture Model with isotropic Gaussians and K components. We construct the following six synthetic two-dimensional datasets: *2-clusters*, a dataset with two Gaussian clusters with isotropic covariance ($K = 2$); *varied*, a datasets with three Gaussian clusters with isotropic covariance but varying standard deviations ($K = 3$); *aniso*, a dataset with three Gaussian with full covariance ($K = 30$); *9-clusters*, a dataset with nine equal Gaussian isotropic clusters arranged in a 3×3 grid ($K = 9$); *2-moons*, a dataset consisting of two half-moon shaped clusters that are slightly intertwined ($K = 50$); *2-circles*, a dataset consisting of an inner and outer circle with different radii ($K = 30$).

We visualize the SPN model distribution of the LL, ADV, and MMD objectives learned on each of the six datasets in Table 5. All objectives can capture the data distribution of *2-clusters*. Increasing the data complexity by intersecting clusters and changing their variance in *varied* leads to minor deviations from the true data distributions for the ADV objective, whereas the MMD objective seems not to put not enough weight on the third cluster with the highest variance. On *aniso*, the ADV objective can model the three cluster modes but is less precise at approximating the non-isotropic cluster shape with multiple smaller Gaussians. The MMD objective does not converge on any mode of *aniso*. Whereas the LL objective can match each cluster in *9-clusters* with one child (although one cluster is weighted too high), the ADV objective distribution only covers some of the outer clusters and the MMD objective cannot distinguish between the single clusters and rather distributes most of its density across the full grid. For the *2-moons* and *2-circles* the LL objective can approximate the non-linear data distribution by enumeration (like filling an arbitrary surface with small circles), whereas the ADV objective mostly collapses on the upper half-moon in *2-moons* and the outer circle in *2-circles*. Similar to the *9-clusters* results, the MMD objective is

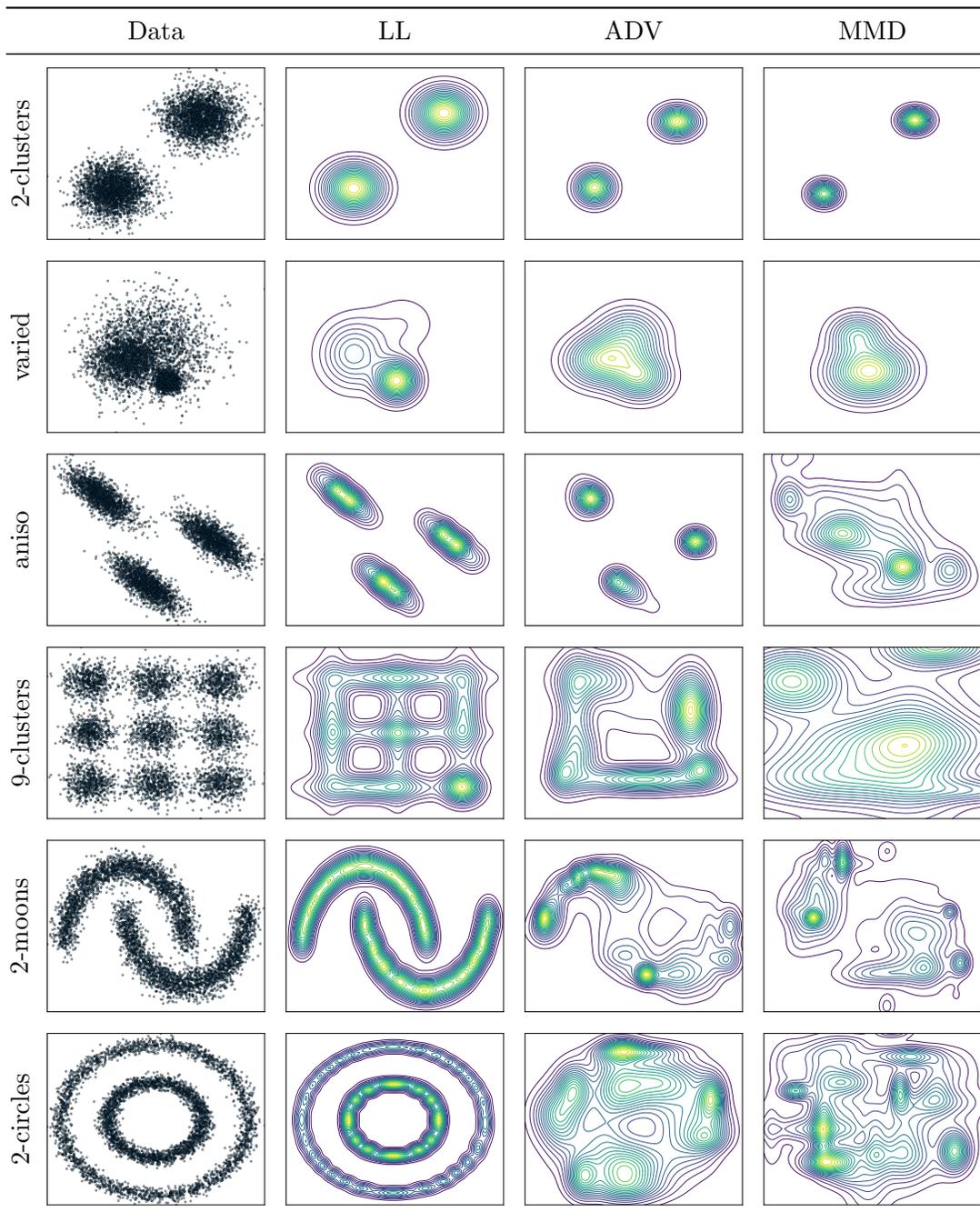


Table 5: SPN distributions on different synthetic two-dimensional datasets with varying difficulty, learned by optimizing the LL, ADV, and MMD objectives. We adapt the number of components, modeled by the SPN, as follows: 2-clusters $K = 2$, varied $K = 3$, aniso $K = 30$, 9-clusters $K = 9$, 2-moons $K = 50$, circles $K = 30$.

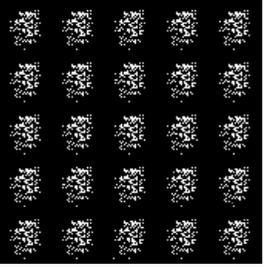
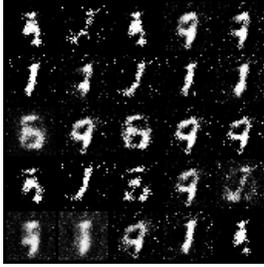
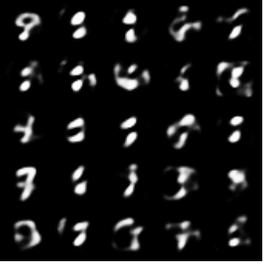
	Stable	Unstable	Mean/Mode Collapse
GAN			
VAE			

Table 6: Stable and failure modes of the vanilla GAN and VAE models. The GAN model can become unstable when the number of steps between discriminator updates is too high and falls into mode collapse when the stochasticity induced by the batch-size is too low. The VAE model can become unstable when the reconstruction error is weighted too high or run into mean collapse issues when the KL-divergence is weighted too high.

not able to distinguish between the inner and outer circle in *2-circles* and puts most of its density only to particular subareas of *2-moons*.

6. Findings

Our results in Section 5 have shown that differentiable sampling in PCs is indeed hard. For now, the initial research hypothesis, that perceptual losses, unlocked by our novel differentiable sampling procedure, will elevate the generative power of PCs and improve their sample quality, could not be fully confirmed. Whereas we were unable to elevate the perceptual quality just yet, the results demonstrated throughout our experiments showcase that the generated samples can have similar sample quality compared to likelihood-based optimized PCs, while at the same time expressing richer contrast, colors, and details. Therefore, it is important to highlight that our results showcase that it is possible to optimize PCs with differentiable sampling and thus allow the use of arbitrary loss functions. Whereas prior to this work, PCs were restricted to likelihood-based optimization, we have now paved the way for PCs to be opened up to a variety of loss formulations that have been built around deep neural networks in recent years.

During the experimental phase, we further observed situations that actually closely relate to the failure modes of the vanilla versions of GANs and VAEs. Training procedures in the ADV setting sometimes either do not converge or quickly collapse such that the generator does not receive any meaningful signal anymore. Although these issues are mostly addressed in modern formulations of these models, it is easy to reproduce those failure modes in their initial versions, as shown in Figure 6, where the GAN model can become unstable when the number of steps between discriminator updates is too high and falls into mode collapse when the stochasticity induced by the batch-size is too low. Similarly, the VAE model can become unstable when the reconstruction error is weighted too high or runs into mean collapse issues when the KL-divergence is weighted too high. Hence, modern GAN and VAE formulations often now work with the help of further less obvious advances and many finicky implementation details, that seem to positively influence the optimization stability. To a similar degree, we found late in the experimental phase, that replacing Gaussian leaves with the Binomial leaves and a Gaussian approximation for differentiable sampling as described in Section 5 vastly improved our results. Consequently, we acknowledge having perhaps underestimated the necessary engineering efforts that have gone into GANs and VAEs in recent years. However, we are cautiously optimistic that similar tricks can be found or adapted from deep neural networks to improve PC optimization with objectives based on differentiable sampling and as a result further improve perceived sample quality.

We additionally want to highlight, that the FID and KID evaluations performed on the CelebA dataset do not reflect the subjective sample quality of Figure 1. We suspect the problem to be a combination of the well-known issues with FID (Borji, 2019) and its susceptibility to a range of factors such as image artifacts, compression rates, and re-scaling algorithms, as well as the fact that the produced samples in all cases are generally not of a quality level at which an FID and KID comparison is meaningful. We therefore raise caution of making comparisons between results based on the reported FID and KID results.

Future Work Whereas the results show that differentiable sampling in PCs is harder than expected, the main takeaway message of this work is that PCs seem to be at a stage where deep neural networks have been in their early days. Our experiments have shown, that PCs with the corresponding objectives sometimes fall into similar pitfalls as their vanilla deep neural counterparts. As was necessary for deep neural networks, we also expect that for each objective a specific and more in-detail investigation is necessary and will result in new optimization settings, loss combinations, and tricks that will help PCs to be more stable and lead to further improvements. Therefore, we suggest further examining how the investigated objectives behave when combined with pre-training the PC via likelihood maximization as a good initialization, as well as training the respective objective jointly with likelihood maximization, similar to how VAEs incorporate both, the reconstruction and the KLD objective. Furthermore, the fact that our finding of a Gaussian approximation of Binomial leaves improves over direct Gaussian leaves highlights the fact, that the choice of leaves during differentiable sampling plays an important role and requires further inspection.

7. Documented Modifications

During our experimental phase, we found that a small subset of the initially planned evaluations were not appropriate or did not contribute additional insight anymore. In the

following, we enumerate evaluation protocol modifications and provide our reasoning for the decision:

Updated abstract We adapted the abstract to reflect the additional insights that were gained after performing the experimental phase.

Added KID score and focus FID and KID on CelebA We further added the KID score as suggested in Li et al. (2017) and focused the FID and KID evaluation on CelebA.

Included synthetic dataset evaluations To further understand the objective optimization behavior, we decided to additionally investigate the objectives in synthetic two-dimensional data settings that are well understood.

Focus on neural discriminator Whereas the ADV objective with a neural discriminator turned out to be harder than expected, we found no setting in which a discriminative SPN as discriminator leads to stable training. Therefore, we decided that this setting needs further in-depth investigation on its own.

Classifier, data efficiency, and convergence speed analysis Given the current room for improvement at the existing scale, we argue that it is somewhat futile to investigate scenarios with even less data or fewer training steps for SPNs. However, during our experiments, we also realized that this is not in disagreement with the conjecture that SPN may or may not be more efficient in direct comparison with their neural counterparts for the problems at hand. In fact, as specified earlier, the encountered challenges seem to be shared with the original VAE and GAN formulations, suggesting that the complexity and data efficiency analysis should be postponed to even further improved implementations, similar to the very recent advances in the respective deep neural generators. In that sense, our question remains open for future analysis, to be revisited in scenarios that have achieved a larger scale first.

Acknowledgments

This work has been supported by the project “safeFBDC - Financial Big Data Cluster” (FKZ: 01MK21002K), funded by the German Federal Ministry for Economic Affairs and Energy as part of the GAIA-x initiative, the Hessian Ministry of Higher Education, Research, Science and the Arts (HMWK) project “The Third Wave of AI”, and the Federal Ministry of Education and Research (BMBF; Competence Center for AI and Labour; “kompAKI”, FKZ02L19C150).

References

- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *arXiv preprint, arXiv:1701.07875*, 2017.
- Yoshua Bengio, Li Yao, Guillaume Alain, and Pascal Vincent. Generalized denoising auto-encoders as generative models. *Conference on Neural Information Processing Systems (NeurIPS)*, 2013.
- Ali Borji. Pros and cons of gan evaluation measures. *Computer Vision and Image Understanding*, 179:41–65, 2019.

- Xi Chen, Diederik P. Kingma, Tim Salimans, Yan Duan, Prafulla Dhariwal, John Schulman, Ilya Sutskever, and Pieter Abbeel. Variational Lossy Autoencoder. *International Conference on Learning Representations (ICLR)*, 2017.
- Adnan Darwiche. A logical approach to factoring belief networks. *Knowledge Representation and Reasoning (KR)*, 2:409–420, 2002.
- Adnan Darwiche. A differential approach to inference in Bayesian networks. *Journal of the ACM*, 50(3):280–305, 2003.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. *arXiv preprint, arXiv:2105.05233*, 2021.
- Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. A survey of data augmentation approaches for NLP. In *Association for Computational Linguistics and the International Joint Conference on Natural Language Processing (ACL-IJCNLP)*, pages 968–988, 2021.
- Robert Gens and Pedro Domingos. Discriminative learning of sum-product networks. In *Conference on Neural Information Processing Systems (NeurIPS)*, volume 25, 2012.
- I. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2014.
- Alex Graves. Practical variational inference for neural networks. *Conference on Neural Information Processing Systems (NeurIPS)*, 2011.
- Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research (JMLR)*, 13:723–773, 2012.
- Ishaan Gulrajani, Kundan Kumar, Ahmed Faruk, Adrien Ali Taiga, Francesco Visin, David Vazquez, and Aaron Courville. PixelVAE: a Latent Variable Model for Natural Images. *International Conference on Learning Representations (ICLR)*, 2017.
- Geoffrey E. Hinton, Peter Dayan, Brendan J. Frey, and Radford M. Neal. The ”wake-sleep” algorithm for unsupervised neural networks. *Science*, 268 5214:1158–61, 1995.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint, arXiv:2006.11239*, 2020.
- Poon Hoifung and Domingos Pedro. Sum-product networks: A new deep architecture. In *Uncertainty in Artificial Intelligence (UAI)*, 2011.
- Huaibo Huang, Zhihang Li, Ran He, Zhenan Sun, and Tieniu Tan. Introvae: Introspective variational autoencoders for photographic image synthesis. *Conference on Neural Information Processing Systems (NeurIPS)*, 2018.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint, arXiv:1611.01144*, 2017.

- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. *International Conference on Learning Representations (ICLR)*, 2013.
- Diederik P. Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *arXiv preprint, arXiv:2107.00630*, 2021.
- Doga Kisa, Guy Van den Broeck, Arthur Choi, and Adnan Darwiche. Probabilistic Sentential Decision Diagrams. In *Knowledge Representation and Reasoning (KR)*, 2014.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, Toronto, Ontario, 2009.
- Anders Boesen Lindbo Larsen, Soren Kaae Sonderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. *International Conference on Machine Learning (ICML)*, 2016.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *IEEE*, 86(11):2278–2324, 1998.
- Chun-Liang Li, Wei-Cheng Chang, Yu Cheng, Yiming Yang, and Barnabás Póczos. Mmd gan: Towards deeper understanding of moment matching network. *arXiv preprint, arXiv:1705.08584*, 2017.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *International Conference on Computer Vision (ICCV)*, 2015.
- Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint, arXiv:1611.00712*, 2017.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *Conference on Neural Information Processing Systems (NeurIPS) Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- Robert Peharz, Steven Lang, Antonio Vergari, Karl Stelzner, Alejandro Molina, Martin Trapp, Guy Van Den Broeck, Kristian Kersting, and Zoubin Ghahramani. Einsum networks: Fast and scalable learning of tractable probabilistic circuits. In *International Conference on Machine Learning (ICML)*, 2020.
- Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz Kaiser, and Geoffrey Hinton. Regularizing neural networks by penalizing confident output distributions. *arXiv preprint, arXiv:1701.06548*, 2017.
- Andrzej Pronobis, Avinash Ranganath, and Rajesh P. N. Rao. LibSPN: A library for learning and inference with Sum-Product Networks and TensorFlow. In *International Conference on Machine Learning (ICML) Workshop on Principled Approaches to Deep Learning*, 2017.

- Tahrima Rahman, Prasanna V. Kothalkar, and Vibhav Gogate. Cutset networks: A simple, tractable, and scalable approach for improving the accuracy of chow-liu trees. In *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKKD)*, 2014.
- Nimish Shah, Laura I. Galindez Olascoaga, Wannes Meert, and Marian Verhelst. Acceleration of probabilistic reasoning through custom processor architecture. *Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2020.
- Connor Shorten and Taghi M. Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6:1–48, 2019.
- Lukas Sommer, Michael Halkenhäuser, Cristian Axenie, and Andreas Koch. Spnc: Accelerating sum-product network inference on cpus and gpus. In *IEEE 32nd International Conference on Application-specific Systems, Architectures and Processors (ASAP)*, pages 53–56, 2021.
- Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint, arXiv:2011.13456*, 2021.
- Christian Szegedy, V. Vanhoucke, S. Ioffe, Jonathon Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schölkopf. Wasserstein auto-encoders. *International Conference on Learning Representations (ICLR)*, 2018.
- George Tucker, Andriy Mnih, Chris J Maddison, John Lawson, and Jascha Sohl-Dickstein. Rebar: Low-variance, unbiased gradient estimates for discrete latent variable models. In *Conference on Neural Information Processing Systems (NeurIPS)*, volume 30. Curran Associates, Inc., 2017.
- Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. It Takes (Only) Two : Adversarial Generator-Encoder Networks. *Association for the Advancement of Artificial Intelligence (AAAI)*, 2018.
- Antonio Vergari, YooJung Choi, Robert Peharz, and Guy Van den Broeck. Probabilistic circuits: Representations, inference, learning and applications, 2020. Tutorial at AAAI 2020.
- Qingsong Wen, Liang Sun, Xiaomin Song, Jing Gao, Xue Wang, and Huan Xu. Time series data augmentation for deep learning: A survey. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2021.
- Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 1992.
- Zijun Zhang, Ruixiang Zhang, Zongpeng Li, Yoshua Bengio, and Liam Paull. Perceptual generative autoencoders. *International Conference on Machine Learning (ICML)*, 2020.

Appendix A. Hyper-Parameter Grid-Search Results

In addition to the reported results and analysis of the different objectives in the main manuscript, we conduct a grid-search over the learning rates (lr) and the stochasticity induced by mini-batch sizes (bs) using the validation set for all datasets and each objective respectively. Specifically, we train an SPN with binary-tree structure and $K = 20$, $R = 10$, $D = 4$, Binomial leaves, the annealed Gumbel-Softmax temperature schedule with the Adam optimizer for 30 epochs schedule with each objective. We visualize the sample quality for all combinations of $lr \in \{0.001, 0.010, 0.100\}$ and $bs \in \{64, 128, 256\}$ for MNIST in Table 7, for SVHN in Table 8, for CIFAR in Table 9, and for CelebA in Table 10. The ADV and MMD setting further expose a separate learning rate hyper-parameter lr_D for the discriminator (ADV) and the adversarially learned kernel encoder (MMD). Whereas it was sufficient to adopt the default learning rate of the vanilla GAN for the ADV discriminator ($lr_D = 0.0002$), the MMD kernel encoder learning rate was crucial to be adapted for each dataset. Therefore, we have extended the MMD grid-search by $lr_D \in \{0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001, 0.00005, 0.00001\}$ and found the following settings to be optimal: MNIST, $lr_D = 0.01$; SVHN, $lr_D = 0.005$; CIFAR, $lr_D = 0.005$; CelebA, $lr_D = 0.001$.

lr	bs	LL	ADV	MMD	PAE
0.001	64				
0.010	64				
0.100	64				
0.001	128				
0.010	128				
0.100	128				
0.001	256				
0.010	256				
0.100	256				

Table 7: Samples of an SPN trained on MNIST using the LL, ADV, MMD, and PAE objectives with a grid-search over learning rate (lr) and batch-size (bs) settings.

lr	bs	LL	ADV	MMD	PAE
0.001	64				
0.010	64				
0.100	64				
0.001	128				
0.010	128				
0.100	128				
0.001	256				
0.010	256				
0.100	256				

Table 8: Samples of an SPN trained on SVHN using the LL, ADV, MMD, and PAE objectives with a grid-search over learning rate (lr) and batch-size (bs) settings.

lr	bs	LL	ADV	MMD	PAE
0.001	64				
0.010	64				
0.100	64				
0.001	128				
0.010	128				
0.100	128				
0.001	256				
0.010	256				
0.100	256				

Table 9: Samples of an SPN trained on CIFAR using the LL, ADV, MMD, and PAE objectives with a grid-search over learning rate (lr) and batch-size (bs) settings.

lr	bs	LL	ADV	MMD	PAE
0.001	64				
0.010	64				
0.100	64				
0.001	128				
0.010	128				
0.100	128				
0.001	256				
0.010	256				
0.100	256				

Table 10: Samples of an SPN trained on CelebA using the LL, ADV, MMD, and PAE objectives with a grid-search over learning rate (lr) and batch-size (bs) settings.