

Probabilistic Circuits That Know What They Don't Know



Fabrizio Ventola^{*1}



Steven Braun^{*1}



Zhongjie Yu¹



Martin Mundt^{1,2}



Kristian Kersting^{1,2,3,4}



TECHNISCHE
UNIVERSITÄT
DARMSTADT

¹Department of Computer Science, TU Darmstadt

²Hessian Center for AI (hessian.AI)

³German Research Center for Artificial Intelligence (DFKI)

⁴Centre for Cognitive Science, TU Darmstadt



hessian.AI

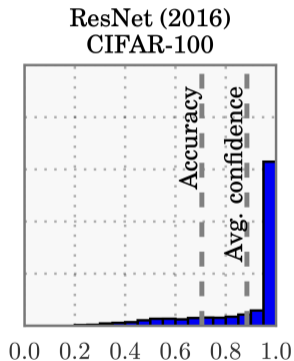
KompA+KI



* indicates equal contribution

Motivation

Many ML models are **overconfident** – often assign high probabilities to (wrong) predictions, even for unseen categories.

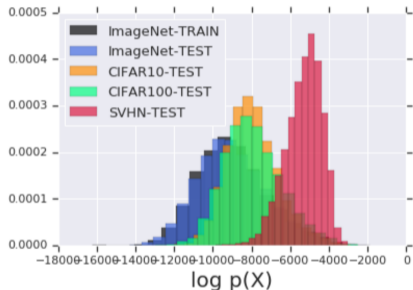


[Amodei et al., arXiv 2016; Guo et al., ICML 2017; Boult et al., AAI 2019; Hendrycks et al., ICLR 2019]

Motivation

- The issue impacts even major deep models like VAEs and normalizing flows.

[Nalisnick et al., ICLR 2019]

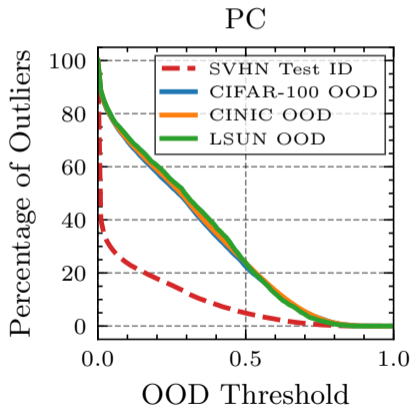


- **Probabilistic circuits (PCs)** are assumed to overcome this — casted as well-calibrated models.

[Peharz et al., UAI 2020; Peharz et al., ICML 2020; Choi et al., UCLA 2020]

Our Contributions (1/2)

Show that PCs suffer from overconfidence and struggle to distinguish **in-distribution from out-of-distribution data** in discriminative setting.



Our Contributions (2/2)

- Introduce **Tractable Dropout Inference (TDI)** – provides tractable model uncertainty estimation.
- Derive a **sampling-free analytical solution for Monte Carlo dropout (MCD)**, a Bayesian method for model uncertainty in neural networks (NNs).
- Demonstrate **TDI's robustness** against distribution shifts and out-of-distribution data in three key scenarios.

Probabilistic Circuits

- We focus on sum-product networks (SPNs), a prominent type of PCs.
- SPNs are known for inference capabilities and representational power.
- Like NNs, SPNs are deep graphs but they explicitly encode a normalized probability distribution.
- SPNs can generate new samples and calculate various exact, tractable probabilistic queries, even with partial evidence.

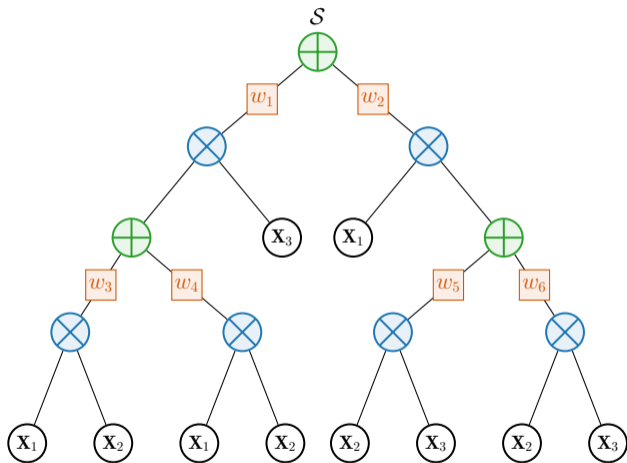
[Poon et al., UAI 2011; Delalleau et al., NeurIPS 2011; Peharz et al., AISTATS 2015]

Probabilistic Circuits

Clear probabilistic semantics!

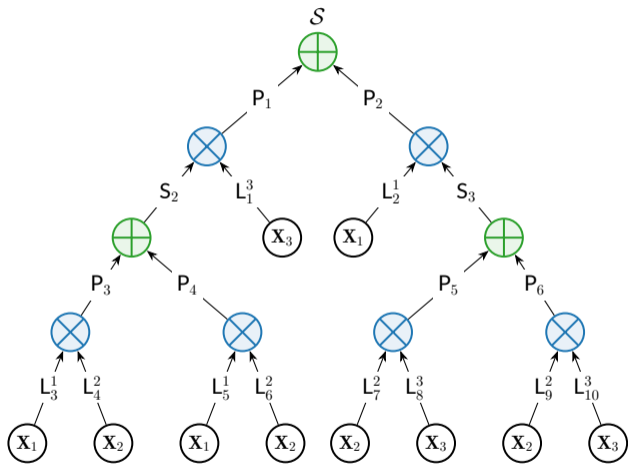
\oplus $S = \sum_i w_i N_i$

\otimes $P = \prod_i N_i$



Probabilistic Circuits

- Instance LL via forward pass
- $\mathcal{S}(\mathbf{x}) = p_{\mathbf{X}}(\mathbf{X} = \mathbf{x})$
- MAP & MPE need top-down pass
- Inference linear in the network size



Remark

- Probabilistic modeling targets uncertainties, but it's crucial to quantify model uncertainty or *epistemic uncertainty*.
- Quantification can tell us how confident the model is in its predictions, including class label distribution $p(Y|\mathbf{X} = \mathbf{x})$.
- Overconfident models can assign near 1 probabilities to out-of-distribution instances or unseen labels.
- Model uncertainty helps us know when the model “*does not know*”, need to take predictions with caution.

[Kendall et al., NeurIPS, 2017; Hüllermeier et al., MLJ 2021]

Model Uncertainty Quantification in Deep Networks

- A common Bayesian way to quantify model uncertainty involves choosing the most likely parameters configuration θ that represents the data \mathcal{D} .

$$\arg \max_{\theta} p(\theta|\mathcal{D}) \quad \text{where} \quad p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

- This approach is usually intractable as it involves integrating over parameters θ .
- Dropout is a popular method in NNs to prevent overfitting and improve generalization by randomly removing connections between layers.

[MacKay, Neural computation 1992; Srivastava et al., JMLR 2014]

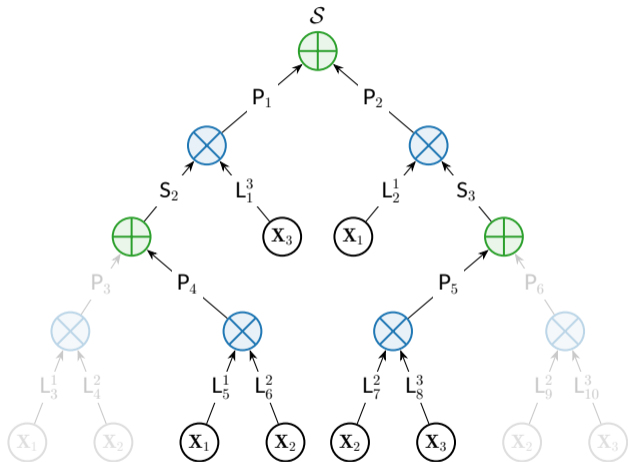
Monte Carlo Dropout

- Gal et al. (ICML 2016) interpret dropout as a Bayesian approximation for model uncertainty, assuming Bernoulli distribution on weights.
- The intractable integration over parameters is approximated by a tractable sum over n parameter sets: $\theta_i \sim p(\theta | \mathcal{D})$.
- With n predictions, first (mean) and second (variance) raw moments serve as prediction and model uncertainty.

[Gal et al., *Dropout as a bayesian approximation: Representing model uncertainty in deep learning*, ICML 2016]

MCD in PCs

- Place Bernoulli distribution at sum node weights
- Calculating two raw moments **requires n evals!**



Can we do better?

MCD needs n stochastic forward-passes of the model ... → **inefficient!**

$\text{Var}[\mathcal{S}]$ is induced by random instantiations of dropout at sum nodes.

Can we find a *closed-form* solution to obtain $\text{Var}[\mathcal{S}]$ in a **single** forward-pass?

Yes! ⇒ **Tractable Dropout Inference (TDI)**

Tractable Dropout Inference (TDI)

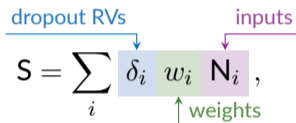
Key idea:

- 1) View node output as function f of inputs and dropout RVs
- 2) Compute $\text{Var}[f] \rightarrow$ **decomposes** into input variances (and more)
- 3) **Propagate** $\text{Var}[f]$ from leaf nodes to root node

TDI only needs a **single** forward evaluation of a model ... \rightarrow **efficient!**

TDI: Basic Idea

Consider sum nodes as linear combinations of their Bernoulli dropout RVs and inputs:

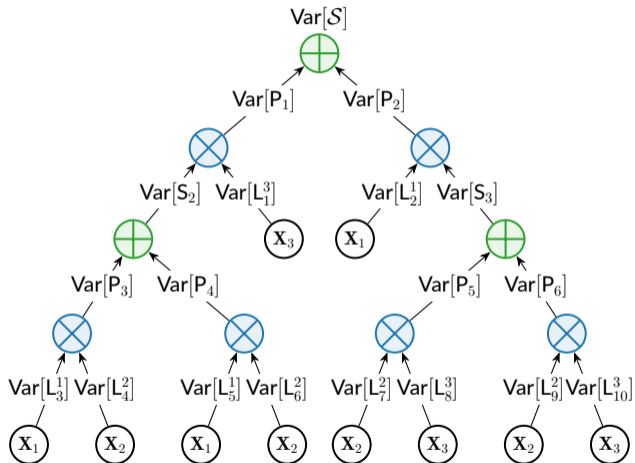
$$S = \sum_i \delta_i w_i N_i,$$


where $\delta_i \sim \text{Bern}(q)$ and $p = 1 - q$ is the dropout probability.

How do we *tractably* compute the variance of a PC with dropout?

⇒ Find a closed-form expression of $\text{Var}[S]$ and perform **variance propagation!**

TDI: Variance Propagation



Closed-form Solutions

Recall: \oplus $S = \sum_i \delta_i w_i N_i$ \otimes $P = \prod_i N_i$

Variance \oplus $\text{Var}[S] = q \sum_i w_i^2 (\text{Var}[N_i] + p \mathbb{E}[N_i]^2) + q^2 \sum_{i \neq j} w_i w_j \text{Cov}[N_i, N_j]$

input variance input expectation input covariance

\otimes $\text{Var}[P] = \prod_i (\text{Var}[N_i] + \mathbb{E}[N_i]^2) - \prod_i \mathbb{E}[N_i]^2$

Expectation \oplus $\mathbb{E}[S] = q \sum_i w_i \mathbb{E}[N_i]$

\otimes $\mathbb{E}[P] = \prod_i \mathbb{E}[N_i]$

Covariance \oplus $\text{Cov}[S^A, S^B] = q^2 \sum_i w_i^A \sum_j w_j^B \text{Cov}[N_i^A, N_j^B]$

\otimes $\text{Cov}[P^A, P^B] = \mathbb{E}\left[\prod_i N_i^A \prod_j N_j^B\right] - \prod_i \mathbb{E}[N_i^A] \prod_j \mathbb{E}[N_j^B]$

‡ not decomposable in general ‡ ↑

Covariance: Three Solutions

⚡ not decomposable in general ⚡

⊗
$$\text{Cov}[\mathbf{P}^A, \mathbf{P}^B] = \mathbb{E}\left[\prod_i \mathbf{N}_i^A \prod_j \mathbf{N}_j^B\right] - \prod_i \mathbb{E}[\mathbf{N}_i^A] \prod_j \mathbb{E}[\mathbf{N}_j^B]$$

We're not at a dead-end with covariance! We provide three solutions:

- a) **Structural Knowledge:** LearnSPN, RAT-SPN, ... → add. independencies
- b) **Covariance Bounds:** Cauchy-Schwarz inequality → lower/upper bound

$$\text{Cov}[\mathbf{N}_i, \mathbf{N}_j] \in \left[-\sqrt{\text{Var}[\mathbf{N}_i] \text{Var}[\mathbf{N}_j]}, +\sqrt{\text{Var}[\mathbf{N}_i] \text{Var}[\mathbf{N}_j]} \right]$$

- c) **Copy-paste Solution:** Graph augmentation → enforce covariance to be zero

Leaf Nodes

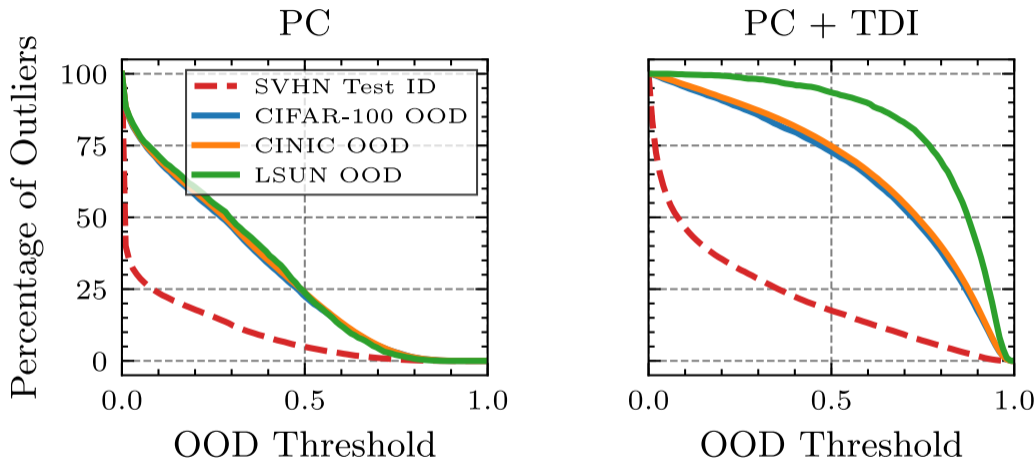
No dropout (in our current framework) \rightarrow leaves become a point estimate!

$$\mathbb{E}[L] = L, \quad \text{Var}[L] = 0, \quad \text{Cov}[L_i, L_j] = 0$$

BUT: This allows to include **prior knowledge** about *aleatoric* and *epistemic* uncertainty by setting $\text{Var}[L] > 0$ and $\text{Cov}[L_i, L_j] \neq 0$ (future work).

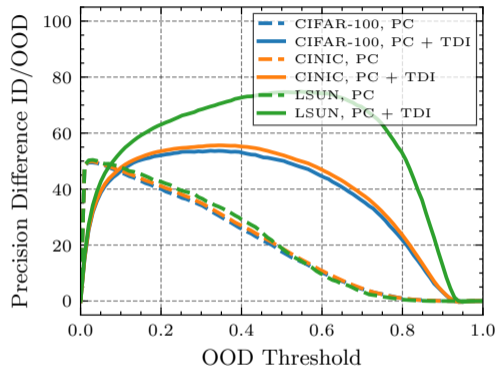
(Note: MCD does *not* allow that)

Experiments: Out-of-Distribution Detection



TDI remarkably improves OOD detection in PCs!

Experiments: Out-of-Distribution Detection

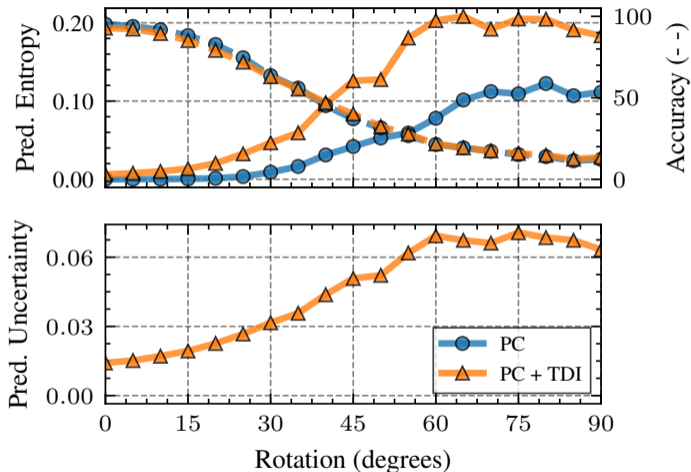


AUC (\uparrow)	CIFAR	CINIC	LSUN
PC	29.3	29.9	30.3
PC + TDI	64.6	66.1	81.8
PC + MCD	68.5	70.0	84.9

TDI is competitive with MCD and **only needs a single instead of $n = 100$ forward passes!**

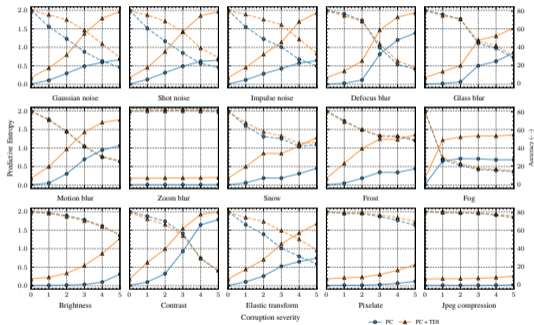
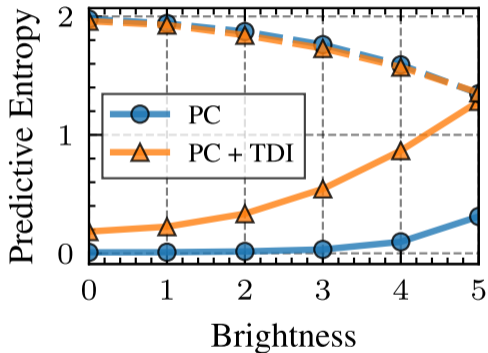
TDI allows PCs to adequately balance ID vs. OOD detection!

Experiments: Data Perturbations (Rotated MNIST)



TDI better captures the distribution shift while retaining predictive accuracy!

Experiments: Data Corruptions (MNIST)



TDI detects distribution shifts and is more robust in predictive accuracy against the corruption!

Future Work

- Influence of **dropout parameter** p ?
- Can we generalize to **arbitrary PC structures**?
- **Density estimation**?
- Can we use **uncertainty during training**?

Conclusion

- Overconfidence in PCs makes ID and OOD data separation challenging
- **Tractable Dropout Inference**: MCD-inspired solution, offers straightforward single-pass uncertainty estimation for PCs which ...
 - enhances PC robustness and helps in detecting distribution shifts
 - removes the computational burden of MCD
 - paves the way to include uncertainty into PC training
 - allows for prior knowledge of epistemic or aleatoric uncertainty

Paper



Code



Backup Slides

Structural Knowledge

The easiest solution is when we know two product nodes P^A and P^B share no common ancestors. This implies $P^A \perp\!\!\!\perp P^B$ and thus $\text{Cov}[P^A, P^B] = 0$.

This is always the case for tree-structured PCs!

For binary tree random and tensorized (RAT) structures, the covariance simplifies to:

$$\begin{aligned}\text{Cov}[P_{l,r}, P_{l',r'}] &= \text{Cov}[S_l^L, S_{l'}^L] \mathbb{E}[S_r^R] \mathbb{E}[S_{r'}^R] \\ &\quad + \text{Cov}[S_r^R, S_{r'}^R] \mathbb{E}[S_l^L] \mathbb{E}[S_{l'}^L] \\ &\quad + \text{Cov}[S_l^L, S_{l'}^L] \text{Cov}[S_r^R, S_{r'}^R] .\end{aligned}$$

In this case, the covariance of two product nodes only depends on the covariance of the input sum nodes of the same graph partition (L or R).

Covariance Bounds

If we don't have structural knowledge, we can still find lower and upper bounds for the covariance using the Cauchy-Schwarz inequality:

$$\begin{aligned} \text{Cov}[N_i, N_j]^2 &\leq \text{Var}[N_i] \text{Var}[N_j] \\ \Leftrightarrow \text{Cov}[N_i, N_j] &\in \left[-\sqrt{\text{Var}[N_i] \text{Var}[N_j]}, \right. \\ &\quad \left. +\sqrt{\text{Var}[N_i] \text{Var}[N_j]} \right] . \end{aligned}$$

Copy-paste Solution

- We can augment the DAG to enforce zero covariance between two nodes, N_A and N_B , by treating their common input as two separate nodes.
- For each node N_C with paths $\text{Path}_A := N_A \rightarrow N_C$ and $\text{Path}_B := N_B \rightarrow N_C$, we “copy” N_C to an equivalent node $N_{C'}$ and replace the original N_C in Path_B with $N_{C'}$.
- This enforces a tree structure on the PC, making the covariance between two inputs of a node N zero.
- In practice, we don't need to modify the DAG. Instead, *we can simply ignore the covariance*.

Classification Uncertainty

In classification, we express class conditionals $p(\mathbf{x} | y_i) = S_i$ with class priors $p(y_i) = c_i$, and find the posterior via Bayes' rule:

$$p(y_i | \mathbf{x}) = \frac{S_i c_i}{\sum_j S_j c_j}.$$

Here, the expectation and variance of the posterior are those of a random variable ratio, $\mathbb{E}\left[\frac{A}{B}\right]$ and $\text{Var}\left[\frac{A}{B}\right]$, with $A = S_i c_i$ and $B = \sum_j S_j c_j$.

While this ratio isn't well-defined, we can approximate it using a second-order Taylor approximation:

$$\begin{aligned}\mathbb{E}\left[\frac{A}{B}\right] &\approx \frac{\mathbb{E}[A]}{\mathbb{E}[B]} - \frac{\text{Cov}[A, B]}{(\mathbb{E}[B])^2} + \frac{\text{Var}[B] \mathbb{E}[A]}{(\mathbb{E}[B])^3} \\ \text{Var}\left[\frac{A}{B}\right] &\approx \frac{\mathbb{E}[A]^2}{\mathbb{E}[B]^2} \left[\frac{\text{Var}[A]}{\mathbb{E}[A]^2} - 2 \frac{\text{Cov}[A, B]}{\mathbb{E}[A] \mathbb{E}[B]} + \frac{\text{Var}[B]}{\mathbb{E}[B]^2} \right].\end{aligned}$$

Classification Uncertainty

We can simplify each component of the previous equations.

The expectations are:

$$\mathbb{E}[A] = \mathbb{E}[S_i c_i] = \mathbb{E}[S_i] c_i$$

$$\mathbb{E}[B] = \mathbb{E}\left[\sum_j S_j c_j\right] = \sum_j \mathbb{E}[S_j] c_j .$$

The variances become:

$$\text{Var}[A] = \text{Var}[S_i] c_i^2$$

$$\text{Var}[B] = \sum_j \text{Var}[S_j] c_j^2 + \sum_{j_1 \neq j_2} \text{Cov}[S_{j_1}, S_{j_2}] c_{j_1} c_{j_2} .$$

Classification Uncertainty

- Covariance between a root node and sum of all root nodes can be deconstructed:

$$\text{Cov}[A, B] = c_i \sum_j c_j \text{Cov}[S_i, S_j] \quad ,$$

- Assumption of statistical independence between A and B is an approximation.
- This approximation has proven effective in practice.

[Seltman N., CMU 2018]

Tractability

- PCs and SPNs are tractable models with linear time queries.
- TDI formulations have at most quadratic space and time complexity.
- Sparse structures like trees allow for linear computational cost.
- In cases needing all covariance combinations, cost is locally quadratic.
- Using the “copy-paste” DAG augmentation, cost can be reduced to linear.
- Full bottom-up pass is tractable and parallelizable.