

# Learning Sum-Product Networks

---

Martin Trapp

# Sum-Product Networks

---

	GANs	VAEs	Flows
Sampling	Y	Y	Y
Density	N	N/Y	Y
Marginals	N	N	?
Conditionals	N	N	?
Moments	N	N	?
MAP	N	N	?

**Table 1:** Robert Peharz, Sum-Product Networks and Deep Learning: A Love Marriage. Talk at ICML, 2019.

- Sum-product networks (SPNs)<sup>1</sup> are a sub-class of so-called probabilistic circuits<sup>2</sup>, that admit tractable probabilistic inference.
- Probabilistic circuits admit many probabilistic inference tasks, such as marginalisation, in linear time in their representation size.

---

<sup>1</sup>H. Poon & P. Domingos: Sum-product networks: A new deep architecture. In UAI, 2011.

<sup>2</sup>Van den Broeck et al.: Tractable probabilistic models: Representations, algorithms, learning and applications. Tutorial at UAI, 2019.

	GANs	VAEs	Flows	SPNs
Sampling	Y	Y	Y	Y
Density	N	N/Y	Y	Y
Marginals	N	N	?	Y
Conditionals	N	N	?	Y
Moments	N	N	?	Y
MAP	N	N	?	N/Y

**Table 2:** Robert Peharz, Sum-Product Networks and Deep Learning: A Love Marriage. Talk at ICML, 2019.

# What is a Sum-Product Network?

Let  $\mathbf{X} = \{X_1, \dots, X_D\}$  be set of random variables.

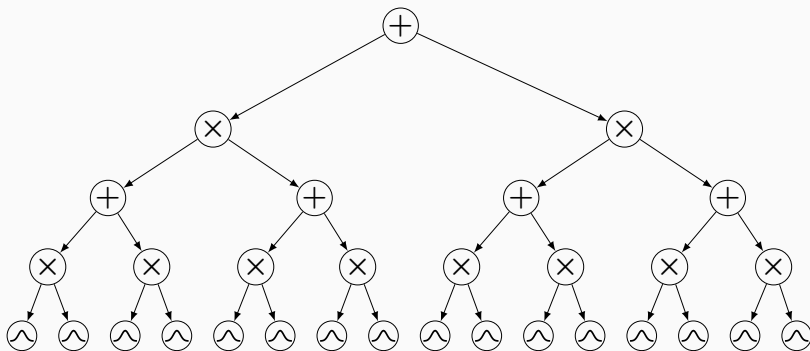
A probabilistic circuit over  $\mathbf{X}$  is a tuple  $\mathcal{S} = (\mathcal{G}, \psi, \theta)$ , where

- $\mathcal{G}$  is a computational graph.
- $\psi$  is a scope function.
- $\theta$  is a set of parameters, e.g. sum-weights and leaf node parameters.

A Sum-Product Network (SPN) is a *smooth* (complete) and *decomposable* probabilistic circuit.

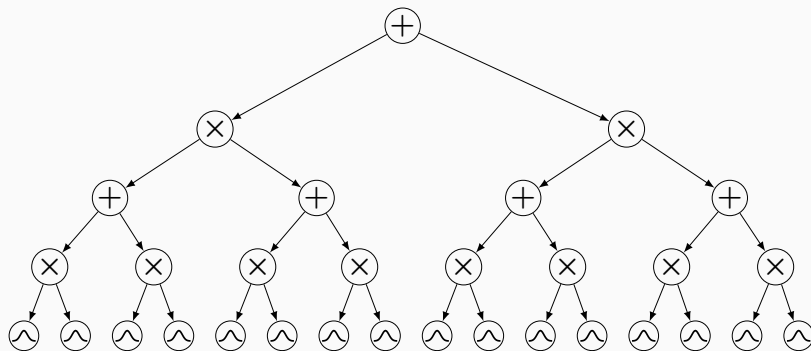
# Computational Graph $\mathcal{G}$

$\mathcal{G}$  is a rooted connected directed acyclic graph (DAG), containing: sum (S), product (P) and leaf nodes (L).



# Computational Graph $\mathcal{G}$

$\mathcal{G}$  is a rooted connected directed acyclic graph (DAG), containing: sum (S), product (P) and leaf nodes (L).



$$S = \sum_{c \in \mathcal{C}(\mathcal{P})} w_{S,c} C(x)$$

$$P = \prod_{c \in \mathcal{C}(\mathcal{P})} C(x)$$

$$L = p(x | \theta_l)$$



## Scope Function $\psi$

The scope function assigning each node  $N$  in a sub-set of  $\mathbf{X}$ ,<sup>3</sup> and has to fulfil the following properties:

1. If  $N$  is the root node, then  $\psi(N) = \mathbf{X}$ .
2. If  $N$  is a sum or product, then  $\psi(N) = \bigcup_{N' \in \text{ch}(N)} \psi(N')$ .

---

<sup>3</sup>This sub-set is often referred to as the scope of a node.

## Scope Function $\psi$

The scope function assigning each node  $N$  in a sub-set of  $\mathbf{X}$ ,<sup>3</sup> and has to fulfil the following properties:

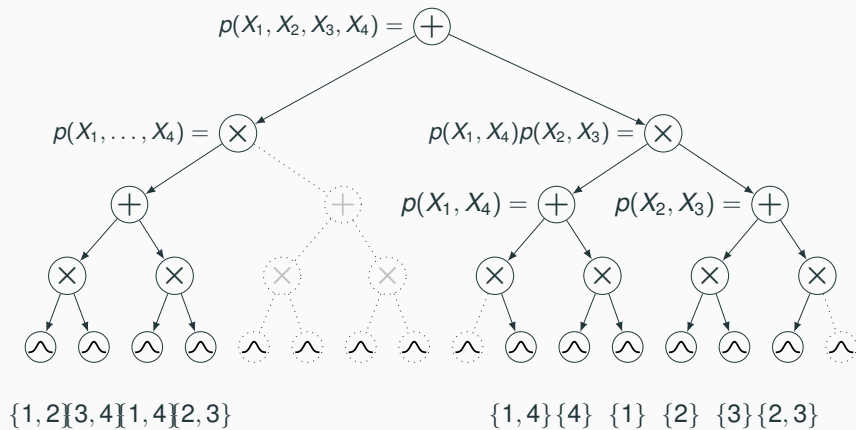
1. If  $N$  is the root node, then  $\psi(N) = \mathbf{X}$ .
2. If  $N$  is a sum or product, then  $\psi(N) = \bigcup_{N' \in \mathbf{ch}(N)} \psi(N')$ .

In case of SPNs we also assume that:

1. For each  $S \in \mathbf{S}$  we have  $\forall N, N' \in \mathbf{ch}(S): \psi(N) = \psi(N')$  (*smoothness*)
2. For each  $P \in \mathbf{P}$  we have  $\forall N, N' \in \mathbf{ch}(P): \psi(N) \cap \psi(N') = \emptyset$  (*decomposability*).

<sup>3</sup>This sub-set is often referred to as the scope of a node.

# Example SPN $\mathcal{S} = (\mathcal{G}, \psi, \theta)$



Note that we define  $L(x) := 1$  for every  $x$  if and only if  $\psi(L) = \emptyset$ .

# Learning Sum-Product Networks

---

# Parameter Learning in SPNs

We can use backprop for parameter learning in SPNs.

More advanced approaches:

- **Expectation Maximisation** [R. Peharz et al.: On the latent variable interpretation of sum-product networks. TPAMI, 2017.]
- **Variational Inference** [H. Zhao et al.: Collapsed variational inference for sum-product networks. In ICML, 2016.]
- **Bayesian moment matching** [A. Rashwan et al.: Online and distributed Bayesian moment matching for parameter learning in SPNs. In AISTATS, 2016.]
- **Safe Semi-Supervised Learning** [M. Trapp et al.: Safe semi-supervised learning of sum-product networks. In UAI, 2017.]

# Challenges in Structure Learning

- The structure has to be smooth and decomposable, i.e., a sparsely connected graph.
- Structure learning has to be efficient.
- How to learn structures that generalise well, many approaches learn deep trees that are prone to overfitting.
- What is a good SPN structure? or What is a good principle to derive an SPN structure?

## Why do we want Bayesian structure learning?

- Occam's razor effect prevents overfitting.
- Works on discrete, continuous and heterogeneous data domains.
- We can use nonparametric formulations, e.g. infinite SPNs, for continual learning.
- Structures can be inferred even in cases of missing values using exact marginalisation.

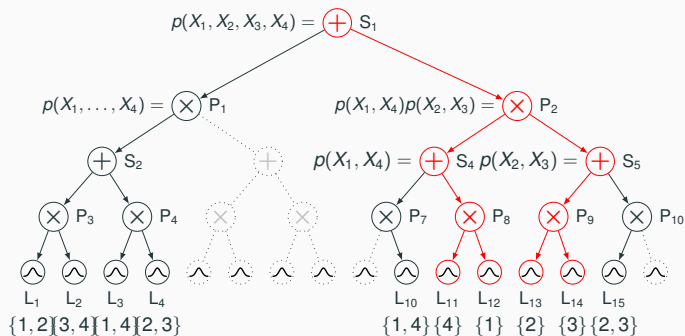
# Bayesian Parameter Learning

The key insight for Bayesian parameter learning is that *sum nodes can be interpreted as latent variables*.



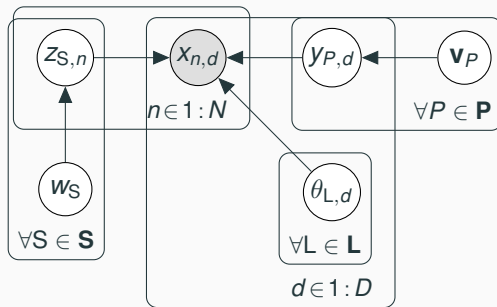
# Bayesian Parameter Learning

The key insight for Bayesian parameter learning is that *sum nodes can be interpreted as latent variables*.



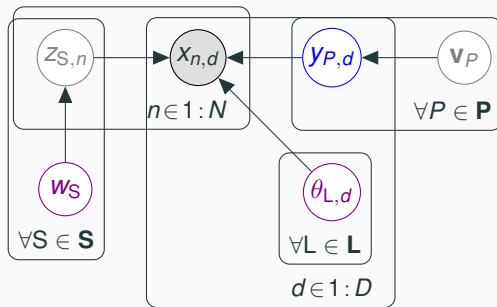
# Bayesian Structure

Generative model for Bayesian learning of SPNs.



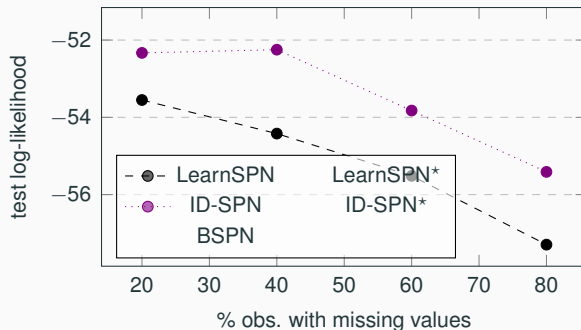
# Bayesian Structure

Posterior inference using ancestral sampling within Gibbs.



# Bayesian Structure - Missing Values Experiment

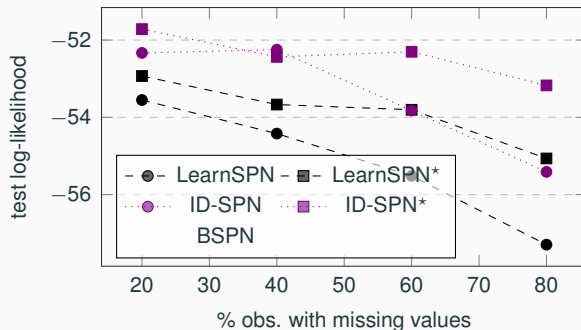
Performance under increasing amount of missing values.



**Figure 1:** Results on EachMovie (D: 500, N: 5526) dataset.

# Bayesian Structure - Missing Values Experiment

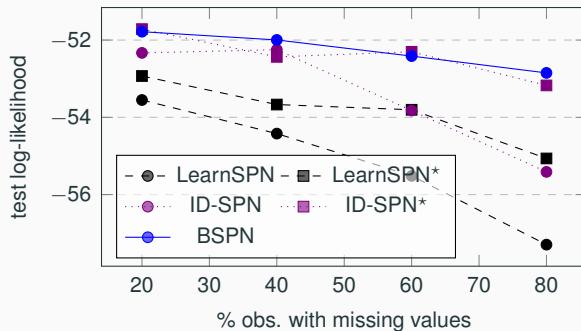
Performance under increasing amount of missing values.



**Figure 1:** Results on EachMovie (D: 500, N: 5526) dataset.

# Bayesian Structure - Missing Values Experiment

Performance under increasing amount of missing values.



**Figure 1:** Results on EachMovie (D: 500, N: 5526) dataset.

# Applications

---

# Existing Applications

Some existing applications of SPNs:

- Computer vision, e.g. image classification, medical image processing, attend-infer-repeat.
- Language processing, e.g. language modelling, bandwidth extension.
- Robotics, e.g. semantic mapping.
- **Non-linear regression**, and many more<sup>4</sup>

---

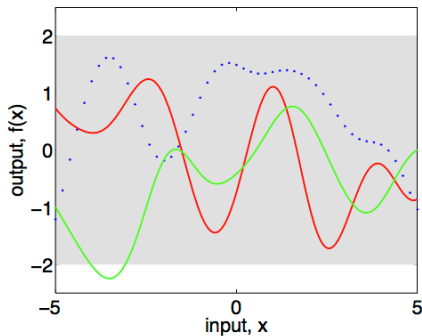
<sup>4</sup><https://github.com/arranger1044/awesome-spn>



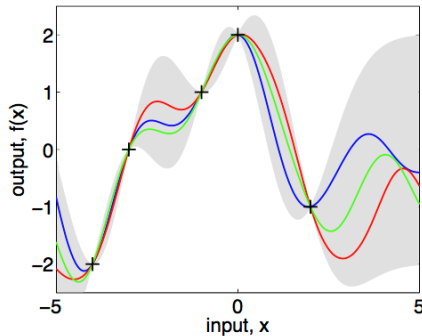
# Gaussian Processes

A Gaussian Process (GP) is a collection of random variables indexed by an arbitrary covariate space  $\mathcal{X}$ , where any finite subset is Gaussian distributed.

A GP is a prior over functions, that admits exact posterior inference.



(a), prior



(b), posterior

A GP is uniquely specified by a *mean-function*  $m : \mathcal{X} \rightarrow \mathbb{R}$  and a *covariance function*  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ .

The posterior predictive distribution (used for predictions) of a GP is Gaussian, i.e.,

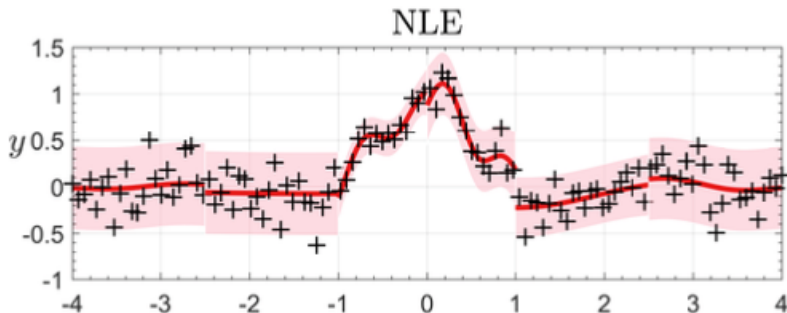
$$p(f^* | \mathbf{f}) = \mathcal{N} \left( k_{\mathbf{x}^*, \mathbf{x}}^T k_{\mathbf{x}, \mathbf{x}}^{-1} \mathbf{f}, k_{\mathbf{x}^*, \mathbf{x}^*} - k_{\mathbf{x}^*, \mathbf{x}} k_{\mathbf{x}, \mathbf{x}}^{-1} k_{\mathbf{x}^*, \mathbf{x}}^T \right) \quad (2)$$

The inversion of  $k_{\mathbf{x}, \mathbf{x}}$  is computed using the Cholesky decomposition of  $k_{\mathbf{x}, \mathbf{x}}$ , which scales  $\mathcal{O}(N^3)$ .

## Local Experts

*Local experts to approximate the GP or approximate the computation of predictions.*

A natural way is to partition  $\mathcal{X}$  into sub-sets  $\mathcal{X}^{(k)}$ ,  $k = 1, \dots, K$ . This is called the naive-local-experts model.



Existing solutions to discontinuities.

1. Product-of-Experts (PoE) / Bayesian Committee Machine (BCM)
  - Instead of partition  $\mathcal{X}$ , partition  $\mathcal{X}$  into sub-sets  $\mathcal{X}^{(k)}$ .
  - Use an algorithm that works only on the sub-sets.
  - **Problem:** Not a stochastic process, results in over-conservative or over-confident estimates.
2. Mixture-of-Experts (MoE)
  - Use a gating network to assign observations to experts instead of hard boundaries.
  - Often intractable (due to the gating network).
3. Impose Continuity Constraints
  - Suffers from inconsistent variances and does not scale.

Why not use a large (finite) mixture of NLEs?

Deep Structured Mixture of GPs (DSMGP)<sup>5</sup>:

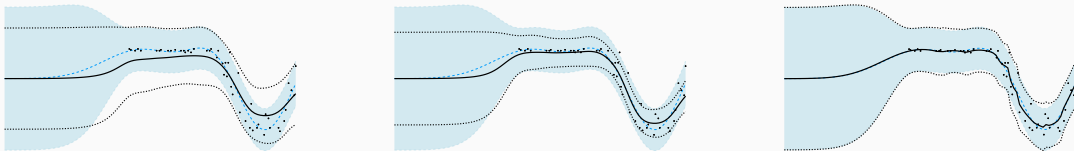
1. A DSMGP is a hierarchically defined convex combination of product measures with Gaussian measures as base measure.
2. DSMGPs perform exact Bayesian model averaging over a large set of NLEs.

---

<sup>5</sup>M. Trapp et al.: Deep structure mixtures of Gaussian Processes. To appear in AISTATS, 2020.

# Deep Structured Mixtures of Gaussian Processes

- DSMGPs are a *sound stochastic process*.
- We can perform *exact posterior inference, efficiently*.
- DSMGPs capture predictive uncertainties consistently better than existing approximations.
- In DSMGPs we can model non-stationary data and perform exact inference over kernel functions.





Thank you for your attention!