

Learning Sum-Product Networks

Martin Trapp

Probabilistic Machine Learning

Probabilistic machine learning: How can we make machines that learn from data using tools from probability theory?

Uncertainty is key to probabilistic machine learning and is expressed in all forms, e.g. noise in the data, uncertainty over the predictions, uncertainty over the model.

Probabilistic Inference

Problem: Probabilistic inference tasks are often intractable for interesting models.

	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

Sum-Product Networks

- Sum-product networks (SPNs)¹ are a sub-class of so-called tractable probabilistic models or probabilistic circuits², that admit tractable probabilistic inference.
- A class of queries Q on a class of models M is tractable, iff for any query $q \in Q$ and model $m \in M$ the computational complexity is at most polynomial.
- SPNs admit many probabilistic inference tasks, such as marginalisation, in linear time in their representation size.

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

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

Probabilistic Inference

	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?

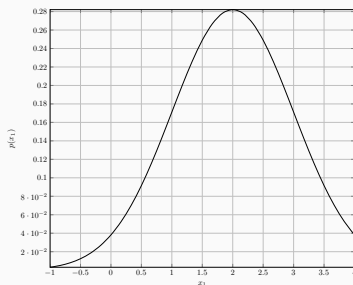
- Let $\mathbf{X} = \{X_1, \dots, X_D\}$ be set of random variables.
- A Probabilistic Circuit (PC) over \mathbf{X} is a tuple $\mathcal{S} = (\mathcal{G}, \psi, \theta)$, where
 - \mathcal{G} is a computational graph.
 - ψ is a scope function.
 - θ is a set of parameters, e.g. sum-weights and leaf node parameters.
- A Sum-Product Network (SPN) is a *smooth* and *decomposable* PC.

Leaves L in \mathcal{G}

Leaf/input nodes are arbitrary (tractable) distributions, e.g. Gaussian, Multinomial, variational autoencoder.



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

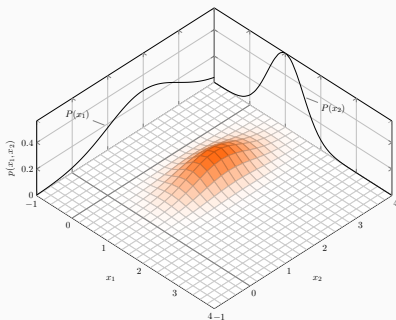


Product Nodes P in \mathcal{G}

Product nodes encode independence assumptions between sets of random variables.



$$P(x) = \prod_{C \in \text{ch}(P)} C(x)$$

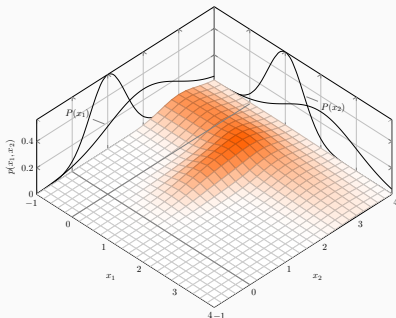


Sum Nodes S in \mathcal{G}

Sum nodes³ replace independence with conditional independence within the network.



$$S(x) = \sum_{C \in \text{ch}(P)} w_{S,C} C(x)$$



³ $\sum_{C \in \text{ch}(S)} w_{S,C} = 1$ and $w_{S,C} \geq 0$.

Computational Graph \mathcal{G}

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

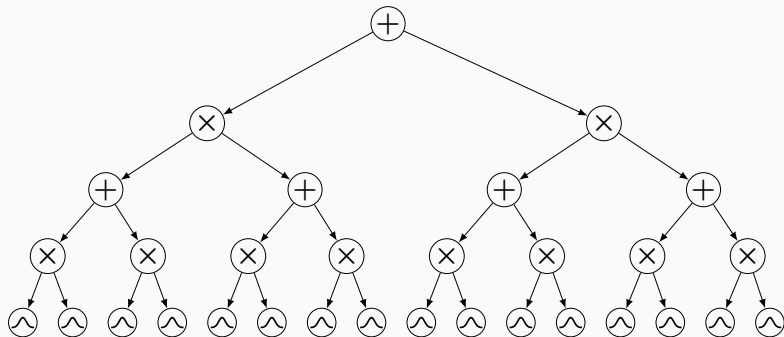


Figure 1: Example of a tree-shaped computational graph.

Scope Function ψ

The scope function assigning each node N in a sub-set of \mathbf{X} ,⁴ and has to fulfil the following properties:

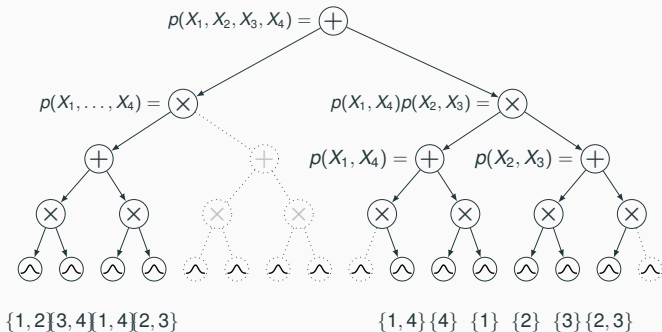
1. If N is the root node, then $\psi(N) = \mathbf{X}$.
2. If N is a sum or product, then
$$\psi(N) = \bigcup_{N' \in \mathbf{ch}(N)} \psi(N').$$

In case of SPNs we also assume that:

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

⁴This 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 $\psi(L) = \emptyset$.

Parameter Learning

Parameter Learning in SPNs

Generative Learning^a

$$\mathcal{L}(\theta | \mathcal{X}) = \sum_{n=1}^N \log \mathcal{S}(\mathbf{x}_n | \theta) - \log \mathcal{S}(* | \theta), \quad \mathbf{x}_n \in \mathbb{R}^D \quad (1)$$

Note that $\mathcal{S}(* | \theta)$ is the partition function which can be evaluated efficiently using a single upward pass.

Discriminative Learning^b

$$\mathcal{L}(\theta, \lambda | \mathcal{X}) = \sum_{n=1}^N \log \mathcal{S}(\mathbf{x}_n, \lambda_n | \theta) - \log \mathcal{S}(\mathbf{x}_n | \theta), \quad \mathbf{x}_n \in \mathbb{R}^D, \lambda_n \in \mathbb{R} \quad (2)$$

^aH. Poon & P. Domingos: Sum-product networks: A new deep architecture. In UAI, 2011.

^bR. Gens & P. Domingos: Discriminative learning of sum-product networks. In NeurIPS, 2012.

Contrastive Pessimistic Likelihood Estimation (CPLE)⁵

- Assume $\mathcal{X} = \{\mathbf{x}_n, \lambda_n\}_{n=1}^N$, $\lambda_n \in \mathcal{R}^K$ to be a small set of labelled data points and $\mathcal{U} = \{\mathbf{u}_m\}_{m=1}^M$ a set of unlabelled data points.
- θ^+ is an SPN trained solely on \mathcal{X}
- We propose soft labels $\mathbf{q}_m \in \Delta_{K-1}$ for all unlabelled data points, e.g. randomly.
- We initialise the semi-supervised SPN θ randomly, then:
 1. Update all soft labels so that the performance of θ^+ on all data points is better than θ .
 2. Update θ so that the performance of θ on all data points is better than θ^+ .
 3. If not converged, move to step 1.

⁵M. Trapp et al.: Safe semi-supervised learning of sum-product networks. In UAI, 2017.

Semi-Supervised Learning

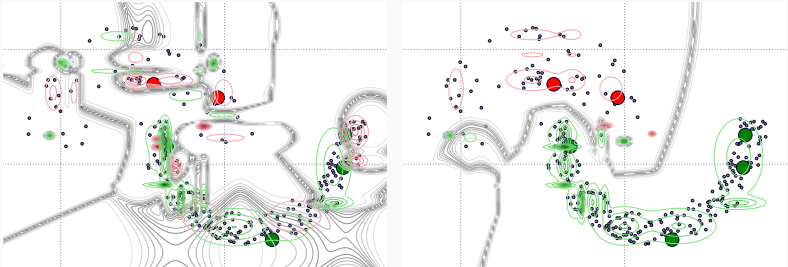


Figure 2: Decision boundary of a semi-supervised SPN at iteration 1 (left) and after convergence (right).

Over-parametrization in SPNs⁶

$$w_k^{(t)} \approx w_k^{(t)} + \rho^{(t)} \nabla_{w_k^{(t)}} + \left[\sum_{l=0}^{L-1} \eta \nabla_{w_{\phi(k,l)}^{[l]}} (w_{\phi(k,l)}^{[l]})^{-1} \right] w_k^{(t)} \quad (3)$$

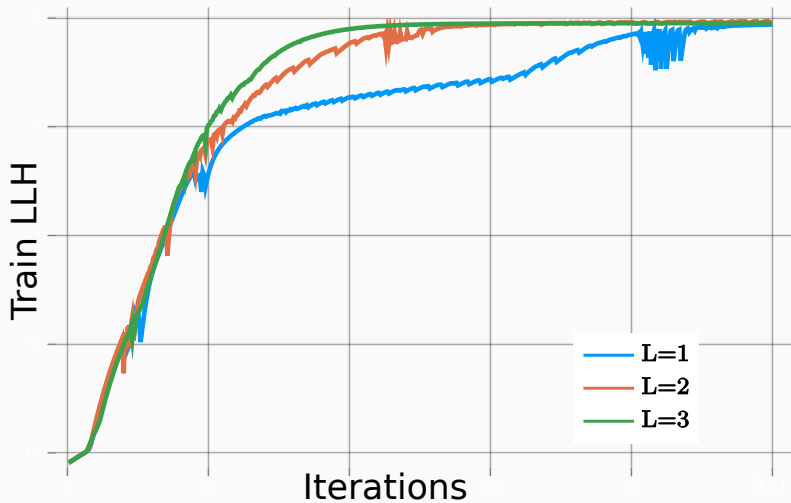
$$= w_k^{(t)} + \rho^{(t)} \nabla_{w_k^{(t)}} + \sum_{\tau=1}^{t-1} \mu^{(t,\tau)} \nabla_{w_k^{(\tau)}} \quad (4)$$

Gradient-based optimisation in deep tree-structured sum-product network with small (fixed) learning rate and near zero initialisation of the weights is equivalent to gradient-based optimisation with adaptive and time-varying **learning rate** and **momentum term**.

⁶M. Trapp et al.: Optimisation of Overparametrized Sum-Product Networks. ICML Workshop on Tractable Probabilistic Models, 2019.

Over-parametrization Experiment

Overparameterization of SPNs on NLPs



Structure Learning

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?⁷

⁷M. Trapp et al.: Bayesian learning of sum-product networks. In NeurIPS, 2019.

General-Purpose Learners (Selection)

- LearnSPN⁸ recursively constructs a tree structure.
- ID-SPN⁹ is a generalisation of LearnSPN with tractable Markov networks as leaves.
- **RAT-SPN**¹⁰ constructs large SPNs with random decompositions.
- **BSPN**⁸ a well-principled framework to learn structures using Bayesian inference.

⁸R. Gens & P. Domingos: Learning the structure of sum-product networks. In ICML, 2013.

⁹A. Rooshenas & D. Lowd: Learning Sum-Product Networks with Direct and Indirect Variable Interactions. In ICML, 2014.

¹⁰R. Peharz et al.: Random sum-product networks: A simple but effective approach to probabilistic deep learning. In UAI, 2019.

Todo

Why care about Bayesian structure learning?

- Competitive results.
- Occam's razor effect prevents overfitting.
- Works on heterogeneous data domains
- We can use nonparametric formulations, e.g. infinite SPNs.
- Can be used in cases of missing values, by using exact marginalisation of missing values.

Bayesian Structure

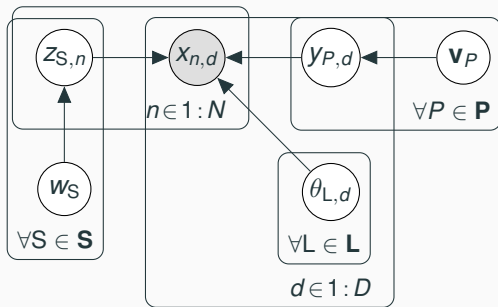


Figure 3: Generative model for Bayesian structure learning.

Bayesian Structure

Posterior inference can be performed using ancestral within Gibbs sampling.

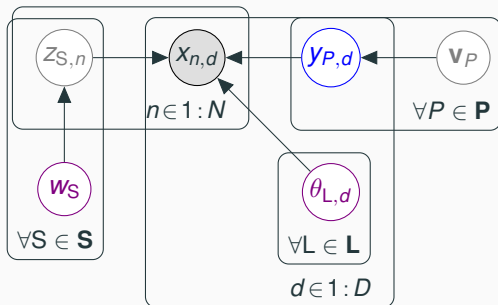
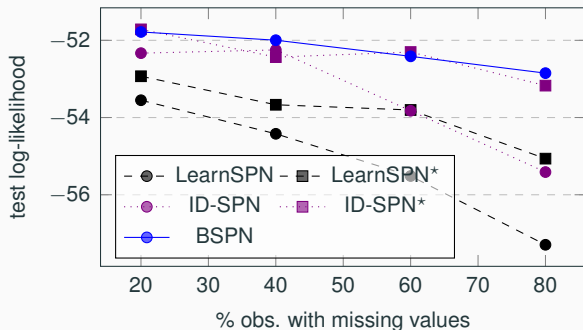


Figure 4: Generative model for Bayesian structure and learning.

Bayesian Structure - Missing Values Experiment

Performance under increasing amount of missing values
(missing completely at random).



Applications

There exist several applications of SPNs in various fields:

- 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.
- Nonlinear regression, and many more¹¹

¹¹<https://github.com/arranger1044/awesome-spn#applications>

We have recently¹² introduced a combination of SPNs with Gaussian processes, which yields an interesting nonparametric regression model that allows exact and efficient posterior inference.

Our model can be understood as an exponentially large mixture over naive-local-experts of Gaussian processes.

¹²M. Trapp et al.: Deep structured mixtures of Gaussian processes. To appear at AISTATS, 2020.

Deep Structured Mixtures of GPs

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

A GP can be understood as a prior over functions.

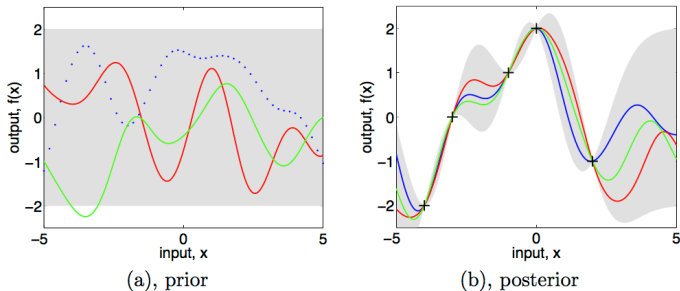


Figure 5: C. E. Rasmussen & C. K. I. Williams, Gaussian Processes for Machine Learning, 2006.