

Tractable probabilistic models for hierarchical data

Tomáš Pevný

Czech Technical University in Prague

February 1, 2024

Problem statement

We assume a set of samples $\{x_i | x_i \in \mathcal{X}\}_{i=1}^n$.

We want to fit a model $p(x|\theta)$ by maximizing likelihood

$$\arg \max_{\theta} \sum_i \log p(x_i|\theta)$$

such that

- ▶ $p(x|\theta)$ is a valid probability distribution
- ▶ $p(x|\theta)$ is tractable,

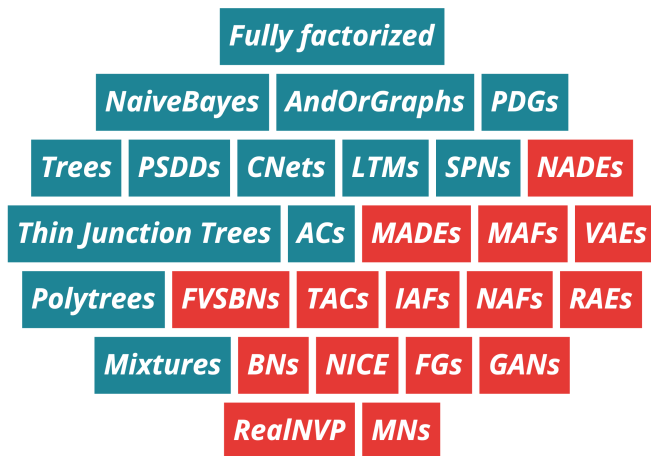
What is tractability?

The model $p(x)$ is tractable with respect to $f \in \mathcal{F}$ if the integral

$$\int_{\mathcal{X}} f(x)p(x|\theta)dx,$$

can be computed in polynomial time with respect to the size of $p(x|\theta)$.

Tractability



Sum-Product networks

- ▶ in **leaf node** $L(x)$

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

- ▶ in **sum node** $S(x)$

$$p(x) = \sum_{N \in \text{ch}(S)} w_N p_N(x_{\psi(N)})$$

- ▶ in **product node** $P(x)$

$$p(x) = \prod_{N \in \text{ch}(P)} p_N(N(x_{\psi(N)})),$$

$\psi(N)$ is a scoping function.

Special cases of SPNs

- ▶ Mixture model

$$p(\vec{x}) = \sum_{i=1}^n w_i p_i(x), \quad \sum_{i=1}^n w_i = 1$$

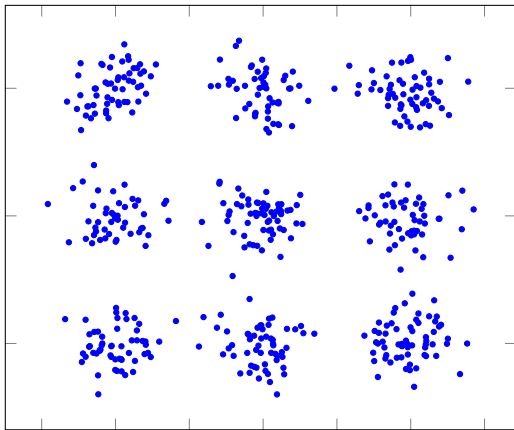
- ▶ Product of marginals (naive bayes)

$$p(\vec{x}) = \prod_{i=1}^n p_i(x_i)$$

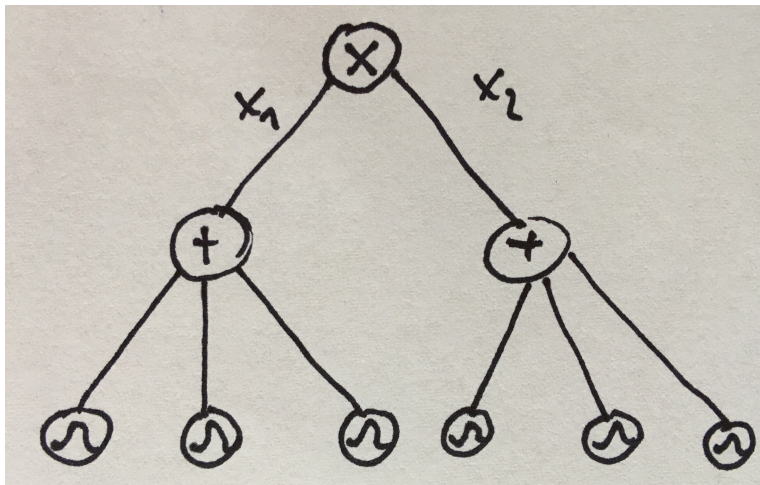
- ▶ Diagonal mixtures

$$p(\vec{x}) = \sum_{i=1}^n w_i \prod_{j=1}^n p_{ij}(x_j)$$

Are SPNs useful?



Are SPNs useful?



Scoping function

Scoping function $\psi(N)$ is used to specify a subspace on which the distribution of a given node is defined.

Example

Let, $\mathcal{X} = \mathbb{R}^n$, but the node L operates just on features (x_1, x_5, x_6) , the scope function $\psi(L) = (x_1, x_5, x_6)$.

Tractability

Tractability of SPNs

Sum-Product networks are tractable

- ▶ if they are *smooth* and *decomposable*
- ▶ and for query f holds

$$f(x) = \prod_{\psi_u \in \cup\{\psi(L)\}} f_u(\psi_u),$$

where

- ▶ $\cup\{\psi(L)\}$ is the set of all possible scopes of leafnodes
- ▶ probability distribution of leafnodes are tractable

Smoothing and Decomposability

- ▶ **Smoothness:** All childs of sumnode S has to have the same scope as S, i.e.

$$\psi(a_1) = \psi(a_2) \quad \forall a_1, a_2 \in \mathbf{ch}[S],$$

and weights $w_i \geq 0$ and $\sum_i w_i = 1$.

- ▶ **Decomposability:** Scopes of childs of productnode P are pairwise disjoint

$$\bigcap_{a \in \mathbf{ch}[P]} \psi(a) = \emptyset,$$

but complete $\psi(P) = \bigcup_{a \in \mathbf{ch}[P]} \psi(a)$.

- ▶ The scope of rootnode R is over all features, $\psi(R) = [P]$.

Recursive computation of integrals

For *smooth* and *decomposable* SPNs,
the integral $I = \int f(x)p(x)dx$ can be computed recursively as
follows:

$$I_u = \begin{cases} \sum_{c \in \mathbf{ch}(u)} w_{u,c} I_c, & \text{for } u \in S, \\ \prod_{c \in \mathbf{ch}(u)} I_c, & \text{for } u \in P, \\ \int f_u(\psi_u) p_u(\psi_u) d\psi_u, & \text{for } u \in L, \end{cases}$$

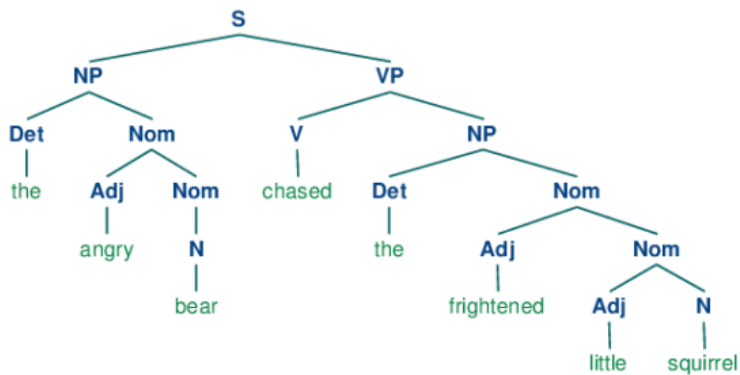
Hierarchically-Structured trees

also called Hierarchical Multi-Instance Learning

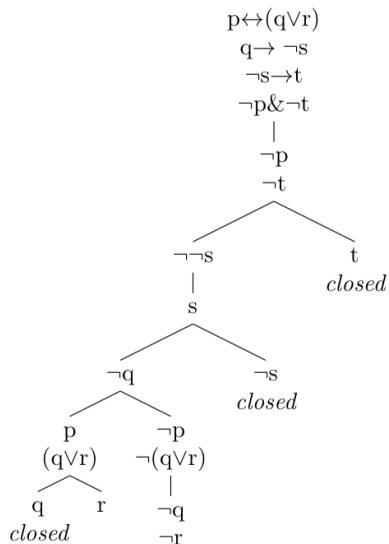
Motivation

```
{
  services: [
    {
      port: 22,
      protocol: tcp
    },
    {
      port: 4070,
      protocol: tcp
    },
    {
      port: 4071,
      protocol: tcp
    },
    {
      port: 5353,
      protocol: udp
    }
  ],
  device_id: 8bb8971c-5983-4baa-9753-f0ac21faf162,
  ip: 192.168.1.80,
  mac: ac:63:be:a5:50:43,
  mdns_services: [_workstation._tcp.local., _ssh._tcp.local., _sftp-ssh._tcp.local.]
}
```

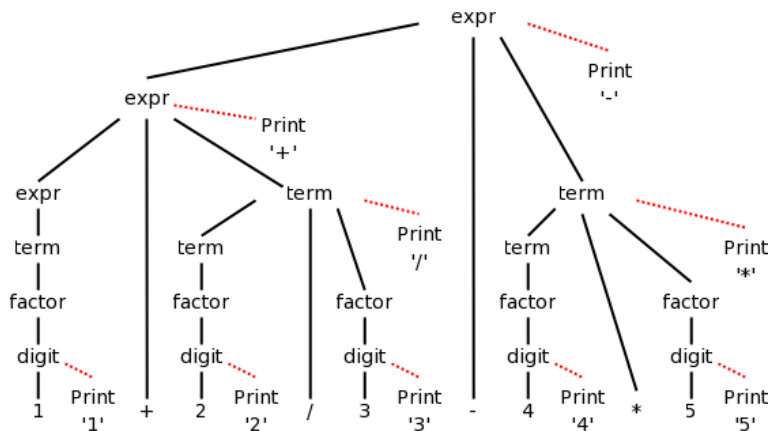
Motivation — semantic tree



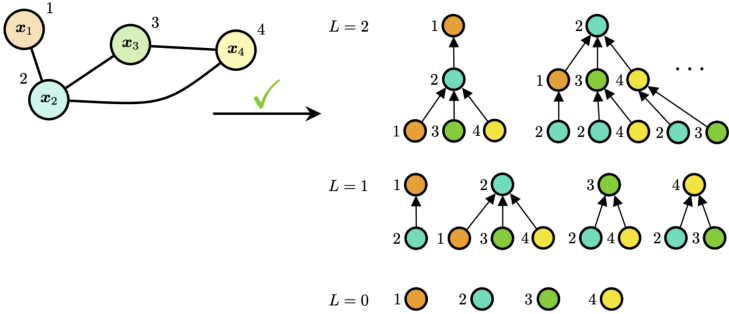
Motivation — logic



Motivation — abstract syntax tree



Motivation — general graph



Types of nodes

```
{
  services: [
    {
      port: 22,
      protocol: tcp
    },
    {
      port: 4070,
      protocol: tcp
    },
    {
      port: 4071,
      protocol: tcp
    },
    {
      port: 5353,
      protocol: udp
    }
  ],
  device_id: 8bb8971c-5983-4baa-9753-f0ac21faf162,
  ip: 192.168.1.80,
  mac: ac:63:be:a5:50:43,
  mdns_services: [_workstation._tcp.local., _ssh._tcp.local., _sftp-ssh._tcp.local.]
}
```

- ▶ Dictionaries
- ▶ Lists (sets)
- ▶ Leafs

Set node

- ▶ Set node computes a probability density over a hyper-space

$$\bar{\mathcal{X}} = \cup_{m=0}^{\infty} \underbrace{\mathcal{X} \times \dots \times \mathcal{X}}_m$$

- ▶ probability density of a *set node* B is computed as

$$p_B(x) = p(m)c^m m! p(x_1, \dots, x_m)$$

where

- ▶ $p(m)$ is a cardinality distribution
- ▶ $p(x)$ is a probability distribution on \mathcal{X}
- ▶ c is a constant for unit normalization

Why the factorial?

- ▶ We need probability distribution on sets $\{x_1, \dots, x_m\}$, but $p(x_1, \dots, x_m)$ is a probability distribution on Cartesian space.
- ▶ We define distribution on sets as

$$p(\{x_1, \dots, x_m\}) = \sum_{\pi \in \text{perm}} p(x_{\pi(1)}, \dots, x_{\pi(m)})$$

- ▶ If p is *exchangeable*, it can be simplified to

$$p(\{x_1, \dots, x_m\}) = m! p^{\text{sym}}(x_1, \dots, x_m)$$

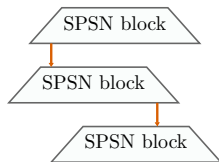
- ▶ *cluster model*

$$p(\{x_1, \dots, x_m\}) = m! \prod_{i=1}^m p(x_i)$$

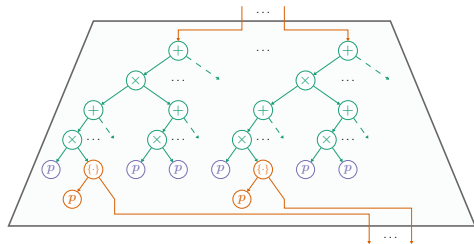
Construction of Sum-Product-Set network



(a) Schema



(b) SPSN



(c) SPSN block

Recursive computation of integrals

For *smooth* and *decomposable* SPSNs,
the integral $\int f(x)p(x)dx$ can be computed recursively as follows:

$$I_u = \begin{cases} \sum_{k=0}^{\infty} p(k) \prod_{i=1}^k I_i, & \text{for } u \in B, \\ \sum_{c \in \text{ch}(u)} w_{u,c} I_c, & \text{for } u \in S, \\ \prod_{c \in \text{ch}(u)} I_c, & \text{for } u \in P, \\ \int f_u(\psi_u) p_u(\psi_u) d\psi_u, & \text{for } u \in L, \end{cases}$$

Experimental comparison

dataset	MLP	GRU	LSTM	HMIL	SPSN
chess	0.41 ± 0.03	0.41 ± 0.05	0.34 ± 0.04	0.39 ± 0.02	0.39 ± 0.03
citeseer	0.69 ± 0.02	0.74 ± 0.01	0.74 ± 0.02	0.75 ± 0.01	0.75 ± 0.01
cora	0.75 ± 0.03	0.86 ± 0.01	0.84 ± 0.01	0.85 ± 0.00	0.86 ± 0.01
genes	0.99 ± 0.01	1.00 ± 0.01	0.98 ± 0.01	1.00 ± 0.01	0.95 ± 0.01
hepatitis	0.86 ± 0.02	0.88 ± 0.01	0.87 ± 0.03	0.88 ± 0.02	0.88 ± 0.02
mutagenesis	0.84 ± 0.02	0.83 ± 0.02	0.82 ± 0.04	0.83 ± 0.00	0.84 ± 0.02
uwcse	0.84 ± 0.02	0.87 ± 0.03	0.85 ± 0.02	0.86 ± 0.03	0.84 ± 0.02
webkp	0.77 ± 0.02	0.82 ± 0.01	0.81 ± 0.02	0.82 ± 0.01	0.81 ± 0.02
rank	3.62	1.62	3.88	1.62	2.38

Conclusion

- ▶ We have extended SPNs to a class of HS-trees (HMIL)
- ▶ It is the first generative model for this class of problems.
- ▶ Seems to deliver similar accuracy to non-tractable models.
- ▶ Is it useful?