The Sum-Product Theorem: A Foundation for Learning Tractable Models

Abram Friesen and Pedro Domingos

{afriesen, pedrod}@cs.uw.edu

Key inference problems in Al

Marginal probability MPE Satisfiability Constraint satisfaction Integration Nonconvex optimization Database querying

 $=\sum \prod P_i(x)$ $= \max_{x \in X} \prod P_i(x)$ $= \bigvee_{x \in X} \bigwedge_{i} N_i(x)$ $= \bigvee \bigwedge C_i(x)$ $F_i(x)$ $= \min_{x \in X} \sum_{i \in X} F_i(x)$ $= \bigcup \bowtie_i Q_i(x)$

These problems all have the same general form: summing a function over a semiring. [Bistarelli et al. (1997); Aji & McEliece (2000); Wilson (2005); Green et al. (2007); Dechter & Mateescu (2007)]

Inference, in general

We are interested in computing summations $\bigoplus_{x \in X} S(x)$ where $(R, \oplus, \otimes, 0, 1)$ is a commutative semiring, $S : X \rightarrow R$ is a function on that semiring, and $X = \{X_1, ..., X_n\}$ is a set of variables.

We refer to *S* as a sum-product function (SPF).



Definition. A (*commutative*) semiring $(R, \oplus, \otimes, 0, 1)$ is a nonempty set R on which operations sum (\oplus) and product (\otimes) are associative and commutative and have identity elements $0, 1 \in R$, such that 1) product (\otimes) distributes over sum (\oplus), and 2) $0 \neq 1$, $a \oplus 0 = a$, $a \otimes 1 = a$, and $a \otimes 0 = 0$ for all $a \in R$.

Definition. A sum-product function (SPF) $S : X \rightarrow R$ over (R, X, Φ) is any of 1) a function: $\phi_j(X_i), \phi_j \in \Phi$ 2) a product of SPFs: $S_1(X_1) \otimes S_2(X_2)$, or $S_1(X_1) \oplus S_2(X_2),$ 3) a sum of SPFs:

We identify and prove the **sum-product theorem**, which states a simple sufficient condition (decomposability) for tractable, high-treewidth, exact inference in any problem with this form.

Based on it, we show how to define and learn tractable, hightreewidth representations for any such problem.

In general, the cost of computing $\bigoplus_X S(X)$ is O(exp(n)).

where *R* is a semiring, $X = \{X_1, ..., X_n\}$ is a set of variables with finite domains, and $\Phi = \{\phi_j\}$ is a set of constant and univariate functions.

Definition. A product node is *decomposable* iff the scopes of its children are disjoint. An SPF is *decomposable* iff all of its product nodes are decomposable.

The sum-product theorem

Decomposability is a simple condition that defines a class of expressive functions for which inference is tractable.

Theorem (sum-product theorem). Every decomposable SPF *S* can be summed in time linear in its size.

Proof: push outer summation to leaves.



Corollary. Every SPF with *n* variables and treewidth bounded by a constant can be summed in time O(n), but not every SPF that can be summed in time O(n) has bounded treewidth.

Learning tractable representations with LearnSPF

From the sum-product theorem, identifying and exploiting decomposability is the key to learning a tractable model.

To decompose, an algorithm must find a partition of *X* into $\{X_1, X_2\}$ such that $\bigoplus_X S(X) \approx (\bigoplus_{X_1} S_1(X_1)) \otimes (\bigoplus_{X_2} S_2(X_2)).$

Based on this, LearnSPF is able to learn tractable, hightreewidth models in any semiring.



- Sum nodes: $\bigoplus_X \bigoplus_i S_i(X_i) = \bigoplus_i (\bigoplus_{X_i} S_i(X_i) \otimes C)$
- Product nodes: $\bigoplus_X \bigotimes_i S_i(X_i) = \bigotimes_i \bigoplus_{X_i} S_i(X_i)$ $\cos t = d^n$ $\cos t = d \cdot n$
- Leaf nodes can be summed in constant time
- Sum *S* by summing each leaf node and then evaluating the remaining nodes bottom-up. Where *C* is a constant and $i \in$ Children(node)



Specific choices for decomposition, clustering and leaf-creation subroutines depend on the domain. Instances with analogous decomposability structure can be clustered by virtually any algorithm, including naive Bayes and *k*-means. Correlation and independence tests can be used to identify decomposability.

The SPT is the most recent step in a long line of work on tractable inference, which includes Darwiche (2001, 2003); Darwiche & Marquis (2002); Bacchus et al. (2002, 2009), Dechter & Mateescu (2007); Poon & Domingos (2011); and Gens & Domingos (2013).

Applications to specific semirings

Domain	Inference task	Semiring	Variables	Leaf functions
Logical inference	SAT #SAT MAX-SAT	$ \begin{aligned} & (\mathbb{B}, \lor, \land, 0, 1) \\ & (\mathbb{N}, +, \times, 0, 1) \\ & (\mathbb{N}_{-\infty}, \max, +, -\infty, 0) \end{aligned} $	Boolean Boolean Boolean	Literals Literals Literals
Constraint satisfaction	CSPs Fuzzy CSPs Weighted CSPs	$(\mathbb{B}, \lor, \land, 0, 1)$ ([0, 1], max, min, 0, 1) ($\mathbb{R}_{+,\infty}$, min, +, ∞ , 0)	Discrete Discrete Discrete	Univariate constraints Univariate constraints Univariate constraints
Probabilistic inference	Marginal MPE	$(\mathbb{R}_+, +, \times, 0, 1)$ $(\mathbb{R}_+, \max, \times, 0, 1)$	Discrete Discrete	Potentials Potentials
Continuous functions	Integration Optimization	$(\mathbb{R}_+, +, \times, 0, 1)$ $(\mathbb{R}_\infty, \min, +, \infty, 0)$	Continuous Continuous	Univariate functions Univariate functions
Relational databases	Unions of CQs Provenance	$egin{aligned} &(\mathbf{U}_m,\cup,\Join,arnothing,1_R)\ &(\mathbb{N}[\mathbf{X}],+, imes,0,1) \end{aligned}$	Sets of tuples Discrete	Unary tuples <i>K</i> -relation tuples

Learning efficiently-optimizable nonconvex functions

SPFs can be extended to continuous (real) variables.

Corollary. The global minimum of a decomposable SPF on the min-sum semiring (i.e., a MSF) can be found in time linear in its size.

RDIS implicitly constructs a decomposable MSF. [Friesen & Domingos (2015)]

Typically, optimizing nonconvex functions is hard, and learning them is even harder; however, LearnSPF provides a method for effectively learning nonconvex functions that can then be efficiently optimized.

Experiment: learning and optimizing a decomposable MSF with

T	ODT	1 •	

Learning tractable knowledge bases

SPFs on the Boolean semiring correspond to negation normal form (NNF), and summation of an NNF *F* is $\bigvee_{x \in X} F(x) = SAT(F)$

Corollary (Darwiche, 2001). The satisfiability of a decomposable NNF is decidable in time linear in its size.

Learned rule sets are typically encoded in large CNF knowledge bases (KBs), making reasoning over them highly intractable. In contrast, LearnSPF provides a method for directly learning large, complex KBs that are encoded in decomposable NNF and therefore support efficient querying.

LearnSPF versus learning and optimizing a nonconvex function.

- Problem is a continuous variant of structured prediction.
- LearnSPF used *K*-means to cluster and correlation to decompose.



• Dataset: 300 train, 50 test, with labels $y^{(i)} = \arg \min_{y \in Y} F_{x^{(i)}}(y)$.

Department of Computer Science and Engineering, University of Washington