# A Minimalistic Approach to Sum-Product Network Learning for Real Applications

**Viktoriya Krakovna**
Department of Statistics, Harvard University
vkrakovna@fas.harvard.edu

**Moshe Looks**
Google
madscience@google.com

## Abstract

Sum-Product Networks (SPNs) are a class of expressive yet tractable hierarchical graphical models. LearnSPN is a structure learning algorithm for SPNs that uses hierarchical co-clustering to simultaneously identifying similar entities and similar features. The original LearnSPN algorithm assumes that all the variables are discrete and there is no missing data. We introduce a practical, simplified version of LearnSPN, MiniSPN, that runs faster and can handle missing data and heterogeneous features common in real applications. We demonstrate the performance of MiniSPN on standard benchmark datasets and on two datasets from Google's Knowledge Graph exhibiting high missingness rates and a mix of discrete and continuous features.

## 1 Introduction

The Sum-Product Network (SPN) [Poon & Domingos, 2011] is a tractable and interpretable deep model. An advantage of SPNs over other graphical models such as Bayesian Networks is that they allow efficient exact inference in linear time with network size. An SPN represents a multivariate probability distribution with a directed acyclic graph consisting of sum nodes (clusters over instances), product nodes (partitions over features), and leaf nodes (univariate distributions over features), as shown in Figure 1.

The standard algorithms for learning SPN structure assume discrete data with no missingness, and mostly test on the same set of benchmark data sets that satisfy these criteria [Rooshenas & Lowd, 2014]. This is not a reasonable assumption when dealing with messy data sets in real applications. The Google Knowledge Graph (KG) is a semantic network of facts, based on Freebase [Bollacker et al., 2008], used to generate Knowledge Panels in Google Search. KG data is quite heterogeneous, with a lot of it missing, since much more is known about some entities in the graph than others. High missingness rates can also worsen the impact of discretizing continuous variables before doing structure learning, which results in losing more of the already scarce covariance information.

Applications like the KG are common, and there is a need for an SPN learning algorithm that can handle this kind of data. In this paper, we present MiniSPN, a simplification of a state-of-the-art SPN learning algorithm that improves its applicability, performance and speed. We demonstrate the performance of MiniSPN on KG data and on standard benchmark data sets.

## 2 Variation on the LearnSPN algorithm

LearnSPN [Gens & Domingos, 2013] is a greedy algorithm that performs co-clustering by recursively partitioning variables into approximately independent sets and partitioning the training data into clusters of similar instances, as shown in Figure 2. The variable and instance partitioning steps are applied to data slices (subsets of instances and variables) produced by previous steps. The variable partition step uses pairwise independence tests on the variables, and the approximately independent sets are the connected components in the resulting dependency graph. The instance clustering step uses a naive Bayes mixture model for the clusters, where the variables in each cluster are assumed independent. The clusters are learned using hard EM with restarts, avoiding overfitting using an exponential prior on the number of clusters. The splitting process continues until the data
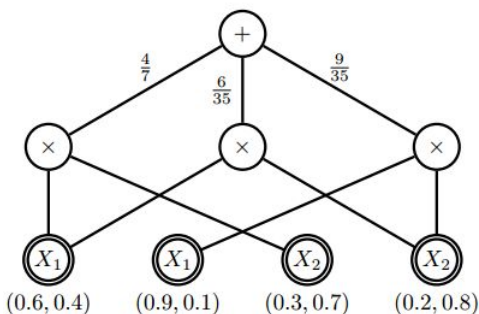
Figure 1: Example of an SPN structure
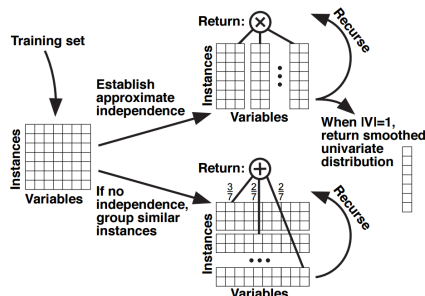(figure from Zhao et al. [2015])



Figure 2: Recursive partitioning process in the LearnSPN algorithm
(figure from Gens & Domingos [2013])

slice has too few instances to test for independence, at which point all the variables in that slice are considered independent. The end result is a tree-structured SPN.

The standard LearnSPN algorithm assumes that all the variables are discrete and there is no missing data. Hyperparameter values for the cluster penalty and the independence test critical value are determined using grid search. The clustering step seems unnecessarily complex, involving a penalty prior, EM restarts, and hyperparameter tuning. It is by far the most complicated part of the algorithm in a way that seems difficult to justify, and likely the most time-consuming due to the restarts and hyperparameter tuning. We propose a variation on LearnSPN called MiniSPN that handles missing data, performs lazy discretization of continuous data in variable partition step, simplifies the model in the instance clustering step, and does not require hyperparameter search.

We simplify the naive Bayes mixture model in the instance clustering step by attempting a split into two clusters at any given point, and eliminating the cluster penalty prior, which results in a more greedy approach than in LearnSPN that does not require restarts or hyperparameter tuning. This seems like a natural choice of simplification - an extension of the greedy approach used at the top level of the LearnSPN algorithm. We compare a partition into univariate leaves to a mixture of two partitions into univariate leaves (generated using hard EM), and the split succeeds if the two-cluster version has higher validation set likelihood. If the split succeeds, we apply it to each of the two resulting data slices, and only move on to a variable partition step after the clustering step fails. The greedy approach is similar to the one used in the SPN-B method [Vergari et al., 2015], which however alternates between variable and instance splits by default, thus building even deeper SPNs.

In the variable partition step, we perform an independence test using the subset of rows where both variables are not missing, and conclude independence if the number of such rows is below threshold. We apply binary binning to each continuous variable, using its median in the data slice as a cutoff.

We compare to the "Pareto" algorithm, previously used for learning SPN models in KG, inspired by the work of Grosse et al. [2012]. It produces a Pareto-optimal set of models, trading off between degrees of freedom and validation set log likelihood score. At each iteration, production rules are randomly applied to add partition and mixture splits to the models in the current model set, and the new models are added to the model set. If a model in the model set has both lower degrees of freedom and higher log likelihood score than another model, the inferior model is removed from the set. The algorithm returns the model from the Pareto model set with the highest validation log likelihood. We also compare to a hybrid method, with the Pareto algorithm initialized by MiniSPN.

## 3 SUMMARY OF EXPERIMENTS

We use two data sets from the Knowledge Graph People collection. In the KG Professions data set, most of the variables are boolean indicators of whether each person belongs to a particular profession. There are 83 boolean variables and 4 continuous variables. In the KG Dates data set, there are 14 continuous variables representing dates of life events for each person and their spouse(s),

Table 1: Average log likelihood and runtime comparison on KG data sets (best performing methods are shown in bold).

| Data set | Test set log likelihood | | | Runtime (seconds) | | |
|---|---|---|---|---|---|---|
| | Pareto | Hybrid | MiniSPN | Pareto | Hybrid | MiniSPN |
| Professions-10K | -10.2 | -6.2 | **-6.09** | 5.3 | 3.7 | **0.4** |
| Professions-100K | -6.61 | -6.53 | **-6.44** | 72 | 131 | **7.2** |
| Dates-10K | -8.66 | **-8.53** | -8.68 | 1.7 | 2.4 | **0.26** |
| Dates-100K | -17.1 | -16.7 | **-16.5** | 29 | 566 | **5.4** |

Table 2: Average log likelihood and runtime comparison on literature data sets (best performing methods are shown in bold).

| Data set | Test set log likelihood | | | | Runtime (seconds) | | | |
|---|---|---|---|---|---|---|---|---|
| | Pareto | Hybrid | MiniSPN | LearnSPN | Pareto | Hybrid | MiniSPN | LearnSPN |
| NLTCS | -6.33 | **-6.03** | -6.12 | -6.1 | **4.8** | 35 | **1.4** | 60 |
| MSNBC | -6.54 | -6.4 | -6.61 | **-6.11** | 61 | 212 | **5.6** | 2400 |
| KDDCup | -2.17 | **-2.13** | -2.14 | -2.21 | 152 | 2080 | **23** | 400 |
| Plants | -17.3 | **-13.1** | **-13.2** | **-13** | 28 | 780 | **11** | 160 |
| Audio | -41.9 | **-39.9** | **-40** | -40.5 | 28 | 556 | **12** | 955 |
| Jester | -54.6 | **-52.9** | **-53** | -53.4 | 13 | 193 | **6.7** | 1190 |
| Netflix | -59.5 | **-56.7** | **-56.8** | -57.3 | 27 | 766 | **14** | 1230 |
| Accidents | -40.4 | -32.5 | -32.6 | **-30.3** | 31 | 1140 | **18** | 330 |
| Retail | -11.1 | **-11** | -11.1 | **-11.09** | 25 | 63 | **7.3** | 100 |
| Pumsb-star | -40.8 | -28.4 | -28.3 | **-25** | 47 | 1100 | **22** | 350 |
| DNA | -98.1 | -91.5 | -93.9 | **-89** | **6.3** | 45 | **3** | 300 |
| Kosarek | -11.2 | **-10.8** | **-10.9** | **-11** | 90 | 537 | **22** | 200 |
| MSWeb | -10.7 | **-9.94** | **-10.1** | -10.26 | 75 | 572 | **34** | 260 |
| Book | -35.1 | **-34.7** | **-34.7** | -36.4 | 83 | 181 | **32** | 350 |
| EachMovie | -55 | **-52.3** | **-52.2** | -52.5 | 62 | 218 | **22** | 220 |
| WebKB | -161 | **-155** | **-155** | -162 | 37 | 169 | **38** | 900 |
| Reuters-52 | -92 | **-85.2** | **-84.7** | -86.5 | 76 | 656 | 95 | 2900 |
| Newsgroup | -156 | **-152** | **-152** | -160.5 | 181 | 1190 | **139** | 28000 |
| BBC | -258 | **-250** | **-249** | **-250** | **33** | 123 | 42 | 900 |
| Ad | -52.3 | -49.5 | -49.2 | **-22** | **58** | 92 | **50** | 300 |

with around 95% of the data missing. We use subsets of 10000 and 100000 instances from each of these data sets, and randomly split the data sets into a training and test set.

On the KG data sets, we compare MiniSPN, Pareto and Hybrid algorithms. We were not able to apply the standard LearnSPN algorithm on these data sets, since they contain missing data. Table 1 shows log likelihood performance on the test set and runtime performance. MiniSPN does better than Pareto, both in terms of log likelihood and runtime. Hybrid performs comparably to MiniSPN, but is usually the slowest of the three.

We use 20 benchmark data sets from the literature (exactly the same ones used in the LearnSPN paper [Gens & Domingos, 2013]) to compare the performance of MiniSPN with the standard LearnSPN algorithm. We are particularly interested in the effect of MiniSPN's simple two-cluster instance split relative to the more complex instance split with the exponential prior and EM restarts used in the standard LearnSPN. Table 2 shows log likelihood performance on the test set and runtime performance. Like on the KG data, we find that MiniSPN uniformly outperforms Pareto, and performs similarly to Hybrid and LearnSPN but runs much faster (on the most time-intensive data set, Newsgroup, MiniSPN takes 2 minutes while LearnSPN takes 8 hours).

## 4    CONCLUSION

Sum-product networks have been receiving increasing attention from researchers due to their expressiveness, efficient inference and interpretability, and many learning algorithms have been developed in the past few years. While recent developments have mostly focused on improving performance on benchmark data sets, our variation on a classical learning algorithm is simple yet has a large impact on usability, by improving speed and making it possible to apply to messy real data sets.

## REFERENCES

Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pp. 1247–1250. ACM, 2008.

Robert Gens and Pedro M. Domingos. Learning the Structure of Sum-Product Networks. In *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*, pp. 873–880, 2013.

Roger B. Grosse, Ruslan Salakhutdinov, William. T. Freeman, and Joshua B. Tenenbaum. Exploiting compositionality to explore a large space of model structures. In *Proceedings of the 28th Conference on Uncertainty in AI (UAI)*, 2012.

Hoifung Poon and Pedro M. Domingos. Sum-Product Networks: A New Deep Architecture. In *UAI 2011, Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence, Barcelona, Spain, July 14-17, 2011*, pp. 337–346, 2011.

Amirmohammad Rooshenas and Daniel Lowd. Learning Sum-Product Networks with Direct and Indirect Variable Interactions. In Tony Jebara and Eric P. Xing (eds.), *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pp. 710–718. JMLR Workshop and Conference Proceedings, 2014.

Antonio Vergari, Nicola Di Mauro, and Floriana Esposito. Simplifying, Regularizing and Strengthening Sum-Product Network Structure Learning. In *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2015, Porto, Portugal, September 7-11, 2015, Proceedings, Part II*, pp. 343–358, 2015.

Han Zhao, Mazen Melibari, and Pascal Poupart. On the Relationship between Sum-Product Networks and Bayesian Networks. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pp. 116–124, 2015.