

The aitoools library

Wieger Wesselink

October 3, 2024

1 Introduction

The aitoools library is a C++ library that contains a few data structures and algorithms related to AI, see <https://github.com/wiegerw/aitools>. It currently supports binary decision trees, random forests, probabilistic circuits, generative forests and some algorithms. It was originally intended as a companion library for the paper *Joints in Random Forests*, see [2], but it was never used as such.

This document contains precise mathematical specifications of the data structures and algorithms that are used in the aitoools library. This is not a tutorial, so it is assumed that the reader is already familiar with the material. In Appendix A some notation used in the algorithms is explained.

2 Graphs

Let $G = (V, E)$ with $E \subseteq V \times V$ be a directed graph. We write $u \rightarrow v$ whenever $(u, v) \in E$. Let \rightarrow^* be the reflexive-transitive closure of \rightarrow . We define

$$\text{pred}(v) = \{u \in V \mid u \rightarrow v\}$$

$$\text{succ}(u) = \{v \in V \mid u \rightarrow v\}$$

$$\text{desc}(u) = \{v \in V \mid u \rightarrow^* v\}$$

A node u is called a leaf or terminal node if $\text{succ}(u) = \emptyset$.

A *topological ordering* of an acyclic graph G is a linear ordering of the vertices V such that for every directed edge $(u, v) \in E$ we have that u comes before v in the ordering. In the context of probabilistic circuits, this ordering is also known as *feedforward order*.

3 Random variables

Let $\mathbf{X} = \{X_1, \dots, X_m\}$ be a set of random variables or features. We assume that each continuous variable X_i assumes values in some compact set $\mathcal{X}_i \subseteq \mathbb{R}$ and each discrete variable X_i assumes values in $\mathcal{X}_i = \{1, \dots, K_i\}$, where K_i is the number of states for X_i . The feature space is denoted as

$$\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_m.$$

Let $\mathcal{D} = \{x_1, \dots, x_n\}$ be a data set for the variables \mathbf{X} . Elements of \mathcal{D} may have missing values, which we denote with the symbol \perp . Hence

$$x_{ij} \in \mathcal{X}_i \cup \{\perp\} \quad (1 \leq i \leq m, 1 \leq j \leq n).$$

In classification problems, there is also an output variable Y (the *class variable*) that assumes values in $\mathcal{Y} = \{1, \dots, K\}$. A data set for a classification problem consists of pairs of inputs and outputs: $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$ with $y_i \in \mathcal{Y}$.

We define the function `project` that projects a data set \mathcal{D} on coordinate i as follows:

$$\text{project}(\mathcal{D}, i) = \{x_i \mid x \in \mathcal{D}\}.$$

We define the function `ncat` that operates on random variables $X_i \in \mathbf{X}$ as follows:

$$\text{ncat}(X_i, \mathcal{D}) = |\text{project}(\mathcal{D}, i)|,$$

In particular we have

$$\text{ncat}(X_i) = \text{ncat}(X_i, \mathcal{X}_i).$$

Furthermore we define the function `counts` that counts the number of samples in class k as:

$$\text{count}(\mathcal{D}, k) = |\{(x_i, y_i) \in \mathcal{D} \mid y_i = k\}|.$$

3.1 The Kendall rank correlation test

Let $(x_1, y_1), \dots, (x_n, y_n)$ be a set of observations of the joint random variables X and Y , such that all the values of (x_i) and (y_i) are unique (the ties are neglected for simplicity). Any pair of observations (x_i, y_i) and (x_j, y_j) , where $i < j$, are said to be concordant if the sort order of (x_i, x_j) and (y_i, y_j) agrees: that is, if either both $x_i > x_j$ and $y_i > y_j$ holds or both $x_i < x_j$ and $y_i < y_j$; otherwise they are said to be discordant.

The Kendall τ coefficient is defined as:

$$\tau = \frac{2}{n(n-1)} \sum_{i < j} \text{sgn}(x_i - x_j) \text{sgn}(y_i - y_j),$$

where

$$\text{sgn}(x) = \begin{cases} -1, & x < 0 \\ 0, & x = 0 \\ 1, & x > 0 \end{cases}$$

Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be a sample of n pairs of observations. The rank correlation coefficient τ of Kendall is defined as

$$\tau = 1 - \frac{2K_n}{\binom{n}{2}},$$

where K_n is the number of inversions: the number of pairs $\{(X_i, Y_i), (X_j, Y_j)\}$ such that $X_i < X_j$ and $Y_i > Y_j$ for $i < j$, $i = 1, \dots, n-1$ and $j = 2, \dots, n$.

4 Probability distributions

4.1 Uniform distribution

A continuous uniform distribution on an interval $[a, b]$ is denoted as $U(a, b)$, and has the following probability density function:

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{for } a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$

and cumulative density function

$$F(x) = \begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } a \leq x \leq b \\ 1 & \text{for } x > b \end{cases}$$

4.2 Normal distribution

A normal or Gaussian distribution $\mathcal{N}(\mu, \sigma)$ is a distribution with the following probability density function:

$$\phi(\mu, \sigma; x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

and cumulative distribution function

$$\Phi(\mu, \sigma; x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt,$$

The parameter μ is the mean, and σ is the standard deviation of the normal distribution. The function Φ can be expressed in terms of the error function as follows:

$$\Phi(\mu, \sigma; x) = \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{(x-\mu)/\sigma}{\sqrt{2}}\right) \right).$$

$$\Phi^{-1}(\mu, \sigma; x) = \mu + \sigma\sqrt{2} \operatorname{erf}^{-1}(2x - 1).$$

Sampling a normal distribution can be done by first randomly drawing a value $p \sim U(0, 1)$ and then transforming it to a value x using

$$x = \mu + \sigma \cdot \Phi^{-1}(0, 1; p) = \mu + \sigma \cdot \sqrt{2} \operatorname{erf}^{-1}(2p - 1).$$

4.3 Standard normal distribution

The standard normal distribution is $\mathcal{N}(0, 1)$. It has the following cumulative distribution function:

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt,$$

which is closely related to the error function erf :

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_{-\infty}^x e^{-t^2} dt.$$

We have

$$\Phi^{-1}(x) = \sqrt{2} \operatorname{erf}^{-1}(2x - 1).$$

4.4 Truncated normal distribution

A truncated normal distribution is obtained from a normal distribution $\mathcal{N}(\mu, \sigma)$ by truncating it to an interval $[a, b]$ and then scaling it appropriately, see also [4]. It has the following probability density function:

$$\psi(\mu, \sigma, a, b; x) = \begin{cases} 0 & \text{if } x < a \\ \frac{\phi(\mu, \sigma; x)}{\Phi(\mu, \sigma; b) - \Phi(\mu, \sigma; a)} & \text{if } a \leq x \leq b \\ 0 & \text{if } x > b, \end{cases}$$

where μ and σ are the parameters of the original distribution. The cumulative density function is given by

$$\Psi(\mu, \sigma, a, b; x) = \begin{cases} 0 & \text{if } x < a \\ \frac{\Phi(\mu, \sigma; x) - \Phi(\mu, \sigma; a)}{\Phi(\mu, \sigma; b) - \Phi(\mu, \sigma; a)} & \text{if } a \leq x \leq b \\ 1 & \text{if } x > b. \end{cases}$$

and the inverse by

$$\Psi^{-1}(\mu, \sigma, a, b; x) = \Phi^{-1}(\mu, \sigma; \Phi(\mu, \sigma; a) + x \cdot (\Phi(\mu, \sigma; b) - \Phi(\mu, \sigma; a)))$$

Sampling a truncated normal distribution can be done by first drawing a value $p \sim U(0, 1)$ and then transforming it to a value x using

$$x = \Psi^{-1}(\mu, \sigma, a, b; p).$$

4.5 Bernoulli distribution

A categorical distribution is a discrete probability distribution that takes the value 1 with probability p and the value 0 with probability $1 - p$. The probability mass function is given by

$$f(x = i) = p_i.$$

4.6 Categorical distribution

A categorical distribution is a discrete probability distribution that models the possible results of a random variable that can take on one of K possible categories, with the probability of each category separately specified. It is characterized by k probabilities $[p_1, \dots, p_k]$ with $p_i \geq 0$ and $\sum_{i=1}^k p_i = 1$. The probability mass function is given by

$$f(x = i) = \begin{cases} p & \text{if } i = 1 \\ 1 - p & \text{if } i = 0 \end{cases}$$

4.7 Multinomial distribution

A multinomial distribution is a discrete probability distribution that models the probability of counts for each side of a k -sided die rolled n times. It is characterized by k probabilities $[p_1, \dots, p_k]$ with $p_i \geq 0$ and $\sum_{i=1}^k p_i = 1$, and a number of trials n . The probability mass function is given by

$$f(x_1, \dots, x_n) = \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}.$$

5 Decision trees

A *decision tree* (DT) is a tree $G = (V, E)$ that is used to partition the feature space \mathcal{X} into a number of disjoint subsets. We denote the root of the tree as $\text{root}(G)$. Each node u of the tree is associated with a subset $\mathcal{X}_u \subseteq \mathcal{X}$ of the feature space. We have $\mathcal{X}_{\text{root}(G)} = \mathcal{X}$. Each non-terminal node u is labeled with a *decision tree classifier split*(u) that partitions \mathcal{X} into disjoint subsets $\{U_1, \dots, U_p\}$, and such that $\mathcal{X}_{v_i} = \mathcal{X}_u \cap U_i$, for $v_i \in \text{succ}(u)$. In particular we have that the leaf nodes of G define a partition of the feature space: $\mathcal{X} = \cup \{\mathcal{X}_v \mid v \in V \wedge \text{succ}(v) = \emptyset\}$.

A decision tree is usually defined over a data set \mathcal{D} . Each node u is associated with the subset $\mathcal{D}_u = \mathcal{X}_u \cap \mathcal{D}$. We define $\text{count}(u) = |\mathcal{D}_u|$.

5.1 Decision tree classifiers

This section gives an overview of the supported decision tree classifiers. Currently, all of them are binary splits, i.e. they partition the feature space in two subsets. For each classifier, a function `partition-number` is defined that assigns a partition to each element x of the feature space. In case of a binary split, the partition numbers are 1 and 2. Furthermore, a function `domain` is defined that specifies for each variable the part of the domain where it can be nonzero. For the root u of the decision tree, we have $\text{domain}(u, X_i) = \mathcal{X}_i$ for each variable $X_i \in \mathbf{X}$.

Single split

A *single split* $\text{SingleSplit}(X_i, v)$ is an axis-aligned split defined by a discrete random variable X_i with domain \mathcal{X}_i and a value $v \in \mathcal{X}_i$. We define

$$\text{partition-number}(\text{SingleSplit}(X_i, v), x) = \begin{cases} 1 & \text{if } x_i = v \\ 2 & \text{otherwise.} \end{cases}$$

$$\begin{aligned} \text{domain}(v_1, X_i) &= \{v\} \\ \text{domain}(v_2, X_i) &= \text{domain}(u, X_i) \setminus \{v\}, \end{aligned}$$

where $\text{succ}(u) = \{v_1, v_2\}$. For other variables X_j ($j \neq i$), we have $\text{domain}(v_1, X_j) = \text{domain}(v_2, X_j) = \text{domain}(u, X_j)$.

Subset split

A *subset split* $\text{SubsetSplit}(X_i, V)$ is an axis-aligned split that is defined by a discrete random variable X_i with domain \mathcal{X}_i and a set of values $V \subseteq \mathcal{X}_i$.

$$\text{partition-number}(\text{SubsetSplit}(X_i, V), x) = \begin{cases} 1 & \text{if } x_i \in V \\ 2 & \text{otherwise.} \end{cases}$$

$$\begin{aligned} \text{domain}(v_1, X_i) &= \text{domain}(u, X_i) \cap V \\ \text{domain}(v_2, X_i) &= \text{domain}(u, X_i) \setminus V, \end{aligned}$$

where $\text{succ}(u) = \{v_1, v_2\}$. For other variables X_j ($j \neq i$), we have $\text{domain}(v_1, X_j) = \text{domain}(v_2, X_j) = \text{domain}(u, X_j)$.

Threshold split

A *threshold split* $\text{ThresholdSplit}(X_i, v)$ is an axis-aligned split defined by a continuous random variable X_i with domain \mathcal{X}_i and value $v \in \mathcal{X}_i$. We define

$$\text{partition-number}(\text{ThresholdSplit}(X_i, v), x) = \begin{cases} 1 & \text{if } x_i < v \\ 2 & \text{otherwise.} \end{cases}$$

$$\begin{aligned} \text{domain}(v_1, X_i) &= \text{domain}(u, X_i) \cap (-\infty, v) \\ \text{domain}(v_2, X_i) &= \text{domain}(u, X_i) \cap [v, +\infty), \end{aligned}$$

where $\text{succ}(u) = \{v_1, v_2\}$. For other variables X_j ($j \neq i$), we have $\text{domain}(v_1, X_j) = \text{domain}(v_2, X_j) = \text{domain}(u, X_j)$.

Applying a splitting criterion to a data set

When constructing a decision tree, a splitting criterion is used to partition a data set. Let \mathcal{D} be a data set, and let u be a decision tree node with scope X_i , splitting criterion split , and with successors $\text{succ}(u) = \{v_1, \dots, v_p\}$. Then

$$\text{apply-split}(\text{split}, \mathcal{D}) = [\{x \in \mathcal{D} \mid \text{partition-number}(\text{split}, x) = 1\}, \dots, \{x \in \mathcal{D} \mid \text{partition-number}(\text{split}, x) = p\}].$$

Indicator functions

Each decision tree classifier *split* associated with node u naturally defines an indicator function for each outgoing edge. We define

$$\text{indicator}(u, j) = \mathbb{1}_{\mathcal{X}_j},$$

where

$$\mathcal{X}_j = \{x \in \mathcal{X} \mid \text{partition-number}(\text{split}, x) = j\}.$$

Missing values

When the data set contains missing values for variable X_i , the following approach is taken. When computing the gain of a split, samples x with $x_i = \perp$ are simply ignored. When partitioning a data set using a split, the samples with $x_i = \perp$ are randomly assigned to one of the partitions.

5.2 Executing a decision tree

Let $x \in \mathcal{X}$ be an arbitrary element of the feature space. The algorithms below determines a leaf node to which x corresponds.

Algorithm 1 Executing a decision tree

Input: A decision tree $G = (V, E)$ and a sample $x \in \mathcal{X}$

Output: A leaf node $v \in V$.

```

1: function PREDICT( $G = (V, E), x$ )
2:    $u := \text{root}(G)$ 
3:   let  $\text{succ}(u) = \{v_1, \dots, v_p\}$ 
4:   while true do
5:     if  $\text{succ}(u) = \emptyset$  then
6:       break
7:      $j := \text{partition-number}(\text{split}(u), x)$   $\triangleright$   $\text{split}(u)$  is the split criterion of node  $u$ 
8:      $u := v_j$ 
9:   return  $u$ 
```

5.3 Learning a decision tree

In this section we describe an algorithm for learning a decision tree for a classification problem with data set $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$ that is defined on random variables $\mathbf{X} = \{X_1, \dots, X_m\}$ and class variable Y .

Impurity measures

There are several *impurity measures* of a data set \mathcal{D} that can be used to select a suitable split. The *Gini index* and the *cross-entropy* are two impurity measures on a data set that are defined as follows (see also [3]):

$$\text{gini-index}(\mathcal{D}) = \sum_{k=1}^K p_k(1 - p_k) = 1 - \sum_{k=1}^K p_k^2$$

$$\text{cross-entropy}(\mathcal{D}) = - \sum_{k=1}^K p_k \log(p_k)$$

where $p_k = \text{count}(\mathcal{D}, k) / |\mathcal{D}|$, i.e. the proportion of samples of class k in \mathcal{D} . Furthermore we define

$$\text{purity}(\mathcal{D}) = \max_{1 \leq k \leq K} p_k.$$

Choosing a best split

Suppose we have a data set \mathcal{D} , and a splitter that splits \mathcal{D} into $\{\mathcal{D}_1, \dots, \mathcal{D}_p\}$. Given an impurity measure imp , the quality of the split is then measured by

$$\text{gain}(\mathcal{D}, \{\mathcal{D}_1, \dots, \mathcal{D}_p\}, imp) = imp(\mathcal{D}) - \sum_{i=1}^p \frac{|\mathcal{D}_i|}{|\mathcal{D}|} imp(\mathcal{D}_i).$$

A splitter with the maximum gain is selected. In practice we use the simplified expression below to determine the maximum:

$$\text{gain}_1(\{\mathcal{D}_1, \dots, \mathcal{D}_p\}, imp) = - \sum_{i=1}^p |\mathcal{D}_i| imp(\mathcal{D}_i).$$

This function gives the same results, but is more efficient to compute.

Algorithm 2 Learning a decision tree

Input: A training set $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$ that is defined on random variables $\mathbf{X} = \{X_1, \dots, X_m\}$ and class variable Y ; the maximum number of split variables M , ($1 \leq M \leq m$); a family of split functions $Splits$; a function $gain$ to measure the quality of a split; a criterion $stop$ for ending the recursion.

Output: A decision tree (V, E) .

```

1: function LEARNDECISIONTREE( $\mathcal{D}, M, Splits, gain, stop$ )
2:   choose  $u_0 \in \mathcal{V}$   $\triangleright u_0$  is a fresh node, taken from a universal set of nodes  $\mathcal{V}$ 
3:    $V := \{u_0\}$ 
4:    $E := \emptyset$ 
5:    $\mathcal{D}_{u_0} := \mathcal{D}$ 
6:    $todo := \{u_0\}$ 
7:   while  $todo \neq \emptyset$  do
8:     choose  $u \in todo$ 
9:      $todo := todo \setminus \{u\}$ 
10:    if  $stop(\mathcal{D}_u)$  then
11:      continue
12:     $\mathbf{Z} := \text{random-sample}(\mathbf{X}, M)$   $\triangleright \mathbf{Z} \subseteq \mathbf{X}$  is a random subset of size  $M$ 
13:     $A := \arg \max_{split \in Splits(\mathcal{D}_u, \mathbf{Z})} gain(\text{apply-split}(split, \mathcal{D}_u))$ 
14:    if  $A = \emptyset$  then continue  $\triangleright A = \emptyset$  means that no suitable split is found
15:    choose  $split \in A$ 
16:    for  $D \in \text{apply-split}(split, \mathcal{D}_u)$  do
17:      choose  $v \in \mathcal{V} \setminus V$   $\triangleright v$  is a fresh node
18:       $V := V \cup \{v\}$ 
19:       $E := E \cup \{(u, v)\}$ 
20:       $todo := todo \cup \{v\}$ 
21:       $\mathcal{D}_v := D$ 
22:  return  $(V, E)$ 

```

Split families

The family of split functions $Splits(\mathcal{D}, \mathbf{Z})$ that is considered can for example be chosen as follows:

$$Splits(\mathcal{D}, \mathbf{Z}) = \{SingleSplit(X_i, v) \mid X_i \in \mathbf{Z} \wedge \text{ncat}(X_i) \leq 5 \wedge v \in \text{project}(\mathcal{D}, i)\} \cup \{ThresholdSplit(X_i, v) \mid X_i \in \mathbf{Z} \wedge \text{ncat}(X_i) > 5 \wedge v \in \text{project}(\mathcal{D}, i)\}.$$

Stop criterion

A possible stop criterion is

$$stop(\mathcal{D}_u) = |\mathcal{D}_u| \leq \text{min-samples-leaf} \vee \text{mis-classification}(\mathcal{D}_u) \leq 0.01 \vee \text{depth}(u) > \text{max-depth}$$

where min-samples-leaf is a user defined constant and $\text{depth}(u)$ is the distance between the root of the decision tree and u .

6 Random forests

A random forest is a union of decision trees $\{G_i = (V_i, E_i)\}_{i \in I}$ that are all defined on the same feature space \mathcal{X} . The trees are disjoint, i.e. $V_i \cap V_j = \emptyset$ for $i \neq j$.

6.1 Learning a random forest

The following algorithm is used to learn a random forest from a training set $\mathcal{D} = \{x_1, \dots, x_n\}$ that is defined on random variables $\mathbf{X} = \{X_1, \dots, X_m\}$.

Algorithm 3 Learning a random forest

Input: A training set $\mathcal{D} = \{x_1, \dots, x_n\}$ that is defined on random variables $\mathbf{X} = \{X_1, \dots, X_m\}$; a function *sample* that draws a sample of \mathcal{D} ; a fraction $p \in [0, 1]$ that determines the size of the sample; the number d of decision trees in the forest; the maximum number of split variables M , ($1 \leq M \leq m$); a family of split functions *Splits*; a function *gain* to measure the quality of a split; a criterion *stop* for ending the recursion.

Output: A random forest

```
1: function LEARNRANDOMFOREST( $\mathcal{D}, \mathbf{X}, \text{sample}, p, d, M, \text{Splits}, \text{gain}, \text{stop}$ )
2:   result :=  $\emptyset$ 
3:   for  $1 \leq i \leq d$  do
4:      $\mathcal{D}' := \text{sample}(\mathcal{D}, p)$   $\triangleright |\mathcal{D}'| \simeq p|\mathcal{D}|$ 
5:     result := result  $\cup$  LEARNDECISIONTREE( $\mathcal{D}', M, \text{Splits}, \text{gain}, \text{stop}$ )
6:   return result
```

7 Probabilistic circuits

A *probabilistic circuit* (PC) is a rooted directed acyclic graph $G = (V, E)$. Each terminal node (also called *distribution node*) represents a multivariate probabilistic distribution. We denote the root of the PC G as $\text{root}(G)$. There are two types of non-terminal nodes: sum nodes and product nodes, and there are many different kinds of terminal nodes, see also [1]. For each node $u \in V$ an evaluation function $\text{evi}(u, x)$ is defined that is used to compute the probability of the node, where $x \in \mathcal{X}$ is a given sample.

7.1 Definitions

We denote the set of random variables that correspond to terminal node v as $\text{scope}(v)$. This definition is extended to a non-terminal node u using

$$\text{scope}(u) = \cup \{ \text{scope}(v) \mid v \in \text{desc}(u) \wedge \text{succ}(v) = \emptyset \}.$$

An *unnormalized distribution* is any nonnegative function $\Phi(x)$ where $\exists x : \Phi(x) > 0$.

For each random variable X_i and each $t \in \mathcal{X}_i$ we define an *indicator variable* $\lambda_{X_i=t}$, which is a function on \mathcal{X} defined as

$$\lambda_{X_i=t}(x) = \begin{cases} 1 & \text{if } x_i = t \\ 0 & \text{otherwise.} \end{cases}$$

Definition 1 (*Network Polynomial*) Let Φ be an unnormalized probability distribution over random variables \mathbf{X} with finitely many states and λ their indicator variables. The network polynomial f_Φ of Φ is defined as

$$f_\Phi(\lambda) := \sum_{x \in \mathcal{X}} \Phi(x) \prod_{\lambda_{X_i=t} \in \lambda} \lambda_{X_i=t}.$$

Definition 2 (*SPN Distribution*) Let S be a Sum-product network over \mathbf{X} . The distribution represented by S is defined as

$$P_S(x) := \frac{\text{evi}(\text{root}(S), x)}{\sum_{x' \in \mathcal{X}} \text{evi}(\text{root}(S), x')}.$$

Definition 3 (*Smoothness a.k.a. completeness*) A sum node u is called smooth if its children have the same scope: $\text{scope}(v) = \text{scope}(w)$ for any $v, w \in \text{succ}(u)$. A PC \mathcal{P} is called smooth if every sum node in \mathcal{P} is smooth.

Definition 4 (*Consistency*) A product node u is called consistent if $\forall v, w \in \text{succ}(u), v \neq w$ it holds that $\lambda_{X=x} \in \text{desc}(v) \Rightarrow \forall x' \neq x : \lambda_{X=x'} \notin \text{desc}(w)$.

Definition 5 (*Decomposability*) A product node u is called decomposable if its children have non-overlapping scopes: $\text{scope}(v) \cap \text{scope}(w) = \emptyset$ for any $v, w \in \text{succ}(u)$ with $v \neq w$. A PC is called decomposable if all its product nodes are decomposable.

Definition 6 (*Shared random variables*) A random variable X is a shared random variable of product node u if there are $v, w \in \text{succ}(u)$ with $v \neq w$ such that $X \in \text{scope}(v) \cap \text{scope}(w)$.

Definition 7 (*Locally normalized*) A sum node u is called locally normalized if the weights of its outgoing edges sum to 1. A PC \mathcal{P} is called locally normalized if every sum node in \mathcal{P} is locally normalized.

Definition 8 (*Deterministic*) A PC is called deterministic if it holds that for each complete sample $x \in \mathcal{X}$, each sum node has at most one non-zero child.

7.2 Nodes in a probabilistic circuit

This section gives an overview of the nodes in a PC that are supported. We use the convention that for every node u the function evi returns the value 1 for missing data (denoted as \perp):

$$\text{evi}(u, \perp) = 1.$$

Sum nodes

A sum node $u = \text{Sum}([w_1, \dots, w_p])$ is a non-terminal node with $[w_1, \dots, w_p]$ the edge weights of the outgoing edges. Let $\text{succ}(u) = \{v_1, \dots, v_p\}$. We have $0 \leq w_i$ ($1 \leq i \leq p$).

$$\text{evi}(u, x) = \sum_{i=1}^m w_i \text{evi}(v_i, x)$$

Sum split nodes

A sum split node $u = \text{SumSplit}([w_1, \dots, w_p], \text{split})$ is an extension of a sum node with a decision tree classifier *split*. Note that this is not a standard node of probabilistic circuits, but it is used by generative forests that are discussed in section 8. The decision tree qualifier *split* is used to ensure that only one of the successors of u can have a non-zero evaluation:

$$\text{evi}(u, x) = w_j \cdot \text{evi}(v_j, x) \quad \text{with } j = \text{partition-number}(\text{split}, x).$$

Product nodes

A product node $u = \text{Product}()$ is a non-terminal node. Let $\text{succ}(u) = \{v_1, \dots, v_p\}$. We have

$$\text{evi}(u, x) = \prod_{i=1}^m \text{evi}(v_i, x)$$

Normal nodes

A normal node $u = \text{Normal}(X_i, \mu, \sigma)$ is a terminal node that models a normal distribution $\mathcal{N}(\mu, \sigma)$ for random variable X_i . We have

$$\text{evi}(u, x) = \phi(\mu, \sigma; x_i).$$

Truncated normal nodes

A truncated normal node $u = \text{TruncatedNormal}(X_i, \mu, \sigma, a, b)$ is a terminal node that models a truncated normal distribution $\mathcal{N}(\mu, \sigma)$ on an interval $[a, b]$ for random variable X_i . We have

$$\text{evi}(u, x) = \psi(\mu, \sigma, a, b; x_i).$$

Categorical nodes

A categorical node $u = \text{Categorical}(X_i, [p_1, \dots, p_k])$ is a terminal node that models a categorical distribution with probabilities $[p_1, \dots, p_k]$ for random variable X_i . It is only defined for integer values $x \in [1, \dots, k]$. We have

$$\text{evi}(u, x) = p_{x_i}$$

Equal nodes

An equal node $u = \text{Equal}(X_i, t)$ is a degenerate case of a categorical distribution with probability 1 for a single integer value $i \in \mathbb{N}$ for random variable X_i . It should only be evaluated for integer values $x \in \mathbb{N}$.

$$\text{evi}(u, x) = \begin{cases} 1 & \text{if } x_i = t \\ 0 & \text{otherwise} \end{cases}$$

Not equal nodes

A not-equal node $u = \text{NotEqual}(X_i, t)$ is a degenerate case of a categorical distribution with probability 1 for integer values not equal to a value $i \in \mathbb{N}$ for random variable X_i . It should only be evaluated for integer values $x \in \mathbb{N}$.

$$\text{evi}(u, x) = \begin{cases} 1 & \text{if } x_i \neq t \\ 0 & \text{otherwise} \end{cases}$$

Less nodes

A less node $u = \text{Less}(X_i, t)$ evaluates to 1 for values less than a given threshold t for random variable X_i .

$$\text{evi}(u, x) = \begin{cases} 1 & \text{if } x_i < t \\ 0 & \text{otherwise} \end{cases}$$

Greater-equal nodes

A greater node $u = \text{GreaterEqual}(X_i, t)$ evaluates to 1 for values greater than or equal to a given threshold t for random variable X_i .

$$\text{evi}(u, x) = \begin{cases} 1 & \text{if } x_i \geq t \\ 0 & \text{otherwise} \end{cases}$$

Subset nodes

A subset node $u = \text{Subset}(X_i, V)$ evaluates to 1 for values in a set V for random variable X_i .

$$\text{evi}(u, x) = \begin{cases} 1 & \text{if } x_i \in V \\ 0 & \text{otherwise} \end{cases}$$

7.3 Inference

Algorithm 4 Iterative computation of an EVI query

Input: A probabilistic circuit $G = (V, E)$ over random variables \mathbf{X} and a complete state $x \in \mathcal{X}$.

Output: A probability.

```

1: function EVI-ITERATIVE( $G, x$ )
2:    $U := \text{topological-ordering}(G)$ 
3:    $c := \{\mapsto\}$   $\triangleright c \subset U \times \mathbb{R}$  is an empty mapping
4:   for  $u \in U$  do
5:     let  $\text{succ}(u) = [v_1, \dots, v_p]$ 
6:     if  $u = \text{Sum}([w_1, \dots, w_p])$  then
7:        $c[u] := \sum_{j=1}^p w_j \cdot c[v_j]$ 
8:     else if  $u = \text{SumSplit}([w_1, \dots, w_p], \text{split})$  then
9:        $j := \text{partition-number}(\text{split}, x)$ 
10:       $c[u] := w_j \cdot c[v_j]$ 
11:     else if  $u = \text{Product}()$  then
12:        $c[u] := \prod_{j=1}^p c[v_j]$ 
13:     else  $\triangleright u$  is a terminal node
14:        $c[u] := \text{evi}(u, x)$ 
15:   return  $c[U[-1]]$   $\triangleright$  The last node  $U[-1]$  is the root of the PC

```

Algorithm 5 Recursive computation of an EVI query

Input: A probabilistic circuit $G = (V, E)$ over random variables \mathbf{X} and a complete state $x \in \mathcal{X}$.

Output: A boolean value.

```

1: function EVI-RECURSIVE( $G, x$ )
2:   return  $\text{evi}(\text{root}(G), x)$ 

```

7.4 Sampling

We start at the root node. If we have a sum node, we pick one of the children by sampling from the categorical distribution defined by the weights of the sum node—intuitively, we will pick the child with largest weight more often—and sample from it. If we have a product node, we sample from all of its children. If we have a distribution node, we just sample from its distribution and write the value into the corresponding variable in our sample.

In this algorithm we start of with an empty vector \mathbf{x} , and write values to it once we reach a leaf. Note that we sample each variable only once, so the set of scopes of the leaves the function reaches will form a partition of the scope.

Algorithm 6 Sampling

Input: The root node u of the PC one wants to sample from, a vector \mathbf{x} where we store the sample.

Output: A sample x .

```

1: function SAMPLE( $u, \mathbf{x}$ )
2:   let  $\text{succ}(u) = [v_1, \dots, v_p]$ 
3:   if  $u = \text{Sum}([w_1, \dots, w_p])$  then ▷ If sum node.
4:      $j \sim \text{Categorical}(\cdot | [w_1, \dots, w_p])$  ▷ Choose a child according to distribution given by the weights.
5:     SAMPLE( $v_j, \mathbf{x}$ ) ▷ Sample from the chosen child.
6:   else if  $u = \text{Product}()$  then ▷ If product node.
7:     for  $v_j \in \text{succ}(u)$  do ▷ Sample from every child.
8:       SAMPLE( $v_j, \mathbf{x}$ )
9:   else ▷ A leaf node with a distribution parameterised by  $\theta$ .
10:     $\mathbf{x}[\text{scope}(u)] \sim \text{Distribution}(\cdot | \theta)$ 

```

7.5 Learning a sum-product network

Algorithm 7 Learning a Probabilistic Circuit with LearnSPN

Input: A training set $\mathcal{D} = \{x_1, \dots, x_n\}$ that is defined on a set random variables \mathcal{X} ; A function *isplit* that splits the set of variables \mathcal{X} into independent sets; A function *cluster* that splits the instances in \mathcal{D} into K clusters.

Output: A Probabilistic Circuit (V, E) .

```

1: function LEARNSPN( $\mathcal{D}, \mathcal{X}, \text{isplit}, \text{cluster}, K$ )
2:   if  $|\mathcal{X}| = 1$  then ▷ we have a single variable in the scope.
3:     return univariate distribution node fitted on  $\mathcal{D}$ .
4:   else
5:      $\mathcal{X}_1, \dots, \mathcal{X}_J = \text{isplit}(\mathcal{D}, \mathcal{X})$ .
6:     if  $J > 1$  then ▷ we found at least one independence relationship.
7:       return  $\prod_j \text{LearnSPN}(\mathcal{D}, \mathcal{X}_j, \text{isplit}, \text{cluster}, K)$  ▷ add a prod. node with one child for each  $\mathcal{X}_j$ .
8:     else
9:        $\mathcal{D}_1, \dots, \mathcal{D}_K = \text{cluster}(\mathcal{D}, \mathcal{X}, K)$ .
10:    return  $\sum_k \text{LearnSPN}(\mathcal{D}_k, \mathcal{X}, \text{isplit}, \text{cluster}, K)$  ▷ add a sum node with one child for each  $\mathcal{D}_k$ .

```

The *isplit* function needs to identify whether there is any independence relationship in a set of variables. This typically means iterating through each possible pair $(X_i, X_j) \in \mathcal{X}$ and running an independence test: Kendall (X_i and X_j continuous), chi-square (X_i and X_j discrete), or Kruskal (one continuous, one discrete). The *cluster* function is usually just K-means with $K=2$.

8 Generative forests

A generative forest (GeF) is a probabilistic circuit that is derived from a random forest. The structure of the underlying forest is preserved, and also the decision tree classifiers (splitters) are encoded in the PC.

8.1 Converting a random forest to a generative forest

A decision tree $G = (V, E)$ defined over a data set \mathcal{D} can be straightforwardly converted into a generative forest G' as follows, see also [2]. The graph $G' = (V', E')$ contains an isomorphic copy of G . Each non-terminal node $u \in V$ with decision tree classifier *split* is transformed into a sum split node $SumSplit([w_1, \dots, w_p], split)$, with $w_j = |\mathcal{D}_{v_j}|/|\mathcal{D}_u|$ for all $v_j \in \text{succ}(u)$. Each leaf node $u \in V$ is transformed into a small PC that fits a distribution to the samples in \mathcal{D}_u . By default u is transformed into a tree $G_u = (V_u, E_u)$ that consists of a product node with an outgoing edge for each random variable X_i as follows:

$$\begin{aligned} V_u &= \{u, v_1, \dots, v_m\} \\ E_u &= \{(u, v_1), \dots, (u, v_m)\} \\ u &= Product() \\ v_j &= \begin{cases} \text{fit-normal}(u, X_i) & \text{if } X_i \text{ is continuous} \\ \text{fit-categorical}(u, X_i, \alpha) & \text{otherwise} \end{cases} \end{aligned}$$

For a continuous variable X_i we define

$$\text{fit-normal}(u, X_i) = \begin{cases} TruncatedNormal(X_i, \text{mean}(D_i), \text{stddev}(D_i), a, b) & \text{if } |D_i| > 0 \\ TruncatedNormal(X_i, 0, 1, a, b) & \text{otherwise,} \end{cases}$$

where $D_i = \{x_i \mid x \in \mathcal{D}_u\} \setminus \{\perp\}$, and where a and b are chosen such that $\text{domain}(u, X_i) \subseteq [a, b]$.

For a categorical variable X_i and a given Laplace smoothing factor α we define

$$\text{fit-categorical}(u, X_i, \alpha) = Categorical(X_i, [p_1, \dots, p_{\text{ncat}(X_i)}]),$$

where

$$p_k = \frac{|\{x \in \mathcal{D}_u \mid x_i = k\}| + \alpha}{|\{x \in \mathcal{D}_u \mid x_i \neq \perp\}| + \alpha \cdot \text{ncat}(X_i)}.$$

If the number of samples \mathcal{D}_u is larges (say ≥ 30), then the node u can be transformed into a sum-product network using the LearnSPN algorithm, see section 7.5.

A random forest $\{G_i = (V_i, E_i)\}_{i=1}^N$ is converted to a generative forest by first converting each decision tree G_i into a PC G'_i , and then add a new root node $r = Sum([\frac{1}{N}, \dots, \frac{1}{N}])$ with outgoing edges $(r, \text{root}(G'_i))$ for $i = 1 \dots N$.

8.2 Converting a generative forest to a probabilistic circuit

A generative forest can be converted to an equivalent PC by expanding the sum split nodes. Node $u = SumSplit([w_1, \dots, w_p], split)$ is replaced by $u = Sum([w_1, \dots, w_p])$. Each edge (u, v_j) is replaced by a subgraph (V_u, E_u) with

$$\begin{aligned} V_u &= \{u, y_j, z_j, v_j\} \\ E_u &= \{(u, y_j), (y_j, v_j), (y_j, z_j)\} \\ y_j &= Product() \\ z_j &= \text{make-indicator-node}(u, j), \end{aligned}$$

where y_j and z_j are fresh nodes, and

<code>make-indicator-node(<i>SingleSplit</i>(X_i, v), 0)</code>	$=$	$Equal(X_i, v)$
<code>make-indicator-node(<i>SingleSplit</i>(X_i, v), 1)</code>	$=$	$NotEqual(X_i, v)$
<code>make-indicator-node(<i>SubsetSplit</i>(X_i, V), 0)</code>	$=$	$Subset(X_i, V)$
<code>make-indicator-node(<i>SubsetSplit</i>(X_i, V), 1)</code>	$=$	$Subset(X_i, V^C)$
<code>make-indicator-node(<i>ThresholdSplit</i>(X_i, v), 0)</code>	$=$	$Less(X_i, v)$
<code>make-indicator-node(<i>ThresholdSplit</i>(X_i, v), 1)</code>	$=$	$GreaterEqual(X_i, v)$

A Pseudocode

This section describes the pseudocode conventions and data structures that are used in this document.

A.1 Attributes

In many algorithms objects are manipulated that have various attributes. For example a node u in a graph might be assigned the label a , and a common way to denote this is using the statement $u.label := a$. We do not allow this kind of object oriented notation. Instead our way to handle this is to introduce a global mapping $label$, and to write $label[u] := a$ instead. This is done to keep the amount of concepts used in the pseudocode as small as possible. In an actual implementation it may be very inefficient to store these attributes in a separate mapping. But we consider it obvious that the implementer has the freedom to choose an efficient way to store a mapping.

A.2 Lists

Let l and m be lists, and i, j natural numbers. We use the following notations:

Expression	Meaning	Precondition
$[]$	The empty list	
$[a, b, c]$	The list with elements a, b and c	
$ l $	The number of elements in l	
$l[i]$	The element at position i	$0 \leq i < l $
$l[-i]$	The element at position $ l - i$	$0 < i \leq l $
$l[i : j]$	The sublist $[l[i], \dots, l[j - 1]]$	$0 \leq i \leq j < l $
$l[i :]$	$l[i : l - 1]$	$0 \leq i \leq l $
$l[: i]$	$l[0 : i]$	$0 \leq i \leq l $
$l ++ m$	The concatenation of l and m	
$a \in l$	$\exists i : 0 \leq i < l : l[i] = a$	
$\text{index}(l, a)$	The smallest value i such that $l[i] = a$	$a \in l$

Table 1: List operations

A.3 Mappings

We define a mapping as a set of (key, value) pairs. Let $m \subseteq V \times W$ be a mapping, and let $v \in V$ and $w \in W$ be two values. We use the following notations:

Expression	Meaning	Precondition
$\{\mapsto\}$	The empty mapping	
$\text{keys}(m)$	$\{v \in V \mid \exists w \in W : (v, w) \in m\}$	
$\text{values}(m)$	$\{w \in W \mid \exists v \in V : (v, w) \in m\}$	
$m[v]$	The value $w \in W$ such that $(v, w) \in m$	$v \in \text{keys}(m)$
$v \in m$	$v \in \text{keys}(m)$	

Table 2: Map operations

References

- [1] YooJung Choi, Antonio Vergari, and Guy Van den Broeck. Probabilistic circuits: A unifying framework for tractable probabilistic models. oct 2020.
- [2] Alvaro H. C. Correia, Robert Peharz, and Cassio de Campos. Joints in random forests, 2020.
- [3] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer New York, 2009.
- [4] J.Burkardt. The truncated normal distribution. Technical report, Florida State University (FSU), Florida, USA, 2014.