

ALC 2019/2020

3rd Project – Decision Tree Construction with CSP

16 November 2019, version 1.0

Overview

The 3rd ALC project is to develop a software tool for constructing the smallest decision tree [2, 1]. In order to solve this problem, students must use solvers for Constraint Satisfaction Problems (CSP). You may assume that the `minizinc` solver is installed on the machines where the project will be tested.¹

Problem Specification

Given a set of samples, the task is to find a decision tree with the fewest nodes that correctly classifies the samples. All features have only values 0 or 1. Further, samples always fall into one of the classes positive (T) and negative (F).

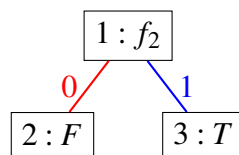


Figure 1: Simple decision tree with the nodes $\{1, 2, 3\}$ deciding only on feature f_2 .

Figure 1 shows a simple decision tree with 3 nodes, where the root node decides on the feature f_2 and responds T if the feature is 1 and responds F if the feature is 0. Such tree correctly classifies for instance the samples positive = $\{(0, 1)\}$, negative = $\{(0, 0)\}$. Observe that the tree can ignore f_1 because its value is always 0 in the given samples.

Project Goals

You are to implement a tool, or optionally a set of tools, invoked with command `proj3`. This set of tools should use an SMT solver to find a smallest decision tree. This means that the tree has the fewest number of nodes (there might be multiple such trees).

¹<https://www.minizinc.org/software.html>

The tool does not take any command-line arguments. The problem instance is to be read from the standard input. Consider an instance file named `samples.smp`. The tool is expected to be executed as follows:

```
./proj3 < samples.smp > solution.txt
```

The tool must write the solution to the standard output, which can then be redirected to a file. The programming languages to be used are only C/C++, Java or Python. The formats of the files used by the tool are described below.

File Formats

You can assume that all input files follow the description provided in this document. There is no need to check if the input file is correct. Additionally, all lines (input or output) must terminate with the end-of-line character.

Input Format

The input file representing a problem instance is a text file that follows the following format:

- One line with one integer K defining the number of features.
- A sequence of lines where each line consists of $K + 1$ characters 0/1, separated by a single space.
- The first K characters on a line determine the values of features in the sample, and the last character determines its class.
- It is guaranteed that there are no conflicting samples, i.e., there are never two lines with the same values of features but different class.

Output Format

The output of the program representing an optimal solution to the problem instance must comply with the following format:

The program produces a binary tree in the following format consisting of 4 types of lines.

- “l i j ” representing that j is a left (value 0) child of i
- “r i j ” representing that j is a right (value 1) child of i
- “c i v ” leaf i responds with the class $v \in \{0, 1\}$
- “a r i ” the feature r is assigned to internal node i

The order of these lines is not important and any line that does not begin with either of the characters l, c, r, a, will be ignored. It is always assumed that node 1 is the root. The tree must be binary, i.e., a node either has no children or it has two (left, right). The tree must be

consistent with the given set of samples. In particular, for any sample s , there must not be a path in the tree consistent with the sample s , where the leaf at the end of the path responds with a different class than the sample s .

The tree must be smallest possible, i.e., there is no tree with fewer nodes that classifies the same set of samples S .

Important: The final version to be submitted for evaluation must comply with the described output. Project submissions that do not comply will be severely penalized, since each incorrect output will be considered as a wrong answer.

Example

The file describing the problem discussed above is described by the following input.

```
2
0 0 0
0 1 1
```

A solution corresponding to Figure 1 would be:

```
l 1 2
r 1 3
a 2 1
c 2 0
c 3 1
```

Another example with still 2 features but requiring 5 nodes is the following.

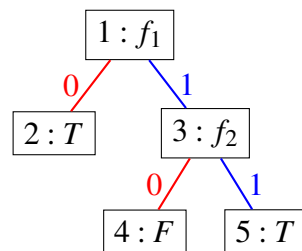
```
2
0 0 1
0 1 1
1 1 1
1 0 0
```

With a possible solution:

```
l 1 2
r 1 3
l 3 4
r 3 5
c 2 1
```

c 4 0
c 5 1
a 1 1
a 2 3

This tree can be visualized as in the following figure.



Additional Information

The project is to be implemented in groups of one or two students.

The project is to be submitted through the course website. Jointly with your source code, you should submit a short text file describing the main characteristics of your project.

The evaluation will be made taking into account correctness given a reasonable amount of CPU time (80%) and efficiency (20%). The input and output formats described in this document must be strictly followed.

Project Dates

- Project published: 16/11/2019.
- Project due: 29/11/2019 at 23:00.

Omissions & Errors

Any detected omissions or errors will be added to future versions of this document. Any required clarifications will be made available through the course's official website.

Versions

16/11/2019, version 1.0: Original version.

References

- [1] N. Narodytska, A. Ignatiev, F. Pereira, and J. Marques-Silva. Learning optimal decision trees with SAT. In J. Lang, editor, *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden.*, pages 1362–1368. ijcai.org, 2018.
- [2] J. R. Quinlan. Induction of decision trees. *Machine learning*, 1(1):81–106, 1986.