

1 Explanation Generation

Algorithm 1: EXPLAINCSP(\mathcal{C}, U, f, I)

```

input   :  $\mathcal{C}$  a CNF  $\mathcal{C}$  over a vocabulary  $V$ 
input   :  $U$  a user vocabulary  $U \subseteq V$ 
input   :  $f$ , a cost function  $f : 2^{\mathcal{G}} \rightarrow \mathbb{N}$  over a CNF  $\mathcal{G}$ 
input   :  $I$ , a partial interpretation over  $U$ 
output  :  $E$ , a sequence of explanation steps as implications  $I_{expl} \implies N_{expl}$ 
1  $E \leftarrow \langle \rangle$ 
2  $I_{end} \leftarrow \text{OPTIMALPROPAGATE}(\mathcal{C} \cup I, U)$  // assignment on variables of U
3 while  $I \neq I_{end}$  do
4    $X \leftarrow \text{BESTSTEP}(\mathcal{C}, f, I_{end}, I)$ 
5    $I_{best} \leftarrow I \cap X$ 
6    $N_{best} \leftarrow \text{OPTIMALPROPAGATE}(\mathcal{C} \cup I_{best}, U) \setminus I$ 
7   add  $\{I_{best} \implies N_{best}\}$  to  $E$ 
8    $I \leftarrow I \cup N_{best}$ 
9 end
10 return  $E$ 

```

Algorithm 2: OPTIMALPROPAGATE(SAT, \mathcal{U}, I)

```

input   : SAT, a SAT solver bootstrapped with a CNF.
input   :  $U$  a user vocabulary  $U \subseteq V$ 
optional:  $I$ , a set of assumption literals.
output  : The projection onto  $U$  of the intersection of all models of  $U$ 
1  $sat?, \mu \leftarrow \text{SAT}(I)$ 
2  $\mu \leftarrow \{x \mid x \in \mu : \text{var}(x) \in U\}$ 
3  $b_i \leftarrow$  a new blocking variable
4 while true do
5    $\mathcal{C} \leftarrow \mathcal{C} \wedge (\neg b_i \bigvee_{x \in \mu} \neg x)$ 
6    $sat?, \mu' \leftarrow \text{SAT}(I \wedge \{b_i\})$ 
7   if  $\neg sat?$  then
8     add clause  $(\neg b_i)$  to SAT solver
9     return  $\mu$ 
10  end
11  $\mu \leftarrow \mu \cap \{x' \mid x' \in \mu' : \text{var}(x') \in U\}$ 
12 end

```

Algorithm 3: BESTSTEP-C-OUS($\mathcal{C}, f, I_{end}, I$)

```

input   :  $\mathcal{C}$ , a CNF.
input   :  $f$ , a cost function  $f : 2^{\mathcal{G}} \rightarrow \mathbb{N}$  over CNF  $\mathcal{G}$ .
input   :  $I_{end}$ , the cautious consequence, the set of literals that hold in all models.
input   :  $I$ , a partial interpretation s.t.  $I \subseteq I_{end}$ .
output  : a single best explanation step
1  $A \leftarrow I \cup (\overline{I_{end}} \setminus \bar{I})$  // Optimal US is subset of A
2 set  $p \triangleq \sum_{l \in \overline{I_{end}}} l = 1$  i.e. exactly one of  $\overline{I_{end}}$  is present in the hitting set
3 return C-OUS( $\mathcal{C}, f, p, A$ )

```

Algorithm 4: c-OUS(\mathcal{C}, f, p, A)

```
input  :  $\mathcal{C}$ , a CNF.
input  :  $f$ , a cost function  $f : 2^{\mathcal{G}} \rightarrow \mathbb{N}$  over CNF  $\mathcal{G}$ .
input  :  $p$ , a predicate  $p : 2^{\mathcal{G}} \rightarrow \{t, f\}$  over CNF  $\mathcal{G}$ .
input  :  $A$ , a set of assumption literals, s.t.  $\mathcal{C} \cup A$  is unsatisfiable.
output : a  $p$ -constrained  $f$ -optimal unsatisfiable subset  $(p, f) - OUS$ .

1  $\mathcal{H} \leftarrow \emptyset$ 
2 while true do
3    $A' \leftarrow \text{CONDOPTHITTINGSET}(f, p, A, \mathcal{H})$ 
4   if  $\neg \text{SAT}(\mathcal{C} \cup A')$  then
5     return  $A'$ 
6   end
7    $A'' \leftarrow \text{GROW}(\mathcal{C}, f, p, A', A)$ 
8   Optional Grow, if the sat solver can provide a provide a good model, we can skip the expensive call
to the grow procedure. Needs to be checked experimentally!
9    $\mathcal{H} \leftarrow \mathcal{H} \cup \{A \setminus A''\}$ 
10   $\mathcal{H} \leftarrow \mathcal{H} \cup \{A \setminus A'\}$ 
11 end
```

2 MIP model

We define a set of user variables U defined over a vocabulary V of the CNF \mathcal{C} as $U \subseteq V$. Given an initial assignment I , where $\text{vars}(I) \subseteq U$, I_{end} is as the cautious consequence (the set of literals that hold in all models) of U .

The Mixed Integer Programming model for computing c-OUSes has many similarities with a set covering problem. The CONDOPTHITTINGSET computes the optimal hitting set over a p -constrained collection of weighted sets \mathcal{H} .

In practice, to ensure that MIP model takes advantage of the incrementality of the problem, namely across different c-OUS calls, the specification is defined on the full set of literals of I_{end} . The constrained optimal hitting set is described by

- $x_l = \{0, 1\}$ is a boolean decision variable if the literal is selected or not.
- $w_l = f(l)$ is the cost assigned to deriving the literal or using the derived literal (∞ otherwise).
- $c_{ij} = \{0, 1\}$ is 1 (0) if the literal l is (not) present in hitting set j .

$$\min_x \sum_{l \in I_{\text{end}} \cup \overline{I_{\text{end}}}} w_l \cdot x_l \quad (1)$$

$$\sum_{l \in I_{\text{end}} \cup \overline{I_{\text{end}}}} x_l \cdot w_{lj} \geq 1, \quad \forall j \in \{1..|hs|\} \quad (2)$$

$$\sum_{l \in \overline{I_{\text{end}}} \setminus \overline{I}} x_l \geq 1, \quad \forall j \in \{1..|hs|\} \quad (3)$$